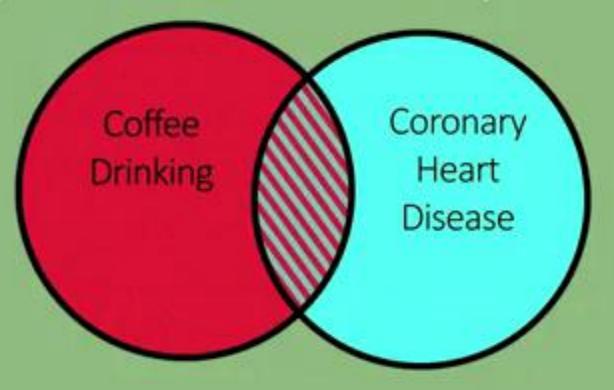
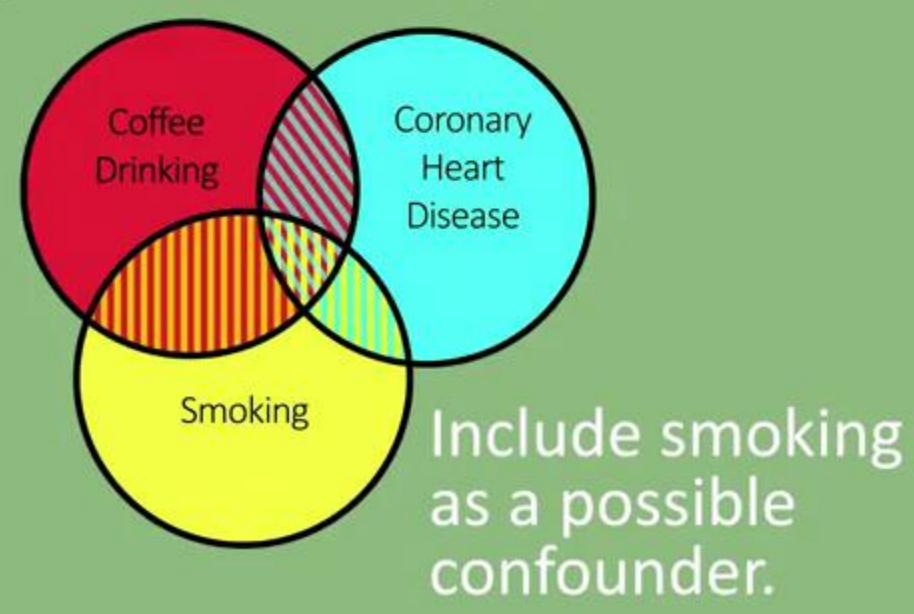
Module 2 Lesson 1 - More on Confounding



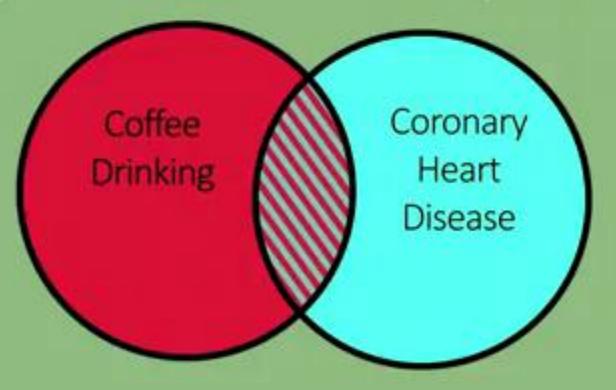
© Creative Commons, 2015

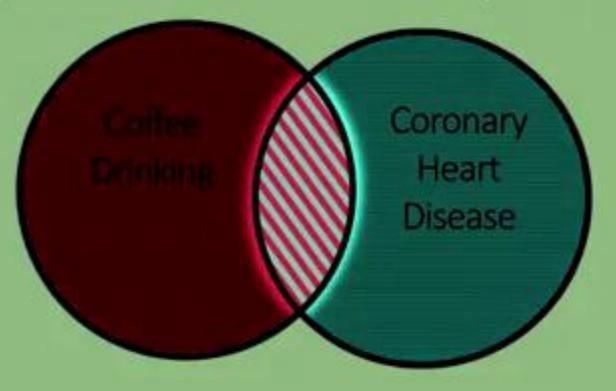




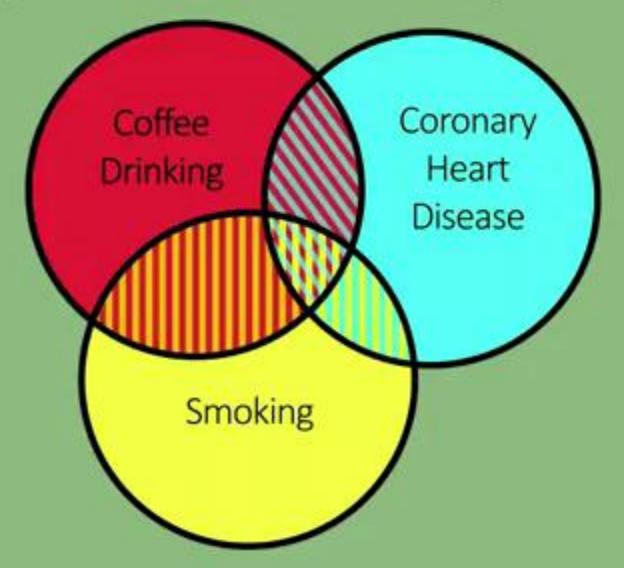


- -Confounder
- -Control Variable
- -Covariate
- -Third Variable
- -Lurking Variable



