

SOC 3811/5811:
BASIC SOCIAL STATISTICS

Multiple Regression

Multiple Regression with k Independent Variables

We've seen how to estimate regression models that include two continuous predictor variables (X_1 and X_2) and a continuous response variable (Y)

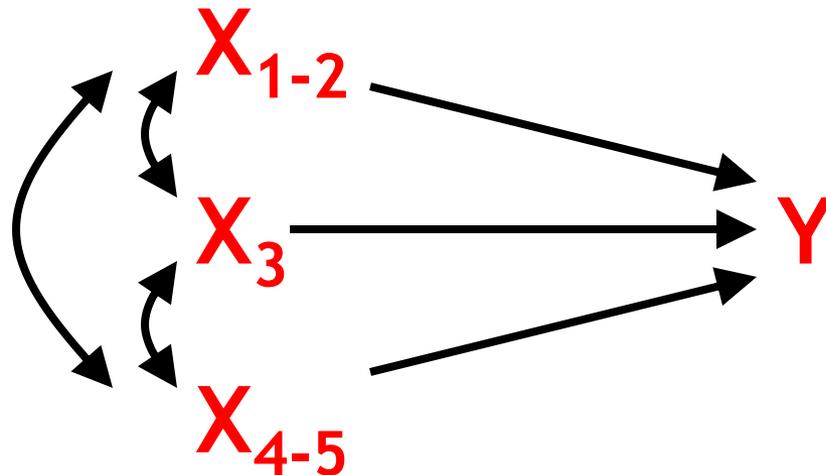
Extension #1: Models with k continuous predictor variables (X_1 through X_k) and a continuous response variable (Today)

Extension #2: Models with k predictor variables—some of which are continuous and some of which are discrete—and a continuous response variable (Bonus Material)

Extension #3: Models with k predictor variables and a discrete response variable (SOC 5811 and 8811)

Multiple Regression with k Independent Variables

Example: How do state educational resources (X_{1-2}), state education policies (X_3), and state economic conditions (X_{4-5}) influence states' high school graduation rates (Y)?



Multiple Regression with k Independent Variables

Example: How do state educational resources (X_{1-2}), state education policies (X_3), and state economic conditions (X_{4-5}) influence states' high school graduation rates (Y)?

(always start by looking at descriptive statistics)

Descriptive Statistics

	N	Mean	Variance
Y = Estimated Completion Rate, 2002	51	72.704	63.446
X1 = Per Pupil Expenditures, 2002	51	67.684	194.999
X2 = Pupil-Teacher Ratios, 2002	51	12.633	11.061
X3 = Carn. Units Required for Graduation, 2002	51	17.980	51.780
X4 = Poverty Rate, 2002	51	11.369	11.667
X5 = Unemployment Rate, 2002	51	4.522	.809
Valid N (listwise)	51		

D.C. is a state! 

Multiple Regression with k Independent Variables

Example: How do state educational resources (X_{1-2}), state education policies (X_3), and state economic conditions (X_{4-5}) influence states' high school graduation rates (Y)?

(always start by looking at descriptive statistics)

Correlations

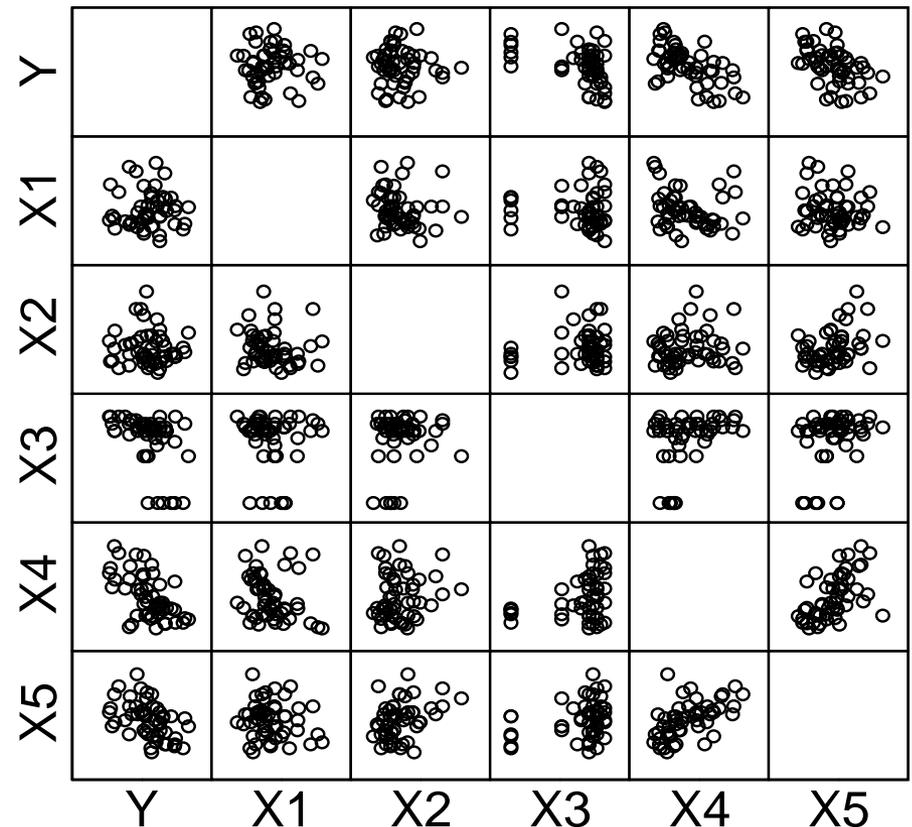
	Y	X1	X2	X3	X4	X5
Estimated Completion Rate, 2002	1.000					
Per Pupil Expenditures, 2002	-.004	1.000				
Pupil-Teacher Ratios, 2002	-.052	-.116	1.000			
Carn. Units Req. for Grad., 2002	-.436	-.057	.125	1.000		
Poverty Rate, 2002	-.574	-.196	.155	.327	1.000	
Unemployment Rate, 2002	-.475	-.121	.359	.289	.549	1.000

Multiple Regression with k Independent Variables

Example: How do state educational resources (X_{1-2}), state education policies (X_3), and state economic conditions (X_{4-5}) influence states' high school graduation rates (Y)?

