



# Linear Regression

Let's learn something!



# Python and Spark

- The Linear Regression section has:
  - Theory Overview Lecture
  - Documentation Example
  - Custom Code Example
  - Consulting Project Exercise



# Python and Spark

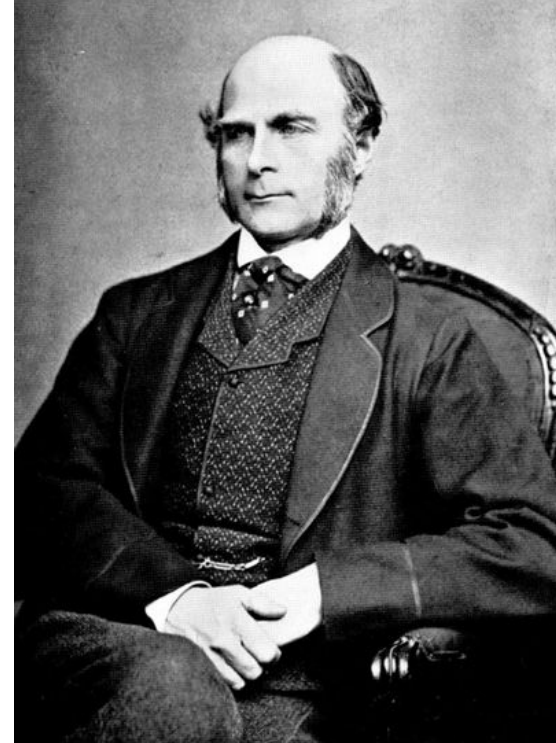
- A quick note, if you've already gone through another course of mine that has Machine Learning topics, you may have already seen this material!
- Feel free to skip the rest of this, the theory of linear regression hasn't changed :)



# Reading Assignment

Chapters 2 & 3 of  
**Introduction to Statistical Learning**  
By Gareth James, et al.

This all started in the 1800s with a guy named **Francis Galton**. Galton was studying the relationship between parents and their children. In particular, he investigated the relationship between the heights of fathers and their sons.

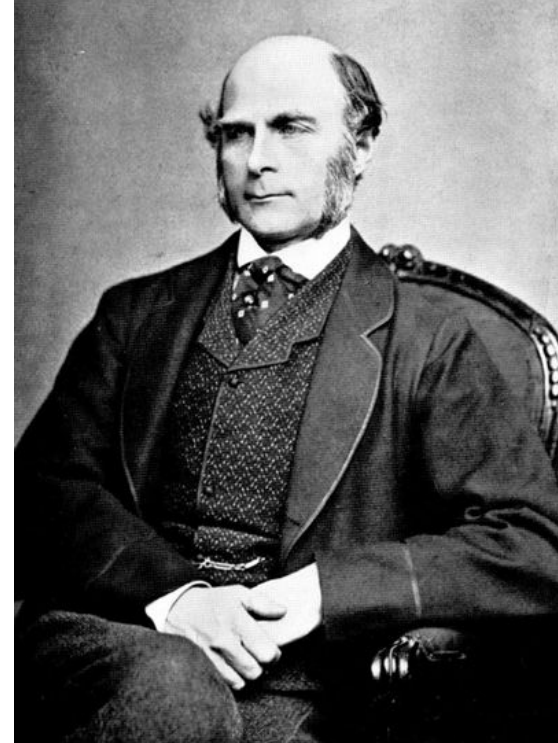




# History

What he discovered was that a man's son tended to be roughly as tall as his father.

However Galton's breakthrough was that the son's height **tended to be closer to the overall average** height of all people.