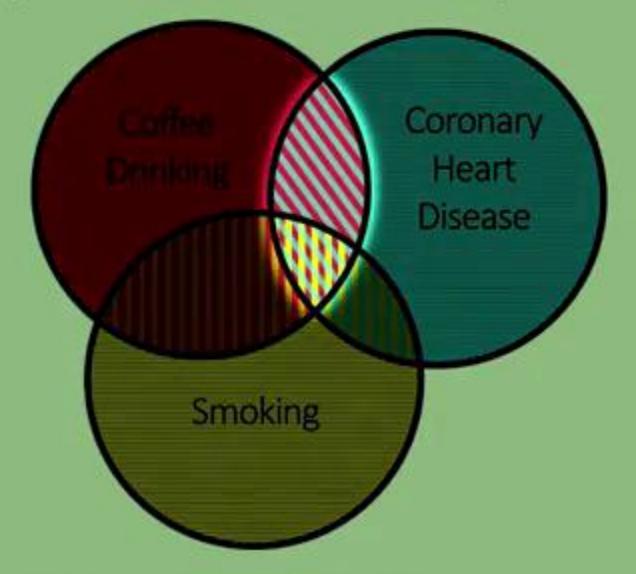


Response Variable



Confounding Variable

Response Variable



Confounding Variable

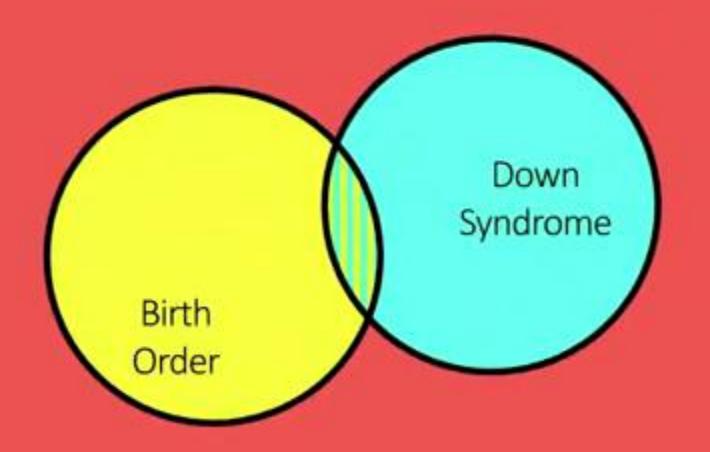
Two Types of Multivariate Model

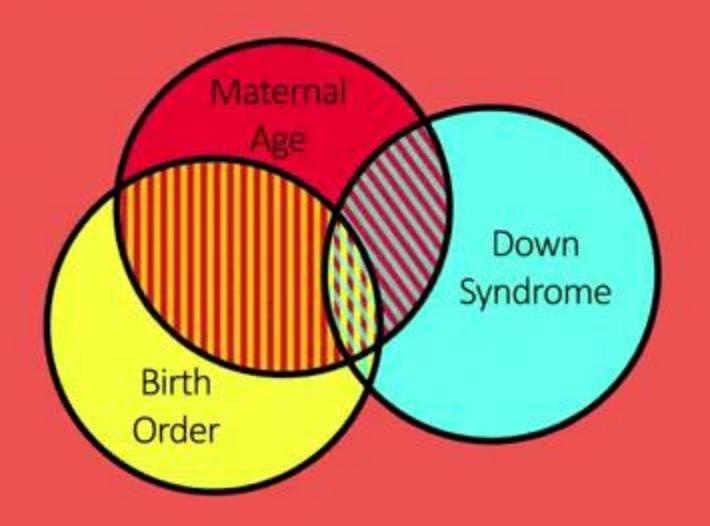
MULTIPLE
REGRESSION

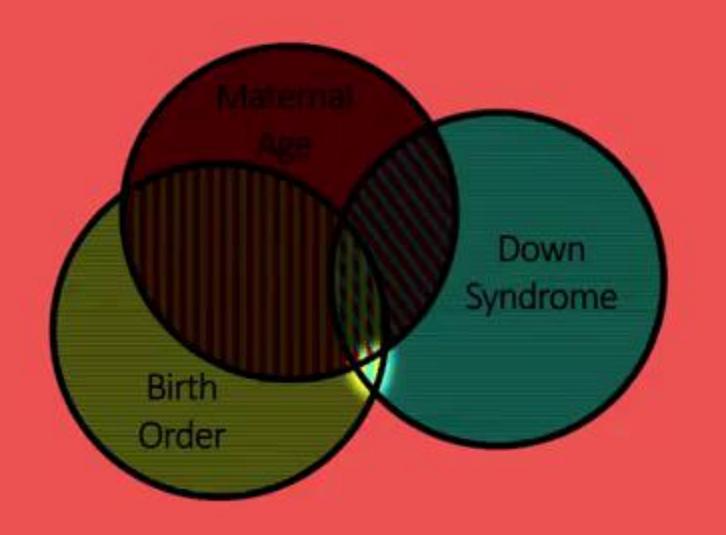
Quantitative Response Variable



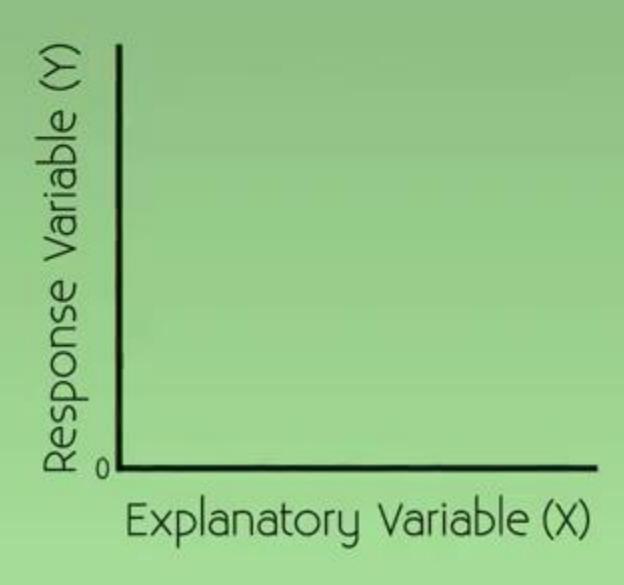
Two Types of Multivariate Model LOGISTIC MULTIPLE REGRESSION REGRESSION Quantitative Binary Response Response Variable Variable

