

NEAS Time Series Project Documentation

PREFACE TO STUDENT PROJECTS ON TIME SERIES

(The attached PDF file has better formatting.)

Updated: February 9, 2006

This forum explains the student projects on time series.

- Some statistical techniques, such as regression analysis, have specified forms. The textbook gives formulas for ordinary least squares estimators and explains how to test hypotheses.
- Time series analysis is as much art as science. Several models may be reasonable, and selecting a model is subjective.

The postings have two guiding principles:

- They show the *range* of student projects that you can do. They give guidance on each type of project so that you can complete it easily.
- They provide project templates for interest rates and daily temperature, with Excel files of the needed data, so you can complete a reasonable project.

Several NEAS faculty members and advisors have contributed to these postings. We have revised some postings for consistency in the recommendations, but you may notice slight differences in the emphasis on one procedure or another.

- Some statisticians emphasize first differences; others explicitly detrend the time series.
- Some statisticians prefer simpler models; others choose the model with the best Box-Pierce Q statistic.
- Some statisticians prefer structural models; others use ARIMA models.

In general, the recommendations emphasize simple, structural models with explicit trend and seasonal adjustments. We emphasize intuition and parsimony over multi-parameter models. The project templates and other suggestions on this discussion forum will help you complete a statistical project that shows the rationale of ARIMA modeling and helps you understand the economic and mathematical foundations of common time series.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

STUDENT PROJECTS: INTRODUCTION

Updated: May 3, 2008

(The attached PDF file has better formatting.)

Jacob: What is the student project?

Rachel: The joint SOA-CAS VEE Administration Committee believes that book learning is only part of a good statistics education. Students in a statistics course must learn to apply the statistical concepts and techniques to real data using statistical software.

Jacob: Is this a project that we design or a project that NEAS designs and we complete?

Rachel: The ideal project is designed by the student. The project can be on any topic.

Jacob: Can we use a project we did at work or in college?

Rachel: You can use a statistical project completed for work as the student project for the course. It should be modified to emphasize the application of the statistical concepts. Your work projects focus on actuarial topics, such as setting rates or estimating reserves. The student project focuses on statistical topics.

Jacob: Designing a new project is difficult. If we don't normally use statistical techniques, how do we get ideas and data for projects? What constitutes a good project?

Rachel: We provide project templates. Some project templates are actuarial applications of statistical concepts to insurance or finance. Other project templates are designed for the VEE requirements.

We give you data sets and explain the procedures to apply. You can use Excel for all the project templates. We provide illustrative spread-sheets with code for statistical techniques that are not available as built-in functions or add-ins.

We focus on the concepts taught in the on-line courses; you don't have to learn other statistical concepts. You can use any statistical software you have; you are not restricted to Excel.

Take heed: Use the project templates as starting points.

- When you first read Chapter 19 of the time series course, or you first use F tests in a regression analysis., the procedures are hard to grasp.
- After reading through several project templates, you will pick up the concepts. Once you understand the student project requirements, the work is straight-forward.

Jacob: What aspects of the student project differ from the homework and final exam?

NEAS Time Series Project Documentation

Rachel: The homework assignments and final exams do not examine real data. Actuarial and statistical work differ in several ways:

- Actuaries use formulas and come to definite answers. A pricing actuary uses a set of procedures and derives a premium rate.
- Statisticians focus on charts, graphs, and plots. The results are stochastic, so the plots are often ambiguous and no single answer is necessarily correct.

These generalizations are often true, though they are not always true.

Jacob: What project templates do we have?

Rachel: We give several project templates on the discussion forum:

- Regression analysis of loss reserving: dummy variables and squares of variables.
- Sports won-loss records: F tests
- ARIMA modeling (interest rates, inflation, daily temperature, and other time series)

Each project template provides data and ideas for dozens of student projects.

Jacob: Do we have other project templates?

Rachel: We add more project templates each semester. The project templates are hard to form: we must ensure that candidates can complete the project, but candidates must do significant work on their own. A poorly constructed project template annoys candidates. Our faculty spends several months on each project template, to ensure that it meets the requirements for VEE credit and can be completed by all candidates who have taken the on-line courses.

- For regression analysis, we have two project templates: dummy variables for regression analysis of loss reserves and F tests for sports won-loss records. The papers from which the projects are drawn are on the actuarial syllabus for CAS Exams 6 and 9 and the statistical skills are much desired by insurers.
- For time series, we use a project template on interest rates. Interest rates are much discussed in the textbook, and they are used by most candidates. The SOA and CAS syllabus for their investment exams have much material on interest rate analysis.

Jacob: Do we have our choice of student project? Can we pick any project template?

Rachel: You have your choice of any project template or any other student project.

Jacob: Can we discuss the project templates on the discussion forum?

Rachel: We encourage discussion. For some items, discussion is the quickest way to learn. You must use the Excel built-in functions to complete the regression analysis and form the residual plots. Some candidates have no problem with Excel; others are less familiar with it. We give guidance in the project templates, but we can't anticipate some problems.

NEAS Time Series Project Documentation

Jacob: How is the student project submitted to NEAS?

Rachel: The student project should be emailed to NEAS as an electronic file. A text portion explains what you have done. The text can be in Microsoft Word, Corel WordPerfect, a text file, or a PDF file. The project template shows what questions you must answer. You have much lee-way; you can substitute a related topic if you want.

An Excel or similar attachment contains the statistical output and the graphs or plots: residual plots for regression analysis and sample autocorrelation functions for time series. In most cases, the built-in functions in Excel or other software forms the plots.

Take heed: The discussion forum explains the contents of the write-up.

- If you have written college papers, the student project is no more complex.
- If you avoided papers during college, read the requirements on the discussion forum. We explain the minimum requirements; you can use any writing style you like.

Jacob: Why do you want electronic files?

Rachel: We must send the files to our faculty members to review each project. Paper copies are harder to keep track of.

Jacob: What if we get stumped and can't figure out how to read a graph?

Rachel: Post your question on this discussion forum. Many statistical procedures are confusing at first, and probably other candidates are equally stumped. We will review the questions and post answers.

Jacob: What if the student project is not satisfactory?

Rachel: We will send you an email explaining what else needs to be done.

Jacob: How long will these student projects take?

Rachel: That depends on the candidate. If you understand the material, the student project is not hard. The project template explains what you do and what decisions you make. If you do not understand the statistics material and can not apply the methods to data, it takes time to work through each step of the project template.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

STUDENT PROJECT: OBJECTIVES AND DATA

Updated: December 26, 2006

(The attached PDF file has better formatting.)

Jacob: What is the objective of the independent student project? How does it differ from the final exam and homework assignments? Do we use regression analysis to test hypotheses? Do we apply the regression concepts to actuarial scenarios?

Rachel: The final exam tests if you understand how to use the regression techniques. The final exam may ask you to calculate t statistics or F statistics or Durbin-Watson statistic or use them to evaluate regression results. The homework assignments apply the concepts to insurance and actuarial scenarios.

Practical regression analyses are more complex. Classical regression analysis assumes the regression parameters (the α and the β 's) are constant and the error term is normally distributed with a constant variance. The time series chapters assume an ARIMA process works well, though we may have to take first or second differences.

The student project addresses real problems. One project template uses regression analysis to estimate loss reserve. If the regression parameters are stable, anyone could do reserve analysis; we would not need actuaries. We show how to analyze scenarios where parameters are not constant.

Another project template uses an F test to compare sports teams. The student project has no correct answer. The results depend on the sport, teams, years, and hypotheses you choose.

The project template on interest rates uses actual inflation and interest rate series. The proper ARIMA process depends on the series, time period, and adjustments you make.

Our faculty reviewing your student project is not looking for a specific result. They check if you apply the statistical techniques reasonably well to actual or simulated data. Two candidates may examine the same issue but use different data and come to opposite conclusions.

DATA

Jacob: Do we choose the topic for the student project, or do we work on assigned data?

Rachel: Ideally, candidates work on topics that interest them. An auto actuary may do a project on territorial rate relativities and a life insurance actuary may do a project on mortality differences by sex. You can design your own project.

NEAS Time Series Project Documentation

Many candidates prefer more structured projects. To ensure the projects are independent, we use one of two methods:

We supply extensive data. For the project template on interest rates, we supply many time series over the past hundred years, and we show how to form combinations of these. You can choose any of the time series on the NEAS web site, and you can form others, such as real interest rates from nominal interest rates and inflation rates. You choose the time periods and the questions to examine. You can use similar time series from dozens of public web sites, and we provide suggestions in many postings.

For the project template on sports won-loss records, we supply statistics for all the teams in four sports: baseball, basketball, hockey, and football. We suggest numerous analysis that you can perform. You select the sport, teams, years, and hypotheses.

For the project template on loss reserving, we suggest two analyses you can do. You simulate data, so every candidate uses different figures. We provide illustrative spread-sheets and a step-by-step guide to all parts of the project template.

The project template gives you a format: the type of data, the type of question, and the type of analysis. You perform the regression analysis to answer the questions.

Jacob: With all this material, it takes hours just to examine the data.

Rachel: Some candidates can not start a project until they see all the possibilities. The project is open-ended: you choose data, hypotheses, and statistical techniques. You show that you can apply statistical techniques to real (or simulated) data to test hypotheses.

The project templates are suggestions; the data on the NEAS web site eliminates time wasted gathering data on your own. The illustrative spread-sheets give you Excel code for the common statistical techniques.

Some candidates ignore all this and submit their own projects. That is perfect. Other candidates want more structured assignments, and they use the project templates. Don't examine all the data; that takes too much time. If you are a basketball fan, do a student project on basketball statistics. Use your home team and replicate the illustrative spread-sheet on the NEAS web site.

NEAS Time Series Project Documentation

STUDENT PROJECT DOCUMENTATION

(The attached PDF file has better formatting.)

The student project has two parts:

- A statistical workbook (Excel spread-sheets, SAS, MINITAB)
- A written document (Word, WordPerfect, PDF)

Many candidates use Excel for the workbook and Word for the text document. You may use any software; you are not restricted to Excel and Word.

- ~ You can use any statistical software package or spread-sheet. Copy the primary output to your Word, WordPerfect, or PDF file and attach the workbook to your email.
- ~ You can use any text document: Word, WordPerfect, PDF or a TXT file. We grade the content of the file, not its format.

Your documentation is in the text file (Word, PDF, WordPerfect). Copy the graphs and exhibits that support your conclusions to the document, and include all statistical results, such as the parameters of your regression equation or ARIMA model, goodness-of-fit tests (F test, p value, R^2 , \bar{R}^2 , Durbin-Watson statistic, Box-Pierce Q statistic, Bartlett's test), and comparison of forecasts to actual values. Refer to other exhibits or regression output in the workbook, and specify where they are found (work-sheet name, tab).

The written document is essential. Many candidates can run the *REGRESSION* add-in and graph the data even without taking a statistics course. The student project shows that you understand what you are doing. Explain what the workbook contains, the logic from the initial hypothesis to the conclusions, why you chose specific statistical tests, significance levels, and procedures, and the implications of your project.

Our faculty does not compare your conclusions with their own views. We check if you use the statistical techniques appropriately and if you understand the material.

Take heed: The course instructor reads the written document, which should include all results of your project, such as

- “The fitted regression line is ...; the R^2 is Z%, indicating that ...; the p -values for the explanatory variables are ... indicating that ...”
- “The fitted ARIMA process is ...; the Box-Pierce Q statistic is ...; with T observations in the time series and K lags, the χ -squared value for a 10% significance level is ...”

If your word processor allows, copy relevant graphs to the document. Form the residual plot or correlogram in Excel (or other software) and copy the graph to Word (or other word processor). Explain in the document what the graph shows, such as

NEAS Time Series Project Documentation

- “The residual plot, showing the mean residual at each calendar year, is shaped as a V, indicating that ...”
- “The correlogram slopes downward slowly. The sample autocorrelation function is positive for the first 20 lags and negative for the next 20 lags, indicating a non-stationary time series with a change in the mean after 20 periods ...”

The written document is a report, not just documentation of your workbook. It should include the topics listed in this posting.

- Do not say: “The course instructor can read the Excel workbook to see what I did.”
- The course instructor examines if you understand the statistical analysis. The project templates, the data on the NEAS web site, and Excel built-in functions enable you to put together an Excel file. The student project shows if you understand the analysis.

The workbook (or other statistical file) has the supporting spread-sheets and exhibits. The written document may say

- The data are from the web site www.datasource.com.
- They are shown as a matrix of dates and values in Sheet 1 in the Excel workbook.
- The regression of (Y variable) on (X variables) uses the *REGRESSION* add-in.
- Residual plots on all the explanatory variables are shown in Charts 1, 2, and 3.
- Correlograms of the time series and its first and second differences are in Charts ...

Statistical analyses that you do not need for the final report should be referenced in the text document, with the exhibits in the workbook.

Illustration: “I examined monthly birth rates for seasonality several ways. Sheet S1 shows average birth rates by month, and Sheet S2 shows the average month-to-month change. No month has significantly different rates. Sheet S3 shows the sample autocorrelations for lags 1 through 24. Lags 12 and 24 are not materially higher than other lags. Sheet S4 shows an ARMA model with an AR(12) coefficient. The p -value for this coefficient is 55%, indicating that it is not significant.”

If you use the statistical analysis, copy the graphs into your text document.

Illustration: “I examined monthly marriage rates for seasonality several ways. Sheet S1 shows average marriage rates by month, and Sheet S2 shows the average month-to-month change. The graphs are reproduced here as Table T1 and T2. June has higher marriage rates than other months. Sheet S3 shows the correlogram for lags 1 through 24, reproduced below as Figure F1. Lags 12 and 24 have high autocorrelations. Sheet S4 shows an ARMA model with an AR(12) coefficient. The R^2 , t values, and p -values are reproduced in Table T4. The p -value for the AR(12) coefficient is 5%, indicating that it is significant.”

Some candidates have writers' block. They do the Excel analysis easily, but they can't write up the results. Don't worry about the quality of your writing.

NEAS Time Series Project Documentation

- We do *not* grade you on writing style, exposition, or grammar.
- For many candidates, English is a second language; we don't expect flawless style.

Take heed: This posting is a guide for candidates who can do the statistical techniques but have trouble with written work. Many student projects submitted in past semesters are excellent: some are humorous; some show keen insight; some are innovative and much impressed our faculty. This outline is *not* a constraint on your write-up.

You have written memos at work. Your supervisor may ask you to change the style, rewrite the memo, make it clearer, expand the reasoning, or shorten the text. The student project is not graded on writing style. We check if

- you understand how to apply the statistical techniques to real data.
- you explain what you have done and how it supports the conclusions.

NEAS Time Series Project Documentation

OUTLINE

We give a *minimum* outline for the written document (write-up). The elements below are recommendations, not requirements. You may write the project in any style you like.

- Don't restrict yourself to this outline. Every student project differs, and your document may have a different structure. Don't force your writing style into the format here.
- Read this outline to see what is expected, and write your student project as you write other memos. If you have trouble with the style, say to yourself: "They are not grading me on the writing style." Explain what you have done and submit the project.

The minimum write-up has the following sections.

Introduction: State the problem and how you deal with it. Use statistical terms, but explain in words what you do.

Illustration: For a project on non-constant regression coefficients, you might write:

Classical regression analysis assumes the regression coefficients are constant. If the regression coefficients change over the range of explanatory variables, the R^2 , t statistic, p -value, and F statistic do not show this. I use residual plots to identify changes in the regression coefficients. To adjust for these changes, I use dummy variables (or the square of an explanatory variable).

For project on modeling interest rates with ARIMA processes, you might write:

Financial economists assume that nominal interest rates move with inflation rates, but real interest rates are more stable. I regressed the three month Treasury bill rate on the CPI and examined how to model the residuals. I used first differences and a seasonal adjustment to create a stationary time series, which I modeled with AR(1), AR(2), and MA(1) processes. I used Bartlett's test, the Box-Pierce Q statistic, and the principle of parsimony to select the AR(1) model.

These are sample paragraphs. Say what your student project does, so that the faculty member reviewing it knows what to expect. Feel free to use any format.

- One candidate used a Jacob/Rachel dialogue – except that Rachel had the questions and Jacob had the answers.
- Many candidates use personal introductions, explaining their interest in the topic, such as "My parents were both divorced from previous marriages, but their present marriage works well. I have often wondered whether divorce rates differ for first vs second marriages or by age at marriage."

Data: You gain most from the student project if you use data that interest you.

NEAS Time Series Project Documentation

- The NEAS web site has many data sets that you can use for the project.
- The world wide web has thousands of web sites with good data. Choose a topic and search the web for data.

If you use data not from the NEAS web site, describe the data. If the student project is a time series, define the series, the period, and the intervals. You might say:

- The time series is the French unemployment rate from 1980 through 2003 at quarterly intervals from the *XYZ.gov* web site.
- The time series is daily temperature in Puerto Rico for 2000-2005, from the *XYZ.com* web site. I adjusted for seasonality by dividing by the average daily temperatures from 1950 through 2005.
- The data are (i) the monthly crime rate in Los Angeles, (ii) immigration rates and gang activity in Los Angeles, and (iii) unemployment rates in Los Angeles. These figures are published on *XYZ.com*.
- I used female mortality rates for ages 35 through 65 based on the *XYZ* mortality table.

Any data set is fine. Do not worry that your data will not show an ARIMA process or will not support your statistical hypothesis. But:

- Use enough data points to perform the statistical techniques. If you have only ten points, you can not fit an ARIMA process and your regression will not be significant.
- Do not choose data that are white noise or random walks.

Illustration: A time series of earthquake frequency may be white noise. If the correlogram shows no significant autocorrelations, choose a different data set.

Take heed: Do not send an email to NEAS asking if you have enough data points. If a topic interests you and has only 50 data points, use it. If the NEAS web site has a larger data set for the same topic, use the larger data set.

Use appropriate periods. For a time series project on outside temperature, use daily readings (or hourly readings), not monthly averages.

Company data: If you use your own company's data, disguise the data by scaling the figures. Multiply the data by a factor or add a constant. The adjustment may cause different regression coefficients or ARIMA parameters; that is fine. You might write:

- I used a paid loss triangle for 12 years. I multiplied the figures by a factor, subtracted a constant, and added a random draw from a normal distribution with a mean of zero.

If you use data from a project template on the NEAS web site, state which data you chose. You might write:

- I used baseball won-loss records for the National League for 1911 through 1960.
- I simulated paid loss triangles for a 15 by 15 array.

NEAS Time Series Project Documentation

- I used overnight LIBOR rates from 1945 through 1975.

If you use outside data, explain its characteristics.

Illustration: If you regress mortality rates on age and sex, say what table the data are from, whether they are population data or insured data, and what years they are from.

If you simulate data, describe the simulation parameters (the α , the β 's, and the σ). You can simulate data of various types for a student project. The simulated paid loss triangles for the project templates on regression analysis for loss reserving are an example. Mahler's *Guide to Regression* has many simulations for actuarial subjects.

Topic: A good student project has a limited focus. Choose a topic, formulate a question, apply the statistical techniques covered in the course, and answer the question.

Illustration: For the project template on regression analysis of loss reserving, you forecast future loss payments. The student project has a specific topic, such as:

- ~ The inflation rate or the development trend may not be constant.
- ~ The variance of the error term may not be constant, but may vary with the calendar year or the development period.

Simulate data with a discrete change in the inflation rate or payment pattern, or a continuous change in the inflation rate or payment pattern, or a variance of the error term that depends on the calendar year or development period.

Show the R^2 , t statistic, p -value, and F statistic assuming a constant inflation rate, payment pattern, or variance. Explain why these values indicate the regression is significant even though the assumption is not correct. The explanation need not be long. You might write:

These statistical tests help decide whether to reject a null hypothesis if the assumptions of classical regression analysis are true. They do not test whether the assumptions are true. If the inflation rate changes from 5% in one part of the sample to 15% in another part, and the null hypothesis is that inflation is zero, the estimated inflation rate may be 10%. The statistical tests reject the null hypothesis that inflation is a constant 0%. We should assume inflation is 10% only if it is constant over the entire sample.

For the time series student project, state the topic clearly. Do not just say "the project deals with interest rates." Your project might

- Fit an ARIMA process to a given time series to see if an autoregressive or moving average process explains the observed values.
- Compare two periods to see if the same ARIMA process is appropriate before and after a critical event.
- Fit an ARIMA process to the residuals of a regression analysis.

NEAS Time Series Project Documentation

You might write:

- I fit ARIMA processes to time series of nominal interest rates and real interest rates, or nominal interest rates divided by the expected inflation rate. Both processes can be fit with ARMA(1,1) processes, but the fit works well only for real interest rates.
- I compare the interest rate time series for two periods to see if the same ARIMA process is appropriate. The rising interest rates in the stagflation era of the 1970's require logarithms and first differences to form a stationary time series; the stable interest rates of the 2000's can be fit with an MA(1) process.

Statistical techniques: State the statistical techniques you use and explain how they relate to the problem in your project. Explain the results of your analysis and the implications.

Illustration: For an F test, explain what the unrestricted and restricted equations are, so that a lay reader would understand your project.

- Write the equations for the unrestricted and restricted regression lines.
- State the null hypothesis, relating it to the restricted regression equation.
- State the degrees of freedom for the F test.
- State the result of the F test, the critical values, and the significance level.

You might write: "I use an F test to see if the two Leagues have the same relation of past experience to future won-loss records. The null hypothesis is that the different relations in the previous section of this project reflect fluctuations, not true differences. The result is significant at 10% level, for which the critical value is ..., but not at the 5% level, for which the critical value is"

Do not just say: "I used an F test to see if two samples are from the same distribution."

Take heed: Explain your charts. Label the indices, and explain the index values.

- If you analyze monthly interest rates over 30 years, the horizontal axis may be a month from 1 to 360. If the interest rates peak at month 125, identify the date of this month.
- If you form residual plots, specify the axes. The horizontal axis is one of the independent variables or the dependent variable; the vertical axis is the residual.

Say what you look at in a graph. For a residual plot, you might write:

- ~ To test if a regression coefficient is constant, I examine the slope of the line connecting the average residuals.
- ~ For conditional heteroscedasticity, I examine the spread of the residuals.

Copy Excel graphs to your write-up and label the axes, so the course instructor follows your arguments. If you leave the graphs and charts in the Excel workbook, it is difficult to follow your reasoning. A course instructor who does not know what graph or chart you refer

NEAS Time Series Project Documentation

to or who does not know what the axes represent may ask for better documentation. It may take an extra 6 to 8 weeks to get VEE credit approved by the SOA.

Explain what the statistical test implies.

- ~ If the residual plot looks like an upside-down V, explain what this means. If you get stuck writing an abstract explanation, give an example.
- ~ If the residuals are more spread out for high values of X than for low values, explain what this means.

Illustration: The residual plot appears like an upside-down V. If the actual inflation rate is declining / rising... (explain the effects). For example, if the actual inflation rate is Z in the first 5 years and Z' in the next five years ...

Corrections and Adjustments: Explain how you correct problems in your data. For a project on non-constant regression coefficients, explain how dummy variables, squares of explanatory variables, or other methods correct for the problem. Show the results of the revised regression equation and explain why it is superior.

Take heed: Not all student projects have corrections. Your project may use an F test to compare two sports teams, or your project may fit an AR(2) process to a time series and conclude that it fits well. But most projects have adjustments.

Many projects have a series of statistical tests or a series of tests and corrections. Fitting ARIMA models uses a sequence of procedures. You first graph the time series, looking at means, trends, drifts, variances, cycles, seasonality. The first corrections are

- Breaking the time series into periods
- Using seasonally adjusted data
- Using structural models, such as real interest rates instead of nominal interest rates

Your student project is a series of adjustments and corrections. You

- Form a correlogram and check if the time series is stationary.
- Compute moving averages to identify trends.
- Adjust by taking differences or logarithms and differences.
- Compute monthly averages or 12 month autocorrelations to identify seasonality.
- Adjust the data to offset seasonality.
- Fit ARIMA models, based on the sample autocorrelation function.

For each step, you may form charts or graphs. Copy the charts or graphs into your Word document and explain how you formed them and what they imply. Don't explain the Excel code, but say what you did.

Illustration: I checked for seasonality two ways:

NEAS Time Series Project Documentation

- I formed monthly averages in three steps:
 - Compute the average monthly sales in each year
 - Divide each month's sales by the average monthly sales to get monthly relativities
 - Take the straight average of the monthly relativities over 20 years
 - The column chart by month (Figure 1) shows high sales in May, June, and July and low sales in December, January, and February.
- I examined the sample autocorrelations for 6, 12, and 24 months:
 - The 12 and 24 months sample autocorrelations are high
 - The 6 month sample autocorrelation is negative
 - The correlogram is Figure 2, with markers at 6, 12, and 24 months

Do not write: "Everything is in the Excel work-sheets." Figuring out what you did from an Excel worksheet is hard. If the work-sheet is not well documented, it may be impossible.

Copy graph, charts, and conclusions table from the work-sheet into your document.

- Copy correlograms and plots into your document and put them next to your comments.
- If you have trouble with cut and paste from Excel into Word, or if you use a text file that does not accept pictures, refer to the Excel chart or graphic.

Don't copy long Excel tables into the written document.

- Refer to the work-sheet or region in your written document.
- Cite the important results, such as an F statistic or a Box-Pierce Q statistic.

STOCHASTICITY AND SIMULATION

Some student projects simulate data. If you simulate data, turn off the stochasticity on your first run. You might write:

- I ran the simulation first with $\sigma = \text{zero}$ to verify that I get the expected results. The regression equation is ..., with an R^2 of 100%.
- I then used $\sigma = 0.01$ to verify that I am adding the error term correctly. The regression equation is ..., with an R^2 of 95%.
- I then used $\sigma = 0.2$ for the student project."

PROJECT LENGTH (TIME AND PAGES)

The student project is a serious task. To cover the topics in this posting, you need several pages of text plus charts, graphs, and tables. The project has no required length, but do not leave out critical sections.

Take heed: If you are unsure whether to include a section, put it in. Instead of saying: "I did not find any seasonality in the time series," write: "I computed monthly averages, as shown in Chart 1, which do not show seasonality." The data are yields on 20 year Treasury bonds, so I did not expect seasonality.

NEAS Time Series Project Documentation

The time needed for the student project varies with your grasp of the statistical concepts, your motivation, and past work with these topics.

- If you understand the concepts, the student project does not take long to complete.
- If you don't understand the concepts, you might spend days wondering what to do.

TOPICS

If you have never before written statistical reports and you are not familiar with internet search engines, you may find the project difficult at first. If you feel lost, use one of the project templates and follow the instructions.

Illustration: You want to do a student project on crime rates: perhaps a time series analysis of murder rates over the past 20 years or a regression analysis project relating crime rates to city governance or city size.

The NEAS web site has no data on crime rates. Use internet search engines (Google, Yahoo, MSN) to find data on crime rates. A few minutes of looking at sites leads you to the FBI home page. Use the site map to find a section called *statistics*, which brings up hundreds of files of crime statistics. You can download many of these files in Excel format. You may spend an hour looking at the different files and picking data you want to analyze.

If you use statistical analysis at work, you can mold your work into a student project.

Illustration: You fit exponential curves to average claim severities to project loss cost trends. Use the same data for a time series student project, with several additions:

- Take logarithms and first differences to make the series stationary.
- Fit an ARIMA process, not an exponential curve.
- Examine seasonality. Average claim severities vary by quarter in most lines.
- Use a structural model. Deflate average claim severities and apply the ARIMA process to the claim severities in real dollars.

The student project will not take long, since you already have the data in Excel files, you know the data characteristics, and you have already worked with these figures. Most insurers do not adjust trend data for seasonality and use very simple exponential trends.

Your analysis should improve the trend projections. You receive VEE credit for the on-line course and time series analysis may be used in your pricing studies.

If you feel lost, review the NEAS project templates and the past student projects on the discussion forum. Project templates for sports scores, interest rates, daily temperature, and other topics are on the discussion forum. The NEAS web site has hundreds of data sets in Excel format, with full explanations of the types of analyses you can do. The postings for the project templates have several forms:

NEAS Time Series Project Documentation

- Step-by-step guides for the more common analyses.
- Jacob / Rachel dialogues for the hypotheses and analyses you can use.
- Discussions of topics used for past student projects. These discussions are broad, with many suggestions. You are not expected to do everything; pick a topic to focus on.

Look at the past student projects on the discussion forum. The discussion forum has the write-ups, not the work-sheets. The past projects give you ideas, which you implement in our own work-sheets. Your student project need not be new. If a past student project interests you, choose other data and do a similar project.

Illustration: Many candidates analyze won-loss records of their home teams. The material on the web site gives dozens of ways to write a student project on sports figures. If you want to do a sports project, read the project templates, see what other candidates have done, select data from the NEAS web site or from other internet sites, and do the analysis.

DISCUSSION FORUMS

The student projects are independent, but you get ideas by discussion with others. Feel free to discuss the student project on the discussion forum. Some items are confusing at first, but they become clear after a bit of discussion.

The Excel built-in functions are hard if you have never used them. Discuss how to use the Excel built-in functions like *RAND* and *NORMSINV*, Excel names, the *REGRESSION* add-in, and the *SOLVER* add-in. The student project does not test Excel; feel free to copy material from any Excel file on the web site.

Use your own data for the student project. Data from another web site are the best. If you choose data from the NEAS web site, use a different time period for a time series or different parameters for a simulation.

Your student project may be similar to the illustration on the web site. If you perform the statistical techniques and apply them to the data you choose, that demonstrates you can use the techniques.

If you are completely lost, begin by reproducing the analysis in a project template. Once you understand what is expected, choose a different set of data and do a similar project.

JOINT PROJECTS

Two or more candidates from one firm may be taking the same on-line course. You can discuss the project with another candidate at your firm. You might both do ARIMA modeling of interest rates or analysis of sports won-loss records. But your projects and your write-ups must be separate.

- You should not have a single project with two authors.
- You should not have the same project with different parameters.

NEAS Time Series Project Documentation

The SOA / CAS wants each candidate to write a student project.

DUE DATES

Actuarial candidates are busy with work, courses, student projects, and study for actuarial exams. Some candidates have young children and many responsibilities.

A candidate might take an on-line course in January/February, study for an actuarial exam in March and April, and do the student project in May and June. Similarly, a candidate might take an on-line course in July/August, study for an actuarial exam in September and October, and do the student project in November and December.

We don't set due dates. You are responsible for effective use of your time. But be mindful of your future work schedule. Some candidates presume they have more time after their exam. But after your exam, you will be thinking about the next exam. If you are like most actuaries, your work schedule keeps getting busier until you retire.

We recommend: Right after the final exam, do the student project. You will finish it quickly when the course material is fresh. Getting VEE credit gives enormous motivation for exam study. A project that takes a few days when you know the material well may drag on for several weeks if you keep postponing it.

QUESTIONS

Questions about the student projects are of several types. For quicker response, send questions to the proper person(s).

Send administrative questions (by email) to the NEAS office: to whom should you send the project, when is the project due, and so forth. (Send the project to the NEAS office; the project has no due date, but you don't receive credit for the course until the project is received and graded.)

Questions about statistical techniques and Excel built-in functions should be posted on the discussion forum. The illustrative work-sheets and step-by-step guides show to use many techniques and built-in functions needed for the student projects. Additional questions may be posted on the discussion forums, such as "Does Excel have a built-in function for generalized linear modeling? What is a reasonable confidence interval for Bartlett's test?"

The NEAS faculty compiles questions into dialogues and general postings. The common questions about student projects are discussed on the discussion forum. Review the forum threads, and you may find your answer.

Some candidates have questions about their project results.

- The correlogram for a time series declines slowly; is the time series stationary?

NEAS Time Series Project Documentation

- The p value for an explanatory variable is 11%; can I use it in the regression analysis?

The student project shows how you deal with these questions. Statistical analyses are subjective because the data are stochastic.

Answer each question based on your understanding of the statistical methods. Explain why you chose a specific answer, and how else you might answer the question. The student project shows if you can deal with real statistical issues. We grade the project by your reasoning, not by whether the answer is right or wrong.

Illustration: A correlogram declines slowly. You might write:

- The correlogram declines to about zero after 8 lags, suggesting an autoregressive process with several parameters.
- The correlogram does not decline rapidly enough for a stationary process. The correlogram of the first differences declines to zero by the second lag.

Some items that occur frequently in actuarial and financial time series are not covered well in the textbook. Candidates repeatedly find correlograms that decline slowly, are positive for the first K lags, and then negative for the next K lags. This correlogram indicates two phases of a time series with different means. We explain the implications of this correlogram in a separate discussion forum posting. Read the postings on the discussion forums; they answer many of the questions on student projects.

NEAS Time Series Project Documentation

GRADING OF STUDENT PROJECTS

Updated: May 2, 2008

Jacob: How are the student projects graded? Are they number graded, letter graded, or pass/fail?

Rachel: The NEAS statistics faculty give a *letter grade* to each student project.

- A project must receive a B- or better for the candidate to receive VEE credit for the course. If your student project receives a B- grade or better, the NEAS office sends you an official transcript. The NEAS office also sends a transcript to the SOA/CAS office that grants VEE credit. You must still fill out a form requesting credit. The form (and instructions) are on the SOA web site.
- Projects that receive less than a B- grade are sent back with comments about what needs to be improved.

Jacob: What are the criteria used by the NEAS faculty to grade the projects? Does a comment mean that the project has an error that lowers the grade below B-?

Rachel: That depends on the type of comment.

- ~ Some comments refer to choices; others correct a mis-understanding of the technique.
- ~ Some statistical techniques are not clearly explained in the textbook. As we review the student projects, we explain the more common errors on the discussion board. An error in using the partial autocorrelation function for ARIMA modeling or in using generalized linear modeling to model automobile insurance rate relativities does not prevent VEE credit. We want you to use statistical techniques, even if they are not well explained in the textbook.
- ~ The grade depends on the overall quality of the student project. The student project need not be perfect. If you set up the statistical test correctly but you have an error in the degrees of freedom, your student project still satisfies the course requirements.
- ~ The grade depends on whether we have discussed this error in the course modules or the discussion threads. Some errors are common and are fully explained on the discussion forum. Using one ARIMA process to model a time series with different attributes in different periods is a common error.

Jacob: What do you mean by choice vs mis-understanding?

Rachel: Some errors are choices that are not ideal. One candidate with a time series of 132 observations used 120 lags for the Box-Pierce Q statistic, presumably because this is the highest degrees of freedom in the table at the back of the textbook. This is not the ideal choice: the last 20% to 30% of the lags are not a good test of a white noise process, since the sample autocorrelation formula may cause low values at high lags. Ideally, we use between 15 and 40 lags for a time series of 132 observations. This is a minor item in

NEAS Time Series Project Documentation

the analysis. If the candidate uses the statistical techniques correctly, we ignore the non-optimal choice of K .

Other errors are more serious. One candidate thought that an AR(2) model regresses the current time series value on the value two periods back instead of on *two* previous values (1 period and 2 periods back). The candidate compared an AR(1) model with an AR(2) model (with just the two period back value) and found the AR(1) model to be superior. This misunderstands autoregressive processes.

Another candidate thought the AR(1) model regresses the current time series value y_t on the time period t , not on the lagged value y_{t-1} . The candidate found that the residuals from this model were not a white noise process; that is, the regression had serial correlation. This would be true for any random walk or autoregressive process; it is a serious error.

Jacob: What do you mean by the *textbook explanation*?

Rachel: The textbook has a good explanation of the sample autocorrelation function but a poor explanation of the partial autocorrelation function. Candidates using Excel have no easy way of forming the partial autocorrelation function. We expect student projects to consider the sample autocorrelations and compare them to the theoretical autocorrelations from the ARIMA process. We do not expect the student projects to make full use of the partial autocorrelation function. (Some student projects using Minitab have made good use of the partial autocorrelation function.)

The most common error in the student projects is the interpretation of the correlogram. The textbook says that a stationary model has a correlogram that declines fairly rapidly to zero, but it does not specify the exact meaning of *declines fairly rapidly to zero*.

Some candidates see a correlogram that declines slowly to zero over its full length, or a correlogram that declines to zero, becomes negative, and then rises slowly back to zero, and assume that these correlogram indicate a stationary series. They do not (except in rare cases), but it is hard to give precise rules. As long as the candidate also examines the first differences, this error does not require a revision of the student project. Several posted comments by the NEAS faculty discuss when a correlogram indicates stationarity.

The second most common error is taking too many differences or using too high an order for the ARIMA process. The textbook often examines both first and second differences. If the first differences form a stationary process, we construct an ARIMA process on the first differences, even if the second differences have lower sample autocorrelations. Several postings discuss the errors in taking too many differences.

The proper order of an ARIMA model is disputed by statisticians. If an ARIMA(2,1,0) model has residuals that pass Bartlett's test and give a reasonably low Box-Pierce Q statistic, we don't use an ARIMA(12,2,8) model that performs slightly better. But the choice between ARIMA(2,1,0) and ARIMA(2,1,1) is less clear. The optimal model depends not only on the goodness-of-fit tests but also on the intuitive reasonableness of each model. Candidates using Minitab may form too many ARIMA models and not examine their reasonableness.

NEAS Time Series Project Documentation

Jacob: What do you mean by the *overall quality of the student project*?

Rachel: Student projects must receive a grade of B- or higher for VEE credit. We encourage candidates to work on innovative projects. Designing an innovative project is hard, and few designs are error-free. Actuarial candidates have many responsibilities and little spare time. We do not require all student projects to receive A's.

Jacob: What do you mean by whether a topic is discussed on the discussion board?

Rachel: As we review the student projects, we explain the most common errors on the discussion forum. We don't expect candidates to read every posting. But once a topic has been discussed and the proper method is clear, we may direct candidates to the discussion forum postings for more guidance. Good independent work is difficult, and we don't expect candidates to produce perfect projects without extensive guidance. But explaining all the statistical techniques clearly enough that candidates can apply them in their projects takes is equally difficult. It may take another few months before the postings are complete.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

TIME SERIES: INDEPENDENT STUDENT PROJECTS

Updated: May 3, 2008

This posting explains the independent student projects for the time series course. It uses examples from the interest rate time series on the discussion forum. Many suggestions here are repeated in other postings.

Take heed: The student project is not a linear assignment, starting at Point A and ending at Point Z. Some candidates grasp the requirement and can start work; others need more time to understand what they must do. You don't have to read everything on the discussion forum. Select items to read: introduction, write-up, project templates, extracts, and so forth.

Background and General Information

Jacob: Is the student project required? What does the student project show?

Rachel: You must submit the student project to receive course credit. The project shows that you can apply the time series concepts to actual data.

Jacob: When is the student project due? Is it like the final exam, which is given on a set date, or like the homework assignments, which must be completed to receive credit?

Jacob: How complex is the student project? Would a loss cost trend analysis for personal auto claims be suitable for a student project?

Rachel: The student project uses time series concepts and techniques. You fit ARIMA models to a time series, using the methods taught in the on-line course.

- Fitting an exponential curve to average claim severities is a linear regression after taking logarithms. It does not use stochastic ARIMA processes or the model building methods in the time series course.
- Modeling loss cost trends with ARIMA processes is a good topic for a student project. See the project template on loss cost trends.

Jacob: Do we do the student project after finishing the course or as we do the modules?

Rachel: The time series modules build upon one another. The last five modules (Chapter 19 of the textbook) apply the concepts to actual data. You won't have a good sense of building ARIMA models until you complete the modules.

Jacob: Should we wait until the last two weeks of the course to do the student project?

Rachel: The student project deals with ARIMA models: specification, diagnostic testing, forecasting, sample autocorrelation functions, Box-Pierce Q statistic, correlograms,

NEAS Time Series Project Documentation

autoregressive and moving average processes, structural models, stationarity, Yule-Walker equations, and integrated models.

- Learn the material over the eight weeks of the course. Focus on the textbook readings, homework assignments, and practice problems.
- During the last two weeks (last six modules) of the course, begin reading the postings on the student projects. The illustrations in Chapter 19 of the textbook are complex, using non-linear regression and partial autocorrelation functions to fit high order ARIMA processes, such as ARIMA(4,2,4) or ARIMA(6,1,8). In practice, we use simpler ARIMA processes for actuarial and financial time series.
- Until the final exam, focus on the practice problems and the statistical techniques.
- Right after the final exam, review the project templates on the NEAS web site, read several of the past student projects posted on the discussion forum, and pick a topic. Spend an hour surfing the web, looking for data on a topic that interests you. Suitable time series on thousands of topics are freely available on the web.
- The NEAS web site has hundreds of time series you can use for the student project and a wide array of project templates, illustrative workbooks, step-by-step guides, and instructions. Your student project may analyze daily temperature for your home town following the project template and illustrative workbook on the web site or it may use a time series that you pick from the internet or your company's data.

