MBC 638

LIVE SESSION WEEK 7

Agenda

Topic	Time	Thursday Section	Sunday Section
Introduction	5 min	9:00-9:05	6:30-6:35
Quiz 2 Recap	15 min	9:05-9:20	6:35-6:50
Highlights from Week 7 Video	65 min	9:20-10:25	6:50-7:55
Review of Upcoming Assignments and Open Question	5 min	10:25-10:30	7:55-8:00

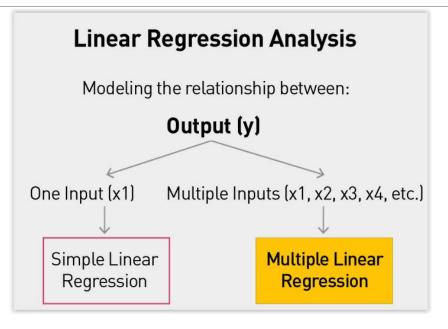
Interim Grade Report

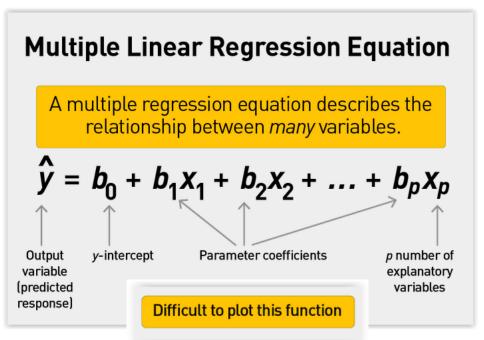
Fall 2017	Points>		10	10	10	5	3	2	10	5	3	2	20	20			
																Highest	
															Points Still	Potential	4
Last Name, First Name		<u>Email</u>	<u>Participation</u>	Prob Def Wkst	Quiz #1	Hwk #1	Hwk #2	Hwk #3	Quiz #2	Hwk #4	Hwk #5	Hwk #6	Project	<u>Final</u>	<u>Availbable</u>	<u>Grade</u>	Rounded
1 Doe	Jane		7	9.5	8	O 4	3	2	9	0	0	0			53	95.5	96.0

Agenda

Topic	Time	Thursday Section	Sunday Section
Introduction	5 min	9:00-9:05	6:30-6:35
Quiz 2 Recap	15 min	9:05-9:20	6:35-6:50
Highlights from Week 7 Video	65 min	9:20-10:25	6:50-7:55
Review of Upcoming Assignments and Open Question	5 min	10:25-10:30	7:55-8:00

Highlights: Video Segment 7.3: Multiple Linear Regression





Highlights: Video Segment 7.4: Multiple Regression Using Excel

- P value helps you identify which input variable(s) are important/useful in describing Y.
- We want Ps lower than our alpha, assume .05, so if p is low Ho must go our implied Ho is, X doesn't describe Y and Ha is, X does describe Y – so if we have a low p value we reject Ho and say X does describe Y.
- Then re-run your regression without the Xs that don't help describe Y to develop a better model.

