

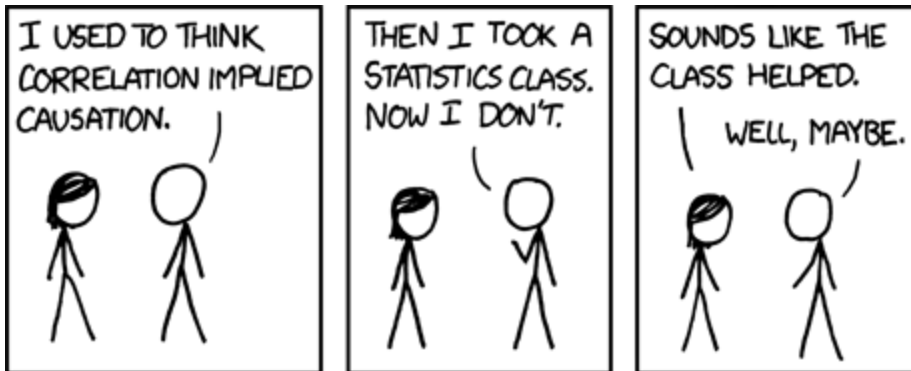
# Regressions

# Cause and Effect

**Correlation:** Two variables are correlated when changes in one variable occur in a pattern corresponding to changes in the other.

# Cause and Effect

Causation: One variable moves, and the second variable changes because of the movement of the first.



# Questioning Causality

When we suspect a causal relationship (that  $x$  causes  $y$ ), it is important to ask ourselves several questions:

1. Is it possible that  $y$  causes  $x$  instead?
2. Is it possible that  $z$  (a new factor that we haven't considered before) is causing both  $x$  and  $y$ ?
3. Could the relationship have been observed by chance?

# Establishing Causality

In order to establish causality, we need to meet several conditions:

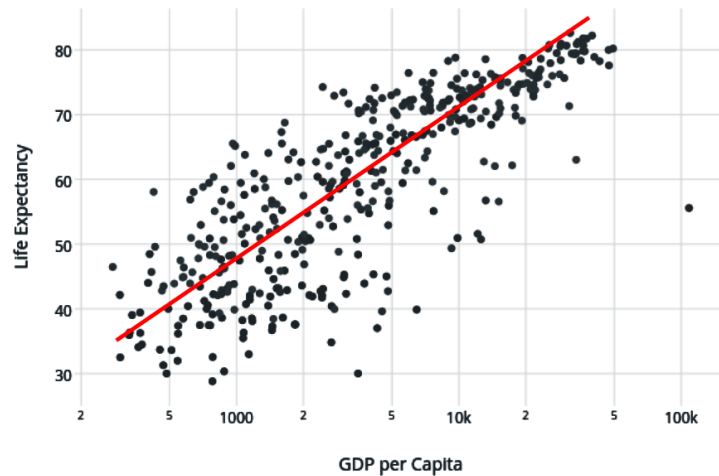
- We can explain (or at least hypothesize) **why**  $x$  causes  $y$
- We can demonstrate that **nothing else is driving the changes** (within reason)
- We can show that there is a **correlation** between  $x$  and  $y$

# Ceteris Paribus

*ceteris paribus* means "all else equal"

# Regression analysis

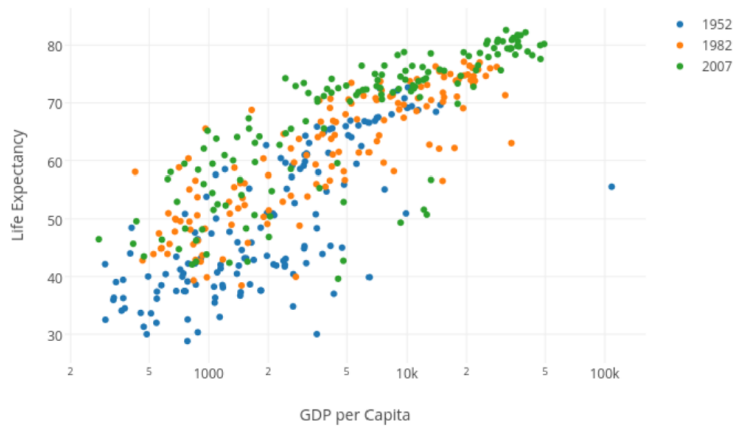
- Allows us to **act as if nothing else were changing**
- Mathematically isolates the effect of each individual **variable** on the outcome of interest
  - Variables are the factors that we want to include in our model



# Regression analysis

- Think about it like a trend line!