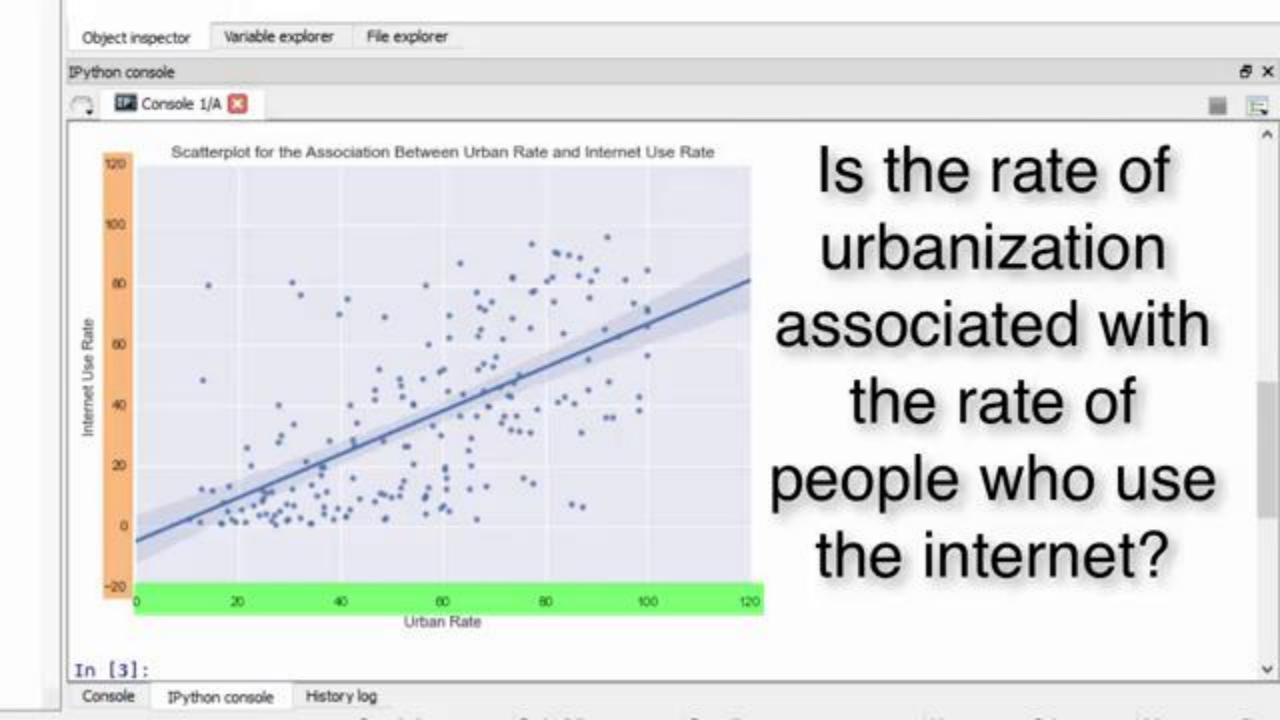


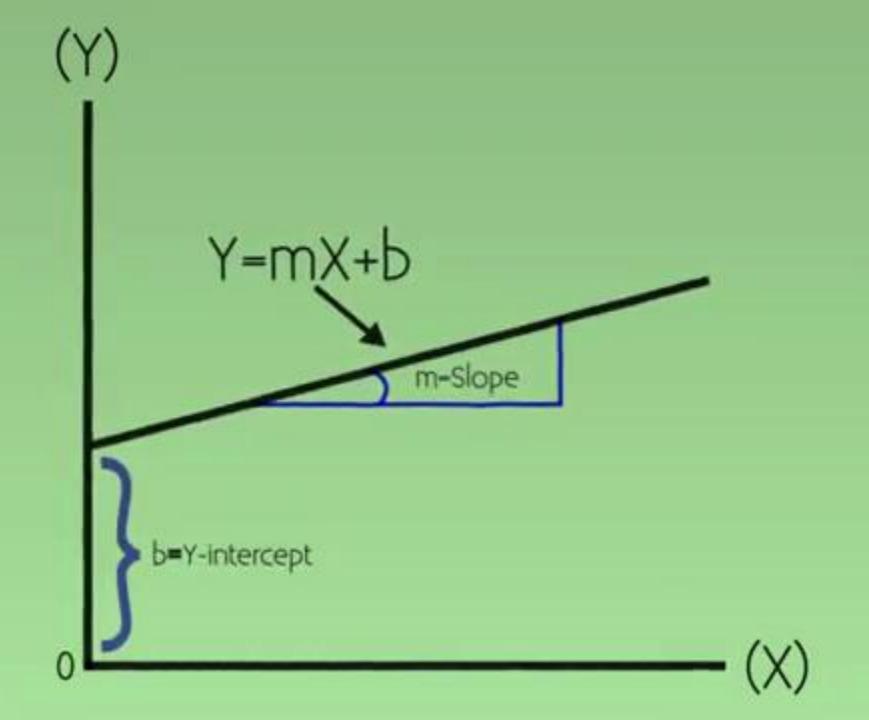












Module 2 Lesson 2 - Testing a Basic Linear Regression Model with Python



© Creative Commons, 2015

```
modelname = smf.ols(formula='QUANT_RESPONSE ~ QUANT_EXPLANATORY',
data=dataframe).fit()
print (modelname.summary())
```

```
GAPMINDER_polynomial_regression_and_diagnostic_plots.py
                                                       NESARC_data_linear_regression_modeling.py
 1# - - coding: utf-8 - --
 3 Created on Wed Oct 14 14:49:10 2015
 5 Mauthor: jml
 8 import pandas
 9 import numpy
10 import statsmodels.api as sm
11 import statsmodels.formula.api as smf
13 data = pandas.read csv('gapminder.csv')
15 # numeric variables that are read into python
16 # from the csv file as strings (objects) with empty cells should be
17 # converted back to numeric format using convert objects function
18 data['internetuserate'] = data['internetuserate'].convert objects(convert numeric=True)
19 data['urbanrate'] = data['urbanrate'].convert_objects(convert_numeric=True)
21 print('OLS regression model for the association between urbanrate and internet use rate')
22 reg1 = smf.ols('internetuserate ~ urbanrate', data=data).fit()
23 print (regl.summary())
```

14

24

	- Married		_
73.1	10.00	Console 1/A	•
4		COLISCIE TAV E	ю.

OLS Regression Results

			3				
Dep. Variabl	e:	internetus	erate	>R-sau	uared:		0.377
Model:			OLS		R-squared:		0.374
Method:		Least Sq	uares	48	atistic:		113.7
Date: Time:		Wed, 04 Nov 2015 15:21:51		Prob (F-statistic): Log-Likelihood:		4.56e-21 -856.14	
Df Residuals	:		188	BIC:			1723.
Df Model:			1				
Covariance Type:		nonr	obust				
	coe	f std err		t	P> t	[95.0% Conf.	Int.]
Intercept	-4.903	7 4.115	-1	1.192	0.235	-13.021	3.213
urbanrate	0.720		10	0.665	0.000	0.587	0.853
Omnibus:		1	0.750	Durb	in-Watson:		2.097
Prob(Omnibus	:):		0.005	Jarqu	ue-Bera (JB):		10.990
Skew:			0.574	Prob	(JB):	6	.00411
Kurtosis:			3.262	Cond.	No.		157.

Warnings: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

In [2]:

03	Console 1/A	×
		-

		OLS Reg	ressi	on Res	ults		
Dep. Variabl Model: Method: Date: Time: No. Observat Df Residuals Df Model: Covariance T	We ions:	Least Squar d, 04 Nov 20 15:21:	es 15 51 90 88	F-stat	red: -squared: istic: F-statistic): kelihood:	4	0.377 0.374 113.7 .56e-21 -856.14 1716. 1723.
	coef	std err		t	P> t	[95.0% Conf	. Int.]
Intercept urbanrate	-4.9037 0.7202	4.115 0.068	-1. 10.		0.235 0.000	-13.021 0.587	3.213 0.853
Omnibus: Prob(Omnibus Skew: Kurtosis:):	10.7 0.6 0.5 3.2	05 74		(T)(E)(O)(2.097 10.990 0.00411 157.

Warnings: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

In [2]:

Console :	L/A	×
-----------	-----	---

		OLS Regr	ession Re	sults		
Dep. Variable Model: Method: Date: Time: No. Observat Df Residuals Df Model: Covariance	tions:	internetuserat OL Least Square led, 04 Nov 201 15:21:5 19 18	S Adj. s F-sta 5 Prob 1 Log-L 0 AIC: 8 BIC:	uared: R-squared: stistic: (F-statistic) ikelihood:		0.377 0.374 113.7 4.56e-21 -856.14 1716. 1723.
	coef	std err	t	P> t	[95.0% Con	f. Int.]
Intercept urbanrate	-4.9037 0.7202	4.115 0.068	-1.192 10.665	0.235	-13.021 0.587	3.213 0.853
Omnibus: Prob(Omnibus Skew: Kurtosis:	3):	10.75 0.00 0.57 3.26	5 Jarqu 4 Prob	LTOTAL OV		2.097 10.990 0.00411 157.