A	E	1	C	D	E	F	G	H	- 1	1	K	L					
				output (y)	inputs (x)								rotar	120	0100117106		
First na	me Lastin	ame	Team	Runs Scored	Hits	Doubles	Triples	Home Runs	RBIs	Walks	Bat Ave	Yankees?					
Ichiro	Suzuk		SEA	111				7 6	6		9 0.35			Coefficients	Standard Error	t Stat	P-value
Delmo			TBD	65				0 13			6 0.28		_	Coefficients		Latur	
Alexis	Rios		TOR	114				7 24			5 0.29		Intercept	-6.677284091	8.358708258	-0.79884	0.4260
Derek	Jeter		NYY	102				4 12			6 0.32		Hits	0.437991874	0.048707908	8.992213	0.00000000
Michae	0.000		TEX	80				1 9			7 0.31		IIIG	0.437331074	0.040707300	0.332213	
Orland		77	LAA	101				1 8			4 0.30		Doubles	0.001881127	0.142779213	0.013175	0.9895
Nick	Marka Sizem		BAL CLE	97				3 23 5 24	11	8 10	0.27		Triples	1.236783363	0.291004023	4.250056	0.000
Brian	Rober		BAL	103		150		5 12			9 0.2	3 - 1					
Robins			NYY	93				7 19			9 0.30		Home Runs	0.757792868	0.174045929	4.353982	0.00002863
Curtis	Grand			- 100							2 0.30		RBIs	-0.203601646	0.075335531	-2.7026	0.0079
Aaron	Hill		TOR	₩ 122 87	177			2 17			1 0.29		100000	0.000540056			
Bobby	Abreu		NYY	123				5 16	10		4 0.28		Walks	0.283549256	0.043817935	6.471078	0.0000
David	DeJes	us	KCR	101	157	29		9 7	5	8 6	4 0.2	5 0	Bat Ave	12.65150623	38.82469554	0.325862	0.7453
Torii	Hunte	er .	MIN	94	172	45		1 28	10	7 4	0.28	7 0	Yankees?	9 194227296	3.330203822	2.76086	0.0067
Adrian	Beltre		SEA	87	164	41		2 26	9	9 3	8 0.27	5 0	rankees:	3.134227200	3.330203022	2.70000	0.000
Maggli	Ordor	nez	DET	117	216	54		0 28	13	9 7	6 0.36	3 0				N	
Jose	Guille	m	SEA	84	172	28		2 23	9	9 4	1 0.2	9 0					-l-h- 0.0F
Justin	Morne	eau	MIN	84	160	31		3 31	11	1 6	4 0.27	1 0					alpha = 0.05

$$\begin{split} \hat{y} &= -6.678 + 0.438 x_{hits} + 0.002 x_{dbles} + 1.237 x_{triples} + 0.758 x_{homerums} \\ &- 0.204 x_{RBIs} + 0.284 x_{walks} + 12.652 x_{batave} + 9.194 x_{yanks} \end{split}$$

Highlights: Video Segment 7.5: Correlation, F Test, and Model Building

Correlation Coefficients

- Multiple correlation coefficient (R)
 - $\circ~$ Measures relationship between observed output y and predicted output y
- Coefficient of determination (R2)
 - Proportion of the variation in response variable that is explained by model
 - $\circ~$ Always increases when another x is added to model
 - $\circ~$ E.g., if model to predict output y has two x inputs, adding a third x will increase ${\bf R}^2$

Adjusted R square is a better measure to look at when trying to determine how good your model is...how much of the variability in Y is explained by your equation

Multiple R is a little different for multiple regression in that it is the relationship between the observed and predicted Y. In simple linear regression this was the correlation between the X and Y.

Correlation Coefficients (cont.)

- Adjusted R²
 - Measure that helps account for too many unnecessary x inputs
 - \circ Often x inputs are added in order to increase \mathbb{R}^2
 - Higher R² makes model seem better, but having more inputs complicates forecast

Adjusted Coefficient of Determination

We measure the goodness of a regression equation using the coefficient of determination $r^2 = SSR/SST$. In multiple regression, we use the same formula for the coefficient of determination (though the letter r is promoted to a capital R).

Multiple Coefficient of Determination R^2

The multiple coefficient of determination is given by:

$$R^2 = SSR/SST \quad 0 \le R^2 \le 1$$

where SSR is the sum of squares regression and SST is the total sum of squares. The multiple coefficient of determination represents the proportion of the variability in the response y that is explained by the multiple regression equation.

Adjusted Coefficient of Determination

Unfortunately, when a new x variable is added to the multiple regression equation, the value of R^2 always increases, even when the variable is not useful for predicting y. So, we need a way to adjust the value of R^2 as a penalty for having too many unhelpful x variables in the equation.

Adjusted Coefficient of Determination R2 adi

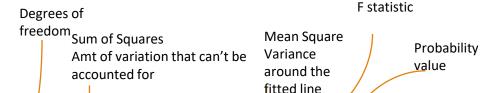
The adjusted coefficient of determination is given by:

$$R_{\text{adj}}^2 = 1 - (1 - R^2) \left(\frac{n - 1}{n - k - 1} \right)$$

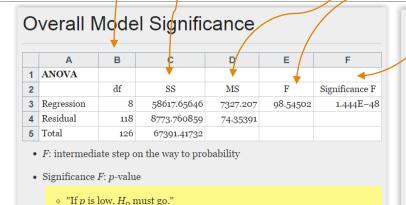
where n is the number of observations, k is the number of x variables, and R^2 is the multiple coefficient of determination.

