

Woody Bendle

#### Advanced methods of analysis & Innovation Today's Topics

Regression Types of Regression / Data Cross Section Time Series Univariate Multivariate

Cluster Analysis K-Means

Innovation

Introduction to the i<sup>3</sup> Innovation Process

Definition:

Regression analysis is a statistical process for estimating the relationships among variables (*data*)

Definition:

Regression analysis is a statistical process for estimating the relationships among variables (*data*)

• How does one thing affect another?

Univariate

Definition:

Regression analysis is a statistical process for estimating the relationships among variables (*data*)

• How does one thing affect another?

Univariate Multivariate

How do several things together affect something else?

Definition:

Regression analysis is a statistical process for estimating the relationships among variables (*data*)

- How does one thing affect another?
- How do several things together affect something else?
- What is the relationship over time?

Univariate Multivariate Time Series

Theoretical

Vs.

Hypothesis: "I think \_\_\_\_\_ affects \_\_\_\_\_ in some way"

I think the amount of rain affects how much the grass grows

A Theoretical

I have a bunch of data, let's dump it into the machine and see what pops out...

#### Theoretical approach

Pros:

You're approaching a problem with a belief construct

Scientific Method

Cons:

You're approaching a problem with a limited belief construct

#### A Theoretical approach

Pros:

Dispassionate

Avoids "Curse of Knowledge"

Can find relationships that challenge convention

Cons:

You're approaching a problem with a limited belief construct

Results can lead to bad / wrong conclusions

- Suspected relationship between two things
- Often shown using X-Y Graph



**Common Univariate Examples** 

**Demand Function in Economics** 



**Common Univariate Examples** 

**Demand Function in Economics** 



**Common Univariate Examples** 

**Demand Function in Economics** 



Y = ƒ(X)

Quantity = f(Price)

Negative or Inverse Relationship

**Common Univariate Examples** 



**Common Univariate Examples** 



**Common Univariate Examples** 

**Cost Function** 



Common Univariate Examples

.

What are some other examples of these types of univariate relationships?

17

Viewing / Visualizing Data



Х

Viewing / Visualizing Data



## Equation for a line



- m = slope of the line = "rise over run"
- b = Y intercept

















Equation for a line Y = mX + b



Х



Exercise:

Y

- Plot (X,Y) coordinates
- Using Line Equation, solve for "m" and "b"

Y	Х
3	10
4	8
5	6
6	4
7	2



Х

# Y = mX + b

Exercise:

Y

- Plot (X,Y) coordinates
- Using Line Equation, solve for "m" and "b"

Y	Х
3	10
4	8
5	6
6	4
7	2



$$b = 8$$

# Real world data are rarely "perfect"



33

Can draw a lot of lines through this scatter plot

How do you determine the "best" line?



Equation for a line Y = mX + b

Regression Equation Y = a + bX

where...

a = Y intercept b = Slope

or...

Y = a + bX

**Ordinary Least Squares (OLS)** 

 $\mathbf{\hat{Y}} = \mathbf{a} + \mathbf{b}\mathbf{X}$ 

A good line is one that minimizes the sum of squared differences between the points and the line


Advanced methods of analysis Univariate Regression Ordinary Least Squares (OLS) Some Equations...

 $\mathbf{\hat{Y}} = \mathbf{a} + \mathbf{b}\mathbf{X}$  $SS_{xy} = S X_i Y_i - \frac{(S X_i) (S Y_i)}{2}$  $b = \frac{SS_{xy}}{SS_{xy}}$  $SS_{xx} = S X_i^2 - \frac{(S X_i)^2}{n}$  $a = \overline{Y} - b\overline{X}$ Y = mean of YX = mean of X37

		х	Y				
	Obs	\$	# Sold		X <sup>2</sup>	ΧΧΥ	
	1	10	4	_	100	40	
	2	5	9		25	45	
	3	6	10		36	60	
	4	3	14		9	42	
	5	6	8		36	48	
	6	6	10		36	60	
	7	8	6		64	48	
	8	5	12		25	60	
	9	7	9		49	63	
n =	10	5	10	_	25	50	
	Sum	61	92	Sum X <sup>2</sup>	405	516	Sum X x Y
	(Sum X) <sup>2</sup>	3,721					
		Mean X	Mean Y				