Do a student project on a topic that interests you. The time series can be sports scores, movie ticket receipts, DVD sales, voting behavior, crime rates, marriages, divorces, births, abortions, immigration, or gas prices. You can choose a project related to your work, such as claim frequency, claim severity, mortality, premium trends, or interest rates.

- Before doing the student project, most candidates fear it will be a burden.
- After completing the project, most candidates say it was the best part of the course. The time series concepts in the textbook are abstract. Using them to analyze a real time series makes them come alive.

Use the discussion forum on student projects for ideas. Read the project templates, review the past student projects, look at the data sets and time series. Explore the internet, using the search engines. You will find hundreds of sites with potential data for a student project. Discuss potential topics with other candidates.

APPLYING TIME SERIES

Jacob: What are ARIMA models used for?

Rachel: ARIMA models are used for many items:

- Business *growth*, product sales, inventory management.
- Regulatory, economic, social, and legal *interventions*.
- Macroeconomic indices: Interest *rates*, inflation, unemployment, exchange rates.
- *Seasonality*, cycles, regimes (eras); see the project template on daily temperature.

NEAS Time Series Project Documentation

- *Residuals* of regression relations (see the project template on structural models)

The textbook emphasizes business growth, product sales, and inventory management. For time series on product sales, use the internet search engines. These data are on industry (trade association) web sites.

The *project templates* on the discussion forum use interest rates, inflation rates, daily temperatures, and insurance statistics. They discuss seasonality, time periods, and residuals.

- Many actuaries deal with interest rates, inflation rates, and similar indices. The project templates on these subjects apply the time series concepts to real actuarial work.
- The macroeconomic indices are publicly available on government web sites. Many indices are on the NEAS web site in Excel sheets that you can download.
- Daily temperatures are available for 1,221 weather stations for 100+ years. The data are in Excel (CSV) format on the NEAS web site that you can download.

Recommendation: If you want to do a student project on macroeconomic indices or daily temperatures, check the data already on the NEAS web site. Check also the suggestions for other student projects on the web site.

INTERVENTIONS

Social scientists use time series analysis to examine changes in the environment. The 1973 Middle East war led to the OPEC oil cartel. We examine if the price of oil follows a different pattern after 1973 than before 1973.

Similarly, Federal Reserve Board policy affects interest rates, minimum wage laws affect unemployment, deregulation affects prices of airline tickets and phone service, insurance regulation affects policy premiums. Examine time series before and after an intervention.

You can examine effects of home DVD sales on movie ticket sales, effects of laws and court decisions on abortions, welfare, and crime rates, effects of sports rule changes on won-loss records, and so forth.

Real research projects examine all material effects on a time series. We don't examine time series of abortion rates before and after the Roe decision because many other social and political changes also affect abortion. If you can easily include other effects in your student project, do so. But don't give up on an interesting project because you are missing data. If you have a time series of abortion rates, compare the pre-Roe vs post-Roe rates.

Take heed: For seasonality, see the project template on daily temperature and the discussion forum posting on seasonality.

Take heed: For cyclical models, see the discussion forum postings on macroeconomic indices and the project templates on real interest rates.

NEAS Time Series Project Documentation

INDEPENDENT STUDENT PROJECT: DESIGN, DATA SETS, AND ANALYSIS

Jacob: The student project applies the time series concepts to an actual time series. Do we apply specific concepts? Do we use a specified time series? How do we know what techniques to use and which data set to apply them to?

Rachel: You can choose the time series, data set, and statistical techniques. You have just taken a course on time series analysis. Use the techniques taught in the course to analyze a time series of your choosing.

Your project can be independent or it can follow the templates on the discussion board.

You will enjoy the project more if you choose a time series that interests you. But open-ended assignments are sometimes baffling. Use the project templates and the student projects posted on the discussion forum two ways:

- The project templates give you ideas. Work through the project templates on interest rates and daily temperature. The project templates give you step-by-step instructions. Download the time series and reproduce the tables, charts, and graphics.
- You can do a student project modeled on a project template. Instead of 90 day Treasury bill yields, use a CPI index or a corporate bond yield.

Jacob: What exactly do we choose? The discussion forum has instructions and step-by-step guides. Do we follow these instructions?

Rachel: Your choice has three levels: design, data sets, and analysis. The instructions and step-by-step guides lead you through setting up and starting a student project. They don't finish the work; no project template completes the project. You decide what your project will cover, and the final form of the project depends on your choices.

You can *design* your own project. The project templates show the type of analysis, but your project can be different. You might use employer data, client data, or publicly available data showing a time series.

- *Pricing:* you might use quarterly premium volume over the past ten years.
- *Investments:* you might use a firm's weekly stock price over the past five years.
- *Social:* you might use U.S. crime rates over the past fifty years.

Take heed: Crime rates are compiled by the FBI and municipal police departments and analyzed by politicians, social scientists, and journalists.

Jacob: Could we use annual profit margins over the past ten years?

Rachel: A time series of ten values is too short for the statistical analysis. The NEAS web site shows monthly interest rates for 55 years, giving enough data to examine the effects of different eras and to test for significance of residuals.

NEAS Time Series Project Documentation

Daily temperatures are ideal for time series analysis. The national weather service provides free data bases showing high and low daily temperatures (and much more weather data) for over a thousand locations for the past hundred years. These data have the time series relations modeled by ARIMA processes:

- Movements of cold and warm fronts cause autoregressive and moving average effects.
- Daily temperature is seasonal.
- Trends in daily temperature are unclear; you can examine possible warming or cooling in your home town over the past hundred years.

A project may also compare similar time series, such as time series models for premium volume (or loss costs) in two states or for interest rates in two periods.

Take heed: Do not worry that the results may not be significant. If you want a particular topic but the only time series you find has just 50 observations, use the data. Your write-up will note that the results are distorted by random fluctuations, but you will enjoy the work.

Jacob: What if we don't have our own data? Do you give us data and a project template?

Rachel: Data are available on hundreds of web sites. But if you have never used ARIMA models, you may not know what time series are best for the student project. We provide data sets and templates on the discussion board.

Recommendation: You enjoy the project more if you choose your own topic. Pick a topic and use Google or another search engine to find data on the web. Add keywords like *history* or *trend* to the search criteria.

Illustration: If you search for *daily temperature*, you get today's weather. If you search for *daily temperature history*, you get historical time series.

Jacob: What time series do the project templates use?

Rachel: We form a time series model for interest rates in the United States. We also post related time series on the web site: inflation, GNP, unemployment, exchange rates, and interest rate futures. Another project template shows ARIMA modeling of daily temperature. We show a long project template on daily temperature and shorter templates on various other topics.

Jacob: People have tried to predict interest rates for years. The interest rate models used by actuaries and economists are complex. Are we supposed to improve on these models?

Rachel: The goal of the student project is not to develop the ideal model.

- The project takes actual interest rates and applies the methods in the textbook. You are learning how to apply the concepts, not developing better interest rate models.

NEAS Time Series Project Documentation

- The project template on daily temperature shows how to adjust for seasonality. We don't use ARIMA processes to forecast the temperature.

Jacob: Interest rates comes in dozens of forms. What time series do we use for the project template on interest rates?

Rachel: We illustrate with 90 day Treasury bills and overnight LIBOR rates, and we provide guidance for other rates. You can pick among several dimensions: long vs short rates, nominal interest rates vs real interest rates, absolute values vs residuals, and spot vs future rates. For example:

- Short rate: three month Treasury bills, overnight LIBOR
- Long rate: twenty year Treasury bonds, Moody's corporate bond rates
- Residuals of three month Treasury bills regressed on the CPI

We provide the time series on the web site along with specific project templates.

Recommendation: If you feel lost, begin with three month Treasury bills, Atlantic City daily temperature, or another project template on the NEAS web site.:

- Compare your results with the examples in the project templates.
- If you can't reproduce the time series charts, post a question on the discussion forum.

After you work through the statistical techniques (correlograms, Box-Pierce Q statistic) ARIMA modeling will make more sense. Select another time series for your student project:

- Another interest rate time series, or another macroeconomic index.
- Daily temperature for another weather station.

Once you are comfortable with the statistical techniques, you can choose any time series and do a good project.

COMPARISON STUDENT PROJECTS

Many student projects compare two or more time series.

- Several types of interest rates, several time periods, or real vs nominal rates.
- High vs low temperatures, daily temperatures in different weather stations or in different time periods.

Do not worry about choosing a correct topic. We require only that the time series is not already a white noise process or a random walk.

- No topic is inherently right or wrong for a student project.
- The keys to a good project are an interesting question on an interesting topic.

NEAS Time Series Project Documentation

An interesting question might be

- How well can we predict Friday's daily temperature with the daily temperatures on Monday through Thursday?
- Does the time series pattern of daily temperature differ on the East Coast (New York) vs the West Coast (San Francisco)?
- Has the pattern of daily temperature changed between 1901-1950 and 1951-2000?

You may be surprised by the answers to the questions above.

NEAS Time Series Project Documentation

STRUCTURAL MODELS VS ARIMA MODELS

Jacob: Interest rates depend on current inflation, expected future inflation, economic growth, Federal Reserve Board policy, and business cycles. Don't we have to consider these variables to forecast interest rates?

Rachel: The textbook discusses this issue several times. The best model for any time series uses all the relevant explanatory variables.

- Sales of pharmaceutical firms depend on the quality of new medications, patents for existing medications, and (perhaps) marketing. To forecast drug sales, we consider outstanding patents and the new medications under development.
- Premium volume for an insurer depends on its rate level compared to other insurers, the quality of its sales force, and special marketing efforts.
- The daily temperature depends on a host of weather variables, not just the temperature on the past few days.

The textbook gives several answers.

Answer 1: Structural models, which rely on other explanatory variables, are complex. They work well in hindsight but have little predictive power if they rely on unknown values, such as future inflation or business trends. The time series is a simple model that may perform nearly as well.

Answer 2: Most effects on a time series are gradual, hard to observe, and hard to quantify. An insurer's premium volume depends on rate level changes, competitors' actions, or the opening of a new sales office. When forecasting the insurer's sales, we may not have this information, and we don't have models linking competitors' actions to the premium volume.

These exogenous causes affect the time series values for both past and future periods. An ARIMA model may be a good forecasting tool.

Answer 3: The absolute level of the time series (such as interest rates) depends on other economic influences. But the remaining effect on interest rates – the residual – may not be random. We estimate a time series of the residuals. Chapter 19 of the textbook shows examples of this. Your student project may focus on residuals; we explain the intuition.

Residuals

Jacob: Suppose we forecast interest rates from inflation rates. How can ARIMA modeling of the residuals help? Don't we assume the residuals of a regression model are random? If the residuals are not random, are the results of the regression analysis still valid?

Rachel: Classical regression analysis assumes the residuals are a white noise process if the explanatory variables completely explain the dependent variable and the regression

NEAS Time Series Project Documentation

equation is correct. But we never know all the explanatory variables, and the true relation is not exactly linear. For a time series, the residuals often have serial correlation.

The regression analysis course uses the Durbin-Watson statistic. This is similar to the autocorrelation of lag 1, but scaled from 0 to 4 instead from -1 to 1.

The Durbin-Watson statistic tests for serial correlation. Ordinary least squares estimators are unbiased even with serial correlation, but the statistical testing is no longer accurate.

Using a time series model on the residuals is more sophisticated, in two ways:

- We examine more complex relations than just the autocorrelation of lag 1.
- We use the relation among the residuals to better forecast the values of the time series.

Illustration: Fisher Effect

Financial economists assume a Fisher effect between interest rates and inflation rates. The nominal interest rate is the inflation rate times the real interest rate, such as 1.02, or plus the real interest rate, such as 2%.

In the long-run, interest rates and inflation rates are correlated. In the short run, the relation is weaker. If inflation rises by 50 basis points, the interest rate may rise by 25 basis points.

Jacob: Does the student project test the Fisher effect?

Rachel: We assume the Fisher effect holds, and we test whether a time series model is appropriate for the residuals. A financial economist may presume that the relation between interest rates and inflation rates has a long-term mean, such as 200 basis points, and a strong autoregressive quality in the short run.

Illustration: Suppose the long-term real interest rate is 2% for an additive model or 1.02 for a multiplicative Fisher model. If the real interest rate is 4% in January 20X6, we don't expect it to stay 4% forever or to become 2% in February 20X6. The real interest rate may drift back to 2% over the next year or two.

Your student project may compare ARIMA models for nominal and real interest rates. If inflation is volatile, an ARIMA process may forecast real interest rates well, but not nominal interest rates.

Take heed: The student project is most interesting if it examines structural relations and it allows you to include more statistical techniques.

Illustration: A time series on Chicago crime rates fulfills the VEE requirements. Adding other pieces makes the project more interesting. Look at population densities, police force policies, gang activity, and similar items.

NEAS Time Series Project Documentation

Seasonality

Interest rates and inflation rates may vary over the year.

- In December, consumers shop for holiday gifts, and the demand for money is high.
- In January, they return to work, and the demand for money falls.

The Federal Reserve Board adjusts the money supply each month to dampen fluctuations in inflation rates and interest rates.

- The adjustments are not perfect, and some weak seasonality remains.
- The seasonality is strongest in 30 day commercial paper rates, weak but noticeable in 90 Treasury bill rates, and smoothed in smoothed in yields of one year Treasury bills and longer securities.

Your student project can examine seasonality in various types of rates. For Treasury securities, the seasonality is strongest in monthly first differences of 90 day Treasury bills.

The student project may also examine the seasonality in nominal interest rates vs real interest rates.

- The seasonality affects inflation rates and nominal interest rates, not real interest rates.
- Use the raw CPI to derive real interest rates.
- The seasonally adjusted CPI does show the seasonal fluctuation.

Your student project can examine residuals of interest rates regressed on other indices.

Illustration: Regress interest rates on GDP (also on the NEAS web site.)

- Economists do not agree on the expected regression coefficients.
- The FED often tries to moderate GDP growth by raising or lowering interest rates.
- It is unclear if the FED can influence interest rates or GDP. The macroeconomic course explains the views of different financial economists.
- Your student project may compare ARIMA models fitted to the nominal interest rates vs fitted to the residuals of this regression.

The NEAS web site shows no index for the *expected* inflation rate. Use the actual CPI as a proxy for the expected inflation rate. This is not an ideal proxy, but we are teaching statistics, not economics.

NEAS Time Series Project Documentation

Four Step Process for ARIMA Modeling: Specification

STATIONARY INTEREST RATE SERIES

The textbook's four step process for modeling time series begins with specifying the model. Form a stationary time series. Several scenarios may occur. For interest rates:

- The time series is not stationary and cannot be converted into a stationary series. Most interest rates series are stationary or homogeneous non-stationary (so their differences are stationary). Economies with run-away inflation have non-stationary nominal interest rates and probably non-stationary real interest rates. ARIMA modeling of inflation rates in Zimbabwe does not work.
- The time series (or its logarithm) is stationary. Most mean reverting time series are stationary. Logarithms change multiplicative models into additive models. For interest rates, we don't usually take logarithms. For dollar denominated items (sales, premium, stock prices), take logarithms. See the discussion forum posting on random walks.
- The first differences of the time series (or of its logarithms) is stationary. If the time series is growing or declining (not mean reverting), take first differences. Economists disagree whether interest rates are mean reverting. Your student project examines the correlograms of the time series and its first differences. If the time series is not stationary but its first differences are, your report should explain why this occurs.
- More complex models using second differences might be needed. Second differences are not used for interest rates. Second differences might be used for a time series of the balance in an account where monthly deposits are made.

If a time series is not growing or declining, it is not seasonal, it is mean reverting, and it has no changes in its mean, it is usually stationary.

- Nominal interest rates may have a upward or downward drift (reflecting inflation) in some periods. Other times they have no drift. If your student project uses nominal interest rates (or any other index that moves with inflation), you may have to use separate periods or take differences.
- It is unclear if real interest rates are mean reverting. Examine the sample autocorrelation function (the correlogram) to determine if a time series is stationary. Use real interest rates for a better model.

A student project may examine differences between short rates, long rates, and residuals of interest rates on inflation rates.

- The short rate fluctuates more than the long rate.
- Their mean reversion is hard to predict.

Many economists assume the interest rate residuals (on inflation) are mean reverting.

- If the time series is not mean reverting, model it as a random walk.

NEAS Time Series Project Documentation

- If interest rates have strong mean reversion, use a AR(1) or AR(2) process.

Deciding whether to model the initial time series or its first differences is critical. Graph the time series. If it has a drift, it is not stationary, so take first differences.

The interest rate graphs on the web site show drifts that vary by time period. Your student project examines if the time series changes.

Illustration: Suppose you examine the daily high temperature in City Z for 1875 - 2005.

- For 1875 – 1940, the daily temperature is recorded every six hours. The high temperature is the 12:00 noon reading.
- For 1941 – 2005, the daily temperature is recorded every hour. The high temperature is usually the 2:00 pm or 3:00 pm reading.

The temperature at 2:00 pm is 3 or 4° higher than at noon. Any single year gives a stationary time series, but the full time series of 131 years is not stationary.

Changes in the average daily temperature from smog, urbanization, ocean currents, or global warming may cause your time series to be non-stationary.

Drifts may be linear or exponential.

- If the drift is linear, take first differences.
- If the drift is exponential, take logarithms and then first differences.

Take heed: The project templates recommend detrending the time series instead of taking differences. If you eliminate the trend with an inflation index, an ARIMA process fits better.

If no drift is evident, we examine if the time series is mean reverting. A mean reverting time series moves up if it is below the mean and down if it is above the mean. A mean reverting time series can be either oscillatory or asymptotic. Many actuarial and financial time series have an exponential drift or are not mean reverting. Interest rates may or may not have a drift. Real interest rates are generally mean reverting; nominal interest rates may or may not be mean reverting.

Jacob: If the time series has no drift and is mean reverting, is it stationary?

Rachel: A time series is not stationary if the mean or variance changes. The project template on interest rates shows how FED policy affects your ARIMA modeling.

Illustration: In the 1970's, the U.S. tried to use inflation to control unemployment. In the early 1980's, the Federal Reserve Board slowed the money supply growth to restrain inflation. The mean nominal interest rate rose and then fell by over 500 basis points in a few years, and the variance of the interest rates increased and then declined significantly.

NEAS Time Series Project Documentation

You have two possible solutions:

- Separate the time series into three eras (regimes).
- Use the residuals of interest rates on inflation.

Using second differences is not correct.

Jacob: How do we find and test for a change in regime?

Rachel: We use exogenous knowledge; we have no specific statistical test. We examine the plot of interest rates and see if a change in the mean, variance, or drift seems likely. We show an illustration on the discussion forum: the mean, variance, and drift of interest rates changes materially among three periods.

If we suspect a change in regime, we use two or more time periods, often with an interlude in between. For interest rates, we show illustrative time periods of 1945-1978, 1979-1982, and 1983-2000. For your student project, you can divide the time period differently. We re-do the time series analysis on each part.

Jacob: If the ARIMA models fit better by part, do we use separate models?

Rachel: We get more stable results in each part even if there are no real changes. Unless the change is large, we don't use separate parts.

Illustration: If we get a mean interest rate of 8% in one period and 7% in the second period, we assume this is random fluctuation and use a single time series. If the means are 11% and 5%, we assume a regime change and use two periods.

We form correlograms for each period. *The correlogram should validate the intuition from looking at the plots.* For the student project, you form the plot, make a guess about the type of time series, and then form a correlogram to validate your inference.

Speed of Decline

If the sample autocorrelations do not decline to zero, or if the decline is *too slow*, the time series is not stationary. We take first differences (and perhaps logarithms) to form a stationary time series.

Jacob: What is *too slow*? Are sample autocorrelations of $\frac{1}{2}$, $\frac{1}{3}$, $\frac{1}{4}$, ..., $\frac{1}{n}$, ..., too slow?

Rachel: *Too slow* means the time series is not autoregressive or moving average.

- If the time series is moving average, the autocorrelation drops to zero after period Q, where Q is the order of the moving average model.
- If the time series is autoregressive, the autocorrelations have a geometrically declining envelop after period P, where P is the order of the model.

NEAS Time Series Project Documentation

The sample autocorrelations in your example show too slow a decline.

Jacob: Should we examine also the second and third differences?

Rachel: We stop taking differences as soon as we get a stationary series. Few time series are homogeneous of order more than one. It is sometimes worthwhile to look at the second differences, but that is not required for the student project.

Jacob: Are the interest rate time series all similar?

Rachel: The time series differ by type of interest rate and by era. The short rate may have a varying mean, the long rate may be a random walk, and the real interest rate (the residual) may be white noise. Your student project can compare the time series.

Length of Periods

Jacob: How long a period do we use? Some interest rates on the web site start in 1945. Do we use all 60 years?

Rachel: The length of the time series depends on two items:

(1) For a short period, the model parameters depend on factors we don't want to include.

- A single year of monthly interest rates is not enough. The short run fluctuations depend on economic factors specific to that year. The ARIMA model will not forecast well.
- Random fluctuations distort the sample autocorrelations for short periods. For a single year of monthly rates, the standard deviation of the sample autocorrelations is $1/\sqrt{12} = 28.87\%$ even if the time series is white noise. Even an autocorrelation of 50% may be random fluctuation. The observed time series doesn't tell us much.
- Seasonal patterns require a series of several years. The student project should test for seasonality, so we use several years.

Jacob: Don't the statistical tests take random fluctuations into account?

Rachel: If the period is too short, the time series is affected by other factors and is not correctly specified. *In-period* goodness-of-fit tests are good: R^2 is high and the Box-Pierce Q statistic is low. But the time series forecasts poorly; we get poor out-of-period results.

(2) For a long period, other factors affect the model. The depression years, World War II, and the modern period have different interest rate models. Putting them together gives a non-stationary time series. Monthly interest rates show one pattern when the Federal Reserve Board focuses on rate stability and a different pattern when the Federal Reserve Board focuses on unemployment.

Choosing Periods

NEAS Time Series Project Documentation

Part of the student project is choosing the proper length of the time series. We illustrate with post World War II years: 1945 and onward. You can choose other periods if you want.

Jacob: How do we use to choose the proper length?

Rachel: We look for stability of the model coefficients. Suppose we use historical data through December 2005 to predict interest rates in 2006. In-period goodness-of-fit tests don't tell us if the parameters are stable, so we use out-of-period tests. We form three time series models:

- Data through December 2002 to predict interest rates in 2003.
- Data through December 2003 to predict interest rates in 2004.
- Data through December 2004 to predict interest rates in 2005.

We examine two things:

- The ARIMA model gives a forecast variance. If the actual interest rates are sufficiently close to the forecasted interest rates, we don't reject the model. *Close to* means within one or two standard deviations.
- The ARIMA model gives standard errors for the time series coefficients: the ϕ_1 , ϕ_2 , and so forth. We examine if the coefficients estimated in different periods are the same. For data from 1985 through 2005, we might use three periods: 1985-1991, 1992-1998, and 1999-2005. Each period has seven years of monthly rates, or 84 data points. The samples are large enough that a change in the time series coefficients indicates the model is not good. If the standard errors of the regression parameters are 5 basis points, a ϕ_1 of 20% in one period and 30% in another period, we use separate ARIMA models. (Standard errors of the regression parameters are covered in the regression analysis course.)

Autocorrelations

Jacob: After choosing between the interest rates and the first differences and selecting a length for the series, what is the next step?

Rachel: We examine the sample autocorrelations and the partial autocorrelations.

Jacob: How many lags do we use? Do we focus on the sample autocorrelation for one lag at a time or for all lags?

Rachel: To specify a time series model, we examine one lag at a time, using the principle of parsimony. For diagnostic testing, we examine the Box-Pierce Q statistic for a large number of lags.

Jacob: Is there an order for the models we try? Or do we try a bunch of ARIMA models to see which fits best?

NEAS Time Series Project Documentation

Rachel: The textbook implies that the sample and partial autocorrelation functions indicate the proper type of model. The textbook authors examine various models, comparing the R^2 and the Box-Pierce Q statistic for each.

In practice, we use only four or five models. Most ARIMA models are AR(1). We use an AR(1) model and examine the sample autocorrelations of the residuals from this model. If these autocorrelations are not statistically significant and the Box-Pierce Q statistic is not significant, we assume the time series can be modeled as AR(1).

Out-of-Sample Tests

Jacob: The AR(1) model may not be optimal, even if the fit is reasonable. Suppose the sample autocorrelations are 50% for lag 1 and 30% for lag 2. If the time series were AR(1) with a ϕ_1 of 50%, the sample autocorrelation of lag 2 would be 25%. Since we get a 30% autocorrelation for lag 2, we assume a small but positive value for ϕ_2 .

Rachel: If the standard error of the sample autocorrelations is 15%, the difference between 25% and 30% is not significant.

Jacob: Even if it is not significant, the AR(2) model fits better. Is there any harm in using the AR(2) model?

Rachel: Using a higher order increases the in-sample goodness-of-fit, but the model may not be better. In the example you gave, most statisticians would bet that even though the AR(2) model has the better in-sample fit, the AR(1) model forecasts better.

Jacob: How do we test the out-of-sample goodness-of-fit?

Rachel: Suppose we use monthly interest rates for 1985 through 2005 to forecast future interest rates. Our final model uses all the historical data. To build the model, we use interest rates for 1985 through 2004. We may specify two or three models. The in-sample goodness-of-fit tests are not decisive. The optimal model is the one that minimizes the mean squared error of the forecast.

Jacob: Do we choose the model with the lowest mean squared error for the forecast?.

Rachel: It is easy to pick the model that minimizes the *observed* squared error; it is hard to pick the model that minimizes the *mean* squared error.

Jacob: What is the difference?

Rachel: Interest rates are stochastic. Suppose two models fit equally well in-sample. We may have several scenarios for next year's interest rates. For some scenarios, one model fits better; for other scenarios, the other model fits better.

Jacob: How do we solve this problem?

NEAS Time Series Project Documentation

Rachel: One solution is to examine various periods, or interest rates in various countries, or various interest rate series (Treasury bills, LIBOR rates, ...).

Alternatively, we rely on intuition and experience with time series models. The intuition for an AR(1) model is strong. Most time series show stability; if this period's value is high, next period's value may also be high. Other relations are rare: negative serial correlation and ARIMA models with $p+q > 2$.

Seasonality

Jacob: How do we include seasonality? The textbook uses various seasonal adjustments.

Rachel: Seasonal correlations should be reasonable. High-renewal annual policies have a 12 month autocorrelation in premium volume. Interest rates depend on the demand for money, which varies seasonally. For any time series, we begin with AR(1) and a seasonal adjustment, such as AR(4) for quarterly data or AR(12) for monthly data.

Illustration: Suppose we model interest rates as a ARIMA(2,1,0) time series, using 20 years of quarterly data. We also model the rates with a four quarter seasonal autocorrelation. The textbook calls this ARIMA(4,1,0). Other statisticians use notation to specify that the non-seasonal part is ARIMA(2,1,0) and the seasonal part is ARIMA(4,1,0), with a *single* autocorrelation for lag 4.

Jacob: How do we test for seasonality?

Rachel: We examine the quarterly sample autocorrelations of lags 4, 8, 12, and 16. Seasonality appears as spikes in the correlogram. To test whether the 4 quarter seasonal lag helps, we examine the Box-Pierce Q statistic for the seasonal vs non-seasonal models.

Jacob: Do we use the time series model with the lowest Box-Pierce Q statistic?

Rachel: We consider also the principle of parsimony. More complex models have lower sample autocorrelations. If the Box-Pierce Q statistic adjusted for degrees of freedom is not lower, we use the simpler model. If it is lower, the choice depends on how much lower.

First Differences and Seasonality

Jacob: How do we combine first differences and seasonality?

Rachel: We have several methods.

(1) We take first differences of the time series and examine the four quarter lag of these first differences. If we examine toy sales, the first differences are positive for third to fourth quarter, negative for fourth to first quarter, and zero elsewhere. Similarly, if the Federal Reserve Board did not adjust the money supply for the higher demand in the fourth quarter, we would expect a large seasonal effect on interest rates and inflation rates.

NEAS Time Series Project Documentation

Jacob: Would interest rates be higher or lower in the fourth quarter of the year?

Rachel: Economists disagree about the effects of monetary policy on interest rates. Keynesian economists assume that a demand for money exceeding the supply of money causes real interest rates to rise, as consumers sell bonds and seek loans. Neo-classical economists assume money is neutral and has no effect on real interest rates, but a demand for money exceeding the supply of money reduces the price level, pushing down nominal interest rates. Your student project may examine the seasonality of interest rates.

(2) We take first differences of the interest rate from the rate lagged twelve monthly periods. The time series is the change in the interest rate from its level 12 months ago.

Illustration: Suppose interest rates are about 150 basis points higher in the fourth quarter of each year. The time series we examine is the 1/20X6 interest rate minus the 1/20X5 interest rate, the 2/20X6 interest rate minus the 2/20X5 interest rate, and so forth.

(3) We de-seasonalize the time series and take first differences from month to month. We first compute factors to de-seasonalize the interest rates, such as subtracting 150 basis points from the fourth quarter rates. We then take month to month first differences.

(4) We take both month to month differences and 12 month differences. This is a second difference, of the form $y_t - y_{t-1} - (y_{t-12} - y_{t-13})$.

(5) We take using first differences and an autoregressive model with non-zero parameters for lags 1 and 12.

Jacob: Which method should we use?

Rachel: Your student project can compare several methods of adjusting for seasonality. Often the simplest method is best. The textbook uses the last method above.

Mean Reversion

Jacob: When do we estimate an autoregressive model from the interest rates themselves and when do we use first differences?

Rachel: If mean reversion is strong, we use the interest rates. If mean reversion is weak (or zero), we use first differences.

Jacob: From 1990 through 2005, interest rates have been low; from 1975 to 1985, interest rates were higher. If there were strong mean reversion, interest rates should have regressed toward the mean in both periods.

Rachel: We infer that the mean changed, not that mean reversion is weak. A student project may compare various periods to see whether different models are appropriate.

NEAS Time Series Project Documentation

Interest Rates on the Web Site

Jacob: What interest rates are on the web site?

Rachel: You can use any interest rate you want, but it is difficult to discuss the methods on the discussion forum unless others are using the same rate. We have selected a few interest rates that are useful for the student projects:

- Three month Treasury bills by month from January 1931 through June 2000. Some years are missing, when auctions were not held.
- Twenty year Treasury bonds by month from April 1953 through December 2005. Some years are missing, when auctions were not held.
- Moody's seasoned AAA bond rates, by day, from January 3, 1983, through January 16, 2005.
- The CPI from January 1913 through December 2005, both (SA) seasonally adjusted and (NSA) not seasonally adjusted. There are many CPI indices; we show U.S. city average, all items.

Whatever series you use, choose a period during which the time series is stationary but a long enough period for statistical testing. We will put up additional series on the web site as time goes on.

Jacob: Can we take another interest rate time series from the internet?

Rachel: You can use any interest rate time series you want. Many financial economists believe that the time series models work better for interest rate *futures* than for the interest rates themselves. An excellent student project may compare a model for 20 year Treasury bonds with a model for Treasury bond futures. You can do the same for LIBOR rates and Eurodollar futures.

Jacob: Is the grading of the student project easier or harder if we choose a different time series or a different test than those in the project templates?

Rachel: The grading is easier. If you follow the project template, you should demonstrate that you understand the concepts by performing the various statistical tests. If you design another project, we are less concerned with performing all the statistical tests. If you select a reasonable item to test and you show how this is done, that is sufficient.

DIFFERENCES

Take heed: Do not take differences unless they are needed.

Jacob: Do the first differences of a mean reverting time series give a better model?

NEAS Time Series Project Documentation

Rachel: Taking first differences of a stationary time series makes the forecasts less robust. Once we have a stationary time series, we use it.

NEAS Time Series Project Documentation

TIME SERIES STUDENT PROJECTS: FACULTY GUIDANCE VS INDEPENDENT WORK

(The attached PDF file has better formatting.)

Updated: May 3, 2008

Jacob: Do many candidates have difficulty with the student projects?

Rachel: Most candidates produce good projects; some of the projects we have received are excellent. Other candidates have more trouble. We give enough guidance that all candidates can complete the student projects, but we do not over-specify the work needed.

Jacob: Why not give a list of steps that we must complete to receive VEE credit?

Rachel: The SOA wants *independent* student projects. If we gave a list of steps for the student project which candidates apply to their data, the SOA might not grant VEE credit.

Our over-riding objective for the on-line courses is to provide the education desired by the SOA and CAS. Our second objective is to provide you with the guidance to fulfill the SOA and CAS requirements.

We provide the level of detail in the instructions and the illustrative work-sheets that best fulfills the SOA and CAS educational desires and provides the guidance needed by you to meet them. As we review and grade the student projects, we see what topics are most confusing and we provide more detailed explanations.

Jacob: Not all candidates do know what is expected of them. How will you deal with this?

Rachel: We post project templates with extensive analysis of selected student projects. We show the type of project you can complete on daily temperature or interest rates; we give hundreds of time series with suitable data; and we give illustrative worksheets with cell formulas and VBA macros that perform the basic statistical techniques. Read these project templates to understand what is expected. We also post sample student projects with comments by the NEAS faculty. Once you see what other candidates have done, you feel more comfortable writing your own project. We suggest additional topics for the projects, such as modeling of corporate bond spreads, FBI crime rates, personal income by state, unemployment rates for teen-agers, and simulation projects.

Jacob: Can we use the same techniques as these projects use? Can we choose another time series and apply the same techniques?

Rachel: The basic structure of these projects is similar. Use the techniques in the on-line course, though the inferences and conclusions must be your own. Once you understand how to form correlograms and use the in-sample goodness-of-fit tests, the mystery of the student projects subsides. The SOA wants to see that you can apply the techniques to real data; it does not require innovative work from candidates taking an introductory course.

NEAS Time Series Project Documentation

TS PROJECT TEMPLATE ON BIRTH RATES

ARIMA models are often used to project population growth and decline. Candidates who work with population statistics have much data that can be used for the student project. They are familiar with these time series and they know the trends and changes that can make a good project.

Actuaries examine the components of population. The population in year T is

the population in year T-1 + births – deaths + immigration – emigration.

Even if births, deaths, immigration, and emigration change, the population time series may seem stable. You might model population as an AR(1) process with low stochasticity, which fits reasonably well. But 95% of the population is the same in successive years, so the time series for the total population doesn't say much.

The first differences are the combined effect of births – deaths + immigration – emigration. These four pieces (births, deaths, immigration, and emigration) can all be used for the time series student project. Each has its own process.

This project template outlines a student project on birth rates. The other three pieces are related to government policy (immigration), wars (deaths and emigration), and health care (deaths). You can do a student project on any of the four pieces, but keep in mind that a change in government policy or a major war can disrupt the ARIMA process.

Birth rates are less affected by government policy, wars, and health care. In *developing countries*, health care affects infant mortality. This effect is smaller in advanced economies, where infant mortality is low.

Birth rates have good statistics. Hospitals keep records of live births and most countries tabulate and publish the results.

Birth rates have several characteristics that are useful for the student project. A change in birth rates affects the *long-term population trend*, so it is a good indicator of population.

(1) Most countries show declining birth rates. Japan and Russia have such low birth rates that their populations are declining. Many Western European nations will have negative population growth in the 21st century (Spain, Italy, Germany). In much of the developing world, birth rates are still high but they are well below earlier rates.

Each nation has its own trend, which you estimate from the data. Negative trends are more challenging to model, since the birth rates do not become negative. The trend is neither linear nor logarithmic; you can compare the two types of trend or choose another trend.

(2) Instead of a simple trend, demographers relate birth rates to girls' education, women's participation in the work force, religious commitment of the population, and other factors.

NEAS Time Series Project Documentation

The low birth rates in Western Europe, China, and Japan are characteristic of secular (atheistic) countries with equal education for both sexes and high female participation in the work force. You may look at birth rates in Japan or Western Europe, fit an ARIMA process, and project future rates. China's birth rates reflect government policy, so the ARIMA model differs before and after the single-child policy.

(3) The U.S. shows a different pattern. Birth rates have not declined much in the twentieth century, and they even increased from 1945 to 1960 (the baby boom generation). The U.S. remains a religious country, and birth rates do not show the same decline evident in secular nations.

In the United States, birth rates are positively correlated with economic conditions. We use the U.S. experience for this project template.

(4) Birth rates show a moving average component, which is less common in the other time series on the NEAS web site. We *test for* moving average components in many time series, but interest rates, daily temperatures, and most economic indices do not seem to have moving average components.

We explain the rationale for the moving average component by comparing babies to durable goods. Moving average components occur in sales of expensive, durable goods. Cars are a good example.

Illustration: Suppose the average car lasts five years and all consumers own one car. We add stochasticity and more realistic assumptions in a moment. In this ideal scenario, a consumer has a 20% probability of buying a car each year.

A car may last more or less than five years, and consumers have different views about when the car wears out, so any consumer's purchase is stochastic.

Suppose a consumer buys a car in 20X0. If the car lasts exactly five years and income stays constant, the consumer has a 0% probability of buying a car in 20X1 through 20X4 and a 100% probability of buying a car in 20X5.

But some cars last only 1 year and some cars last 9 years. The probability of buying a car may be 2%, 5%, 10%, 18%, 30%, 18%, 10%, 5%, 2% in the next nine years.

We now add economic conditions. Cars are expensive. Consumers are less likely to buy cars in recessions, when income is low, and more likely to buy cars in prosperous years, when income is high. This is not true for most consumer purchases, such as food, clothing, or books, but it is true for cars, homes, and other durable goods.

We may model car purchases as $20\% + \beta \times \text{GDP growth}$. *GDP growth*, as used here, is the residual from the mean, and β is a positive coefficient. Suppose the mean GDP growth is 2% and β is 3. GDP growth may be -1% in recessions and $+5\%$ in prosperous years, so the residuals are -3% and $+3\%$.

NEAS Time Series Project Documentation

- ~ In recessions, the percentage of consumers buying new cars is $20\% + 3 \times (-3\%) = 11\%$.
- ~ In prosperous years, the percentage of consumers buying new cars is $20\% + 3 \times (+3\%) = 29\%$.

We combine the economic conditions and the probabilities for a consumer.

- ~ If the economy enters a recession in the fifth year (20X5), the 30% probability for a particular consumer may drop to 15%. With more old cars on the road, the demand for new cars rises. Even if the economy stays in a recession, the probability of buying a new car in the sixth year (20X6) may rise to 25%.
- ~ If the economy is prosperous in the fourth year (20X4), the 18% probability may rise to 35%. With fewer old cars on the road, the demand for new cars falls. Even if the economy stays prosperous, the probability of buying a new car in the fifth year (20X5) may drop to 25%.

We summarize the sales of new cars as:

Sales volume depends on economic conditions. If consumers buy many cars one year (autoregressive process), or more cars than expected one year (moving average process), they are likely to buy fewer cars the new year.

BABIES

Children are the most expensive and most durable goods. Many women (couples) want two or three years (or more) in between children. Even if they would prefer children half a year apart, biology decrees otherwise.

Illustration: A young, childless couple in 20X0 may want three children. The probability of having a child in 20X1 may be 30%. If the couple does not have a child in 20X1, they may increase their efforts or consult a fertility clinic, and the probability of having a child in 20X2 may rise to 50%. If the couple has a child in 20X1, they may wait a year before trying to have another child, and the probability of having a child in 20X2 may drop to 10%.

This phenomenon may lead to a moving average process. If more women than expected have babies in year X, fewer women than expected may have babies in year X+1.

Children are extremely expensive. Mothers generally lose at least a year of work, and many mothers leave the work force entirely. For middle-income couples, the cost of raising a child may be several hundred thousand dollars.

Recessions cause a dip in the birth rates; prosperous years raise the birth rates. More precisely, they affect the conception rates, for which the birth rates are a proxy. The birth rate in Year T reflects economic conditions in year T-1.

ABSOLUTE VALUES VS DEVIATIONS

NEAS Time Series Project Documentation

The relation between birth rates and macroeconomic conditions makes the ARIMA process fit better (and may make the student project more interesting). You may do a student project on birth rates using the absolute rates if you prefer. If you use U.S. data, you may divide the time series into different periods. If you use other countries, the time periods may depend on government policies (as the one child policy in China) or other social conditions, such as the dramatic twentieth century fall in birth rates in secular countries.

Using residuals from a regression analysis on explanatory variables often gives a better fitting ARIMA model and a better understanding of the time series dynamics.

STEP-BY-STEP GUIDE FOR A STRUCTURAL MODEL

For a simple structural model, use the following steps. This regression model is too simplistic for actual demographic analysis, but it gives a good student project.

Step #1: GDP Growth Rates

Form GDP growth rates by taking logarithms and then first differences of GDP. This gives a series of GDP growth rates that you can use as the explanatory variable for birth rates.

