

LECTURE 17: MULTIVARIABLE REGRESSIONS I

- I. What Determines a House's Price?
 - a. Open Data Set 6 to help us answer this question. You'll see pricing data for homes based on when they were built, how big each home is, how far it is from the city center, and how many days it was on the market before being sold.
 - i. I don't remember where I got this data from. I'm pretty confident it's real but I doubt it's for our area.
 - b. Suppose you're researching how home prices change as you get closer to a city's downtown area. You'd suspect that homes should get cheaper as you go further from the city.
 - c. Here's a regression output (n=100) with miles from city center causing housing prices:

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	508457.7719	57163.32439	8.894825	3.02137E-14	395019.02	621896.5288	395019	621896.5
miles	-5517.997542	11504.88918	-0.47962	0.632564886	-28349.08	17313.08059	-28349.1	17313.08

- d. While the coefficient is negative (as expected: more miles means a lower price) the result is ***not statistically significant***. Location, location, location...doesn't matter?
- e. That can't be right. And it's not. The problem with this analysis is as homes get farther out, they get bigger.
 - i. We asked the question, "If you buy a home farther from the city, what happens to the price?"
 - ii. We need to ask: "If you buy an ***identical*** home farther from the city, what happens to the price?"
- f. While it's hard to get data so we can compare "identical" homes, we can get data on one of the big variables here: size. Both size (in square feet) and distance from city center (in miles) matter for housing prices. So we turn to a multivariate regression.
- g. Excluding an important variable can distort the regression analysis, resulting in ***omitted variable bias***. It's when a variable that's correlated with the dependent variable and at least one independent variable is not included in the regression.

- i. In our example, size was correlated both distance and price. Without size, we got a distorted understanding of what was going on. We were missing an important control.

II. Basics

- a. A multivariate regression has more than one explanatory variable.

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \varepsilon_i$$

- i. Note that the explanatory variables now have a subscript to differentiate them from other explanatory variables.
- b. You want to do a multiple regression because you think multiple variables matter.
 - i. Example: Life expectancy depends on both diet and exercise.
 - ii. Example: Sales depends on price, the unemployment rate, advertising, and so on.
- c. When you interpret a particular beta value, it is now the change in the dependent variable for every unit change in the corresponding explanatory variable, **holding all other explanatory variables constant**.
- d. Here's the housing regression, now with size and location predicting price (remember when I suggested you use labels? This is why):

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	78.70913492	82608.5323	0.000953	0.999241735	-163876.4	164033.7786	-163876	164033.8
sqft	236.9646693	32.00587304	7.403787	4.84964E-11	173.44187	300.4874676	173.4419	300.4875
miles	-23792.47937	9567.375458	-2.48683	0.014596448	-42781.07	-4803.887471	-42781.1	-4803.89

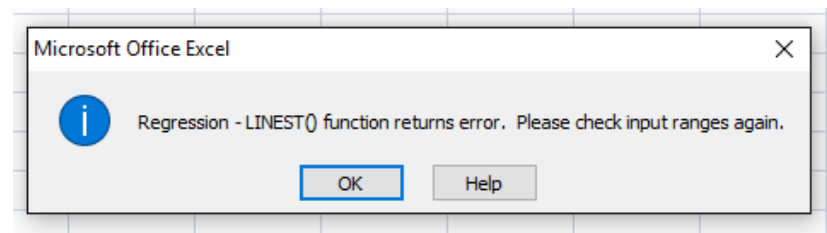
- i. Now distance (and size!) are statistically significant.
- ii. For every additional square foot a house has the price increases by \$236.96, holding distance from the city center constant.
- iii. For every additional mile the house is from the city center the price decreases by \$23,792.48, holding size constant.

III. Preparing Your Data

- a. Excel requires that all explanatory variables for a regression are next to each other. Suppose, for example, I'm interested in how age of 1st marriage, population density, and median age affect the murder rate.
 - i. The easiest way to do this is to right click the column with the variable you're interested in, select "Cut", right click the column of another variable you're interested in, and select "Insert Cut Cells." Like this:

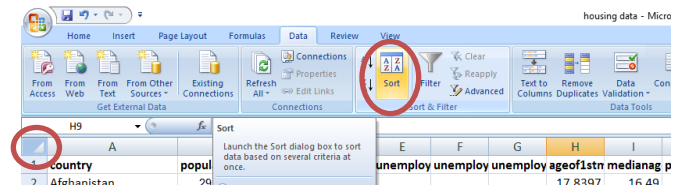
	country	population	unemployment	age of 1st median age	murder rate
547	6.49	77.92	3732	58.1237	17.8397
547	532	57.909	3267	39.3066	23.3265
599	015	23.3662	8921	17.3029	29.6
391	718	18.5	2651	38.2979	16.718
482	9.12	28.362	8496	29.5455	188.991
387	24.3474	24.3474	5377	49.1104	13.93
297	31.142	59.8532	2.35225	56.3553	22.986
508	3.23984	3.23984	40.294	31.142	102.85
547	25.4	12.513	36.417	11.1111	561.311
397	3.99286	63.0269	-32.219	36.553	2.635
547	2.03918	58.0308	48.231	39.953	16.3782
397	16.1394	8.21798	40.352	27.08	23.8863
152	5.05683	25	24.7	28.054	19.5109
271	11.1306	21.3037	26.024	27.617	23.182
458	5.07472	0.24164	23.88	22.544	62.6047

- ii. Now all my explanatory variables are next to each other.
- b. Excel requires that all variables have no blank observations. If you get this error message:

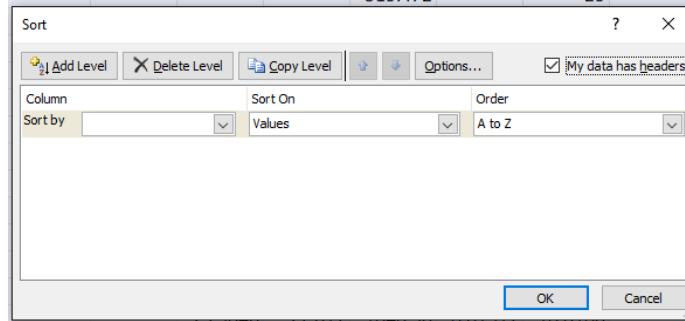


It means you are trying to run a regression using variables with missing values.

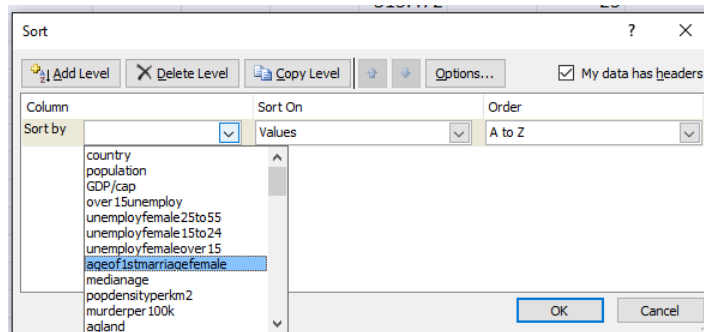
- i. Normally, a program would just ignore those observations.
- ii. But Excel is kind of dumb. You have to delete them. So let me first show you a quick way to do that.
- c. The Sort function is in the Data tab. Highlight the whole Excel sheet (by clicking in the upper right-hand corner of sheet) and select Data.



- d. You'll get something that looks like this:



- i. Make sure to select “My data has headers.” It’ll make this a lot easier.
- e. In the dropdown menu, select ageof1stmariagefemale. Then press OK.



- f. Excel will reorder the data based on that variable. This means all the blank values end up in the same place: at the end. This makes it a lot easier to find and delete all the observations with blank values.

1	country	population	GDP/cap	over15un	unemploy	unemploy	unemploy	ageof1stn	medianag	popdensit	murderpe	agland	aidgivena:	aidrecieve	alca
173	Sweden	9,041,000	31995	7.8	6.3	21.4	7.6	32.405	40.103	20.149	0.89648	7.8374		0.94	
174	Jamaica	2,651,000	7132					33.2029	25.467	242.715	42.5942	43.121			0.44
175	Martinique	396,000	14627.1					33.2695	36.709	361.279	4.5				
176	American Samoa	185,000	9617.82							315.472		25			
177	Andorra	67,000	39002.4							170.459	0.72651	38.2979			
178	Angola	15,941,000	3533						16.718	13.329	48.2062	46.194			1.59
179	Anguilla	12,000	19478.9							150.198					
180	Antigua and Barbuda	81,000	14579							188.991	7.48496	29.5455			0.95
181	Aruba	99,000	26762.7						36.417	561.311		11.1111			
182	Bermuda	64,000	69916.8							1210.83		20			
183	Bosnia and Herzegovina	3,907,000	6506						37.332	73.857	1.82565	42.2941			4.7
184	British Virgin Islands	22,000	44961.2							145.781					
185	Caribbean								28.033	173.12					
186	Cayman Islands	45,000	48632.3							199.182		8.33333			
187	Channel Islands	149,000							40.05	762.303		36.8421			
188	Congo, Dem. Rep.	57,549,000	330						26.079	25.194	45.129	9.90274			26.35
189	Congo, Rep.	3,999,000	3621						38.943	9.99	23.4758	30.8697			31.69
190	Cook Islands		18482.1							80.864	0.84027				
191	Cuba	11,269,000	7407.24						5.488	100.966	5.78205	62.3544			
192	Djibouti		1964						20.077	34.696	3.77898	73.3822			9.81
193	Dominica	79,000	8576							89.822	9.84239	30.6667			7.64
194	Equatorial Guinea	504,000	11999						18.583	21.704	28.2999	11.5508			1.09
195	Faeroe Islands	47,000	39495.7							34.738		2.15827			
196	Falkland Islands (Malvinas)	3,000	26101.6							0.244					
197	Gibraltar	28,000	40734.2							5119					