### OLS – Solving our example by hand

6

9

A = a + bX $= \frac{SS_{xy}}{SS_{xx}}$  $SS_{xy} = X_iY_i - \frac{(X_i)(Y_i)}{n}$ 

$$SS_{xy} = 516 - \frac{61 \times 92}{10} = -45.2$$

		х	Y				
	Obs	\$	# Sold		X <sup>2</sup>	X x Y	
	1	10	4		100	40	
	2	5	9		25	45	
	3	6	10		36	60	
	4	3	14		9	42	
	5	6	8		36	48	
	6	6	10		36	60	
	7	8	6		64	48	
	8	5	12		25	60	
	9	7	9		49	63	
n =	10	5	10		25	50	
	Sum	61	92	Sum X <sup>2</sup>	405	516	Sum X x Y
	(Sum X) <sup>2</sup>	3,721					

OLS – So	lving ou	ır exam	nple k	oy hand
----------	----------	---------	--------	---------

Mean X	Mean Y
6	9

$\mathbf{\hat{Y}} = \mathbf{a} + \mathbf{b}\mathbf{X}$
$=\frac{SS_{xy}}{SS_{xx}}$
$SS_{xy} = X_i Y_i - \frac{(X_i) (Y_i)}{n}$
$SS_{xy} = 516 - \frac{61 \times 92}{10} = -45.2$
$SS_{xx} = X_i^2 - \frac{(X_i)^2}{n}$
$SS_{xx} = 405 - \frac{3,721}{10} = 32.9$
$b = \frac{-45.2}{32.9} = -1.37$

		Х	Y				
_	Obs	\$	# Sold		X <sup>2</sup>	ΧΧΥ	
	1	10	4	-	100	40	
	2	5	9		25	45	
	3	6	10		36	60	
	4	3	14		9	42	
	5	6	8		36	48	
	6	6	10		36	60	
	7	8	6		64	48	
	8	5	12		25	60	
	9	7	9		49	63	
n =	10	5	10		25	50	
-	Sum	61	92	Sum X <sup>2</sup>	405	516	Sum X x Y
	(Sum X) <sup>2</sup>	3,721					
		Mean X	Mean Y				

### OLS – Solving our example by hand

6

9

<b>^</b> Y = a + bX
$= \frac{SS_{xy}}{SS_{xx}}$
$SS_{xy} = X_iY_i - \frac{(X_i)(Y_i)}{n}$
$SS_{xy} = 516 - \frac{61 \times 92}{10} = -45.2$
$SS_{xx} = X_i^2 - \frac{(X_i)^2}{n}$
$SS_{xx} = 405 - \frac{3,721}{10} = 32.9$
$b = \frac{-45.2}{32.9} = -1.37$
$=\overline{Y} - \overline{X}$
a = 9 - (-1.37 x 6) = <b>17.6</b>
$\hat{Y} = 17.6 - 1.37 \times X^{40}$



- **Univariate Regression**
- **Ordinary Least Squares (OLS)**
- **Using Excel**

## Data Analysis (Regression)



					<u>^</u>	. – .			
Regression St	atistics	Y = 17.6 - 1.37 * X							
Multiple R	0.93				7	7			
R Square	0.87								
Adjusted R Square	0.85								
Standard Error	1.09								
Observations	10.00								
ANOVA						_			
	df	SS	MS	F	Significance F				
Regression	1	62.10	62.10	52.29	0.00	_			
Residual	8	9.50	1.19						
Total	9	71.60				_			
						-			
	Coefficients St	andard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%	
Intercept	17.58	1.21	14.54	0.00	14.79	20.37	14.79	20.37	
X Variable 1	-1.37	0.19	-7.23	0.00	-1.81	-0.94	-1.81	-0.94	

- **Univariate Regression**
- Ordinary Least Squares (OLS)
- Assessing the Model

		R Square
EXCEL SUMMARY O	UTPUT	
Regression St	atistics	<u>^</u>
Multiple R	0.93	
R Square	0.87	$ _{P^2} - S(r_i - r_i)$
Adjusted R Square	0.85	$R^{-} = \overline{C/V}$
Standard Error	1.09	$  S(Y_i - Y_i) $
Observations	10.00	

$$R^{2} = \frac{S(\hat{Y_{i}} - \overline{Y})^{2}}{S(Y_{i} - \overline{Y})^{2}}$$

