



Discussion `printf("C")`

Today we will learn:

Today we will learn:



Lina

Regression

But what is Linear ReGURAssion?



In simplified words:

It is a way to capture the trend line of dataset

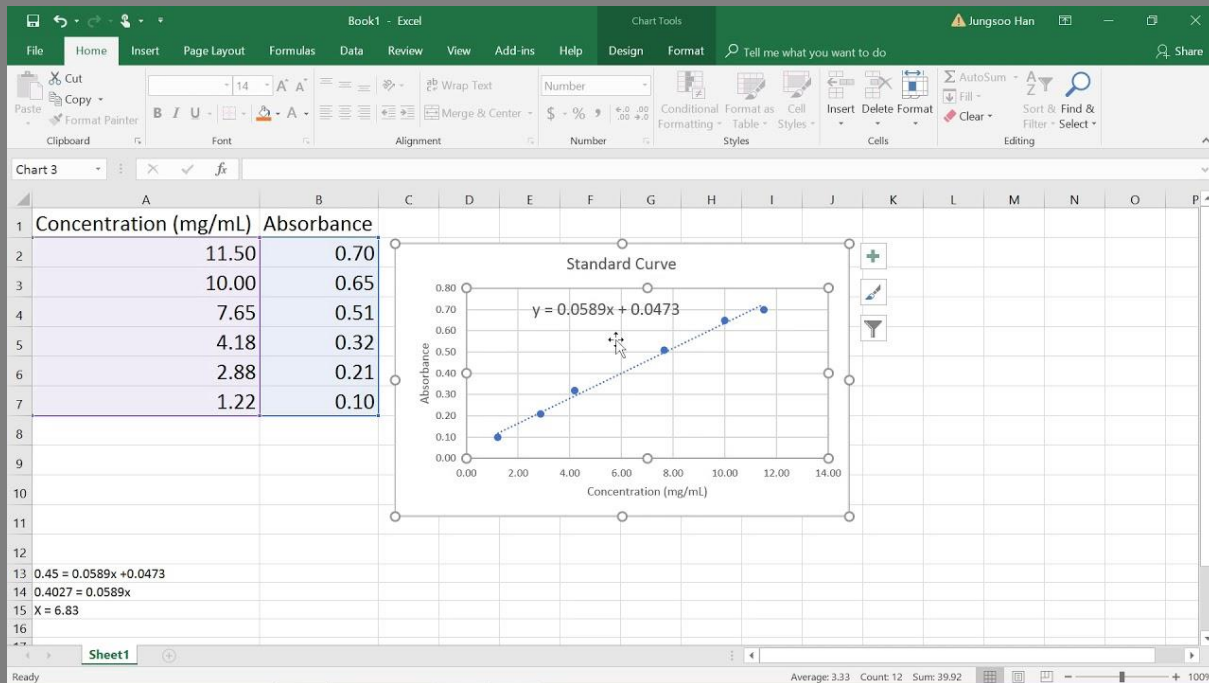
In longer explanation:

In statistics, linear regression is a linear approach to modelling the relationship between a scalar response and one or more explanatory variables (also known as dependent and independent variables) [wiki]

Let's look inside the math behind
linear rEVANGELION



Remember this?



This is him now

Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal: minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

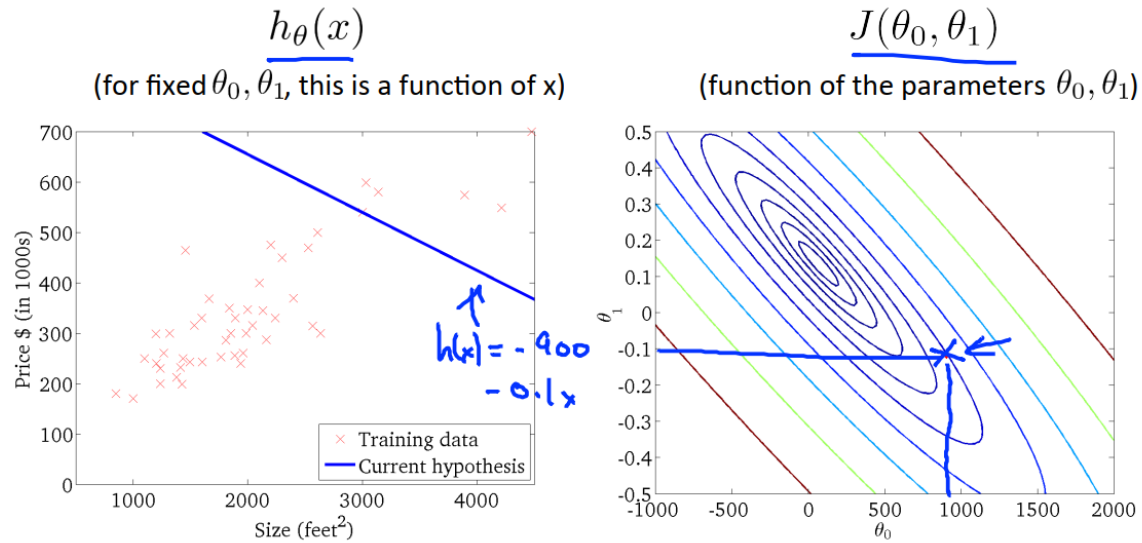
Andrew Ng

Credit: Andrew Ng

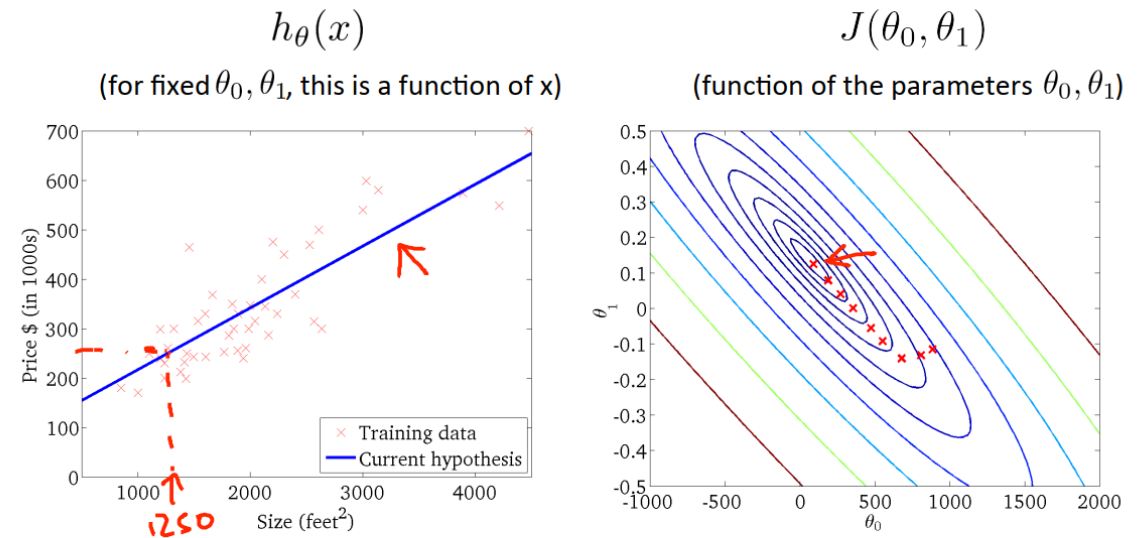
Feels old yet?

