



Chapter 11

Multiple Regression

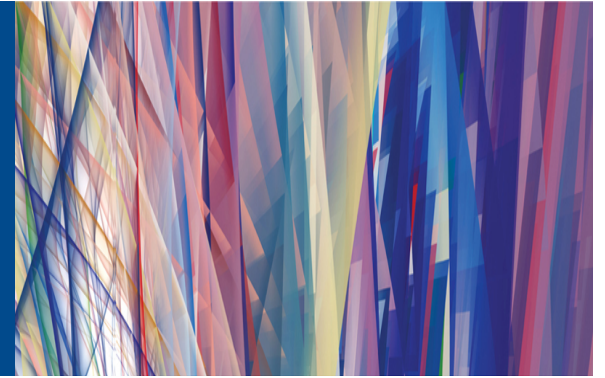
Introduction to the Practice of
STATISTICS SEVENTH
EDITION

Moore / McCabe / Craig

Lecture Presentation Slides

Chapter 11

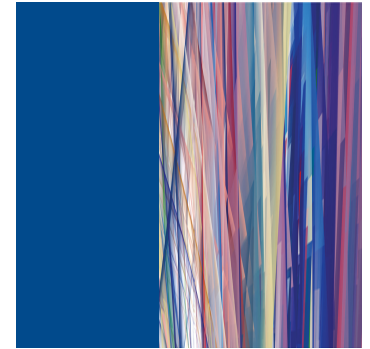
Multiple Regression



11.1 Inference for Multiple Regression

11.2 A Case Study

11.1 Inference for Multiple Regression



- Population Multiple Regression Model
- Data for Multiple Regression
- Multiple Linear Regression Model
- Confidence Intervals and Significance Tests
- Squared Multiple Correlation R^2

Population Multiple Regression Equation



Up to this point we have considered, in detail, the linear regression model in one explanatory variable x .

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Usually more complex linear models are needed in practical situations.

There are many problems in which a knowledge of more than one explanatory variable is necessary in order to obtain a better understanding and better prediction of a particular response.

In multiple regression, the response variable y depends on p explanatory variables, x_1, x_2, \dots, x_p

$$\mu_y = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

Data for Multiple Regression



The data for a simple linear regression problem consist of n observations (x_i, y_i) of the two variables.

Data for multiple linear regression consist of the value of a response variable y and p explanatory variables (x_1, x_2, \dots, x_p) on n cases.

We write the data and enter them into software in the form:

Case	Variables				
	x_1	x_2	...	x_p	y
1	x_{11}	x_{12}	...	x_{1p}	y_1
2	x_{21}	x_{22}	...	x_{2p}	y_2
...
n	x_{n1}	x_{n2}	...	x_{np}	y_n

Multiple Linear Regression Model



The **statistical model for multiple linear regression** is

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \varepsilon_i$$

for $i = 1, 2, \dots, n$.

The **mean response** μ_y is a linear function of the explanatory variables:

$$\mu_y = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip}$$

The **deviations** ε_i are independent and Normally distributed $N(0, \sigma)$.

The parameters of the model are $\beta_0, \beta_1 \dots \beta_p$, and s .

Estimation of the Parameters



Select a random sample of n individuals for which $p + 1$ variables are measured ($x_1 \dots, x_p, y$). The least-squares regression method minimizes the sum of squared deviations $e_i = (y_i - \hat{y}_i)$ to express y as a linear function of the p explanatory variables:

$$\hat{y}_i = b_0 + b_1 x_{i1} + \dots + b_p x_{ip}$$

As with simple linear regression, the constant b_0 is the y -intercept.

- The regression coefficients (b_1, \dots, b_p) reflect the unique association of each independent variable with the y variable. They are analogous to the slope in simple regression.
- The parameter s^2 measures the variability of the responses about the population regression equation. The estimator is:

$$s^2 = \frac{\sum e_i^2}{n - p - 1} = \frac{\sum (y_i - \hat{y}_i)^2}{n - p - 1}$$

Confidence Interval for β_j



Estimating the regression parameters $\beta_0, \dots, \beta_j, \dots, \beta_p$ is a case of one-sample inference with unknown population variance.