Also called the coefficient of determination. Reflects the percent of the variance (dispersion) in the data that is explained by the model

Higher is better

ANOVA								
	df	SS	MS	F	Significance F			
Regression	1	62.10	62.10	52.29	0.00	-		
Residual	8	9.50	1.19					
Total	9	71.60				_		
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	17.58	1.21	14.54	0.00	14.79	20.37	14.79	20.37
X Variable 1	-1.37	0.19	-7.23	0.00	-1.81	-0.94	-1.81	-0.94

**Univariate Regression** 

**Ordinary Least Squares (OLS)** 

Assessing the Model

#### **EXCEL SUMMARY OUTPUT**

Regression Statistics						
Multiple R	0.93					
R Square	0.87					
Adjusted R Square	0.85					
Standard Error	1.09					
Observations	10.00					

#### Adjusted R Square

Adjusts the R<sup>2</sup> value to control for the number of variables included in the model

Adding variables will increase R<sup>2</sup>

	df	SS	MS	F	Significance F			
Regression	1	62.10	62.10	52.29	0.00			
Residual	8	9.50	1.19					
Total	9	71.60				_		
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	17.58	1.21	14.54	0.00	14.79	20.37	14.79	20.37
X Variable 1	-1.37	0.19	-7.23	0.00	-1.81	-0.94	-1.81	-0.94

**Ordinary Least Squares (OLS)** 

Assessing the Model

Regression Statistics					
Multiple R	0.93				
R Square	0.87				
Adjusted R Square	0.85				
Standard Error	1.09				
Observations	10.00				

**EXCEL SUMMARY OUTPUT** 

#### **Model Standard Error**



# Is the standard deviation of the model's residual errors

Is a measure of the amount of "noise" in the data

Upper 95.0%

20.37

-0.94

						_	
	df	SS	MS	F	Significance F		
Regression	1	62.10	62.10	52.29	0.00		
Residual	8	9.50	1.19				
Total	9	71.60				-	
	Coefficients S	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%
Intercept	17.58	1.21	14.54	0.00	14.79	20.37	14.79
X Variable 1	-1.37	0.19	-7.23	0.00	-1.81	-0.94	-1.81

Ordinary Least Squares (OLS)

### Assessing the Model

Regression Statistics						
Multiple R	0.93					
R Square	0.87					
Adjusted R Square	0.85					
Standard Error	1.09					
Observations	10.00					

**EXCEL SUMMARY OUTPUT** 

### **Coefficient Standard Error**



Is essentially a measure of the signal-to-noise

The larger the coefficient standard error, the worse the signal-to-noise, meaning, less precise

	df	SS	MS	F	Significance F
Regression	1	62.10	62.10	52.29	0.00
Residual	8	9.50	1.19		
Total	9	71.60			

	Coefficients Sta	ndard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	17.58	1.21	14.54	0.00	14.79	20.37	14.79	20.37
X Variable 1	-1.37	0.19	-7.23	0.00	-1.81	-0.94	-1.81	-0.94

Advanced methods of analysis Univariate Regression Ordinary Least Squares (OLS)

## Assessing the Model

EXCEL SUMMARY OUTPUT		t Statistic	Is the Coefficient divided by the standard error
Regression Sta	tistics		
Multiple R	0.93		Test of the hypothesis that the "true"
R Square	0.87		value of the coefficient is zero
Adjusted R Square	0.85		
Standard Error	1.09		Rule of thumb, looking for values with
Observations	10.00		absolute value $> 1.69$

	df	SS	MS	F	Significance F
Regression	1	62.10	62.10	52.29	0.00
Residual	8	9.50	1.19		
Total	9	71.60			

	Coefficients	Standard Error 🤇	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	17.58	1.21	14.54	0.00	14.79	20.37	14.79	20.37
X Variable 1	-1.37	0.19	-7.23	0.00	-1.81	-0.94	-1.81	-0.94

**Univariate Regression** 

**Ordinary Least Squares (OLS)** 

## Assessing the Model

EXCEL SUMMARY OU	ITPUT	P-Value	Also known as the significance level
Regression Sta	tistics		The probability that the result isn't reliable
Multiple R	0.93		
R Square	0.87		(1 - P-Value) is the % of times you'd you'd
Adjusted R Square	0.85		likely see this result in similar situations
Standard Error	1.09		incery see this result in similar situations
Observations	10.00		<b>.</b>
			Commonly looking for P-Values < 0.10