But how do you get from this

To this



Andrew Ng

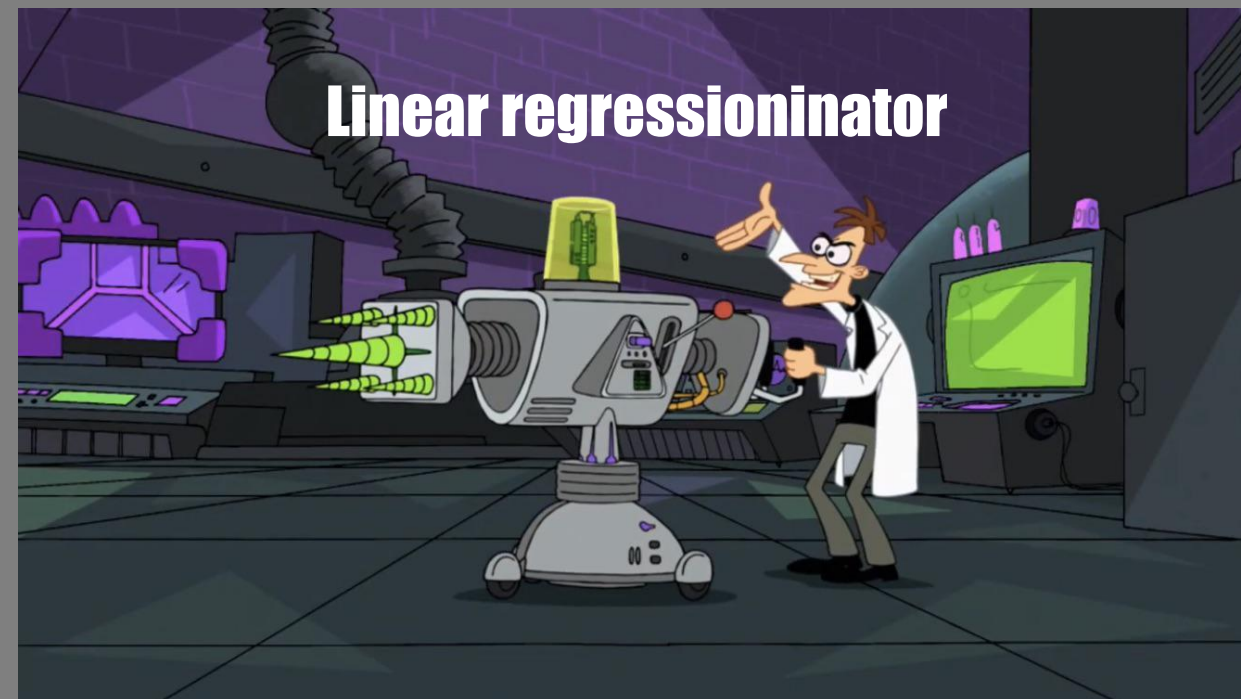
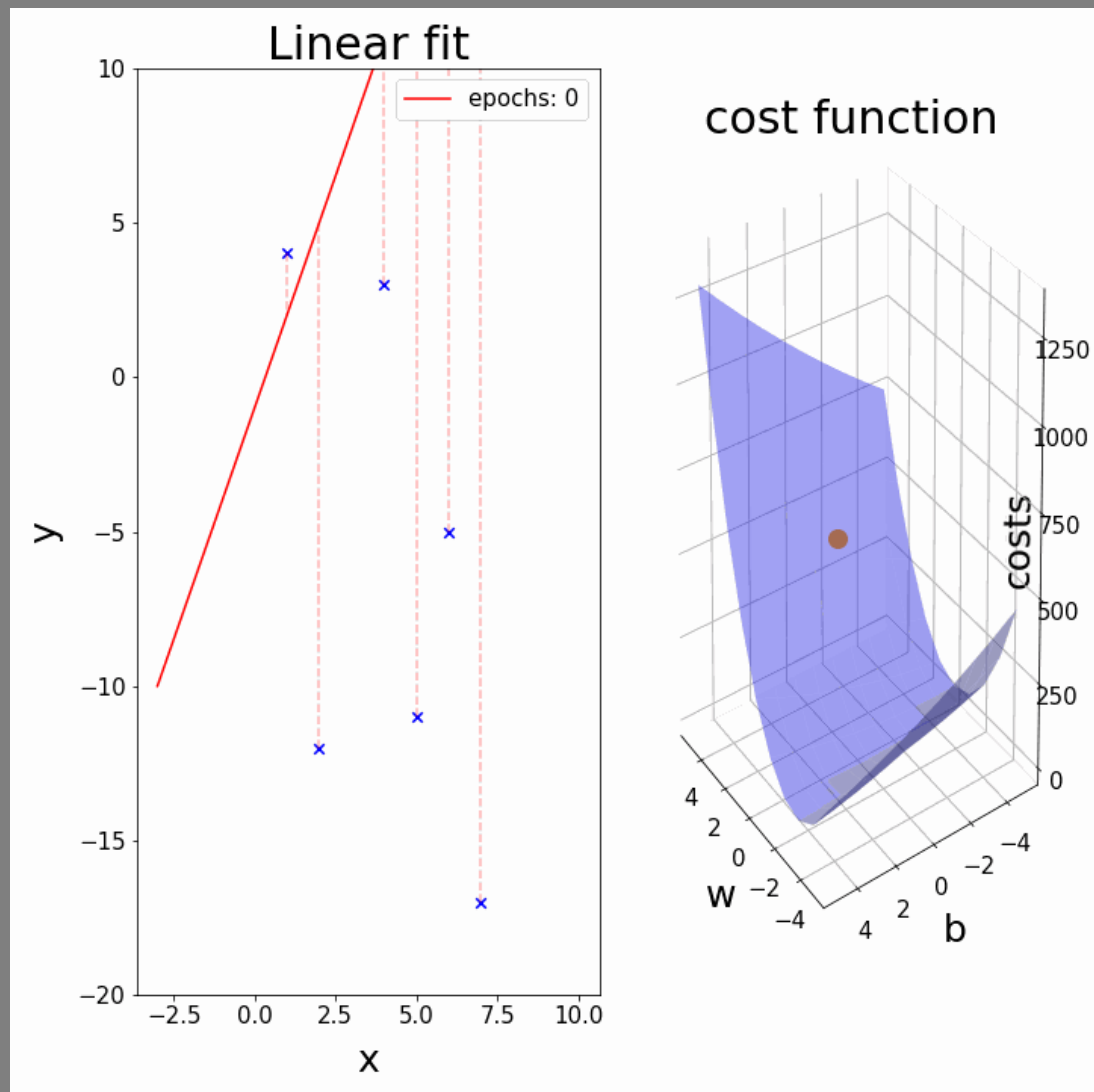