We rely on the t distribution, with **$n - p - 1$ degrees of freedom**.

A **level C confidence interval for β_j** is:

$$b_j \pm t^* SE_{b_j}$$

where SE_{b_j} is the standard error of b_j and t^* is the t critical for the $t(n - p - 1)$ distribution with area C between $-t^*$ and $+t^*$.

Significance Test for β_j



To test the hypothesis $H_0: \beta_j = 0$ versus a one- or two-sided alternative,

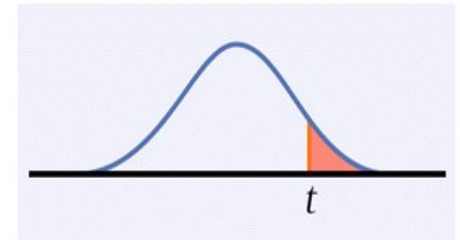
we calculate the t statistic **$t = b_j / \text{SE}_{b_j}$**

which has the $t (n - p - 1)$ distribution to find the p -value of the test.

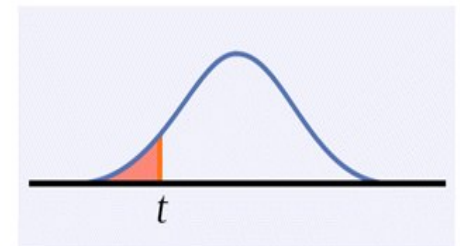
Note: Software typically provides two-sided p -values.

Important: this tests the significance of ONE variable AFTER adjusting for the effects of all others!!!

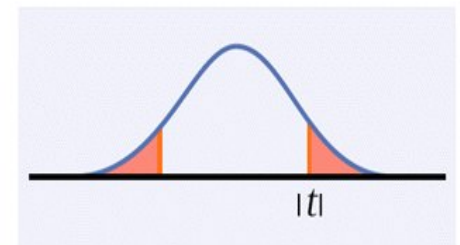
$$H_a: \beta_j > 0 \text{ is } P(T \geq t)$$



$$H_a: \beta_j < 0 \text{ is } P(T \leq t)$$



$$H_a: \beta_j \neq 0 \text{ is } 2P(T \geq |t|)$$



ANOVA F -test for Multiple Regression

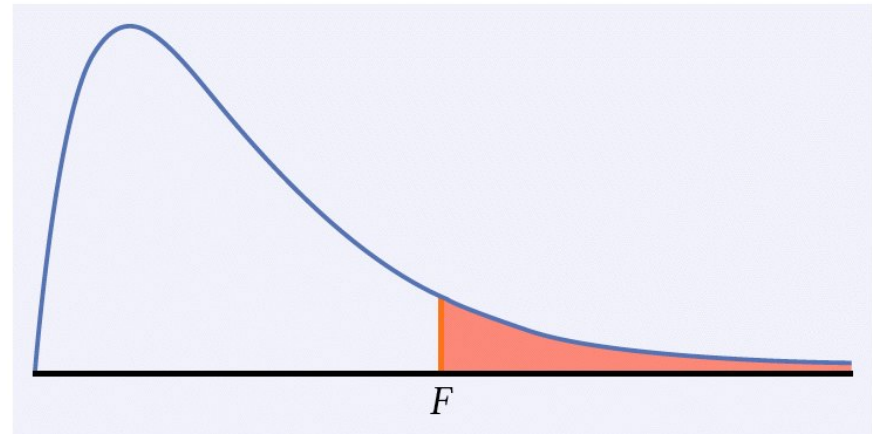


For a multiple linear relationship the ANOVA tests the hypotheses

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \quad \text{versus} \quad H_A: \text{at least one } \beta \neq 0$$

by computing the F statistic: **$F = \text{MSM} / \text{MSE}$**

When H_0 is true, F follows the $F(p, n - p - 1)$ distribution. The p -value is $P(F \geq f)$.