A	Ν	0	V	A
<i>'</i> '		$\sim$	•	•

ANOVA						_		
	df	SS	MS	F	Significance F	_		
Regression	1	62.10	62.10	52.29	0.00			
Residual	8	9.50	1.19					
Total	9	71.60				-		
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	17.58	1.21	14.54	0.00	14.79	20.37	14.79	20.37
X Variable 1	-1.37	0.19	-7.23	0.00	-1.81	-0.94	-1.81	-0.94

# Data Considerations / Challenges

# Advanced methods of analysis Data Transformations

We noted earlier that in the "real world" data are seldom "perfect"

OLS likes data to be as linear as possible...

Sometimes, it is appropriate to "transform" your data to make it more linear...



# Advanced methods of analysis Data Transformations

We noted earlier that in the "real world" data are seldom "perfect"

OLS likes data to be as linear as possible...

Sometimes, it is appropriate to "transform" your data to make it more linear...



## Advanced methods of analysis Data Transformations

Exercise

Input and Plot the three data sets below and perform the appropriate data transformation.

1				
Х	Y			
10.00	1			
5.00	2			
3.33	3			
2.50	4			
2.00	5			
1.67	6			
1.43	7			
1.25	8			
1.11	9			
1.00	10			
0.91	11			
0.83	12			
0.77	13			
0.71	14			
0.67	15			
0.63	16			
0.59	17			
0.56	18			
0.53	19			
0.50	20			

	2
х	Y
1.00	1
1.41	2
1.73	3
2.00	4
2.24	5
2.45	6
2.65	7
2.83	8
3.00	9
3.16	10
3.32	11
3.46	12
3.61	13
3.74	14
3.87	15
4.00	16
4.12	17
4.24	18
4.36	19
4.47	20

3			
х	Y		
1	1		
4	2		
9	3		
16	4		
25	5		
36	6		
49	7		
64	8		
81	9		
100	10		
121	11		
144	12		
169	13		
196	14		
225	15		
256	16		
289	17		
324	18		
361	19		
400	20		

- **Univariate Regression**
- **Ordinary Least Squares (OLS)**
- Output using a different Statistical Package (E-Views)

Dependent Variable: QUANTITY Method: Least Squares Date: 07/20/14 Time: 14:44 Sample: 1 10 Included observations: 10

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C PRICE	17.58055 -1.373860	1.209153 0.190000	14.53956 -7.230843	0.0000 0.0001
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.867297 0.850709 1.089812 9.501520 -13.93372 52.28509 0.000090	Mean depende S.D. dependen Akaike info crite Schwarz criterie Hannan-Quinn Durbin-Watson	nt var t var erion on criter. stat	9.200000 2.820559 3.186744 3.247261 3.120357 2.619326

	Price	Quantity
	Х	Y
Obs	\$	# Sold
1	10	4
2	5	9
3	6	10
4	3	14
5	6	8
6	6	10
7	8	6
8	5	12
9	7	9
10	5	10

### **Univariate Regression**

## Problems to watch out for - Heteroskedasticity

INCOME = f (AGE) (+)

Dependent Variable: INCOME Method: Least Squares Date: 07/20/14 Time: 17:20 Sample: 1 16 Included observations: 16

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C AGE	15063.69 762.1433	11228.48 242.0928	1.341561 3.148146	0.2011 0.0071
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.414491 0.372669 11117.67 1.73E+09 -170.6954 9.910821 0.007118	Mean depende S.D. dependen Akaike info crite Schwarz criterio Hannan-Quinn Durbin-Watson	nt var t var erion on criter. stat	49312.50 14036.71 21.58693 21.68350 21.59187 0.435679

Model results "look" pretty good

Descent R<sup>2</sup>

T-Stat > 1.69

P-Value < 0.10

But....