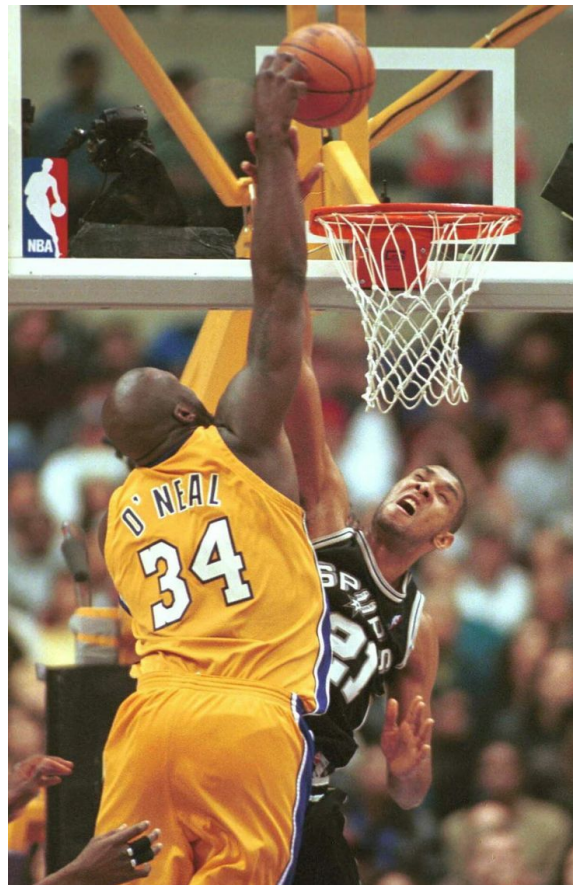




## Example

Let's take **Shaquille O'Neal** as an example. Shaq is really tall: 7ft 1in (2.2 meters).

If Shaq has a son, chances are he'll be pretty tall too. However, Shaq is such an anomaly that there is also a very good chance that his son will be **not be as tall as Shaq**.

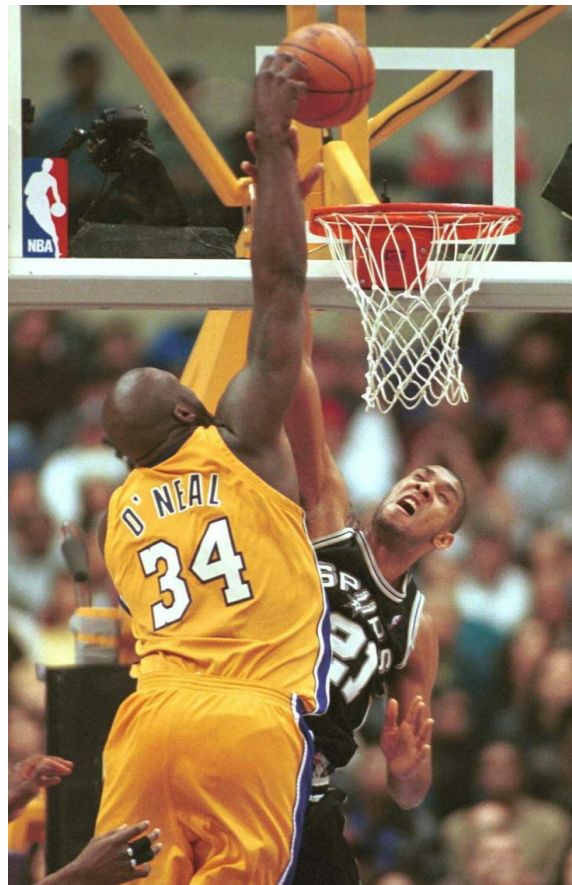




## Example

Turns out this is the case:  
Shaq's son is pretty tall (6 ft 7 in), but not nearly as tall as his dad.

Galton called this phenomenon **regression**, as in "A father's son's height tends to regress (or drift towards) the mean (average) height."