Bivariate Scatterplots

(always start by looking at bivariate plots)



Multiple Regression with k Independent Variables

The population regression equation:

$$Y_i = \alpha + \sum_{j=1}^k \beta_j X_{ji} + \varepsilon_i$$

The population prediction equation:

$$\hat{Y}_i = \alpha + \sum_{j=1}^k \beta_j X_{ji}$$

The sample regression equation:

$$Y_i = a + \sum_{j=1}^k b_j X_{ji} + e_i$$

The sample prediction equation:

$$\hat{Y}_i = a + \sum_{j=1}^k b_j X_{ji}$$

Prediction Equations

Model with ONE Independent Variable

$$\hat{Y}_i = a + b_1 X_{1i}$$

Model with TWO Independent Variables

$$\hat{Y}_i = a + b_1 X_{1i} + b_2 X_{2i}$$

Model with k Independent Variables

$$\hat{Y}_i = a + b_1 X_{1i} + b_2 X_{2i} + b_3 X_{3i} + \dots + b_k X_{ki}$$

$$\hat{Y}_i = a + \sum_{j=1}^k b_j X_{ji}$$

Multiple Regression with k Independent Variables

The ordinary least squares (OLS) method is used to estimate a and b_1 through b_k ... again, this method minimizes the sum of the squared prediction errors

The computational formulas for a and b_1 through b_k are complex when $k > 2$

As before, they are based on the correlation among the response and predictor variables and on their means and variances

Interpreting Multiple Regression Coefficients

How are a and b_k interpreted in the equation:

$$\hat{Y}_i = a + \sum_{j=1}^k b_{kj} X_{ji}$$

Intercept a :

The intercept, a , equals the predicted value of Y when each of the k predictor variables (X_1 through X_k) equal 0

Multiple regression coefficient (or slope) b_k :

Multiple regression coefficient b_k represents the expected change in Y associated with a one unit increase in X_k , controlling for all other predictors in the model

Interpreting Multiple Regression Coefficients

Example: How do state educational resources (X_{1-2}), state education policies (X_3), and state economic conditions (X_{4-5}) influence states' high school graduation rates (Y)?

$$\begin{array}{ll} \hat{Y}_i = 98.728 & \hat{Y}_i = 98.728 \\ -0.064X_{1i} & -0.064(\text{Pupil - Teacher Ratio}_i) \\ +0.274X_{2i} & +0.274(\text{Per - Pupil Expenditure}_i) \\ -0.285X_{3i} & -0.285(\text{Carnegie Units}_i) \\ -0.933X_{4i} & -0.933(\text{Poverty Rate}_i) \\ -2.086X_{5i} & -2.086(\text{Unemployment Rate}_i) \end{array}$$

Coefficient of Determination

As before, we can use R^2 to express the proportion of variation in Y that is accounted for by the predictor variables

Because, at worst, a predictor variable can explain none of the variation in Y , it follows that the addition of more predictor variables to the model will either leave R^2 unchanged or increase it

Worksheet

Below are the means and standard deviations of scores on the second exam (“Exam”), the total points earned on problem sets (“ProblemSets”), in-class worksheet scores (“InClass”), and lab worksheet scores (“InLab”); the latter three are measured at the time the 2nd exam was taken

Variable	Obs	Mean	Std. Dev.	Min	Max
Exam	155	78.96129	11.871	0	90
ProblemSets	155	166.7452	37.3687	0	200
InClass	155	69.09677	11.69125	18	78
InLab	155	35.84516	6.896397	0	40

Worksheet

Below are the correlations between these variables

	Exam	Problems	InClass	InLab
Exam	1.0000			
ProblemSets	0.5028	1.0000		
InClass	0.2897	0.6574	1.0000	
InLab	0.3308	0.6246	0.5801	1.0000

Worksheet

Here are the results of a regression of exam score on the other variables. Interpret the intercept, the slopes, & R^2

Number of obs = 155

R-squared = 0.2572

Adj R-squared = 0.2425

Exam	Coef.	Std. Err.
ProblemSets	.1676782	.0322029
InClass	-.0886353	.098682
InLab	.0891642	.1614245
_cons	53.93007	5.308435

Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

Do the k predictor variables collectively explain any of the variation in Y ?

We use $R^2_{Y \bullet X_1 \dots X_k}$ to estimate $\rho^2_{Y \bullet X_1 \dots X_k}$

As before, another way to express $R^2_{Y \bullet X_1 \dots X_k}$ is:

$$R^2_{Y \bullet X_1 \dots X_k} = \frac{SS_{\text{REGRESSION}}}{SS_{\text{TOTAL}}}$$

where $SS_{\text{REGRESSION}} = SS_{\text{TOTAL}} - SS_{\text{ERROR}}$

Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

Hypothesis Testing in 6 Steps

1. State the null (H_0) and alternative (H_1) hypotheses
2. Check that the sample data conform to basic assumptions; if they do not, then do not go any further
3. Choose an α probability level ... that is, a probability associated with incorrectly rejecting the null hypothesis
4. Determine the “critical value” ... that is, how large the test statistic must be in order to reject the null hypothesis at the given α level
5. Calculate the test statistic ... F
6. Compare the test statistic to the critical value

Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

State the null (H_0) and alternative (H_1) hypotheses

$$H_0: \rho^2_{Y \bullet X_1 \dots X_k} = 0$$

$$H_1: \rho^2_{Y \bullet X_1 \dots X_k} > 0$$

This is a one-sided test (with no $<$) because $\rho^2_{Y \bullet X_1 \dots X_k}$ cannot possibly be less than zero

Failing to reject the null means failing to reject the hypothesis that the k predictor variables collectively explain none of the variation in Y

Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

Check that the sample data conform to basic assumptions;
if they do not, then do not go any further

The assumptions of the regression model described earlier
must hold for hypothesis tests about $\rho^2_{Y \bullet X_1 \dots X_k}$ to be valid

Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

Choose an α probability level ... that is, a probability associated with incorrectly rejecting the null hypothesis

Let's choose $\alpha=0.05$

Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

Determine the “critical value” ... that is, how large the test statistic must be in order to reject the null hypothesis at the given α level

