



Baron Schwartz ✅ @xaprb · 15. Nov.

When you're fundraising, it's AI

When you're hiring, it's ML

When you're implementing, it's linear regression

When you're debugging, it's printf()

🌐 Original (Englisch) übersetzen



71



4,4 Tsd.



10 Tsd.



Discussion print("C")

Today we will learn:

Today we will learn:



Regression

Lina

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]

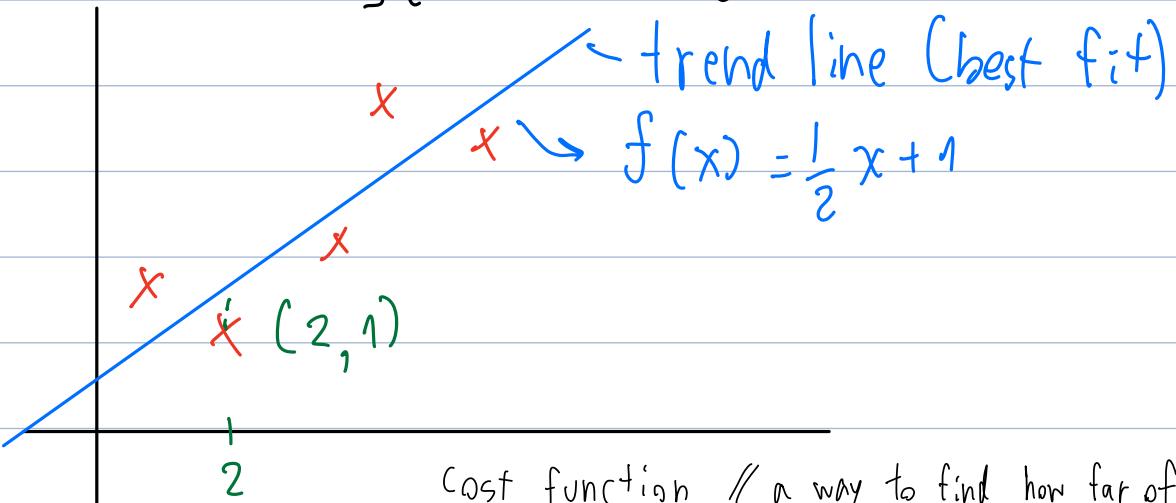
Let's look inside the math behind
linear rEVANGELION



[2D] : point : (x, y)

line : $y = ax + b$

$$f(x) = ax + b$$



Cost function // a way to find how far off//
 (loss) \uparrow builds on top
 // distance //

$$(a - b)^2$$

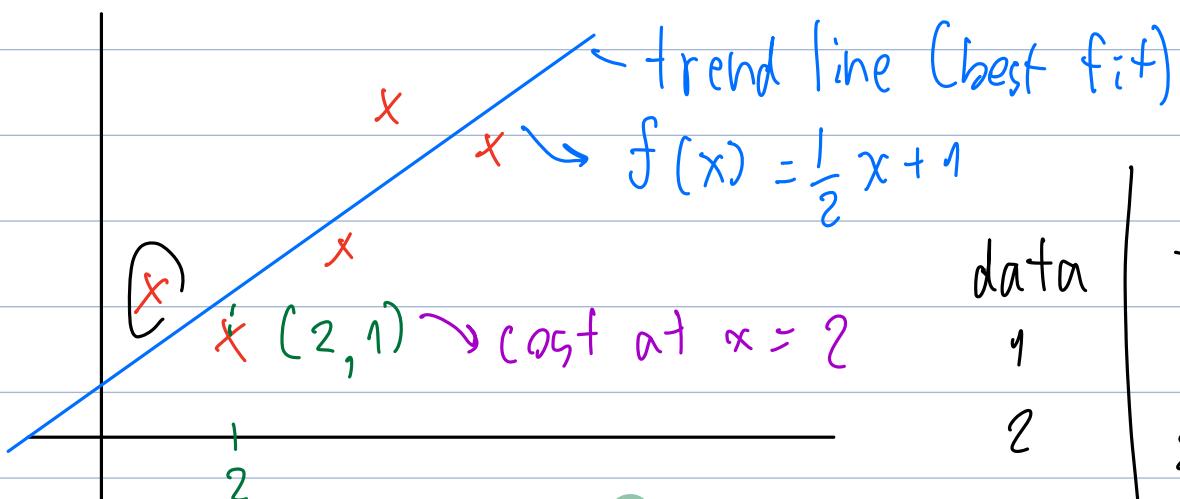
Cost function $(a, b) = (\text{predicted} - \text{observed})^2$

$$\begin{aligned} //y = ax + b// &= (\hat{y} - y)^2 \\ &= ((ax + b) - y)^2 \end{aligned}$$

↑ observed

Ex Find cost at $x = 2$

$$\begin{aligned} \text{Cost}(a, b) &= (\text{predicted} - \text{observed})^2 \\ &= \left(\left(\frac{1}{2}x + 1 \right) - 1 \right)^2 \end{aligned}$$



data	x	y
1	1	2.5
2	2	1
3	:	:
4	:	:

$$\text{Cost}(a, b) = (\text{predicted} - \text{observed})^2$$

$$= \left(\left(\frac{1}{2}x + 1 \right) - 1 \right)^2$$

$$= \left[\left(\frac{1}{2}(2) + 1 \right) - 1 \right]^2$$

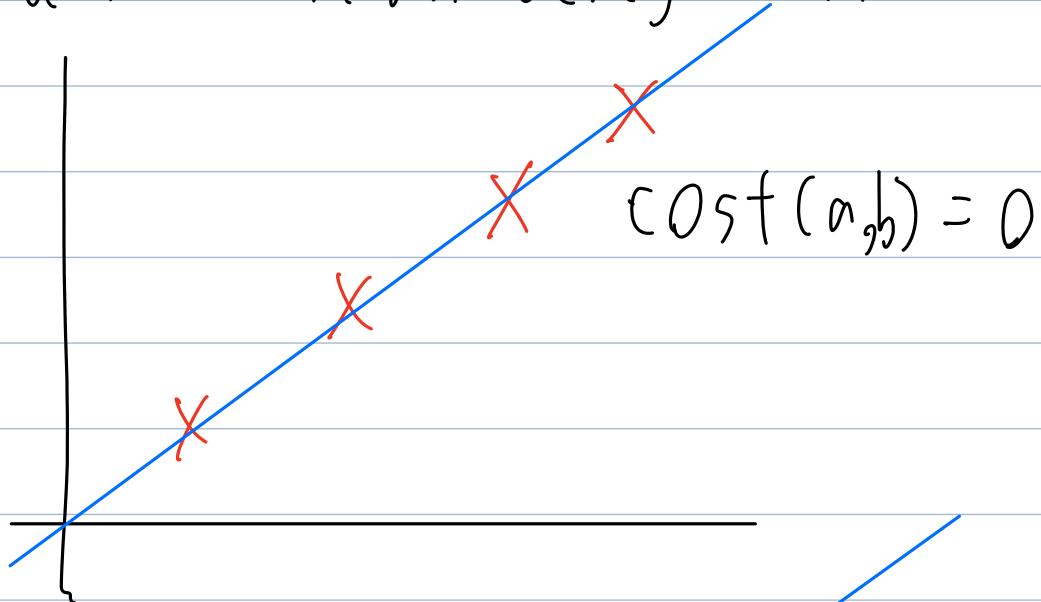
$$= (2 - 1)^2 = 1$$

$$\text{cost}(a, b) = \sum_{i=1}^n (\text{predicted}_i - \text{observed}_i)^2$$

Note: ppl use MSE (Mean-square-error)