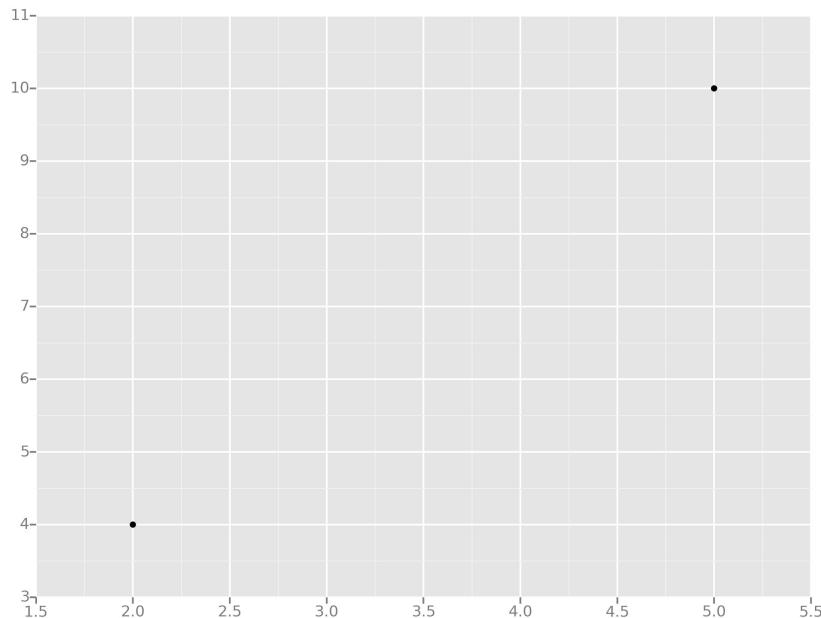






## Example

Let's take the simplest possible example: calculating a regression with only 2 data points.

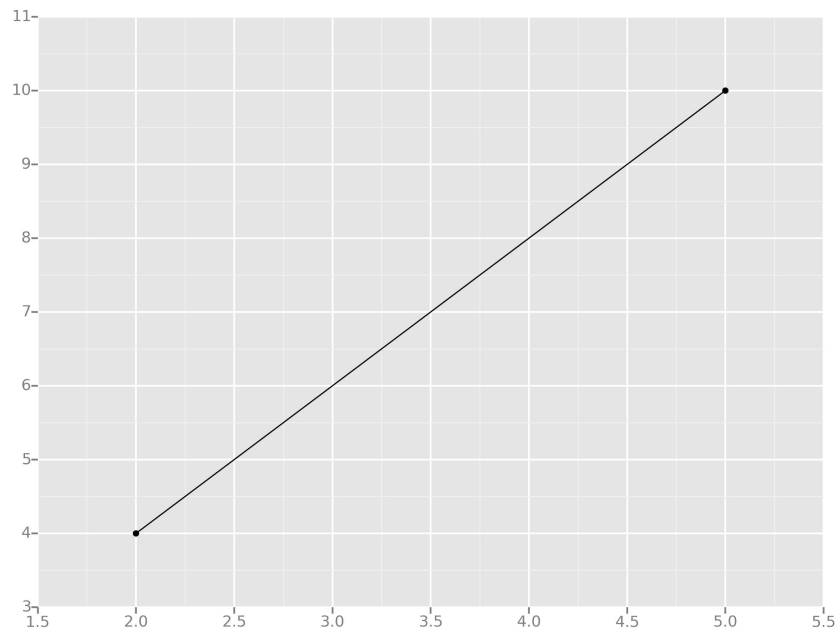




## Example

All we're trying to do when we calculate our regression line is draw a line that's as close to every dot as possible.

For classic linear regression, or "Least Squares Method", you only measure the closeness in the "up and down" direction