Highlights: Video Segment 7.5: Correlation, F Test, and Model Building

	Α		В	
1	8 Input	<u>Variables</u>		
2			Audtinia D	
3	SUMMARY OUTPUT		•	sure of the relationship of the
4	Regressio	n Statistics O	bserved output	of Y and the predicted of Y
5	Multiple R	0.9326355		
6	R ²	0.0090009		ortion of the variation explained by
7	Adjusted R ²	0.8609824 ti	ne model, alway	s increases when you add more Xs
8	Standard Error	8.6228711	alternational D. Construction	- d
9	Observations	127	•	ed = accounts for excess Xs, a better ariability explained by the model



Highlights: Video Segment 7.5: Correlation, F Test, and Model Building



ANOVA= analysis of variance

F provides a p value, our measure of goodness

Ho: coefficient of variables =0

Ha: At least something in my model is a good x input and should have a value, at least one variable helps forecast Y.

Must have a low F to have a valid model.

Hypothesis Test

•
$$H_0$$
: $\beta_1 = \beta_2 = ... = \beta_p = 0$

• H_a : at least one $\beta_j \neq 0$

$$\hat{y} = -6.678 + 0.438x_{hits} + 0.002x_{dbles} + 1.237x_{triples} + 0.758x_{homeruns} -0.204x_{RBIs} + 0.284x_{walks} + 12.652x_{batave} + 9.194x_{vanks}$$

- Significant p-value does not mean all x's have significant influence on y
 - o Does mean one or more x's has significant influence
 - Refine model to eliminate less useful variables
- Failing to reject H_0 : no evidence that any coefficient $\beta_i \neq 0$
- To assess model, look at ANOVA
 - Low F statistic indicates good model

F Test for Multiple Regression

The multiple regression model is an extension of the model from Section 13.1, and approximates the relationship between *y* and the collection of *x* variables.

Multiple Regression Model

The **population multiple regression equation** is defined as:

$$y = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

where β_1 , β_2 , ..., β_k are the parameters of the population regression equation, k is the number of x variables, and ε is the error term that follows a normal distribution with mean 0 and constant variance.

The population parameters are unknown, so we must perform inference to learn about them. We begin by asking: *Is our multiple regression useful?* To answer this, we perform the *F* test for the overall significance of the multiple regression.

F Test for Multiple Regression

The hypotheses for the *F* test are:

$$H_0$$
: $\beta_1 = \beta_2 = \dots = \beta_k = 0$

 H_a : At least one of the β 's \neq 0.

The *F* test is not valid if there is strong evidence that the regression assumptions have been violated.

F Test for Multiple Regression

If the conditions for the regression model are met

Step 1: State the hypotheses and the rejection rule.

Step 2: Find the *F* statistic and the *p*-value. (Located in the ANOVA table of computer output.)

Step 3: State the conclusion and the interpretation.

Highlights: Video Segment 7.5: Correlation, F Test, and Model Building

Model Building

- Nonlinear regression models exist
 - o Remember: practical, graphical, statistical
 - Plot data whenever possible
- Curved relationship: x^2 variable(s); possibly quadratic function(s)
- Reciprocals: 1/x
- Interaction terms: $x_1 \times x_2$
- Too many variables → pare down
 - Best models: fewer x inputs but same predictive value

Highlights: Video Segment 7.6:Just Correlation

	Runs Scored	Hits	Doubles	Triples	Home Runs	RBIs	Walks	Bat Ave	Yankees?
Runs Scored	♣ 1								
Hits	0.843722818	1							Look through the table to
Doubles	0.672563335	0.788082	1						
Triples	0.32311721	0.277987	0.123538	1					identify high correlation
Home Runs	0.521110368	0.346261	0.424893	-0.17987	1				variables.
RBIs	0.665750495	0.672397	0.682796	-0.06029	0.808734	1			
Walks	0.630325512	0.384619	0.352119	-0.00888	0.604142	0.544677	1		
Bat Ave	0.568341926	0.727978	0.540431	0.166983	0.124469	0.417129	0.185354	1	
Yankees?	0.34809467	0.263619	0.189084	0.098172	0.172827	0.272852	0.203198	0.199359	1

Multiple R is the observed to predicted Y relationship.