Step #2: Regression

Regress the birth rates on a lagged GDP growth rate. If the birth rates are calendar year figures, use the GDP growth rate of the previous calendar year. Even if you have monthly birth rates, the decision to have a baby reflects economic conditions of about a year before the birth.

Step #3: Residuals

Form graphs of the absolute birth rates and the residuals of the birth rates. If the residuals are smoother, fit an ARIMA process to the residuals. The relation of birth rates and GDP growth is not linear, so the residuals are not a white noise process.

You can examine the trends in the absolute rates and the residuals. If you use monthly birth rate figures, you can examine seasonality.

Step #4: ARIMA Modeling

Use the time series techniques to fit an ARIMA process to the absolute birth rates and the birth rate residuals. Use correlograms, regression analysis, the Yule-Walker equations, or other techniques.

Use simple ARIMA processes. The residuals of the birth rates may be a stationary time series, and a simple ARIMA process may fit well and forecast well. It is harder to fit the absolute birth rates to a simple ARIMA process, since the trend in birth rates depends on the time period.

NEAS Time Series Project Documentation

- ~ Choose a time period with a homogeneous trend, take first differences, form a correlogram, and fit an ARIMA process to the first differences.
- ~ If your birth rates time series does not have a strong trend, you may not need to take first differences.

Step #5: Testing for White Noise

For each ARIMA process, test if the residuals are a white noise process, using the Box-Pierce Q statistic and Bartlett's test. We don't expect any ARIMA process to fit perfectly, since many other things affect birth rates.

These are the residuals of the ARIMA fitting, or the actual birth rates (or actual residuals from the regression analysis on GDP growth) minus the ARIMA estimates.

Step #6: Comparison

If you use both absolute birth rate and the regression analysis, judge whether the regression on the GDP growth rate improves the ARIMA fitting.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

PROJECT TEMPLATE ON ARIMA MODELING OF UNEMPLOYMENT RATES

Unemployment rates are a good time series for a student project because

- They vary by age, sex, region, and ethnic group.
- The relative unemployment rates by group have varied over the past decades.
- Each group's unemployment rate series has its own seasonality.
- Each group's unemployment rate depends on different macroeconomic factors.

The Excel file on the discussion forum shows unemployment rates

- seasonally adjusted and not seasonally adjusted
- for male vs female and teenager vs adult.

You can find unemployment rates by state, by other age groups, and by ethnic group on the internet. Demographers and social scientists study numerous unemployment rate relations that are suitable for student projects, such as

- Effects of state minimum wage laws on teen-age unemployment rates.
- Effects of Hispanic immigration on African-American unemployment rates.
- Effects of GDP changes on male vs female unemployment rates.

The unemployment rates are by month, so you can examine seasonality.

Illustration: School lets out in June for summer vacation.

- Adult male unemployment rates do not change much from May to June.
- Teen-age unemployment rates may double from May to June, as young people look for summer jobs.
- Female adult unemployment rates may reflect school-age children with vacation in the summer.

Your student project may examine seasonally adjusted and non-seasonally adjusted rates for adults vs teen-agers and male vs female.

Compare the seasonal adjustment actually used vs a simple seasonal adjustment using the procedure in the textbook. The actual adjustment should give a smoother series and a better ARIMA fit, but you can test these relations.

Retail firms have more job positions in December, with a slight seasonal effect on the unemployment rate.

Economists presume that unemployment rates are correlated with other macroeconomic indices, such as GDP and (perhaps) inflation.

NEAS Time Series Project Documentation

Unemployment rates should be inversely correlated with GDP. Regress unemployment rates on GDP and fit an ARIMA process to the residuals. The regression may differ for males vs females and adults vs teen-agers.

The relation of unemployment rates and inflation is debated.

- The macroeconomics on-line course assumes no relation.
- A previous generation of economists assumed a Phillips curve and a strong relation.

You can regress the unemployment rate on a lagged inflation rate. Use seasonally adjusted figures for both rates, or you will get a spurious relation.

Take heed: Some economists regress unemployment on unexpected inflation. One month LIBOR minus the CPI change is a measure of unexpected inflation. But don't feel bound by the economics. Any structural model is fine for the student project; we do not grade you on the economic reasoning. We examine if you properly regress one index on another and fit an ARIMA process to the residuals.

The teen unemployment rate is sometimes assumed to reflect minimum wage laws. The U.S. pattern in the minimum wage laws is

- Periods of no change in the nominal minimum wage law \Rightarrow a steady decrease in the real minimum wage.
- Federal minimum wage law legislation \Rightarrow a sudden jump in the real minimum wage.

Find the dates of minimum wage law changes by using the internet search engines (look for "minimum wage law"). See if the revised law causes a change in the teen-ager vs adult seasonally adjusted unemployment rate pattern.

Take heed: The data on the NEAS discussion forum are countrywide figures, with labor force participation combined with job finding and separation rates. Actual unemployment rates differ by industry, and careful econometric studies examine numerous influences on employment in each state, industry, age group, ethnic group, and other class dimension. You won't see clear effects, such as a rise in the unemployment rate when the minimum wage increases. We grade the student project by your use of the ARIMA techniques, not the validity of your conclusions.

Recommendations: Candidates who want a student project on unemployment rates that is particularly relevant to modern issues may compare unemployment rates in Ireland and Germany from 1946 to now.

- The Post-WW2 "German Miracle" rebuild the country in the 1950's – 1970's, with low unemployment and high GDP growth. By the 1980's, the German welfare state was beginning to stifle the economy. The pressure of unification caused high unemployment and slow growth. Compare German macroeconomic indices before and after the fall of the Berlin Wall.

NEAS Time Series Project Documentation

- Ireland had a poor economy with high unemployment and a low standard of living for three decades after World War 2. Liberalization of the economy in the 1980s's made Ireland the fastest growing European economy for the past 20 years, with low unemployment. Examine Ireland's macroeconomic indices, choose two periods that differ sharply, and fit an ARIMA process to each one.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

SEASONALITY: PROJECT TEMPLATE

(The attached PDF file has better formatting.)

Many candidates examine seasonality in their time series student projects. Your project can focus on seasonality, have a section on seasonality, or include a test for seasonality.

This project template summarizes seasonality issues that can be used in student projects. Most of this project template explains the types of seasonality and adjustment methods. The last page outlines a student project comparing the adjustment methods for CPI.

Take heed: We do *not* require the optimal seasonality adjustment for a student project. The sophisticated seasonal adjustments actually used for most economic time series are not covered in this course. We explain here the methods discussed in the textbook, and we indicate when each is appropriate.

Scope: Some student projects focus on seasonality, examining methods of identifying and correcting for seasonality or comparing two time series for seasonality. An analysis of ARIMA models with and without seasonal adjustments can form a good student project.

We examine seasonality in almost every time series. Even time series with no obvious seasonality, such as stock prices, over-night interest rates, inflation, claim severities, and business profits, may have seasonal patterns.

Recommendation: Include a section on seasonality in your student project. Form graphs of monthly averages (or quarterly averages) and 12 month sample autocorrelations, and explain whether they show seasonality.

Data: For a student project focusing on seasonality, choose data with a seasonal pattern.

- ~ The seasonality of interest rates is offset by the FED's monetary policy.
- ~ The global interest rate patterns of the past two decades have weaker seasonality, since many parts of the world do not share America's holiday shopping.

Interest rates are not a suitable time series for a student project on seasonality. But a good choice for a project on seasonality is the CPI indices on the NEAS web site.

- The non-seasonally adjusted index shows seasonality.
- The seasonally adjusted index does not show much seasonality.

Your student project might look at the following time series:

- The two CPI indices on the NEAS web site.
- The first differences of each CPI index.
- The non-seasonally adjusted index with a seasonality adjustment.

NEAS Time Series Project Documentation

- The first differences of this seasonally adjusted time series.

The data collection is minimal, since all the data are on the NEAS web site.

Weather: Projects on daily temperature and other weather patterns deal with seasonality. For daily temperature and rainfall, de-seasonalize the data. The project template on daily temperature discusses the procedures in more detail. Hourly temperature forms a good project on seasonality, since it overlays two patterns: 24 hour pattern and 365 day pattern.

Take heed: An excellent student project topic is rainfall in Los Angeles or Chicago or any other large city with heavy smog from weekday traffic. Some people say that the build-up of smog particles during weekdays and the dispersion of smog during week-ends causes a seven day seasonality in rainfall. You might search for smog levels in various cities: many large cities provide daily for hourly smog reports. You can see if smog has a seven day seasonality and fit an ARIMA process to a regression of rainfall on smog levels.

NEAS Time Series Project Documentation

ARIMA MODELING FOR CPI INDICES

The CPI index is not stationary. Your student project can fit ARIMA processes to the CPI, its first differences, or the first differences of its logarithms. You can use several methods to de-seasonalize the index.

Take heed: The CPI index has an exponential trend. Instead of first differences, take logarithms and first differences. Alternatively, take ratios of CPI index values in successive periods and logarithms of these ratios. The ratios have a lognormal distribution, and the logarithms of the ratios have a normal distribution.

Your student project explains how to test for stationarity, identify the type of trend, and derive the stationary time series.

- ~ Examine the seasonality in each time series: the CPI indices, your adjusted CPI index, and the first differences of each. You can also use the seasonally adjusted CPI, to see if that gives a better ARIMA fit.
- ~ Fit an ARIMA process to each time series.
- ~ Compare the ARIMA process for the seasonally adjusted index with the ARIMA process for the index with a simple seasonality adjustment.

The seasonal adjustments actually used for the CPI is more sophisticated than the simple seasonality adjustment in the textbook. An ARIMA process should fit better to the seasonally adjusted CPI than to the other indices. Compare the goodness-of-fit with the Box-Pierce Q statistic and Bartlett's test.

This is a simple student project, using data solely from the NEAS web site and the time series techniques in the textbook. For a time series of CPI price level figures, we expect an ARIMA (1,2,0) or an ARIMA (2,2,0) process.

- One difference converts the price level to the inflation rate. Don't forget to use logarithms if you see an exponential trend.
- If the inflation rate is a random walk, we take differences to form a white noise process.
- If the inflation rate has a trend, we take differences to eliminate the trend.

INFLATION AND THE MONEY SUPPLY

For a more sophisticated student project, you can examine the non-seasonally adjusted index divided by the money supply. The FED adjusts the money supply to remove the seasonality in interest rates. The change in the money supply causes a proportional change in the price level, though the effect is not immediate. Because the economic effects are not immediate and we have only rough monthly figures, you won't get strong relations.

INTERVALS

NEAS Time Series Project Documentation

Much seasonality in the time series used in student projects relates to end of the year shopping, summer vacations, and weather. With monthly data, you can identify spikes in January, August, December, or any other relevant month.

- For end of the year shopping, quarterly data are fine. Many business time series have quarterly figures, which differentiate higher sales in November/December from lower sales in January through March. For inflation, weather, and other business patterns, quarterly data doesn't provide fine enough information.
- For daily temperature, use days. See the project template on daily temperature for an example of fitting an ARIMA process to a highly seasonal time series.

Many time series have seasonal patterns: vacation travel in August, student job seekers in May through July, and weddings in June. Graph the monthly figures for your time series.

Insurance time series like claim frequency and severity have strong seasonal patterns in many lines of business. An analysis of seasonality for loss cost trends in personal auto or workers' compensation is valuable to insurers. Your student project will benefit your employer in addition to getting VEE credit.

GRAPHS

Graph the time series. Many monetary time series have an inflationary trend, so the seasonality may be more evident in the first differences than in the original series. If the trend is exponential, take logarithms and first differences. Form graphs of each series: the initial time series, its logarithms, and the first differences.

If the data are stochastic, fluctuations in the time series may obscure the seasonal pattern. We have two methods of identifying seasonality in stochastic data, corresponding to the two method of correcting for seasonality: multi-year averages and correlograms.

Illustration: We examine home sales by month in a small state. The values are stochastic, and home sales depend on economic conditions and interest rates, so the seasonality in any year is unclear. We use a 20 year average or a correlogram of 20 years of data to see the seasonality.

METHODS

Seasonality has several forms. The proper method of dealing with seasonality depends on the type of seasonality. Your student project can focus on the optimal method of correcting for seasonality. We contrast three scenarios, covering most time series.

Scenario #1: The current value depends on the time of year, not on the value at the same date one year ago.

Illustration: The expected daily temperature depends on the time of the year. It may be 25° on February 15 and 95° on August 15. It does not depend on the daily temperature one

NEAS Time Series Project Documentation

year ago. If the daily temperature is 45° on February 15, 20X8, we expect it to be 25° on February 15, 20X9, not 45° .

If we do not deseasonalize the data, we get high autoregressive (ϕ) parameters and sub-optimal ARIMA models. The illustration below shows the rationale.

Illustration: Suppose the proper model is an AR(1) process, where μ is the average daily temperature for that day, ϕ_1 is 20%, the average daily temperature varies greatly through the year (from 25° F to 95° F), and σ is high.

The δ parameter varies from 20° F to 76° F, depending on the time of year. Instead of a δ parameter, the ARIMA process might use a centered moving average of twenty values.

An ARIMA model may have $\phi_1 = 20\%$ and $\phi_2 = \phi_3 = \dots = \phi_{10} = \phi_{356} = \phi_{357} = \dots = \phi_{366} = 4\%$. We actually want $\delta = 80\% \times$ expected temperature for that day, which we estimate as a centered moving average of 20 surrounding days.

A better ARIMA model may have $\phi_1 = 20\%$ and values of 1% for 80 other lags, using 20 day moving averages from four years. An even better model may have $\phi_1 = 20\%$ and values of 0.1% for 800 other lags, using 20 day moving averages from 40 years.

This is wrong. Don't use autoregressive parameters to de-seasonalize the data.

- Separate the seasonality and the autoregressive process by deseasonalizing the data.
- Then fit an AR(1) or AR(2) process.

Scenario #2: Autoregressive Seasonality Parameter

The current value depends on the value N periods back, not on the time of the year.

Illustration: We model monthly auto insurance policy counts with an ARIMA process. The auto policy has a six month term, and 90% of policyholders renew their policies.

We don't expect policy counts to be higher or lower in any month. (This is not strictly true. There is a slight seasonality in auto insurance policy counts stemming from high auto sales in the fall and higher home sales in the summer, but the effect on policy counts is small.) But if policy counts were high in month j , they are high in month $j+6$.

Illustration: Auto insurance policy counts have no seasonality. Insurer ABC had a sales competition in September 20X6, with double commissions to agents with 40 or more new policies. Policy counts in September 20X6 were 20% higher than in other months.

With six month policies, renewals in March 20X7 will be 20% higher than in other months. This is not inherent seasonality: policy counts in March 20X6 were no higher than usual.

NEAS Time Series Project Documentation

The policy counts have no (significant) seasonality. The long-term monthly averages are about equal. The simple seasonal adjustment has no effect. But the autoregressive parameter for a lag of six months is high. We model the seasonality by the ARIMA process.

PRINCIPLES FOR COMPANY SALES VOLUME, BY TYPE OF PRODUCT:

For products without renewal purchases, we expect an autoregressive process with a high ϕ_1 parameter or a moving average process with a negative θ_1 parameter, depending on (i) the type of product and (ii) company sales vs industry sales. See the module postings for the modules on autoregressive and moving average processes.

For non-insurance products with high consumer loyalty, we have two scenarios:

- For short-duration products, like cigarettes and beer, high consumer loyalty raises the ϕ_1 parameter in a time series of monthly sales. We can't distinguish repeat purchases from autoregressive effects.
- For long-duration products, like autos and mobile phones, consumer loyalty is lower and the lag between purchases varies widely. If autos last three to ten years, a 50% consumer loyalty to a given model slightly raises the ϕ_j parameter for lags of three to ten years. But stochasticity in sales overwhelms these effects.

Property-casualty insurance has high renewal rates and 6 or 12 month terms. The ARIMA process for auto policy counts may have a ϕ_6 coefficient of 85% and a ϕ_1 coefficient of 5%.

The difference in the type of seasonality is seen in multi-year averages and correlograms spanning several years.

- If a high 12 month autocorrelation stems from weather patterns or holiday sales, the multi-year average shows the pattern more strongly.
- If the 12 month autocorrelation reflects consumer loyalty, the effect wears off over time.

If the insurer has an 85% renewal rate (retention rate), the renewal effect after ten years is $0.85^{20} = 3.88\%$. Even with a 95% renewal rate, the effect after ten years is $0.95^{20} = 35.85\%$. But if policy counts vary much from month to month by random fluctuations, the ϕ_6 autoregressive parameter may be 80% to 90%.

The correlogram shows the seasonality clearly. The correlogram shows the short-term 6 month sample autocorrelation over the past ten years. The average of 20 half-years eliminates the random fluctuations. The sample autocorrelations of lags 5 and 7 are close to zero, and the sample autocorrelation of lag 6 is 85%.

The difference in the type of seasonality is seen in the correlogram:

- For daily temperature, the 120 month sample autocorrelation is about the same as the 12 month sample autocorrelation.
- For policy counts, the 120 month sample autocorrelation is near zero.

NEAS Time Series Project Documentation

Scenario #3: Combination

Many seasonal time series are a combination of these two patterns.

Illustration: A sporting goods store sells both summer and winter sports equipment: surf-boards and bicycles in the summer and ski equipment in the winter.

The relative volume of winter vs summer goods varies over the years.

- A cold and rainy summer may depress sales of surf-boards.
- A cold winter with heavy snow may raise sales of ski equipment.

An ARIMA process of monthly sales is distorted by the seasonality.

- If we know the relative volume of summer vs winter goods, we use the relative volume to adjust for seasonality in sales.
- Our best estimate of the relative volume is last year's relative volume. But last year's monthly sales volume is affected by random fluctuations that do not repeat.

Illustration: Suppose the store gets much of its revenue from sales of bicycles in May, June, and July. If the weather last year was exceptionally rainy in May but clear in June, sales may have been low in May and high in June.

The optimal adjustment for seasonality may be a combination of a long-term average by month, and last year's sales by month.

Illustration: We examine a time series of unemployment rates for workers age 18 to 25 by month in Boston. Boston is a college town, with many students who look for jobs in May and June, so unemployment rates are high. In November and December, retail stores hire young workers for high holiday sales, so unemployment rates are low. The effects vary from year to year, depending on minimum wage laws and the types of employers looking for staff each year.

- Stores selling children's toys or video equipment have high demand for sales clerks in November and December.
- Firms offering part-time programming work hire college students in summers.

The stochasticity of unemployment rates affects the optimal seasonality adjustment.

- If the unemployment rates are not stochastic, we use last year's monthly relativities to adjust for seasonality.
- If the unemployment rates are highly stochastic, we use a multi-year average to adjust for seasonality.

NEAS Time Series Project Documentation

For many time series, the proper method to adjust for seasonality is unclear. Your student project can adjust a time series with both methods and compare the residuals.

Illustration: Begin with the monthly CPI with no seasonal adjustment. Use two models:

1. Seasonally adjust the data using average monthly relativities. Take logarithms and first differences, and fit an ARIMA process.
2. Take logarithms and first differences, and fit an ARIMA process with a 12 month seasonal term.

Choose *corresponding* processes to compare the two models.

- If the first model is an AR(2) process, the second model may be an AR(2) process with a ϕ_{12} parameter as well.
- If the optimal model differs for the two methods, choose one model. If the first model is an AR(2) process, and the second model is an AR(1) process with a ϕ_{12} parameter, use the AR(2) process for both methods.

Use Bartlett's test and the Box-Pierce Q statistic to select the adjustment for seasonality. One adjustment may lead to an ARIMA process with a lower Box-Pierce Q statistic than the other adjustment.

NEAS Time Series Project Documentation

TIME SERIES STUDENT PROJECTS: TIME SERIES TECHNIQUES

(The attached PDF file has better formatting.)

Updated: May 1, 2008

The SOA requires independent student projects for the regression analysis and time series courses.

- Statistics is not book knowledge alone. Knowing the characteristics of ARIMA models does not make one a statistician.
- Candidates must show they can apply the statistical techniques to empirical data to receive VEE credit for the statistics courses.

The student projects use

- Correlograms, Durbin-Watson statistics, Box-Pierce Q statistics, and Bartlett's test to specify, estimate, and diagnose ARIMA models.
- Multiple linear regression, residual plots, F tests, and dummy variables, plus
- Analyses of serial correlation, heteroscedasticity, and multicollinearity to formulate and test hypotheses.

The on-line courses do not require candidates to buy a statistical package, such as SAS or Minitab, which has these techniques. Buying a statistics package would add several hundred dollars to your cost. You may use Excel or similar spread-sheet software.

Excel is a powerful modeling package that handles number crunching, statistics, data base queries, and graphing. Its built-in functions and add-in features (data analysis and solver) provide most of the functions needed for the statistics courses. VBA macros can handle all the remaining statistical techniques you are likely to use. Most candidates use Excel at work, and they can immediately apply the techniques to data sets from public web sites.

Experienced statisticians often use *R*, a free on-line statistics package, which two of our faculty advisors use in their university courses. *R* is not user-friendly and has a steep learning curve. If you work with statistics in your actuarial jobs, the three to four weeks needed to become familiar with *R* are worthwhile. For most candidates, a good knowledge of Excel, its add-ins, and its VBA capabilities provide sufficient statistical software.

We do not provide VBA macros that do all the work for the student projects. ARIMA modeling is highly subjective. Based on the attributes of the time series, you may specify periods, adjust for seasonality, take first differences, or otherwise re-form the data. The SOA wants to ensure that candidates receiving VEE credit for statistics can apply the techniques to real data, not just run VBA macros that NEAS provides.

NEAS Time Series Project Documentation

Take heed: For some items, a VBA macro is much more efficient. Forming a correlogram from 50,000 records of daily temperature using simple cell formula may run slowly on your machine. If you are familiar with VBA, do the computations in a macro.

We provide a VBA macro for sample autocorrelations (correlograms) and the Box-Pierce Q statistic. Learn first the cell formulas, so you can modify the statistical tests for your student project. The explanation of the VBA macro is in the project template for weather student projects (daily temperature), which use time series of 40,000+ values.

NEAS Time Series Project Documentation

STATISTICAL TECHNIQUES

The time series student project requires linear regression, correlograms, residuals, and the Box-Pierce Q statistic. We provide Excel work-sheets with the time series techniques.

We use the *REGRESSION* section of the *DATA ANALYSIS* add-in for the linear regression and the residuals. We explain this add-in for the regression analysis student project.

Excel has a *CORREL* built-in function (correlation) but no sample autocorrelation function. This posting explains the difference between them. We do not require candidates to code the cell formulas for the sample autocorrelation function. The correlogram sheet on the illustrative work provide the cell formulas.

- Read this documentation so you understand the cell formulas.
- Copy the cell formulas to the work-sheet with your time series.
- You must replace the parameter showing the number of observations.

We have kept the cell formulas simple. If you have never used the *OFFSET* built-in function, read this documentation carefully.

Copy code from these worksheets for your student project. This course is statistics, not programming. You must know enough Excel to use the techniques, form graphs, and copy the output to your written document. If you had to write the Excel functions yourself, you would spend hours writing code. This is not an efficient use of your time.

If you know VBA or if you have worked with Excel data bases, use macros and data base techniques to simplify your project work. You may also use SAS or MINITAB or any other package. You do not have to write or create any functions from scratch.

ARIMA uses partial autocorrelation function and nonlinear regression for processes with a moving average component. We do *not* provide functions or macros for these tasks. Nonlinear regression is not covered in the regression analysis or time series course.

- If the objective of the student project were to find the optimal ARIMA model for a time series, you need nonlinear regression to optimize the moving average part.
- The student project shows that you can apply the course concepts to real data. The courses do not cover nonlinear regression, so we do not use it in the student project.

The course does not require knowledge of Excel. You can use another spread-sheet facility or a statistical package (SAS, Minitab, R) for the student project. You must know some package. We assume that all candidates know how to use Excel or similar packages.

Take heed: This documentation does not fully explain Excel built-in functions or VBA macros. If the cell formulas or the VBA macro is not clear to you, use the on-line help or Excel manuals for more explanation. Then post a question on the discussion forum.

NEAS Time Series Project Documentation

SAMPLE AUTOCORRELATIONS AND CORRELOGRAMS

The time series illustrative worksheets use sample autocorrelation functions, correlograms, the Durbin-Watson statistic, and the Box-Pierce Q statistic. None of these are built into Excel, but all of them are easily coded.

Take heed: The Excel built-in function *CORREL* computes the correlation of two random variables. The sample autocorrelations are not the same as the correlations used by Excel's *CORREL* built-in function. The differences are subtle, but they are important for statistical testing. The illustrative worksheet show the correlogram in three steps to clarify the differences between statistical correlations and sample autocorrelations.

- ~ The *CORREL* built-in function shows the shape of the correlogram, but it does *not adjust for degrees of freedom* and has *random fluctuations* that distort the statistical tests.
- ~ We can adjust the *CORREL* built-in function for degrees of freedom.
- ~ We use the *SUMPRODUCT* built-in function and the exact formula. This corrects for degrees of freedom and the random fluctuations.

The student project should use the exact formula. The two simpler versions help you grasp the intuition: why the *correl* built-in function is biased and how to correct.

Take heed: The sample autocorrelation formula still has a slight bias. The bias is too small to affect your results, and you need not concern yourself with it.

For your student project, copy the code using the *SUMPRODUCT* built-in function and the exact formula. Change the code to adjust for the number of observations.

- The code is general, and you can copy the code from one cell to any other cell. It uses Excel's *OFFSET* built-in function and a combination of relative and absolute references.
- If you have never used Excel's *OFFSET* built-in function and combinations of relative and absolute references, review the on-line help facility for these items. You can complete the student project even if you do not know how these items work, but you are likely to make errors. Ten minutes reading the on-line help may save you hours.

The correlogram is a chart. Most candidates do not need instructions for Excel charts. Use the chart wizard (if necessary) to create correlograms.

- We do not require a fancy chart.
- A line chart or a bar chart using the Excel chart wizard is sufficient.
- Create labels using the chart wizard. Edit the labels to make sure they are clear.
- If you have never used charts in Excel, you can document the chart in your write-up.

If you know how, add features to make your chart clear. The guidance here is a minimum.

NEAS Time Series Project Documentation

Illustration: To show seasonality in a correlogram, use markers or arrows for the seasonal autocorrelations. If you can't do so, explain in the write-up how the correlogram shows the seasonality.

NEAS Time Series Project Documentation

THE ILLUSTRATIVE WORKSHEET

Illustration: We use the last 42 months (3½ years) of 90 day Treasury bill interest rates on the NEAS web site: January 1997 – June 2000. This example shows how to use the Excel *CORREL*, *SUMPRODUCT*, and *OFFSET* built-in functions and the chart wizard.

- Copy Treasury bill interest rates for January 1997 through June 2000 from the spreadsheet on the NEAS web site to a blank spreadsheet.
- Place the interest rates in cells B11:B52. Leave one column on the left for line labels and ten rows on top for headers and documentation.

Enter Jan 1997 in cell A11 and Feb 1997 in cell A12. Select both cells and drag down to cell A52. Excel's *AUTOFILL* puts the proper months in these cells. Column A documents the work.

Take heed: The NEAS faculty members review the student projects. The relevant graphs and charts should be copied to your write-up, with references to the Excel spreadsheet where they are formed. If the reviewer can not understand what the worksheet does, we send it back to you for better documentation. This adds 4 - 6 weeks to the grading process. For your own sake, keep your work-sheets clean, and provide a clear write-up.

Format the cells with two decimal places. The formatting aligns the rates at the decimal place and makes them easier to read.

Enter the values 1 and 2 in cells C11 and C12. Select both cells, place the cursor at the lower right corner of cell C12, and drag down to cell C52. This forms a column with the value 1 through 42, using Excel's *AUTOFILL* procedure. Use this index, not the month names, for Excel's *OFFSET* function.

Place column headers in cells B10, C10, D9 and D10, as

- ~ B10: interest rate
- ~ C10: lag
- ~ D9: sample
- ~ D10: autocorrelation

(If you know how, join cells D9 and D10. Joining the cells makes the column headers look better; it is not necessary.)

We create correlograms three ways.

- Use the *third* method for the student project.
- The first two methods explain why the Excel built-in *CORREL* function does not give the sample autocorrelations. Follow the steps here, and you will avoid errors in your code.

NEAS Time Series Project Documentation

Note: Degrees of freedom are discussed in the regression analysis course. You do not have to understand the statistical theory in this posting. But you must copy the cell function for the sample autocorrelation to your work-sheet, and you must change the parameter for the number of observations in your time series.

Method 1: Built-in *CORREL* Function

Begin by entering the full correlation formula separately for two cells, so you see how the lag affects the cell formula. Then use the *OFFSET* function to simplify.

In Cell D11, enter the code `=CORREL(B11:B51,B12:B52)`. This is the correlation with a one period lag: the correlation of the first 41 values with the 41 values lagged one month. Use the first ten rows for comments: document your steps as you do them, and then copy the documentation to your write-up.

In Cell D12, enter the code `=CORREL(B11:B50,B13:B52)`. This is the correlation with a two period lag: the correlation of the first 40 values with the 40 values lagged 2 months.

Format the values in column D with four decimal places.

This procedure uses separate code for each cell. The final method writes the code once.

Take heed: Write the code by hand or select the cells. Write the name of the function or use the function wizard to select the function.

WRITING THE CODE ONCE: THE OFFSET BUILT-IN FUNCTION

A student project on daily temperature may have 50,000 observations. A student project on over-night LIBOR rates may have 3,500 observations. We can't write separate code for each cell.

Take heed: If you are not familiar with the *OFFSET* function, jot down on scrap paper the values in cells D11 and D12. Compare these two values with the values from the final version to make sure you have not made an error.

Erase the formula in cell D12. Change the formula in cell D11 to

`=CORREL(B$11:B51,B12:B$52)`.

The dollar sign makes the row number absolute. The formula asks for the correlation of the 41 values starting cell B11 with the 41 values ending cell B52

Copy this formula from cell D11 to cell D12. We get `=CORREL(B$11:B52,B13:B$52)`. This asks for the correlation of the 42 values starting in cell B11 with the 40 values ending in cell

NEAS Time Series Project Documentation

B52. This is incorrect. We want the formula `=CORREL(B$11:B50,B13:B$52)`: the correlation of the 40 values starting in cell B11 with the 40 values ending in cell B52.

We use the `OFFSET` function. Write the formula in cell D11 as

$$=CORREL(OFFSET(B$11,0,0,42-C11,1),B12:B$52)$$

The `OFFSET` function has five parameters. The formula here selects a range that

1. Begins in cell B\$11 with an offset of 0 rows and 0 columns
2. Has a height of $42 - C11$ rows and a width of 1 column.

The height of $42 - C11$ is $42 - 1 = 41$ for the correlation of lag 1. (Column C has the lags.) The formula is a relative reference. It changes to $42 - C12$ for the next row = $42 - 2 = 40$.

Copy this formula from cell D11 to cells D12 through D49. We don't use the last three cells:

- The correlation of lag $N-2$ (lag 40) is 1 or -1 . This is the correlation of (x_1, x_2) with (y_1, y_2) . The `CORREL` function gives a value, but this figure has no meaning. [Use the definition of the correlation from Module 1 of this course to find this value.]
- The autocorrelation of lag $N-1$ (lag 41) has a division by zero. This is the correlation of a scalar x with a scalar y .
- The correlations of lag N (lag 42) and higher are undefined. There are no values in the time series with lags this great.

The relative cell references adjust to the proper values for each correlation. Examine the formulas in the first several cells, so you see how the formulas change.

Compare your values with those on the illustrative worksheet. If they differ, review this documentation to find the error.

We use the `OFFSET` function for several of the statistical techniques. Excel has several alternatives to this function as well as VBA code that replicates it.

Take heed: This sample autocorrelation function has the proper shape, but the values differ slightly from those in the exact function. We form a correlogram and explain the problems with this sample autocorrelation function.

The cells in this column use the Excel built-in functions. Make sure you understand the `OFFSET` function and the correlation used for each lag.

Use the exact formula for your student project (Method #3). We use several columns for Method #3, but if you proceed through the steps here, you should understand the logic. Copy the code for Method #3 from the illustrative worksheet to your student project.

NEAS Time Series Project Documentation

CORRELOGRAM

The illustrative work-sheet has three methods to form correlograms. Use the third method for your student project. The first two methods explain the logic of correlograms.

- ~ One correlogram is formed directly from the Excel *CORREL* built-in function. This is the easiest correlogram to form, but it does not adjust for degrees of freedom and it has large random fluctuations at late periods. It is biased and can not be used for Bartlett's test or the Box-Pierce Q statistic.
- ~ A second correlogram adjusts for degree of freedom. It uses the *CORREL* built-in function and multiplies by $(N - k) / N$. It removes most of the bias in the first method. These first two correlograms show the difference between the sample autocorrelation and the correlation. Use the third correlogram for the student project.
- ~ A third correlogram uses the formula needed for the Box-Pierce Q statistic. The Excel work-sheet uses the *SUMPRODUCT* built-in function, not the *CORREL* built-in function.

Take heed: When you copy the cell functions, be sure to adjust the number of observations in the time series.

Form correlograms from the sample autocorrelation function. You may form half a dozen correlograms in your student project, corresponding to different versions of the time series.

The illustrative worksheet uses Excel's chart wizard. You may prefer to form charts directly, without the chart wizard. These instructions are for candidates not familiar with Excel.

Select cells C11:D49. Click on the *CHART WIZARD* and choose a *LINE GRAPH*. (You may also use bar graphs for the correlogram. Use whatever seems clearest.)

On the second wizard menu, the data range should be D11:D49, not C11:D49. Make this change manually. Alternatively, select cells D11:D49. Cells C11:C49 are your x-axis, which are used in other charts.

You can re-format any parts of the graph in the chart wizard or after making the chart.

- Eliminate the *LEGEND* on the right hand side, or rename it as sample autocorrelations.
- Give a title to the correlogram, such as *Correlogram of Interest Rate Time Series*.
- Label the axes, such as *Month Lag* for the horizontal axis and *Sample Autocorrelation* for the vertical axis.

Recommendation: These instructions form a simple chart. If you use Excel regularly, add titles, legends, markers, and other documentation to your chart before copying it to Word.

NEAS Time Series Project Documentation

DEGREES OF FREEDOM

The degrees of freedom affects the denominator of the correlation formula. A higher lag means fewer data points and fewer degrees of freedom.

Copy the lags from Column D to Column F. This is not essential, but new Excel users may find it easier to format the chart if the lags are next to the autocorrelations.

Recommendation: Learn Excel's Chart options. Your student project uses many charts and graphs, and graphics are equally useful for other actuarial reports. An hour spent learning to form charts will save you many hours of later work.

Enter the formula $=D11*(42-F11)/42$ in Cell G11. Copy cell G11 to Cells G12:G49.

Each cell in Column G (the revised correlations) is the Column D value $\times (N - k) / N$. This correlogram does not have the large random fluctuations at high lags and the decay is closer to a straight line. The *shape* of the correlogram remains the same.

NEAS Time Series Project Documentation

EXACT SAMPLE AUTOCORRELATIONS

The exact sample autocorrelation function differs slightly from the correlations above.

Take heed: The illustrative work-sheet uses simple Excel functions. The cells formulas can be made more efficient, but the formulas here are easier to understand.

Replace the interest rates by their deviations.

- Place the average interest rate in cell B9. Cell I11 has the formula =B11-B\$9.
- Copy cell I11 to the rest of this column: cells I12:I52.

Use the *SUMPRODUCT* built-in function for Column J, not the *CORREL* built-in function. Use the *OFFSET* function in the same fashion as for the *CORREL* built-in function.

- Cell J11 has =SUMPRODUCT(OFFSET(I\$11,0,0,42-C11,1),I12:I\$52).
- Copy this formula to cells J12:J52.

Take heed: It is easy to make an error with the *OFFSET* function. To avoid errors, compute the *SUMPRODUCT* function for two or three cells by explicitly referencing all the cells, and compare this value with the results using the *OFFSET* function.

Enter the cell formula =I11^2 in Cell K11. Copy this formula to Cells K:12:K52. Enter the cell formula =SUM(K11:K52) in Cell K9.

Column K has squares of the interest rate deviations. Cell K9 has the sum of the squares.

Enter the formula =J11/K\$9 in Cell L11.

The sample autocorrelations in column L are the *SUMPRODUCT* in Column J divided by the sum of the squares in cell K9.

Form a correlogram from the sample autocorrelations in Column L. This correlogram is smoother than the other correlograms.

- The correlogram using the *CORREL* built-in function has fewer terms in the denominator, so it is more distorted by random fluctuations.
- The adjustment for degrees of freedom adjusts the magnitude of the sample autocorrelations, but keeps the distortions from random fluctuations.
- The exact formula has more terms in the denominator, and it is less distorted by random fluctuations.

Summary: This worksheet explains how to form the correlogram and why it differs from ordinary correlations. For your student project, use the code for the exact autocorrelation formula. You do not have to explain the code in your student project.

NEAS Time Series Project Documentation

Take heed: This code is written for new Excel users. Experienced Excel users may prefer other cell formulas. You may write a VBA macro to form a correlogram from any sequence of data points. If you know VBA, this macro will save time and prevents typos.

The SOA wants candidates to show they can work with the statistical techniques. We don't automate the steps of the student project with macros, so that you work with the figures.

NEAS Time Series Project Documentation

USE OF THE CORRELOGRAM

Your student project forms a correlogram for each time series. Know how to interpret the sample autocorrelation function and the correlogram.

Take heed: The sample autocorrelation function and the correlogram have many uses. Statisticians differ on the implications of a sample autocorrelation function. We explain how we might interpret the sample autocorrelations of these 42 time series observations:

- ~ The statistical techniques are used in combination. We use these 42 points for the correlogram, an AR(1) process, the Durbin-Watson statistic, and the Box-Pierce Q statistic. Some implications are ambiguous.
- ~ The analysis in this workbook is incomplete. We explain what else to do for the student project, but we don't show every part. Using the data of 42 months, we should examine first differences, second differences, an AR(2) model, and perhaps an ARMA(1,1) process. An ARMA(1,1) process is not easy to fit using standard Excel functions, so this process is not required for the student project.
- ~ The textbook tries several complex models for interest rates, such as ARIMA(8,1,4). The student project does not require you to construct complex ARIMA models.
- ~ Statisticians differ in their interpretations of the statistical techniques, results, and plots. We give several explanations for the results here.
- ~ The results differ by time period. You may compare ARIMA models for one time period or one ARIMA model for two or three time periods. Dividing these 42 observations into two time series gives a different result.

Recommendation: When you begin the student project, do not worry about ARMA(1,1) processes. You can fit them several ways:

- Use a statistical package, such as SAS, MINITAB, or "R"
- Use the Yule-Walker equations to get an estimate of θ_1 .
- Use the *SOLVER* add-in and have Excel iterate for the solution.

After fitting AR(1) and AR(2) processes (or ARIMA(1,1,0) and ARIMA(2,1,0) processes), fit an MA(1) or ARIMA(0,1,1) process, using the Yule-Walker equations for an estimate of θ_1 . If you find this easy, fit an ARMA(1,1) or ARIMA(1,1,1) process. If this is not easy, do not worry about the ARMA(1,1) and ARIMA(1,1,1) processes. The student project gives you experience with using statistical techniques. New Excel users are not expected to code the sophisticated formulas needed for these processes.

NEAS Time Series Project Documentation

STATIONARITY

The correlogram drops to zero by a lag of 7 months. Similar correlograms appear in many student projects. Some candidates mistakenly infer that the time series is a stationary AR(1) process. Know the principles:

- The sample autocorrelations from a stationary AR(1) process decline geometrically to zero. It stays zero at subsequent lags, with random fluctuations that depend on the length of the time series.
- The sample autocorrelation here drops steadily to -47% by lag 14. The sample autocorrelation does not remain close to zero until lag 22.

Take heed: Some correlograms have the following form:

- Decline steadily to zero by lag N .
- Continue declining to a minimum by lag $2N$.
- Rise to zero by lag $3N$.

A stationary AR(1) process does not do this. A stationary process has autocorrelations close to zero after a few lags.

Illustration: An AR(1) process with ϕ_1 of 50% has an autocorrelation of $0.5^{10} = 0.10\%$, or a tenth of a percent after ten lags. With random fluctuations of perhaps 5% at each lag, we do not see a pattern in the autocorrelations after 4 or 5 lags.

Large ϕ_1 parameters do not give the pattern in this correlogram.

- An AR(1) process with ϕ_1 of 95% has an autocorrelation of $0.95^{14} = 48.77\%$, not -47% , at 14 lags.
- An AR(1) process with ϕ_1 of -95% has an oscillatory pattern about the mean, which is not the pattern in this time series.

The pattern here reflects a change in the mean or trend of the time series. A change in the trend is a change in the mean of the first differences.