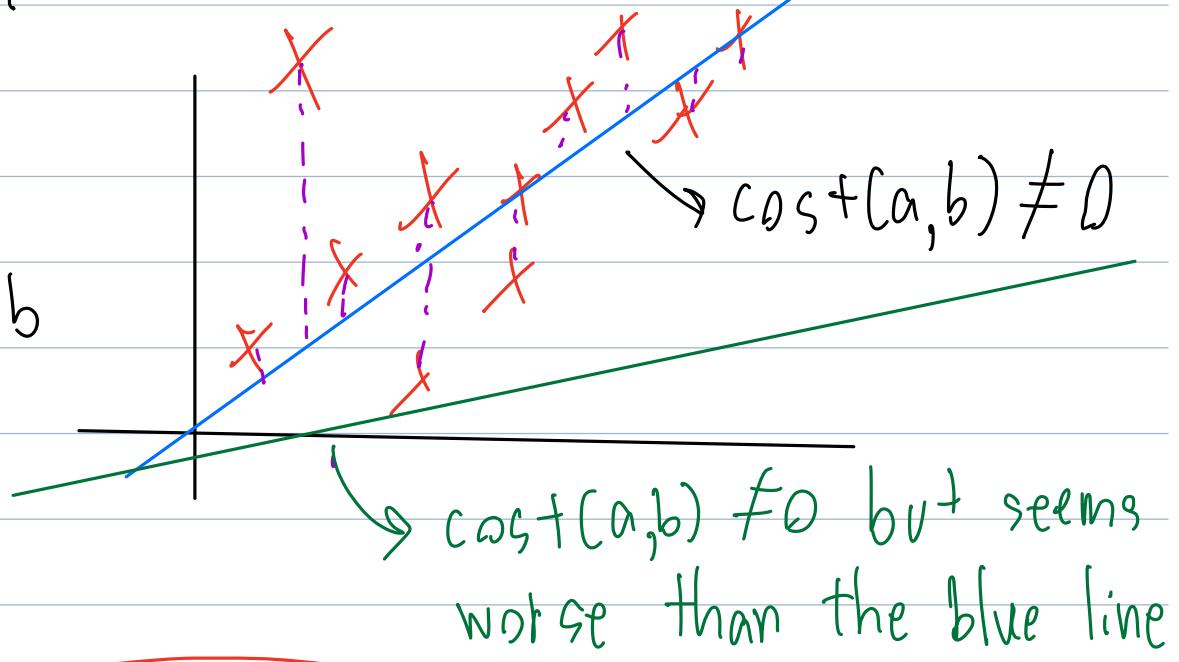
$$\text{MSE} = \text{Cost}(a, b) = \frac{1}{2n} \sum_{i=1}^n (\text{predicted}_i - \text{observed}_i)^2$$

- Ideal data & trendline (regression)



- Reality

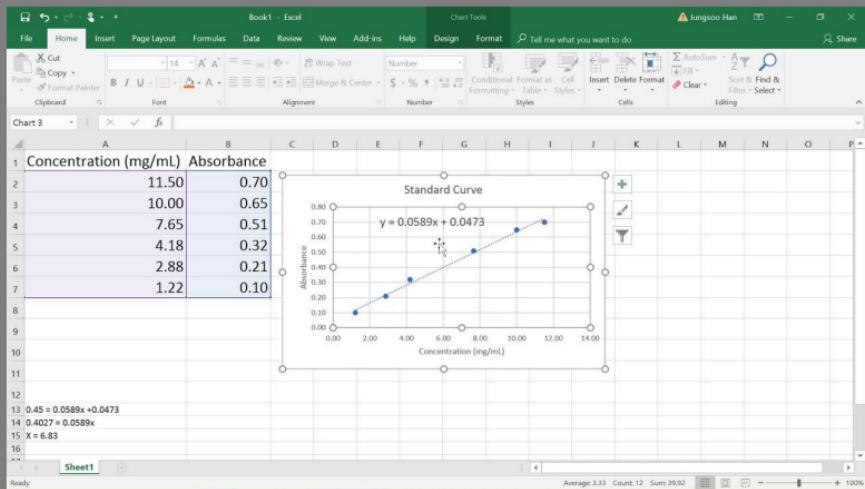
$$f(x) = ax + b$$



Goal: We wanna find a line (a, b) s.t. $\text{cost}(a, b)$ give us a value as close to zero

s.t. \rightarrow such that

Remember this?



This is him now

$f(x) = b + ax$ in example
 $h_{\theta}(x) = \theta_0 + \theta_1 x$

θ_0, θ_1 b, a predicted vs observed

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

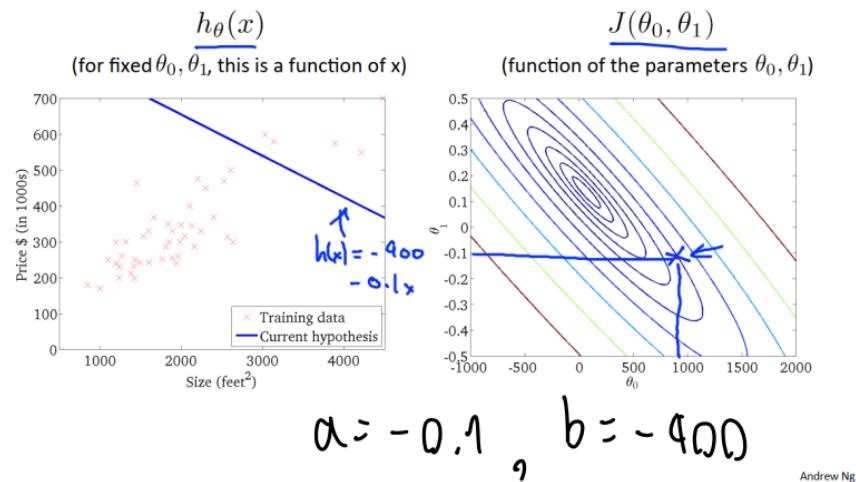
Goal: $\downarrow \min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ cost function

Andrew Ng

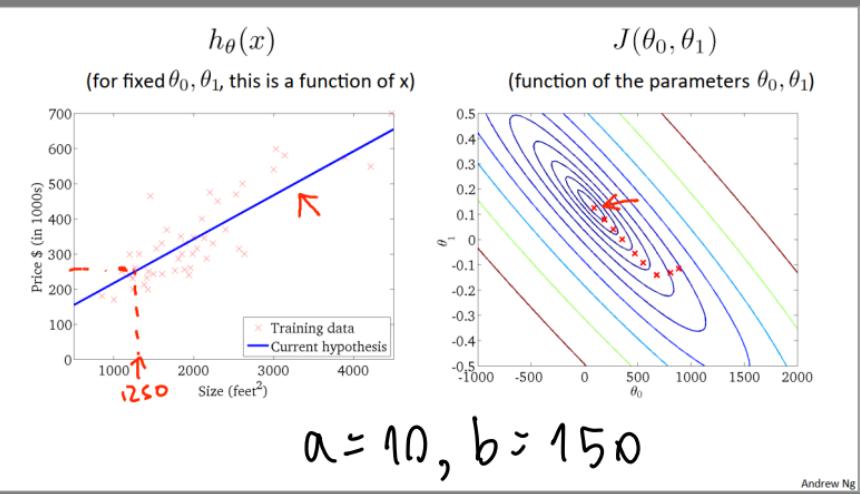
Feels old yet?