The hypothesis test for $\rho^2_{Y \bullet X_1 \dots X_k}$ is an F test with $df_{\text{NUM}}=k$ and $df_{\text{DENOM}}=N-k-1$

In our example, we want $F_{5,45}$ for $\alpha=0.05$ which is close to 2.45

We will thus reject H_0 if our F statistic exceeds 2.45

Critical Values of F ($\alpha=0.05$)



DENOMINATOR Degrees of Freedom

		NUMERATOR Degrees of Freedom																	
		1	2	3	4	5	6	7	8	9	10	15	20	30	40	50	100	200	∞
1	161.45	199.50	215.71	224.58	230.16	233.99	236.77	238.88	240.54	241.88	245.95	248.01	250.10	251.14	251.77	253.04	253.68	254.31	
2	18.51	19.00	19.16	19.25	19.30	19.33	19.35	19.37	19.38	19.40	19.43	19.45	19.46	19.47	19.48	19.49	19.49	19.50	
3	10.13	9.55	9.28	9.12	9.01	8.94	8.89	8.85	8.81	8.79	8.70	8.66	8.62	8.59	8.58	8.55	8.54	8.53	
4	7.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04	6.00	5.96	5.86	5.80	5.75	5.72	5.70	5.66	5.65	5.63	
5	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82	4.77	4.74	4.62	4.56	4.50	4.46	4.44	4.41	4.39	4.36	
6	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15	4.10	4.06	3.94	3.87	3.81	3.77	3.75	3.71	3.69	3.67	
7	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73	3.68	3.64	3.51	3.44	3.38	3.34	3.32	3.27	3.25	3.23	
8	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44	3.39	3.35	3.22	3.15	3.08	3.04	3.02	2.97	2.95	2.93	
9	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23	3.18	3.14	3.01	2.94	2.86	2.83	2.80	2.76	2.73	2.71	
10	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07	3.02	2.98	2.85	2.77	2.70	2.66	2.64	2.59	2.56	2.54	
11	4.84	3.98	3.59	3.36	3.20	3.09	3.01	2.95	2.90	2.85	2.72	2.65	2.57	2.53	2.51	2.46	2.43	2.40	
12	4.75	3.89	3.49	3.26	3.11	3.00	2.91	2.85	2.80	2.75	2.62	2.54	2.47	2.43	2.40	2.35	2.32	2.30	
13	4.67	3.81	3.41	3.18	3.03	2.92	2.83	2.77	2.71	2.67	2.53	2.46	2.38	2.34	2.31	2.26	2.23	2.21	
14	4.60	3.74	3.34	3.11	2.96	2.85	2.76	2.70	2.65	2.60	2.46	2.39	2.31	2.27	2.24	2.19	2.16	2.13	
15	4.54	3.68	3.29	3.06	2.90	2.79	2.71	2.64	2.59	2.54	2.40	2.33	2.25	2.20	2.18	2.12	2.10	2.07	
16	4.49	3.63	3.24	3.01	2.85	2.74	2.66	2.59	2.54	2.49	2.35	2.28	2.19	2.15	2.12	2.07	2.04	2.01	
17	4.45	3.59	3.20	2.96	2.81	2.70	2.61	2.55	2.49	2.45	2.31	2.23	2.15	2.10	2.08	2.02	1.99	1.96	
18	4.41	3.55	3.16	2.93	2.77	2.66	2.58	2.51	2.46	2.41	2.27	2.19	2.11	2.06	2.04	1.98	1.95	1.92	
19	4.38	3.52	3.13	2.90	2.74	2.63	2.54	2.48	2.42	2.38	2.23	2.16	2.07	2.03	2.00	1.94	1.91	1.88	
20	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45	2.39	2.35	2.20	2.12	2.04	1.99	1.97	1.91	1.88	1.84	
21	4.32	3.47	3.07	2.84	2.68	2.57	2.49	2.42	2.37	2.32	2.18	2.10	2.01	1.96	1.94	1.88	1.84	1.81	
22	4.30	3.44	3.05	2.82	2.66	2.55	2.46	2.40	2.34	2.30	2.15	2.07	1.98	1.94	1.91	1.85	1.82	1.78	
23	4.28	3.42	3.03	2.80	2.64	2.53	2.44	2.37	2.32	2.27	2.13	2.05	1.96	1.91	1.88	1.82	1.79	1.76	
24	4.26	3.40	3.01	2.78	2.62	2.51	2.42	2.36	2.30	2.25	2.11	2.03	1.94	1.89	1.86	1.80	1.77	1.73	
25	4.24	3.39	2.99	2.76	2.60	2.49	2.40	2.34	2.28	2.24	2.09	2.01	1.92	1.87	1.84	1.78	1.75	1.71	
26	4.23	3.37	2.98	2.74	2.59	2.47	2.39	2.32	2.27	2.22	2.07	1.99	1.90	1.85	1.82	1.76	1.73	1.69	
27	4.21	3.35	2.96	2.73	2.57	2.46	2.37	2.31	2.25	2.20	2.06	1.97	1.88	1.84	1.81	1.74	1.71	1.67	
28	4.20	3.34	2.95	2.71	2.56	2.45	2.36	2.29	2.24	2.19	2.04	1.96	1.87	1.82	1.79	1.73	1.69	1.65	
29	4.18	3.33	2.93	2.70	2.55	2.43	2.35	2.28	2.22	2.18	2.03	1.94	1.85	1.81	1.77	1.71	1.67	1.64	
30	4.17	3.32	2.92	2.69	2.53	2.42	2.33	2.27	2.21	2.16	2.01	1.93	1.84	1.79	1.76	1.70	1.66	1.62	
31	4.16	3.30	2.91	2.68	2.52	2.41	2.32	2.25	2.20	2.15	2.00	1.92	1.83	1.78	1.75	1.68	1.65	1.61	
32	4.15	3.29	2.90	2.67	2.51	2.40	2.31	2.24	2.19	2.14	1.99	1.91	1.82	1.77	1.74	1.67	1.63	1.59	
33	4.14	3.28	2.89	2.66	2.50	2.39	2.30	2.23	2.18	2.13	1.98	1.90	1.81	1.76	1.72	1.66	1.62	1.58	
34	4.13	3.28	2.88	2.65	2.49	2.38	2.29	2.23	2.17	2.12	1.97	1.89	1.80	1.75	1.71	1.65	1.61	1.57	
35	4.12	3.27	2.87	2.64	2.49	2.37	2.29	2.22	2.16	2.11	1.96	1.88	1.79	1.74	1.70	1.63	1.60	1.56	
40	4.08	3.23	2.84	2.61	2.45	2.34	2.25	2.18	2.12	2.08	1.92	1.84	1.74	1.69	1.66	1.59	1.55	1.51	
50	4.03	3.18	2.79	2.56	2.40	2.29	2.20	2.13	2.07	2.03	1.87	1.78	1.69	1.63	1.60	1.52	1.48	1.44	
75	3.97	3.12	2.73	2.49	2.34	2.22	2.13	2.06	2.01	1.96	1.80	1.71	1.61	1.55	1.52	1.44	1.39	1.34	
100	3.94	3.09	2.70	2.46	2.31	2.19	2.10	2.03	1.97	1.93	1.77	1.68	1.57	1.52	1.48	1.39	1.34	1.28	
200	3.89	3.04	2.65	2.42	2.26	2.14	2.06	1.98	1.93	1.88	1.72	1.62	1.52	1.46	1.41	1.32	1.26	1.19	
∞	3.84	3.00	2.60	2.37	2.21	2.10	2.01	1.94	1.88	1.83	1.67	1.57	1.46	1.39	1.35	1.24	1.17	1.00	



Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

Calculate the test statistic

The F statistic when there are k predictors is

$$F_{k, N-k-1} = \frac{SS_{\text{REGRESSION}}/k}{SS_{\text{ERROR}}/N-k-1} = \frac{MS_{\text{REGRESSION}}}{MS_{\text{ERROR}}}$$

Computationally:

$$SS_{\text{TOTAL}} = (s_Y^2)(N-1)$$

$$SS_{\text{REGRESSION}} = (R^2_{Y \bullet X_1 \dots X_k})(SS_{\text{TOTAL}})$$

$$SS_{\text{ERROR}} = SS_{\text{TOTAL}} - SS_{\text{REGRESSION}}$$

Testing Hypotheses about $\rho^2_{Y \bullet X_1 \dots X_k}$

Calculate the test statistic

In our example:

$$SS_{\text{TOTAL}} = (s_Y^2)(N-1) = (63.446)(51-1) = 3,172.3$$

$$SS_{\text{REGRESSION}} = (R^2_{Y \bullet X_1 \dots X_k})(SS_{\text{TOTAL}}) = (0.448)(3,172.3) = 1,421.19$$

$$SS_{\text{ERROR}} = SS_{\text{TOTAL}} - SS_{\text{REGRESSION}} = 3,172.3 - 1,421.19 = 1,751.11$$

SO

$$F_{5,45} = \frac{SS_{\text{REGRESSION}}/5}{SS_{\text{ERROR}}/45} = \frac{1,421.19/5}{1,751.11/45} = 7.30$$

Testing Hypotheses about $\rho^2_{Y \cdot X_1 \dots X_k}$

Compare the test statistic to the critical value

If the test statistic is as large or larger than the critical value, then reject H_0

If the test statistic is less than the critical value, then do not reject H_0

We can restate the hypotheses:

$H_0: \rho^2_{Y \cdot X_1 \dots X_k} = 0 \rightarrow$ Fail to reject H_0 if $F \leq 2.45$

$H_1: \rho^2_{Y \cdot X_1 \dots X_k} > 0 \rightarrow$ Reject H_0 if $F > 2.45$

Since $F=7.30$, we reject H_0 ... so it appears that in the population the k predictors (in combination) account for some of the variability in Y

Worksheet

Here are the results of a regression of exam score on the other variables. Test the hypothesis that $\rho^2=0$; use $\alpha=0.05$

Number of obs = 155

R-squared = 0.2572

Adj R-squared = 0.2425

Exam	Coef.
ProblemSets	.1676782
InClass	-.0886353
InLab	.0891642
_cons	53.93007

Variable	Obs	Mean	Std. Dev.
Exam	155	78.96129	11.871
ProblemSets	155	166.7452	37.3687
InClass	155	69.09677	11.69125
InLab	155	35.84516	6.896397

Testing Hypotheses about β_k

Can we conclude that β_k is different from 0?

We use b_k to estimate β_k

In the model with $k=2$ predictors the variance of the sampling distribution of slopes b_1 and b_2 were

$$s_{b_1}^2 = \frac{MS_{\text{ERROR}}}{(s_{X_1}^2)(N-1)(1-R_{X_1 \cdot X_2}^2)} \quad s_{b_2}^2 = \frac{MS_{\text{ERROR}}}{(s_{X_2}^2)(N-1)(1-R_{X_2 \cdot X_1}^2)}$$

In the model with k predictor variables the variances of the sampling distributions of b_k is

$$s_{b_k}^2 = \frac{MS_{\text{ERROR}}}{(s_{X_k}^2)(N-1)(1-R_{X_k \cdot X_1 \dots X_{k-1}}^2)}$$

Testing Hypotheses about β_k

Hypothesis Testing in 6 Steps

1. State the null (H_0) and alternative (H_1) hypotheses
2. Check that the sample data conform to basic assumptions; if they do not, then do not go any further
3. Choose an α probability level ... that is, a probability associated with incorrectly rejecting the null hypothesis
4. Determine the “critical value” ... that is, how large the test statistic must be in order to reject the null hypothesis at the given α level
5. Calculate the test statistic ... t
6. Compare the test statistic to the critical value

Testing Hypotheses about β_k

State the null (H_0) and alternative (H_1) hypotheses

$$H_0: \beta_k = 0$$

$$H_1: \beta_k \neq 0$$

This is a two-sided tests (although it needn't be)

Failing to reject H_0 means failing to reject the hypothesis that there is no net association between Y and X_k (holding constant the other predictors in the model)

Testing Hypotheses about β_k

Check that the sample data conform to basic assumptions; if they do not, then do not go any further

The assumptions of the regression model described earlier must hold for hypothesis tests about β_k to be valid

Testing Hypotheses about β_k

Choose an α probability level ... that is, a probability associated with incorrectly rejecting the null hypothesis

Let's choose $\alpha=0.05$

Testing Hypotheses about β_k

Determine the “critical value” ... that is, how large the test statistic must be in order to reject the null hypothesis at the given α level

The hypothesis test for β_k is a t tests with $N-k-1$ degrees of freedom (because MS_{ERROR} has $N-k-1$ degrees of freedom)

In our example, we want t_{45} for $\alpha=0.05$ which is close to 2.021 (because $N-k-1$ is 45 and thus close to 40)

For hypothesis tests about β_k we will thus reject H_0 if our t statistic exceeds 2.021 in absolute value

Table entry for p and C is the critical value t^* with probability p lying to its right and probability C lying between $-t^*$ and t^* .