Exercise 1.2: Change in Mean of Time Series

We simulate a time series of 100 observations.

- Observations 1 – 50 are random draws from a normal distribution with a mean of 2 and a standard deviation of 1.
- Observations 51 – 100 are random draws from a normal distribution with a mean of -2 and a standard deviation of 1.

A. Graph the time series.

NEAS Time Series Project Documentation

- B. For each era, form a correlogram. The sample autocorrelations have a normal distribution, with a mean of zero and a standard deviation of $1/\sqrt{100} = 0.100$.
- C. Form the correlogram from the entire series. The mean is zero, and the standard deviation is about 2. The distribution is not normal.
- D. Consider the sample autocorrelation of lag 1. The values for the pairs (1,2), (2,3), ..., (49, 50), (51, 52), (52, 53), ..., (99,100) are highly positive. The value for the pair (50, 51) is highly negative. The average for these 99 pairs is highly positive.
- E. Consider the sample autocorrelation of lag 2. The values for the pairs (1,3), (2,4), ..., (48, 50), (51, 53), (52, 54), ..., (98,100) are highly positive. The values for the pairs (49, 51), (50, 52) are highly negative. The average for these 98 pairs is highly positive, but lower than the average for the 99 pairs of lag 1.

NEAS Time Series Project Documentation

Graphs, Intuition, and Correlograms (Sample Autocorrelations)

The student project shows you can apply statistical techniques to empirical data and interpret the results to model a time series. Your written report shows three elements:

- The correlogram for the initial time series and any adjusted time series (first and second differences, seasonal averages, moving averages, logarithms).
- Graphs of the initial and adjusted time series.
- The reasoning linking the time series and its correlogram, explaining why the time series has the correlogram and what the correlogram implies for ARIMA modeling.

In this example, the student project would show

- The graph of the time series with a different mean or trend in two periods. A different trend in two periods causes a different mean of the first differences.
- The correlogram of the initial series or the first differences showing the pattern described here.
- The possible reason(s) for the change, if you are aware of them. The change in interest rates is explained in other postings.

You may not know why the time series changed. The student project does not require you to research the change.

Illustration: A project on crime rates in City Z may say: “The frequency of violent crimes increased from 1971 to 1992 and decreased from 1992 to 2006. The first differences show means of +0.06% in the first period and –0.02% in the second period. The correlogram shows decreasing sample autocorrelations for the first 14 lags, reaching a minimum of –22%, and then increasing toward zero in subsequent lags. I chose this time series to examine if the 1992 mayoral election on a crime fighting platform had any effect on the time series.

NEAS Time Series Project Documentation

OSCILLATION VS CHANGE IN MEAN OR TREND

{This sub-section comments on the shape of this correlogram. We include this sub-section so you see how to analyze a correlogram. Several other postings on this discussion forum discuss the shapes of correlograms.}

In this illustration, the interest rates decline in the first half of the series and rise in the second half. The decline and rise are not smooth because interest rates are stochastic, but a three month moving average shows a clear pattern. The sample autocorrelations for lags 8 through 26 are less than zero. This is a common pattern, which we model by dividing the time series into two periods and taking first differences.

Take heed: The sample autocorrelation functions in your student project may differ from the one in this correlogram.

- For a stationary AR(1) process, the sample autocorrelations begin at 1 and decline to 0.
- For a random walk, the sample autocorrelations begin at 1 and stay high.

Do not confuse this with an oscillatory model, where the sample autocorrelations alternate about the mean. The sample autocorrelations here decline and then rise.

- If this occurs repeatedly, the correlogram is cyclical; the authors say *sinusoidal*.
- This pattern occurs once here. It is a change in the trend, not an oscillatory model.

HIGH VALUES AT END

The *CORREL* function gives large values for the long lags, where the correlation is based on few values. The high correlations are random fluctuations.

Excel's *CORREL* built-in function shows the shape of the sample autocorrelation function.

- ~ An autoregressive model shows geometric decay
- ~ A moving average model shows a sharp drop
- ~ A white noise process shows small fluctuations

All stochastic time series show the white noise process about the expected values.

- ~ An autoregressive model shows random fluctuations about the geometric decay
- ~ A moving average model shows random fluctuations about the sharp drop
- ~ A white noise process shows small fluctuations about the mean of zero

The sample autocorrelation function has two changes.

NEAS Time Series Project Documentation

- ~ The correlation here is not adjusted for degrees of freedom; it has the same number of terms in the numerator as the denominator. The sample autocorrelation of lag k has k more terms in the denominator than in the numerator.
- ~ The correlation divides the sum of the cross-products by the product of the standard deviations of each series. The sample autocorrelation divides the sum of the cross-products, which has $N - k$ terms, by the sum of the squares of the elements, which has N terms. This second adjustment smooths much of the random fluctuations.

NEAS Time Series Project Documentation

Jacob: For the student project, do we just explain this reasoning or do we use graphs?

Rachel: Suppose we want to test whether the time series is an AR(1) model with a coefficient of 95%. See the worksheet with the AR(1) model for this coefficient. We form the autocorrelation function and the associated correlogram for an AR(1) process with a coefficient of 95% and compare that autocorrelation with the one in this worksheet. The two correlograms look different, indicating the time series is not an AR(1) process. The student project also uses other statistical techniques, as we explain below.

Jacob: The correlogram goes to zero eventually; does that mean the series is stationary?

Rachel: We have only 42 observations. The sample autocorrelations at high lags, such as the last 25% or 30% of the lags, may be small even for a non-stationary series.

Jacob: How would we deal with this series in the student project?

Rachel: We have several methods. They are used in combination. Statisticians have their preferred methods; no method is necessarily right or wrong.

- ~ Examine first and second differences of the observations. The pattern of this sample autocorrelation function suggests second differences may be stable. For the student project, take first and second differences. Examine the correlograms, explain what the correlograms imply, and explain in English what the first and second differences mean. We almost always examine the first differences of interest rates. Economists differ of whether we should use first or second differences; your student project can decide.
- ~ Deflate the interest rate time series (to real interest rates), adjust for business activity (GDP) and interest rate cycles (if any exist). The textbook speaks of structural models, or fitting an ARIMA process to the residuals of a regression analysis on another index.
- ~ Use higher order ARIMA processes. The textbook authors examine the correlogram and infer the *maximum* order of the ARIMA process. If the sample autocorrelations are near zero after 4 lags, they infer a maximum order of 4 and test various models. Most statisticians have the opposite perspective. They begin with an AR(1) model and work up to more complex models *only if needed*. We recommend the later approach for the student project. Your student project can decide if an AR(1), AR(2), or higher order model is appropriate.
- ~ We divide the time series into segments and fit ARIMA processes to each segment. These interest rates decline for about a year and a half and then increase. We might fit different random walk models to each segment. Your student project can decide if we should use one ARIMA process for the entire time series or different processes for different parts.

Take heed: Structural models and homogeneous time periods (segments) are preferred methods, but they require knowledge of the time series. Other postings on this discussion

NEAS Time Series Project Documentation

forum explain that first or second differences and ARIMA processes with more parameters may complicate a simple time series. Do not worry that your model may not be optimal. We review if you use the statistical techniques correctly, not if you form the optimal model.

NEAS Time Series Project Documentation

EXCEL REGRESSION BUILT-IN FUNCTION

The regression analysis on-line course shows how to fit regression lines with ordinary least squares estimators. For the student project, use the Excel REGRESSION built-in function.

Jacob: Where is the Excel REGRESSION built-in function?

Rachel: Choose the TOOLS menu from the menu bar. From the menu, choose DATA ANALYSIS. You may have to include the DATA ANALYSIS add-in to your version of Excel. From DATA ANALYSIS, choose REGRESSION.

Jacob: Simpler Excel built-in functions determine a linear trend line using regression. Can we use those built-in functions?

Rachel: We need the Excel REGRESSION add-in to get the table of residuals.

Jacob: How do we include the DATA ANALYSIS add-in?

Rachel: Check to see if the add-in is already installed. Some actuarial departments have the add-in installed. If the add-in is not installed, choose ADD-INS... from the tools menu. From the menu that appears, choose ANALYSIS TOOLPAK. To work with VBA, include also the ANALYSIS TOOLPAK VBA.

Jacob: What does the ANALYSIS TOOLPAK VBA give that the plain add-in doesn't have?

Rachel: You need the VBA version to invoke the add-in from VBA code. Most candidates do not need this facility.

Your version of Excel may differ. If you can't find the REGRESSION built-in function, post a question on the discussion forum, listing your version of Excel and of windows.

NEAS Time Series Project Documentation

DURBIN-WATSON STATISTIC

The Durbin-Watson statistic is covered in the regression analysis on-line course. It is a simple test of serial correlation. Coding and using the Durbin-Watson statistic is a good prelude to the Box-Pierce Q statistic. The Durbin-Watson statistic by itself is not a valid statistical measure for a lagged regression, as we use for autoregressive processes. Use it as a prelude to the Box-Pierce Q statistic.

Jacob: How do we form the Durbin-Watson statistic? Is there an Excel built-in function?

Rachel: Excel has no built-in function for this; we write the formula.

Jacob: The formula uses the residuals. How do we calculate the residuals? We can do this from the equations in the textbook, but it would take a while. Is there a simple method?

Rachel: The Excel *REGRESSION* add-in calculates the residuals. The add-in computes the ordinary least squares estimators and the residuals. Copy the formula for the Durbin-Watson statistic and the Box-Pierce Q statistic from the illustrative spreadsheet on the NEAS web site and use it with the residual output from the Excel *REGRESSION* add-in.

Form an AR(1) model from the last 3½ years of Treasury bill interest rates: January 1997 through June 2000.

Jacob: Must the time series be stationary?

Rachel: Even if the time series is not stationary, we can form the Durbin-Watson statistic for the residuals from an AR(1) model.

Start with the *REGRESSION* add-in. Copy the January 1997 – June 2000 Treasury bill rates to a new worksheet. Place these in cells B11:B52 and also in cells C12:C53. Column B is the Y values and Column C is the X values. We don't use the values in B11 or C53.

On the illustrative worksheet, we eliminate rows 53 and 11, getting rid of the original cell B11 and cell C53. This gives a matrix of B11:C51 for the regression analysis.

Jacob: If we use an AR(2) model, can we still use the *REGRESSION* add-in?

Rachel: To use an AR(2) model, place these rates in cells D13:D54 as well. Column B is the Y values, Column C is the X_1 values, and Column D is the X_2 values. We don't use the values in B11, B12, C12, C53, D53, or D54. We have 40 triplets (observations).

Jacob: What do we choose on the *REGRESSION* menu for the AR(1) model?

Rachel: The dependent variable is in cells B11:B51, after eliminating the original cell B11. The independent variable is in cells C11:C51.

NEAS Time Series Project Documentation

Ask for *RESIDUALS*. You don't need the *STANDARDIZED RESIDUALS* since the interest rates are all about the same size in this illustration. If interest rates change greatly over the time series, we would examine *STANDARDIZED RESIDUALS*.

Take heed: For time series analysis, we use first differences to eliminate trends. The time series values should not change materially. Use *RESIDUALS* for your analysis.

You can place the output on the same sheet or a new sheet. In a new worksheet, the residual output is in Columns A, B, and C. We place the output on the same worksheet starting in cell A61, so the output is in Columns A, B, and C, rows 85 through 125.

Take heed: The Excel default is a new worksheet. To use the same worksheet, over-ride the default and enter the upper-left cell of the output region on the *REGRESSION* screen.

The residual output shows the observation number, the fitted Y value, and the residual.

Take heed: The observation number is not the X value. For some analyses, like conditional heteroscedasticity, you may copy it the X values.

In column D, place the square of the residual. For cell D85, write $=C85^2$. Copy this formula to cells D86:D125.

In column E, place the difference of successive residuals. For cell E86, write $=D86-D85$. Copy this formula to cells E87:E125. Column E has one less figure than Column D has; this reflects the degrees of freedom in the regression.

In Column F, place the square of the difference of successive residuals. For cell F86, write $=E86^2$. Copy this formula to cells F87:F125.

Use Excel's quick sum function to get totals for Columns D and F. Place the cursor in cell D126 and click on the quick sum icon. Do the same for cell F126.

The Durbin-Watson statistic is the sum in cell F126 divided by the sum in cell D126. Place the Durbin-Watson statistic in cell G126 as the formula $=F126/D126$.

Jacob: What do we expect to find?

Rachel: A Durbin-Watson statistic of 2 indicates no serial correlation. This example gives a Durbin-Watson statistic of 2.10, which is not significantly different from 2.

- The correlation of lag 1 on the residuals, using the *CORREL* built-in function, is -0.06289 (cell C127), which is not significantly different from zero.
- The autocorrelation of lag 1 on the residuals, using the exact formula, is -0.06195 (cell C128).

NEAS Time Series Project Documentation

Other time series show serial correlation. Your student project should explain the meaning of your results.

Jacob: The slope coefficient is 95%. Is the t statistic significantly different from one?

Rachel: Be careful that you read the regression output correctly. This regression shows $\beta = 0.95$, its standard error = 0.07, and the t statistic is 13.66.

- This t statistic means we reject the null hypothesis that $\beta = 0$, not the null hypothesis that $\beta = 1$.
- For the null hypothesis of $\beta = 1$, the t statistic is $(0.95315 - 1) / 0.06977 = -0.480$. We do *not* reject the null hypothesis that $\beta = 1$.

NEAS Time Series Project Documentation

BOX-PIERCE Q STATISTIC

Your student project uses the Box-Pierce Q statistic

- To select among ARIMA models and
- To verify that a particular model fits the empirical data.

The textbook explanation of the Box-Pierce Q statistic is abstract.

- Use 15 to 40 sample autocorrelations.
- Ignore the first several sample autocorrelations.

The textbook does not specify precisely how many sample autocorrelations to use and how many to ignore. You understand the use of this technique after seeing its application.

The dialogue below explains how to use the Box-Pierce Q statistic. The illustrative Excel worksheet provides the code. Be sure you understand how to use the technique before applying it to the time series in your student project.

Jacob: The Box-Pierce Q statistic tests whether the time series is a white noise process. How do we use the Box-Pierce Q statistic for ARIMA processes, which are not white noise?

Rachel: We illustrate with the file of residuals from the AR(1) process.

Jacob: Do we apply the Box-Pierce Q statistic to the values of the time series (such as the interest rates) or to the residuals from an ARIMA model?

Rachel: If the time series values (interest rates) themselves show no pattern, we test if they are a white noise process. As an example, form the Box-Pierce Q statistic for the sample autocorrelations from the interest rates themselves on the *CORRELOGRAM* worksheet.

- The sample autocorrelations do not remain close to zero until lag 22.
- The Box-Pierce Q statistic is too high to reflect a white noise process.

Interest rates have random walks, trends, or mean reversion; we do not expect the interest rates themselves to be a white noise process. If an ARIMA model fits well, the residuals from the model may be a white noise process.

Take heed: Your student project may use the Box-Pierce Q statistic several ways. You may use the Box-Pierce Q statistic to identify a white noise process or a random walk.

- A time series may be a white noise process. Weekly rainfall in a tropical rain forest may show no seasonality or cycles. The weekly rainfall may be a white noise process.
- A time series may be a random walk, and its first differences are a white noise process. Stock prices are often assumed to be a random walk.

NEAS Time Series Project Documentation

If your chosen time series is a white noise process or a random walk, pick another time series for your student project. Use the Box-Pierce Q statistic to confirm white noise.

Sometimes changing the length of the periods gives a better model.

Take heed: Using periods that are too long causes may cause a time series to seem like white noise. Daily rainfall in a tropical rain forest may have strong autoregressive process (if it rains much on Monday, it may rain much on Tuesday as well) or a moving average process (if the rainfall was more than expected on Monday, it may be less than expected on Tuesday). If the effects die out after a few days, the weekly rainfall is white noise. Choose periods that show autoregressive or moving average processes.

Take heed: Use the Box-Pierce Q statistic to verify that your ARIMA model fits.

Jacob: If the residuals are a white noise process, is the ARIMA model correct? If the residuals pass the Box-Pierce Q statistic, have we fit the proper model?

Rachel: An ARIMA process is stochastic. The forecasts are never exact, since each value is a random variable. If the residuals are close to a white noise process, the ARIMA process *may* be a suitable model of the time series. We compare alternative models; we don't solve for an exact answer.

Take heed: The student project has no exact solution. No ARIMA process fits perfectly, and even the best fit may change from one year to the next. Model fitting combines science and art. You have wide discretion in choosing a model.

- One statistician may prefer an AR(2) process and another statistician may prefer an ARIMA(1,1,0) process.
- One statistician may fit an ARIMA(2,1,1) process to the entire time series and another statistician may fit ARIMA(1,1,0) processes separately to two periods.

We examine if you use the statistical tools properly and you proceed correctly through the model-fitting process. We do not judge if you came to the right answer.

NEAS Time Series Project Documentation

CODING THE BOX-PIERCE Q STATISTIC

Jacob: How do we form the Box-Pierce Q statistic in Excel? Does Excel have a built-in function, or do we code the formula ourselves?

Rachel: We provide the cell formulas, which you may copy to your student project. You select the discretionary items:

- How many values of K to use, and how many initial values to ignore.
- What significance level to use.

Take heed: The number of observations in the time series (T) is given. The Box-Pierce Q statistic uses K sample autocorrelations. We use 15 to 40 sample autocorrelations. The illustrative worksheet shows the Box-Pierce Q statistic for values of K from 1 to 40.

The choice of K depends on the number of time series observations. With more observations, we can use a higher value for K.

Illustration: With only 42 observations for the time series in the illustrative worksheet, we use 15 to 20 sample autocorrelations for the Box-Pierce Q statistic.

If the time series has many elements, we can choose a high value of K. But if the Box-Pierce Q statistic suggests the process is (or is not) white noise for $K = 40$, the indication probably won't change for $K = 80$.

Illustration: A time series of daily temperature from 1890 through 2005 has $116 \times 365.25 = 42,369$ observations. We can use as high a value for K as we like.

- If you use the Box-Pierce Q statistic to ensure that the time series is correctly de-seasonalized and then fit with an ARIMA process, you might choose $K = 365$ (assuming you have a table with the appropriate χ^2 values).
- If you first de-seasonalize the time series and then fit the ARIMA process, $K = 40$ or 50 is large enough.

Take heed: Be sure to estimate the sample autocorrelation function exactly. The Box-Pierce Q statistic is an approximation, but it is a good approximation with the exact formula. Using the wrong degrees of freedom for the sample autocorrelation function won't change your choice of ARIMA process. But if your time series has few observations, the Box-Pierce Q statistic will not be accurate.

Columns H – K show the autocorrelations formed with the *SUMPRODUCT* built-in function.

- ~ Column H shows the sum of the cross-products of lag k .
- ~ Column I is the square of Column H.
- ~ Column J is the sum of the first K terms in Column I.

NEAS Time Series Project Documentation

- ~ Column K divides Column J \times 41 (number of observations) by the sum of all 41 squared residuals.
- Place the formula =SUMPRODUCT(OFFSET(C\$85,0,0,41-A85,1),C86:C\$125) in Cell H85. Copy the formula to cells H86:H122. We don't use the last three lags, so we don't copy the formula to the last three cells. The *OFFSET* built-in function ensures that the *SUMPRODUCT* uses the correct autocorrelations. Be sure that your relative and absolute references are correct and that you use the proper number of time series observations.
 - Place the formula =H85^2 in Cell I85. Copy the formula to Cells I85:I122.
 - Place the formula =SUM(I\$85:I85) in Cell J85. Note the combination of relative and absolute references for a downward sum. Copy the formula to Cells J86:J122.
 - Place the formula =J85*41/D\$126^2 in Cell K85. Copy the formula to Cells K86:K122. The figures should increase down the column at a decreasing rate.

Jacob: We have 42 interest rates; why do we have only 41 residuals?

Rachel: We are using an AR(1) model, so we have $42 - 1 = 41$ residuals. The degrees of freedom for the Box-Pierce Q statistic is $K - p - q$.

Jacob: To what do we compare the figures in Column K?

Rachel: The Box-Pierce Q statistic is distributed approximately like a χ -squared distribution with $K - p - q$ degrees of freedom. If we use the sum of the first 15 terms, we compare to a χ -squared distribution with 14 degrees of freedom. If we use the sum of the first 20 terms, we compare to a χ -squared distribution with 19 degrees of freedom.

Jacob: How do we get these figures?

Rachel: Table 2 on page 604 of the textbook has the χ -squared distribution. At a 10% significance level, the critical χ -squared value is 21.06 for 14 degrees of freedom and about 27.2 for 19 degrees of freedom.

Take heed: For your student project, use the Excel built-in functions.

- We placed the formula =CHIINV(0.1,14) in Cell L99. This is the critical χ -squared value for 14 degrees of freedom at a 10% significance level. You can make the formula relative by writing =CHIINV(0.1,A99 - 1). You can copy the formula to every cell of column L. The illustrative worksheet shows the relative formula for two cells. You don't have to look up critical χ -squared values in the textbook chart.
- We placed the formula =CHIDIST(21.06414,14) in Cell M99.

Take heed: Explain in the write-up the procedure you used.

Jacob: The figures in the worksheet are lower. Can we conclude that (i) the residuals are a white noise process, (ii) the interest rates are an AR(1) process, and (iii) we have solved the student project?

NEAS Time Series Project Documentation

Rachel: The Box-Pierce Q statistic says that we can not reject the null hypothesis that the residuals are a white noise process. As the textbook authors note in their examples, many ARIMA models may be close enough fits that we can not reject the null hypothesis that the residuals are a white noise process. We compare the ARIMA processes and decide which one seems best.

A random walk is not stationary, but its residuals are a white noise process. This series is not stationary, so we take first differences and examine an ARIMA(0,1,0) process. We have not shown this on the illustrative workbook so that you can set up your own sheet.

NEAS Time Series Project Documentation

BOX-PIERCE FORMULA, T, K, AND DEGREES OF FREEDOM

Jacob: What is the formula for the Box-Pierce Q statistic?

Rachel: The Box-Pierce Q statistic is $Q = T \times \sum_1^K \hat{r}_k^2$.

Jacob: This formula seems strange. If we have twice as many observations T, won't the Box-Pierce Q statistic be twice as high? If the sample autocorrelations are about 10%, then twice as many of them gives a Box-Pierce Q statistic that is twice as great.

Rachel: If the sample autocorrelations are about 10% for all lags, the time series is not a white noise process. The Box-Pierce Q statistic tests whether the time series is a white noise process; it is not used for non-stationary time series.

If we have twice as many observations, the square of each sample autocorrelation for a white noise process is about half as great, on average. The *autocorrelations* for a white noise process are identically zero. A non-zero *sample* autocorrelation stems from random fluctuation. More observations (a higher T) smooths the random fluctuations.

Jacob: If we use twice as many lags, will the Box-Pierce Q statistic be twice as great?

Rachel: The increase tracks a χ -squared distribution with $K - p - q$ degrees of freedom. The increase is not proportional.

DOF's	1%	10%	90%	99%
10	2.56	4.87	15.99	23.21
20	8.26	12.44	28.41	37.57
30	14.95	20.60	40.26	50.89
40	22.16	29.05	51.81	63.69

Jacob: Why is the increase not proportional?

Rachel: The expected value of the Box-Pierce Q statistic increases proportionally for a white noise process. The probability that the Box-Pierce Q statistic lies outside a Z% confidence interval increases more than proportionally for $Z < 50\%$ and less than proportionally for $Z > 50\%$.

NEAS Time Series Project Documentation

QUALITY OF FIT: NON-STATIONARY RANDOM WALK VS STATIONARY WHITE NOISE PROCESS

Jacob: When I regress the *interest rates* on the *lagged interest rates* for the time period in my student project, using an AR(1) model on the monthly interest rates themselves, I get a β coefficient of one, indicating that the model is not stationary. The t statistic for β is high, the p -value is low, and the R^2 for the regression is high. The fit seems excellent.

When I regress the *first differences* of the interest rates on the *lagged first differences*, using an AR(1) model on the first differences, I get a β coefficient of zero, indicating that the model is stationary. But the t statistic for β is low and not significant, the p -value is high, and the R^2 for the regression is low. The fit seems poor.

I had expected the opposite results.

- If $\beta \approx 1$, the time series is a non-stationary random walk; we do not use it for forecasts.
- If $\beta \approx 0$, the time series is a stationary white noise process; we use it for forecasts.

Why does the random walk have a good fit and the white noise process have a poor fit?

Rachel: These are expected results. Consider what each test implies.

The t statistic tests the null hypothesis that $\beta = 0$. If the time series is a random walk, β is 1, not zero. We reject the null hypothesis. When we take first differences, $\beta = 0$; that is exactly correct. We do not reject the null hypothesis.

To test if the time series is a random walk, the null hypothesis should be $\beta_0 = 1$. This gives a t statistic close to zero, and we do not reject the null hypothesis.

The R^2 says how much of the variance is explained by the β coefficient. If the time series is a random walk with a low standard error (low σ), the R^2 is high. Values of 98% or 99% are reasonable for a perfect random walk with many observations.

If the β coefficient is zero for the AR(1) model of first differences, we expect the R^2 to be about zero. This says that the β coefficient doesn't explain anything. If the β is zero, it doesn't explain anything.

Forecasts: A random walk is not stationary, but we still use it for forecasts. For a random walk with a drift of zero, the L-period forecast is the current value. For a random walk with a drift of k , the L-period forecast is the current value + $L \times k$.

Jacob: For a perfect random walk, with $\beta = 1$ and a drift of zero, I assume the R^2 depends on the stochasticity. If σ is high, the R^2 should be low; if σ is low, the R^2 should be high.

NEAS Time Series Project Documentation

Rachel: That is true for most regression equations, where the values of X are not stochastic. For the AR(1) model, the X values are the Y values lagged one period. The dispersion of the X values varies directly with σ .

If σ is twice as large, σ^2 is four times as large, the sum of squared deviations of the X values ($\sum x_i^2$) is four times as large, the variance of $\hat{\beta}$ does not change, and the t statistic does not change. The R^2 is the σ^2 divided by the total sum of squares (TSS) of the Y values. The X values are also the Y values, so the R^2 does not change.

Jacob: That is counter-intuitive. Suppose the starting interest rate is 8%. If σ is 1%, which is high stochasticity for monthly interest rates, I presume the R^2 will be low. If σ is 0.01%, which is low stochasticity, I presume the R^2 will be high.

Rachel: To conceive of the relations, assume the starting interest rate is zero. The deviation is the actual value minus the mean, so we might as well start with a mean of zero.

- > If $\sigma = 1\%$, we have much random fluctuation in the Y values. This makes the X values more dispersed, which lowers the variance of the ordinary least squares estimator.
- > If $\sigma = 0.01\%$, we have little random fluctuation in the Y values. The X values are less dispersed, which raises the variance of the ordinary least squares estimator.

The two effects offset each other.

Jacob: Shouldn't the degree of stochasticity affect the R^2 ?

Rachel: The change from $\sigma = 1\%$ to $\sigma = 0.01\%$ as a change in the units of measurement. Nothing about the regression has changed.

NEAS Time Series Project Documentation

Student Project Variants

Updated: February 9, 2006

Jacob: The other postings give the considerations and techniques for the student project. What are possible topics for the student project?

Rachel: The student projects are of several types: fitting a model, comparing two or more time series, or applying a technique.

Fitting ARIMA Models

Jacob: What does a student project on fitting ARIMA models consist of?

Rachel: One type of student project fits an ARIMA model to a time series. We first graph the time series and select a time period that seems stationary or homogeneous non-stationary. We suggest eras in other postings.

The student project considers four types of ARIMA models: AR(1), MA(1), AR(2), and ARMA(1,1). For each of these, we consider both the interest rates themselves and the first differences, and we consider seasonality. This gives $4 \times 2 \times 2 = 16$ potential models.

Jacob: Do we form 16 models?

Rachel: We test for stationarity and seasonality to determine if we use first differences and a 12 month seasonal lag. We then try an AR(1) model and use diagnostic testing to see if it fits well. We can if it fits well, we are done; if it does not fit well, we try one of the other three models. If the student project shows you know how to fit a model, you need not repeat the fitting work, but you should explain how you would fit the other models.

Comparisons

Jacob: What would a student project doing comparisons consist of?

Rachel: You can compare a short rate and a long rate, showing the differences in ARIMA parameters, seasonality, and autocorrelations. You can compare nominal interest rates and real interest rates, or interest rates and interest rate futures. You can compare the three eras in the post-World War II period.

Techniques

Jacob: What would a student project on a technique consist of?

NEAS Time Series Project Documentation

Rachel: We discuss five ways of adjusting for seasonality in another posting. You can compare these five methods.

Jacob: When doing a student project on comparisons, do we just show the techniques?

Rachel: The comparison must show the goodness-of-fit tests using each technique. For each technique, we form the correlogram of the residuals and see which technique leads to white noise. Sometimes more than one technique works well; sometimes no technique works well.

NEAS Time Series Project Documentation

Time Series Graphics

Updated: February 9, 2006

Jacob: Why do we graph the time series? Why don't we start with the sample and partial autocorrelations and specify a model?

Rachel: Some items are easier to see in a graph than from numbers. We can separate a time series into eras, detect seasonality, observe changes in drift, volatility, and mean, and see cycles more easily in graphs than in tables of figures.

Jacob: What types of graphs do we form?

Rachel: We use four types of graphs:

We start with scatter plots of interest rates by month. These shows changes in mean, volatility, and drift; help differentiate eras; and show cycles.

Stochasticity overwhelms minor changes. We use moving average charts to eliminate stochasticity and better observe trends.

High stochasticity obscures weak seasonality. We use bar charts to eliminate the stochasticity and better detect seasonality.

Scatter plots do not show sample autocorrelations. We use correlograms to observe the sample autocorrelations.

Jacob: Can we do a student project showing these graphs for various types and time periods of interest rates?

Rachel: The graphs are a part of any student project; they are a tool, not the content of the student project. The student project should deal with some question, such as "What is a reasonable time series model for these interest rates?"

NEAS Time Series Project Documentation

Time Series Student Projects: Hypothesis Testing and Distributions

Updated: April 7, 2006

Jacob: What do the Durbin-Watson statistic, Bartlett's test, and the Box-Pierce Q statistic evaluate?

Rachel: These tests evaluate if the sample autocorrelations of the residuals are statistically different from zero.

Jacob: How do these tests work? I understand the formulas, but I don't get the intuition.

Rachel: To test the null hypothesis that the sample autocorrelation is zero, we need several assumptions about the distribution of sample autocorrelations.

Jacob: Regression analysis uses t statistics and p -values to test hypotheses that a value is zero. Do we make assumptions about distributions?

Rachel: Hypothesis testing in regression analysis generally assumes that the residuals are normally distributed with a constant variance.

Jacob: How do the assumptions affect the hypothesis test?

Rachel: Suppose the *sample autocorrelation* of lag 1 is 8% and we want to test the hypothesis that the autocorrelation is actually zero. The null hypothesis assumes that the sample autocorrelation has a normal distribution with a mean of zero, but we don't know the variance or the standard deviation (the square root of the variance) of the distribution.

If the standard deviation is 8%, the probability of the sample autocorrelation being greater in absolute value than 8% is about one third.

If the standard deviation is 4%, the probability of the sample autocorrelation being greater in absolute value than 8% is about 5%.

Jacob: How do we estimate the variance of this distribution?

Rachel: Bartlett's theorem says that the sample autocorrelations of a white noise process have a normal distribution with a variance of $1/T$, where T is the number of observations. If the residuals are a white noise process, their sample autocorrelations are normally distributed with a standard deviation of $1/\sqrt{T}$.

Jacob: How does Bartlett's test work?

NEAS Time Series Project Documentation

Rachel: If the sample autocorrelations have a normal distribution with a mean of zero and a standard deviation of $\hat{\rho}$, then the probability of an observed sample autocorrelation being greater in absolute value

~ than $1.96 \hat{\rho}$ is 5%.

~ than $1.45 \hat{\rho}$ is 10%.

Jacob: Suppose we find that a sample autocorrelation exceeds $1.96 \hat{\rho}$. Do we infer that the sample autocorrelations are not a white noise process?

Rachel: That depends on the lag, the type of data, and the other sample autocorrelations. Suppose $T = 400$, $1/\sqrt{T} = 5\%$, $\sigma = 5\%$, $1.96 \hat{\rho} \approx 10\%$, the data are monthly interest rates, and the sample autocorrelation of lag k is 12.5%.

~ If $k = 1$, we avoid drawing any inference. If the interest rates are correlations with the time period, the residuals may appear to be serially correlated, even if they are not. It is common for autocorrelations of lag 1 to be more than zero even if the autocorrelation is zero. The Durbin-Watson statistic would say that the data are inconclusive. This is like Bayesian estimation. The *prior* distribution for the autocorrelation of lag 1 has a high probability of being greater than zero because of other reasons. An observed value of 2.5 standard deviations is inconclusive.

~ If $k = 10$ and the sample autocorrelations of lags 1 through 9 are less than 10%, we presume the high sample autocorrelation for lag 10 is the expected random fluctuation in a normal distribution. The null hypothesis of a zero mean assumes that 5% of the sample autocorrelations have absolute values greater than 10%. One high sample autocorrelation out of ten is not unexpected.

~ If $k = 12$ and the sample autocorrelations of lags 1 through 11 are less than 10%, we suspect the high sample autocorrelation for lag 12 reflects annual seasonality.

Jacob: Bartlett's test sounds subjective; is this a problem?

Rachel: A skilled statistician prefers a subjective test that relies on our intuition about the time series. Seasonal correlations are more common than non-seasonal correlations. We may have subjective views on the types of autocorrelations that are most likely, such as positive vs negative autocorrelations.

Jacob: If we assume a normal distribution with a standard deviation of $\hat{\rho}$, the probability of a positive sample autocorrelation is the same as the probability of a

NEAS Time Series Project Documentation

negative sample autocorrelation. Why should the sign of the sample autocorrelation affect our inference?

Rachel: If we model monthly sales data and find a sample autocorrelation of 15% for a lag of 12 months, we presume this is annual seasonality. If the sample autocorrelation for a lag of 12 months is -15% , we may attribute this to random fluctuation.

~ If the actual autocorrelation is zero, the probability of positive vs negative sample autocorrelations is the same.

~ In practice, the incidence of positive autocorrelation is not the same as the incidence of negative autocorrelation. An AR(1) model with a negative parameter and an oscillating pattern is often a spurious result.

Jacob: Do we examine the Durbin-Watson statistic, Bartlett's test, and the Box-Pierce Q statistic on the residuals or the sample autocorrelations of the residuals?

Rachel: The Durbin-Watson statistic is applied to the residuals; the Durbin-Watson statistic calculates the autocorrelation of lag 1. The Durbin-Watson statistic $\approx 2 - 2 \times$ the sample autocorrelation of lag 1.

~ For perfect positive autocorrelation, the Durbin-Watson statistic = 0.

~ For perfect negative autocorrelation, the Durbin-Watson statistic = 4.

Bartlett's test and the Box-Pierce Q statistic are applied to the sample autocorrelations of the residuals.

Jacob: How does the Durbin-Watson statistic differ from Bartlett's test and the Box-Pierce Q statistic?

Rachel: Bartlett's test and the Box-Pierce Q statistic evaluate whether the residuals have a normal distribution with a variance of $1/T$, where T is the number of observations. They examine sample autocorrelations of various lags. The Durbin-Watson statistic examines if the sample autocorrelations of lag 1 are statistically different from zero.

Jacob: Are these tests equally strict?

Rachel: The Durbin-Watson statistic considers other factors that affect the observed sample autocorrelation of lag 1. Bartlett's test and the Box-Pierce Q statistic should be applied to many sample autocorrelations, not just the first one. Both these tests acknowledge that the sample autocorrelations for the first

NEAS Time Series Project Documentation

several lags many not follow the normal distribution or the χ -squared distribution in these tests.

NEAS Time Series Project Documentation

Time Series Projects: White Noise Process

Updated: April 7, 2006

Jacob: If the time series itself is a white noise process, have we made an error?

Rachel: Some time series are white noise processes.

~ If the time series is the number of earthquakes each *year* in the U.S., a white noise process is reasonable.

~ If the time series is the number of hurricanes each *month* in the Gulf Coast or the number of tornadoes in the mid-West, we don't expect a white noise process, since hurricanes and tornadoes have strong seasonality and possibly long-term cycles.

Jacob: What about interest rates? Might they be a white noise process?

Rachel: Interest rates themselves are not a white noise process. But if interest rates are a random walk, the first differences are a white noise process. Random walks are common in financial and actuarial work, so an ARIMA(0,1,0) model is reasonable.

Jacob: Are the interest rates on the NEAS web site a random walk that should be modeled as ARIMA(0,1,0)?

Rachel: For all eras combined, the interest rates are clearly not ARIMA(0,1,0). In the first era, the drift is positive, and in the third era, the drift is negative. But the interest rates in sub-periods may be a random walk, depending on the time period.

Jacob: If the observed interest rates have a drift, is the process not a random walk?

Rachel: Random walks can have drifts. Stock prices are often presumed to be random walks, but they have strong drifts.

Jacob: Do financial economists assume interest rates are a random walk?

Rachel: That depends on the time period, the country, and the actions of the central bank. One reason we use interest rates for the student project is that they are hard to model.

TIME SERIES PROJECTS MOVING AVERAGE MODELS

Updated: April 7, 2006

Jacob: How do we back into a moving average parameter?

Rachel: For an MA(1) model, the autocorrelation of lag 1 is $\rho_1 = \frac{-\theta_1}{(1 + \theta_1^2)}$. If we have

enough points, the sample autocorrelation is a good estimator of the autocorrelation, and we can back into the θ_1 parameter.

Jacob: Can you give an example of this?

Rachel: Suppose the sample autocorrelation of lag 1 is 50% and sample autocorrelations of lags greater than 1 are close to zero. We presume the time series is an MA(1) process and we compute the θ_1 parameter as $1 + \theta_1^2 = -2 \times \theta_1 \Rightarrow \theta_1 = -1$.

Jacob: Can you explain this model intuitively?

Rachel: Suppose the mean is zero. If $\theta_1 = -1$ for an MA(1) process, $y_t = \epsilon_t + \epsilon_{t-1}$. Since ϵ_t and ϵ_{t-1} are identically distributed random variables, $\rho(y_t, \epsilon_t) = \rho(y_t, \epsilon_{t-1}) = 1/2$.

Jacob: If we have only 20 observations; can we still back into the moving average parameter?

Rachel: A time series with only 20 observations is common; an example might be five years of quarterly numbers. We can back into the moving average parameter, but our estimate is not efficient. In the example above, θ_1 parameters between -1.2 and -0.8 give autocorrelations of lag 1 of about 50%. Even with hundreds of observations, we can't get an exact estimate of θ_1 . With only 20 observations, a θ_1 between -1.5 and -0.5 may give a sample autocorrelation of 50%. We must be careful with inferences from small data sets.

Jacob: Is the same true for the regression analysis used to estimate autoregressive models?

Rachel: The concept is the same, though the estimates may be more accurate. We should always compute the standard error of the estimator. Sometimes the textbook gives a formula for the standard error. These standard errors are not in the chapters for the time series course, so they are not required for the student projects.

Jacob: Are we supposed to examine moving average processes for the student project.

Rachel: Some candidates are confused by this.

NEAS Time Series Project Documentation

- ~ If the sample autocorrelation indicates an MA(1) model or an ARIMA(0,1,1) model, you should back into the MA(1) parameter and test if the residuals are a white noise process.
- ~ You do *not* have to test ARMA(p,q) models or ARIMA(p,d,q) models for which both p and q are equal to or greater than 1. You may comment if the sample autocorrelation function suggests a moving average component, but estimating the parameters requires nonlinear regression.

Jacob: If we have a statistical package like Minitab, should we examine moving average models?

Rachel: Statisticians differ on the value of moving average models. The textbook authors believe that all reasonable ARIMA models should be examined. Other statisticians examine simple autoregressive models and use moving average models only if a strong intuitive reason exists for these models.

Jacob: Why is it hard to identify moving average models? From the description in the textbook, they should be easier to identify than autoregressive models.

- ~ A moving average model has a high sample autocorrelation followed by sample autocorrelations of zero. This should be easy to spot.
- ~ An autoregressive model has an exponentially declining envelop about its sample autocorrelations; this should hard to spot.

Rachel: In the extreme scenarios, a moving average model can be spotted. If the sample autocorrelation is 50% for lag 1 and 0% for all other lags, the model is moving average. Even in this simple scenario, specifying the moving average parameter may be difficult in small data sets. θ_1 could range from -80% to -120% and give sample autocorrelations of 50% for lag 1.

In practice, the moving average component is usually weaker than the autoregressive component. If the sample autocorrelation of lag 1 is 80% and the following sample autocorrelations have an exponentially declining pattern of 50%, 40%, 32%, ..., it is hard to decide if this model is AR(1) or ARMA(1,1).