Credit: Andrew Ng

But how do you get from this



To this



Credit: Andrew Ng

Introducing: Gradient Descent



Allow us to introduce ourselves,

predicted : $f(x) = ax + b$

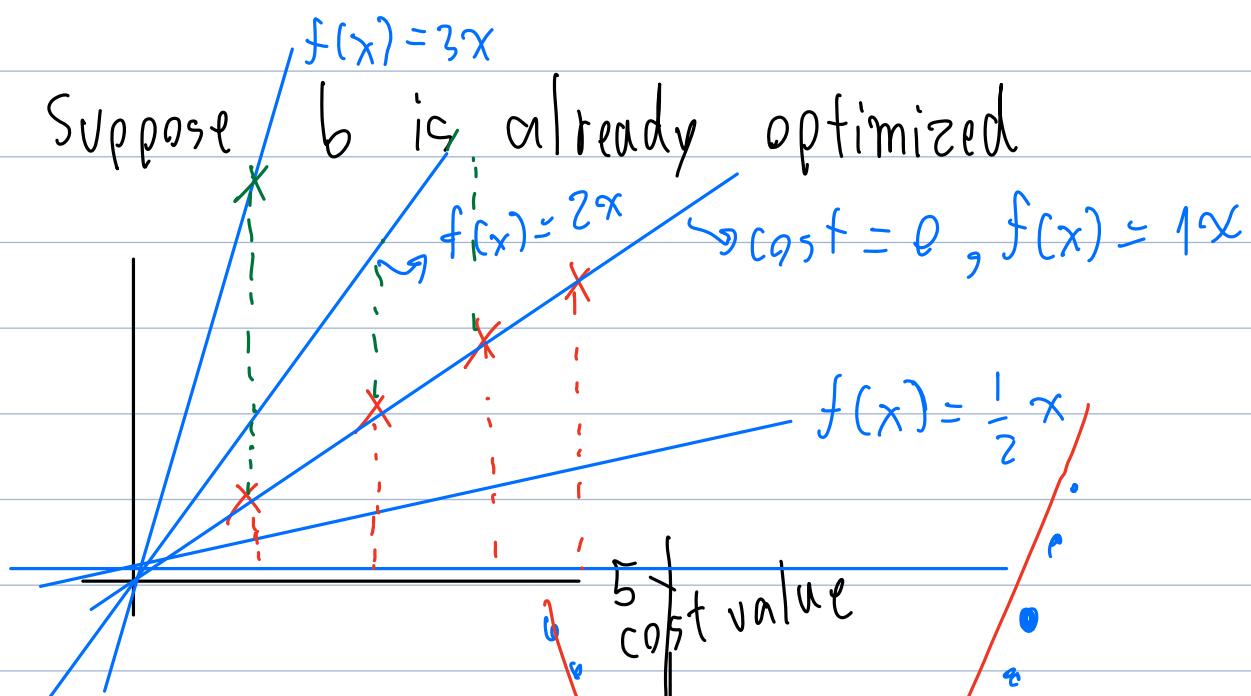
as lowest
as possible

Goal : change a, b s.t. $\text{cost}(a, b) \approx 0$

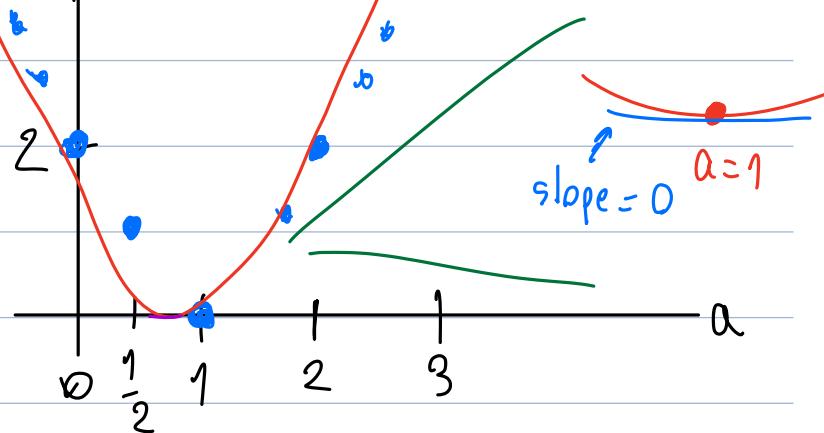
$$a_{(\text{new})} = a_{(\text{old})} - \text{somevalue}$$