Andrew Ng

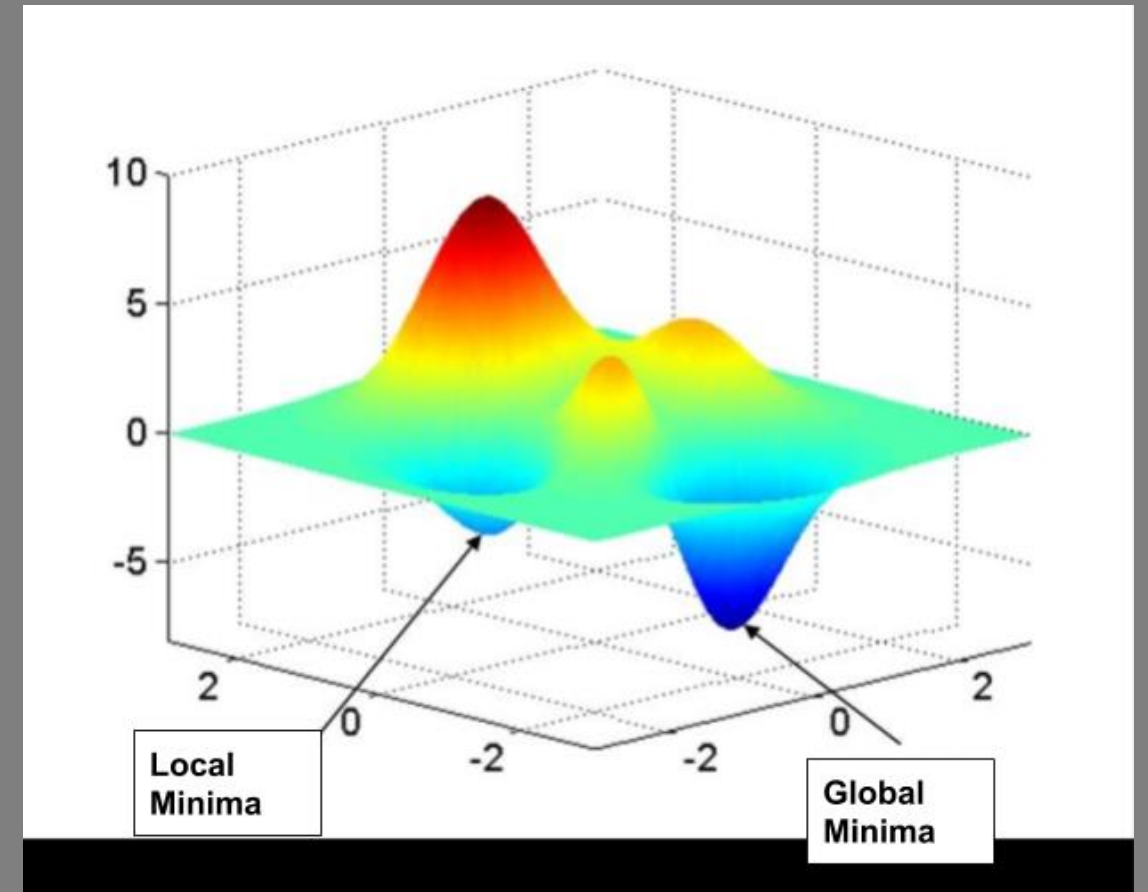
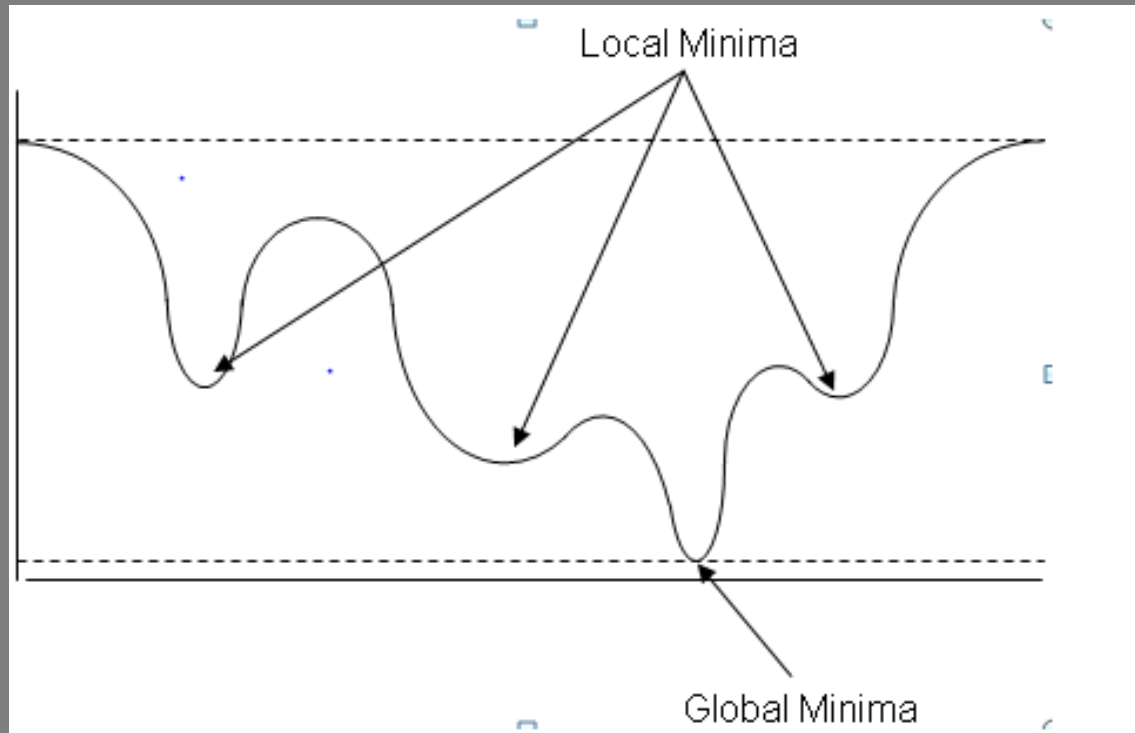
Credit: Andrew Ng