- ~ If the exponential decay is exactly 50%, 40%, 32%, and so forth, we might assume the sample autocorrelation of lag 1 is exact as well, and the model is ARMA(1).
- ~ If the exponential decay is roughly equal to 50%, 40%, 32%, and so forth, we might assume the sample autocorrelation of lag 1 is also not exact, and the model is AR(1).

Jacob: Are there statistical tests to decide this question?

Rachel: The decision depends on the intuition for a moving average component. If we have no reason to assume a moving average component, we prefer the AR(1) model. If we think a moving average component is likely, we may prefer the ARMA(1,1) model.

NEAS Time Series Project Documentation

OSCILLATING SERIES AND NEGATIVE AUTOCORRELATIONS

Updated: April 19, 2006

Jacob: Random walks, white noise processes, and mean reverting AR(1) models are often appropriate. What about oscillating series (that is, negative autoregressive parameters) and negative autocorrelations? Negative autoregressive parameters and negative autocorrelations seem unusual.

Rachel: They are indeed unusual. As the textbook states about the Durbin-Watson statistic, positive serial correlation is common but negative serial correlation is rare.

Jacob: How might negative autocorrelations occur?

Rachel: True negative autocorrelations occur in multi-period cycles. But multi-period cycles are uncommon, and they can rarely be modeled. Property-casualty underwriting cycles average 6 to 8 years in length, but no one has yet modeled them successfully. If we could predict the turning points or the severity of an underwriting cycle, the cycle would probably not follow the predictions.

Oscillating series and negative autocorrelations often stem from errors in the modeling process. These errors have occurred in some student projects for interest rate modeling. We review them here for candidates working on student projects.

- ~ Random fluctuations that *look like* oscillations (but are not)
- ~ Measurement error with short time intervals, a stable series, and a steady trend
- ~ Taking second or third differences when first or second differences are white noise
- ~ Combining time periods with upward and downward drifts
- ~ Modeling trend as a random walk

Jacob: If we observe an oscillating pattern or negative autocorrelations, what do we do?

Rachel: Consider first if the pattern is just random fluctuation. A white noise process has random autocorrelations, half of which are negative. A line graph connecting a white noise process might appear to be a cycle, since we tend to see patterns even in random data. Some candidates assume this is an oscillating series (which it is not), and deduce that an AR(1) model with a negative ϕ_1 is needed.

If more than half the sample autocorrelations are within one standard deviation from zero, you may be seeing patterns in a white noise process. Use Bartlett's test and the Box-Pierce Q statistic to check for a white noise process before you assume an oscillating series.

Jacob: What if the process is not white noise? What if the time series truly oscillates?

NEAS Time Series Project Documentation

Rachel: Check if you are taking frequent differences of a time series with a low drift. This occurs with daily differences of the Moody's investment yield rate.

DAILY VS MONTHLY RATES

Jacob: How does the time interval affect the observed patterns? What time intervals should we use for interest rates? Should we use daily changes or monthly changes in the interest rate? If we have daily figures, why not use them? Doesn't that improve the model?

Rachel: We consider interest rate volatility and the accuracy of measurement.

- ~ If rates are volatile and change frequently, we use shorter time intervals.
- ~ If interest rates are rounded and are not volatile, too frequent time intervals may cause spurious patterns.

The NEAS web site shows *daily* values for Moody's long-term corporate bond yield, since Moody's announces the rate each business day. You may combine rates into a monthly average, or you may use the rate on the first day of each month or week. Fitting an ARIMA model to the daily rates (or their first differences) may create spurious effects.

Suppose interest rates are a random walk with a drift of 0.5 basis points a day. Tomorrow's rate is forecasted as today's rate plus 0.005%. Interest rates increase by 2½ basis points a week, or $52 \times 0.025\% = 1.30\%$ a year. This is a reasonable interest rate process.

A random walk is not stationary. The first differences are a stationary white noise process. The mean of the white noise process is the drift of the random walk, or 0.005%.

If the interest rate is 8.002% on day 1, the expected rates for days 2, 3, 4, ... are

{8.007, 8.012, 8.017, 8.022, 8.027, 8.032, 8.037, 8.042, 8.047, 8.0452 ...}

We measure interest rates in basis points, or 0.01%. With two decimal place accuracy, the expected rates are

{8.01%, 8.01%, 8.02%, 8.02%, 8.03%, 8.03%, 8.04%, 8.04%, 8.05%, 8.05% ...}

Because of the rounding, the first differences are an oscillating series:

{0.00, 0.01, 0.00, 0.01, 0.00, 0.01, 0.00, 0.01, 0.00, ... }

We model this series as AR(1) with $\delta = 0.01$ and $\phi_1 = -1$, not as a white noise process. This gives a mean of $0.01\% / (1 - (-1)) = 0.005\%$, which is correct, but the oscillating pattern and the autoregressive parameter are not correct. This is a *spurious* relation; the true interest rate first differences are a white noise process with no oscillation.

Jacob: Would we see this in practice?

NEAS Time Series Project Documentation

Rachel: The actual Moody's long-term corporate bond yield has declined by about 25 basis points a year over the past two decades. This may create spurious effects if we fit daily first differences to an ARIMA model.

Jacob: Can't we spot this effect without much trouble?

Rachel: The effect of measurement errors is not always easy to spot. Suppose interest rates are a random walk with a drift of 0.4 basis points per day. The expected rate increases by 2 basis points a week, or 1% a year. The first differences are a white noise process with a mean of 0.004%. If the interest rate is 8.00% on day 1, the forecasts for days 2, 3, 4, ... are

{8.004, 8.008, 8.012, 8.016, 8.02, 8.024, 8.028, 8.032, 8.036, 8.040 ...}

With two decimal place accuracy, the forecasts are

{8.00%, 8.01%, 8.01%, 8.02%, 8.02%, 8.02%, 8.03%, 8.03%, 8.04%, 8.04% ...}

The first differences are an oscillating series:

{0.00, 0.01, 0.00, 0.01, 0.00, 0.00, 0.01, 0.00, 0.01, 0.00, ... }

The sample autocorrelations are negative for lags 1 and 3, positive for lags 2 and 4, and +1 for lag 5. We might model this series as an AR(5) process instead of white noise. The spurious process is not always simple year to year oscillation.

Jacob: Can we spot this effect from the consistent pattern?

Rachel: When we combine measurement error and stochasticity, the patterns are hard to spot. As an exercise, add a random white noise process to the pattern above. As the stochasticity increases, it becomes harder to see the pattern.

Jacob: Don't the Durbin-Watson statistic and the correlogram smooth the stochasticity and uncover the true patterns?

Rachel: If the sample autocorrelations are caused by real factors, the Durbin-Watson statistic and the correlogram smooth the random fluctuations and uncover the true pattern. If the sample autocorrelations are caused by measurement error, the Durbin-Watson statistic and the correlogram have the same effect: they smooth the random fluctuations and uncover the measurement error, which now looks like a true pattern. The statistical tests do not differentiate between measurement error and true patterns.

ONE DIFFERENCE TOO MANY

Jacob: You mention that taking too many differences can cause spurious effects. What does this mean? If first differences are a white noise process, aren't second differences also a white noise process?

NEAS Time Series Project Documentation

Rachel: If interest rates are a random walk with a drift of zero, the *first* differences are a white noise process with a mean of zero. We examine the time series process for the *second* differences.

Suppose the *second* difference in period j is positive, such as $+Z$. We estimate the likely *second* difference in period $j+1$.

- ~ The expected *first* difference in periods j and $j-1$ is zero, since this is a white noise process with a mean of zero.
- ~ The *second* difference in period j is positive, so the *first* difference may be high in period j or low in period $j-1$.
- ~ The pattern in a white noise process is symmetric, so our estimate of the *first* difference is $-\frac{1}{2}Z$ in period $j-1$ and $+\frac{1}{2}Z$ in period j .
- ~ The *expected first* difference in period $j+1$ is zero, so the *expected second* difference is $-\frac{1}{2}Z$.
- ~ The autocorrelation for lag 1 is negative.
- ~ This is AR(1) with $\phi_1 < 0$, or an ARIMA(1,2,0) process with $\phi_1 < 0$

Jacob: If taking differences can lead to spurious patterns, why do the textbook authors take differences so often?

Rachel: Stochasticity makes it hard to see the pattern in the first differences. The authors *examine* the second differences. They use them only if the first differences are not stationary and the second differences are more stable. If the first differences are a stationary time series, we don't use second differences.

LINEAR TREND AND RANDOM WALKS

Jacob: You mention trends vs random walks as another cause of spurious correlations. What does this mean?

Rachel: Stock prices, interest rates, commodity prices, business sales, and many other time series may be geometric or additive Brownian motion; that is, random walks with or without first taking logarithms. Some actuarial series, such as average claim severity, may be better modeled as linear or exponential trends.

Jacob: Average claim severity should be like stock prices and commodity prices. What causes this difference?

Rachel: The difference stems from sampling error. The prices for stocks and commodities are known with certainty, since the markets are liquid for these items. The expected average claim severity is estimated from a small sample of claims that occur.

- ~ If the claim severity trend is 10%, the expected claim severity in 20X8 is 10% greater than the expected claim severity in 20X7.

NEAS Time Series Project Documentation

~ If we don't know the expected severity in 20X7 with certainty, because our sample size is too small, we may estimate the expected claim severity in 20X8 from a trend line fit to several past years.

If we model a trend as a random walk, the residuals have an oscillating pattern.

Jacob: Isn't a random walk just a trend process with stochasticity?

Rachel: A random walk is different from trend. Consider a linear *trend of zero* vs a random walk with a *drift of zero*. The linear trend of zero is white noise, not a random walk.

We relate this to the previous discussion. If we presume the process is a random walk, we take first differences to get a white noise process. But since the process is actually white noise, taking first differences gives an oscillating AR(1) process with $\phi_1 < 0$.

Jacob: What if the series has a non-zero trend? Does that change the autocorrelation?

Rachel: If the interest rates have a trend of 10 basis points a period and the random walk has a drift of 10 basis points a period, the logic remains the same.

Jacob: The previous discussion shows that first differences of a white noise process are a random walk with a negative ϕ_1 . If we model a linear trend as a random walk, what happens?

Rachel: Consider first a linear trend of zero, which is a white noise process. Suppose the expected interest rate is 8% each month, but we model it as a random walk with a drift of zero. Let the actual interest rate be 8% in January 20X6.

The forecasted interest rate is 8% for February 20X6 as well. Interest rates are stochastic, so the actual interest rate may be higher or lower. Suppose the actual interest rate is 9%. The stochastic term for February 20X6 is +1%.

If we model the interest rate as an AR(1) process with no drift, the forecast for March 20X6 is 9%. The expected interest rate is actually 8%, so the expected error term in our model is -1%. The residuals have a negative serial correlation.

NEGATIVE SERIAL CORRELATION

Jacob: What should we do if we find negative serial correlation?

Rachel: Examine whether each of these spurious effects may be occurring. Long-term interest rates change slowly. If we model daily changes in twenty year rates with two decimal point precision, we expect spurious effects from measurement error.

Check whether you are taking too many differences. Some candidates presume that taking differences can't hurt. The opposite is true. If a time series is a white noise process, taking differences gives spurious results.

NEAS Time Series Project Documentation

Jacob: Suppose the sample autocorrelations are high. Shouldn't we take differences?

Rachel: We *examine* the differences; we don't always use them. Even if the time series is a white noise process, 10% of the time its sample autocorrelation function gives a Box-Pierce Q statistic above the 90% critical value.

Jacob: What if the differences give an AR(1) process with a negative ϕ_1 ?

Rachel: If the *differences* give an AR(1) process with a negative ϕ_1 , we go back and re-examine the sample autocorrelation function of the time series.

NEGATIVE SAMPLE AUTOCORRELATIONS

Jacob: What if the sample autocorrelations are negative for many lags in a row? This is not statistical fluctuation or measurement error, and it occurs often in the student projects on interest rates. What causes this?

Rachel: If the sample autocorrelation function is *positive* for many lags in a row, such as lags 1 through 30, we assume the interest rates follow a random walk. We sometimes find that the sample autocorrelations are *positive* for lags 1 through 20 and then *negative* for lags 21 through 40. The negative autocorrelation usually stems from a positive drift in the first half of the series and a negative drift in the second half of the series (or vice versa).

Jacob: Why does the drift change? Isn't our goal to model this change in the drift?

Rachel: Two types of causes lead to the change in drift; one is real and one is spurious.

Real: Some business cycles have automatic corrective effects. For example, some economists believe that prosperity builds up forces leading to recession, and recession builds up forces leading to prosperity. Similarly, rising interest rates build up forces leading to declining interest rates, and vice versa. If this is true, we should model the alternating interest rates eras as an interest rate cycle.

Spurious: Other economists do not believe interest rates have cycles. But changes in regulation caused by changes in regulators give the *appearance* of interest rate cycles.

Illustration: Suppose economist Y believes the Federal Reserve Board should raise interest rates to strengthen the economy, and economist Z believes the Federal Reserve Board should reduce interest rates to strengthen the economy. If the chairmanship of the FED changes from economist Y to economist Z, the rising pattern of interest rates may shift to a declining pattern.

Jacob: Don't we model this by an ARIMA process?

Rachel: We model endogenous effects by ARIMA processes, such as prosperous years creating forces that bring a recession. We do not model exogenous changes, such as a

NEAS Time Series Project Documentation

new FED chairman. Instead, we separate the time series into two eras and model each one by a simple ARIMA process.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

INTUITION, PARSIMONY, AND STATISTICAL TESTS: FITTING ARIMA MODELS

Updated: April 19, 2006

Jacob: How do we judge which ARIMA model is best?

Rachel: We want to forecast future values. The best model is the one that gives the lowest mean squared error of the forecasts.

GOODNESS-OF-FIT TESTS

Jacob: Why do we emphasize the Durbin-Watson statistic, Bartlett's test, and Box-Pierce Q statistic? These are in-sample goodness-of-fit tests; they do not test forecast accuracy. Shouldn't we use out-of-sample tests? Should we test which ARIMA model best forecasts interest rates over the next 12 months?

Rachel: Out-of-sample tests are important, but they don't always work well.

Illustration: Suppose interest rates follow a random walk with a drift of zero and a high stochasticity. The interest rate at December 20X7 is 8%. We compare three models:

- ~ A random walk with a drift of +5 basis points a month.
- ~ A random walk with a drift of zero basis points a month.
- ~ A random walk with a drift of -5 basis points a month.

Consider three scenarios, all of which are reasonable.

If the stochastic error term in January 20X8 is +30 basis points, the average forecast for the rest of 20X8 is 8.30%, and Model A gives the best forecasts.

If the stochastic error term in January 20X8 is -30 basis points, the average forecast for the rest of 20X8 is 7.70%, and Model C gives the best forecasts.

If the stochastic error term in January 20X8 is 0 basis points, the average forecast for the rest of 20X8 is 8.00%, and Model B gives the best forecasts.

Jacob: What does this illustration show?

Rachel: The actual interest rates for the 12 months of 20X8 are not independent. They are one realization of an interest rate path, not 12 independent realizations of future interest rates. The actual interest rates are one scenario out of many, and they do not confirm or disprove a model.

The in-sample goodness-of-fit tests are often good predictors of forecast accuracy. A model that has a low Box-Pierce Q statistic for its residuals is likely to have a low mean squared error for its forecasts. This is not *always* true, but it is reasonable.

NEAS Time Series Project Documentation

Jacob: For regression analysis, we use an R^2 test; why don't we use it for time series?

Rachel: For many ARIMA models, we take first differences to get a process that is close to white noise. The fitted process has a slope coefficient close to zero and a low R^2 . If the slope coefficient is zero, the R^2 is zero. We do not use R^2 to compare different types of ARIMA models, and we do not use R^2 when the slope coefficient is close to zero.

INTUITION

Jacob: How does intuition affect the choice of the optimal ARIMA model? We emphasize intuition, but the concept seems so vague. Why not just use the goodness-of-fit tests?

Rachel: Suppose we examine medical malpractice average claim severity over 20 quarters (5 years). We find an upward trend, which we can model as either a linear trend or an exponential trend. We might say: Let us consider two ARIMA models:

- ~ We take first differences to get a linear trend
- ~ We take logarithms and first differences to get an exponential trend

We use the model with the better in-sample fit to forecast future values. Is this what you mean by goodness-of-fit tests?

Jacob: Yes; if the trend is exponential, the exponential trend should fit better.

Rachel: Actuaries do not do this. If a \$10,000 claim this year will grow to an \$11,000 claim next year, we assume that a \$20,000 claim this year will grow to a \$22,000 claim next year, not a \$21,000 claim. We assume intuitively that the exponential trend is correct and will have the lower mean squared forecast error. If the linear trend has the better in-sample fit, we attribute that to random fluctuation. Strong intuition out-weighs the statistical tests.

Jacob: In general, which is more important: in-sample tests or intuition?

Rachel: That depends on the stochasticity of the time series values. If the stochasticity is low, we emphasize the in-sample goodness-of-fit tests. If the stochasticity is high, we emphasize our intuition about the reasonableness of each ARIMA model.

Jacob: How does this relate to interest rates? Do we know the structure of interest rates?

Rachel: We don't know the structure, but we assume that today's interest rate is a good predictor of tomorrow's interest rate.

- ~ Some financial economists assume we can't make better estimates.
- ~ Some financial economists believe that the drift in interest rates is likely to persist.
- ~ Some financial economists believe that interest rates are mean reverting.

NEAS Time Series Project Documentation

These views justify random walks or AR(1) models. We begin with an AR(1) model, either of the interest rates themselves or of their first differences.

Jacob: Isn't the assumption based on the empirical data? Why not fit the empirical data to as many ARIMA models as possible and see which one fits best? That will tell us which assumption is most reasonable.

Rachel: ARIMA modeling has two parts: we specify the model and estimate parameters.

- ~ The intuition is most important for the specification of the model.
- ~ The immediate past data are most important for estimating the parameters.

Illustration: Financial economists assume that stock prices in an efficient market follow a random walk. This assumption is based on arbitrage arguments, not on the past stock prices for any given stock. We use past data to estimate the drift and volatility of the stock price, not to choose among ARIMA models.

The same is true of interest rates. A financial economist might form an ARIMA model of twenty year Treasury bond rates in 2006 from observed Treasury bond rates in 2000-2005. The economist's intuition about the structure of the model may be based on long-term risk-free rates in many time periods and many countries, and the parameters are based on the twenty year bond rates in 2000-2005.

Jacob: What exactly does intuition mean? We say that ARIMA models should be intuitive; what makes one model more intuitive (more reasonable) than another?

Rachel: Intuitive means: "If we had no data about the particular item we are forecasting, what type of model would we use?" This is similar to a Bayesian prior: what do we expect before we examine the data?

Jacob: What ARIMA models are most intuitive (most reasonable)?

Rachel: The answer depends on the time series. Financial economists dealing with efficient markets use random walks. Actuaries dealing with random fluctuations use higher order autoregressive models. Marketing personnel emphasize seasonality. Economists dealing with industry demand for durable goods use moving average models.

Jacob: Can you explain the difference between financial economists and actuaries? (*Financial economists dealing with efficient markets use random walks. Actuaries dealing with random fluctuations use higher order autoregressive models.*)

Rachel: The stock market is efficient. The consensus among financial economists is that an ARIMA(0,1,0) model is optimal for stock prices. Economists examine past data to judge the drift and volatility of stock price movements, not to specify the type of model. We don't use stock prices for the student projects because the arbitrage arguments for an ARIMA(0,1,0) process are so strong that any other model is hard to justify.

NEAS Time Series Project Documentation

Jacob: Why do actuaries often use higher order ARIMA models?

Rachel: Financial economists modeling interest rates do not worry about sampling error. The interest rate may be rounded to the nearest basis point, and fluctuations in supply and demand may cause another 2 or 3 basis point errors, but this is minor. Actuarial statistics have great random fluctuations. Last year's personal auto claim frequency may be the best forecast for this year's claim frequency, but last year's frequency may be distorted by random fluctuations. Instead of a random walk, the actuary may use a three year weighted average, such as a 20%-30%-50% weighting. For low frequently / high severity items, the actuary may use a high order autoregressive process.

Jacob: What about interest rates? Are they like stock prices?

Rachel: The intuition for interest rates is less clear. Financial economists disagree about the optimal model, so interest rates are ideal for student projects.

Jacob: Why are interest rates different from stock prices? For stock prices, we assume the first differences are a white noise process (after taking logarithms). If the long-term average change in stock prices is 1% a month, our forecast for the February stock price change is 1%. If stock prices rose 2% in January, we don't change our forecast for February.

Interest rates have a long-term average movement of zero. The average risk-free interest rate over the past decade is about the same as it was 50 years ago and 100 years ago. Just because interest rates rise (or fall) in one month, why should we assume they will rise (or fall) in the next month?

Rachel: Compare interest rates and the weather. The long-term average temperature on May 1 in Chicago doesn't change much from decade to decade. The simplest forecast for the May 1 temperature is the average May 1 temperature over the past 100 years. If this average temperature is 45°, our time series is $45^\circ + \epsilon$, which is a white noise process.

An ARIMA model does better than the white noise forecast. If the temperature on April 31 was 55°, we might forecast 50° for May 1. If the temperature on April 31 was 35°, we might forecast 40° for May 1. This is an AR(1) model with $\phi_1 = 50\%$.

Jacob: This justifies an AR(1) model, which is a reasonable model for both interest rates and local temperature. Is there any justification for a higher order model?

Rachel: If the temperature on April 31 was 45°, the AR(1) model forecasts 45° for May 1. Perhaps we can do better.

- ~ If the temperature on April 30 was 40°, we might presume a warm front is moving in, and our forecast for May 1 may be higher than 45°.
- ~ If the temperature on April 30 was 50°, we might presume a cold front is moving in, and our forecast for May 1 may be lower than 45°.

NEAS Time Series Project Documentation

We might form an AR(2) model to represent this process. Interest rates may be similar.

- ~ The average interest rate may be 8% over the past century. The simplest model is a white noise process, with a mean of 8%.
- ~ Interest rates change slowly, so we may give 80% weight to last year's interest rate and 20% weight to the long-term average of 8%. This is an AR(1) model with $\phi_1 = 80\%$.
- ~ Interest rates are not traded in an efficient market. If interest rates were 8.4% last month, 8.3% the previous month, and 8.2% the month before, we might presume that interest rates are rising. Numerous economic reasons may cause this: the economy may be recovering from a recession with high demand for investment funds or a budget deficit may create an over-supply of Treasury securities. An AR(2) process may forecast more accurately than an AR(1) process.

Jacob: If the cause of the rising interest rate affects our forecasts, shouldn't we use an econometric model, not an ARIMA model?

Rachel: ARIMA models are proxies for the true economic, financial, or other explanation. An economist who knows why interest rates are rising may make better forecasts than a statistician who uses ARIMA models. But economists understand interest rate movements imperfectly, and statisticians using ARIMA models may make reasonable forecasts.

PARSIMONY

Jacob: How does parsimony fit into this?

Rachel: First differences, second differences, AR(1), AR(2), and MA(1) models have intuitive rationales. They reflect economic influences on interest rates. If these influences had certain effects in past years, we assume they have similar effects in future years.

We have no intuition for the effects of residuals or interest rates from 3 or 4 periods back. We assume the past relations are random fluctuations. If we include them in the model, they are more likely to raise the mean squared error than to lower it.

Jacob: It sounds like parsimony is a part of intuition; is that true?

Rachel: Yes. Parsimony does not mean that we prefer a linear trend to an exponential trend. We base the choice on the intuition for linear vs exponential trends. Parsimony means: if we have no explanation for a moving average or autoregressive factor, leave it out.

INTUITION VS DATA

Jacob: If the intuition is so important, why do we focus on graphs and correlograms?

Rachel: We rarely understand the causes of interest rate movements. Many economic factors affect interest rates, and the possible ARIMA process are varied. This makes interest rates ideal for the student projects.

NEAS Time Series Project Documentation

An example is seasonality. Some candidates say there is no justification for seasonality in interest rates, so they treat any patterns as random fluctuations. A financial economist would observe the high demand for money in December and wonder why the seasonality is not stronger. An experienced economist would explain that interest rates were highly seasonal 80 years ago, but the Federal Reserve Board now adjusts the supply of money to mitigate the fluctuations in interest rates.

The same is true of many economic indices. Inflation, unemployment, GDP growth, and similar items are seasonal. A student project on any of these items might explore how the correlogram allows us to measure the seasonality. We are expanding the project templates to include other subjects. We have a conservative perspective: we introduce a new project template when candidates are comfortable with the current templates. Too rapid expansion leads to unnecessary confusion.

NEAS Time Series Project Documentation

TS: Note on Positive and Negative Autocorrelations

Many student projects have correlograms of the first differences showing

- ~ Slowly declining positive sample autocorrelations for the first N lags from Z to 0.
- ~ Slowly declining negative sample autocorrelations for the next N' lags from 0 to Z' .
- ~ Slowly rising negative sample autocorrelations for the next N'' lags from Z' to 0.

Often, $N' \approx \frac{1}{2} N$ and $Z' \approx -\frac{1}{2} Z$. The magnitude of N'' varies.

The correlogram has a clear pattern, but this pattern is not discussed in the textbook. Many candidates infer that the time series is not stationary, and they take second differences. They fit ARIMA processes to the second differences. Their models do not fit well and they forecast poorly.

This pattern of sample autocorrelations indicates two time series with different means. We should divide the full time series into two time periods. Taking second differences hides the problem and creates a poor model.

This error occurs most frequently in student projects on GNP (Gross National Product) and CPI (inflation). This posting explains this pattern and why it occurs.

Suppose GNP (or GDP) grows 3% a year for 15 years and then 1% a year for the next 15 years. We examine a time series by calendar quarter of annualized GDP.

This example reflects actual U.S. experience. GDP grew at 3.5% a year in the 1950's and 1960's, drifted down to about 1.5% a year by the second half of the 1970's, and then rose back to 2.5% or 3% by the late 1980's, where it remained for much of the 1990's. Actual GDP varies by quarter, begin low in recessions and high in prosperous years.

Economists assume that long-term GDP growth reflects numerous conditions, including marginal tax rates, international trade and currency restrictions, government policies, level of development, education and labor force participation of women, and a host of other items. The time series course emphasizes that we should first use a structural model, such as regression analysis, to explain GDP growth from these exogenous conditions. We then use an ARIMA process to explain the residuals.

The student project does not expect a sophisticated structural model. We encourage you to form simple regression models, but these are not required for the time series course. For GDP growth, economists do not agree on the proper explanatory variables, so a structural model is difficult to form.

For this illustration, we do not use a structural model and we use the simplified pattern of 3% and 1% per annum with low stochasticity. Stochasticity and business cycles affect the correlograms, but the pattern described here shows up in the actual results for U.S. GDP.

NEAS Time Series Project Documentation

We examine GDP by quarter for 30 years, or quarters 1 through 120. To simplify the mathematics, we add one more value at the beginning (at quarter 0), so we have 60 values of 3% growth and 60 values of 1% growth.

We model GDP by an ARIMA process. GDP is not stationary, since it grows exponentially by improvements in worker productivity. We use real GDP, not nominal GDP, so inflation does not affect the ARIMA modeling.

Productivity growth is multiplicative (exponential), not additive. We take logarithms of real GDP to convert the exponential trend to a linear trend. We take first differences to eliminate the trend. We get a time series of 120 values: 60 of +3% and 60 of +1%.

We first convert this time series into its deviations from the mean of 2%. We get 60 values of +1% and 60 values of -1%. To simplify the mathematics, we use units of 1 percentage point, so we have 60 values of +1 and 60 values of -1.

We form the correlogram for lags of 1 through 119. The denominator of the sample autocorrelation function is $120 \times 1^2 = 120$.

The numerator of the sample autocorrelation function depends on the lag.

~ *Lag 1*: 59 values of $+1 \times +1 = +1$; 1 value of $+1 \times -1 = -1$; 59 values of $-1 \times -1 = +1$; the sum of $59 - 1 + 59 = 117$.

~ *Lag 2*: 58 values of $+1 \times +1 = +1$; 2 values of $+1 \times -1 = -1$; 58 values of $-1 \times -1 = +1$; the sum of $58 - 2 + 59 = 114$.

The sample autocorrelation decreases by 3 for each additional lag. Using the N and Z in the first lines of this posting: $N = 39$ (≈ 40) and $Z = 117/120 \approx 1$. In practice, variations in the GDP growth rates and stochasticity of GDP growth each quarter reduce Z and N. Even in clear instances of different trends for two perspective, Z is usually about 40%.

~ *Lag 40*: 20 values of $+1 \times +1 = +1$; 40 values of $+1 \times -1 = -1$; 20 values of $-1 \times -1 = +1$; the sum of $20 - 40 + 20 = 0$.

~ *Lag 41*: 19 values of $+1 \times +1 = +1$; 41 values of $+1 \times -1 = -1$; 19 values of $-1 \times -1 = +1$; the sum of $19 - 41 + 19 = -3$.

~ *Lag 59*: 1 values of $+1 \times +1 = +1$; 59 values of $+1 \times -1 = -1$; 1 value of $-1 \times -1 = +1$; the sum of $1 - 59 + 1 = -57$.

~ *Lag 60*: 0 values of $+1 \times +1 = +1$; 60 values of $+1 \times -1 = -1$; 0 values of $-1 \times -1 = +1$; the sum of $0 - 60 + 0 = -60$.

~ *Lag 119*: 0 values of $+1 \times +1 = +1$; 1 value of $+1 \times -1 = -1$; 0 values of $-1 \times -1 = +1$; the sum of $0 - 1 + 0 = -1$.

NEAS Time Series Project Documentation

Taking second differences appears to solve the problem. We get a time series of 59 values of zero, one value of -2 , and 59 values of zero. With a stochastic time series and a slowly changing mean, the second differences may even form a stationary time series which passes Bartlett's test and the Box-Pierce Q statistic for a white noise process.

But this time series gives incorrect forecasts, since it assumes the mean first difference is $+2\%$. The true GDP growth rate in the second period is $+1\%$. We forecast 1% for the immediately following quarters, not 2% .

The proper analysis is that the time series process changed after the first 15 years, from a 3% trend to a 1% trend. The forecasts should use a 1% trend, not a 2% trend.

The textbook is not always clear about these issues, since it is hard for the statistician to judge if the time series has changed. The textbook assumes that you have examined the time series and the explanatory variables that affect the process. You have chosen a time period for which the time series is homogeneous (whether it is stationary or non-stationary). If the time series has changed because of some exogenous explanatory variable, you are using the residuals from a structural model, and these residuals form a homogeneous time series.

The time series section of this textbook does not much discuss exogenous variables that make the time series heterogeneous. Other chapters of the textbook cover this topic, such as the chapter on dummy variables that is covered in the regression analysis course. The authors assume you have applied the needed statistical techniques to make the time series homogeneous.

The student projects do not require you to use residuals from a structural model or to make the time series homogeneous by other means. We simply select a time period in which the time series is homogeneous.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

BUILDING ARIMA MODELS: A STEP BY STEP GUIDE

Updated: April 22, 2008

(The attached PDF file has better formatting.)

The postings on the discussion forums provide guidance for your student project. We describe what each posting covers, and we suggest an order for your initial review.

The student projects are independent projects. The NEAS web site has hundreds of data sets and various project templates that you may use for the student project. You may use any time series with enough observations, as long as it is not a random walk or white noise.

This posting is a step-by-step guide to ARIMA modeling. Separate postings

- explain the requirements for the student project and the learning objectives
- outline the written documentation that accompanies the statistical work
- review the statistical techniques for student projects on time series analysis
- compare common ARIMA processes and suggest which ones to explore
- document the illustrative worksheets for the project template on interest rates
- clarify the balance between faculty guidance and independent work
- answer questions from candidates about the time series student projects

The instructions summarize questions and answers in past semesters. They provide more guidance than most candidates need. You have wide latitude: you may choose the data, topics, and statistical procedures, and write a student project that differs from the project templates on the discussion forums.

Your student project applies statistical techniques to real data. It is not a full statistical study to select the optimal ARIMA model. You need not compare all ARIMA processes or test all structural models. We review if your project properly uses modeling techniques, not if your solution is optimal.

This posting is discursive. It discusses each step with many examples. It refers to other discussion forum postings for further explanation. It is not a cookbook to a rigid sequence of statistical techniques

Some candidates want to know an exact order for the student project, and they are frustrated by the subjectivity of ARIMA modeling. As the course textbook says, ARIMA modeling is a second-best alternative. We search for patterns in the data. We don't find the true causes of the time series patterns, but the ARIMA pattern helps our forecasts.

Your student project is successful if it sheds light on the pattern over time of the data. Read the suggestions in this step-by-step guide and apply them to your time series.

ARIMA MODELING

Chapter 19 of the textbook has five examples of ARIMA models. The text assumes you

- are familiar with nonlinear regression and partial autocorrelation functions
- have the needed statistical software
- understand well the time series that you are modeling.

The text gives an outline of ARIMA models: specification, diagnosis, and validation.

The time series on-line course assumes

- you know linear regression and sample autocorrelation functions
- you have Excel, but not other statistical software
- you have not worked with the data except for the student project

The ARIMA models for the student project can be built with basic Excel functions. The illustrative work-sheets on the discussion forum provide code for the statistical techniques. If you understand the concepts, you can complete the student project without difficulty.

This step-by-step guide to building ARIMA models is geared to candidates who have taken the time series on-line course but who have never worked with ARIMA models.

- ~ We use plain language for most of the explanation.
- ~ When we use a statistical procedure, we explain the steps you need.

A step-by-step guide is not a cookbook. You make decisions at each step. We give enough guidance that you won't get stuck, but you form the model.

- We do specify exactly what to do at each step. The empirical data are stochastic, and they do not fit any model exactly. Differences from the model may reflect random fluctuation or a poor model.
- We explain what to consider at each step. Your analysis depends on the empirical relations and your statistical judgment.

The student project is an educational process, not a consulting job.

- Working through the techniques helps you learn them.
- We examine if you can apply techniques to actual data, not if your solution is correct.

Read this step-by-step guide when you begin the ARIMA modeling part of the time series course, so you know the techniques you must master.

Illustration: The Box-Pierce Q statistic seems vague in the textbook, but it is a basic tool to validate ARIMA models. Validation is subjective, since ARIMA models rarely fit the data

NEAS Time Series Project Documentation

exactly. The student project uses the Box-Pierce Q statistic, Bartlett's test, intuition, and parsimony to select an ARIMA model.

See also: The illustrative work-sheets provide templates for correlograms, Durbin-Watson statistic, Box-Pierce Q statistic, Yule-Walker equations, and statistical tools. You need not have statistical software or Excel expertise to complete the student project.

Take heed: Statisticians differ in their approaches.

- The text uses complex ARIMA processes, often with six to eight parameters.
- For the student project, use simpler processes, with two or three parameters.

The project demonstrates that you understand the concepts, not that you can form a model with many parameters.

NEAS Time Series Project Documentation

Step #1: *KNOWLEDGE*

The student project assumes you know time series analysis and regression techniques.

- ARIMA modeling seems complex at first, but the modeling sequence is logical.
- Each step requires a decision; a limited set of decisions leads to many combinations.

Below are topics in the time series on-line course that are used in the student project.

1. Identify *stationary* time series. *Most time series that actuaries use are homogeneous non-stationary.* A time series with a trend or drift, or a random walk even with no drift, is not stationary. Examine the correlogram, check for unit roots, and graph your results.
2. Form a stationary time series: take first differences, using logarithms if appropriate; divide the time series into periods; and correct for seasonality.
3. The *autocorrelation function* of an ARIMA process is not the *sample autocorrelation function* of a time series. Use the *sample autocorrelation function to specify the model* and the *autocorrelation function to validate the model*.
4. The textbook discusses partial autocorrelation functions as well. You do not have the statistical software for the partial autocorrelation function, and the discussion in the text is weak. You need not use the partial autocorrelation function for the student project.
5. The pattern of the correlogram reflects the type of ARIMA process. Focus on *geometric decay vs sudden drops*. Low ARIMA parameters and high standard errors of short time series obscure the pattern. A correlogram is hard to analyze if stochasticity is high. If a pattern is not obvious, explain the pros and cons of alternative models. We judge if you understand the reasoning, not if you choose the optimal model
6. Understand the *intuition* for moving average vs autoregressive models. Much practical statistical work is subjective. A non-intuitive ARIMA process with good in-sample fits may have poor out-of-sample results. Explain in lay terms what each parameter implies.
7. Chapter 19 builds ARIMA models for several time series. The student project does the same with less complex time series and ARIMA processes.

As you work through your project, review the course modules for these topics.

NEAS Time Series Project Documentation

Step #2: *TOOLS*

You need statistical software for the student project. The illustrative Excel worksheets

- Have cell formulas and examples for statistical functions not built into Excel.
- Have VBA macros and custom functions that simplify your work.

The items below are provided on the illustrative work-sheets.

- *Sample autocorrelation function*: Excel has a *CORREL* built-in function but no built-in function for the sample autocorrelation. An illustrative worksheet explains the difference and gives the cell formulas, using the Excel *SUMPRODUCT* and *OFFSET* built-in functions.
- The *SUMPRODUCT* with *OFFSET* is slow for large time series, such as 40,000 or 50,000 days. We provide a VBA macro that is quicker and simpler.
- *Correlogram*: Use the chart wizard to construct correlograms. Label your axes so our faculty can review your work. Copy and paste the correlograms into your write-up.
- *Durbin-Watson statistic* is taught in the regression analysis course and used in the time series student projects. An illustrative worksheet gives you the code.
- Use the *Box-Pierce Q statistic* to test the model, using the sample cell formulas in the illustrative worksheet. Compare the Q statistic with the critical values for the χ -squared distribution, using Excel's built-in function or the tables in the textbook.
- Use linear regression to form autoregressive models. Use the Excel *REGRESSION* add-in and use the residual output for the Box-Pierce Q statistic.

Take heed: We explain the code in the illustrative worksheets and the rationale for each procedure. You may copy the cell formulas from the illustrative worksheets.

The cell formulas in the illustrative worksheets are simple. Experienced Excel users will find *SOLVER* and VBA to be more efficient tools. Nothing in the student project requires more advanced Excel knowledge than in the simple cell formulas.

Your write-up states the results in your worksheet. *Our faculty can not figure out what you have done from the Excel workbook alone*. State the techniques you use and the results. Explain what the results imply and how you test them for significance (when appropriate).

You can use any statistical software or any spread-sheet package. If you use SAS at work, you can save time by using it for the student project. If you know VBA and Excel built-in functions, they can save you time as well.

SAS, MINITAB, and "R" have all the built-in functions you might use in a student project. You may use these software packages; they are not required.

NEAS Time Series Project Documentation

Step #3: CHOOSE THE TIME SERIES

You can use any time series you want. We show ARIMA models for interest rates and daily temperature, and structural models for various macroeconomic indices.

Many project templates suggest a variety of student projects. Use the project templates to generate ideas for your own student project.

Choose a topic that interests you. The web has dozens of sites with statistics on almost any topic.

- If climate change intrigues you, do a project on daily temperature or rainfall.
- If you are a sports fan, do a project on won-loss ratios of your home team.
- If you like music, do a project on monthly DVD sales by genre.

We suggest numerous topics for your student project.

Illustration: Interest rates have a hundred flavors. We show sample work with 90 day Treasury bill rates, and extracts from student projects on other rates. You can choose

- short rates (three month bills, over-night rates) vs long rates (twenty year bonds).
- private rates (Moody's corporate bond index, the bank prime rate) or government rates (Treasury securities, bank discount rates).
- spot rates, forward rates, futures rates, or other derivatives.
- nominal interest rates or real interest rates (structural models)

The type of rate should reflect the analysis.

- For interest rate seasonality, use short rates, not long rates, such as over-night LIBOR.
- For the relation of interest rates and budget deficits, use real interest rates.
- For the relation of interest rates and recessions, use corporate spreads.

The NEAS discussion forum has many interest rate time series. Read the project templates and the other discussion forum postings. Feel free to choose another type of rate, such as risk-free rates in other currencies.

Take heed: It is often easier to model real interest rates, the residuals of interest rates on inflation rates, corporate spreads, or the spread between long and short rates with ARIMA processes. Spend an extra half hour setting up the data; you save hours in your analysis.

Suggestion: You read dozens of theories about interest rates and other macroeconomic indices: real interest rates are higher or lower in recessions, higher or lower when the U.S. runs a deficit vs a surplus, and so forth. Choose a hypothesis, form a structural model, and fit an ARIMA process to the residuals.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

Step #4: *STRUCTURAL VS TIME SERIES MODELS*

If the time series is a by-product of *clear and easily accessible* causes, we use regression analysis. For a stochastic time series with its own internal logic, we use ARIMA models.