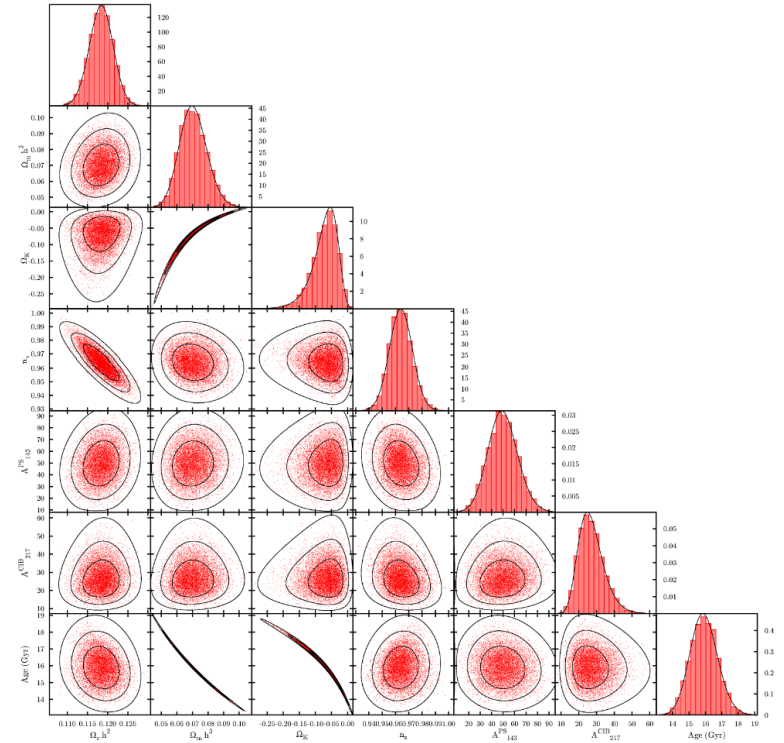


# Regression analysis

Whoops! What if there is another variable?

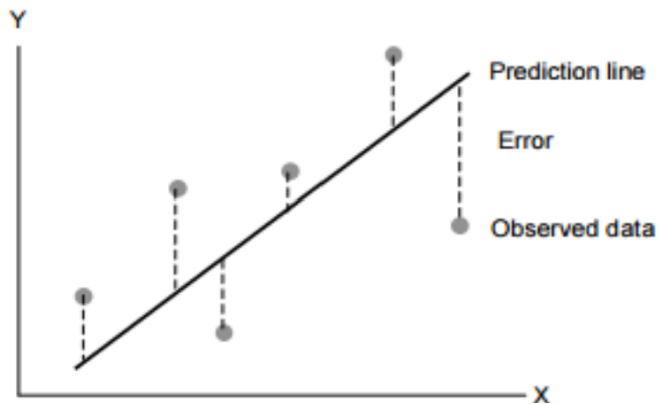
# Regression analysis

Or lots of variables??