A significant p -value doesn't mean that all p explanatory variables have a significant influence on y —only that at least one does.

ANOVA Table for Multiple Regression



Source	Sum of squares SS	df	Mean square MS	F	P -value
Model	$\sum(\hat{y}_i - \bar{y})^2$	p	SSM/DFM	MSM/MSE	Tail area above F
Error	$\sum(y_i - \hat{y}_i)^2$	$n - p - 1$	SSE/DFE		
Total	$\sum(y_i - \bar{y})^2$	$n - 1$			

$$SST = SSM + SSE$$

$$DFT = DFM + DFE$$

Squared Multiple Correlation R^2



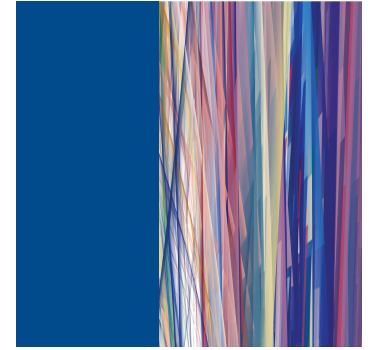
Just as with simple linear regression, **R^2 , the squared multiple correlation**, is the proportion of the variation in the response variable y that is explained by the model.

$$R^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{SSModel}{SSTotal}$$

In the particular case of multiple linear regression, the model is all p explanatory variables taken together.

The square root of R^2 , called the **multiple correlation coefficient**, is the correlation between the observations and the predicted values.

11.2 A Case Study



- Preliminary Analysis
- Relationships Between Pairs of Variables
- Regression on High School Grades
- Interpretation of the Results
- Refining the Results
- Regression Using All Variables
- Test for a Collection of Regression Coefficients

Preliminary Analysis



The data on 224 first-year computer science majors at a large university in a given year for each student include:

- * Cumulative GPA after two semesters at the university (y, response variable)
- * SAT math score (SATM, x1, explanatory variable)
- * SAT verbal score (SATV, x2, explanatory variable)
- * Average high school grade in math (HSM, x3, explanatory variable)
- * Average high school grade in science (HSS, x4, explanatory variable)
- * Average high school grade in English (HSE, x5, explanatory variable)

Here are the summary statistics for these data given by software **SAS**:

Variable	N	Mean	Std Dev	Minimum	Maximum
GPA	224	2.6352232	0.7793949	0.1200000	4.0000000
SATM	224	595.2857143	86.4014437	300.0000000	800.0000000
SATV	224	504.5491071	92.6104591	285.0000000	760.0000000
HSM	224	8.3214286	1.6387367	2.0000000	10.0000000
HSS	224	8.0892857	1.6996627	3.0000000	10.0000000
HSE	224	8.0937500	1.5078736	3.0000000	10.0000000

Relationships Between Pairs of Variables



The first step in multiple linear regression is to study all pair-wise relationships between the $p + 1$ variables. Here is the SAS output for all pair-wise correlation analyses (value of r and two-sided p -value of $H_0: \rho = 0$).

Pearson Correlation Coefficients / Prob > |R| under Ho: Rho=0 / N = 224

	GPA	SATM	SATV	HSM	HSS	HSE
GPA	1.00000 0.0	0.25171 0.0001	0.11449 0.0873	0.43650 0.0001	0.32943 0.0001	0.28900 0.0001
SATM	0.25171 0.0001	1.00000 0.0	0.46394 0.0001	0.45351 0.0001	0.24048 0.0003	0.10828 0.1060
SATV	0.11449 0.0873	0.46394 0.0001	1.00000 0.0	0.22112 0.0009	0.26170 0.0001	0.24371 0.0002
HSM	0.43650 0.0001	0.45351 0.0001	0.22112 0.0009	1.00000 0.0	0.57569 0.0001	0.44689 0.0001
HSS	0.32943 0.0001	0.24048 0.0003	0.26170 0.0001	0.57569 0.0001	1.00000 0.0	0.57937 0.0001
HSE	0.28900 0.0001	0.10828 0.1060	0.24371 0.0002	0.44689 0.0001	0.57937 0.0001	1.00000 0.0