Warnings: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

In [2]:

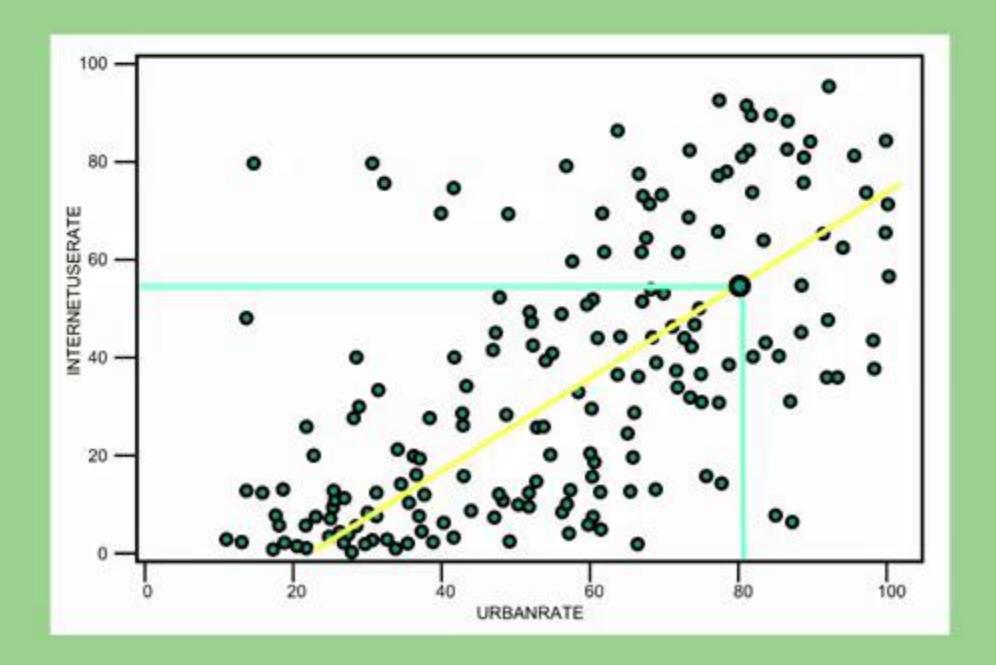
Dependent Variable Population Population Slope Coefficient $y = \beta_0 + \beta_1 x$

Explanatory Variable Variables Related in Some Meaningful Way

$$y = \beta_0 + \beta_1 x$$

$$\hat{y} = \beta_0 + \beta_1 x$$

$$\beta_{0} = -4.90$$
 $\beta_{1} = 0.72$
URBANRATE = 80



Module 2 Lesson 3 - Categorical Explanatory Variables



© Creative Commons, 2015



MAJOR DEPRESSION binary categorical explanatory variable

"Is having major depression associated with an increased number of nicotine dependence symptoms?"

NICOTINE DEPENDENCE SYMPTONS quantitative response variable

NDSymptoms = 2.19 + 1.36 (majordeplife)

Without Depression:

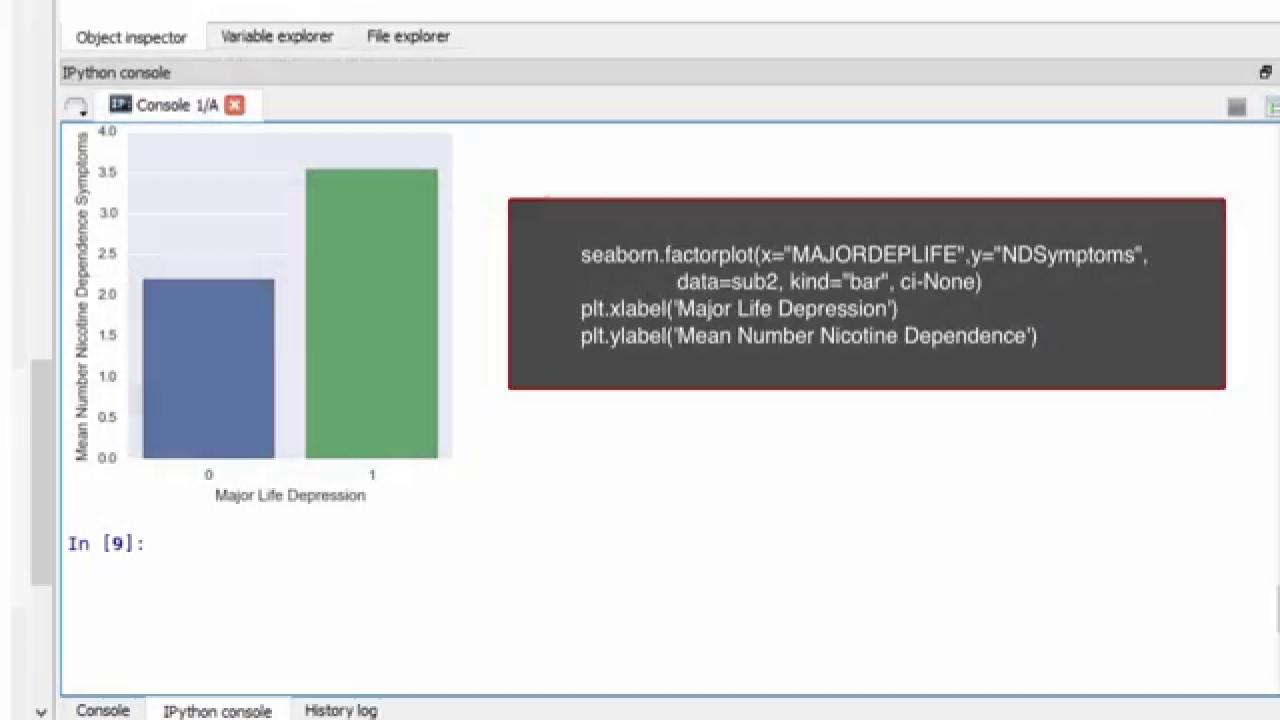
NDSymptoms = 2.19 + 1.36 (0)

NDSymptoms = 2.19

With Depression:

NDSymptoms = 2.19 + 1.36 (1)

NDSymptoms = 3.55



It is possible to include several explanatory or predictor variables to evaluate the independent contribution of multiple explanatory variables in predicting our response variable and to evaluate whether specific variables confound the relationship between our explanatory variable of interest and our response variable.