Adjusted R Square accounts for the fact that you have possibly too many Xs, discounts the Xs that are not relevant to the model-the high p value Xs.

Overall, how well do my Xs describe my output.

SUMMARY OUTPUT		
Regression St	atistics	
Multiple R	0.932635474	÷
R Square	0.869808928	
Adjusted R Square	0.860982415	
Standard Error	8.622871076	
Observations	127	

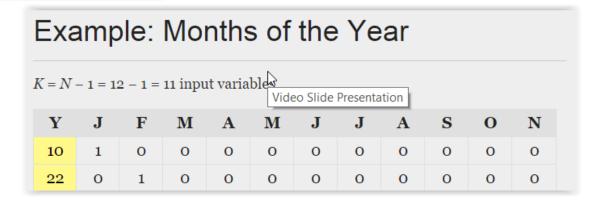
	Coefficients	andard Err	t Stat	P-value	Lower 95%	Upper 95%	ower 95.09	pper 95.09
Intercept	-6.677284091	8.358708	-0.79884	0.425987	-23.2298	9.875234	-23.2298	9.875234
Hits	0.437991874	0.048708	8.992213	4.81E-15	0.341537	0.534447	0.341537	0.534447
Doubles	0.001881127	0.142779	0.013175	0.98951	-0.28086	0.284623	-0.28086	0.284623
Triples	1.236783363	0.291004	4.250056	4.29E-05	0.660516	1.813051	0.660516	1.813051
Home Runs	0.757792868	0.174046	4.353982	2.86E-05	0.413135	1.102451	0.413135	1.102451
RBIs	-0.203601646	0.075336	-2.7026	0.007895	-0.35279	-0.05442	-0.35279	-0.05442
Walks	0.283549256	0.043818	6.471078	2.31E-09	0.196778	0.370321	0.196778	0.370321
Bat Ave	12.65150623	38.8247	0.325862	0.745106	-64.232	89.53497	-64.232	89.53497
Yankees?	9.194227286	3.330204	2,76086	0.006687	2.599517	15.78894	2.599517	15.78894

This can be used to look through to identify high p-value Xs to potentially eliminate.

Highlights: Video Segment 7.7: Categorical Input Variables

Categorical Input Variables (x's)

- Discrete, categorical x variable may be influencing output y
 - o E.g., x_{vankees}: o or 1
 - $\circ~$ Put categorical x variables in regression by assigning each variable 0 or 1



Highlights: Video Segment 7.7: Categorical Input Variables

Example: Power Tools

What drives the price of a particular power tool?

- Output (y) = price of the tool (continuous data)
- Inputs (x) =
 - o Product brands (discrete data)
 - Types of accessories (discrete data)
 - Weight of the tool (continuous data)

Example: Power Tools: Data

	Brand 1	Brand 2	Brand 3	Access	sories	
Price	Sears	Toshiba	Dremel	Basic	Xtra	Weight
у	X1	X2	х3	X4	X5	Х6
20	1	0	0	1	О	10
40	O	1	0	1	О	30
60	O	0	1	О	1	32
30	1	0	0	1	O	12
35	0	1	0	1	0	11

1 less than N is needed for regression

		Brand 1	Duanda	Accessories	
			Brand 2		
	Price	Sears	Toshiba	Basic	Weight
	y	X_1	X ₂	x_4	x ₆
ď	20	1	0	1	10
	40 60	0	1	1	30
	60	0	0	0	32
	30	1	0	1	12
	35	0	1	1	11
	6-	^	^	^	~0