Regression on High School Grades

For simplicity, let's first run a multiple linear regression using **only the three high school grade averages**:

Dependent Variable: GPA

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F
Model	3	27.71233	9.23744	18.861	0.0001
Error	220	107.75046	0.48977		
C Total	223	135.46279			

P-value very significant

Root MSE	0.69984	R-Square	0.2046
Dep Mean	2.63522	Adj R-sq	0.1937
C.V.	26.55711		

R^2 is fairly small (20%)

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T
INTERCEP	1	0.589877	0.29424324	2.005	0.0462
HSM	1	0.168567	0.03549214	4.749	0.0001
HSS	1	0.034316	0.03755888	0.914	0.3619
HSE	1	0.045102	0.03869585	1.166	0.2451

HSM significant

HSS, HSE not

Regression on High School Grades



Dependent Variable: GPA

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F	P-value very significant
Model	3	27.71233	9.23744	18.861	0.0001	
Error	220	107.75046	0.48977			
C Total	223	135.46279				
Root MSE	0.69984	R-Square	0.2046			R^2 is fairly small (20%)
Dep Mean	2.63522	Adj R-sq	0.1937			
C.V.	26.55711					

The ANOVA for the multiple linear regression using only HSM, HSS, and HSE is significant. At least one of the regression coefficients is significantly different from zero.

R^2 is fairly small (0.205) → only about 20% of the variations in cumulative GPA can be explained by these high school scores. (*Remember, a small p-value does not imply a large effect.*)

Interpretation of the Results



Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T	
INTERCEP	1	0.589877	0.29424324	2.005	0.0462	
HSM	1	0.168567	0.03549214	4.749	0.0001	HSM significant
HSS	1	0.034316	0.03755888	0.914	0.3619	
HSE	1	0.045102	0.03869585	1.166	0.2451	HSS, HSE not

The tests of hypotheses for each b within the multiple linear regression reach significance for HSM only.

We found a significant correlation between HSS and GPA when analyzed by themselves, so why is b_{HSS} not significant in the multiple regression equation? Well, HSS and HSM are also significantly correlated.

When all three high school averages are used together in the multiple regression analysis, only HSM contributes significantly to our ability to predict GPA.

Refining the Model

We now **drop** the least significant variable from the previous model: HSS.

Dependent Variable: GPA

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F
Model	2	27.30349	13.65175	27.894	0.0001
Error	221	108.15930	0.48941		
C Total	223	135.46279			
Root MSE		0.69958	R-Square	0.2016	
Dep Mean		2.63522	Adj R-sq	0.1943	
C.V.		26.54718			

P-value very significant

R² is small (20%)

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T
INTERCEP	1	0.624228	0.29172204	2.140	0.0335
HSM	1	0.182654	0.03195581	5.716	0.0001
HSE	1	0.060670	0.03472914	1.747	0.0820

HSM significant

HSE not

The conclusions are about the same. But notice that the actual regression coefficients have changed.

predicted GPA = .590 + .169HSM + .045HSE + .034HSS

predicted GPA = .624 + .183HSM + .061HSE

Refining the Model



Let's run a multiple linear regression with the **two SAT scores only**.

Dependent Variable: GPA

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F
Model	2	8.58384	4.29192	7.476	0.0007
Error	221	126.87895	0.57411		
C Total	223	135.46279			

P-value very significant

Root MSE	0.75770	R-Square	0.0634
Dep Mean	2.63522	Adj R-sq	0.0549
C.V.	28.75287		

R² is very small (6%)

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T
INTERCEP	1	1.288677	0.37603684	3.427	0.0007
SATM	1	0.002283	0.00066291	3.444	0.0007
SATV	1	-0.000024562	0.00061847	-0.040	0.9684

SATM significant
SATV not

The ANOVA test for β_{SATM} and β_{SATV} is significant. At least one is not zero. R^2 is really small (0.06). Only 6% of GPA variation is explained by these tests. When taken together, only SATM is a significant predictor of GPA (P 0.0007).