$$b_{(\text{new})} = b_{(\text{old})} - \text{somevalue}$$

Ex Suppose b is already optimized

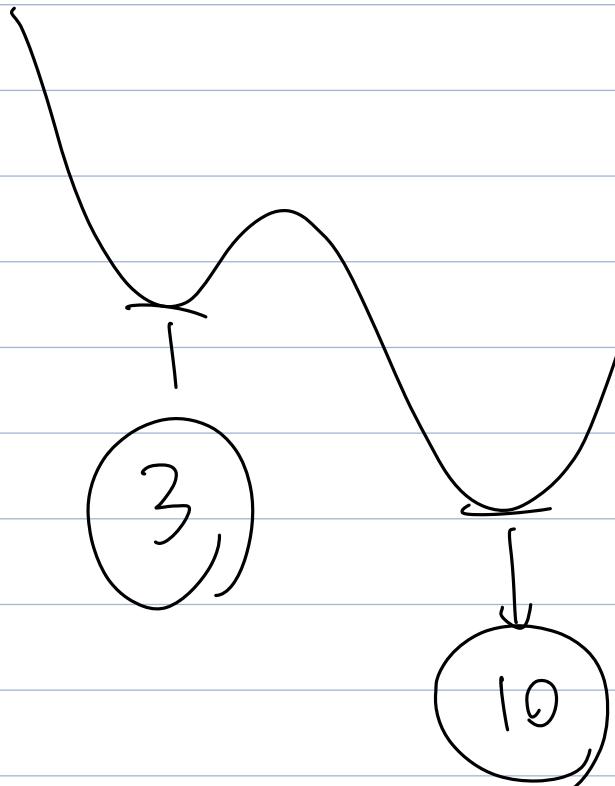


cost graph:
of slope $\rightarrow a$



If you want to find a minimum of a graph
use derivative

Derivative \rightarrow finding slope of a graph



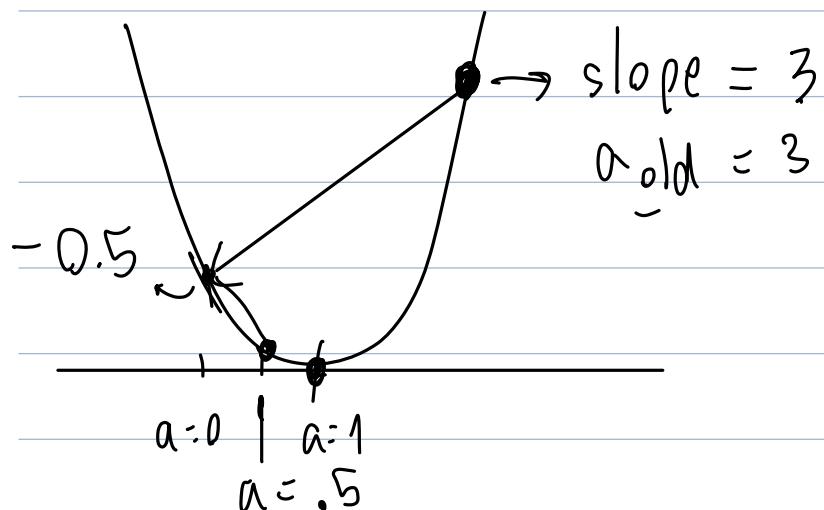
\rightarrow take derivative
at every point,
then find
which give you
 $= 0$

$$\text{slope} = \frac{dy}{da} \stackrel{\Delta}{=} 0$$

$$x^n \rightarrow n(x)^{n-1}$$

$\cos(x) \rightarrow -\sin(x)$

$$\text{// Some value} = \lambda \text{slope}$$



$$a_{\text{new}} = 0 - (1)(-0.5)$$
$$= 0.5$$

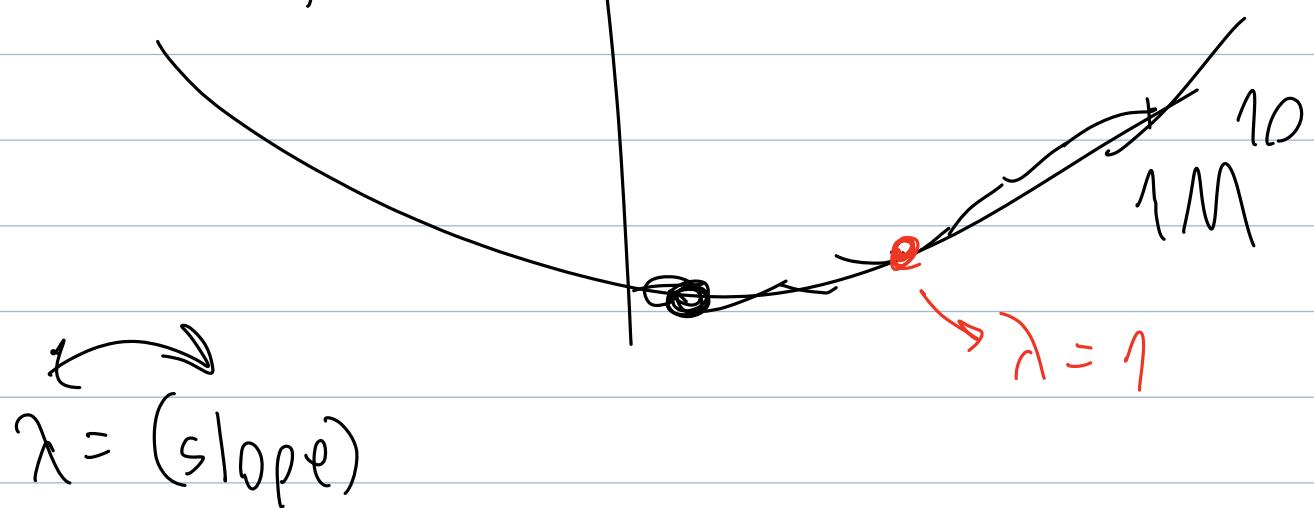
Generalized : Update every parameters

$$a_{\text{new}} = a_{\text{old}} - \lambda \text{slope}(a)$$

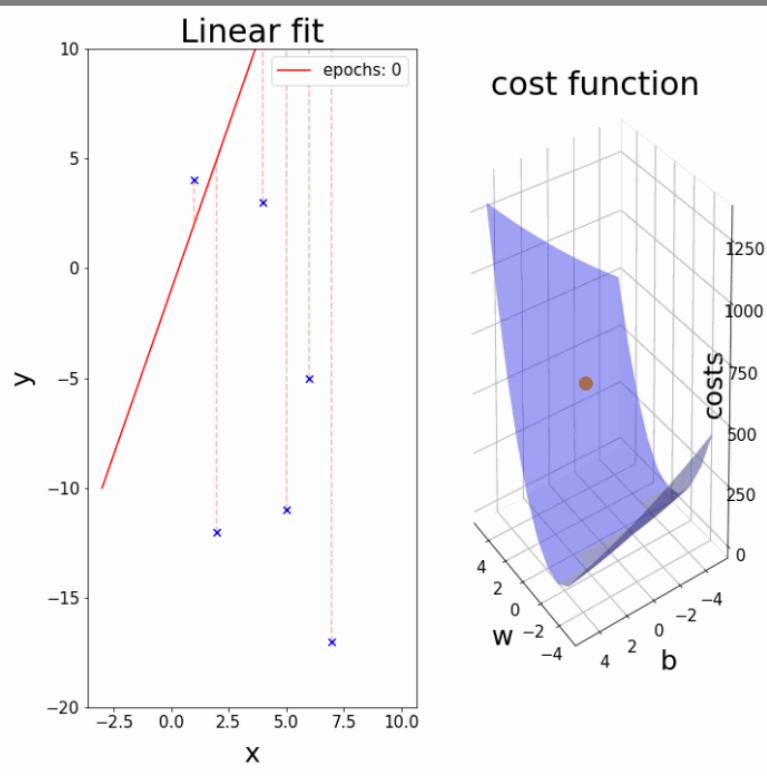
$$b_{\text{new}} = b_{\text{old}} - \lambda \text{slope}(b)$$