- i. *This is an incredibly useful function for your everyday understanding of data.* It makes it easy to, for example, find the

largest values or put all observations of the same category next to each other.

- ii. Remember when we used RMP data and I had you analyze ratings for different disciplines? I got all the disciplines next to each other using the Sort function.
 - iii. You can also add “levels” (see window on the previous page). This will tell Excel to sort within categories. For example, I could have it sort by discipline and then sort by number of ratings. Within each discipline, the professor with the most (or, if I choose, fewest) ratings would be listed first.
- g. Highlight rows starting in 176 all the way down to 237. Right click and select Delete.

1	country	population	GDP/cap	over15un	unemploy	unemploy	unemploy	ageof1stn	medianag	popdensit	murderpe	agland	aidgivena	aidreceiv	alca
213	Marshall Islands	62,000	6206							313.37	1.74394	77.7778			31.42
214	Mayotte	160,265	9617.82						18.549	466.479		54.0541			
215	Melanesia								20.219	14.556					
216	Micronesia, Fed. Sts.	110,000	5508						19.784	155.862	0.78562	31.4286			41.54
217	Montserrat	4,000	11579.6							55.176					
218	Nauru	14,000	6933.94							481.476	12.579				
219	Niue	1,000	5630.64							6.323	1.01416				
220	Northern Mariana Islan	81,000	9617.82							172.847		6.52174			
221	Palau	20,000	13012							43.845	0.84483	10.8696			
222	Palau		67							10					
223	Palau	5,000	2566.03							38.639					
224	Saint Kitts and Nevis	43,000	13677							188.268	11.3276	19.2308			0.64
225	Samoa	6,000	6859.54							25.401					
226	San Marino	8,000	41590							495.787		16.6667			
227	Seychelles	1,000	14202							181.613	3.22677	8.69565			2.19
228	Sri Lanka	8,000	932.962						17.645	13.101	1.85063	70.7384			
229	Sri Lanka	9,000	7234						26.143	3.051	10.5511	0.46154			2.53
230	Sri Lanka	8,000	4059						20.566	103.259	3.06892	75.3282			0.28
231	Sri Lanka	7,000	2203						16.731	66.667	16.2515	25.8911			26.82
232	Sri Lanka	1,000	889.433							101.083					
233	Sri Lanka	8,000	31209.1							70.995		1.05263			
234	Sri Lanka	0,000	4978.91							375.462	1.92412	33.3333			
235	Vietnam	84,238,000	2142						25.572	253.473	4.34193	32.4249			3.66
236	Wallis et Futuna	15,000	3612.17							74.58					
237	Western Sahara	341,000							24.221	1.656					

- h. Repeat this process for each variable that you care about (including your dependent variable) and you’re ready to run the regression.

IV. Dummy variables

- a. A common control is a *dummy variable*—a variable that’s either zero (for “no”) or one (for “yes”).

- i. These variables are binomial: gender (male or female), employment (working or not working), immigration status (legal or illegal).