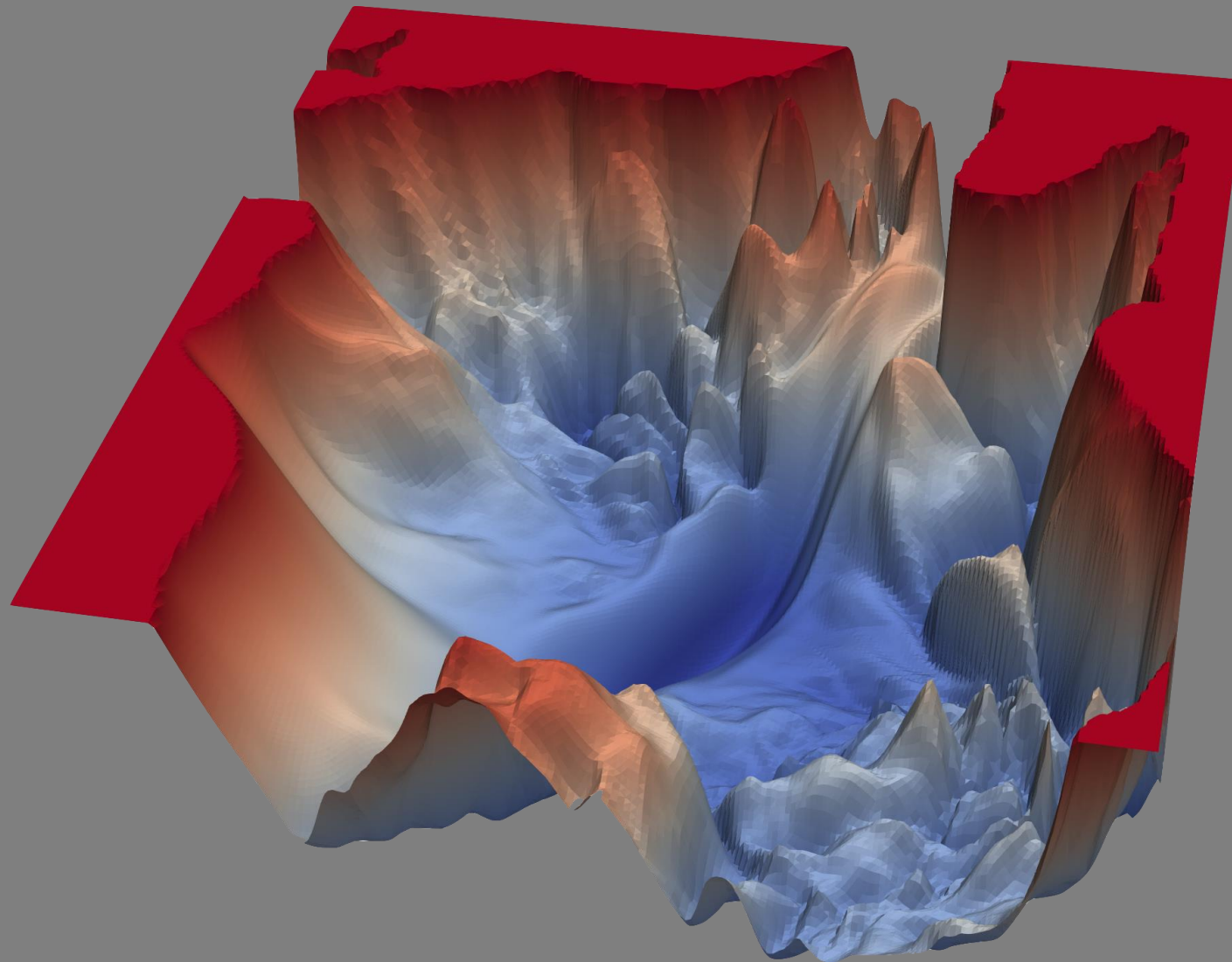
Introducing: Gradient Descent





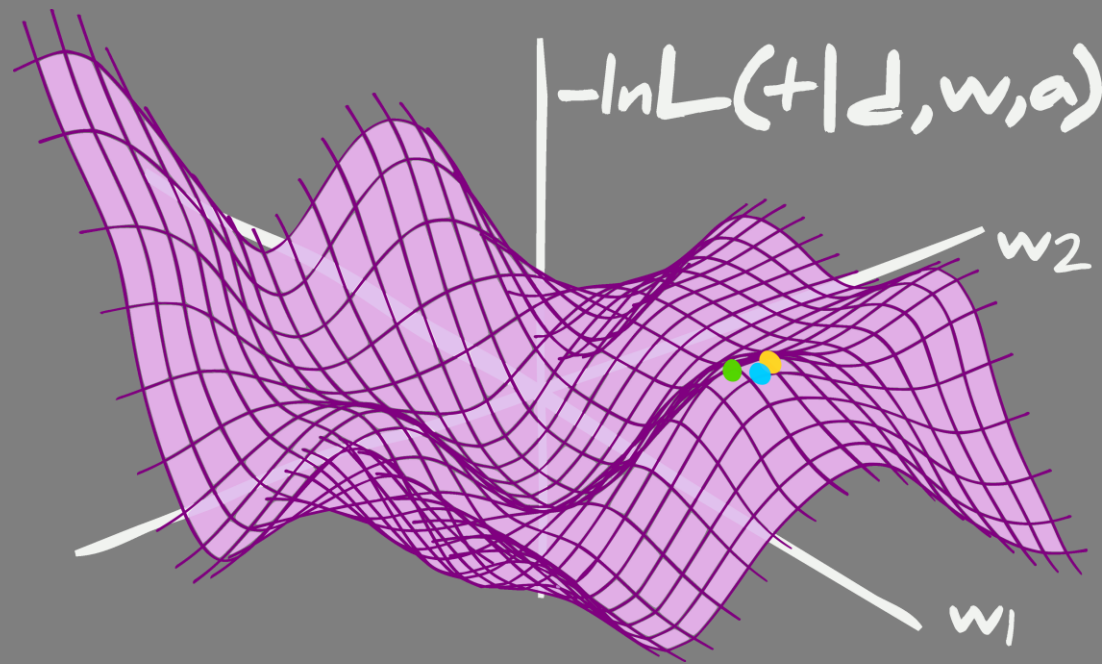
Issue: what if we cannot get the lowest cost function





A Complicated Loss Landscape *Image Credits:*
<https://www.cs.umd.edu/~tomg/projects/landscapes/>

- Adjust learning rate (difficult to find the right learning rate)
- Repeat the trial while changing starting point every time (Stochastic Gradient Descent)



Linear Regression with multiple features



Simple
Linear
Regression

$$y = b_0 + b_1 * x_1$$

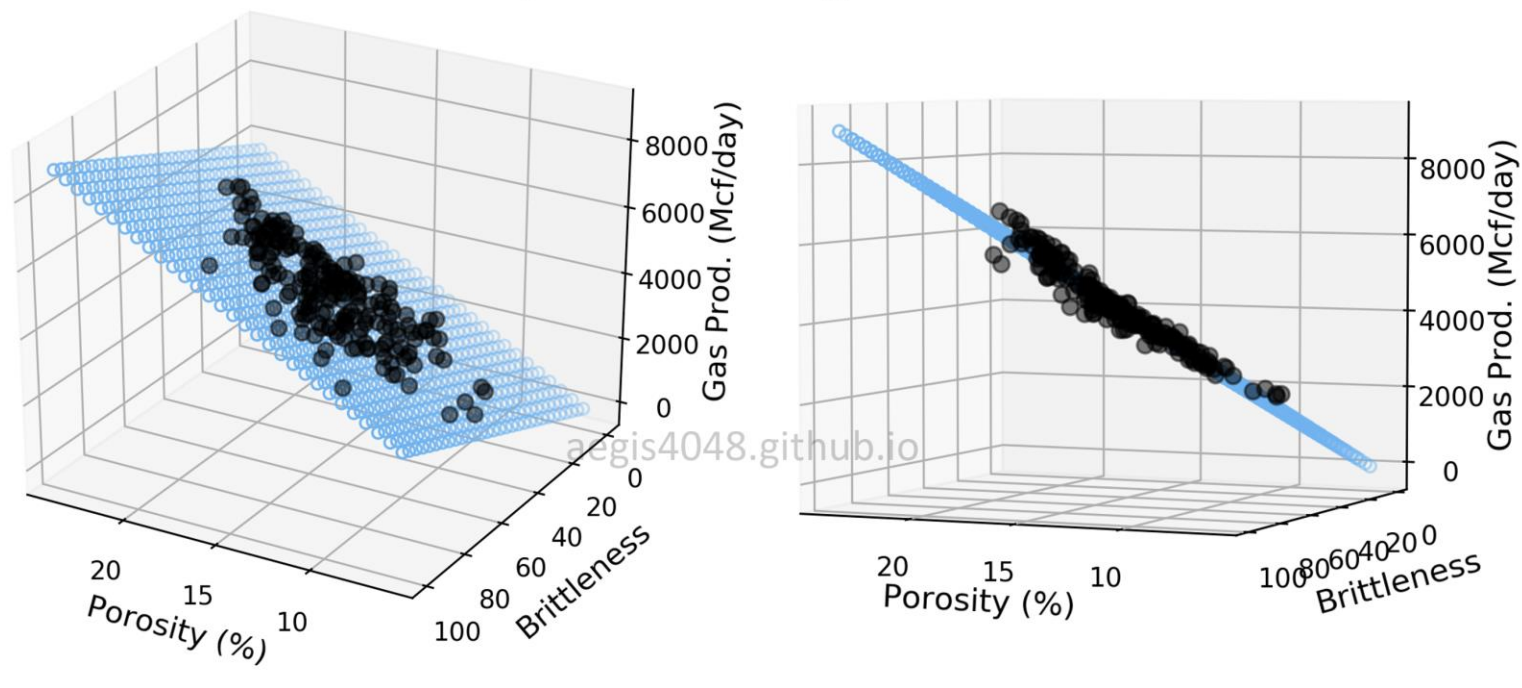
Multiple
Linear
Regression

Dependent variable (DV) Independent variables (IVs)

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

Constant Coefficients

3D multiple linear regression model



If we have this....