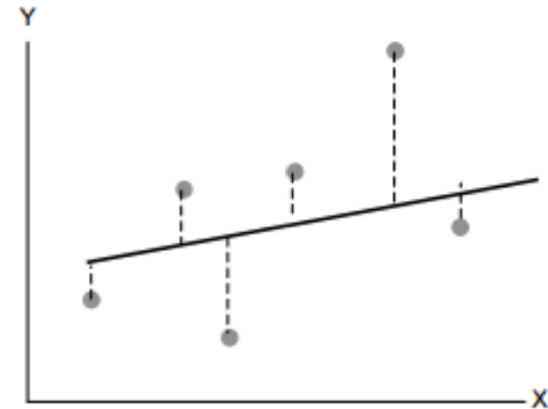


# Minimize Errors and Best Fit Lines

Best Fit



Something Else



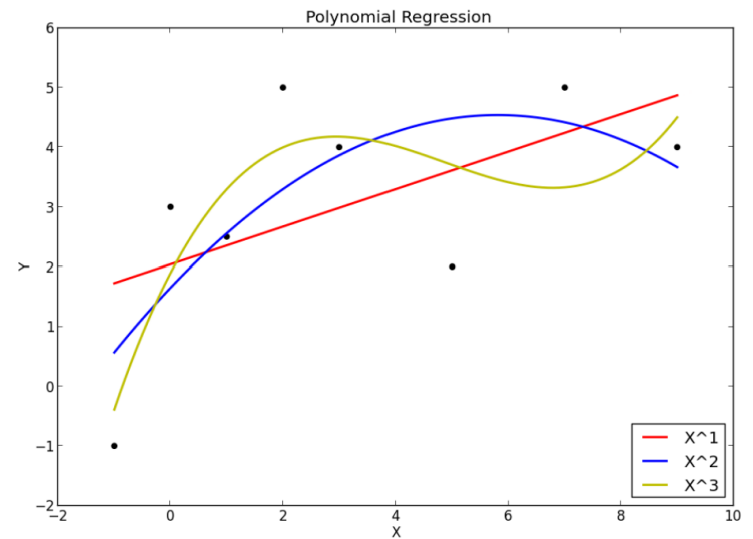
# Minimize Errors and Best Fit Lines

Try it by hand!



# Why LINEAR regression?

- Faster
- More honest



# OLS in Python

```
import pandas as pd
import statsmodels.formula.api as smf

data = pd.read_csv(
    "https://github.com/dustywhite7/pythonMikkeli/raw/master/exampleData/fishWeight.csv")

reg = smf.ols("Weight ~ Length1", data=data)

reg = reg.fit()

print(reg.summary())
```

In [5]: `reg.summary()`

Out[5]: OLS Regression Results

<b>Dep. Variable:</b>	Weight	<b>R-squared:</b>	0.839			
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.837			
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	815.3			
<b>Date:</b>	Tue, 09 Jun 2020	<b>Prob (F-statistic):</b>	4.75e-64			
<b>Time:</b>	20:09:35	<b>Log-Likelihood:</b>	-1015.1			
<b>No. Observations:</b>	159	<b>AIC:</b>	2034.			
<b>Df Residuals:</b>	157	<b>BIC:</b>	2040.			
<b>Df Model:</b>	1					
<b>Covariance Type:</b>	nonrobust					
	<b>coef</b>	<b>std err</b>	<b>t</b>	<b>P&gt; t </b>	<b>[0.025</b>	<b>0.975]</b>
<b>Intercept</b>	-462.3751	32.243	-14.340	0.000	-526.061	-398.690
<b>Length1</b>	32.7922	1.148	28.554	0.000	30.524	35.061
<b>Omnibus:</b>	9.385	<b>Durbin-Watson:</b>	0.369			
<b>Prob(Omnibus):</b>	0.009	<b>Jarque-Bera (JB):</b>	9.768			
<b>Skew:</b>	-0.489	<b>Prob(JB):</b>	0.00757			
<b>Kurtosis:</b>	3.721	<b>Cond. No.</b>	79.2			

Warnings:

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

# Regression Equations

$$\text{dependent} \sim x_1 + x_2 + x_3 + \dots$$

We can force variables to be categorical:

$$\text{dependent} \sim x_1 + x_2 + C(x_3) + \dots$$

Here, we make `x3` categorical



# Regression Equations

$$\text{dependent} \sim x_1 + x_2 + x_3 + \dots$$

We can use arithmetic transformations:

$$\text{dependent} \sim x_1 + I(x_2^{**2}) + x_3 + \dots$$

Here, we square `x2`

# When OLS Fails

OLS is an inappropriate model whenever you have a binary or discrete dependent variable (think "yes" or "no" questions)

In this case, you should use Logistic Regression instead. More details can be found in the class notes on Mimir/Github.

# Implementing Logistic Regressions

```
formula = "y ~ all_of_the_xs"  
reg = smf.logit(formula, data)  
reg = reg.fit()  
reg.summary()
```

**Lab Time!**