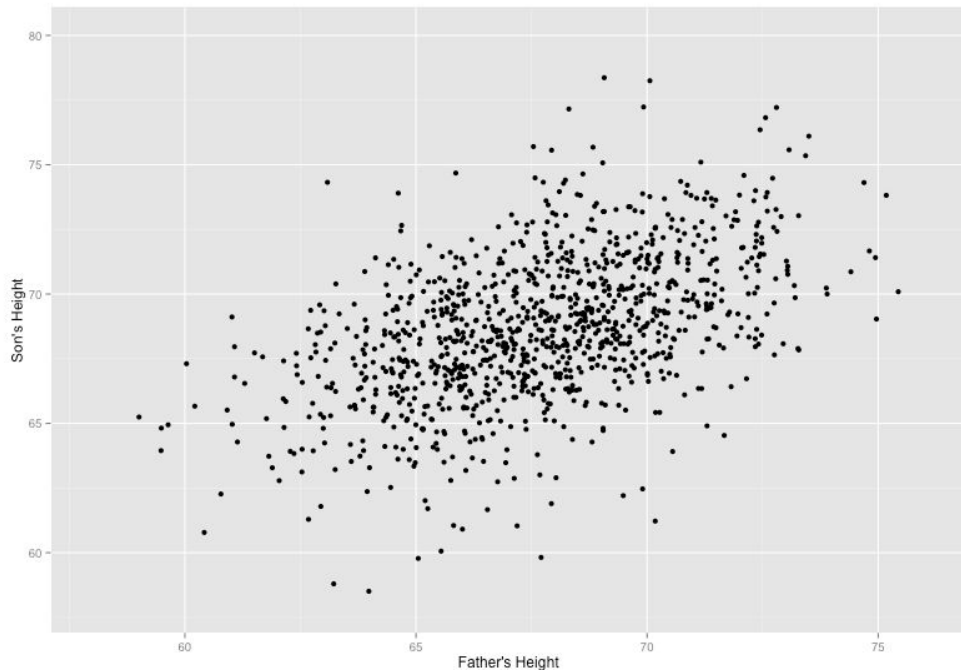




## Example

Now wouldn't it be great if we could apply this same concept to a graph with more than just two data points?

By doing this, we could take multiple men and their son's heights and do things like tell a man how tall we expect his son to be...before he even has a son!

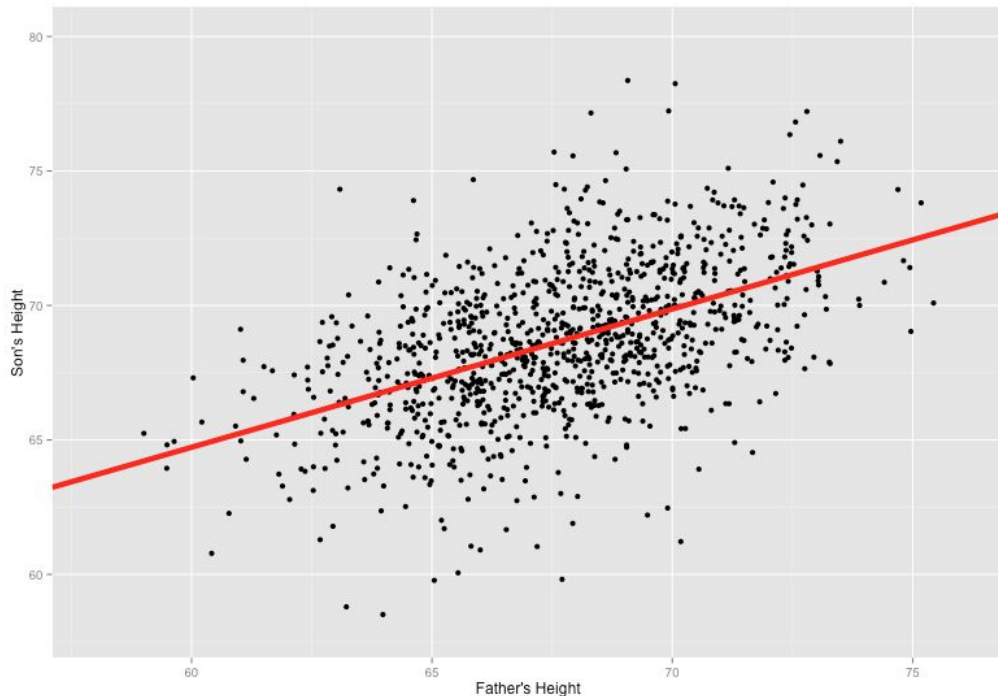




## Example

Our goal with linear regression is to **minimize the vertical distance** between all the data points and our line.