Module 2 Lesson 4 - Linear Regression Assumptions



© Creative Commons, 2015

Linear Regression Assumptions

Normality

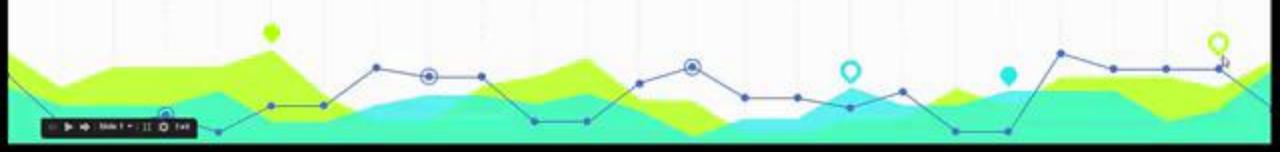
Linearity

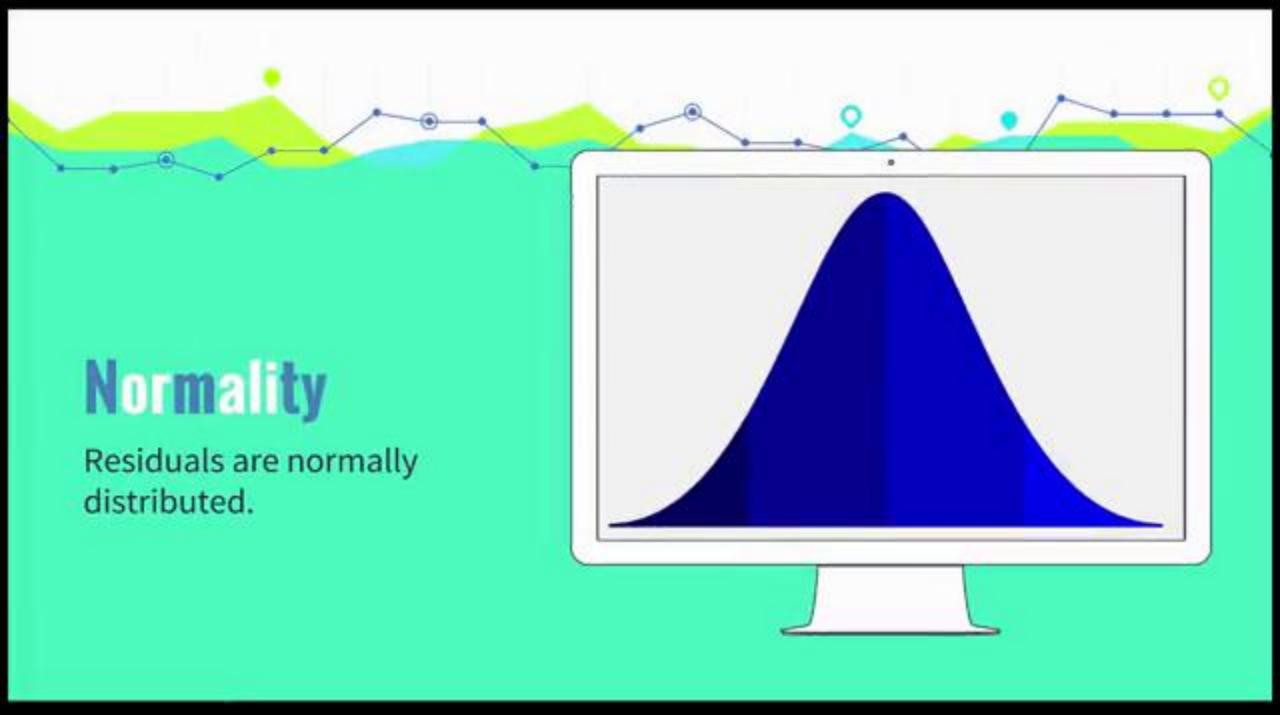
Homoscedasticity

Independence

Multicollinearity

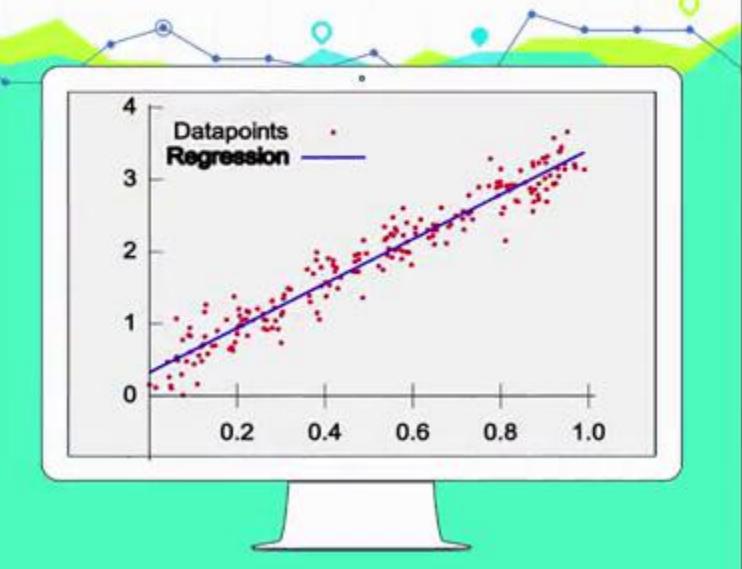
Outliers





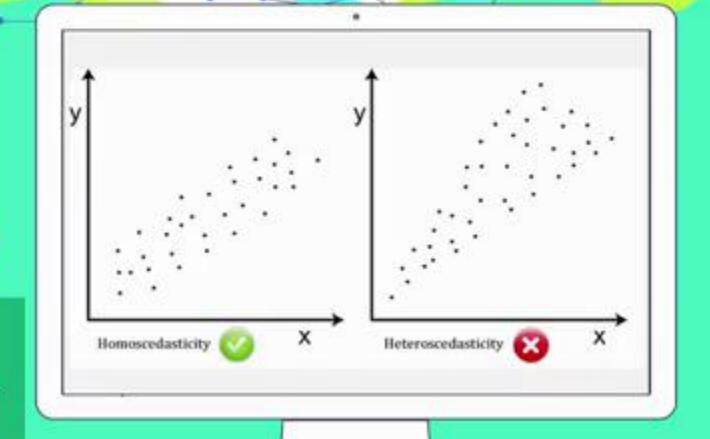
Linearty

Associations between explanatory variables and response variable are linear.



Homoscedasticity

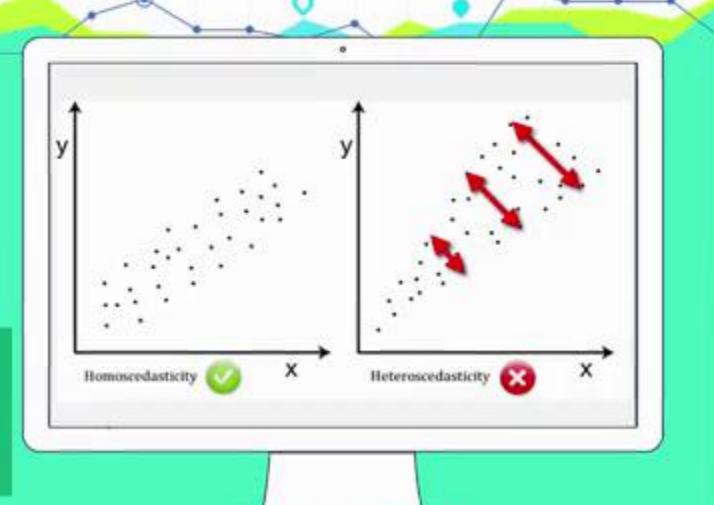
The variability in the response variable is the same at all levels of the explanatory variable



Residual = observed score (y) - predicted score (\hat{y})