**Univariate Regression** 

Problems to watch out for - Heteroskedasticity

Error terms are not randomly distributed

Error terms are correlated the dependent variable (not independent)

INCOME = f (AGE) (+)



Scatter plot will typically look like a cone or funnel

Residual Error increase as Y (INCOME) and X (AGE) increases Not independent from each other

# **Time Series**

**Dice Exercise** 

Take 3 Dice Roll the dice 15 different times at 10 second intervals Create a data series in sequential order (1,2...15) of or our rolls Each data point is the sum of the values of the 3 dice 1<sup>st</sup> roll = Observation 1, 2<sup>nd</sup> roll = Observation 2, etc

Using Excel, plot our results

Discussion

## Data series over time



- Trends
- Patterns
- Changes / Shocks in either trends or patterns



Seasonality

Periodic, repetitive, and predictable movement around a trend



Seasonality

#### Trend

Long-term movement when things like seasonality have been filtered out





Observed breaks or deviations in seasonality and / or trend

### Time Series Data – US Ecommerce Sales

Shock

	Ecommerce			Recession
	Sales	Q4 Dummy	TIME	Dummy
Q1-2007	30,322	0	1	0
Q2-2007	31,493	0	2	0
Q3-2007	32,248	0	3	0
Q4-2007	42,063	1	4	1
Q1-2008	34,169	0	5	1
Q2-2008	34,170	0	6	1
Q3-2008	33,396	0	7	1
Q4-2008	39,498	1	8	1
Q1-2009	32,185	0	9	1
Q2-2009	32,789	0	10	1
Q3-2009	34,317	0	11	0
Q4-2009	45,617	1	12	0
Q1-2010	36,842	0	13	0
Q2-2010	38,230	0	14	0
Q3-2010	39,809	0	15	0
Q4-2010	54,014	1	16	0
Q1-2011	43,841	0	17	0
Q2-2011	44,948	0	18	0
Q3-2011	45,545	0	19	0
Q4-2011	63,549	1	20	0
Q1-2012	50,589	0	21	0
Q2-2012	51,285	0	22	0
Q3-2012	52,643	0	23	0
Q4-2012	72,361	1	24	0
Q1-2013	58,215	0	25	0
Q2-2013	60,498	0	26	0
Q3-2013	61,857	0	27	0
Q4-2013	83,709	1	28	0
Q1-2014	66,917	0	29	0

Seasonality	Created a "dummy variable" to control for the spike in Q4 each year		
Trend	Simple linear time variable that increments by one each period		

Created a "dummy variable" to control for the recession between Q4-2007 and Q2-2009



Let's do some regressions and discuss our results

Dependent Variable: ECOM\_SALES Method: Least Squares Date: 07/22/14 Time: 17:37 Sample: 2007Q1 2014Q1 Included observations: 29

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C Q4_DUMMY	43014.00 14244.71	2715.331 5526.788	15.84116 2.577395	0.0000 0.0157
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.197455 0.167731 12736.03 4.38E+09 -314.2266 6.642963 0.015736	Mean dependen S.D. dependen Akaike info crite Schwarz criterio Hannan-Quinn Durbin-Watson	nt var t var erion on criter. stat	46452.38 13960.54 21.80873 21.90303 21.83826 0.082249

Coefficient = 14,244 MM

Meaning that all things equal, Q4 is typically \$14.2B larger than the other quarters

t-Stat = 2.577 and P-Value = 0.01

OK results – definitely significant

R2 = 0.19

Not great, but we're only controlling for one piece of information, so not bad

Assessing the residual errors

Dependent Variable: ECOM\_SALES Method: Least Squares Date: 07/22/14 Time: 17:37 Sample: 2007Q1 2014Q1 Included observations: 29

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C Q4_DUMMY	43014.00 14244.71	2715.331 5526.788	15.84116 2.577395	0.0000 0.0157
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.197455 0.167731 12736.03 4.38E+09 -314.2266 6.642963 0.015736	Mean depender S.D. dependen Akaike info crite Schwarz criterio Hannan-Quinn Durbin-Watson	nt var t var erion on criter. stat	46452.38 13960.54 21.80873 21.90303 21.83826 0.082249



There is a time trend in our residual error (blue line) at the bottom of the graph

Sometimes called a "drift"

We want the blue line to essentially bounce (randomly) around the "Zero Line" with roughly equal amount of points above and below the zero line (Zero Mean)