Regression Using All Variables

Dependent Variable: GPA

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Prob>F	P-value very significant
Model	5	28.64364	5.72873	11.691	0.0001	
Error	218	106.81914	0.49000			
C Total	223	135.46279				

Root MSE 0.70000 R-Square 0.2115 **R² fairly small (21%)**
 Dep Mean 2.63522 Adj R-sq 0.1934
 C.V. 26.56311

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	T for H0: Parameter=0	Prob > T	HSM significant
INTERCEP	1	0.326719	0.39999643	0.817	0.4149	
SATM	1	0.000944	0.00068566	1.376	0.1702	
SATV	1	-0.000408	0.00059189	-0.689	0.4915	
HSM	1	0.145961	0.03926097	3.718	0.0003	
HSS	1	0.035905	0.03779841	0.950	0.3432	
HSE	1	0.055293	0.03956869	1.397	0.1637	

The overall test is significant, but only the average high school math score (HSM) makes a significant contribution in this model to predicting the cumulative GPA. This conclusion applies to computer majors at this large university.

Test for a Collection of Regression Coefficients



Regression Statistics						
Multiple R	0.459837234					
R Square	0.211450282					
Adjusted R Square	0.193364279					
Standard Error	0.699997195					
Observations	224					
ANOVA						
	df	SS	MS	F	Significance F	
Regression	5	28.64364489	5.728729	11.69138	5.06E-10	
Residual	218	106.8191439	0.489996			
Total	223	135.4627888				
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	upper 95%
intercept	0.326718739	0.399996431	0.816804	0.414932	-0.461636967	1.115074446
HSM	0.14596108	0.039260974	3.717714	0.000256	0.068581358	0.223340801
HSS	0.03590532	0.037798412	0.949916	0.343207	-0.03859183	0.11040247
HSE	0.055292581	0.039568691	1.397382	0.163719	-0.022693622	0.133278785
SATM	0.000843593	0.000685657	1.376187	0.170176	-0.000407774	0.002294959
SATV	0.00040785	0.000591893	-0.68906	0.491518	-0.001574415	0.00075816

Test for a Collection of Regression Coefficients



The regression equation is

$$\text{GPA} = 0.327 + 0.146 \text{ HSM} + 0.0359 \text{ HSS} + 0.0553 \text{ HSE} + 0.000944 \text{ SATM} - 0.000408 \text{ SATV}$$

Predictor	Coef	StDev	T	P
Constant	0.3267	0.4000	0.82	0.415
HSM	0.14596	0.03926	3.72	0.000
HSS	0.03591	0.03780	0.95	0.343
HSE	0.05529	0.03957	1.40	0.164
SATM	0.0009436	0.0006857	1.38	0.170
SATV	-0.0004078	0.0005919	-0.69	0.492

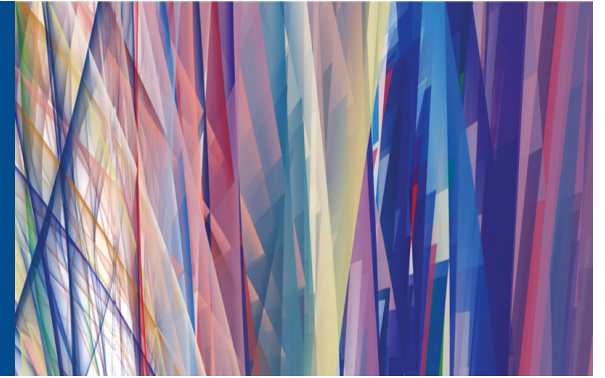
S = 0.7000 R-Sq = 21.1% R-Sq(adj) = 19.3%

Analysis of Variance

Source	DF	SS	MS	F	P
Regression	5	28.6436	5.7287	11.69	0.000
Error	218	106.8191	0.4900		
Total	223	135.4628			

Chapter 11

Multiple Regression



11.1 Inference for Multiple Regression

11.2 A Case Study