$$\lambda \leq 1$$

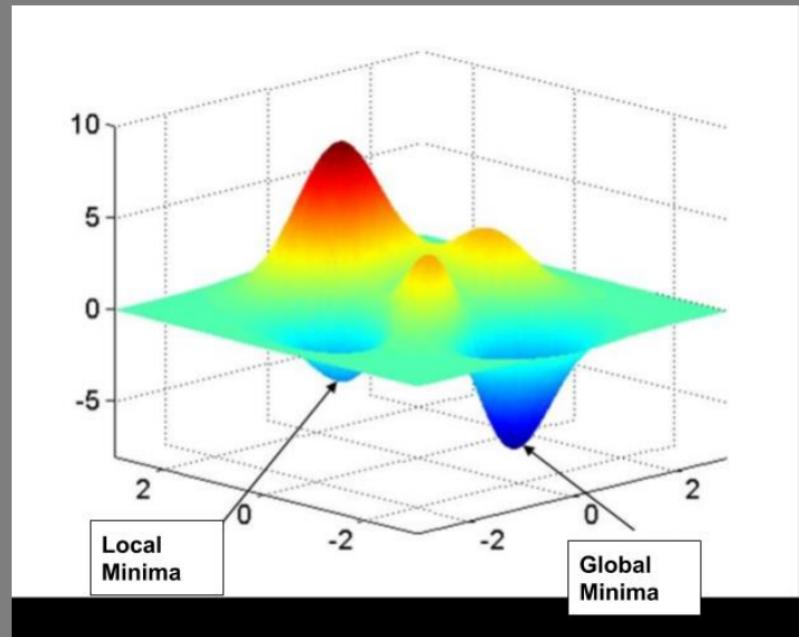
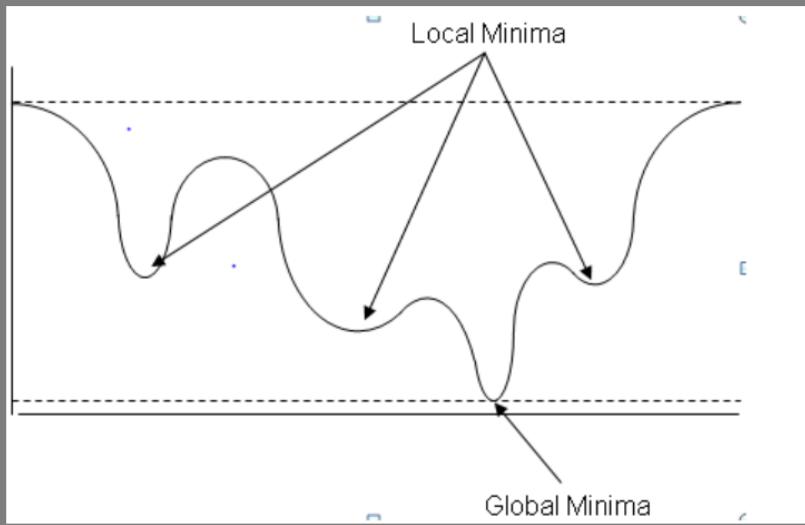
$$\lambda = 500,000$$

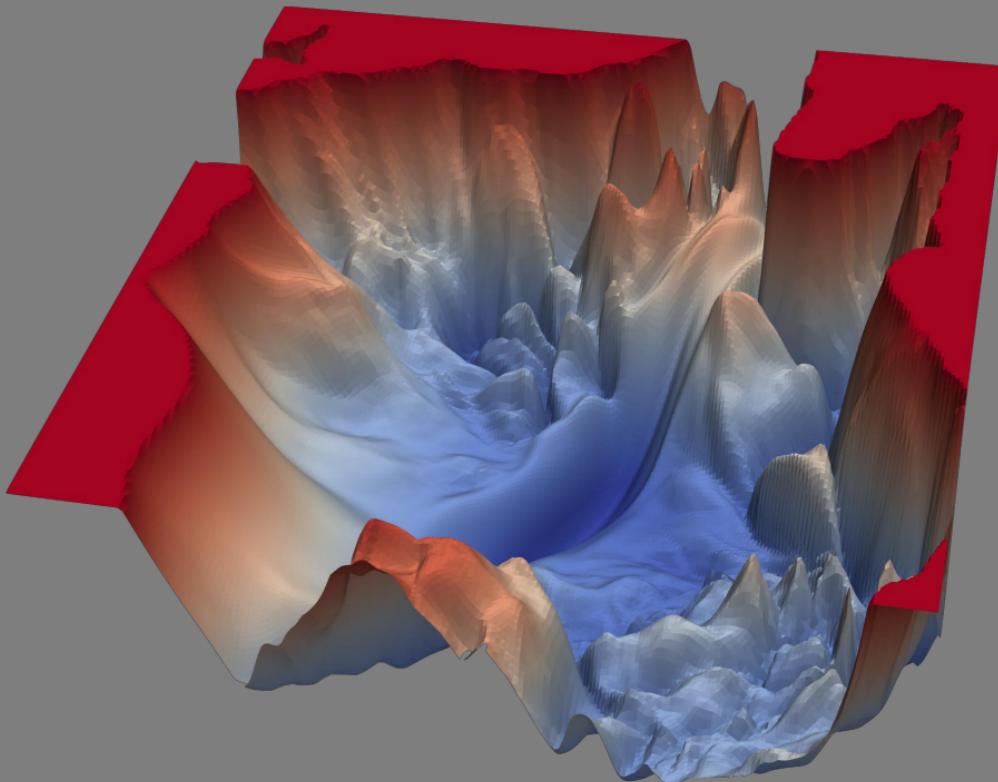


$$a_{\text{old}} - \lambda \cdot \text{slope} \rightarrow a_{\text{old}} - \text{slope} \cdot \text{slope}$$



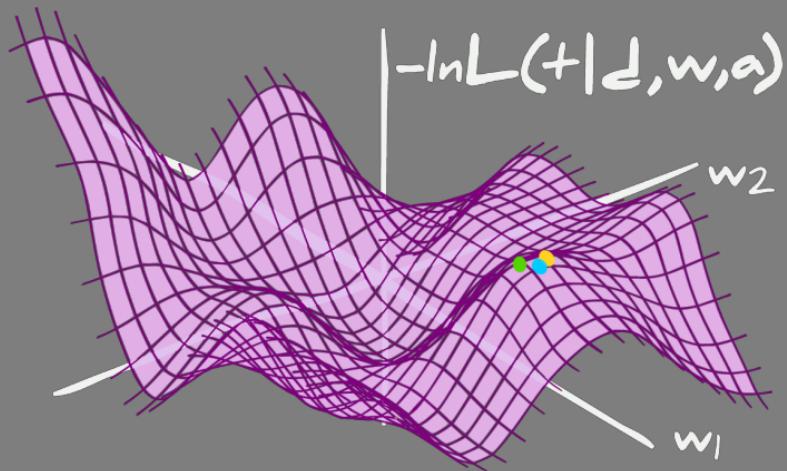
Issue: what if we cannot get the lowest cost function





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 REYgression with multiple features