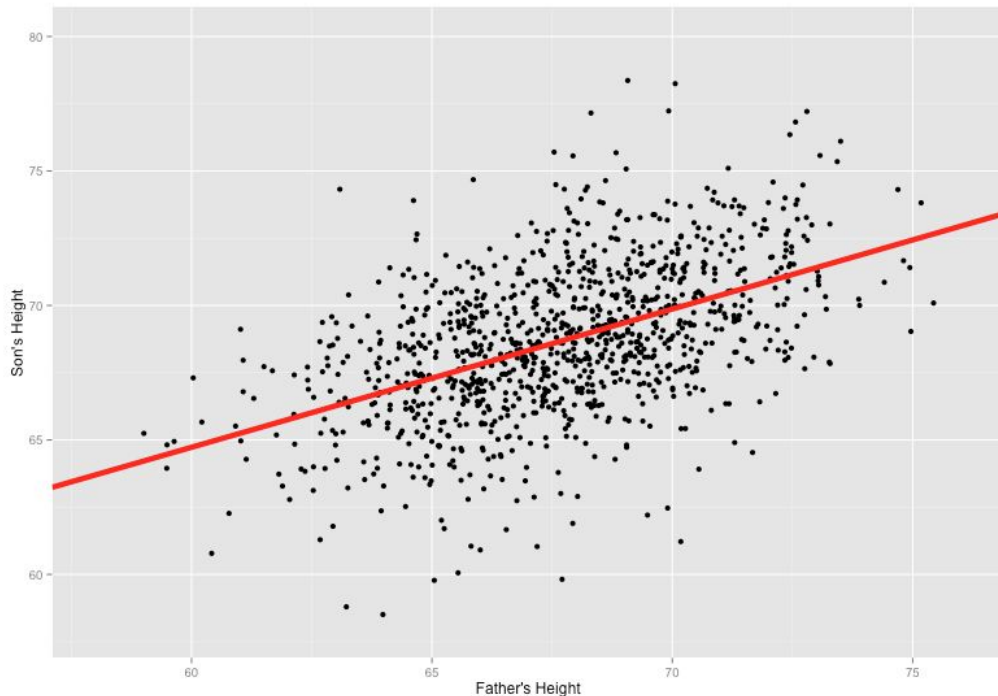
So in determining the **best line**, we are attempting to minimize the distance between **all** the points and their distance to our line.





## Example

There are lots of different ways to minimize this, (sum of squared errors, sum of absolute errors, etc), but all these methods have a general goal of minimizing this distance.

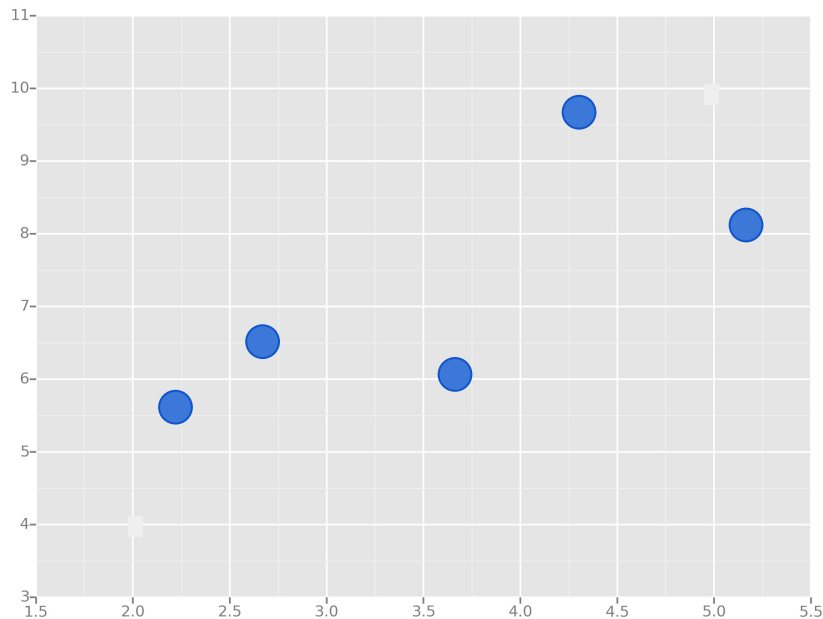




## Example

For example, one of the most popular methods is the least squares method.

Here we have blue data points along an x and y axis.