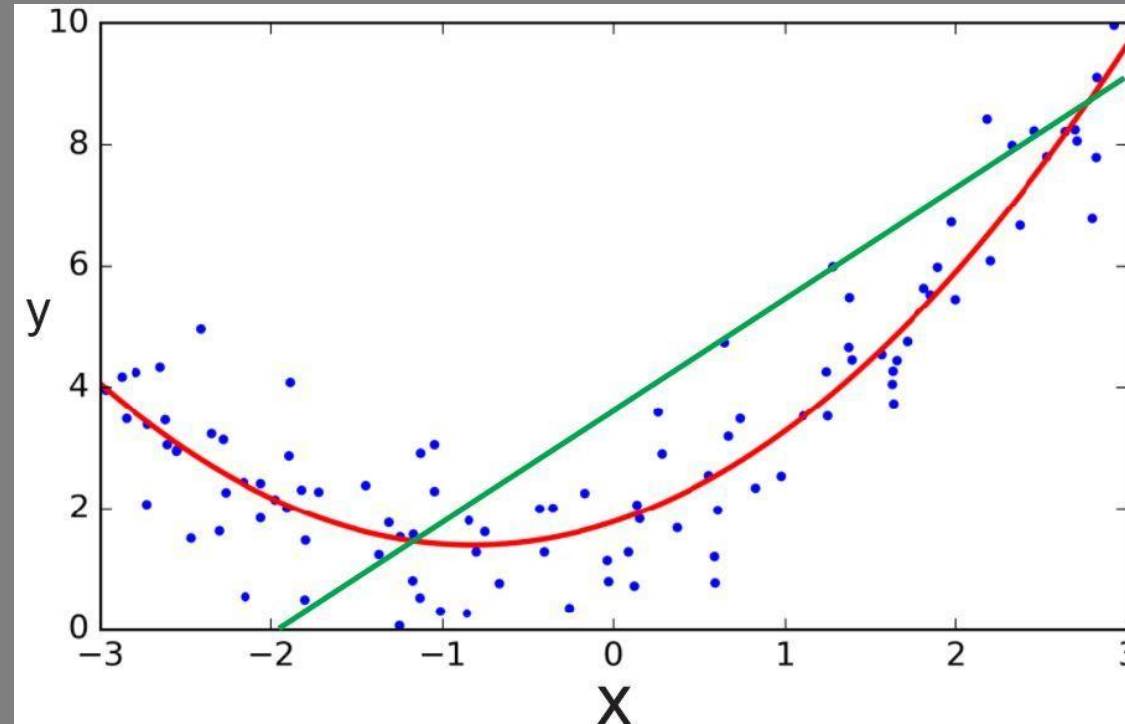
Multiple
Linear
Regression

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

We can also have this!

Polynomial
Linear
Regression

$$y = b_0 + b_1x_1 + b_2x_1^2 + \dots + b_nx_1^n$$



Same cost(loss) function and gradient descent

repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

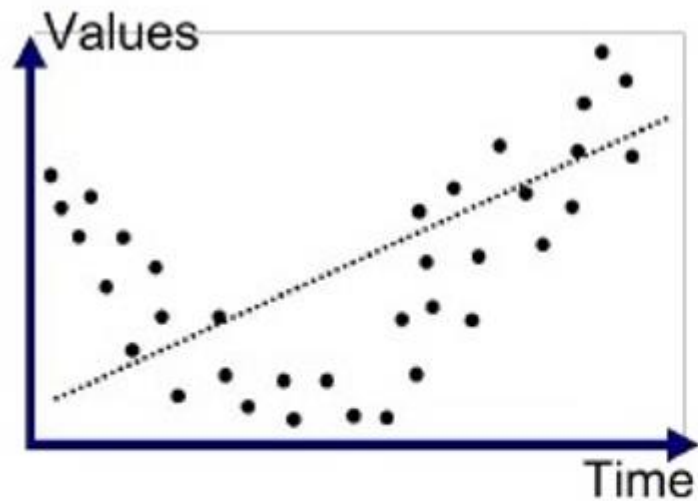
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}$$

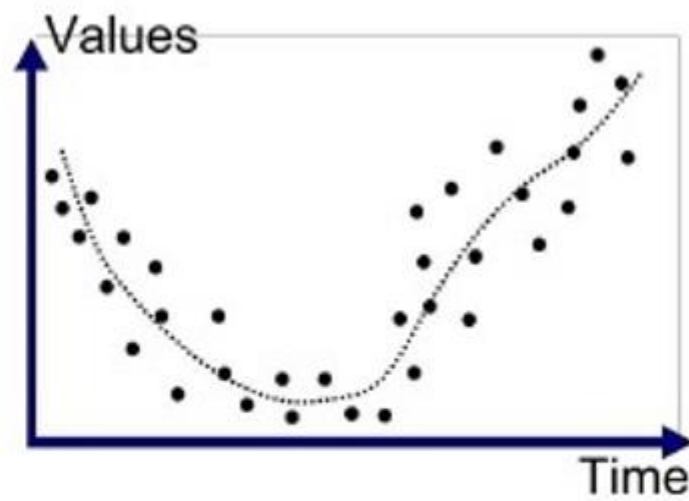
...