Homoscedasticity

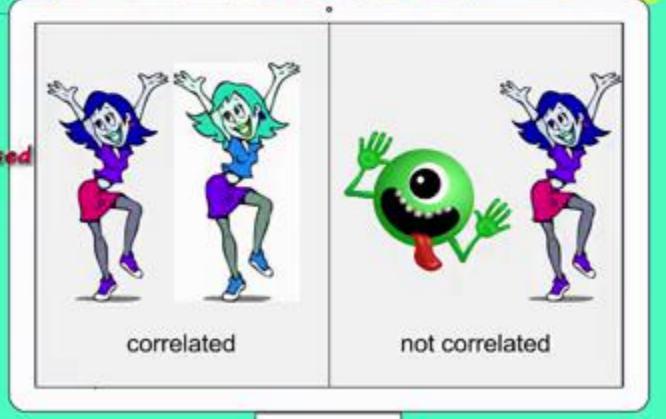
The variability in the response variable is the same at all levels of the explanatory variable



Clustered data and repeated measures data have correlated observations

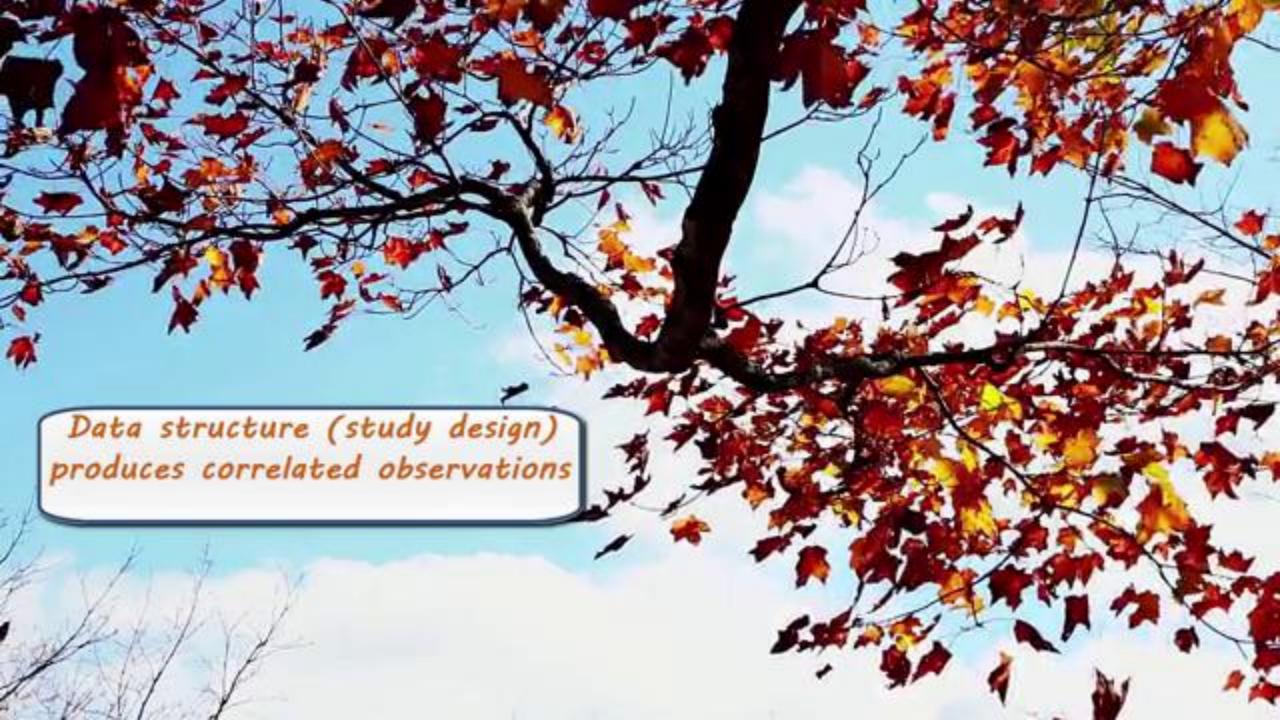
Independence

Observations are not correlated with each other.









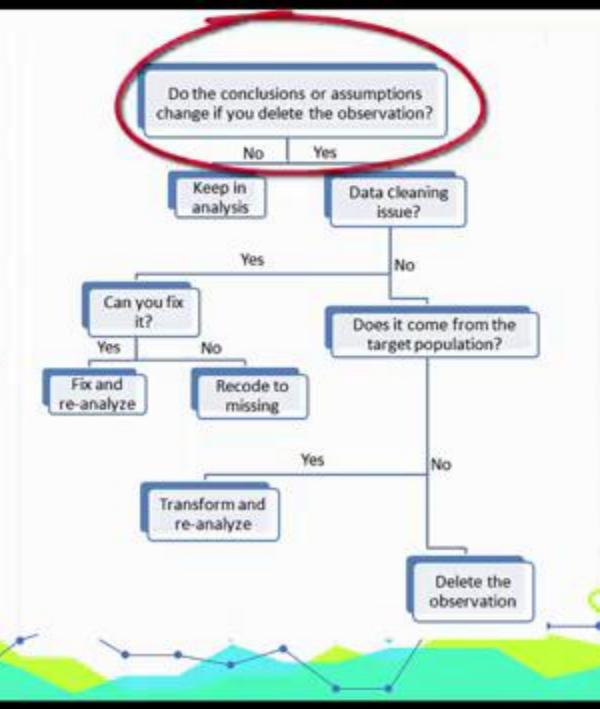
Outliers One or more observations has extreme values on variable(s) relative to other observations.

Outliers One or more observations has extreme values on variable(s) relative to other observations.



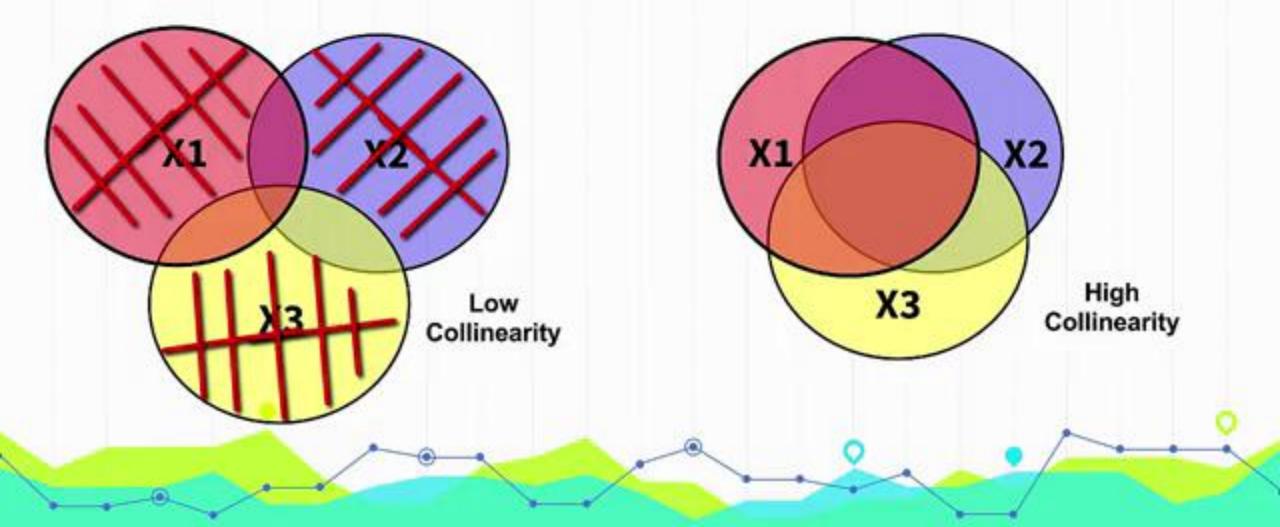


A: TRY USING THIS DECISION FLOWCHART



Multicollinearity

Explanatory variables are highly correlated with each other



Multicollinearity

Signs:

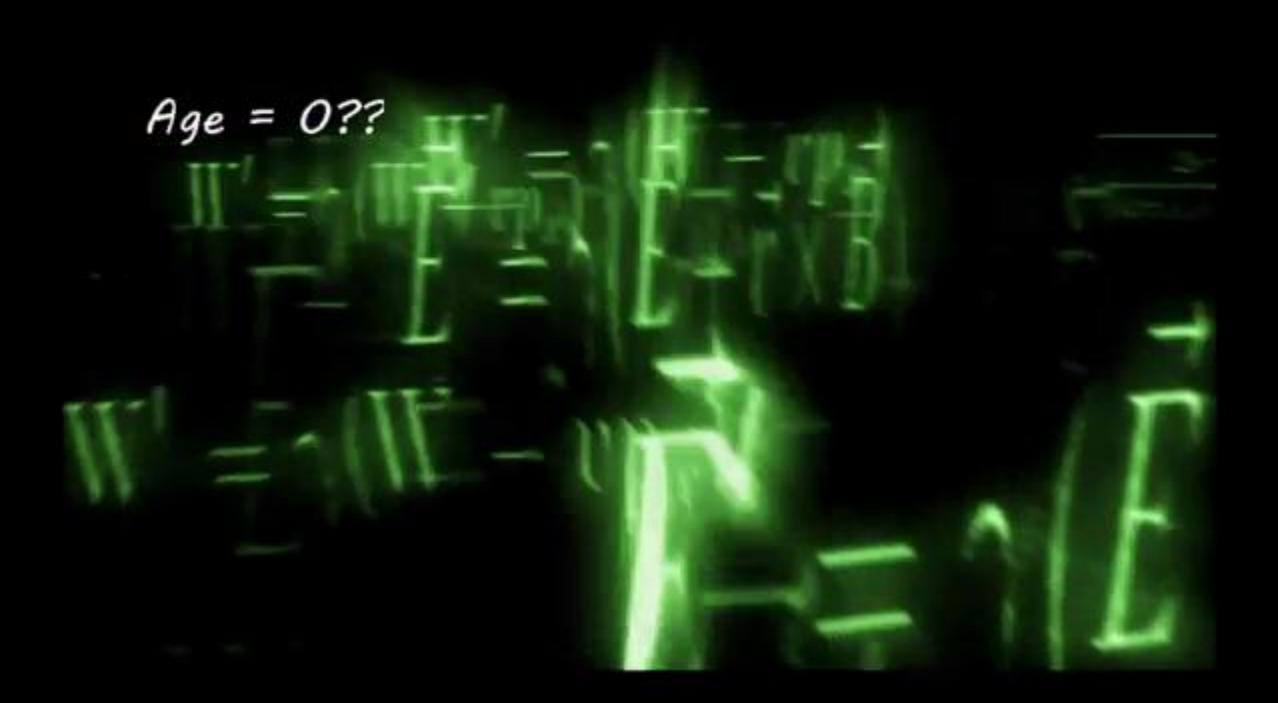
- 1) highly associated explanatory variable not significant
- 2) negative regression coefficient that should be positive
- 3) taking out an explanatory variable drastically changes results

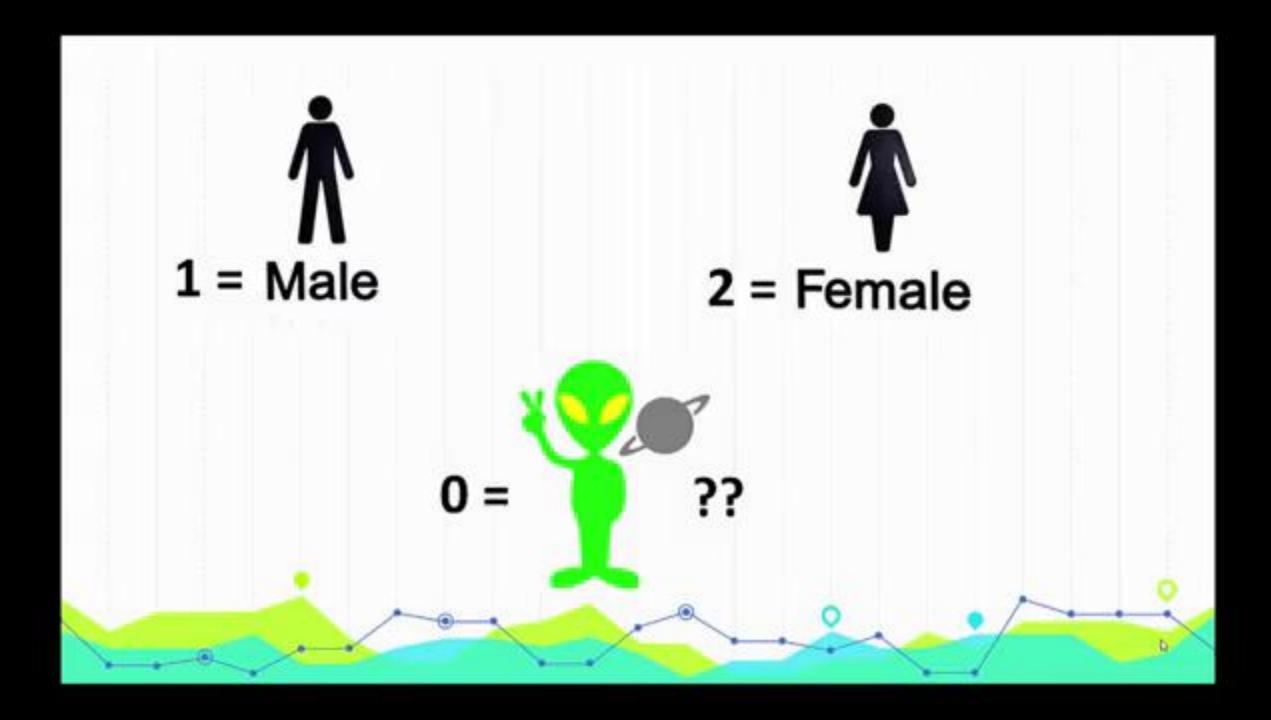
Module 2 Lesson 5 - Centering Explanatory Variables