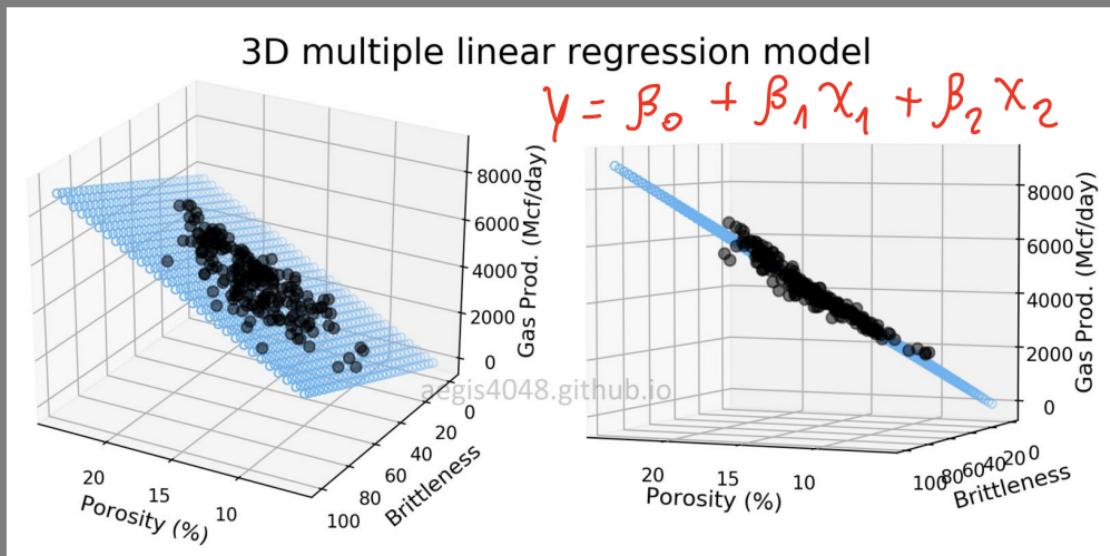
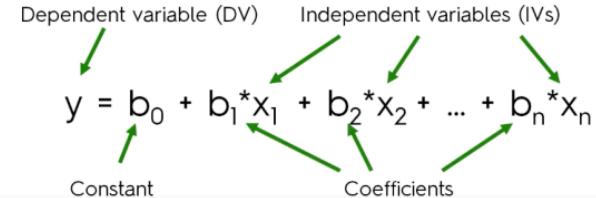


Simple Linear Regression

$$y = b_0 + b_1 * x_1$$

$$b_0 + \alpha = \beta_1$$

Multiple Linear Regression



If we have this...

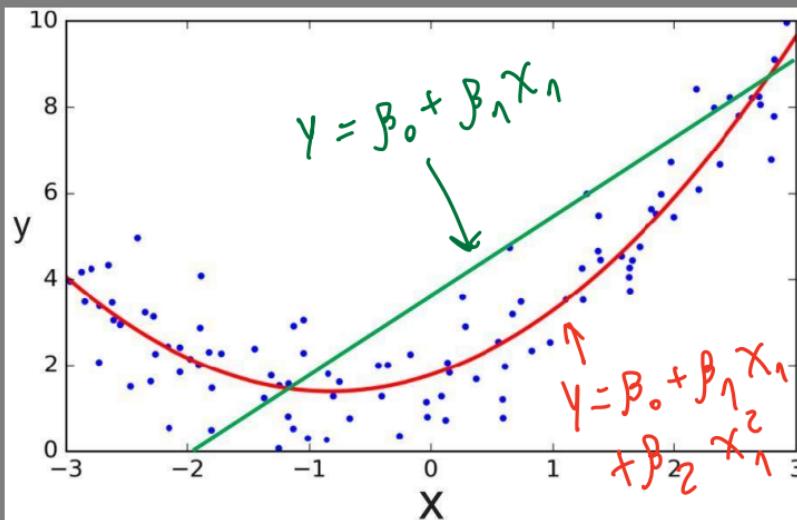
Multiple
Linear
Regression

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

We can also have this!

Polynomial
Linear
Regression

$$y = b_0 + b_1 \underline{x_1} + b_2 \underline{x_1^2} + \dots + b_n \underline{x_1^n}$$



$$\uparrow \\ b_2 x_1 \cdot x_1$$

previous $\rightarrow a_{\text{new}} = a_{\text{old}} - \lambda \text{slope}, b_{\text{new}} = b_{\text{old}} - \lambda \text{slope}$

Same cost(loss) function and gradient descent

this is loss

repeat until convergence: {

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

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

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

...

$$\} \quad 2 \cdot \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_i^{(i)}$$

$$\text{loss} \rightarrow \text{MSE}$$

$$\frac{1}{m} \cdot \sum_i^m (\text{pred}_i - \text{observed}_i)^2$$

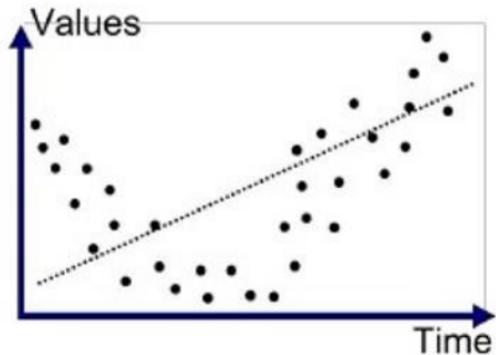
$$x^n \rightarrow n(x)^{n-1}$$

$$f(x) = ax + b$$

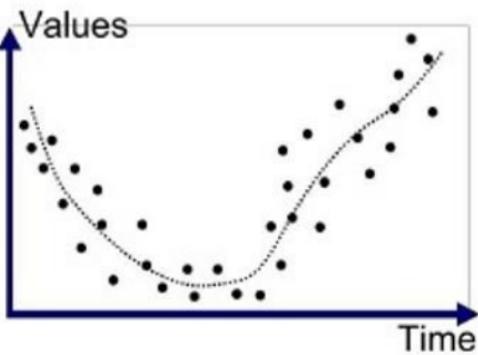
$$\frac{1}{m} \sum_i^m (\text{pred}_i - \text{observed}_i)^2$$

derivative wr.t. a := $2 \cdot \frac{1}{m} \sum_i^m (\text{pred}_i - \text{observed}_i) x$