Company	West?	Midwest?	Northeast?
Red Sun	1	0	0
Yellow Sun	0	0	0
Blue Sun	0	0	1
Green Sun	1	0	0
Orange Sun	0	1	0
Purple Sun	0	0	0
Black Sun	0	0	1
White Sun	0	0	0
Grey Sun	0	1	0

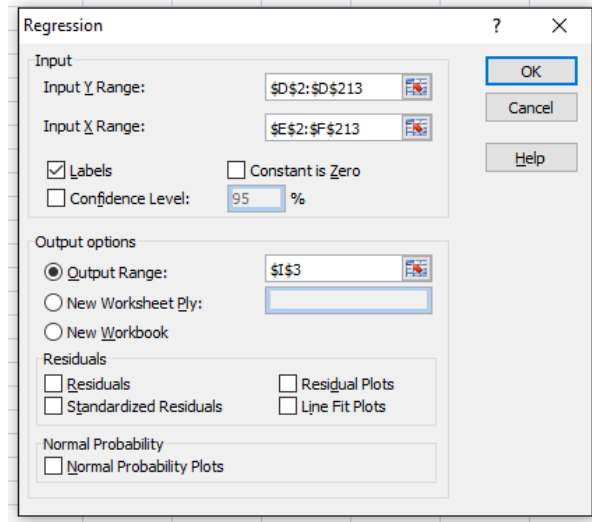
- ii. You can use multiple dummies for a variable with a few categories (White? Black? Asian? Hispanic?). For example, here’s hypothetical data where each observation is a

U.S. company. The dummy variable is the region of country where the company's headquarters are.

- iii. You typically want to have a number of dummies equal to one minus the number of categories. If the dummy is "Female?" then you know $1=F$ and $0=M$. Adding "Male?" is redundant. Note on the table of the hypothetical firms, there is no dummy variable for the South. That's because if a U.S. firm doesn't have their HQ in any of the other regions, it must have it in the South. That's where Yellow Sun, Purple Sun, and White Sun have their HQs.
 - iv. The only time you don't want to have one fewer dummy variables than categories is when the categories aren't mutually exclusive. A firm can't have their HQ in two different regions. But a student can have more than one major, a person can identify as multiple races, a rug can have several different colors in it, etc.
- b. You interpret the variable as you would when there's a single variable: examine the coefficient. Again, you're holding the other variables constant.

V. More Output from Excel

- a. Let's go back to the Rate My Professor ratings data in Data Set 5. Recall we explored how a professor's easiness can predict his or her quality.
- b. Rate My Professor also asks students to indicate if the professor is attractive or not (hot or not). I've set this up as a dummy variable: 1 means the professor is rated as hot and 0 means the professor is rated as not hot.
- c. If a professor becomes "hot," is it possible that results in a better quality? We need a plausible causation story (remember: regressions are all about causation). Perhaps students pay more attention and are more likely to attend class if the professor is attractive. That means students learn more and the class is more enjoyable, encouraging students to think the professor is a better educator.
- d. To run a regression with multiple explanatory variables, you just highlight multiple columns for the X range rather than just one column. I do below, highlighting the E and F columns:



- i. This is why all your dependent variables have to be next to each other: so you can create a continuous box.
- e. Here is the full output:

Regression Statistics	
Multiple R	0.806814419
R Square	0.650949506
Adjusted R Square	0.647593251
Standard Error	0.518875933
Observations	211

ANOVA					
	df	SS	MS	F	Significance F
Regression	2	104.435809	52.2179	193.9512	2.88655E-48
Residual	208	56.0003047	0.269232		
Total	210	160.4361137			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	5.756302906	0.153848027	37.41551	2.51E-94	5.453001574	6.059604238	5.453001574	6.059604238
DIFFICULTY	-0.754138639	0.049301502	-15.2965	6.84E-36	-0.85133333	-0.656943949	-0.85133333	-0.656943949
Hot?	0.552312714	0.086112722	6.413834	9.39E-10	0.382547109	0.722078319	0.382547109	0.722078319

These are the items we will focus on. The rest we've already discussed or don't matter for our purposes. Well, expect observations but it's obvious what that is.

- f. If a professor simply becomes “Hot” (going from a 0 to a 1), his or her rating increases by about 0.55, holding their DIFFICULTY rating constant. Note this is the most a professor could get out of this variable because there're only two values this variable can be.

VI. Interpretation

- a. *Explained (Regression) Sum of Squares (ESS)*—the squared vertical difference between the average and the predicted value of the dependent variable. This difference is taken for each observation and then added together.