}

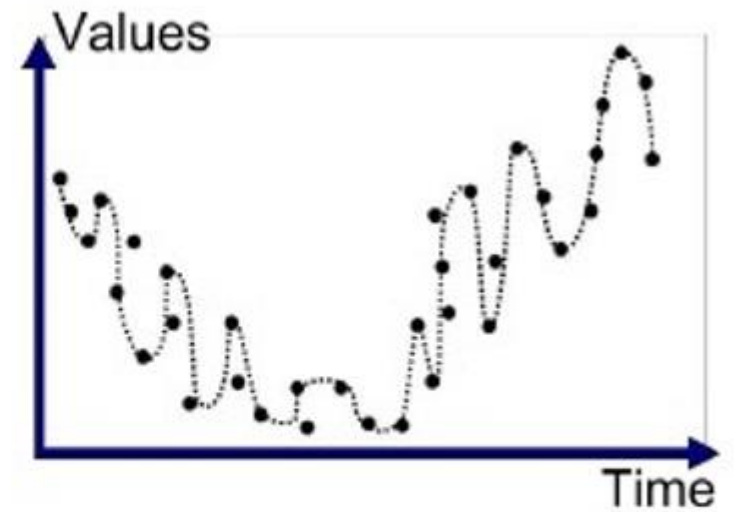
The good, the bad, and the ugly



Underfitted

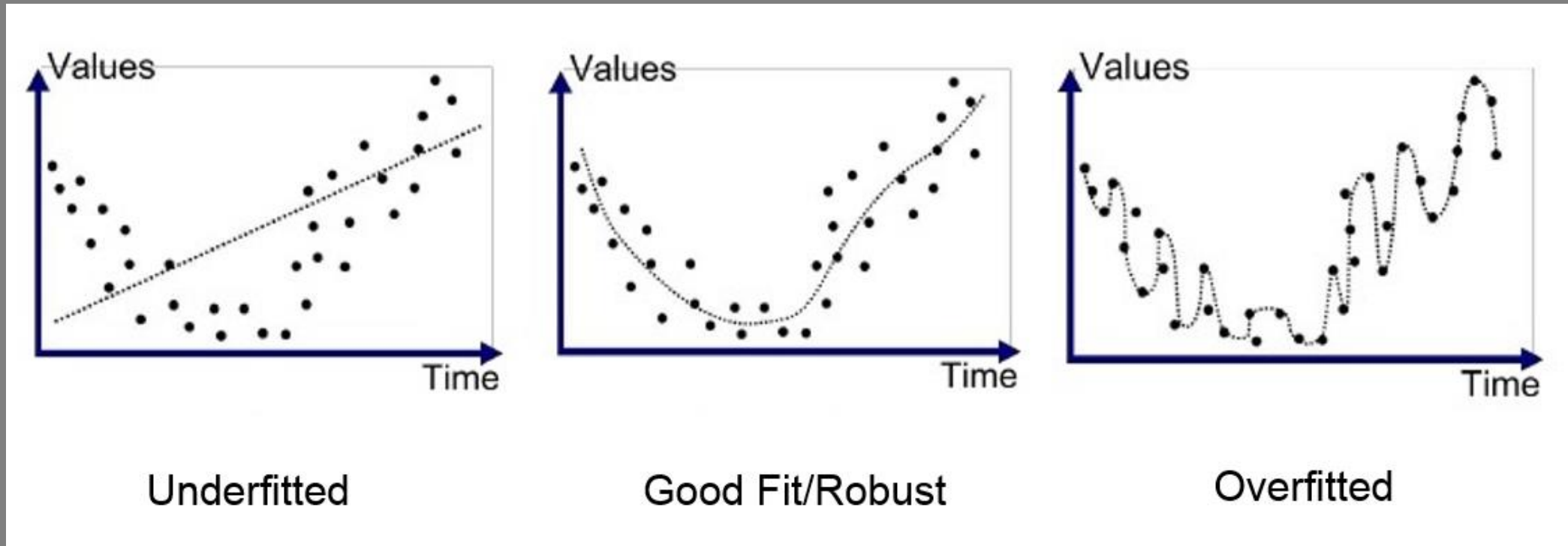


Good Fit/Robust



Overfitted

Terminology and Tips



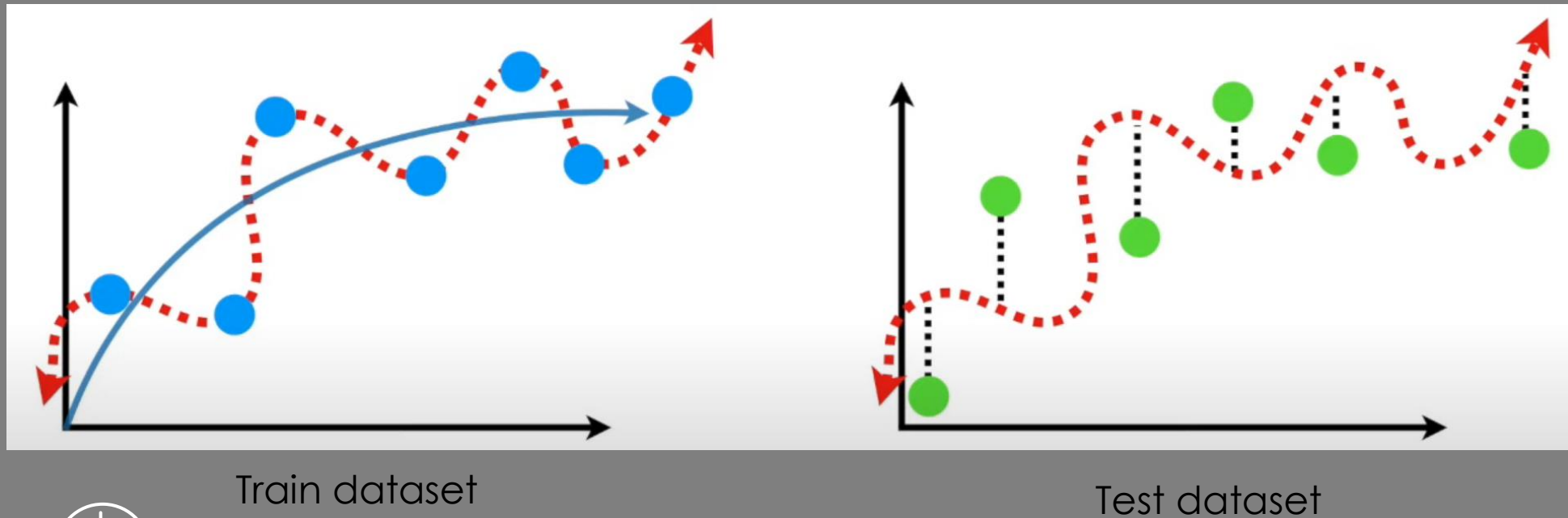
We use the term “bias” and “variance” as another way to explain how well the trend line (or plane if 3D) captures data