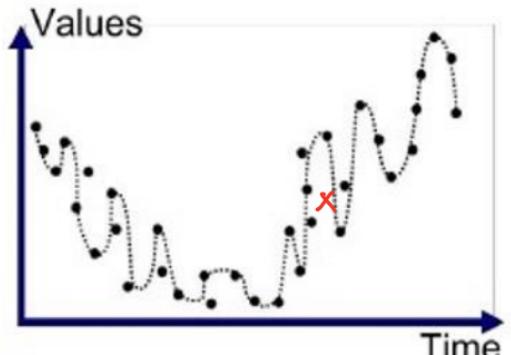
The good, the bad, and the ugly



Underfitted

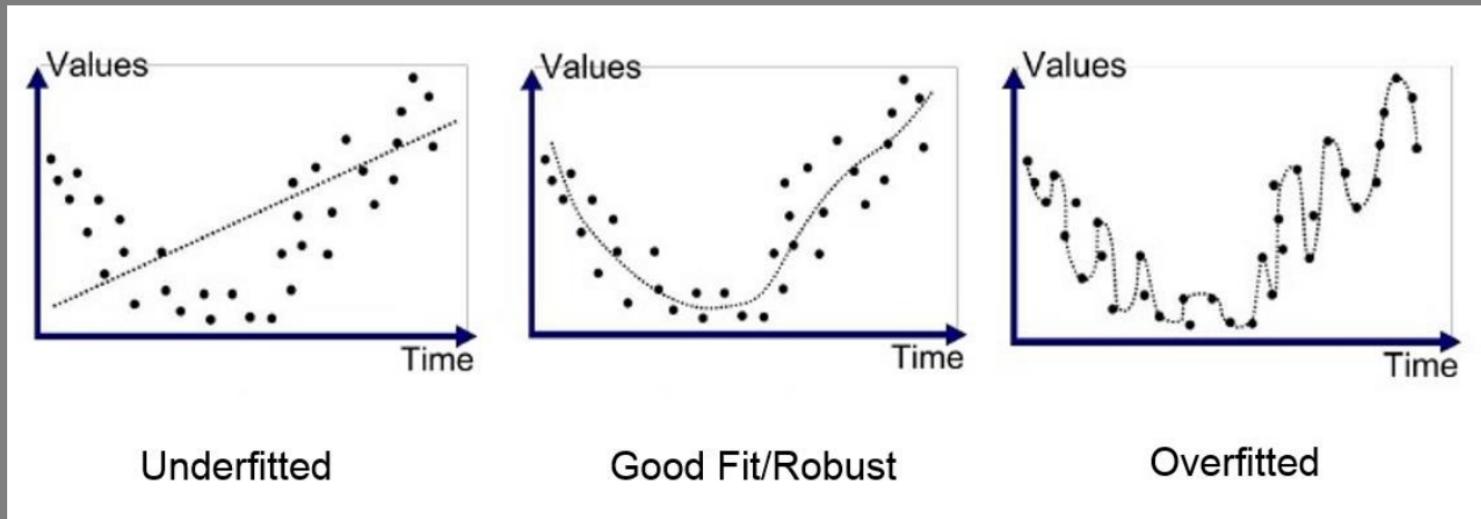


Good Fit/R robust



Overfitted

Terminology and Tips



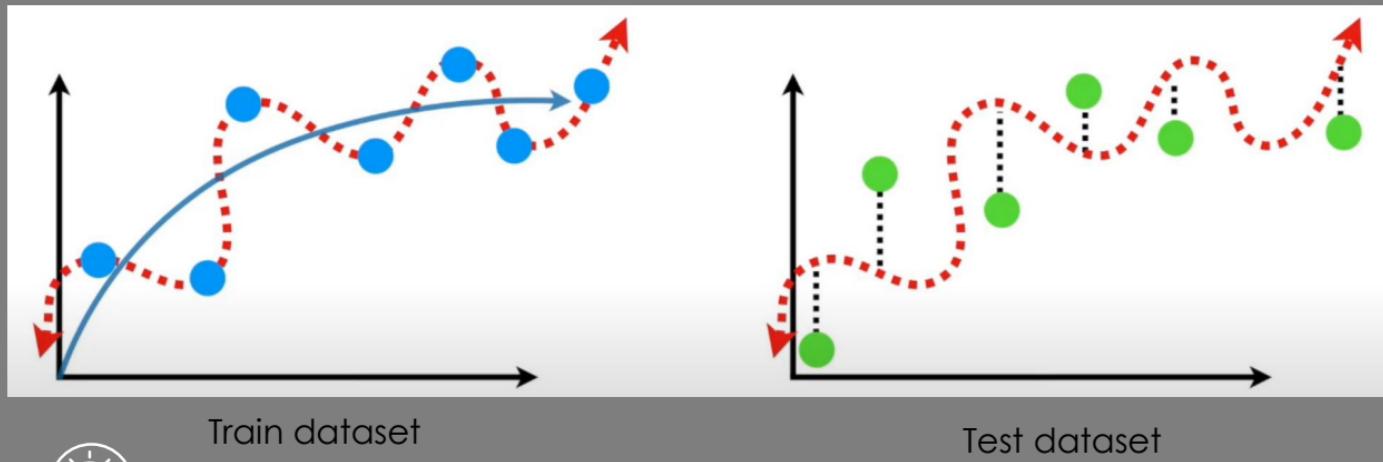
We use the term “bias” and “variance” as another way to explain how well the trend line(or plane if 3D) captures data