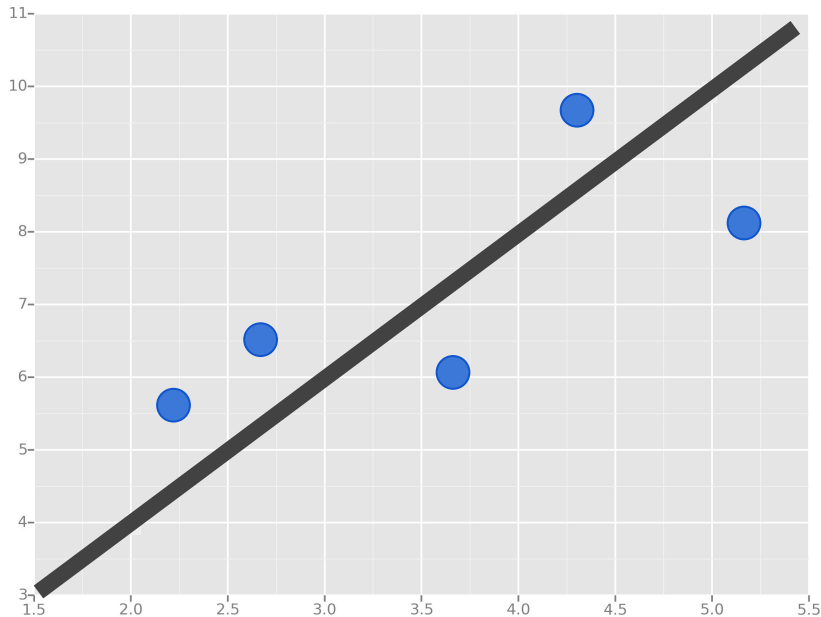




## Example

Now we want to fit a linear regression line.

The question is, how do we decide which line is the best fitting one?

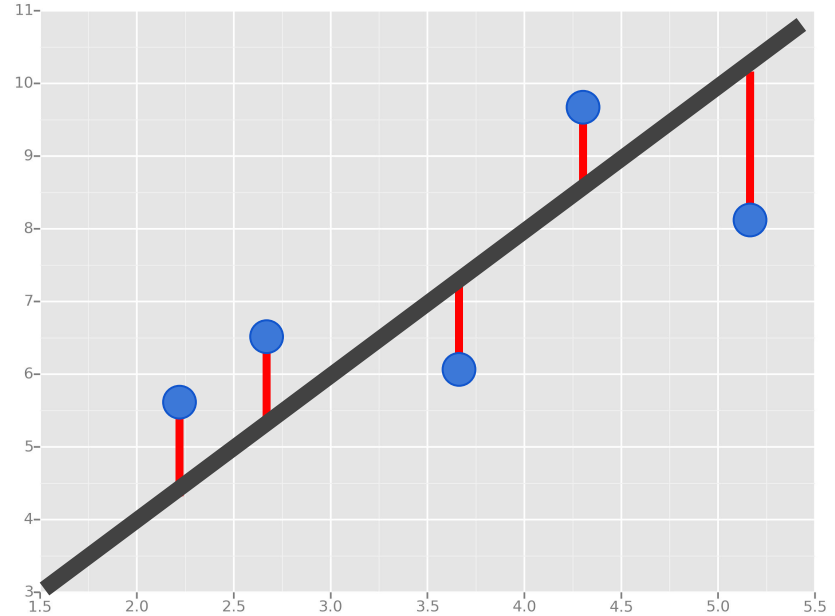




## Example

We'll use the Least Squares Method, which is fitted by minimizing the ***sum of squares of the residuals***.

The residuals for an observation is the difference between the observation (the y-value) and the fitted line.







# Python and Spark

- Now let's get an idea of how to implement this idea with PySpark!
- Will ease into all of this by checking out the simpler documentation example first!



# **Linear Regression Documentation Example**



# Python and Spark

- While it may not be super interesting, understanding how to read and re-apply documentation examples will rapidly speed up your own learning process!
- Let's walk through the Linear Regression Documentation Example!



# Python and Spark

- In the browser we'll have:
  - Documentation Page
  - Data from Documentation
  - Linear\_Regression\_Example.ipynb
  - New Untitled Notebook



# Evaluating Regression



# Python and Spark

- Let's take a quick break to discuss evaluating Regression Models
- Not just Linear Regression, but any model that attempts to predict continuous values (unlike categorical values, which is classification)



# Python and Spark

- You may have heard of some evaluation metrics like accuracy or recall.
- These sort of metrics aren't useful for regression problems, we need metrics designed for **continuous** values!