Bias: Inability to capture the true relationship

Variance: The difference in cost function between train dataset and test dataset

Bias: Inability to capture the true relationship

Variance: The difference in cost function between train dataset and test dataset



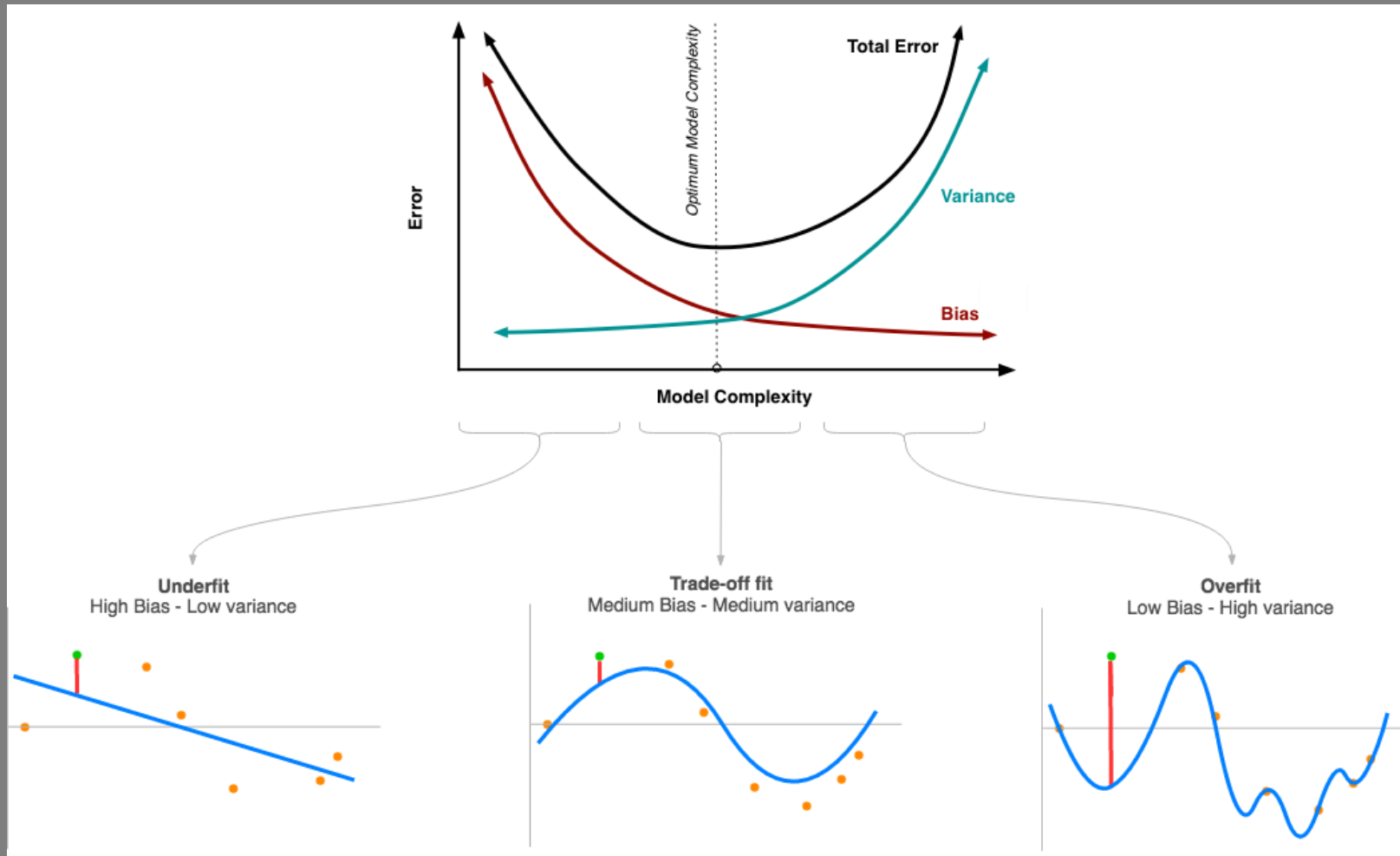
Ideally: low bias and low variance

Bias: Inability to capture the true relationship

Variance: The difference in cost function between train dataset and test dataset

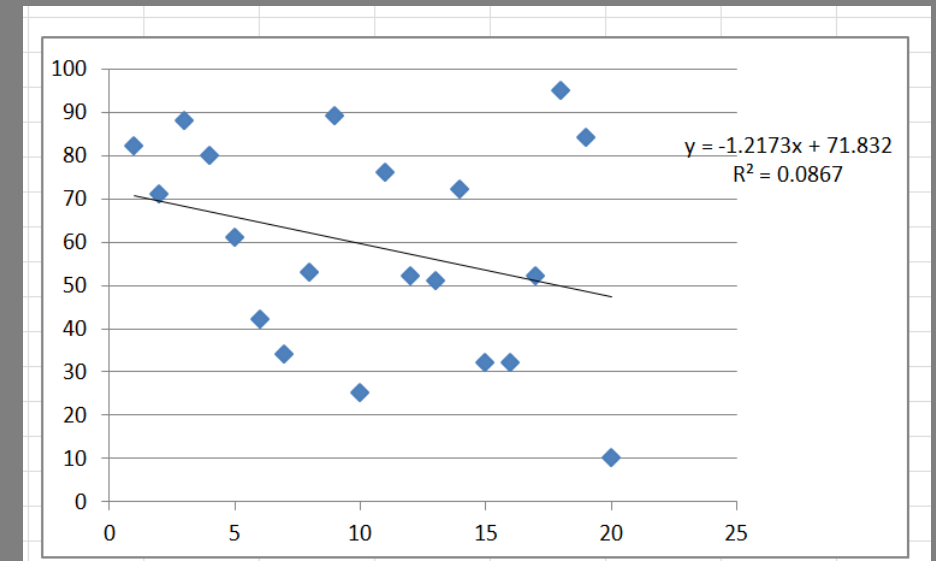
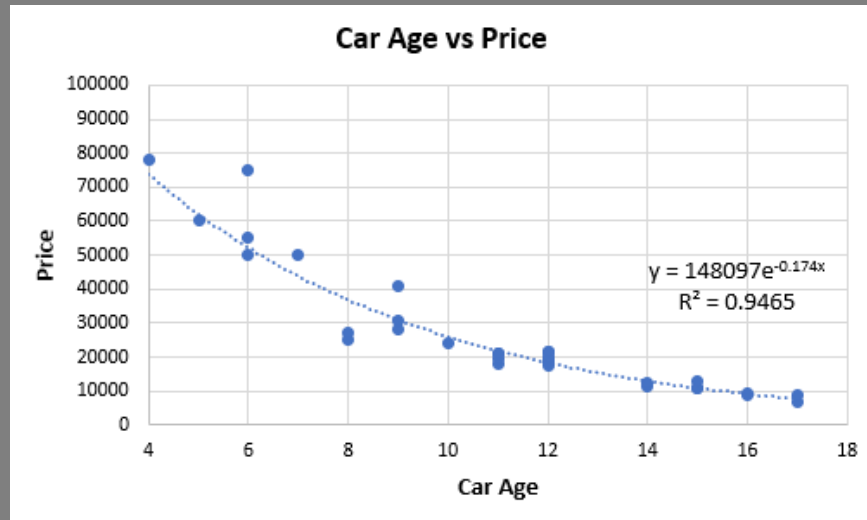


Ideally: low bias and low variance

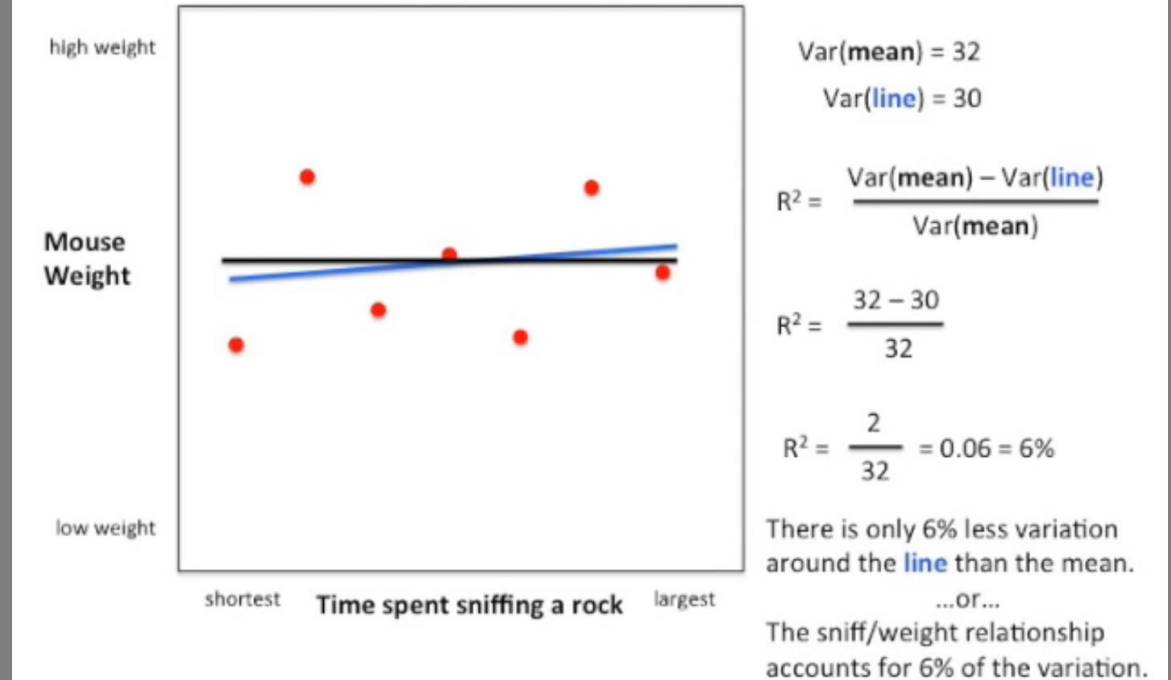
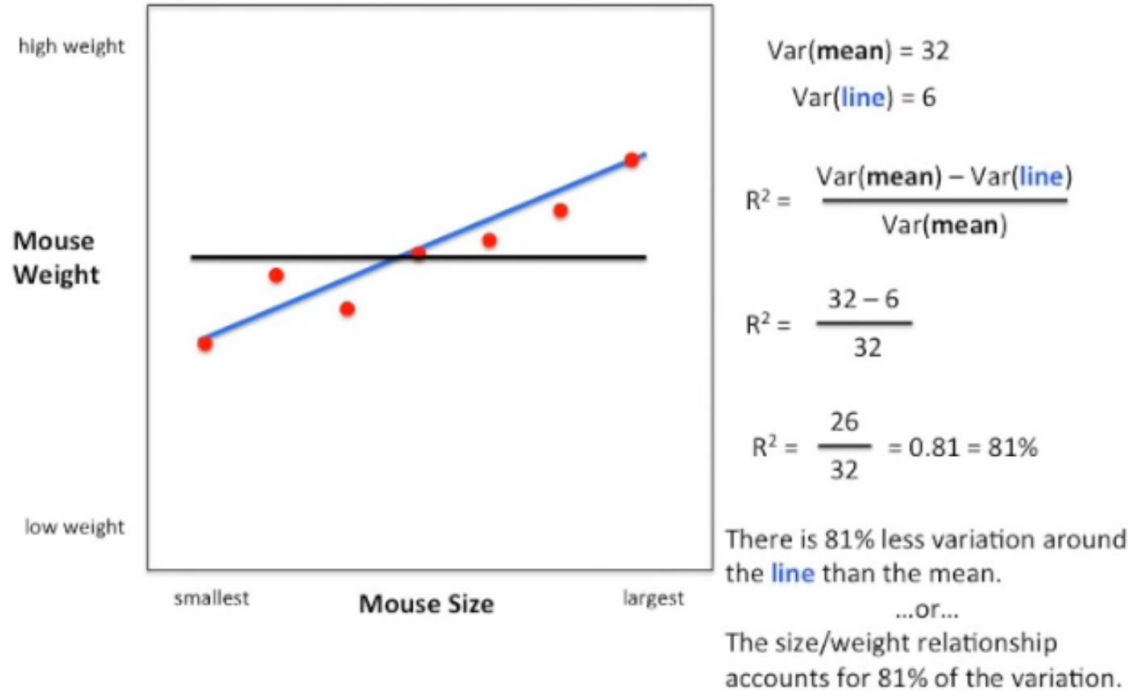