#### Looking at only the time variable (ignoring the Q4 spikes)

Dependent Variable: ECOM\_SALES Method: Least Squares Date: 07/22/14 Time: 17:53 Sample: 2007Q1 2014Q1 Included observations: 29

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C TIME	24973.52 1431.924	2639.852 153.6981	9.460195 9.316470	0.0000 0.0000
R-squared	0.762735	Mean dependent var		46452.38
Adjusted R-squared	0.753947	S.D. dependent var		13960.54
S.E. of regression	6924.950	Akaike info criterion		20.59012
Sum squared resid	1.29E+09	Schwarz criterion		20.68442
Log likelihood	-296.5568	Hannan-Quinn criter.		20.61965
F-statistic	86.79662	Durbin-Watson stat		1.987252
Prob(F-statistic)	0.000000			



Coefficient = 1,431.9 MM

Meaning Ecommerce sales are growing by ~ \$1.4B each quarter (between 2007 – present)

t-Stat = 9.31 and P-Value = 0.00

OK Strong – definitely significant

R2 = 0.76

Pretty good!

Removed most of the drift in the residual error

See these periodic "shocks" which we know always happen in Q4

Looks like something was happening from 2007 – mid-2009

#### Looking at only the time variable (ignoring the Q4 spikes)

Dependent Variable: ECOM\_SALES Method: Least Squares Date: 07/22/14 Time: 17:53 Sample: 2007Q1 2014Q1 Included observations: 29

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C TIME	24973.52 1431.924	2639.852 153.6981	9.460195 9.316470	0.0000 0.0000
R-squared	0.762735	Mean dependent var		46452.38
Adjusted R-squared	0.753947	S.D. dependent var		13960.54
S.E. of regression	6924.950	Akaike info criterion		20.59012
Sum squared resid	1.29E+09	Schwarz criterion		20.68442
Log likelihood	-296.5568	Hannan-Quinn criter.		20.61965
F-statistic	86.79662	Durbin-Watson stat		1.987252
Prob(F-statistic)	0.000000			



Coefficient = 1,431.9 MM

Meaning Ecommerce sales are growing by ~ \$1.4B each quarter (between 2007 – present)

t-Stat = 9.31 and P-Value = 0.00

OK Strong – definitely significant

R2 = 0.76

Pretty good!

Removed most of the drift in the residual error

See these periodic "shocks" which we know always happen in Q4

Looks like something was happening from 2007 – mid-2009

Now also bringing in the Q4 Spikes into the model

Dependent Variable: ECOM\_SALES Method: Least Squares Date: 07/22/14 Time: 17:59 Sample: 2007Q1 2014Q1

Included observations: 29

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C TIME Q4_DUMMY	22619.21 1389.119 12413.60	1676.178 95.60087 1869.169	13.49452 14.53040 6.641243	0.0000 0.0000 0.0000
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(E-statistic)	0.912006 0.905237 4297.547 4.80E+08 -282.1740 134.7377 0.000000	Mean dependent var S.D. dependent var Akaike info criterion Schwarz criterion Hannan-Quinn criter. Durbin-Watson stat		46452.38 13960.54 19.66717 19.80862 19.71147 0.593286



Coefficients changed a little from our prior specifications

Meaning Ecommerce sales are growing by ~ \$1.4B each quarter (between 2007 – present)

t-Stats and P-Values are very strong definitely significant

R2 = 0.91

Getting better!

Removed most of the drift in the residual error – but it is still present

We are capturing the Q4 Seasonal Sales Spikes nicely

Looks like something was happening from 2007 – mid-2009

Recap

Analysis of data series over time

- Trends
- Patterns
- Changes / Shocks in either trends or patterns

Time Trend

Seasonality Dummy Variable

Semi- A Theoretical Approach We've Addressed What – Not Why

## Advanced methods of analysis Time Series Regression

Plot the data

Over time and vs other data with possible relationships What patterns / relationships do you observe? Are there recurring patterns

Run a regression with "time" as the independent variable What patterns / trends (if any) do you see in the residuals

Discussion

# Advanced methods of analysis Multivariate Time Series Regression

ACME Widget Company

- 52 Weeks of Data
- Units Sold
- Net Price of Widgets Sold
- Supply of Widgets on Hand
- Knowledge of Promotional Event Timing
- Major Weather Event Mid Year Caused Significant Disruption
- 1) How much does Price affect Units Sold?
- 2) Does my supply on hand impact sales?
- 3) Are my Promotional Events working?
- 4) How much did the weather event cost me?
- 5) Anything else I need to be aware of?