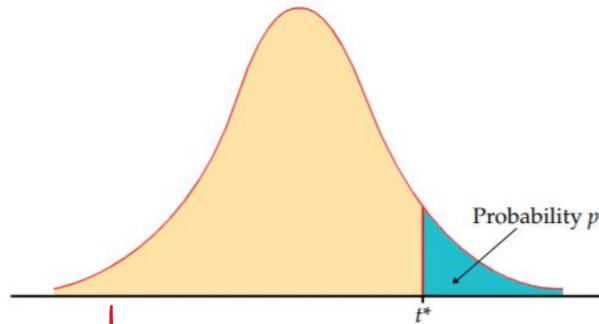


TABLE D

t distribution critical values

df	One-tail probability p											
	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001	.0005
1	1.000	1.376	1.963	3.078	6.314	12.71	15.89	31.82	63.66	127.3	318.3	636.6
2	0.816	1.061	1.386	1.886	2.920	4.303	4.849	6.965	9.925	14.09	22.33	31.60
3	0.765	0.978	1.250	1.638	2.353	3.182	3.482	4.541	5.841	7.453	10.21	12.92
4	0.741	0.941	1.190	1.533	2.132	2.776	2.999	3.747	4.604	5.598	7.173	8.610
5	0.727	0.920	1.156	1.476	2.015	2.571	2.757	3.365	4.032	4.773	5.893	6.869
6	0.718	0.906	1.134	1.440	1.943	2.447	2.612	3.143	3.707	4.317	5.208	5.959
7	0.711	0.896	1.119	1.415	1.895	2.365	2.517	2.998	3.499	4.029	4.785	5.408
8	0.706	0.889	1.108	1.397	1.860	2.306	2.449	2.896	3.355	3.833	4.501	5.041
9	0.703	0.883	1.100	1.383	1.833	2.262	2.398	2.821	3.250	3.690	4.297	4.781
10	0.700	0.879	1.093	1.372	1.812	2.228	2.359	2.764	3.169	3.581	4.144	4.587
11	0.697	0.876	1.088	1.363	1.796	2.201	2.328	2.718	3.106	3.497	4.025	4.437
12	0.695	0.873	1.083	1.356	1.782	2.179	2.303	2.681	3.055	3.428	3.930	4.318
13	0.694	0.870	1.079	1.350	1.771	2.160	2.282	2.650	3.012	3.372	3.852	4.221
14	0.692	0.868	1.076	1.345	1.761	2.145	2.264	2.624	2.977	3.326	3.787	4.140
15	0.691	0.866	1.074	1.341	1.753	2.131	2.249	2.602	2.947	3.286	3.733	4.073
16	0.690	0.865	1.071	1.337	1.746	2.120	2.235	2.583	2.921	3.252	3.686	4.015
17	0.689	0.863	1.069	1.333	1.740	2.110	2.224	2.567	2.898	3.222	3.646	3.965
18	0.688	0.862	1.067	1.330	1.734	2.101	2.214	2.552	2.878	3.197	3.611	3.922
19	0.688	0.861	1.066	1.328	1.729	2.093	2.205	2.539	2.861	3.174	3.579	3.883
20	0.687	0.860	1.064	1.325	1.725	2.086	2.197	2.528	2.845	3.153	3.552	3.850
21	0.686	0.859	1.063	1.323	1.721	2.080	2.189	2.518	2.831	3.135	3.527	3.819
22	0.686	0.858	1.061	1.321	1.717	2.074	2.183	2.508	2.819	3.119	3.505	3.792
23	0.685	0.858	1.060	1.319	1.714	2.069	2.177	2.500	2.807	3.104	3.485	3.768
24	0.685	0.857	1.059	1.318	1.711	2.064	2.172	2.492	2.797	3.091	3.467	3.745
25	0.684	0.856	1.058	1.316	1.708	2.060	2.167	2.485	2.787	3.078	3.450	3.725
26	0.684	0.856	1.058	1.315	1.706	2.056	2.162	2.479	2.779	3.067	3.435	3.707
27	0.684	0.855	1.057	1.314	1.703	2.052	2.158	2.473	2.771	3.057	3.421	3.690
28	0.683	0.855	1.056	1.313	1.701	2.048	2.154	2.467	2.763	3.047	3.408	3.674
29	0.683	0.854	1.055	1.311	1.699	2.045	2.150	2.462	2.756	3.038	3.396	3.659
30	0.683	0.854	1.055	1.310	1.697	2.042	2.147	2.457	2.750	3.030	3.385	3.646
40	0.681	0.851	1.050	1.303	1.684	2.021	2.123	2.423	2.704	2.971	3.307	3.551
50	0.679	0.849	1.047	1.299	1.676	2.009	2.109	2.403	2.678	2.937	3.261	3.496
60	0.679	0.848	1.045	1.296	1.671	2.000	2.099	2.390	2.660	2.915	3.232	3.460
80	0.678	0.846	1.043	1.292	1.664	1.990	2.088	2.374	2.639	2.887	3.195	3.416
100	0.677	0.845	1.042	1.290	1.660	1.984	2.081	2.364	2.626	2.871	3.174	3.390
1000	0.675	0.842	1.037	1.282	1.646	1.962	2.056	2.330	2.581	2.813	3.098	3.300
z^*	0.674	0.841	1.036	1.282	1.645	1.960	2.054	2.326	2.576	2.807	3.091	3.291

Testing Hypotheses about β_k

Calculate the test statistic

The t statistic for β_k is

$$t_{N-k-1} = \frac{b_k - 0}{s_{b_1}} = \frac{b_k - 0}{\sqrt{\frac{MS_{\text{ERROR}}}{(s_{X_k}^2)(N-1)(1-R_{X_k \cdot X_1 \dots X_{k-1}}^2)}}$$

Testing Hypotheses about β_k

Calculate the test statistic

In our example...

