

Name: Tong Ren
Registration ID: 39899381
EMAIL: tongren1984@gmail.com
Regression Analysis Project
Winter 2012

Introduction

Concrete is the most important material in civil engineering of construction and maintenance of the physical and naturally built environment. The project analyzes the effectiveness of the relationships between ingredients and concrete compressive strength. The Excel Regression Add-In was used to run regression analysis.

Data

Data Source:

Website http://archive.ics.uci.edu/ml/machine-learning-databases/concrete/compressive/Concrete_Data.xls

Description: The data are related to concrete compressive strength depending on age and ingredients. The ingredients include cement, blast furnace slag, fly ash, water, superplasticizer, coarse aggregate and fine aggregate. The dependent variable (denoted by Y) is concrete compressive strength.

Number of cases: 1030

Variables

Dependent Variable: Concrete compressive strength -- quantitative -- MPa

Independent Variables:

1. Cement (component 1) -- quantitative – kg/m³
2. Blast Furnace Slag (component 2) – quantitative – kg/m³
3. Fly Ash (component 3) -- quantitative – kg/m³
4. Water (component 4) -- quantitative – kg/m³
5. Superplasticizer (component 5) – quantitative – kg/m³
6. Coarse Aggregate (component 6) – quantitative – kg/m³
7. Fine Aggregate (component 7) – quantitative – kg/m³
8. Age -- quantitative -- Day (1~365)

Equation and Variables

Where the equation and variables for the full model are:

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + \beta_7 X_7 + \beta_8 X_8 + \varepsilon_i$$

Where: Y = Concrete compressive strength -- quantitative -- MPa

α = Intercept

β_i = Least squares coefficients

X_1 = Cement (component 1) -- quantitative – kg/m³

X_2 = Blast Furnace Slag (component 2) – quantitative – kg/m³

X_3 = Fly Ash (component 3) -- quantitative – kg/m³

X_4 = Water (component 4) -- quantitative – kg/m³

X_5 = Superplasticizer (component 5) – quantitative – kg/m³

X_6 = Coarse Aggregate (component 6) – quantitative – kg/m³

X_7 = Fine Aggregate (component 7) – quantitative – kg/m³

X_8 = Age -- quantitative -- Day (1~365)

Hypothesis

The null hypothesis is that all least squares coefficients are zero: $\beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = \beta_6 = \beta_7 = \beta_8 = 0$

Data Analysis

Results from the Regression Analysis

The regression was performed for some combinations of the above stated variables. Results are illustrated below separately:

Model 1 (Using All Explanatory Variables)

Regression Statistics	
Multiple R	0.7845156
R Square	0.6154647
Adjusted R Square	0.6124517
Standard Error	10.399849
Observations	1030

ANOVA	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	8	176744.87	22093.109	204.26914	6.76E-206
Residual	1021	110428.16	108.15686		
Total	1029	287173.03			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	-23.163756	26.588421	-0.8711971	0.3838513	-75.337951	29.010439	-75.337951	29.010439
Cement	0.11978	0.0084894	14.1100	1.963E-	0.10312	0.13644	0.103126	0.136443

	53		32	41	67	39	7	9
Blast Furnace Slag	0.10384 72	0.0101362	10.2451 69	1.633E- 23	0.08395 71	0.12373 74	0.083957 1	0.123737 4
Fly Ash	0.08794 31	0.0125851	6.98786 6	5.03E-1 2	0.06324 74	0.11263 87	0.063247 4	0.112638 7
Water	-0.15029 79	0.0401793	-3.7406 843	0.00019 37	-0.2291 413	-0.0714 546	-0.229141 3	-0.071454 6
Superplasticizer	0.29068 69	0.0934599	3.11028 6	0.00192 09	0.10729 15	0.47408 23	0.107291 5	0.474082 3
Coarse Aggregate	0.01803 02	0.0093942	1.91928 66	0.05522 65	-0.0004 04	0.03646 44	-0.000404	0.036464 4
Fine Aggregate	0.02015 45	0.0107027	1.88312 1	0.05996 8	-0.0008 473	0.04115 62	-0.000847 3	0.041156 2
Age	0.11422 56	0.0054275	21.0457 41	5.841E- 82	0.10357 53	0.12487 59	0.103575 3	0.124875 9

The corresponding regression equation is:

$$Y = -23.1638 + 0.1198X_1 + 0.1038X_2 + 0.08794X_3 - 0.1503X_4 + 0.2907X_5 + 0.01803X_6 + 0.02015X_7 + 0.1142X_8$$

The adjusted R^2 statistic of 0.61245 means that about 61% of variation in the dependent variable, i.e. concrete compressive strength, can be explained by these eight variables which suggest a very reasonable relationship of these explanatory variables to concrete compressive strength.

The T-statistic for variable X_8 (age) is the highest (21.0457) suggesting a great deal of reliance on this variable.

A high p-value means low relevance of an explanatory variable to the dependent variable. Variable X_7 (fine aggregate) has a comparatively very high p-value and hence seemingly irrelevant.