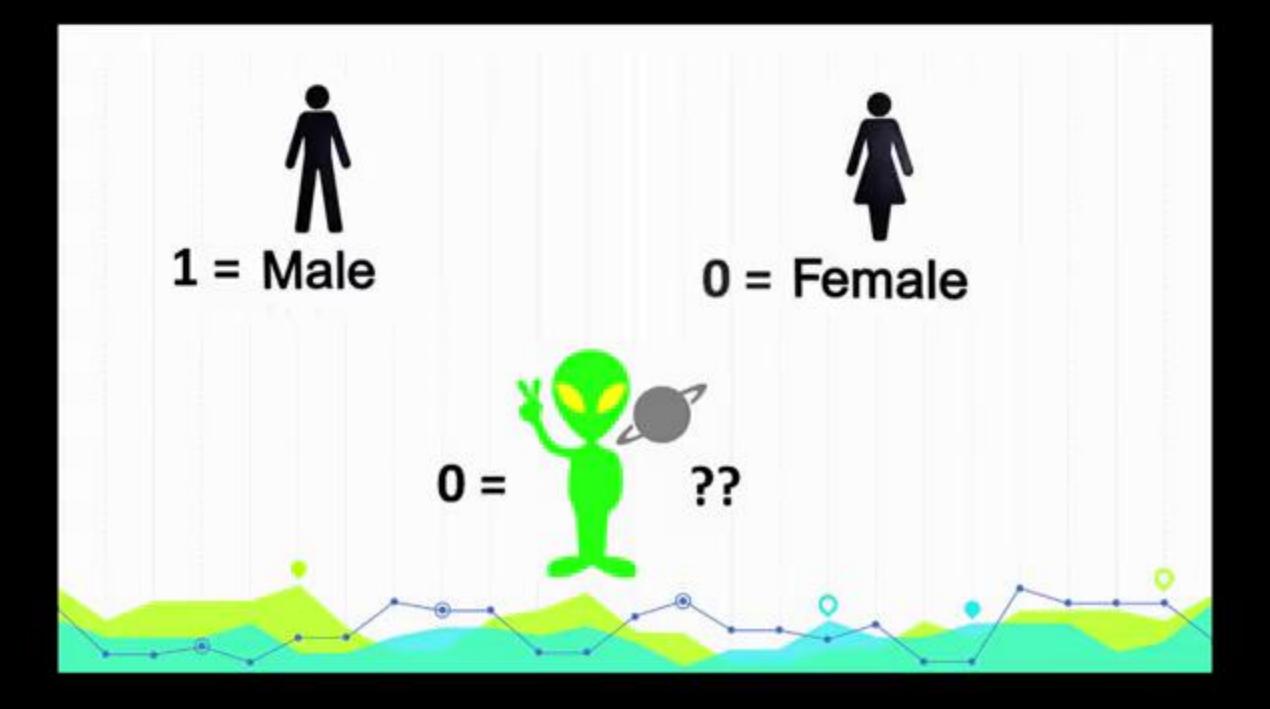


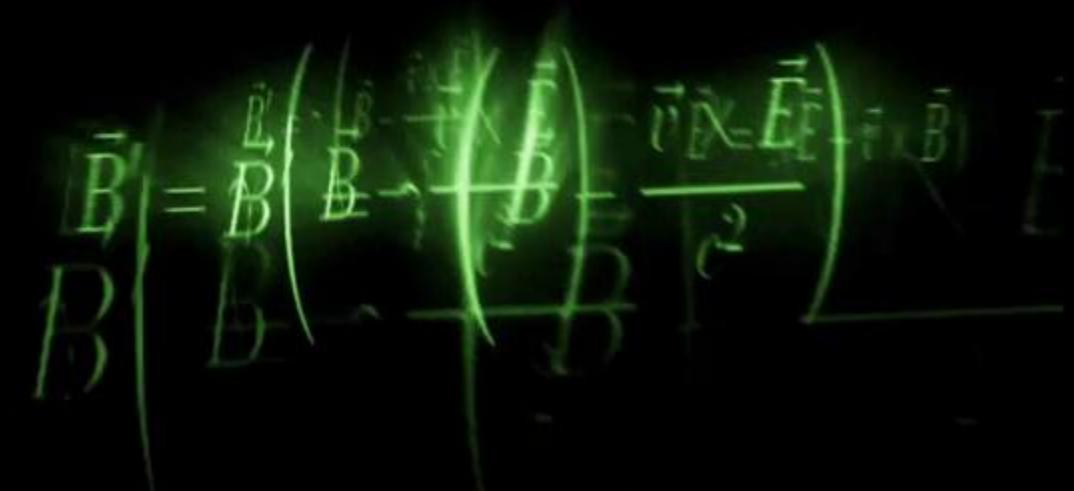
© Creative Commons, 2015



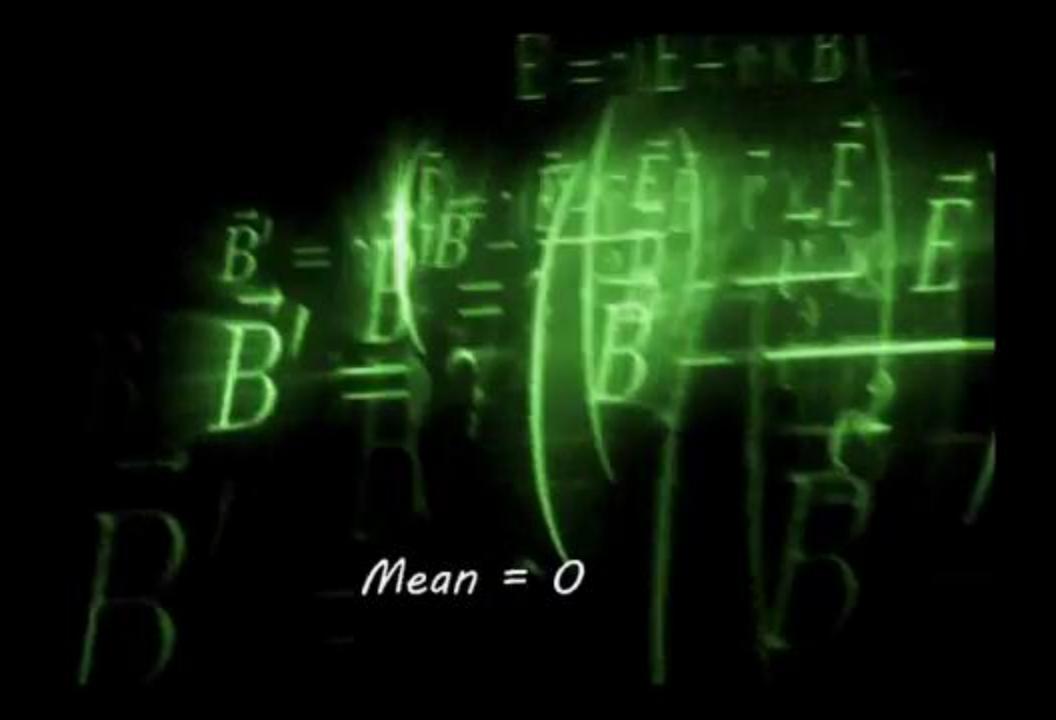








Centering = subtracting the mean of a variable from the value of the variable



Do not center the response variable