Highlights: Video Segment 7.7: Categorical Input Variables

	Brand 1	Brand 2	Accessories	
Price	Sears	Toshiba	Basic	Weight
v	X ₁	X ₂	X_{d}	X6

Example: Power Tools: Regression Equation

	Coefficients	Std Error	t Stat	P-Value	Lower 95%
Intercept	48.3628029	14.17717	3.411315	0.019016	11.91922486
X1	-21.377046	9.757484	-2.19084	0.080012	-46.45945872
X2	-12.478716	7.311039	-1.70683	0.148563	-31.27233932
X4	-4.2239686	8.881025	-0.47562	0.6544	-27.0533713
x6	0.2848723	0.391725	0.727225	0.499698	-0.722089489

$$\hat{y} = 48.36 - 21.38x_1 - 12.48x_2 - 4.22x_4 + 0.28x_6$$

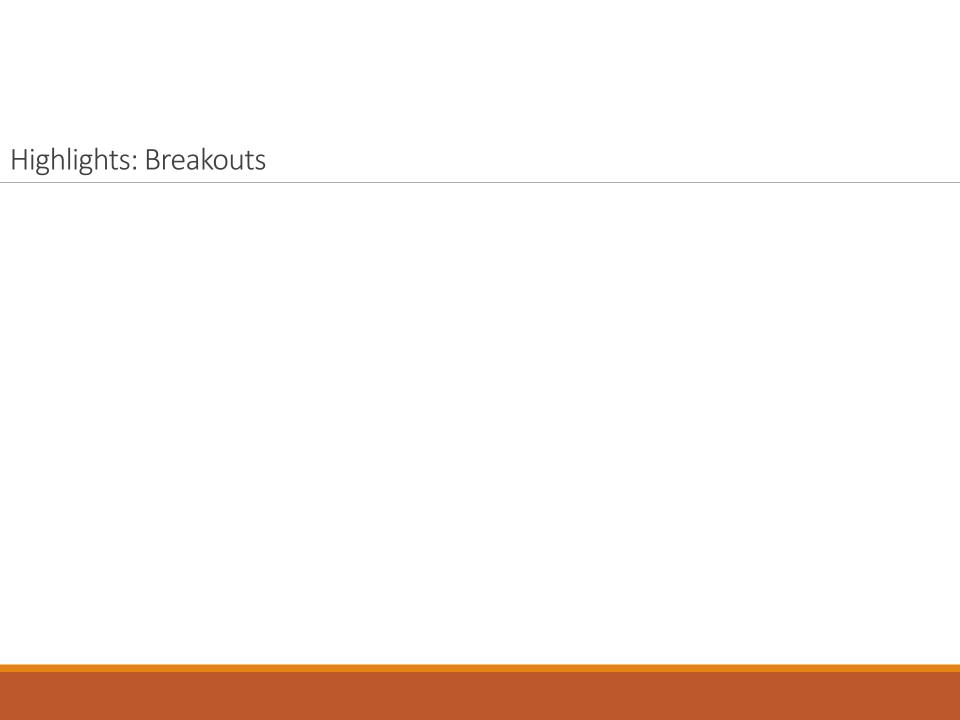
Estimating Price (y)

$$\hat{y} = 48.36 - 21.38x_1 - 12.48x_2 - 4.22x_4 + 0.28x_6$$

Tool characteristics:

- Dremel (Brand 3): $x_1 = 0$ and $x_2 = 0$
- Extra accessories: $x_4 = 0$
- Weight 25 lb: $x_6 = 25$