Coefficients^a

	Unstandardized Coefficients		Std. Coef.	t	Sig.
	B	Std. Error	Beta		
(Constant)	98.73	7.21		13.7	.00
Per Pupil Expenditures, 2002	-.06	.06	-.11	-.99	.33
Pupil-Teacher Ratios, 2002	.27	.29	.11	.96	.34
Carn. Units Req. for Grad., 2002	-.29	.13	-.26	-2.17	.03
Poverty Rate, 2002	-.93	.32	-.40	-2.91	.01
Unemployment Rate, 2002	-2.09	1.25	-.24	-1.66	.10

a. Dependent Variable: Estimated Completion Rate, 2002

Testing Hypotheses about β_k

Compare the test statistic to the critical value

If the test statistic is as large or larger than the critical value, then reject H_0

If the test statistic is less than the critical value, then do not reject H_0

We can restate the hypotheses:

$$H_0: \beta_k = 0$$

$$H_1: \beta_k \neq 0$$

Our values of t exceed our critical value t^* (2.021) for two of the five hypothesis tests about b_k , so we reject the null hypothesis in those instances

Worksheet

Here are the results of the regression of exam score on the other variables. Test the hypothesis that $\beta_{\text{ProblemSets}}=0$; use $\alpha=0.05$

Source	SS	df	MS
Model	5582.60155	3	1860.86718
Residual	16119.1662	151	106.749445
Total	21701.7677	154	140.92057

Number of obs = 155
F(3, 151) = 17.43
Prob > F = 0.0000
R-squared = 0.2572
Adj R-squared = 0.2425
Root MSE = 10.332

Exam	Coef.	Std. Err.
ProblemSets	.1676782	.0322029
InClass	-.0886353	.098682
InLab	.0891642	.1614245
_cons	53.93007	5.308435

BONUS TOPIC #1

Comparing Nested Equations

We can use hypothesis tests about specific slopes, b_k , to assess whether particular X variables add to the predictive power of the regression model

If we reject H_0 for β_k , we are concluding that X_k is significantly associated with Y net of the other covariates in the model ... and thus that our ability to predict Y is improved by including X_k in the model

Sometimes, however, we are interested in assessing the contribution of theoretically derived groups of X variables to our ability to predict Y

Comparing Nested Equations

In the example we've been considering, there are three groups of predictors of states' graduation rates (Y):

- state educational resources (X_1 and X_2)

- state education policies (X_3)

- state economic conditions (X_4 and X_5)

For example, do state educational resources significantly add to the predictive power of the model (as compared to a model that does not include these predictors)?

Comparing Nested Equations

The nested equations below allow us to making useful statements about the contributions of groups of variables to our ability to predict state graduation rates

	Model 1	Model 2	Model 3	Model 4
Pup.-T Ratio	-0.01	—	—	-0.06
P.P. Exp.	-0.13	—	—	0.27
Carn. Units	—	-0.48**	—	-0.29*
Poverty	—	—	-1.04**	-0.93**
Unempl.	—	—	-2.03	-2.09
Constant	74.68**	81.38**	93.75**	98.73**
R ²	0.003	0.190	0.366	0.448
F (df ₁ ,df ₂)	0.066 (2,48)	11.478**(1,49)	13.831**(2,48)	7.301**(5,45)

* = $p < 0.05$; ** = $p < 0.01$

Comparing Nested Equations

In this example we can say that Models 1, 2, and 3 are nested within Model 4 ... they contain a subset of the predictors included in Model 4

	Model 1	Model 2	Model 3	Model 4
Pup.-T Ratio	-0.01	—	—	-0.06
P.P. Exp.	-0.13	—	—	0.27
Carn. Units	—	-0.48**	—	-0.29*
Poverty	—	—	-1.04**	-0.93**
Unempl.	—	—	-2.03	-2.09
Constant	74.68**	81.38**	93.75**	98.73**
R ²	0.003	0.190	0.366	0.448
F (df ₁ ,df ₂)	0.066 (2,48)	11.478**(1,49)	13.831**(2,48)	7.301**(5,45)

* = $p < 0.05$; ** = $p < 0.01$

Comparing Nested Equations

The F tests for each individual model indicate whether the X variables in that model improve our ability to predict Y relative to the null model (with no X variables)

	Model 1	Model 2	Model 3	Model 4
Pup.-T Ratio	-0.01	—	—	-0.06
P.P. Exp.	-0.13	—	—	0.27
Carn. Units	—	-0.48**	—	-0.29*
Poverty	—	—	-1.04**	-0.93**
Unempl.	—	—	-2.03	-2.09
Constant	74.68**	81.38**	93.75**	98.73**
R ²	0.003	0.190	0.366	0.448
F (df ₁ ,df ₂)	0.066 (2,48)	11.478**(1,49)	13.831**(2,48)	7.301**(5,45)

* = $p < 0.05$; ** = $p < 0.01$

Comparing Nested Equations

A different question is whether a particular subset of X variables adds to our ability to predict Y relative to a model that contains a different subset of X variables

	Model 1	Model 2	Model 3	Model 4
Pup.-T Ratio	-0.01	—	—	-0.06
P.P. Exp.	-0.13	—	—	0.27
Carn. Units	—	-0.48**	—	-0.29*
Poverty	—	—	-1.04**	-0.93**
Unempl.	—	—	-2.03	-2.09
Constant	74.68**	81.38**	93.75**	98.73**
R ²	0.003	0.190	0.366	0.448
F (df ₁ ,df ₂)	0.066 (2,48)	11.478**(1,49)	13.831**(2,48)	7.301**(5,45)

* = p < 0.05 ; ** = p < 0.01

Comparing Nested Equations

Call the model that contains the full set of X variables the complete model; it has k_2 independent variables

Call a model that contains a subset of those X variables the reduced model; it contains k_1 independent variables

Question: Does the addition of the $k_2 - k_1$ new predictor variables in the complete model improve our ability to predict Y (relative to the reduced model)?