- b. *Residual Sum of Squares (RSS)*—The squared vertical difference between the observed value and the predicted value. This difference is taken for each observation and then added together.
- c. *Total Sum of Squares (TSS)*—ESS + RSS
- d. R^2 —ESS/TSS, or the percent of deviation that our regression explains. There is no threshold for a “good” R^2 .
 - i. We are explaining 65% of the distance between a rating’s observed value and the average rating.
 - ii. R^2 is sometimes also called the “coefficient of determination.”
- e. *Adjusted R^2* —The R^2 value adjusted for the number of explanatory variables.
 - i. A weakness of R^2 is that it adding additional explanatory variables causes it to increase, regardless of the quality of explanatory variables. This is a problem because having many explanations for something is the same as having few.
 - ii. Adjusted R^2 penalizes the researcher for adding explanations, especially if it’s large relative to the number of observations. The equation is:

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

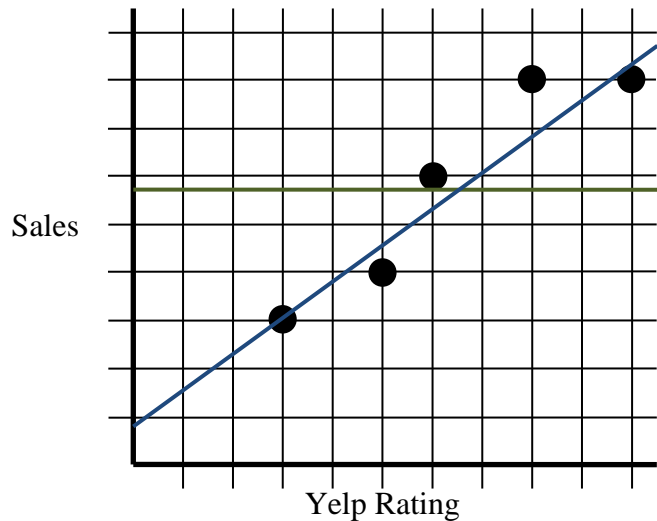
Where n is the number of observations and k is the number of explanatory variables, excluding the intercept.

- f. *F*—The ratio between the explained and unexplained variance. Like R^2 , it’s used for evaluating the model as a whole. And like the t distribution, the F distribution is a family of distributions. Significance level depends on degrees of freedom.
 - i. Higher values of F indicate a model with more explanatory power. Because the shape of the F distribution is known (its exact shape changes based on the number of observations and number of explanatory variables), it is possible to determine critical values.
- g. *Significance F*—this is the p-value for the F stat and uses the same criteria. If the value is very small, the model is quite good.

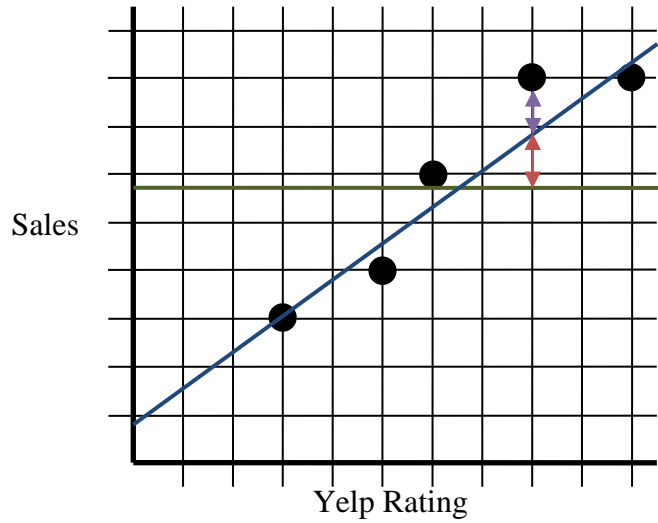
VII. Thinking About Regressions

- a. Suppose you have sales data on various Chinese restaurants. If you pick a restaurant at random, what do you suppose that restaurant’s sales are?

- b. Your best guess would be average sales. Obviously, your guess probably won't be right but based on how little information you have, there's no better guess.
- c. Now suppose you know that restaurant you chose has 4 out of 5 stars on Yelp, the popular review site. How do you adjust your expected sales? It should go up, right?
- d. Regressions are about how you can explain why an observation's value is different from the average (that's why causation is so important).



- e. The green line is the average sales. The blue line is the regression line. Note that we get a much better estimation of sales if we employ something we know that has predictive power (Yelp ratings) than if we just guessed based on the average.
 - i. Indeed, of the five observations, four give us a much better estimate of sales than the average (one is spot on!). Only one observation—the middle one—does using the line rather than the average worsen the guess. And it's not that much worse.



- f. The red line is that observation's contribution to ESS; it's the part of the deviation the regression line can explain.
- g. The purple line is that observation's contribution to RSS; it's the part of the deviation the regression line can't explain.
- h. I write "contribution" in each of these cases because ESS and RSS are the *sum* of squares. It's the result (after squaring it) from all the observations.