Illustration: Unemployment rates depend on economic, demographic, and legislation, such as hiring practices, restraints on firing, unemployment benefits, and minimum wages.

- If the legislature raises the minimum wage, teen-age employment drops.
- If the state raises unemployment benefits or mandates employer provided health insurance, unemployment rises.
- During recessions, unemployment rises.

The macroeconomics on-line course reviews these effects. Barro's textbook is an excellent source of ideas for time series and regression analysis student projects. Government web sites have extensive data on economic variables.

Take heed: The residuals from structural models are easier to fit with ARIMA processes, and the time series are more meaningful.

We use regression analysis if the explanatory factors change frequently. We use ARIMA models for the *residuals* of the regression.

Illustration: Inflation rates affect interest rates. We may regress interest rates on inflation and use an ARIMA model to forecast the residuals. We provide both interest rates and inflation rates on the web site so you can model either nominal interest rates or real interest rates. Chapter 19 of the textbook has a similar example.

Take heed: Many macroeconomic indices are functions of other indices. Use differences, periods, or structural models.

Illustration: Nominal interest rates are a function of inflation. Inflation is not mean reverting, so nominal interest rates are not stationary. Using first differences and dividing the time series into periods creates stationary time series, but your student project will be better if you use real interest rates and adjust for economic activity.

Illustration: Suppose you want to model interest rates.

- Divide the one month LIBOR by the CPI for the previous period to get the real LIBOR.
- Real GDP is the detrended GDP.
- Regress the one month LIBOR on real GDP.

Your student project may show a sequence of models.

NEAS Time Series Project Documentation

- Fit nominal interest rates by taking first differences. Use the mean squared error over the next 12 months to estimate goodness-of-fit.
- Convert to real interest rates and fit a new ARIMA process. It is more difficult to decide if first differences are needed. Re-compute the mean squared error.
- Use a regression on real GDP.

Your student project may explain the advantages and drawbacks of a structural model. You may have a good model of real interest rates, but if you don't know future inflation and real GDP, you can't forecast future interest rates.

You may determine real interest rates three ways:

- Rate A minus Rate B.
- Rate A divided by Rate B.
- Rate A regressed on rate B.

The regression is the best procedure, since it combines additive and multiplicative models. The *REGRESSION* add-in does the regression and gives the residuals. But you can use any of the three methods.

The regression can be done using different inflation rates as the explanatory variable. Choose one and explain the rationale. This is a statistics course, not an economics course. You are not graded on the choice of the inflation rate.

We use residuals for real interest rates, maturity spreads, and corporate spreads. The definitions below use *Rate A minus Rate B*, but you can use any of the definitions above.

- Real interest rate = nominal interest rate minus expected inflation.
- Maturity spread = long risk-free rate minus short risk-free rate.
- Corporate spread = corporate bond rate minus Treasury bond rate.

The NEAS web site has many time series in Excel format. Form the time series you want. Use simplifications, even if they are not perfectly accurate.

Illustration: Use last month's actual inflation as a proxy for expected inflation. The ratio of the CPI in the two previous months as the expected inflation for the current month.

Don't worry that your time series is not perfect. Construct the time series you want and fit an ARIMA process. But be consistent. To analyze corporate spreads, use the average rate in each month, or the corporate bond rate at the start or middle of each month.

Take heed: Do not worry that a time series on the residuals of a regression is too complex for the student project. The opposite is true. Many macroeconomic and demographic time series are too complex to fit with an ARIMA process. The residuals of a regression analysis are easier to fit, and they are more likely to have an intuitive relation.

NEAS Time Series Project Documentation

The project templates discuss structural model for many of the time series on the NEAS discussion forum. For some models, you must find an appropriate explanatory variable on the internet. Spend half an hour or an hour looking for time series on the internet that fit the hypothesis you want to examine. We do not grade your success in finding the right explanatory variables for a structural model.

NEAS Time Series Project Documentation

Step #5: *TIME PERIODS*

We fit an ARIMA process to model a time series over a given period. If the mean, drift, or variance of the time series changes because of *external causes*, we use different models for the different parts of the time series. Statisticians speak of *interventions*, or exogenous events that change the ARIMA process.

- If changes in the time series are random fluctuations, we use a single process to model the underlying structure. If we change models each year, we can't forecast future rates.
- If the time series itself changes, we need separate models. Forcing a single model to cover all years gives an ARIMA model that does not fit well in any period.

A time series may follow different ARIMA process in different periods. If exogenous factors change the mean, variance, or drift of the time series, the time series is not stationary and can not be modeled by a single ARIMA process. Examples:

- We model an insurer's premium volume in 1985 - 2005. For 1985 - 1995, the insurer has a monopoly; in 1996, the market becomes competitive. Premium *volume* may be high when the insurer has a monopoly and lower when the insurer competes. The *variance* of the premium volume may be low when the insurer has a monopoly and high when it competes.
- We model airline passenger volume before and after deregulation. Greater competition and lower fares after deregulation raise the industry's passenger volume. The variance of any carrier's passenger volume increases: some new carriers rapidly gain market share and some established carriers fail.
- We model sales, profitability, and cash flow of firms differently in their start-up phases and their mature phases.
- Oil prices have different time series for the pre-OPEC era (before 1973) and the OPEC era (1973 onwards).

Interest rates have both types of changes.

- ~ Rates are stochastic, varying from month to month. The ARIMA process models these fluctuations.
- ~ Federal Reserve Board policy (monetary policy), federal budget deficits (fiscal policy), domestic and foreign capital investment, economic growth, and perhaps trade balances affect the mean, drift, and variance of interest rates. An ARIMA process for the 1960's may not be a good model for the 1980's.

It is not always easy to identify external causes. Even simple questions, such as warming vs cooling of the earth, are much debated.

We use two methods of dividing a time series into periods:

- We examine the means, drifts, and variances of the time series itself.

NEAS Time Series Project Documentation

- We examine the exogenous events that might change the ARIMA process.

Illustration: An actuary examines a time series of personal auto written premium from 1980 to 2008. An acquisition of a personal auto subsidiary in 1995 writing in different states changes the time series. We use separate models for 1980-1994 and 1996-2008.

Examine the time series you choose. You may divide it into two or three periods based on the attributes of the time series, even if you have no explanation.

We don't expect you to know post-World War II Federal Reserve Board policy or other events that affect interest rates. For the student project, examine the means, drifts, and variances of the time series. Select periods that have reasonably stable attributes.

We provide some basic information about post-World War II Federal Reserve Board policy to explain the differences observed in the time series. Just as actuaries examine policy provisions, distribution systems, and market competition to set optimal rates, a statistician should know the attributes of the time series to fit an ARIMA process. But the student project focuses on the statistics, just as the SOA and CAS exams focus on the actuarial procedures. You are not graded on your knowledge of economics.

- From the end of World War II (1945) through the mid-1970's, the U.S. economy expanded briskly. Government officials worried about Depression-era deflation, not the mild inflation of an expanding economy. Inflation was thought to be an antidote to unemployment, which had been high during the Depression. The federal government and the Federal Reserve Board believed that mild inflation was beneficial, in that it restrained unemployment and did not hamper economic prosperity.

This presumed relation of inflation and unemployment was an error, but it was the prevailing macroeconomic policy in the 1960's and 1970's. Interest rates had a steady upward trend; the time series is not stationary.

- From the late 1970's through early 1980's, inflation and interest rates were high and volatile, resulting from (i) the mistaken macroeconomic policies of these times and perhaps (ii) the supply shocks of OPEC oil price increases.
- Paul Volcker, who became chairman of the FED in 1981, adopted a monetarist perspective (Milton Friedman's views). The money supply grew at a steady, slow rate. Interest rates and inflation rates declined steadily (downward drift). Greenspan continued Volcker's policies. The interest rate patterns in the first and second periods should not recur if the FED continues a monetarist policy.

No single ARIMA process is an appropriate model for all three periods. *For the first part of the student project, you select appropriate periods.*

- Global daily temperature has long periods (thousands of years). We have had ice ages and warm eras; different models are appropriate for each. A student project on daily

NEAS Time Series Project Documentation

temperature over the past 130 years may examine if a trend exists and if the trend rate has changed. Many web sites on global warming have information about changes in the daily temperature. A student project may compare the daily temperature time series in a period of no trend vs a period of trend.

- If the time series is unemployment rates, the time periods depend on legislation for hiring practices, restraints on firings, unemployment benefits, and minimum wages. U.S. law stayed relatively constant over the past fifty years; European law provided increasing liberal benefits. You might compare U.S. and European unemployment.
- Interest and inflation rates depend on monetary and fiscal policies. In Europe, policies changed with entry into the European Union. Use European, Asian, or Latin American rates to make your student project different.

We use separate models for each period; we don't take differences and use a single model. We use two or three *interest rate eras* because of different Federal Reserve Board policies.

For the student project, you can rely on internal characteristics of the time series. Inspect the graphs for the time series and choose the time periods. We show *illustrative* graphs for three interest rate periods:

- Period 1: January 1945 – December 1978
- Period 2: January 1979 – December 1982
- Period 3: January 1983 – June 2000

These are *illustrative* periods. For your project, examine the data and choose periods.

The periods need not be contiguous. You can leave gaps. You can choose January 1945 – December 1978 as Period 1 and July 1979 – June 1982 as Period 2. This leaves an 18 month gap between the periods. During the gap, policies are changing, and no ARIMA process may properly model interest rates.

If the drift is steady for several years and then reverses for several years, a single ARIMA process doesn't work. Taking second differences is not proper, since the first differences form stationary time series in adjoining periods. Instead, you can

- Use separate ARIMA processes for the two periods.
- Form *real interest rates* and see if a single ARIMA process works for both periods.

The second method is a structural model: form a regression and use the ARIMA process on the *residuals*. Use this method if an exogenous factor affects the time series.

Take heed: Shifts in daily temperature are not well understood.

A shift in the mean *suggests* an exogenous intervention. Unemployment rates may be 6% for several decades followed by a rise to 11% for a decade, as in France and Germany. The cause may be higher unemployment benefits and restrictions on work terminations. Unemployment rates have since declined in Europe, and you may have three periods.

NEAS Time Series Project Documentation

The discussion board graphs three month Treasury bills and suggests rough interest rate eras. For the student project, graph the time series and choose time periods.

- Do not just copy the three periods on the discussion board. Examine the graph of your time series and explain whether two or more time periods are needed.
- A project *comparing* time periods may use two or more periods. A project *fitting* an ARIMA process may focus on one period.
- You can leave gaps between periods or ignore some periods. A student project can compare ARIMA processes for the first and third periods on the discussion board.

A time series with different means in different segments is not stationary. Separating the time periods is necessary to create stationary time series, but it is generally not enough. For several reasons, a time series may not be stationary.

- A time series with an upward or downward drift is not stationary. A moving average graph of the time series reveals most drifts.
- A random walk (an autoregressive time series with a unit root) is not stationary.
 - The graph doesn't show whether the time series is a random walk.
 - The correlogram shows sample autocorrelations that do not decline rapidly.
 - Fitting an AR(1) process shows a β of about one (a unit root).

A time series can have a drift and also be a random walk. In both scenarios, we test for stationarity and take first differences; see the later steps of this guide.

Take heed: As an alternative to first differences, you may detrend the time series. If daily temperature increases 0.03% a year, detrend the time series.

COMMON ERRORS

As you construct ARIMA models, check if periods are needed. If you divide a time series into two periods and the ARIMA process is similar for both, you don't need two periods.

Illustration: Daily temperature may have cycles or long-term trends, but the ARIMA process may be similar to each period.

An unusual pattern in a correlogram may reflect a changing trend or drift. Suppose the correlogram shows sample autocorrelations

- declining from 25% to zero over the first 20 lags
- declining from zero to -15% from lags 20 to 30
- rising back to zero by lag 40.

This pattern indicates two periods with different drifts. Many time series have this pattern. Taking second differences to eliminate the pattern loses the information in the time series.

Second differences: First differences may remove stable trends in the time series. If the first differences are not stationary, we have three alternatives:

- If the trend is exponential, take logarithms and then first differences. *Do not take second differences instead of logarithms.* The resulting process is not stationary, but it might *seem* stationary because the trend is small.
- If the first differences have a stable trend, take second differences. This is uncommon. See if you can explain why this occurs.

Illustration: If we invest \$1,000 in a common stock portfolio, the value of the portfolio is a non-stationary random walk. Taking logarithms and first differences creates a stationary time series. If we invest \$1,000 in a common stock portfolio each month, we take first differences, subtract \$1,000, then logarithms, and then second differences.

- Graph the first differences. If the graph shows different means by period, we separate into two time periods. If one part of the time series has an upward trend and another part has a downward trend, taking first differences may not make it stationary.

Illustration: If the time series is $\{1, 2, \dots, 99, 100, 99, 98, \dots, 2, 1\}$, the first differences are $\{+1, +1, \dots, +1, +1, -1, -1, \dots, -1, -1\}$. The first differences have a mean of +1 in the first half of the time series and -1 in the second half. This is not a stationary process.

Be careful about taking second differences. If the time series comprises two eras with different means, variances, or drifts, taking higher order differences to make a stationary series obscures the true relation. If different models are appropriate for one part of a time series versus another, taking first and second differences obscures the problem.

Recommendation: Comparing two eras of a time series makes a good student project. Divide the time series into two parts and fit two ARIMA processes. An intuitive break is best, such as movie ticket sales before and after home DVD players. If the explanatory variables are not obvious, separate the eras by their means, drifts, or variances.

NEAS Time Series Project Documentation

Step #6: *ROBUST MODELS*

A robust model doesn't change much if we make small changes in the scenario. Examine if the periods chosen create robust models.

Illustration: If a 20 year period has a drift of +2% per annum, each 10 year sub-period should have a drift of about 2% per annum. If the first ten years have a drift of +5% and the second ten years have a drift of -1%, the overall drift of +2% is not meaningful.

Illustration: For a period of a few months and volatility of 0.1% a month, even a time series with no drift may show a small drift. Do not mistake volatility for drift.

The three interest rate periods on the discussion forum illustrate this concept.

- The observed interest rate drifts are +0.02%, -0.03%, and -0.01% per month for the three periods.
- The interest rate volatility is much higher in the second period.

A statistician would say that the first period has an upward drift, the second period is volatile, and the third period has a downward drift.

- The drift in the first and third periods is stable. If we divide the periods in half, we get the same drifts for each half.
- The second period is volatile and short. The absolute value of the drift is high, but it is *not robust*. Changing the period by a few months changes the drift.

To measure the drift, consider also the volatility of the rates and the length of the period.

The +0.02% and -0.01% drifts in the first and third periods reflect FED policy. The -0.03% drift in the middle period is an artifact of the short time period and the high volatility.

NEAS Time Series Project Documentation

Step #7: SCALING, INTERVAL LENGTH, STOCHASTICITY, AND MOVING AVERAGES

Choose an appropriate interval length. Some time series specify the interval length. Others allow you to choose the intervals. For stock prices, you may use daily, weekly, or monthly intervals. One might reason: Shorter intervals give more observations.

- ~ A monthly Treasury bill rate give 12 observations a year.
- ~ A daily corporate bond index gives 242 to 250 observations a year (business days).

More data points increase the accuracy of the analysis, so a Box-Pierce Q statistic that is not significant with 12 observations may be significant with 250 observations. But intervals that are too short hide the relations. Daily intervals may complicate the analysis.

Illustration: Take first differences of the interest rates and graph the results. The horizontal axis is the month and the vertical axis is the change in the interest rate. The graph looks like white noise. It is hard to see the upward or downward trends.

Monthly data in a stable series do not show the drifts well. The average monthly drifts are

- Period 1: +0.0215% \approx +0.02%
- Period 2: -0.0283% \approx -0.03%
- Period 3: -0.0099% \approx -0.01%

The monthly drifts are too small to see, unless we use narrow markers for the vertical axis. The *annual* drifts of 0.258%, -0.340%, and -0.119% are clear.

If we change just the scale of the vertical axis to use 0.01% as the marker, the interest rate stochasticity overwhelms the drift. We must change the scale *and* use a 12 month moving average to reduce the stochasticity. To observe the drift in your time series, examine

- a line graph of moving averages, which eliminates the stochasticity
- the 12 month first differences (the year to year change in the monthly rate)

Even monthly intervals are short. Monthly intervals give enough data to test hypotheses. But monthly intervals might make a stationary time series seem like a random walk.

Illustration: Suppose *annual* interest rates are a stationary AR(1) time series with $\phi_1 = 80\%$ and $\delta = 2\%$, so the mean interest rate is $2\% / (1 - 80\%) = 10\%$. *Monthly* interest rates might have an AR(1) process with $\phi_1 = 98\%$ and $\delta = 0.2\%$, so the mean interest rate is still 10%. This looks like a random walk, but it is not. Similarly, the daily corporate bond spread looks like a random walk, but it is mean reverting process using longer periods.

If the interest rate per annum is 12% in January 20X7, the forecast for January 20X8 is $80\% \times 12\% + 2\% = 11.60\%$. Using the monthly model, we get

NEAS Time Series Project Documentation

$$\text{February 20X7: } 98\% \times 12\% + .2\% = 11.96\%$$

Daily or weekly intervals of annual interest rates obscure the process. One time series on the web site is daily values of the Moody's 30 year corporate bond rate. A time series with daily intervals obscures the process. You choose monthly values as the average value in the month or as the first value in the month. The monthly time series is easier to work with. With the monthly values, use the techniques in this guide to fit a model.

Take heed: The first steps are preparation for your analysis. Your write-up should explain the periods and intervals. Show that you understand the change in the ARIMA process.

- If you choose a robust time series with proper periods and intervals, your ARIMA analysis proceeds smoothly.
- If you don't choose reasonable periods, you may waste much time in your analysis.

Take heed: Contrast overnight LIBOR and corporate bond spreads:

- Overnight LIBOR is a one day rate, and it changes rapidly; use daily periods.
- The corporate bond spread is a twenty year rate; use monthly periods.

Step #8: SEASONALITY

Examine your time series for seasonality, even if you do not expect seasonality. The write-up should explain what you examined, what you found, and the adjustments you made.

- Daily temperature and rainfall have smooth seasonality.
- Children's toys (high sales in December) or group health insurance, reinsurance, and workers' compensation (high sales in January) have strong, discrete seasonality.
- If stochasticity obscures the seasonality in the graph, use monthly or quarterly averages over several years. The hurricane season may not be clear in any one year, but a 20 year average by month shows the pattern. The textbook has several ways of identifying seasonality, such as dividing monthly rates by a 12 month centered moving average.
- Correlograms identify even weak seasonality. GDP, unemployment rates, and inflation have weak seasonality. Check the 12 month sample autocorrelations.
- The hypothesized relations must make sense. Don't follow numbers blindly. A high sample autocorrelation for a lag of 7 months is random fluctuation, not seasonality.
- Annual figures smooth seasonality. Daily and weekly rainfall is seasonal; semi-annual rainfall may not be seasonal.
- Many macroeconomic indices are adjusted for seasonality. CPI (inflation indices), price levels, unemployment rates, and Gross National Product are seasonally adjusted.

Take heed: If a time series has two peaks at opposite ends of the year, annual seasonality may appear as semi-annual seasonality.

Illustration: A quarterly series may have positive autocorrelations at lags 2, 4, and 8, and negative autocorrelations at lags 1 and 3. If the autocorrelation at lag 4 is greater than the autocorrelation at lag 2, this is annual seasonality. An autoregressive parameter of lag 4 may correct all the autocorrelations.

We correct for seasonality several ways, depending on the time series:

- seasonally adjust the data
- use seasonal differences
- use a seasonal lag in the ARIMA model

Illustration: Youth unemployment is highest in the summer, when school is not in session. Farm and construction work is high in the summer and low in the winter. We seasonally adjust unemployment rates to identify trends, cycles, and other effects.

Seasonal adjustments are covered in the chapter on non-stochastic time series. They are used whenever the value of a series depends on the time of the year, not the value of the series one year back.

Illustration: Suppose we model daily temperature with an ARIMA process. We seasonally adjust the data and then fit the ARIMA process. We *don't* use a 365 day lag, and we *don't*

NEAS Time Series Project Documentation

model the year-to-year changes in the daily temperature. See the project template on daily temperature for explanation. A student project may find the optimal method to seasonally adjust the data.

- For seasonal items, such as textbooks, camping equipment, heating oil, and wedding dresses, we examine growth by the year-to-year change in monthly sales.
- If a figure depends on the value 12 months back, ARIMA seasonal lags are best.

Illustration: We model personal auto written premium by month for a direct writer. The policy renewal rate is 90%, so the value 12 months ago is the proper base. We use ARIMA models with a ϕ_{12} of about 90%.

A student project on seasonality may have the following steps.

DAILY AVERAGES

Examine average values by day of the year.

- If interest rates have no trend or cycles, compute the average interest rate over the entire period for each day of the year.
- If interest rates have a trend or a cycle, simple averages conflate trend and seasonality.

Distinguish trend from seasonality:

- Compute 365 day centered moving averages. This eliminates seasonality and random fluctuations, leaving trend and cycles.
- To eliminate trend, convert interest rates to their deviations from a 365 day centered moving average.

Take heed: Overnight LIBOR shows business days only, or about 242 to 250 days a year. Use a centered moving average of the 365 calendar days, which cover a variable number of business days.

Take heed: Use the *COUNTIF* and *SUMIF* built-in functions to compute daily averages. The illustrative work-sheet for the project template on daily temperature uses these functions.

- Leap years cause an extra day every fourth year.
- Overlap of holidays with weekends may cause more business days some years.
- Missing values may cause fewer days.

Graph the moving averages. You see a decline followed by a rise in the overnight LIBOR for the period on the discussion forum Excel work-book. A student project might examine the relation of these movements to other macroeconomic indices. The centered moving average smooths the trends.

NEAS Time Series Project Documentation

Take heed: The daily temperature over the past 130 years may have weak trends or cycles (depending on the weather station). See the project template on daily temperature for methods of dealing with temperature trends.

Subtract the centered moving averages from the observed values. This eliminates trends and cycles, leaving seasonality and fluctuations. Each observed value is a deviation from the average in the surrounding year.

Compute long-term averages. This reduces the random fluctuation and leaves seasonality.

Take heed: Daily temperature has a high error term. The daily temperature may fluctuate $\pm 20^\circ$ because of unexpected weather. The daily temperature may be 30° one day and 65° two days later. Even a 130 year average daily temperature shows random fluctuations. Interest rates fluctuate less, and fluctuations are gradual. Graph the results as a line chart.

- If the line is smooth, the random fluctuations don't distort the long-term averages.
- If the line is jagged, replace each value by its centered moving average for 3 days or 5 days or 7 days or some other period. Use judgment to select the proper period. See the project template for daily temperature for an example.

INTEREST RATES AND SEASONALITY

For interest rates, seasonality is much weaker now than in the past and may be seen only in short rates. A student project using rates from the past few decades or rates with durations of one year or more need not discuss seasonality.

Interest rates are seasonal because they depend on the supply and demand for money.

The demand for money varies over the year. It is high for holiday shopping and low in January and February. If the money supply is held constant, interest rates are seasonal.

Eighty years ago, interest rates were seasonal. Now the Federal Reserve Board varies the supply of money to offset the demand for money.

The Federal Reserve Board is not perfect, and some seasonality remains. See if overnight LIBOR has any seasonality. You may examine whether a seasonal autoregressive term improves the model for overnight LIBOR.

Overnight, one week, two week, and one month LIBOR might be seasonal. A rate for one year or longer is not seasonal.

The graphs don't show much seasonality even for short LIBOR rates. But the correlogram may show a high 12 month sample autocorrelation. The pattern may be even clearer in the first differences.

NEAS Time Series Project Documentation

Recommendation: Decomposing LIBOR, insurance claim costs, and other actuarial items into long-term trends, cycles, seasonality, and stochasticity makes a good student project that can be valuable to your employer.

Take heed: If the time series is a random walk with weak mean reversion or seasonality, the sample autocorrelations show a slow decline for the first 11 months and a slight spike in the twelfth sample autocorrelation.

Stochasticity obscures weak seasonality. Annual seasonality may also cause high sample autocorrelations at 6 months, which further disrupts the pattern. Use the correlogram for the first differences to identify weak seasonality.

Illustration: The 1945-1978 period in the interest rates illustrative worksheet has a 17% correlation for the 12 month lag and lower correlations for the other lags. We have $34 \times 12 = 374$ observations in the first period. We subtract 1 for the first differences and 12 for the 12 month lag to get 361 observations. A sample autocorrelation higher than $2 \times 1/\sqrt{361} = 10.5\%$ is significant at a 5% level. A 17% sample autocorrelation is significant; most other sample autocorrelations are below 10.5% in absolute value.

Take heed: Don't expect all other sample autocorrelations to be below the critical value. By chance, one or two may be higher.

Many student projects use a correlogram, graph the sample autocorrelations, and conclude that seasonality is not important. Always examine seasonality; you may be surprised.

Illustration: Many candidates are not aware that claim frequency and severity are seasonal in many lines of business. Workers' compensation, group health insurance, reinsurance, are often written on January 1, so written premium is also seasonal.

If you find a significant 12 month sample autocorrelation, try an ARIMA(12,1,0) model: ϕ_1 and ϕ_{12} are non-zero, and the other coefficients are zero. Estimate the parameters with multiple linear regression.

If interest rates are a random walk with annual seasonality, we expect

- ϕ_{12} is low for a non-seasonal product and high for a seasonal product.
- ϕ_1 is low for a white noise process and high for a random walk.

NEAS Time Series Project Documentation

Step #9: *TRENDS*

A time series with a trend or drift is not stationary. We distinguish the two terms.

- A regression line has a trend.
- A random walk and other autoregressive processes have drifts.

Illustration: Suppose inflation has a drift of 5% per annum. If inflation is 3% in 20X8, we might expect it to be 4% in 20X9: an AR(1) process of the first differences with $\phi_1 = 50\%$.

Illustration: Suppose average claim severity has a trend of 5% per annum. If claim severity increases 3% in 20X8, we might expect it to increase 6% in 20X9. We assume that random fluctuations causes claim severity to be below trend in 20X8, so the increase in 20X9 is higher than trend.

Take heed: Trends in marriage rates, divorce rates, abortion rates, crime rates may change direction. A student project may examine the ARIMA process before and after the change in trend. Graph the data, separate into periods, and fit models to each period.

Stochasticity and seasonality obscure trends. A student project on climate change can be wonderful. Extensive data can be found on public web sites, and the implications are hotly debated. But high weather stochasticity overwhelms small trends. You may fit an ARIMA process to long-term weather indices to see if trends are real.

Illustration: To see trends in home sales, we use 12 month moving averages to remove the seasonality and dampen the stochasticity. Households buy homes in the summer months more than in winter months. A 2% annual trend in real (inflation-adjusted) home sales is obscured by the seasonality and random fluctuations.

Recommendation: Recent economic changes show the difficulty of identifying trends. Economists disagree if the U.S. economy is heading toward recession, if credit problems are a correction of lax lending practices, or if banks gave loans to weaker borrowers to meet federal non-discrimination requirements. A student project on home sales or mortgage rates might examine these issues.

LINEAR VS EXPONENTIAL TRENDS

To distinguish among linear, exponential, and other trends, graph the data.

- A linear trend appears as a straight line.
- An exponential trend appears as a convex (concave upward) curve.

Checking if a trend is linear or exponential is not easy. A non-linear trend is not necessarily exponential. It may be

NEAS Time Series Project Documentation

- A linear trend whose slope coefficient changes over time.
- A linear trend with much random fluctuation.
- A non-linear and non-exponential trend.

Decide if a trend is linear or exponential two ways:

1. Compare the trend of the time series to the trend in the logarithm of the time series.
2. Decide *intuitively* whether a linear or exponential trend makes more sense.

Illustration: If \$100 rises to \$110, \$200 should rise to \$220. The relation is multiplicative and the trend is exponential. But if Greenland's temporary rises from 1° Celsius to 2° Celsius in the 20th century, it might rise to 3° Celsius in the 21st century: a linear trend.

To adjust a time series for trend:

- For a linear trend, take first differences.
- For an exponential trend, take logarithms and then take first differences.

For stock prices, financial analysts take logarithms of ratios, which are the first differences of the logarithms. Either method is fine for the student project.

Summary

The initial steps in ARIMA modeling include: separate the time series in homogenous periods, de-seasonalize the data, and adjust for trends.

- If the time series differs in two time periods, separate the periods.
- If the time series is seasonal, de-seasonalize the data, take a seasonal difference, or use a seasonal lag in the ARIMA model.
- If the time series has a linear trend, take first differences.
- If the time series has an exponential trend, take logarithms and first differences.

NEAS Time Series Project Documentation

Step #10: STATIONARITY

Convert the time series to stationary form. Be careful not to take differences unless they are appropriate.

Illustration: The first and third interest rate periods in the interest rate project template have upward or downward drifts. They are not stationary, but their first differences may be stationary. Your student project tests for stationarity and fits an appropriate model.

The middle period is more complex. The interest rates in the middle period are volatile, and the drift may be random fluctuation. Even if the drift is zero and the volatility is constant, the time series could be white noise or a random walk. White noise is stationary and a random walk is not stationary. Your student project may test if the process is white noise, a random walk, or something else.

To *test* if a series is stationary, use sample autocorrelations, unit roots, and correlograms.

(1) Regress the time series on the same values one period back. This is an AR(1) model, which is the most common ARIMA process. If ϕ_1 (the β of the regression equation) is more than 1 or less than -1 , the time series is not stationary. We see this in the graph as well.

- If $\phi_1 > 1$, the time series grows continually. Random fluctuations may cause any single value to be smaller (in absolute value) than the preceding one, but the growth is clear over long periods. To correct this, take (logarithms and) first differences.
- If $\phi_1 < -1$, the time series grows continually and oscillates. Random fluctuations may obscure the exact process, but the oscillations are evident. This type of process is rare.

(2) If ϕ_1 is ≈ 1 , the time series is a random walk and is not stationary. Because of random fluctuations, the ordinary least squares estimator of the parameter is never exactly one.

- If we estimate ϕ_1 as 0.95 in a time series of 40 observations, we assume it is one and the time series is a non-stationary random walk.
- If we estimate ϕ_1 as 0.80 in a time series of 400 observations, we assume it is less than one and the time series is a stationary AR(1) process.

We rely on judgment, not on hard rules: t statistics and p -values for the null hypothesis that $\phi_1 = 1$ help us decide, but we don't have rigid rules.

Illustration: In the regression for the AR(1) model on the interest rate project template, the first and third periods have a ϕ_1 of 0.99, and the second period has a ϕ_1 of 0.85. [These are the values in the illustrative worksheet. You may choose a different time series and different periods. You will get different parameters and perhaps different conclusions.]

NEAS Time Series Project Documentation

The volatility is higher in the middle period than in the first or third periods. All three periods may be random walks, but the volatility is so high in the middle period and the length of the period is so short (48 months) that the slope coefficient has a high variance.

(3) The correlogram tests if the time series is stationary. If the autocorrelations do not *decline rapidly*, the time series is not stationary. *Rapid decline* means *at least geometric decline*. The time series is stochastic, and it may be hard to judge if the decline is rapid.

Checking if an interest rate time series is stationary is not easy. Annual interest rates are moving averages of 12 monthly forecasts. Since 11 of the 12 months are the same in adjoining periods, we get a $\phi_1 \approx 1$ in an AR(1) process.

Take heed: Be careful when you examine the stationarity of long duration interest rates. Overnight LIBOR fluctuates rapidly; Moody's 30 year corporate bond rate is steady.

If the time series or its first differences is stationary with a $\phi_1 < 1$, we have a *possible* AR(1) model. We consider three more items:

- Is the model correct? (Are the residuals close to a white noise process?)
- Is the model optimal? (Do other ARIMA processes fit equally well or better?)
- Does the model forecast well? (Do future values fall within a confidence interval?)

COMMON ERRORS: FIRST DIFFERENCES

A random walk is not stationary and an AR(1) process with a high ϕ_1 (but less than one) is stationary. Time series are stochastic, and it hard to distinguish the two scenarios.

- Not taking first differences of the random walk leaves a non-stationary series.
- Taking first differences of the AR(1) process is an error. We lose information about the time series, making it harder to forecast future values.

Take heed: Some candidates reason: the correlogram does not decline to zero quickly. The first differences form a more clearly stationary time series, which is easier to model.

Don't take first differences simply to make the time series easier to model. Take first differences only if the time series is not stationary. First differences lose information.

Illustration: We use monthly interest rates to get enough data points to test hypotheses. But monthly interest rates might make a stationary time series seem like a random walk.

Suppose *annual* interest rates are a stationary AR(1) time series with $\phi_1 = 80\%$ and $\delta = 2\%$, so the mean interest rate is 10%. *Monthly* interest rates might have an AR(1) process with $\phi_1 = 98\%$ and $\delta = 0.2\%$, so the mean interest rate is still 10%. This looks like a random walk, but it is not.

NEAS Time Series Project Documentation

If the interest rate per annum is 12% in January 20X7, the forecast for January 20X8 is $80\% \times 12\% + 2\% = 11.60\%$. Using the monthly model, we get

$$\text{February 20X7: } 98\% \times 12\% + .2\% = 11.96\%$$

Repeating this 12 times gives a forecast for January 20X8 of about 11.6%.

Take heed: If the correlogram shows geometric decline, the time series is stationary, even if the decline is slow.

COMMON ERRORS: MOVING AVERAGES

Use moving averages to find trends, not to test for stationarity. Moving averages are often used to remove seasonality. The moving averages obscure seasonality; they don't test for seasonality. Test for seasonality by examining monthly figures, and adjust for seasonality by one of the other methods in this course.

Don't use 12 month moving averages to form better correlograms. 11 of 12 months are the same in adjoining periods, so the sample autocorrelation function is high.

12 month interest rates are moving averages of 12 monthly forecasts, so be careful when you examine the stationarity of long rates. The correlogram examines the autocorrelation of long range forecasts. The forecasts change slowly; an interval of one month might show a random walk even if the true process is AR(1) with a high ϕ_1 parameter.

Take heed: The NEAS web site shows daily estimates of the investment grade corporate bond rate. A daily correlogram has such short intervals that patterns are hard to observe. You may use monthly intervals to evaluate the correlogram.

Take heed: The length of the time series (number of observations) does not determine the proper length of intervals.

- A time series of daily temperature over 100 years may have 36,524 observations.
- The daily temperature changes so quickly that intervals longer than 1 day do not show ARIMA processes.
- You may use hourly time series, with 24 hour seasonality overlaid on 365 day patterns.

Use quarterly or annual intervals to test if 90 day or 1 year Treasury bills are stationary. This is fine for the first period on the NEAS web site, which has 34 years. The middle period has only 4 years, and we can not use annual rates.

Take heed: If possible, detrend the time series, eliminate cycles and inflation, adjust for seasonality, and use methods besides taking differences to make a time series stationary. Taking differences is the proper adjustment only for random walks. The textbook does not make this clear.

NEAS Time Series Project Documentation

Step #11: *TESTING FOR WHITE NOISE*

The objective of ARIMA modeling is to forecast future values. Time series are stochastic, so forecasts do not exactly equal the future values. Ideally, the ARIMA process eliminates everything but the white noise of random fluctuations (the error term).

Check for white noise two places in your student project:

- ~ Once the time series is stationary, check if it is white noise.
- ~ After fitting an ARIMA process, check if the residuals are white noise.

Some statisticians consider white noise an ARIMA(0,0,0) process. If the first differences are white noise, the series is an ARIMA(0,1,0) process. Other statisticians say that white noise doesn't require an ARIMA model.

Use three tests for white noise: Durbin-Watson statistic, Bartlett's test, and Box-Pierce Q statistic. If the model is correct and the residuals are white noise:

- The regression of the series on the series lagged one period has no serial correlation, so the Durbin-Watson statistic is ≈ 2 .
- The sample autocorrelation of the residuals is normally distributed with a standard deviation of $1/\sqrt{T}$. Test this by examining percentiles. (Excel has a built-in function to test the percentiles. The function is not explained in the textbook, but you may use it.)
- The Box-Pierce Q statistic has a χ -squared distribution with the appropriate degrees of freedom.

If you have taken the regression analysis course, you can check the significance of the test using the Durbin-Watson table in the textbook. Keep in mind two items:

- ~ We are using a lagged value as the independent variable in the regression. The critical values for significance are not proper in this scenario. We use the Durbin-Watson statistic to help examine the time series, but we do not draw firm conclusions.
- ~ If the independent variable itself has a high autocorrelation, the Durbin-Watson statistic overstates the correlation of the residuals. The Durbin-Watson statistic may give wrong conclusions for time series modeling, so be careful with your analysis. The textbook mentions the problem, and recommends the Box-Pierce Q statistic instead.

For the student project, you may use the Durbin-Watson statistic. We want to see if you understand how to use the tool and what the results mean. We know that the test is not accurate for time series, and we do not require that you comment on this.

The Durbin-Watson statistic differs from Bartlett's test and the Box-Pierce Q statistic:

- The Durbin-Watson statistic uses autocorrelations of lag 1. Bartlett's test and the Box-Pierce Q statistic use autocorrelations of many lags. Bartlett's test and the Box-Pierce

NEAS Time Series Project Documentation

Q statistic are more robust, but if the sample autocorrelation of lag 1 is close to zero, the time series is probably a white noise process.

- The tests are scaled differently. White noise has a Durbin-Watson statistic of 2, sample autocorrelations that vary normally about zero (Bartlett's test), and a Box-Pierce Q statistic that is lower than the relevant χ -squared statistic.
- The progression of X values affects the autocorrelation of lag 1, so hypothesis testing is more complex for the Durbin-Watson statistic. The regression analysis module on the Durbin-Watson statistic explains how the correlation of the X values obscures the sample autocorrelation of the residuals.

If the Durbin-Watson statistic is 2, the process does not have an autoregressive coefficient of lag 1. It may have moving average coefficients or a seasonal autoregressive coefficient.

- Most ARIMA process have an autoregressive coefficient of lag 1. If the Durbin-Watson statistic is less than 1.600 to 1.700 (depending on the length of the time series), we examine AR(1) and AR(2) processes.
- If the Durbin-Watson statistic is between 1.800 and 2.200, the time series may be a white noise process. We examine Bartlett's test and the Box-Pierce Q statistic.

Bartlett's test and the Box-Pierce Q statistic examine more sample autocorrelations, such as the first 20 values. If they are close to zero, the time series is probably white noise. The *standard deviation* of the sample autocorrelations depends on the length of the time series. We have decades of interest rates, so we use them for the project templates. If you use ten years of annual premium volume for your student project, the data are too sparse for the statistical tests.

We check the percentage of sample autocorrelations with absolute values above various levels. We use judgment to evaluate the significance. This is a strong test. If you are familiar with Excel, you can use built-in functions for most of the work.

Review the discussion forum posting on time series techniques. We provide cell formulas and functions needed to complete the student project.

NEAS Time Series Project Documentation

Step #12: ARIMA MODELS

Once the time series is stationary but not white noise, specify an ARIMA process. The most common processes are AR(1), AR(2), MA(1), and ARMA(1,1). Each process has seasonal versions, versions with first differences, and versions with logarithms.

Illustration: An AR(1) process might be ARIMA(1,1,0), AR(12), where the ϕ_{12} parameter is for seasonality, ARIMA (12,1,0), or logarithmic versions of these.

Use simple processes for the student project. We review the student projects to see if you use statistical techniques properly. If you specify and test AR(1), MA(1), and ARIMA(1,1,0) models, and you explain what each model implies, you have completed the student project.

For each process, select parameters.

- For AR(1) and AR(2) processes, fit the model with linear regression.
- For MA(1) processes, use the Yule-Walker equations.

You may use other statistical software to fit the models. You may also use Excel *SOLVER* built-in function to fit a model.

Some time series are white noise after taking differences and logarithms, adjusting for seasonality, and regressing on economic or financial variables. If you begin with a time series that is not a simple random walk, and you take differences and logarithms, adjust for seasonality, regress on economic or financial variables, and convert your time series to a white noise process, your student project is fine.

Illustration: You use interest rates, daily temperatures, unemployment rates, inflation rates, sports won-loss records, sales volume, baby names, claim severity, or claim frequency, and you obtain a white noise process after the adjustments mentioned above:

- Check if a time series with fewer differences is stationary. Some candidates assume a time series with lower sample autocorrelations is better. They use second differences if that gives lower sample autocorrelations than the initial time series or first differences.
- If the time series with fewer differences is not stationary, write up the student project and turn it in. Do not think you erred because your result is a white noise process.

If you start with a white noise process or a random walk, you don't use ARIMA modeling.