# Advanced methods of analysis Multivariate Time Series Regression

ACME Widget Company

Answer questions by estimating a weekly demand model

#### Steps:

- 1) Look at (plot) the data
- 2) Look at Correlation Matrix
- 3) Iteratively build the Weekly Demand Model
- 4) Assess Output with each iteration
- 5) Review the residual errors of each iteration

UNITS = 
$$f$$
 (Price, etc...)

**Multivariate Time Series Regression** 


# Advanced methods of analysis Multivariate Time Series Regression



#### Independent Variable Correlation Matrix

	PRICE	SUPPLY	PROMO_EVENT	BAD_WEATHER
PRICE	1.00	-0.44	-0.15	-0.18
SUPPLY	-0.44	1.00	-0.01	0.03
PROMO_EVENT	-0.15	-0.01	1.00	-0.07
BAD_WEATHER	-0.18	0.03	-0.07	1.00

#### Advanced methods of analysis Multivariate Time Series Regression

#### **Final Model Output**

Dependent Variable: UNITS\_SOLD Method: Least Squares Date: 07/23/14 Time: 12:59 Sample: 1/01/2014 12/24/2014 Included observations: 52

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C PRICE SUPPLY PROMO_EVENT BAD_WEATHER TIME	901.3092 -18.47527 0.101883 175.7728 -161.6135 5.536277	215.9840 2.480796 0.015892 34.08523 56.72894 0.793918	4.173037 -7.447315 6.411067 5.156860 -2.848872 6.973357	0.0001 0.0000 0.0000 0.0000 0.0065 0.0000
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.816102 0.796113 76.34355 268103.5 -296.0298 40.82783 0.000000	Mean depende S.D. dependen Akaike info crit Schwarz criteri Hannan-Quinn Durbin-Watson	ent var it var erion on criter. i stat	729.3269 169.0745 11.61653 11.84167 11.70285 1.372811

#### Residual Error Structure





# **Multivariate Regression**

Regression analysis with more than one variable

 $\mathsf{Y}=f(\mathsf{X}_1,\,\mathsf{X}_2,\,\ldots)$ 

### Quantity = f(Price, Income, Price of Substitutes, etc.)

Can be "cross section" or "time series" data

#### Revisiting our model that had a Heteroskedasticity problem

INCOME = f (AGE) (+)

Dependent Variable: INCOME Method: Least Squares Date: 07/20/14 Time: 17:20 Sample: 1 16 Included observations: 16

Variable	Coefficient	Std. Error	t-Statistic	Prob.	\$80,000							
C AGE	15063.69 762.1433	11228.48 242.0928	1.341561 3.148146	0.2011 0.0071	\$70,000 - \$60,000 - E S \$50,000 -					-		
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.414491 0.372669 11117.67 1.73E+09 -170.6954 9.910821 0.007118	Mean depende S.D. dependen Akaike info crite Schwarz criterie Hannan-Quinn Durbin-Watson	nt var t var erion on criter. stat	49312.50 14036.71 21.58693 21.68350 21.59187 0.435679	LI \$40,000 - \$30,000 - \$20,000 - \$10,000 - \$- 0	10	20	30 AGE	40	50	60	

Is often a sign that something important is missing from the model,7

#### Revisiting our model that had a Heteroskedasti

Revisiting our model that had a Heteroskedasticity problem

INCOME = f (AGE, Education, Gender) (+) (+) (+ / -)

Obs	INCOME	AGE	GENDER	EDUCATION
1	30,000	29	0	14
2	32,000	27	1	16
3	33,000	35	0	16
4	37,000	49	0	12
5	40,000	33	1	14
6	42,000	56	0	12
7	44,000	52	0	18
8	46,000	35	1	12
9	48,000	62	0	14
10	50,000	37	1	12
11	55,000	41	1	14
12	59,000	45	1	16
13	61,000	42	1	20
14	65,000	53	1	18
15	72,000	59	1	18
16	75,000	64	1	14

Dataset has 16 people in it

Gender Variable is either a "1" or "0" 1 = Male 0 = Female This is called a "dummy" variable

Education is number of years of school completed

- 12 = High School
- 14 = Associates Degree
- 16 = Bachelors Degree
- 18 = Masters Degree
- 20 = PhD or Professional Degree (i.e., Law)

INCOME = f (AGE, Education, Gender) (+) (+) (+ / -)



Plot out the data

Each univariate plot and excel trend line regression suggests there is a positive relationship between the independent and dependent variables

What about the relationship between the independent variables?