(If $k_2 - k_1$ equals one, then we are just adding one new X variable and we can answer this question with a hypothesis test about the coefficient for that variable)

Comparing Nested Equations

In general, when the reduced model is nested within the complete model we test the hypothesis that the new additional variables in the complete model add to the predictive power of the null model

We make this comparison using an F statistic that is based on the change in R^2 between the reduced model (R^2_1) and the complete model (R^2_2)

$$F_{(k_2 - k_1), (N - k_2 - 1)} = \frac{(R^2_2 - R^2_1) / (k_2 - k_1)}{(1 - R^2_2) / (N - k_2 - 1)}$$

Comparing Nested Equations

In our example: Does the addition of the state educational resource variables and the state education policies variables (in Model 4) add to the predictive power of Model 3?

If we use $\alpha=0.05$, the critical value F^* has $df_{\text{NUM}}=k_2-k_1=5-2=3$ and $df_{\text{DENOM}}=N-k_2-1=51-5-1=45$, so $F^*=2.84$

The F statistic equals $F_{3,45} = \frac{(0.448 - 0.366)/(5 - 2)}{(1 - 0.448)/(51 - 5 - 1)} = 2.228$

We fail to reject H_0 ... the addition of these 3 variables does not add to the predictive power of Model 3

Comparing Nested Equations

We thus conclude that Model 4 fits no better than Model 3, even though we observe a statistically significant coefficient (for “Carnegie Units”) in Models 2 and 4

	Model 1	Model 2	Model 3	Model 4
Pup.-T Ratio	-0.01	—	—	-0.06
P.P. Exp.	-0.13	—	—	0.27
Carn. Units	—	-0.48**	—	-0.29*
Poverty	—	—	-1.04**	-0.93**
Unempl.	—	—	-2.03	-2.09
Constant	74.68**	81.38**	93.75**	98.73**
R ²	0.003	0.190	0.366	0.448
F (df ₁ ,df ₂)	0.066 (2,48)	11.478**(1,49)	13.831**(2,48)	7.301**(5,45)

* = $p < 0.05$; ** = $p < 0.01$

BONUS TOPIC #2

Regression vs. ANOVA

The regression techniques we have used thus far require continuous predictor variables

It would be wrong—technically and conceptually—to simply enter nominal or ordinal variables as predictor variables, since it is wrong to compute means and standard deviations for these variables

For a discrete variable X that has j categories, we can construct j dummy variables—each of which has possible values 0 and 1—and each of which indicates whether an individual falls into a particular category of X

Regression vs. ANOVA

For example, X might indicate father's education, which is a discrete measure of whether fathers (a) did not finish high school, (b) finished high school but went no further, or (c) completed at least some college

From this we can construct $j=3$ three dummy variables:

	X_1	X_2	X_3
Father: < H.S.	1	0	0
Father: = H.S.	0	1	0
Father: > H.S.	0	0	1

Notice that knowing the value of $j-1$ of the X_j values allows you to infer the value of the j^{th} X value

Regression vs. ANOVA

If we then regress child's education (Y) on j—1 of these dummy variables we observe:

$$\hat{Y}_i = 11.776 + 1.697X_2 + 3.098X_3$$

Compare these results to the means of Y by level of father's education:

Mean of Child's Education (Y)

Father: < H.S.	11.776	$\hat{Y}_i = 11.776 + 1.697(0) + 3.098(0)$
Father: = H.S.	13.473	$\hat{Y}_i = 11.776 + 1.697(1) + 3.098(0)$
Father: > H.S.	14.871	$\hat{Y}_i = 11.776 + 1.697(0) + 3.098(1)$

Regression vs. ANOVA

If we then regress child's education (Y) on $j-1$ of these dummy variables we observe: $\hat{Y}_i = 11.776 + 1.697X_2 + 3.098X_3$

This regression model is exactly equivalent to an ANOVA in which we investigate the association between discrete variable X and continuous variable Y

The F statistic for this regression model is identical to the F statistic for the ANOVA relating Y to X

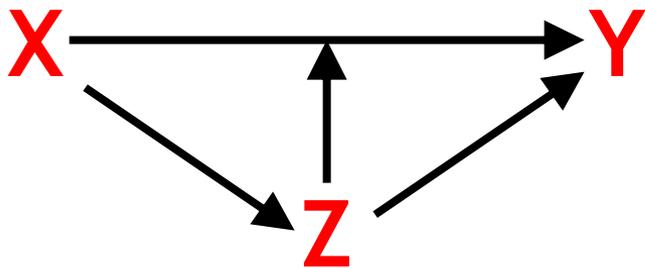
However, the regression framework gives us the ability to control for additional independent variables ... Doing so is known as ANalysis of COVariance (ANCOVA)

BONUS TOPIC #3

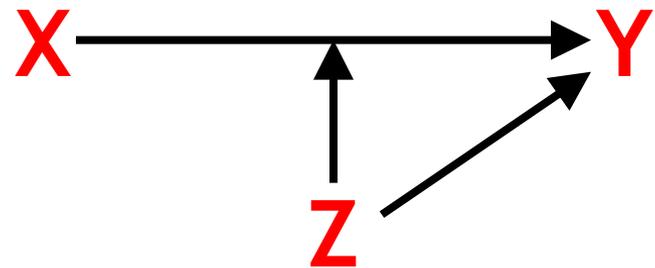
Interaction Effects

In some cases the association between X and Y may be different across levels of Z (or “conditioned by” Z)

In these cases we say that there is an **interaction** between X and Z ... Z is known as a **moderating** variable



or

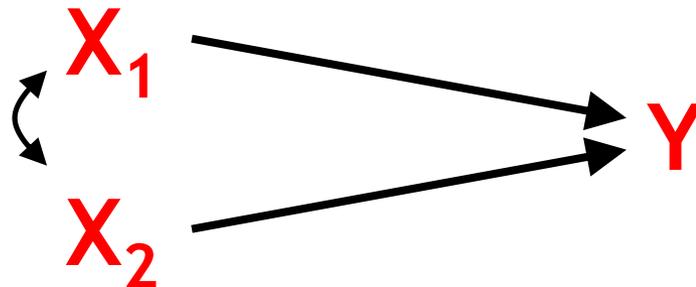


Interaction Effects

How do we model this in the regression context?