- Do not use earthquake frequency (a white noise process) and test for white noise. We are not asking if you can form a white noise process.
- Do not use daily stock prices (a random walk), take logarithms and first differences, and say the result is white noise. This does not show that you can use ARIMA processes.

NEAS Time Series Project Documentation

You don't use the statistical techniques in the time series course for the two series above. But both series can be used if you analyze seasonality, cycles, drifts, and trends.

You can begin with daily stock prices and model weekly or annual seasonality. Economists refer to these as the Monday effect and the January effect. Statisticians have spent years modeling these two effects, and we don't know what causes them. Many financial papers analyze these patterns, and they are good topics for student projects.

Illustration: Monday effect

Take logarithms and first differences of daily stock prices. Index the data so that Mondays are $1 \bmod 5$, Tuesdays are $2 \bmod 5$, ..., and Fridays are $0 \bmod 5$. Use an AR(5) model with values for ϕ_1 and ϕ_5 , which you can fit easily with Excel.

Take heed: Holidays (New Year's, Fourth of July, Thanksgiving) don't have stock prices. To obtain entries for the time series, use the geometric average of the adjoining days.

Illustration: January effect

Take logarithms and first differences of monthly stock prices. Use an AR(12) model with values for ϕ_1 and ϕ_{12} . The monthly stock price may be the average in the month or the value on the 15th of the month.

Illustration: Natural catastrophes

Hurricanes have possible trends and cycles. You can fit an ARIMA model to a hundred year history of hurricane frequency.

Take heed: You don't need complex ARIMA models for the student project. Your student project should show that you understand how a moving average model differs from an autoregressive model and that you know how to test for each model. Use simpler models. If the simple models do not pass the Box-Pierce Q statistic, explain in your write-up what else a statistician might look at.

Illustration: Suppose your student project examines new home sales in Boston, and no simple ARIMA process gives white noise residuals. Your student project may say:

"New home sales are affected by economic conditions and mortgage rates. I regressed new home sales on GDP, but the indices I used were rough. Ideally, we should examine the residuals of new home sales in Boston regressed on per capita real personal income in Boston and on new home mortgage rates. In addition, I looked only at AR(1), AR(2), and MA(1) processes. A more complete analysis would look at ARIMA processes with more parameters."

Step #13: Correlograms

To choose among ARIMA processes, examine the correlograms. The illustrative spreadsheet on the NEAS web site gives the sample autocorrelation function.

- Autoregressive models have geometrically declining sample autocorrelations and spikes in the partial autocorrelation function.
- Moving average models have spikes in the sample autocorrelation function and declining partial autocorrelations.

Use the partial autocorrelation function if you have more sophisticated statistical software. To determine partial autocorrelations, we use nonlinear regression, which we do not cover in the statistics courses. Use the sample autocorrelations in the correlogram.

The autocorrelations from an AR(1) model have a geometric decline beginning with the first lag. The *sample* autocorrelations are the relations in a sample of observations. Ten years of monthly interest rates give 120 observations. The sample autocorrelations are stochastic and do not show a perfect geometric decline.

- If the sample autocorrelations have a reasonably rapid decline (but don't drop to zero immediately), AR(1) is usually the best model. You may test other processes, but the AR(1) solution is fine. Unless the number of observations is high and ϕ_1 is close to 1 (or -1), you can't confirm that the decay is geometric, since stochasticity overwhelms the expected results.
- If the sample autocorrelations are close to zero after the first lag, the indicated model depends on the first sample autocorrelation and the number of elements.
 - If the sample autocorrelation for the first lag is high or negative, the model may be MA(1).
 - If the sample autocorrelation for the first lag is positive but low, such as 15%, the model is probably AR(1). The indicated sample autocorrelation for the second lag of an AR(1) process is $15\%^2 = 2.25\%$, which is overwhelmed by stochasticity. Even if the sample autocorrelation for the first lag is 30% or 40%, the stochasticity is large enough in small or medium-size samples to obscure the sample autocorrelations.
- If the sample autocorrelation for the first lag is high or negative, and the remaining sample autocorrelations have a geometric decline, the model may be ARMA(1,1).
- If the sample autocorrelations for the first two lags are high and geometrically declining afterward, the model may have an AR(2) term and perhaps MA terms of order 1 or 2.

NEAS Time Series Project Documentation

Step #14: *STRUCTURAL MODELS*

Some candidates presume that structural models are more complex, so the student project takes more time. The opposite is true: the regression eliminates much of the variability in the original time series, and the ARIMA fitting is easier.

Illustration: Moody's investment grade corporate bond yield shows fluctuations, cycles, and trends. The corporate bond spread (after subtracting the long-term Treasury bond rate) is an ideal time series for ARIMA modeling.

Regress the corporate bond spread on the GDP growth rate. The residuals should be a stationary time series, which you might fit as white noise, AR(1), MA(1), or ARMA(1,1).

If you model the corporate bond yield itself

- The fit is less good and you may have to separate the time series into periods.
- You spend more time analyzing correlograms and choosing among alternative models.

By modeling the residuals of the corporate bond spread regressed on the GDP growth rate, you spend an extra hour getting the time series, but the rest of your project is quick.

NEAS Time Series Project Documentation

Step #15: *ORDER OF MODELS*

The textbook uses correlograms and in-sample tests to select a model. But ARIMA modeling is imprecise, since other factors may affect the time series values.

For the student project, use a sequence of models. First check trends and seasonality.

- De-trend the values. Use first differences and logarithms to determine the type of trend. You get a better ARIMA fit if you de-trend the time series with an inflation index and you de-seasonalize the data than if you take differences.
- Add seasonal lags or de-seasonalize the time series if needed.

Form a correlogram to see if an AR(1) or MA(1) model is appropriate.

Fit an AR(1) model with ordinary least squares for the autoregressive parameter and the seasonal parameter, if any. Compute the residuals from the AR(1) model and check the Durbin-Watson statistic, Barlett's test, and the Box-Pierce Q statistic.

If these tests are not significant, the residuals may be a white noise process and an AR(1) model is reasonable.

Do the same fitting for an AR(2) model, and state whether the better fit justifies the extra parameter. *If the AR(2) model is much better than the AR(1) model*, use it for the diagnostic testing; otherwise, we use the principle of parsimony and stick with AR(1).

Use the Yule-Walker equations to estimate θ_1 for an MA(1) process. Estimate the residuals from this model and apply the in-sample goodness-of-fit tests.

Take heed: Statistical software uses nonlinear regression to estimate an MA(1) process. For the student project, use the Yule-Walker equations. The estimated model is not the best possible, but it is close. We examine whether you correctly form the residuals from the MA(1) model for the Box-Pierce Q statistic and Bartlett's test.

Jacob: When do we use higher order autoregressive models?

Rachel: Some statisticians routinely use high order ARIMA processes; others do not. If you use the simple model correctly, and you explain what each model implies, you don't need more complex models.

Jacob: What ARMA(1,1), as well as its ARIMA and seasonal variants?

Examine the correlogram to see if ARMA(1,1) is indicated. If it is, explain how we see this from the correlogram. You do not have to fit the model.

NEAS Time Series Project Documentation

MA(1), AR(2), AND ARMA(1,1) MODELS

(The attached PDF file has better formatting.)

Updated: May 2, 2008

The textbook explains the statistical properties of ARIMA processes. This posting gives examples of MA(1), AR(2), and ARMA(1,1) processes in actuarial work, with *intuitive* explanations of when these models are used.

We discuss loss cost trends for stochastic lines of business and short underwriting cycles.

Loss Cost Trends

Autoregressive process: Suppose inflation is highly variable, and we have information about the expected trend. If we observe an 8% trend one year, we forecast an 8% trend the next year as well. This is an autoregressive process with $\phi_1 = 100\%$.

Jacob: Do we expect $\phi_1 = 100\%$?

Rachel: Some Western countries target an inflation rate, such as 6% per annum.

- If inflation one year is more than 6%, the central bank may tighten the money supply. We expect inflation the next year to fall toward 6%.
- If inflation one year is less than 6%, the central bank may loosen the money supply. We expect inflation the next year to rise toward 6%.

The time series may be a stationary autoregressive process with $\phi_1 < 100\%$.

Jacob: When might we see a moving average process?

Rachel: Suppose inflation is stable at 10% per annum and we model commercial fire loss cost trends for a small insurer. We expect the trend series to be +10%, +10%,

Your student project may take logarithms and first differences of average claim severities.

- The average claim severities may be \$10,000, $\$10,000 \times 1.10$, $\$10,000 \times 1.10^2$, ...
- Taking logarithms $\Rightarrow \ln(\$10,000)$, $\ln(\$10,000) + \ln(1.10)$, $\ln(\$10,000) + 2\ln(1.10)$, ...
- The first differences are $\ln(1.10)$, $\ln(1.10)$, ..., $\approx 10\%$, 10% , ...

If losses are not stochastic, the time series is a simple exponential trend. Random loss fluctuations add a moving average part to this time series.

If one year has many large claims (by chance), the observed trend that year may be +12%. The expected trend the next year is +8%, not +10%. This is an MA(1) model with a θ_1 parameter of 100%.

NEAS Time Series Project Documentation

Take heed: The trend is exponential, not linear. The 10% trend is a 10% linear trend after taking logarithms.

Jacob: Do we really expect $\theta_1 = 100\%$?

Rachel: Suppose the expected loss cost trend is +10%, and we observe +12% from 20X7 to 20X8. We consider several scenarios:

- The 20X7 average claim severity was the expected value, and the 20X8 average claim severity was 2% higher than expected. We expect the trend from 20X8 to 20X9 to be 2% less than 10%, or 8%.
- The 20X8 average claim severity was the expected value, and the 20X7 average claim severity was 2% lower than expected. We expect the trend from 20X8 to 20X9 to be 10%.

If both scenarios are equally likely, θ_1 may be 50%.

Jacob: These are AR(1) and MA(1) models. When would an ARMA(1,1) model be used?

Rachel: We combine the two stochastic pieces.

- For inflation, we use an AR(1) model with a positive parameter.
- For the stochasticity of fire losses, we add an MA(1) component.

The AR(1) model reflects the stochasticity of inflation. The MA(1) component adds the sampling error of the observed fire losses. The relative size of the components depends on the relative importance of these two forms of stochasticity.

Take heed: The project template on loss cost trends uses AR(1), MA(1), and ARMA(1,1) processes, along with seasonality and structural elements. The optimal trend model is rarely an exponential fit. Many items affect the loss cost trend. Actuaries have long sought more accurate methods of projecting trends. The ARIMA processes reviewed in the time series course may improve your trend models.

Underwriting Cycles

Underwriting cycles (and other business cycles) may be modeled by AR(2), ARMA(1,1), ARMA(2,1), and similar models.

Jacob: How do cycles differ from seasonality? They both show an oscillating pattern: higher in some months, quarters, or years, and lower in other months, quarters, or years.

Rachel: They differ in two ways. We compare

- Seasonal toy sales, which have higher fourth quarter than second quarter figures.

NEAS Time Series Project Documentation

- Cyclical auto insurance profits, which may have a four year oscillating pattern.

We assume a four year underwriting cycle to make the comparison clearer.

Seasonality says that fourth quarter sales are higher than average. If second quarter toy sales are low one year, we presume that the true average toy sales are low this year, so fourth quarter sales may also be lower than usual. A cycle says the opposite: if 20X5 insurance profits are low, the cycle may be more severe, and 20X7 profits may be higher than usual.

Insurance industry projections often have a cyclical perspective. After September 2001, many insurers expected higher premiums and strong profits in 2002. Some economists say such projections are not compatible with a competitive market. We do not seek the economic rationale for a cycle, but you may do a student project on insurance industry profits by line of business over the past 60 years to see if a cycle is reasonable.

Seasonality has fixed times; toy sales are high in December and perhaps in November. Cycles are not fixed. A four year cycle may extend into a five year cycle. If a cycle stops one year for exogenous reasons, we expect the cycle to renew with a 1 year lag in the future. We model cycles with ARIMA processes, not with curves of fixed periods.

Illustration: Suppose auto insurance has a four year cycle. If profits are low in 20X5, we expect average profits in 20X6 and high profits in 20X7. If profits are again low in 20X6, we expect average profits in 20X7 and high profits in 20X8.

Jacob: Do we use autoregressive or moving average models for seasonality?

Rachel: We use autoregressive, not moving average models. Autoregressive models deal with expected changes, such as higher than average December sales. Moving average models deal with *unexpected* changes. In the cycle illustration above, we expected average profits in 20X6 but actual profits were low.

Jacob: How do we model a cycle? Suppose an industry has profit cycles, where annual profits relative to long-run mean are 0, +Z, 0, -Z. This looks the same as seasonality; do we use $\text{profit}(t) = \text{profit}(t-4) + \epsilon$?

Rachel: This has two problems:

- If profit at bottom of cycle was unusually poor, such as $-2Z$, we expect profit at upturn to be unusually good, such as $+2Z$. So we use $-\text{profit}(t-2)$. This is an autoregressive model.
- If the cycle is delayed one period by an exogenous force, we expect the cycle to continue afterwards with a one period delay. If we have an extra year of zero profits in year 6 of this time series, we expect: 0, +Z, 0, -Z, 0, 0, +Z, 0, -Z.

We might use an ARMA(2,1) process with $\phi_2 = -1$ and $\theta_1 = -1$.

NEAS Time Series Project Documentation

- If the cycle is on schedule, the AR process shows the expected profits.
- If the cycle has a delay, such as the two periods with zero profits, the residual gets the cycle back on track.

Jacob: This illustration seems simplistic. Cycles don't have regular four year periodicity. A cycle may be five, six, or seven years. How do we deal with this uncertainty?

Rachel: We may have two assumptions:

- This year's profits will be like last year's profits, were there no cycle.
- If last year's profits were higher (lower) than those of the year before, we assume the cycle is moving up (down), and this year's profits may be even higher (lower).

We can model this as an AR(2) process. The first assumption is an AR(1) model. If the average profits are 12%, the model may be $y_t = 6\% + \frac{1}{2} y_{t-1} + \epsilon$. The second assumption is an AR(2) process, such as $y_t = \frac{1}{4} (y_{t-1} - y_{t-2}) + \epsilon$. Putting the two together gives

$$y_t = 6\% + \frac{3}{4} y_{t-1} - \frac{1}{4} y_{t-2} + \epsilon.$$

Jacob: This is still too simple; it doesn't catch turning points of the cycle.

Rachel: More sophisticated ARIMA models may be used for more complex cycles. But the more sophisticated model is not necessarily better. Cycles vary so much that we don't have good models.

Jacob: Can we do a student project analyzing property-casualty underwriting cycles?

Rachel: The data are so haphazard that no ARIMA model fits well and you won't see how to apply the statistical techniques. The project templates for interest rates show how to apply the statistical techniques and arrive at proper inferences.

Take heed: Researchers do not agree on the existence, shape, or periods of underwriting cycles. You can collect insurance industry profit margins by line of business and fit ARIMA processes. You can use simple processes; less you have time series software, you can't fit complex ARIMA processes. Use multiple linear regression to estimate and AR(1) and AR(2) process, and the Yule-Walker equations for MA(1) and ARMA(1,1) processes.

NEAS Time Series Project Documentation

TIME SERIES STUDENT PROJECTS: HYPOTHESIS TESTING AND DISTRIBUTIONS

(The attached PDF file has better formatting.)

Updated: May 2, 2008

Jacob: What do the Durbin-Watson statistic, Bartlett's test, and the Box-Pierce Q statistic evaluate?

Rachel: These tests evaluate if the sample autocorrelations of the residuals are statistically different from zero.

Jacob: How do these tests work? I understand the formulas, but I don't get the intuition.

Rachel: To test the null hypothesis that the autocorrelations are zero, we need several assumptions about the distribution of sample autocorrelations. We explain the terms.

- The autocorrelations of the residuals are not observed. If the ARIMA process is a good model for the time series, the autocorrelations of the residuals have a mean of zero.
- The sample autocorrelation of the residuals are observed values that differ from zero. We estimate the standard error of the sample autocorrelations.

Jacob: Regression analysis uses t statistics and p -values to test hypotheses that a value is zero. Do we make assumptions about distributions?

Rachel: Hypothesis testing in regression analysis generally assumes that the residuals are normally distributed with a constant variance.

Jacob: How do the assumptions affect the hypothesis test?

Rachel: Suppose the *sample autocorrelation* of lag 1 is 8% and we want to test the hypothesis that the autocorrelation is actually zero. The null hypothesis assumes that the sample autocorrelation has a normal distribution with a mean of zero, but we don't know the variance or the standard deviation (the square root of the variance) of the distribution.

- If the standard deviation is 8%, the probability of the sample autocorrelation being greater in absolute value than 8% is about one third.
- If the standard deviation is 4%, the probability of the sample autocorrelation being greater in absolute value than 8% is about 5%.

Jacob: How do we estimate the variance of this distribution?

Rachel: Bartlett's theorem says that the sample autocorrelations of a white noise process have a normal distribution with a variance of $1/T$, where T is the number of observations. If the residuals are a white noise process, their sample autocorrelations are normally distributed with a standard deviation of $1/\sqrt{T}$.

NEAS Time Series Project Documentation

Jacob: How does Bartlett's test work?

Rachel: If the sample autocorrelations have a normal distribution with a mean of zero and a standard deviation of σ , then the probability of an observed sample autocorrelation being greater in absolute value

- ~ than 1.96σ is 5%.
- ~ than 1.45σ is 10%.

Jacob: Suppose we find that a sample autocorrelation exceeds 1.96σ . Do we infer that the sample autocorrelations are not a white noise process?

Rachel: That depends on the lag, the type of data, and the other sample autocorrelations. Suppose $T = 400$, $1/\sqrt{T} = 5\%$, $\sigma = 5\%$, $1.96\sigma \approx 10\%$, the data are monthly interest rates, and the sample autocorrelation of lag k is 12.5%.

If $k = 1$, we avoid drawing any inference. If the interest rates are correlated with the time period, the residuals may appear to be serially correlated, even if they are not. It is common for sample autocorrelations of lag 1 to be more than zero even if the autocorrelation is zero. The Durbin-Watson statistic may be inconclusive. This is like Bayesian estimation. The *prior* distribution for the autocorrelation of lag 1 has a high probability of being greater than zero because of other reasons. An observed value of 2.5 standard deviations is inconclusive.

If we observe 10 lags ($k = 10$) and the sample autocorrelations of 9 of these lags are less than 10%, we presume the high sample autocorrelation the tenth lag is the expected random fluctuation in a normal distribution. The null hypothesis of a zero mean assumes that 5% of the sample autocorrelations have absolute values greater than 10%. One high sample autocorrelation out of ten is not unexpected.

If $k = 12$ and the sample autocorrelations of lags 1 through 11 are less than 10%, we suspect the high sample autocorrelation for lag 12 reflects annual seasonality.

Jacob: Bartlett's test sounds subjective; is this a problem?

Rachel: A skilled statistician prefers a subjective test that relies on our intuition about the time series. Seasonal correlations are more common than non-seasonal correlations. We may have subjective views on the types of autocorrelations that are most likely, such as positive vs negative autocorrelations.

Jacob: If we assume a normal distribution with a standard deviation of σ , the probability of a positive sample autocorrelation is the same as the probability of a negative sample autocorrelation. Why should the sign of the sample autocorrelation affect our inference?

Rachel: If we model monthly sales data and find a sample autocorrelation of 15% for a lag of 12 months, we presume this is annual seasonality. If the sample autocorrelation for a lag of 12 months is -15% , we may attribute this to random fluctuation.

NEAS Time Series Project Documentation

- ~ If the actual autocorrelation is zero, the probability of positive vs negative sample autocorrelations is the same.
- ~ In practice, the incidence of positive autocorrelation is not the same as the incidence of negative autocorrelation. An AR(1) model with a negative parameter for one lag is often a spurious result.

Jacob: Do we examine the Durbin-Watson statistic, Bartlett's test, and the Box-Pierce Q statistic on the residuals or the sample autocorrelations of the residuals?

Rachel: The Durbin-Watson statistic is applied to the residuals; the Durbin-Watson statistic calculates the autocorrelation of lag 1. The Durbin-Watson statistic $\approx 2 - 2 \times$ the sample autocorrelation of lag 1.

- ~ For perfect positive autocorrelation, the Durbin-Watson statistic = 0.
- ~ For perfect negative autocorrelation, the Durbin-Watson statistic = 4.

Bartlett's test and the Box-Pierce Q statistic are applied to the sample autocorrelations of the residuals.

Jacob: How does the Durbin-Watson statistic differ from Bartlett's test and the Box-Pierce Q statistic?

Rachel: Bartlett's test and the Box-Pierce Q statistic evaluate whether the residuals have a normal distribution with a variance of $1/T$, where T is the number of observations. They examine sample autocorrelations of various lags. The Durbin-Watson statistic examines if the sample autocorrelations of lag 1 are statistically different from zero.

Jacob: Are these tests equally strict?

Rachel: The Durbin-Watson statistic considers other factors that affect the observed sample autocorrelation of lag 1. The Durbin-Watson statistic may be distorted in a lagged regression, so we don't use this statistic for hypothesis testing of ARIMA models.

Take heed: You may use the Durbin-Watson statistic in your student project, but be aware that it might a serial correlation even if none exists.

Bartlett's test and the Box-Pierce Q statistic are applied to many sample autocorrelations, not just the first one. Both tests acknowledge that sample autocorrelations for the first several lags many not follow the normal distribution or the χ -squared distribution.

NEAS Time Series Project Documentation

NEAS Time Series Project Documentation

TIME SERIES PROJECT OSCILLATING, MEAN REVERTING, AND CYCLICAL

(The attached PDF file has better formatting.)

Updated: April 7, 2006

Jacob: What is the difference between oscillating, mean reverting, and cyclical?

Rachel: A non-oscillating mean reverting model moves part way toward the mean. An oscillating process moves toward the mean and over-shoots it.

- An AR(1) model is mean reverting if $-1 < \phi_1 < 1$.
- An AR(1) model is oscillating if $\phi_1 < 0$.

An AR(1) model is stationary if and only if it is mean reverting.

Jacob: Can an AR(1) model be mean reverting, stationary, and oscillatory?

Rachel: Yes. A mean reverting process is stationary. It is also oscillatory if $-1 < \phi_1 < 0$.

Jacob: If a time series has a mean of μ , is oscillatory, and $y_t > \mu$, is $y_{t+1} < \mu$?

Rachel: You should say: "If a time series has a mean of μ , is oscillatory, and $y_t > \mu$, then the *expected value* of $y_{t+1} < \mu$." The time series is stochastic. We forecast the expected values, not the actual values.

Jacob: If a time series has a mean of μ , is mean reverting but not oscillatory, and $y_t > \mu$, is the expected value of y_{t+1} between μ and y_t ? Is the probability that $\mu < y_{t+1} < y_t$ more than 50%?

Rachel: The first sentence is correct: If a time series has a mean of μ , is mean reverting but not oscillatory, and $y_t > \mu$, then the expected value of y_{t+1} is between μ and y_t . The second sentence should say: "The probability that $\mu < y_{t+1}$ is more than 50% and the probability that $y_{t+1} < y_t$ is more than 50%."

Jacob: Can we make more exact probabilistic statements? If a time series has a mean of μ , is mean reverting but not oscillatory, and $y_t > \mu$, then the expected value of y_{t+1} is between μ and y_t . Can we make probabilistic statements like: "The probability $y_{t+1} > y_t$ is less than 5%"?

Rachel: We must know the value of y_t and the standard deviation of the time series. If y_t is 1% and σ is 25%, the probability $y_{t+1} > y_t$ is close to 50%.

Jacob: Is cyclical another term for oscillating?

NEAS Time Series Project Documentation

Rachel: Oscillating means the terms alternate about the mean. The autoregressive parameter ϕ_1 is negative. A cyclical pattern, also called a sinusoidal pattern, looks like a sine wave. The autoregressive parameter ϕ_1 is positive, but one or more higher order autoregressive parameters are negative.

Jacob: What is the intuition for an oscillatory process? Do we encounter ARIMA processes with negative sample autocorrelations of lag 1?

Rachel: We may have moving average processes with a negative sample autocorrelations of lag 1. We explain these in other discussion forum postings. They are may occur for loss cost trends in small samples, cycles, and industry wide sales of durable goods.

Take heed: A common cause of negative sample autocorrelations of lag 1 is in the first differences of stationary processes. The common example is

- A time series is a white noise process (or a weak autoregressive process) overlaid on a trend or a long-term cycle.
- The candidate takes first differences to eliminate the trend. This makes the time series stationary, but the first differences have a negative sample autocorrelation of lag 1.
- The ideal solution is to detrend the time series or otherwise offset the long-term cycle. If possible, use a structural model to eliminate the trend.
- If the first differences have negative sample autocorrelations first two or three lags that become zero by the fourth or fifth lag, return to the original time series and detrend the series. You have an easily explained series, that the first differences hide.

NEAS Time Series Project Documentation

PROJECT TEMPLATE ON INTEREST RATES AND OTHER ECONOMIC TIME SERIES

This project template illustrates ARIMA modeling for interest rates, inflation, unemployment rates, and other macroeconomic indices. It uses three month Treasury bills from the NEAS web site as a sample time series, and it explains the applications to other interest rates and economic indices.

This project template explains how to fit ARIMA processes to financial and economic time series. It provides guidance for many student projects, not just those on interest rates.

Illustration: You may use overnight LIBOR rates or corporate bond spreads for your student project. Adapt this project template to your time series.

The illustrative worksheets on sample autocorrelation functions, correlograms, Durbin-Watson statistic, and Box-Pierce Q statistic use the interest rate observations in this project template. Some of the discussion forum postings on ARIMA modeling use illustrations from this project template. Even if you choose another topic for your student project, this project template clarifies many of the statistical techniques you use.

Take heed: The project templates are difficult to compose. Candidates begin their student projects at various levels of expertise.

- Some understand the course material and are looking for good topics.
- Others feel lost and need guidance through all steps of the student project.

The student projects must be independent work. We can't say: "Do Step A, then Step B, and so forth," or the SOA would not grant credit for the on-line course. The on-line courses teach the statistics material so that you can apply them to real data.

The project templates on the discussion forum lead you through statistical analysis.

- If you know the statistics material, you may find some of the information to be repetitive. We repeat some items in different postings, since it takes time to absorb the concepts.
- If you are lost, you may find that the instructions seem incomplete. If you have trouble grasping the concepts, copy the data sets to an Excel workbook and recreate the analyses described in the discussion forum postings.

If you are worried about successfully completing the student project, take heart. Many candidates feel overwhelmed at first, and they think the student projects will be a drain on their time. The first year, the student projects were indeed difficult for some candidates. But our faculty has since created an excellent array of project templates and other postings that provide guidance for all levels. You need an hour or two to get oriented, and you will soon be working through statistical analysis of real data.

NEAS Time Series Project Documentation

LEVELS OF GUIDANCE

The project templates describe various methods. The ideal method depends on your data, your hypotheses, and your assumptions.

- The project template tells you to graph the interest rates and examine the graph and the correlogram for stationarity, seasonality, trends, and other patterns. Interest rates come in thousands of varieties, and we do not know the patterns in your time series.
- Another part of the project template suggests taking first differences to make the time series stationary. If the initial time series is a random walk, always take differences. If the initial time series is not stationary but is not a random walk, the proper method is not clear. Examine the correlogram of the first differences, but use one of the other methods if possible.
- A third part implies that taking first differences is the wrong approach. Instead, convert the nominal interest rates to real interest rates. Converting a nominal time series to real terms often makes it easier to fit an ARIMA process.
- A fourth part implies that you should divide the time series into three segments and fit different ARIMA process to each one.
- A fifth part tells you to regress the interest rates on other macroeconomic indices (GDP, inflation, money supply) and fit the ARIMA process to the residuals.

Question: How can we complete a student project with these instructions? *Answer:*

1. The true cause of interest rate movements is not well known. The textbook proposes several models, and financial economists propose many more. We are not hiding the answer from you; we don't know the pattern of interest rates any better than others do.
2. The student project teaches you the statistical techniques. Many methods may be used to make a time series stationary. You may use two or three methods and see which one produces the best model. We are not telling you what to do. We are explaining the statistical techniques that you may apply to the data.
3. The ARIMA process is not the true explanation of interest rates. It is a proxy that often does well for short periods. The proper analysis varies with the objective.
 - An exact financial model may regress real interest rates on GDP
 - A simple short-term proxy may use first differences of nominal interest rates.
4. If the statistician knows the expected inflation rates in coming months, projections of real interest rates are fine. If the statistician does not know expected inflation, the nominal interest rates must be projected.

Jacob: The textbook fits an ARIMA(8,1,4) process to interest rates. Do we fit complex ARIMA processes in the student project?

NEAS Time Series Project Documentation

Rachel: For the student projects, we focus on simple processes. If you understand an ARMA(1,1) process, you understand the more complex processes as well. Focus on the basic ARIMA processes and statistical tests.

Jacob: For the final exam, we had to learn formulas. We don't have statistical software to compute sample autocorrelations, correlograms, the Box-Pierce Q statistic, and other items used for ARIMA fitting. Do we have to code each of these items?

Rachel: This is a course in statistics, not programming. We provide illustrative Excel worksheets that do the number crunching. You copy and paste cell formulas, and you use the Excel built-in functions, but do not spend time on arithmetic.

Illustration: We provide interest rates, cell formulas, and VBA macros for the statistical tests. You form correlograms, interpret the graphs, and fit a good model to the time series.

Take heed: Excel 2007 differs from previous versions in your access to add-ins. To get the *REGRESSION* add-in or the *SOLVER* add-in for previous versions of Excel, click the *TOOLS* menu and choose the add-in. In Excel 2007, click the *DATA* menu and choose the add-in from the *ANALYSIS* portion of the ribbon. The add-ins work the same way. (The add-ins are produced by other firms, not by Microsoft. They did not change with Windows Vista.) Some of the discussion forum postings have the pre-2007 Excel procedure. Nothing substantive have changed in the statistical add-ins, but the series to clicks to get the add-in differs.

NEAS Time Series Project Documentation

Step #1: Choose the interest rate time series

We illustrate with three month Treasury bills, for several reasons:

- The data are readily available. We give you an Excel spread-sheet with the time series in a column; you have no need to gather data.
- The data extend over many years. The ARIMA model differs by sub-period.
- The interpretation of interest rate time series is unclear. You have many ways to continue the analysis here. No two candidates will come to the same conclusion.

Take heed: We encourage you to search the internet for other time series. Data are easier to gather than many candidates think.

We provide several other time series of interest rates: LIBOR rates at various maturities for the U.S. dollar and some other currencies; other Treasury security rates; corporate bond yields; municipal bond yields; bank prime rates; and others.

Hundreds of other interest rate time series are available on the world wide web. Use the internet search engines (Google, MSN, Yahoo) for “interest rate” or “interest rate history.”

Take heed: If you find a web site with good time series data, post a message on the discussion forum. Other candidates appreciate the help.

NEAS Time Series Project Documentation

Step #2: Correct Missing or Erroneous Elements

The time series starts in January 1931 and continues through December 1935. Three month Treasury bills were auctioned (sold at auction) once a month in these five years. The next auction was in February 1941, leaving 61 months with no auction. We have several ways of dealing with missing periods.

- If the time series is 40 years of data, one year of missing data, then 40 years of data, and if the two periods have similar processes, we might ignore the 12 missing months.
- If one observation is missing, interpolate between the two surrounding values.
- Economic conditions differ by decade. Treasury bill rates decline from 3.25% at the end of 1931 to 0.09% at the end of 1935. During World War II, rates were affected by a war-time budget deficit. Financial economists often analyze post-World War II indices. Select an appropriate period, based on the data available and the time series process.

Take heed: The project template on daily temperature uses the first two methods above.

- Weather data are often missing on some days: interpolate for the missing values.
- The daily temperature series may have several months missing: ignore these values.

Take heed: Twenty minutes of checking for errors may save you hours of wasted work.

- Most time series on the discussion forum do not have errors. The data have already been checked and corrected. But the time series may have missing periods.
- If you gather data from other web sites, check the data. A missing value may be coded as -999. If you have 100 values of interest rates with an average of 5% and you include a value of -999, this one value skews the time series.
- If you use in-house data for loss cost trends, a coding error may skew your data. If a value seems erroneous, you save time by eliminating it or correcting it.

No method is right or wrong. The question is whether Treasury bill rates followed the same pattern in the early 1930's and early 1940's as in later years.

If we use post-World War II data, we have 666 months of Treasury bill rates. The first month is January 1945, and the last month is June 2000.

Recommendation: The choice of time series is arbitrary. Start with a long time series, such as 125 years of daily temperature or 2,500 observations of overnight LIBOR rates. As you complete the project, you may find that the time series process changes mid-way through. Sometimes the reason is clear; other times you see the change but can't explain it. This project template explains methods to make the time series stationary.

NEAS Time Series Project Documentation

Step #3: Graph the Time Series

Graphing the time series helps you see patterns. You form many graphs and charts for your student project. General advice on graphing:

- Use line graphs for long time series. Set your chart default to line graph (instead of bar graph). Bar graphs are used for small samples, such as a regression analysis of voter types. For time series, use line charts.
- Use deviations from the mean to see patterns. If you are not certain of a trend, graph the deviations from the mean.
- Use graphs of centered moving averages to smooth random fluctuations. If you are not certain of a trend, graph the centered moving averages.
- Label your axes, and use titles and legends, so the course instructor can read the graph. If you use an index for the observations, state the index values on the graph. If your time series uses a month index from 1 to 666, write “Jan 1945 = Month 1.”
- Callouts, arrows, autoshapes, and text boxes help readers understand your graphs.
- You may graph the time series, the first differences, or the sample autocorrelations.

The type of graph depends on the item you analyze.

- To identify trends in a seasonal process, you may graph 12 month moving averages.
- To identify seasonality, you may graph monthly averages over many years.

Take heed: Explain the pattern in your time series and the implications for your analysis.

Illustration: The graph of 3 month Treasury bill rates for January 1945 – June 2000 has three patterns:

- increasing through 1979, with either cycles or random movements
- volatile for about four years, with no clear pattern
- decreasing through the end of the time series, with cycles or random movements

Take heed: The three periods above are not precise. For your student project, give your impression of the time series pattern. The graph uses call-outs to identify possible changes in the time series. The worksheet uses slightly different periods for correlograms. Examine the data and decide what periods seem homogeneous.

This pattern means the time series is not stationary. Corroborate your analysis with the correlogram. A correlogram starting from February 1941 or January 1945 shows a long series of about 200 positive sample autocorrelations.

You have several ways of proceeding. We list them from easiest to preferred methods.

Take heed: Statisticians differ on the optimal method.

NEAS Time Series Project Documentation

- The course textbook uses an ARIMA(8,1,4) process for interest rates.
- Other statisticians would consider this process an error: it has no intuitive rationale and the slight improvement in the in-sample fit does not offset the added complexity.

Any statistical method below is fine for the student project.

NEAS Time Series Project Documentation

Step #4: *DIFFERENCES*

Take first differences to make the time series stationary. *Always take first differences if the time series is a random walk.* If the process is multiplicative (instead of additive), first take logarithms and then take first differences.

Many project templates discuss random walks, white noise, and mean reversion. A random walk can be hard to distinguish from a stationary AR(1) process with a high ϕ_1 .

The correlogram for the first differences of 3 month Treasury bills is hard to interpret. Some statisticians would say that the long string of positive sample autocorrelations in the original time series disappears from the first differences. The sample autocorrelations for the first ten lags are sometimes positive and sometimes negative. These statisticians say we can fit an ARIMA process to the first differences, not the original time series. The ARIMA(8,1,4) fit in the textbook takes this view.

But interest rates in developed countries are usually not a random walk. An economist might see three distinct patterns in the time series of 3 month Treasury bills.

If interest rates are not a random walk, first differences are not ideal. We mention methods to form a stationary time series that can be modeled by a lower order ARIMA process.

Take heed: The time series for three month Treasury bills has interest rates to two decimal places, such as 5.50%. If the interest rate does not change between months, the first difference is zero. The time series for LIBOR rates has more decimal places. The first differences of LIBOR rates look more like a normal distribution.

Take heed: You can check whether taking first differences is appropriate.

- If the sample autocorrelations of the first differences are a white noise process, taking first differences is correct.
- If the sample autocorrelations of the first differences are negative for first two or three lags, taking first differences is not correct.

The discussion forum posting on *time series simulations* explains these relations.

The sample autocorrelations of first differences of three month Treasury bills do not give a conclusive answer. We show the sample autocorrelations and the correlogram. With 665 observations, the standard deviation of a white noise process is $1/\sqrt{665} \approx 3.88\%$. The sample autocorrelations for the first 20 lags have many values greater than two standard deviations, but no clear pattern. Fitting an ARIMA model to the first differences is not easy.

Take heed: Each time series differs. Your student project may use other rates, such as LIBOR rates or bank prime rates. You need not find the correct ARIMA process or even a good fit. But you should explain your analysis: why you chose a particular approach.

NEAS Time Series Project Documentation

UNIT ROOTS

The textbook discussion of unit roots is concise. The homework assignments and final exam problems do not emphasize this topic, but it is important for the student project.

Interest rates may be mean reverting or random walks.

- A mean reverting time series is stationary and can be modeled by an ARMA process.
- A random walk is not stationary. It is modeled by an ARIMA process.

If your time series looks like a random walk, use three tests:

- Use a one period lagged regression and check for a unit root.
- Form a correlogram and examine the decay in the sample autocorrelation function.
- Take first differences and check the Box-Pierce Q statistic.

Take heed: Many economic and financial indices may be random walks. Interest rates, inflation rates, GDP growth, unemployment rates, and various similar indices look like random walks within moderate bounds.

Unit Root: Regress the time series on the same values one period back. This is an AR(1) model, which is the most common ARIMA process. If ϕ_1 (the β of the regression equation) is more than 1 or less than -1 , the time series is not stationary. We see this in the graph.

- If $\phi_1 > 1$, the time series grows continually. Random fluctuations may cause any single value to be smaller (in absolute value) than the preceding one, but the growth is clear over long periods. To correct this, we take logarithms and first differences.
- If $\phi_1 < -1$, the time series grows continually and oscillates. Random fluctuations may obscure the exact process, but the oscillations are evident. This type of process is rare.

If ϕ_1 is ≈ 1 , the time series is a random walk and is not stationary. Because of random fluctuations, the ordinary least squares estimator of the parameter is never exactly one.

- If we estimate ϕ_1 as 0.95 in a time series of 40 observations, we assume it is one and the time series is a non-stationary random walk.
- If we estimate ϕ_1 as 0.80 in a time series of 400 observations, we assume it is less than one and the time series is a stationary AR(1) process.

NEAS Time Series Project Documentation

Step #5: *PERIODS*

If an exogenous *intervention* causes the different patterns, separate the time series into two or more periods.

- Statutory, regulatory, and judicial interventions are common exogenous factors. The time series project templates have many examples.
- Federal Reserve Board policy on interest rates affected the Treasury bill time series.

If you are not aware of exogenous factors, base the periods on the time series pattern.

Illustration: Suppose interest rates increase for the first ten years and decrease for the second ten years.

- The time series itself is not mean reverting, so it is not stationary.
- The first differences are positive for the first ten year and negative for the last ten years. The mean first difference is not stable, so the first differences are not stationary.
- In this simple process, the second differences may be zero for all periods (except at ten years). Actual time series are less distinct. The interest rate path may be a parabola, and even the second differences are not stationary.

If the time series can be divided into two or three homogeneous processes, fit a separate ARIMA process to each period. Periods often provide interesting student projects. Distinct changes often occur, such as a new government in a country or a new CEO of a firm. You can analyze GDP growth under two governments or sales growth under two CEO's. Some past student projects posted on the discussion board analyze time series in two periods.

Illustration: Examine the graph of three month Treasury bill rates for January 1945 through June 2000. Note the general upward trend for 1945 through 1979, the high volatility for 1980 through 1983, and the downward trend for 1984 through 2000.

For your student project, do the following:

- Graph your time series and examine the means, trend, and volatility in different periods.
- Select periods based on either exogenous information (a change in policy, legislation, or economic environment) or the observed means, trends, and volatility.
- Separate into periods if the differences seem material, not just random fluctuation.
- Fit ARIMA processes to each period. Examine if the processes in adjoining periods are materially different.

Take heed: Structural models and de-trending are better than separating into periods if the time series depends on some other index. You may not be able to fit an ARIMA process to nominal interest rates, but perhaps you can model real interest rates. Unemployment rates may have a cyclical pattern, but the residuals of unemployment rates on GDP growth may be an AR(1) process.

NEAS Time Series Project Documentation

Take heed: Some statisticians avoid dividing time series into periods. They say the goal of time series analysis is to forecast turning points. Periods of stability are easy to forecast. Starting a new period at each turning point does not accomplish the goal. The authors of the course textbook generally prefer to model long time series with higher order terms.