Model 2 (Using Seven Explanatory Variables; dropping X_7)

Regression Statistics	
Multiple R	0.7836639
R Square	0.6141292
Adjusted R Square	0.6114862
Standard Error	10.412796
Observations	1030

ANOVA	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	7	176361.33	25194.476	232.36495	2.49E-206
Residual	1022	110811.7	108.42632		

Total	1029	287173.03			
-------	------	-----------	--	--	--

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	24.06271	8.8425583	2.7212385	0.0066143	6.7110656	41.414355	6.7110656	41.414355
Cement	0.1060683	0.0043655	24.29676	2.72E-103	0.0975019	0.1146348	0.0975019	0.1146348
Blast Furnace Slag	0.087453	0.0051978	16.825022	2.925E-56	0.0772534	0.0976525	0.0772534	0.0976525
Fly Ash	0.0693289	0.0077992	8.8891876	2.731E-18	0.0540246	0.0846333	0.0540246	0.0846333
Water	-0.2112656	0.023824	-8.8677524	3.266E-18	-0.2580152	-0.164516	-0.2580152	-0.164516
Superplasticizer	0.2627082	0.0923862	2.8435866	0.0045498	0.0814198	0.4439965	0.0814198	0.4439965
Coarse Aggregate	0.0033576	0.0052545	0.6389878	0.522974	-0.0069533	0.0136685	-0.0069533	0.0136685
Age	0.1133479	0.0054142	20.935399	2.873E-81	0.1027238	0.1239721	0.1027238	0.1239721

The corresponding regression equation is:

$$Y = 24.0627 + 0.1061X_1 + 0.08745X_2 + 0.06933X_3 - 0.2113X_4 + 0.2627X_5 + 0.003358X_6 + 0.1133X_8$$

The adjusted R^2 statistic of 0.61149 means that about 61% of variation in the dependent variable, i.e. concrete compressive strength, can be explained by these seven variables which suggest a very reasonable relationship of these explanatory variables to concrete compressive strength.

The T-statistic for variable X_8 (age) is the highest (20.9354) suggesting a great deal of reliance on this variable.

A high p-value means low relevance of an explanatory variable to the dependent variable. Variable X_6 (coarse aggregate) has a comparatively very high p-value and hence seemingly irrelevant.

Model 3 (Using Six Explanatory Variables; dropping X_6)

Regression Statistics	
Multiple R	0.7835656
R Square	0.613975
Adjusted R Square	0.6117109
Standard Error	10.409784
Observations	1030

ANOVA					
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	6	176317.06	29386.177	271.18124	1.78E-207
Residual	1023	110855.97	108.3636		
Total	1029	287173.03			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	29.0302 24	4.2124756	6.89148 77	9.639E- 12	20.7641 44	37.2963 04	20.76414 4	37.29630 4
Cement	0.10542 75	0.0042475	24.8207 84	7.82E-1 07	0.09709 26	0.11376 24	0.097092 6	0.113762 4
Blast Furnace Slag	0.08649 36	0.0049748	17.3863 16	1.61E-5 9	0.07673 16	0.09625 56	0.076731 6	0.096255 6
Fly Ash	0.06870 84	0.0077363	8.88129 01	2.913E- 18	0.05352 75	0.08388 92	0.053527 5	0.083889 2
Water	-0.21829 23	0.0211282	-10.331 809	7.193E- 24	-0.2597 519	-0.1768 328	-0.259751 9	-0.176832 8
Superplasticizer	0.23900 25	0.0845858	2.82556 38	0.00481 16	0.07302 11	0.40498 4	0.073021 1	0.404984
Age	0.11349 48	0.0054077	20.9874 98	1.298E- 81	0.10288 33	0.12410 63	0.102883 3	0.124106 3

The corresponding regression equation is:

$$Y = 29.0302 + 0.1054X_1 + 0.08649X_2 + 0.06871X_3 - 0.2183X_4 + 0.2390X_5 + 0.1135X_8$$

The adjusted R² of the full model is 61.40% demonstrating that the model a good indicator of concrete compressive strength based on those six variables. The p-values for the remaining explanatory variables are very close to zero, allowing us to reject the null hypothesis.

Conclusion

I began the regression analysis with full model of eight variables: cement, blast furnace slag, fly ash, water, superplasticizer, coarse aggregate, fine aggregate and age. I eliminated coarse aggregate, fine aggregate because they did not appear to be a good fit with the model based on the results of eight-variable regression analysis and six-variable regression analysis.

$$Y = 29.0302 + 0.1054X_1 + 0.08649X_2 + 0.06871X_3 - 0.2183X_4 + 0.2390X_5 + 0.1135X_6$$

Where: Y = Concrete compressive strength -- quantitative -- MPa

X₁ = Cement (component 1) -- quantitative – kg/m³

X₂ = Blast Furnace Slag (component 2) – quantitative – kg/m³

X_3 = Fly Ash (component 3) -- quantitative – kg/m³

X_4 = Water (component 4) -- quantitative – kg/m³

X_5 = Superplasticizer (component 5) – quantitative – kg/m³

X_6 = Age -- quantitative -- Day (1~365)

Based on three regression analysis performed, I determined that the equation with six explanatory variables provides the best fit. F-value also is the highest in this model reflecting the fact that the six-variable model is a best fit. (F-values: six variable model = 271.18124; seven variable model = 232.36495; eight variable model = 204.26914)