Imagine that we regress a continuous measure of education (Y) on a continuous measure of father's education (X_1) and a dummy variable for gender (X_2) that equals 0 for women and 1 for men

That model looks like this:

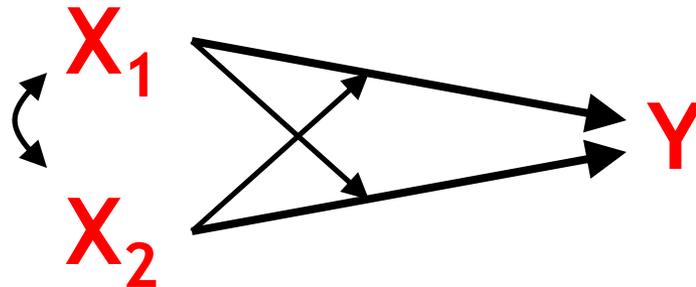


Interaction Effects

To allow for an “interaction effect” — such that the effect of X_1 on Y varies across levels of X_2 and the effect of X_2 on Y varies across levels of X_1 — we add an interaction term

An interaction term between X_1 and X_2 is simply a new variable created by multiplying X_1 by X_2

If we add the interaction term to the model that includes X_1 and X_2 as predictors, we have:



Interaction Effects

The prediction equation for this new model is

$$\hat{Y}_i = a + b_1 \text{Father's Educ.}_i + b_2 \text{Male}_i + b_3 \text{Inter.}_i$$

where “Inter” is the variable that equals $X_1 \times X_2$

If we estimate this model we get

$$\hat{Y}_i = 9.37 + (0.34)\text{Father's Educ.}_i + (0.05)\text{Male}_i + (0.02)\text{Inter.}_i$$

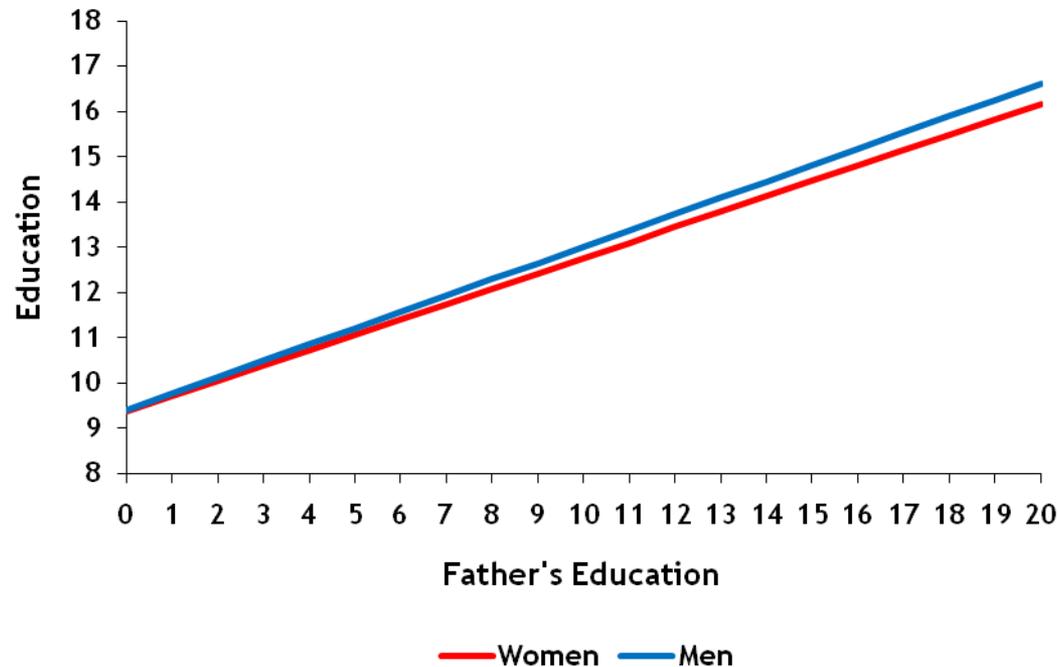
What is the “effect” of a one unit increase in father’s education? What is the “effect” of being male?

What is the predicted value of Y for a man whose father completed 10 years of school? What about for a woman whose father completed 10 years of school?

Interaction Effects

It is almost always easiest to think about interaction effects if you make a graph of predicted values

Below are predicted values of education (Y) by father's education (X_1) and gender (X_2)



Interaction Effects

How do we know whether the interaction term improves the power of the model to predict Y ? We may theorize that an interaction effect exists ... but how do we test the hypothesis that this is true in the data?

Option #1: Look at the statistical significance of the coefficient for the interaction term

Option #2: Treat the model w/o the interaction term as a reduced model that is nested within a full model that does include the interaction term ... conduct an F test .

Interaction Effects

How do we know whether the interaction term improves the power of the model to predict Y ? We may theorize that an interaction effect exists ... but how do we test the hypothesis that this is true in the data?

Option #1: In our example, the test statistic t for the interaction term is 2.48 ... we would reject H_0 at $\alpha=0.05$

Option #2: In our example, the test statistic $F_{1,27356}$ for the improvement in fit of the full model relative to the reduced model is 6.15 ... we would reject H_0 at $\alpha=0.05$

Want More?

David Lane's Books

http://onlinestatbook.com/2/regression/multiple_regression.html

Dallal's Book (see "Simple Linear Regression" section)

<http://www.jerrydallal.com/LHSP/LHSP.htm>

(Look under "multiple linear regression")

Biddle's Book:

http://www.biddle.com/documents/bcg_comp_chapter4.pdf

Another good overview:

<http://www.amstat.org/publications/jse/v16n3/datasets.kuiper.html>