Other statisticians believe that modeling a time series composed of different processes degrades the forecasts. The interest rate process depends on monetary policy. If Federal Reserve Board policy changes, the interest rate process changes.

Periods are useful if the ARIMA process changes, not if it simply turns. Many industries show profit cycles, like the underwriting cycles in insurance.

- Corporate bond spreads may also have cycles, perhaps related to business cycles.
- If such cycles exist, the goal of ARIMA modeling is to model them, not to start a new process every time the cycle turns.

Take heed: Your student project may compare one ARIMA process for the entire time series vs separate ARIMA processes for each period.

For your student project, use periods cautiously. They are appropriate if the process is distinct and homogeneous for each time period. We summarize below the FED policies in the three sub-periods.

Post-World War II Federal Reserve Board policy explains the changes in the Treasury bill time series. Just as a pricing actuary knows the policy provisions, type of insurer, and extent of market competition to set optimal rates, a statistician should know the attributes of the time series to fit an ARIMA process. But the student project focuses on the statistics, just as the SOA and CAS exams focus on the actuarial procedures. You are not required to know the economic and financial effects on the time series, but this knowledge may help you fit an ARIMA process.

- From the end of World War II (1945) through the mid-1970's, the U.S. economy expanded briskly. Government officials worried about Depression-era deflation, not the mild inflation of an expanding economy. Inflation was thought to be an antidote to unemployment, which had been high during the Depression (almost one third of the labor force in 1932). The federal government and the Federal Reserve Board believed that mild inflation was beneficial, in that it restrained unemployment and did not hamper economic prosperity.

This presumed relation of inflation and unemployment was an error, but it was the prevailing macroeconomic policy in the 1960's and 1970's. To reduce unemployment, the FED used expansionary monetary policy. Inflation and interest rates had steady upward trends. The time series of interest rates is not stationary, though real interest rates and the first differences of nominal interest rates may be stationary. You select the end-point of this period: the periods flow into one another; they are not distinct.

NEAS Time Series Project Documentation

Take heed: These comments apply to other interest rate time series as well, such as Treasury bonds or one year Treasury bills.

- From the late 1970's through the early 1980's, inflation and interest rates were high and volatile. See the large swings in interest rates in the graph on the illustrative worksheet. The high and volatile rate reflect (i) the mistaken macroeconomic policies of these times and (ii) the supply shocks of OPEC oil price increases (1973 and 1979).

Take heed: A volatile time series can be stationary, but it is hard for the statistician to distinguish trends from random fluctuations in a short, highly volatile time series.

- Paul Volcker became chairman of the Federal Reserve Board in 1981 and adopted a monetarist perspective (Milton Friedman's views). The money supply grew at a steady, slow rate. Interest rates and inflation declined. Greenspan continued Volcker's policy.

Take heed: We do not give Volcker credit or deny it to him. Some economic historians say Reagan was lucky to be President when Volcker chaired the FED. Others credit Reagan with stopping inflation, and Volcker was lucky to chair the FED in these years.

Most likely, a single ARIMA process is not an appropriate model for all three periods.

- You form a stationary time series several ways, as described in this project template.
- You can also combine methods, such choosing a sub-period, detrending, taking first differences, and forming a structural model.

If you select appropriate periods, state what periods you use and justify your choice.

Illustration: An analysis of three month Treasury bill rates might say:

- I graphed the rates, to see if the time series has the same process in all periods. I used rates for 1982 - 2000, which seem to have a downward drift. I used real interest rates to offset the decline in inflation during these years.
- I took first differences of the rates, to see if one ARIMA process could model all the post-World War II rates. I then excluded years 1979-1982 to see if excluding years with high volatility improved the fit.

Take heed: Some student projects compare two time periods. You might compare abortion rates before and after Row v Wade. Your write-up might say:

I fit ARIMA models to U.S. abortion rates for the year before and after Row v Wade. I found that the drift differs in the two time periods, and separate ARIMA models are needed.

Take heed: You may not have the exogenous knowledge to select proper time periods. You may select time periods based on the mean, drift, or volatility of the observations.

NEAS Time Series Project Documentation

Illustration: Your student project may say: “Overnight LIBOR rates for 1982 - 2006 decline for several years and then increase. (You would specify the years; see the project template on LIBOR rates.) My student project examines if the time series are really different:

- I took first differences of LIBOR rates in each period. The mean first difference differs by period, but each period can be fit by an ARIMA(2,1,0) process.
- I used real interest rates LIBOR divided by expected inflation. Each period can be fit by an ARMA(2,0) process.

Take heed: In short time series, drifts are hard to distinguish from random fluctuation.

Illustration: The years 1979-1982 is a short period of volatile interest rates. A difference of 2 or 3 months changes the observed drift. The drift is *not robust*. If you find a drift in the time series, consider also the volatility of the values and the length of the period.

Illustration: If a 20 year period has a drift of +2% per annum, each 10 year sub-period should have a drift of about 2% per annum. If the first ten years have a drift of +5% and the second ten years have a drift of -1%, the overall drift of +2% is not robust.

Interest rates seem to show drifts for short periods, such as several months of increasing rates followed by several months of decreasing rates. This may reflect an ARIMA(1,1,0) process and random fluctuations.

Illustration: Suppose interest rates are an ARIMA(1,1,0) process with $\delta = 0$ and $\phi_1 = 80\%$.

- If interest rates increase 100 basis points in January 20X6 because of random fluctuations, they are expected to increase 80, 64, 51, 41, and 33 basis points in each of the next five months.
- If they then decrease 100 basis points in July 20X6 because of random fluctuations, they are expected to decrease 80, 64, 51, 41, and 33 basis points in each of the next five months.

Illustration: For a period of 1 month and a volatility of 0.1% per month, even a time series with no drift may show a drift of *about* 0.1% (either positive or negative)..

- The drifts of +0.02% and -0.01% in the first and third periods reflect FED policy.
- The observed drift in the middle period reflects the short time period and high volatility.

NEAS Time Series Project Documentation

SUB-PERIODS AND STATIONARITY

{Taking differences may convert a homogeneous time series to a stationary time series. Separating a time series into sub-periods may help several ways.}

Jacob: Are the time series stationary in each sub-period?

Rachel: The first and third periods, with upward or downward drifts, are not stationary, but their first differences may be stationary.

- You may compare first differences for each sub-period vs for the entire time series.
- You may compare real interest rates for each sub-period vs for the entire time series.

When analyzing sub-periods for an economic or financial time series:

- Examine the raw time series, using first and second differences.
- Detrend the time series or use real interest rates or real dollars.
- Use a structural model by regressing the time series on other indices.

{*Take heed:* These are suggestions for time series analysis. Decide how to fit an ARIMA process. We explain de-trending (real interest rates) and structural models below. They are not required for the student project, but they are good topics.}

NEAS Time Series Project Documentation

Step #6: *DE-TREND*

A time series may combine several patterns reflecting several explanatory variables.

Illustration: Corporate bond rates for auto manufacturers combine expected inflation, real interest rates, business cycles, and default probabilities for the issuing firms.

It is easier to fit an ARIMA process to each piece separately than to the combination.

Illustration: Expected inflation may be an ARIMA(1,1,0) process, real interest rates may be an ARMA(1,1) process, and business cycles may be an oscillatory process.

- It is often easier to model a time series in real dollars than in nominal dollars.
- The same is true for other time series that are functions of a changing measure.

Illustration: It is easier to fit an ARIMA process to GDP per capita than to a country's total GDP. Population growth or decline is like inflation or deflation. We fit separate ARIMA process to (i) population growth and (ii) GDP per capita.

If the time series is in dollars (or other currency), convert it to real dollars. Divide the dollars by the CPI. We provide several CPI indices on the discussion forum for deflating.

Take heed: Decomposing a time series and fitting ARIMA processes to its parts is a good student project. Your project may compare ARIMA processes for nominal vs real interest rates.

Recommendation: Choose a time series that is composed of two or more elements. To keep your student project manageable, use two pieces:

- Nominal interest rates = real interest rates + expected inflation
- Corporate interest rates = risk-free rates + default spreads
- Insurance premium = new business premium + renewal business premium
- Auto sales = new car sales + used car sales

If you want a more adventurous project, decompose the time series into 3 or 4 parts.

Compare an ARIMA process fitted to the combined time series vs ARIMA processes fitted to each part. See which model forecasts better.

NEAS Time Series Project Documentation

ILLUSTRATION: NOMINAL AND REAL INTEREST RATES

To convert nominal interest rates to real interest rates, divide by the inflation rate in the previous month as a proxy for expected inflation.

Illustration: The nominal interest rate is 8% on June 1, 20X8. The CPI is 130 on 6/1/20X8 and 125 on 5/1/20X8. The real interest rate on June 1, 20X8, is

$$1.08 / (130 / 125) - 1 = 3.85\%$$

For your student project, do the following:

- Copy a CPI index from the discussion forum to the work-sheet with your interest rate time series. You have a choice of seasonally adjusted or not seasonally adjusted CPI.
- Compute the inflation in each month as the ratio of two CPI figures.
- The interest rate is annualized. To annualize the inflation rate, raise it to the 12th power.
- Choose a proxy for expected inflation. You might use the actual inflation the previous month or a moving average of inflation rates in the previous several months.
- Divide (1 + nominal interest rate) by (1 + expected inflation rate).

Intuition: The objective of ARIMA modeling is to separate trends, cycles, seasonality, and stochasticity. Time series often overlay a stationary ARIMA process on a trend, drift, or cycle. Inflation may combine a trend with seasonality, and interest rates may be a mean reverting pattern overlaid on a business cycle. ARIMA modeling separates each part.

The procedures differ for each part. We may use

- First differences for the trend
- An autoregressive process for the mean reversion
- A structural model for the business cycle
- A sine pattern for the seasonality.

Take heed: Pick an appropriate inflation index. For health insurance loss costs, you may use medical CPI, not total CPI. Explain if you use seasonally adjusted or non-seasonally adjusted CPI.

- If your dollar-denominated time series has the same seasonal pattern as the CPI, use non-seasonally adjusted CPI to de-trend.
- If your dollar-denominated time series is not seasonal, use seasonally adjusted CPI to de-trend.

Take heed: De-trending, adjusting for seasonality, and fitting the proper ARIMA process is a good student project for insurance loss cost trends. Don't presume that actuaries have already developed optimal trend models. On the contrary: actuaries are just now beginning

NEAS Time Series Project Documentation

to use sophisticated ARIMA processes for loss cost trends. The first CAS paper on ARIMA modeling of loss cost trends was published recently.

Use spreads in the same manner as de-trending. Real interest rates are the spread over expected inflation. Corporate bond rates are the spread over risk-free rates.

Illustration: Instead of analyzing the Moody's AAA Bond rate, analyze the corporate bond spread to Treasury bonds. You may add a structural component (such as GDP growth) as a proxy for default probabilities.

Intuition: Corporate bond rates are a mix of several items:

- The term structure of interest rates, which can be modeled by Treasury securities.
- Expected inflation, which can be modeled by CPI changes and money growth.
- Real interest rates, which is best modeled by Treasuries or LIBOR rates.
- Default expectations, which can be modeled by GDP growth.

Your student project need not model all aspects of the corporate bond yield. Explain how the corporate bond spread to Treasury securities eliminates duration, risk-free rates, and expected inflation, leaving a time series governing by default expectations and business cycles. Regressing the corporate bond spreads on GDP growth may eliminate the business elements, leaving a stationary time series that can be fit to an ARIMA process.

Note: An early use of ARIMA processes was to model economic cycles in the U.S. and Great Britain. You have probably studied the lay version in a college economics course: prosperous years cause consumer over-confidence (or some other item) that leads to an over-heating economy and a recession. The CAS Exam 5 syllabus has a reading on underwriting cycles with a similar ARIMA perspective. The on-line NEAS macroeconomics course has a different perspective on business cycles (called real business cycle theory).

Take heed: The real interest rates formed on the illustrative work-sheet do not form a stationary time series. The real interest rate should be between 0.5% and 4%. A figure outside this range may mean the estimated expected inflation rate is not correct.

- Last month's inflation is sometimes greater than the nominal interest rate, implying that investors believe inflation will fall.
- Inflation may be low one month for exceptional reasons (perhaps oil prices fell because of a peace treaty in the Middle East), but investors presume future inflation will be high.

For the real interest rates on the illustrative worksheet, you still take differences, divide the time series into homogeneous periods, or de-trend the time series to fit an ARIMA process. Your student project may explain whether using real interest rates improves the ARIMA fit.

NEAS Time Series Project Documentation

Step #7: *STRUCTURAL MODELS*

ARIMA models are proxies for the true explanatory model. Treasury bill rates are affected by expected inflation, economic growth, consumer confidence, political stability, demand for money, other investment opportunities, and similar factors.

If we knew all the influences on Treasury bill rates and values of all explanatory variables, we could form a model $R = f(X_1, X_2, X_3, \dots) + \epsilon$. The model might also have lagged terms, such as economic growth last month or two months ago.

Constructing a complete model may not be possible, since we do not know the influences on Treasury bill rates. An ARIMA model is a proxy.

- The explanatory variables are also stochastic time series.
- The ARIMA process puts all the explanatory variables into one time series.

Illustration: Suppose nominal interest rates are a function of real interest rates, expected inflation, economic activity, and demand for money. Each of the explanatory variables can be modeled as an autoregressive or moving average process. Ideally, we might forecast each explanatory variable and then derive the indicated Treasury bill rate each month.

Structural models are a compromise between the full model and a simple ARIMA process.

Illustration: Suppose real interest rates depend on economic growth, consumer confidence, political stability, demand for money, other investment opportunities, and similar factors. Economic growth can be measured by GDP growth and it has a large effect on real interest rates. The other explanatory variables are hard to measure and have less effect on real interest rates.

We regress the real interest rate on real GDP growth. We fit the residuals of the regression to an ARIMA process. These residuals reflect the effects of the other explanatory variables.

A diffuse pattern of real interest rates may be a simpler ARIMA process after regressing on real GDP growth. Your student project may examine the effects of each item.

Illustration: Suppose your time series is the 20 year corporate bond AAA rate.

- Model the nominal corporate bond rate. Take first (and perhaps second) differences to get a stationary time series. The goodness-of-fit tests may be poor, since the ARIMA model combines several influences on corporate bond rates.
- Model the *real* corporate bond rate (offset with the change in the CPI). The in-sample fit may improve, but it may be difficult to forecast future interest rates unless you can forecast future CPI values.
- Model the corporate bond *spread* to *three month Treasury bills*. The time series may be easier to fit, and you may not need to take differences. But duration effects are

NEAS Time Series Project Documentation

overlaid on business cycle effects, and the fit may not be good. [The duration effect is the slope of the term structure of interest rates.]

- Model the corporate bond spread to 20 year Treasury bonds. Corporate bonds and Treasury bonds have about the same duration, so the term structure of interest rates does not affect the spread. With just spread cycles and default expectations remaining, the ARIMA fit may improve.
- Regress the corporate bond spread on a measure of business activity, such as GDP growth. Economists do not agree about the influence of GDP growth on interest rates. In general, as GDP increases, bond defaults decrease, so the corporate bond spread narrows. Your time series of spread cycles may resemble a simpler ARIMA process.

Ideally, GDP growth is annualized, seasonally adjusted, converted to real terms, shown by month, and perhaps lagged. We explain each of these adjustments in the discussion forum posting on structural models. They are important for the economic analysis, not for the student project. If the explanatory variable is not in the form you want, explain the desired form in the student project write-up but do the analysis with the data you have.

Structural models are best if we have suitable explanatory variables. The structural model shows that you can decompose a time series into its pieces and use statistical tests to see if the decomposition gives a better ARIMA model. Even if you are not persuaded that a structural model is correct, you may examine it for your student project.

Illustration: Some economists say that higher GDP growth raises real interest rates, since firms wish to borrow (and invest) more in prosperous year and demand for cash rises. You may not be persuaded by the economic reasoning, and you would like to test it.

- Your regression analysis student project may examine several macroeconomic indices that may affect real interest rates, such as GDP growth, budget deficits, foreign interest rates, and employment rates.
- Your time series student project may fit ARIMA processes to the residuals of each regression line.

{Structural models are discussed in a separate discussion forum posting, and an example is in a separate illustrative worksheet.}

NEAS Time Series Project Documentation

Step #8: CORRELOGRAMS

Terms: Compute sample autocorrelations of the time series. The student project fits an ARIMA process to the observed time series. We test the fit several ways:

- The sample autocorrelations of the time series should have the same pattern as the implied autocorrelations of the ARIMA process.
- The residuals of the time series from the fitted ARIMA process should be random, so the sample autocorrelations of the residuals should be those of a white noise process.

Terms: The sample autocorrelation depends on the lag, such as 30% for a lag of 1, and 20% for a lag of 2. The sample autocorrelation as a function of the lag is the sample autocorrelation function.

The graph of the sample autocorrelation function is the correlogram. When you first start ARIMA fitting, the graph (correlogram) is easier to read than a table of autocorrelations. But you use the table of sample autocorrelations for the statistical tests.

The raw time series may not be stationary. Explain the pattern of sample autocorrelations. The correlogram should confirm the pattern you see in the graph.

Your student project proceeds in two directions.

- The sample autocorrelations reflect trends, seasonality, cycles, and other patterns in the time series. Adjust the data to remove trends, seasonality, and cycles, and compute again the sample autocorrelations.
- After fitting an ARIMA process to the time series, determine the sample autocorrelation function of the residuals. The residuals should be a white noise process, which we test by the pattern of the sample autocorrelation function.

Illustration: Inflation is seasonal, so short interest rates, such as overnight LIBOR, may also show seasonality. If you believe the rates are seasonal, you may de-seasonalize the data and re-compute the sample autocorrelations.

Take heed: We provide several time series of LIBOR rates. These are excellent time series for student projects, since the short rates (overnight, one week, and two weeks) have daily observations. We also show LIBOR rates in other currencies, such as the Euro, and we show inflation and foreign currency exchange rates. These indices are related and they are all stochastic, so you can form various structural models.

Take heed: Examining the seasonality of overnight LIBOR rates is not easy, since trends, seasonality, cycles, and random fluctuations are overlaid on each other.

NEAS Time Series Project Documentation

- Convert each rate to its difference from a centered moving average of one year. The LIBOR time series shows rates for business days only, so a year is somewhat less than 250 days. Use ratios for the differences or first take logarithms.
- Compute long-term averages for a given day. Use all the years on the NEAS web site.
- Graph these averages to see if they show a seasonal pattern.
- Confirm your results with a correlogram. If LIBOR rates are seasonal, the correlogram should start high (reflecting the autoregressive process), decline to about 125 days (half a business year), and then rise to a local maximum at slightly less than 250 days.

You learn how to test for seasonality in a complex time series.

Illustration: The daily temperature, after adjusting for seasonality, may have a trend. Even if the trend is obscured by random fluctuations, it may be seen in the correlogram. After detrending the temperature, the correlogram may indicate a stationary time series. You then fit an ARIMA process so that the residuals show a white noise process.

We discussed earlier segmenting a time series into periods. The periods are indicated if the full time series is not stationary but each period is stationary.

Illustration: Treasury bill rates may show patterns, such as increasing or decreasing rates. You may divide the time series into periods, each of which has its own pattern.

- The first differences in each period may be stationary processes with different means.
- De-trended interest rates, real interest rates, or the interest rates regressed on another index may be distinct ARIMA processes in each period.

NEAS Time Series Project Documentation

Step #9: Fit a Model

The model depends on the sample autocorrelation function. For the student project, look at four models: AR(1), AR(2), MA(1), and ARMA(1,1), along with their first differences: ARIMA(1,1,0), ARIMA(2,1,0), ARIMA(0,1,1), and ARIMA(1,1,1). You may fit more complex models if you believe they are needed, but we do not require more complex models for the student project.

Take heed: If the sample autocorrelation function (the correlogram) indicates that an AR(1) or AR(2) model is appropriate, and if the fitted model passes the Box-Pierce Q statistic, you need not fit MA(1) or ARMA(1,1) processes.

Take heed: Fit the ARMA models or the ARIMA models, not both.

- If the time series is stationary, use the ARMA models, not the ARIMA models.
- If the time series is not stationary, use the ARIMA models, not the ARMA models.

The model fitting is described in separate discussion forum postings.

- For the AR(1) and AR(2) processes, fit a lagged regression.
- For the MA(1) and ARMA(1,1) processes, use the Yule-Walker equations.

After fitting the model:

- Compare the sample autocorrelation function with the autocorrelation function.
- Use Bartlett's test and the Box-Pierce Q statistic to test the goodness-of-fit.
- For forecasts and evaluate the out-of-sample goodness-of-fit.

Each of these is described in separate postings.

Take heed: Document your work as you do it. When you finish the student project, edit your documentation using the guidelines on the discussion forum. The documentation is the write-up for the student project.

NEAS Time Series Project Documentation

Step #10: *AUTOCORRELATIONS AND SAMPLE AUTOCORRELATIONS*

After you fit an ARIMA process, compare the implied autocorrelations from that process with the sample autocorrelations of the time series.

- The correlogram shows the pattern of the observed time series.
- The actual time series is stochastic, so the correlogram is not smooth.
- You fit an ARIMA process, which has a smooth autocorrelation function.
- Form a graph overlaying this smooth function on the correlogram.
- If the smooth function is close to the sample autocorrelations, the fit is good.

Question: Why do we compare the correlogram to the autocorrelation function? Why not compare the observed time series to the pattern of the ARIMA process?

Answer: An ARMA process becomes a straight line at the mean after a few lags, and an ARIMA process becomes a diagonal line with a slope equal to the drift after a few lags. The forecasts from any ARMA process look the same.

Question: Why not compare the graph of the one month forecasts from the ARIMA process to the actual time series?

Answer: This comparison is the sum of squared deviations, which we use to choose the best model. The algebra gives a figure that varies by ARIMA process. On a graph, it is hard to see which process fits best. An AR(1) process looks about the same whether $\phi_1 = 40\%$ or 60%.

The graph of the autocorrelation function looks different for each ARIMA process. It is a good marker of the ARIMA process, and it is easy to compare with the correlogram.

Take heed: The discussion forum posting on *time series simulations* has correlograms for several ARIMA processes. The correlogram reveals the time series process. Know what the correlogram of each ARIMA process looks like, so you can identify a reasonable model for your time series.

Take heed: The discussion forum posting on *time series techniques* and the attached Excel work-sheet gives cell formulas and a VBA macro for sample autocorrelations. That posting uses the interest rates in this project template as the illustration.

Review the illustrative worksheet on time series techniques. Make sure you understand the cell formulas and the correlograms.

The VBA macro is optional. We have not made the macro a screen based facility, so you don't treat the macro as a black box. Open the macro in the VBE (the editor) and choose the length of the time series and the number of lags. The defaults are the entire time series

NEAS Time Series Project Documentation

and all lags. Choose a smaller number of lags if the time series has so many observations that the number crunching takes too much time (e.g., more than 10,000 observations).

Choose a different length for the time series if you want the correlogram for a sub-period.

Illustration: Suppose the time series has 240 observations of monthly interest rates. You want to form correlograms separately for the first 10 years and the second 10 years.

- Place the cursor on the first observation. Choose a length of 120 to form a correlogram from the first 120 observations.
- Place the cursor on the 121st observation. The default is to the end of the time series. This forms a correlogram from the second 120 observations.

If you don't want to use macros, you can use the cell formulas. The macro is more efficient. If you have a slow machine and a time series with tens of thousands of observations, you must use the macro. For interest rates, you can use cell formulas.

NEAS Time Series Project Documentation

Step #11: *LAGGED REGRESSIONS*

For the AR(1) and AR(2) processes, fit a lagged regression. Copy the stationary time series to a second column and use Excel's *REGRESSION* add-in.

Take heed: If the original time series is not stationary and you took first differences, do the lagged regression on the stationary first differences.

Take heed: Excel comes with the analysis tool-pack, but you must load it on your copy.

Illustration: The time series values are in Cells A11:A200 and we fit an AR(2) process.

- Copy Cells A11:A200 to Cells B12:B201 and also to Cells C13:C202. *Take heed:* Column B has the values lagged one period; Column C has values lagged 2 periods.
- Invoke Excel's *REGRESSION* add-in. For Excel versions before 2007, click the *TOOLS* menu \Rightarrow *DATA ANALYSIS* \Rightarrow *REGRESSION*. For Excel 2007, click *DATA* \Rightarrow *ANALYSIS* \Rightarrow *DATA ANALYSIS* \Rightarrow *REGRESSION*. The sequence of clicks may vary for your version of Excel.
- On the regression screen, select Cells A13:A200 as the Y values and Cells B13:C200 as the X values. Do not include headers, since the values in Cells A13:C13 are time series observations, not labels for the variables. Label the output of the regression.
- Select *SHOW RESIDUALS* (not standardized residuals). If the time series is stationary, the mean does not vary, so the standardized residuals should be almost the same as the ordinary residuals.
- Place the output on a new worksheet. Name the new worksheet as the ARIMA process, such as AR(2). This work-sheet will also show the correlogram of the residuals, their Box-Pierce Q statistic, and a comparison with the theoretical autocorrelations.

GOODNESS-OF-FIT

The goodness-of-fit measures used in regression analysis, such as the R^2 of the regression and the significance of the ordinary least squares estimators, are not the most important measures for ARIMA fitting.

- If you have taken the regression analysis course, comment on these items.
- If you have not taken the regression analysis course, you may ignore these items and comment just on Bartlett's test and the Box-Pierce Q statistic.

You can use the \bar{R}^2 (the R^2 adjusted for degrees of freedom) to compare AR(1) and AR(2) processes. If the adjusted R^2 does not increase from AR(1) to AR(2), use AR(1). The R^2 (or adjusted R^2) won't help you decide between AR(1) and MA(1).

Check that the ordinary least squares estimators are significant. For a lagged regression with two periods – an AR(2) model – the correlation of the explanatory variables affects the t values.

NEAS Time Series Project Documentation

Form the sample autocorrelations of the residuals, and use the Durbin-Watson statistic, Bartlett's test, and the Box-Pierce Q statistic. The illustrative worksheet on time series techniques provides cell formulas and a VBA macro. The cell formulas and macro do all the arithmetic. Your write-up should explain how you use and interpret these tests.

NEAS Time Series Project Documentation

Step #12: MOVING AVERAGE MODELS

For MA(1) and ARMA(1,1) processes, the textbook uses nonlinear regression to estimate the ARIMA parameters and Yule-Walker equations for initial estimates of the parameters.

- Nonlinear regression is not covered in the on-line courses.
- The Yule-Walker estimates are close enough for the student project.
- Excel's *SOLVER* add-in makes it simple to use the Yule-Walker equations.

Take heed: You can solve the Yule-Walker equations by hand; you don't need solver. For MA(1) and ARMA(1,1) processes, solve a quadratic equation for θ_1 .

Illustration: Suppose the sample autocorrelations in a time series of 10,000 observations are 45% for lag 1 and 10% for lag 2. The sharp drop from lag 1 to lag 2 suggests a moving average parameter, so we test an ARMA(1,1) model. Your write-up explains why you use an ARMA(1,1) process:

- With 10,000 observations, the standard error of the sample autocorrelations of a white noise process is 1%.
- For an AR(1) process, $45\%^2 = 20.25\% > 10\%$.
- Even if the sample autocorrelations are off by one standard deviation, $46\%^2 = 21.16\% > 11\%$.

Examine both an AR(2) process and an ARMA(1,1) process.

- For an AR(2) process, we fit the autoregressive parameters with linear regression. An AR(2) process has $\phi_1 > 0$ and $\phi_2 < 0$, as might occur in a cyclical time series.
- For MA(1) and ARMA(1,1) process, we use the Yule-Walker equations.

For an MA(1) process,

$$\theta_1 = \frac{-1 \pm \sqrt{1 - 4\rho^2}}{2\rho} = (-1 + \sqrt{1 - 4 \times 0.45^2}) / (2 \times 0.45) = -0.627$$

An MA(1) process would have a zero autocorrelation of lag 2. The MA(1) process fits no better than the AR(1) process. Once you have fit all four models, use Bartlett's test and the Box-Pierce Q statistic to see which fits best. The fit depends on all observations; we can't decide which process fits best from the first two sample autocorrelations alone.

Equations 17.58 and 17.59 on page 536 gives the autocorrelations for lags 1 and 2 of an ARMA(1,1) process.

$$\rho_1 = \frac{(1 - \phi_1\theta_1)(\phi_1 - \theta_1)}{1 - 2\phi_1\theta_1 + \theta_1^2}$$

$$\rho_k = \phi_1\rho_{k-1} \quad \text{for } k \geq 2$$

We know the sample autocorrelations. We solve the equations for the ARIMA parameters.

- Estimate ϕ_1 as $10\% / 45\% = 2/9 = 22.22\%$.
- Estimate $Z = \theta_1$ as $45\% = [(1 - Z \times 2/9) \times (2/9 - Z)] / (1 + Z^2 - 2 \times Z \times 2/9)$

We solve the quadratic equation for $Z = \theta_1 = -0.290855735 \approx -0.291$.

We estimate δ from the ARIMA parameters and the mean of the time series.

- For the MA(1) process, δ is the mean of the time series.
- For the ARMA(1,1) process, δ is the mean of the time series times the complement of the autoregressive parameter.

Take heed: Instead of solving the quadratic equation by hand, we can use Excel's *SOLVER* add-in. Excel's *GOAL SEEK* can also do the arithmetic.

Choose cells for θ_1 and ϕ_1 , such as Cells B2 and B3. Name these cells one of three ways:

- Use the *DEFINE NAMES* dialogue box. (Select Cell B3, type *phis1* in the dialogue box, and press *ADD*.)
- Use the name box on the left end of the formula toolbar. (Select Cell B3, type *phis1* in the name box, and press enter.)
- Place the strings "theta1" and "phis1" in Cells A2 and A3. Use the *CREATE NAMES* dialogue box to assign these names to the values in Cells B2 and B3.

Take heed: You don't have to name the cells. You can refer to the cells with absolute references, such as $\$B\3 . But your cell formulas become hard to understand.

- Choose cells for the sample autocorrelations of lags 1 and 2, such as Cells C2 and C3.
- Name these cells "rhos1" and "rhos2".

Type cell formulas into Cells C2 and C3.

- For Cell B3, type "= $\text{rhos1} / \text{rhos2}$ ".
- For Cell C2, type "= $((1 - \text{phis1} * \text{theta1}) * (\text{phis1} - \text{theta1})) / (1 - 2 * \text{phis1} * \text{theta1} + \text{theta1}^2)$ ".

Use solver to find θ_1 . We show the figures in the illustrative worksheet *YULEWALKER*.

Take heed: In Excel versions before 2007, we use names *phi1*, *rho1*, and *rho2*, not *phis1*, *rhos1*, *rhos2*. Excel 2007 changes the names *phi1*, *rho1*, and *rho2* to *phi1_*, *rho1_*, and

NEAS Time Series Project Documentation

rho2_, since phi1, rho1, and rho2 can refer to cells in Excel 2007. Typing underscores can be a hassle, so use names with 4 letters and a digit. “phis1” stands for phi-sub-1.

Take heed: If *SOLVER* has a poor starting figure for θ_1 , it may not find a solution.

- If Cell B2 has 1.000 when you invoke *SOLVER*, no solution is found.
- Start with a reasonable value for θ_1 . A good starting value is $\phi_1 - \rho_1 = 0.2222 - 0.45 = -2278 = -22.78\%$.

{This step-by-step guide fits MA(1) and ARMA(1,1) processes with Yule-Walker equations. With better statistical software, you can solve a nonlinear regression, as discussed in the textbook (not on the time series syllabus). Alternatively, you can fit an MA(1) or ARMA(1,1) process with Excel's *SOLVER* add-in. We show the *SOLVER* technique in a separate discussion forum posting.}

NEAS Time Series Project Documentation

Step #13: *GOODNESS-OF-FIT TESTS*

Test the goodness-of-fit of the moving average models using Bartlett's test and the Box-Pierce Q statistic, just as for autoregressive processes.

- Excel's *REGRESSION* add-in gives residuals. The VBA macro in the illustrative worksheet gives the sample autocorrelations and the Box-Pierce Q statistic.
- Excel does not have an add-in for the residuals of a moving average process.
- Compute the residuals manually. The cell formulas are simple, and the process shows you understand the goodness-of-fit test.

Illustration: Suppose the time series values are in Cells A11:A200.

- Column B is the forecasts. In Cell B11, type “=A11.”
- Column C is the residuals. In Cell C11, type “=A11-B11.”

Take heed: Column B is the forecast *made in the previous period*.

- Cell A12 has the time series value for Period 2.
- Cell B12 has the forecast for Period 2 made in Period 1.
- Cell C12 has the residual for Period 2.

Name three variables as delta, theta1, and phis1, using the methods described earlier.

Take heed: You don't have to place these values in the work-sheet. You can name constants in the “define names” dialogue box.

- Place these value in work-sheet cells and name the cells if these values change as you analyze the data. This is the more common scenario.
- Use named constants if these values do not change.

Take heed: For the MA(1) process, the phis1 variable is zero.

- In Cell B12, type “=delta + A11 × phis1 – C11 × theta1.”
- Copy Cell B12 to Cells B13:B200.
- Copy Cell C12 to Cells C13:C200.

Excel shows residuals for 190 observations. The residual for the first observation is zero (by definition), so we do not use it. We set B11 equal to A11 because we have no values for A10 and C10.

For the other 189 residuals, form the sample autocorrelations. Use Bartlett's test and the Box-Pierce Q statistic to test the quality of the ARIMA fit.

NEAS Time Series Project Documentation

The VBA macro does most of the number crunching. Place the cursor in Cell C12 and run the macro. The macro also computes the values for the Box-Pierce Q statistic. Bartlett's test is subjective; you must count the sample autocorrelations above a critical value.

Take heed: This procedure slightly overstates the residuals by assuming the first residual is zero. The distortion is not material.

NEAS Time Series Project Documentation

Step #14: *FORECASTS*.

The ARIMA model helps forecast future interest rates. A student project that fits an ARIMA process may have a section testing the fit with an out-of-sample test.

Take heed: A forecasting section is not required for all student projects. It is useful, since it shows that you understand the purpose of ARIMA modeling. If your student project does a good analysis of other topics, you need not include forecasting.

To test the quality of the forecasts, leave out the last N observations and forecast them.

Illustration: For a time series of three month Treasury bills ending June 2000, you might

- use observations through December 1999 and forecast the next six months.
- use observations through December 1998 and forecast the next 18 months.

A separate discussion forum posting covers forecasts and out-of-sample goodness-of-fit tests. For your student project, leave out the last year of observations from the data used to fit the ARIMA model. Do the ARIMA fitting using the techniques in this project template.

Once you have fit two or more models, compare their forecasts of the final observations. Keep in mind that an out-of-sample fit is easily distorted by random fluctuations. The ARIMA modeling and forecasting completes your student project.

Take heed: The project template discusses many techniques that you can apply to the time series. You need not use all the techniques. The illustrative worksheet shows some of the techniques, not all of them.

NEAS Time Series Project Documentation

TS STUDENT PROJECTS: TOPICS

(The attached PDF file has better formatting.)

Some candidates wonder what they can model with ARIMA processes. The NEAS web site has project templates for financial series (interest rates, CPI, GDP, unemployment, gas prices), daily temperature, sports won-loss records, and other topics

These project templates have suggestions for time series analysis and ARIMA fitting. You have hundreds of time series to choose from. You can take ideas from the NEAS postings or past student projects and apply them to your time series.

Illustration: The weather service publishes daily temperature histories for hundreds of locations. Pick a location, such as your home town, and a time period, such as 1905-2005, and fit an ARIMA process.

For a more imaginative student project, compare two adjoining cities with different climates, such as San Francisco and Oakland, or compare a weather station before and after it became a large city.

We have posted Excel data bases of daily temperatures by location. Because we provide so much guidance for this student project, include illustrative worksheets with cell formulas and VBA macros, we expect a good student project on this topic.

Illustration: In past years, we accepted student projects looking at monthly temperature, since it was not easy to get daily temperature for long periods. Now you have all the daily temperature time series you need; don't use monthly temperature.

Take heed: You can do a student project on many other weather topics, such as hourly temperature, rainfall, smog in Los Angeles (or China), by getting other information from weather web sites.

Sales: You can use sales for any industry, including insurance. The web has thousands of sites with business data. A monthly time series shows seasonality, trends, and often business cycles.

Take heed: We show several time series with industry sales (autos, homes, electronics, furniture, department stores, and so forth). You can use any of these time series, a time series from other web sites, or a time series from a private firm.

Birth rates, baby names, marriage rates, divorce rates, mortality rates, and other life events are good time series for the student projects. These topics are fascinating, and the data on the internet are extensive. We show some sample files, such as Excel workbooks from the FBI crime web site.

NEAS Time Series Project Documentation

Take heed: We encourage you to post ideas, links to web sites, and Excel attachments with possible data. These ideas may stimulate student projects by you or other candidates.

Note: Some candidates send us an excess work-book and ask if they can use the data for a student project. We do not answer these questions:

- The answer is almost always yes, if you properly use statistical techniques to analyze the data. We don't want to say "yes" and then receive a student project that does not do the needed work.
- If the answer is "no" because the data have too few observations, the problem is obvious. Don't try to fit 10 years of profits to an ARIMA process.
- If the data contain missing observations or other problems, your write-up explains the problems and your solutions.

You can also use data from your actuarial work, such as personal auto average claim severities. You have these data in Excel work-books, you understand the attributes of the data, and you may already have done much of the work. The student project may be easy to complete. It is fine to use data that you have already collected for other purposes.

Imagination and web surfing give dozens of ideas. Candidates have done student projects on sunspots, airplane crashes, gas prices, border crossing, and other topics.

Illustration: The elections in 2008 will generate hundreds of ideas for student projects, both regression analysis and time series. You can examine the effects of explanatory variables on voting behavior (age, sex, wealthy, religion, ethnic group, education, and so forth).

We post past student projects for many reasons: to illustrate how to apply the statistical techniques, how to adjust for trend, seasonality, cycles, and changes between time periods, and to point out common errors that make your analysis less efficient. The past student projects stimulate ideas for other candidates. Even if the analysis in the student project is not that good, we may post the first paragraph to give ideas to other candidates.

- ~ Some candidates worry that if they pick their own time series, the ARIMA modeling might not work well.
- ~ They worry that if they do not fit an ARIMA process that passes Bartlett's test and the Box-Pierce Q statistic, their student project does not pass.
- ~ They worry that if they choose their own data, they might overlook important effects in their models.
- ~ They worry that if they pick their own data and no ARIMA process fits well, NEAS will ask them to do the student project again.

The student projects are not intended to produce the true models. Actual statistical work requires extensive analysis of exogenous factors that influence the time series. The statistician spends several days reviewing the relations of the time series and other data.

NEAS Time Series Project Documentation

Do not hesitate to choose other time series or regression data sets. We do not reject student projects because we don't like the choice of topic. If you use statistical techniques properly, your project is fine.

Illustration: A student project on sunspots will not determine the real sunspot process. But sunspots are a good choice for a time series student project. They are a stationary time series or a homogeneous non-stationary time series with several attributes that can be fit to an ARIMA process. The student project selects the time intervals (day, month, or year), analyzes the stationarity or the order of homogeneity, and fits an ARIMA process.

For short time intervals, such as one week, sunspots are clearly autoregressive, since the same sunspot continues for longer than a week. Even over longer periods, sunspots have periods of high or low frequency. The type of ARIMA process is unclear, and you examine graphs, correlograms, trends, cycles, differences, and the autocorrelations.

Some astronomers believe that sunspots have long-term cycles. It is hard to model multi-year cycles with simple ARIMA processes, and we do not expect you to construct an AR(1) or AR(2) model which fits exactly. Your student project would focus on graphing the time series and the correlograms to identify the cycle.

Some candidates want confirmation from NEAS that a particular time series can be used for the student project. Use the following guides:

If the time series is a white noise process, choose something else. Few time series are pure white noise. You can recognize a white noise process intuitively: random draws from a distribution. If the time series seems like a white noise process, use Bartlett's test or the Box-Pierce Q statistic to test for white noise.

A random walk has first differences that are a white noise process, so random walks are not good for the student project. We don't recommend stock prices for the student project.

Take heed: We don't rule out random walks. You may compare the time series of stock prices from two firms in the same industry, such as two insurers. Your time series might compare the correlation of the two time series.

If the time series has few points or discrete values, choose something else. The time series values do not have normal distributions, and hypothesis testing can not be used (in the form taught in this course). Don't choose earthquake frequency by year in California. This is a white noise process with rare and discrete values.

If the time series has many points and is not a white noise process, it is fine.

NEAS Time Series Project Documentation