 \hat{y} =48.36-21.38(0)-12.48(0)-4.22(0)+0.28(25)=\$55.36



Highlights: Video Segment 7.8: Test your knowledge

SUMMARY OUTPUT													
5													
Regression Statistics Multiple R 0.91478924		247											
	0.914789247				wnei	when only using hand to predict foot 0.79							
R Square Adjusted R Square	0.825586	839367				0.7	-						
Standard Error	0.789031					0.	.0						
Observations	0.789031	32					-						
Observations		32					-						
ANOVA													
	df		SS	MS		F		Significa	nce F				
Regression		2	93	46	74.3	8694761	_	3.8275					
Residual		29	18	0.6									
Total		31	111										
	Coefficie	nts	dard L	t Stat	t P-value			Lower	95%	Upper 95%	ver 95.0%		
Intercept	4.379011	.892	3	1.5	0.15	566205	58	-1.7648	76316	10.5229001	-1.76	10.523	
M/F	1.096222	729	0.5	2.4	0.02	242955	59	0.1666	20861	2.025824597	0.167	2.0258	
Hand	1.090436	031	0.2	6.4	5.3	9068E-0	<mark>)7</mark>	0.7418	11647	1.439060416	0.742	1.4391	
			4 07	2.4.0	004	4 0004	11/2						
Regression Equation		y =	-4.3/	9+1.0	96X +	1.0904	IX2						
			_	, , , , , ,									
my actual hand size in inches				7.625									
my actual hand size i	n cms		19.	3675								<u> </u>	
								Using Simple Linear Regression - only H			nly Hand		
Use formula to predict foot size y			e y'=4.379+1.096X + 1.0904X2						y' = (1.39X 19.3675)52866				
									predicted foot size		26.3	392165000	
predicted foot size		25.497322						actual foot in cms			25.4		
actual foot in cms		25.4						residual		-0.9	992165000		
residual =Ac		=Ac	Actual foot size - predicted foot size										
			5.4-25.497322										
			0.097322000 So the multiple regression when considering gende							ler creat	es a		
							•	_		better predic			

Agenda

Topic	Time	Thursday Section	Sunday Section
Introduction	5 min	9:00-9:05	6:30-6:35
Quiz 2 Recap	15 min	9:05-9:20	6:35-6:50
Highlights from Week 7 Video	65 min	9:20-10:25	6:50-7:55
Review of Upcoming Assignments and Open Question	5 min	10:25-10:30	7:55-8:00

Review of Upcoming Assignments: Thursday

- 1. HMWK #4, Learning Curve on Chapter 4, in LaunchPad is due until Sunday, 11/18, midnight EST.
- 2. Understanding Variation Book –Live Class #8 is focused on the material from this book
- 3. Optional Learning Opportunity: 8.8 Relate Control Charts to Your Project
- 4. Projects.... Should begin Improve in the next week or so, such that you have time to collect "after" data

	November 2018						
	Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
Week #7	11	12	13	14	15	16	17
	Quiz #2 DUE				Live Class #7		
Week #8	18	19	20	21	22	23	24
	Homework #4 Due: 1. CH 4 Learning Curve Reminder: Start reading Understanding Variation	Live Class #8 RESCHEDULE 9PM EST			Live Class #8 RESCHEDULE		
Week #9	25	26	27	28	29	30	1
	Homework #5 Due: 1. Problems 1-10 pg 114-116 in Understanding Variation				Live Class #9		
	December 2018						
	Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
Week #10	2	3	4	5	6	7	8
	Homework #6 Due: 1. Time Series Problem posted in Excel				Live Class #10		
Week #11	9	10	11	12	13	14	15
		Process Improvement Project DUE			Live Class #11: FINAL EXAM		

Review of Upcoming Assignments: Sunday

- 1. HMWK #4, Learning Curve on Chapter 4, in LaunchPad is due until Wednesday, 11/21, midnight EST.
- 2. Understanding Variation Book –Live Class #8 is focused on the material from this book
- 3. Optional Learning Opportunity: 8.8 Relate Control Charts to Your Project
- 4. Projects.... Should begin Improve in the next week or so, such that you have time to collect "after" data

	November 2018						
	Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
Week #7	18	19	20	21	22	23	24
	Live Class #7			Homework #4 Due: 1. CH 4 Learning Curve Reminder: Start reading Understanding Variation			
Week #8	25	26	27	28	29	30	1
	Live Class #8			Homework #5 Due: 1. Problems 1-10 pg 114-116 in Understanding Variation			·
	December 2018						
	Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
Week #9	2	3	4	5	6	7	8
	Live Class #9			Homework #6 Due: 1. Time Series Problem posted in Excel			
Week #10	9	10	11	12	13	14	15
	Live Class #10				<u>Process</u> <u>Improvement</u> <u>Project DUE</u>		
Week #11	16	17	18	19	20	21	22
	Live Class #11: FINAL EXAM						