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



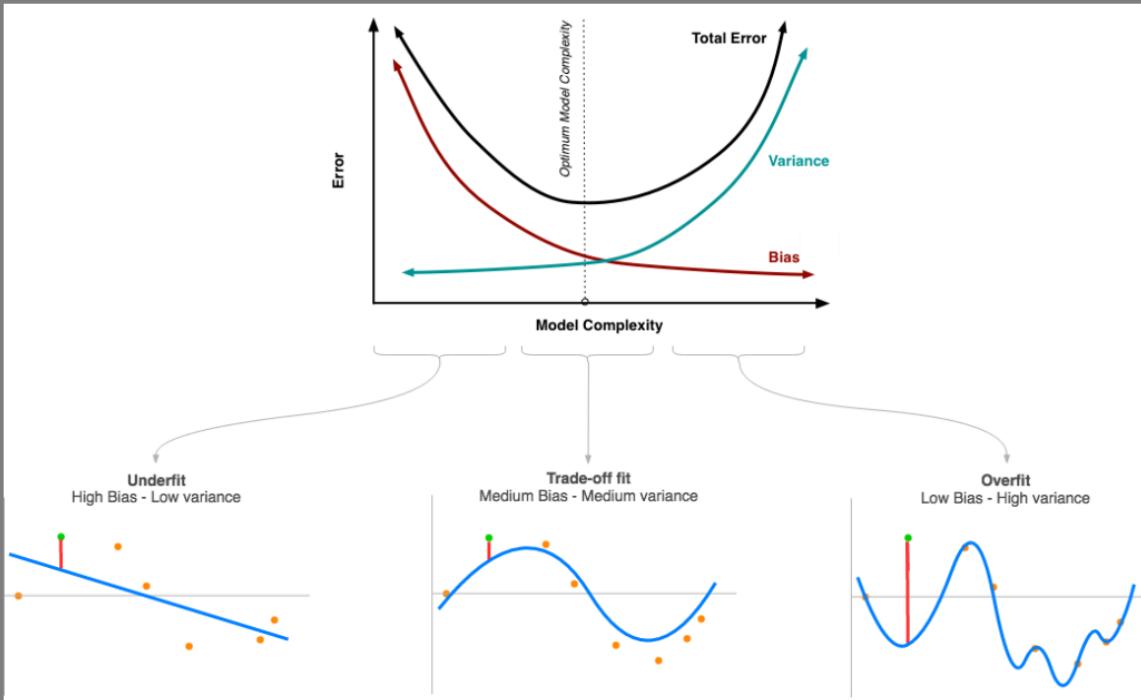
Ideally: low bias and low variance

Bias: Inability to capture the true relationship

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

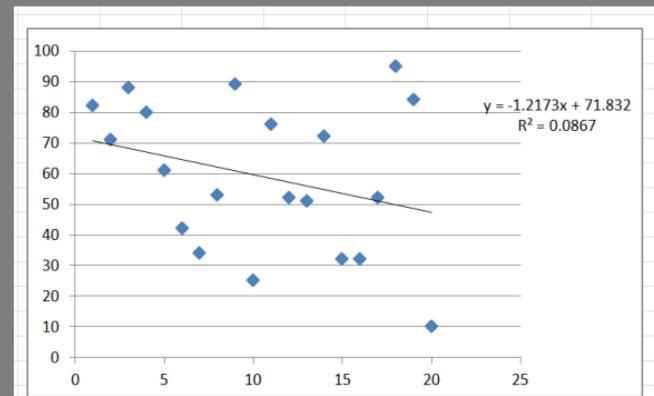
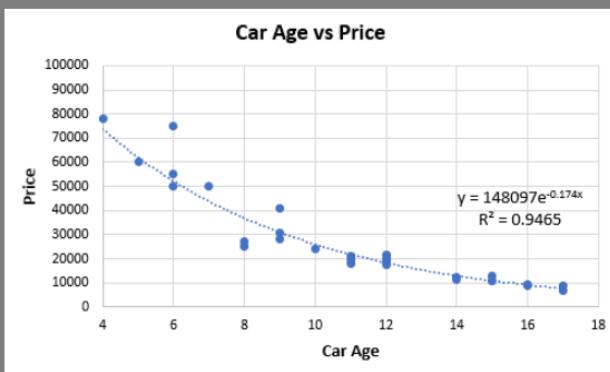


Ideally: low bias and low variance

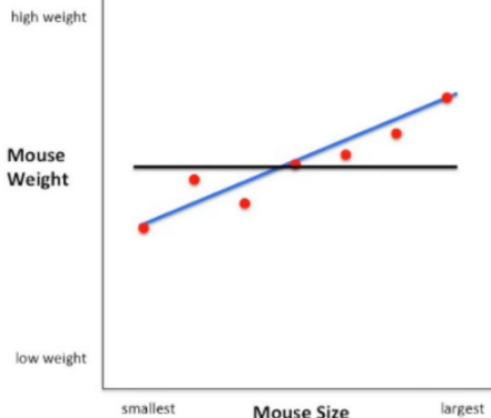


R-squared

- 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



$$\text{Var}(\text{mean}) = 32$$

$$\text{Var}(\text{line}) = 6$$

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{line})}{\text{Var}(\text{mean})}$$

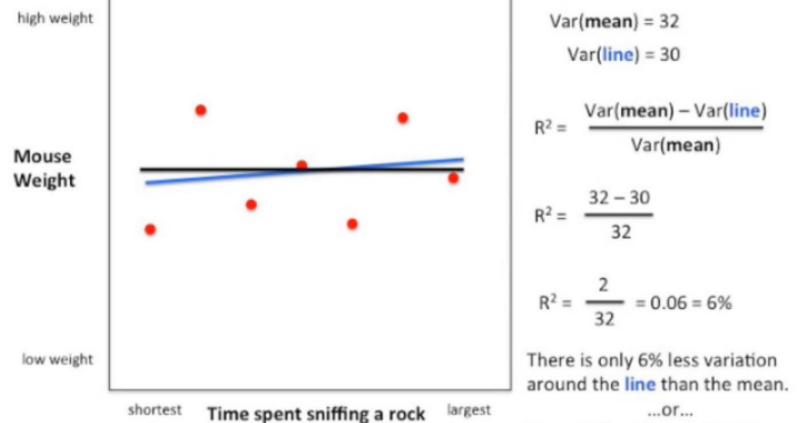
$$R^2 = \frac{32 - 6}{32}$$

$$R^2 = \frac{26}{32} = 0.81 = 81\%$$

There is 81% less variation around the **line** than the **mean**.

...or...

The size/weight relationship accounts for 81% of the variation.



$$\text{Var}(\text{mean}) = 32$$

$$\text{Var}(\text{line}) = 30$$

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{line})}{\text{Var}(\text{mean})}$$

$$R^2 = \frac{32 - 30}{32}$$

$$R^2 = \frac{2}{32} = 0.06 = 6\%$$

There is only 6% less variation around the **line** than the **mean**.

...or...

The sniff/weight relationship accounts for 6% of the variation.

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

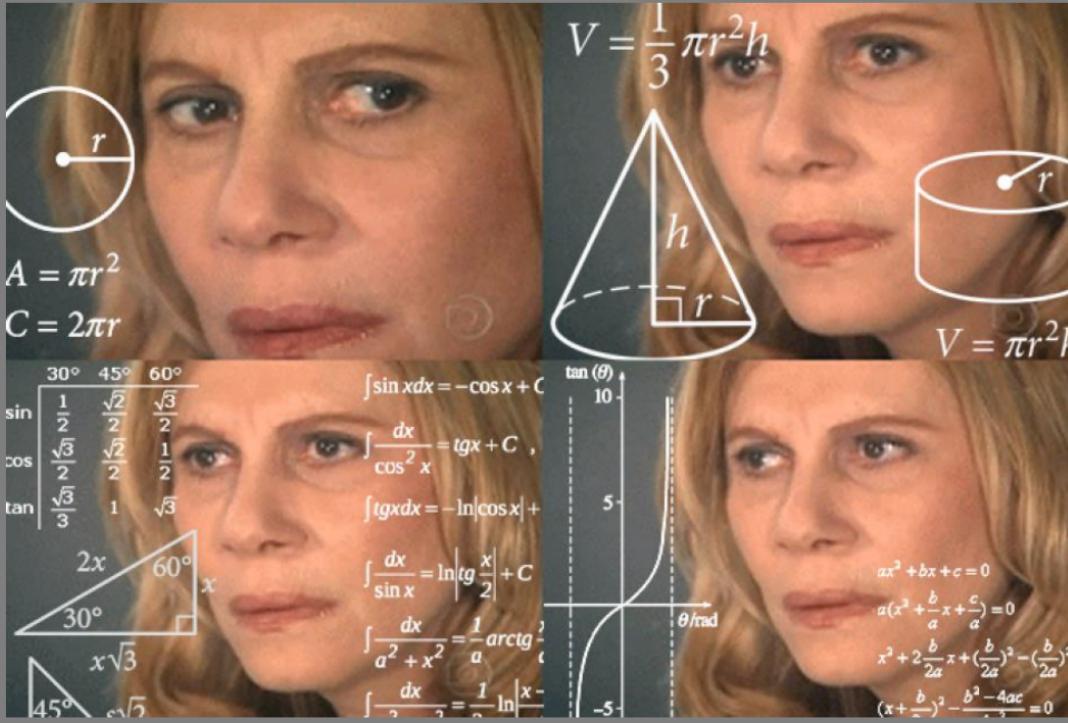
```
[16]: correlation = airbnb_housing.corr()  
correlation["price"].sort_values(ascending=False)
```

```
[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>