R-squared

- R-squared is a goodness-of-fit measure for linear regression models
- It tells how much two variables are correlated. $R^2 = 1$ means two variables are perfectly correlated. $R^2 = 0$ means that two variables are not correlated



R-squared



BUT! Be careful! High R-squared can also mean the model overfits

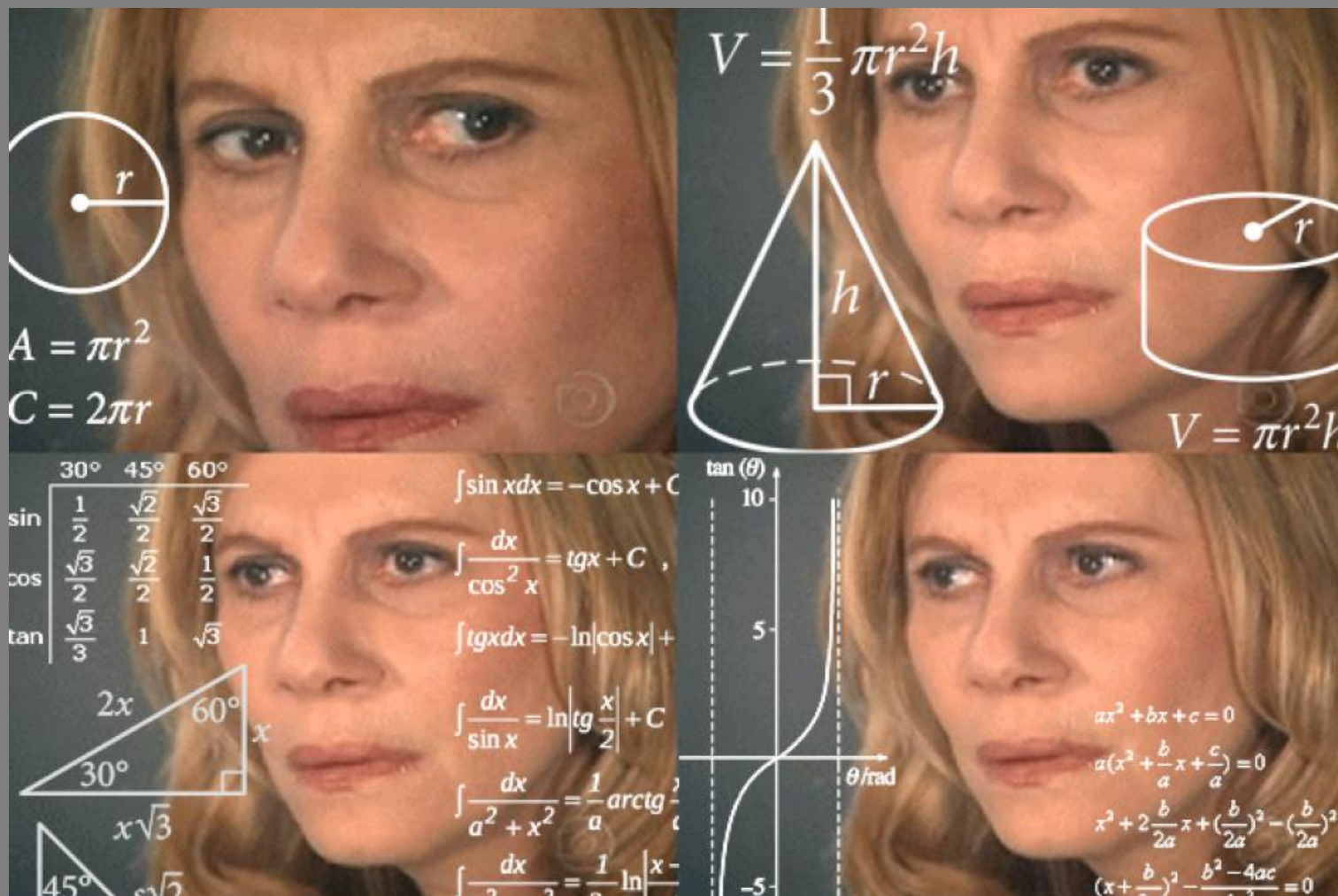
```
[16]: correlation = airbnb_housing.corr()  
      correlation["price"].sort_values(ascending=False)
```

```
[16]: price                1.000000  
      availability_365      0.081829  
      calculated_host_listings_count  0.057472  
      minimum_nights       0.042799  
      latitude             0.033939  
      host_id              0.015309  
      id                  0.010619  
      reviews_per_month   -0.030608  
      number_of_reviews   -0.047954  
      longitude           -0.150019  
      Name: price, dtype: float64
```

Correlation between housing price and other features. Note: this is correlation (R), not R^2 .

Basically: when working with dataset, consider features correlation. It's up to you to drop a certain feature if you believe it does not contribute to the prediction model

// R-squared is for the predicting line correlation to all feature, while R is for correlation between 2 features



But you know, I learned something today



- Linear regression can find the trend line so we can make a prediction
- The model should not be too overfitted or underfitted
- R-square scored is a way to find model's accuracy



<https://colab.research.google.com/drive/1YxHMiHZnEiwnHRFJh2SL83Bj-6eHJzQR?usp=sharing>