INCOME = f (AGE, Education, Gender) (+) (+) (+ / -)

Look at the Correlation Matrix

	AGE	EDUCATION	GENDER
AGE	1.00	0.11	-0.15
EDUCATION	0.11	1.00	0.21
GENDER	-0.15	0.21	1.00

Correlation is the measure the strength of association (relationship) between two variables

A high correlation coefficient indicates two variables are not (or are not sufficiently) independent from one another

- High correlation of the independent variables gives rise to Multicolinearity

Loose rule of thumb to test for model inclusion is  $\sim < 0.75$ 

INCOME = f (AGE, Education, Gender) (+) (+) (+ / -)

Dependent Variable: INCOME Method: Least Squares Date: 07/21/14 Time: 13:41 Sample: 1 16 Included observations: 16

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C AGE EDUCATION GENDER	-13864.13 860.5635 856.3410 18655.91	7989.759 101.1059 479.2658 2438.227	-1.735238 8.511506 1.786777 7.651426	0.1083 0.0000 0.0992 0.0000
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.916257 0.895322 4541.444 2.47E+08 -155.1376 43.76531 0.000001	Mean dependent var S.D. dependent var Akaike info criterion Schwarz criterion Hannan-Quinn criter. Durbin-Watson stat		49312.50 14036.71 19.89220 20.08534 19.90209 0.681418

#### **Regression Checklist**

- R2 and Adj R2 are high and "close" to one another
- Coefficients all have the expected signs
  - t-Statistics are all > 1.69 (Absolute value)
  - P-Values are all 0.10 or smaller

What about the residual errors?

INCOME = f (AGE, Education, Gender) (+) (+) (+ / -)

#### Evaluating the model residual error structure

OBS	INCOME	Estimate	Residual	
1	30,000	23,081	6,919	
2	32,000	41,728	-9,728	
3	33,000	29,957	3,043	
4	37,000	38,580	-1,580	
5	40,000	45,179	-5,179	
6	42,000	44,604	-2,604	
7	44,000	46,299	-2,299	
8	46,000	45,188	812	
9	48,000	51,480	-3,480	
10	50,000	46,909	3,091	
11	55,000	52,064	2,936	
12	59,000	57,219	1,781	
13	61,000	58,062	2,938	
14	65,000	65,816	-816	
15	72,000	70,979	1,021	
16	75,000	71,857	3,143	
				_

0.00 Avg



Correl Est vs. Resid 0.00000

#### A word on Multicollinearity

We noted earlier that this is a situation where two or more of the independent variables are highly correlated with one another

 Imagine two people in a hallway yelling something really important at the top of their lungs – hard to understand what is being said

Let's introduce another variable (average time in years in their current job)

Dependent Variable: INCOME Method: Least Squares Date: 07/21/14 Time: 14:33 Sample: 1 16 Included observations: 16

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C AGE YRS IN JOB EDUCATION GENDER	-88501.46 -868.7334 <del>&lt;</del> 23030.39 809.4847 18579.48	75280.57 1737.290 23097.61 481.6798 2440.022	-1.175622 -0.500051 0.997090 1.680545 7.614472	0.2646 0.6269 0.3401 0.1210 0.0000
R-squared Adjusted R-squared S.E. of regression Sum squared resid Log likelihood F-statistic Prob(F-statistic)	0.923199 0.895271 4542.544 2.27E+08 -154.4453 33.05663 0.000005	Mean depende S.D. dependen Akaike info crite Schwarz criteri Hannan-Quinn Durbin-Watson	nt var t var erion on criter. stat	49312.50 14036.71 19.93067 20.17210 19.94303 0.924386

- Notice that the sign on AGE changed
- Also AGE is no longer significant Std. Error (Coefficient) went way up!
- AGE and YRS IN JOB are 0.998 correlated

Rotate variables in and out to see which specification yields best overall fit