# Python and Spark

- Let's discuss some of the most common evaluation metrics for regression:
  - Mean Absolute Error
  - Mean Squared Error
  - Root Mean Square Error
  - R Squared Values





# Python and Spark

- Mean Absolute Error (MAE)
  - This is the mean of the absolute value of errors.
  - Easy to understand, just average

er

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$



# Python and Spark

- Mean Squared Error (MSE)
  - This is the mean of the squared errors.
  - Larger errors are noted more than with

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

popular.



# Python and Spark

- Root Mean Absolute Error (RMSE)
  - This is the root of the mean of the squared errors.
  - Most popular (has same units as y)

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$



# Python and Spark

- R Squared Values (also known as the coefficient of determination)
- Not quite an error metric, more of a statistical measure of your regression model.



# Python and Spark

- By itself,  $R^2$ , won't tell you the “whole story”.
- In a basic sense it is a measure of how much variance your model accounts for.
- Between 0-1 (0% to 100%)
- There are also different ways of obtaining  $R^2$ , such as adjusted R squared.



# Python and Spark

- A full analysis and explanation of interpreting R Squared is outside the scope of this course, but check out the resource links for more information on this statistical topic.



# Python and Spark

- Just keep in mind that R Squared can enhance your understanding of a model, or help you compare models, but it shouldn't be your sole source of evaluating a model!
- Let's continue on to our code example!



# Linear Regression

# Custom Code Example

Let's learn something!





# Python and Spark

- Now let's walk through a more realistic example of Linear Regression!
- We'll examine some Ecommerce Customer Data for a company's website and mobile app.



# Python and Spark

- The main idea will be trying to predict a customer's total amount expenditure (a continuous money value).
- We'll also discuss how to convert realistic data into a format accepted by Spark's MLlib!



# Python and Spark

- The relevant notebook is called:
  - Linear\_Regression\_Code\_Alone
- Let's get started!



# Python and Spark

- Then we want to see if we can build a regression model that will predict the customer's yearly spend on the company's product.
- Let's get started!



# Linear Regression Consulting Project



# Python and Spark

- Now let's set you loose on your first consulting project!
- You've been contracted by Hyundai Heavy Industries to help them build a predictive model for some ships.



# Python and Spark

- You've been flown out to their HQ in Ulsan, South Korea!





# Python and Spark

- It's one of the world's largest manufacturers of large ships, including cruise liners!







# Python and Spark

- They need your help to give accurate estimates of how many crew members a ship will require.
- They are currently selling ships to some new customers and want you to create a model and use it to predict how many crew members the ships will need.



# Python and Spark

- They provided you data with these features:
  - Ship Name
  - Cruise Line
  - Age (as of 2013)
  - Tonnage (1000s of tons)
  - passengers (100s)
  - Length (100s of feet)
  - Cabins (100s)
  - Passenger Density
  - Crew (100s)



## Python and Spark

- Your job is to create a regression model that will help predict how many crew members will be needed for future ships.
- In other words, use the features you think will be useful to predict the value in the Crew column.



# Python and Spark

- The client also mentioned that they have found that particular cruise lines will differ in acceptable crew counts, so it is most likely an important feature to include in your analysis!



## Python and Spark

- The cruise line value is a string however!
- We haven't covered exactly how to convert strings to numbers with Python and Spark (yet)
- Try to see if you can discover how to use **StringIndexer** from the documentation!



# Python and Spark

- As in any real world project, there are no “100% correct” answers here.
- Just try your best to build the model!
- You can optionally see if you can figure out **StringIndexer** on your own (we’ll cover it more formally in future lectures)



# Python and Spark

- This is a pretty open consulting project, feel free to either try it on your own, or go straight to the “solutions” to do a code-along!
- There can be more than one way to create this model!



# Python and Spark

- Best of luck with the project, all of the necessary information for the project can be found in the files:
- `Linear_Regression_Consulting_Project.ipynb`
- `cruise_line_info.csv`





# Python and Spark

Best of luck!



# **Linear Regression Consulting Project Example Solution**