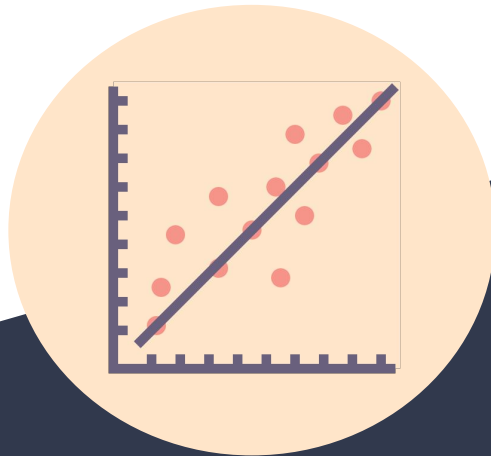


# ECE ING4

# MACHINE LEARNING

Jeremy Cohen

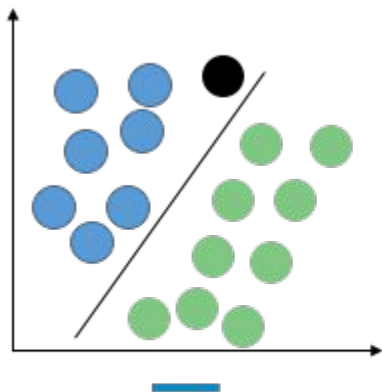


**Linear Regression**

# Week 1 Review

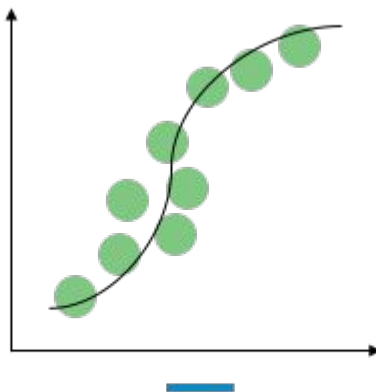
# Types of Machine Learning

## CLASSIFICATION



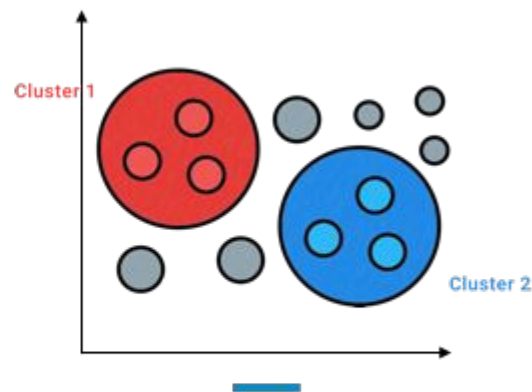
- Supervised
- Output is a discrete number (0,1,2, ...), (SPAM/NOT SPAM), ...

## REGRESSION



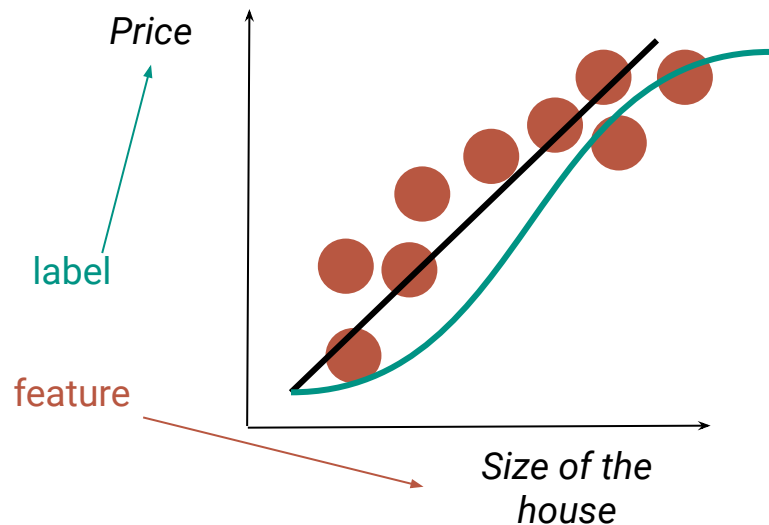
- Supervised
- Output is a continuous number

## CLUSTERING



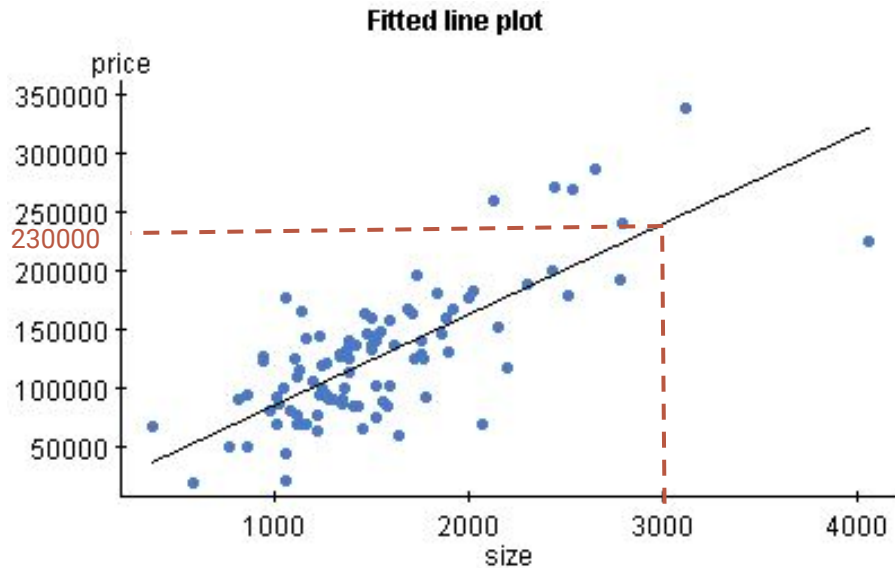
- Unsupervised
- Outputs are clusters

# Regression



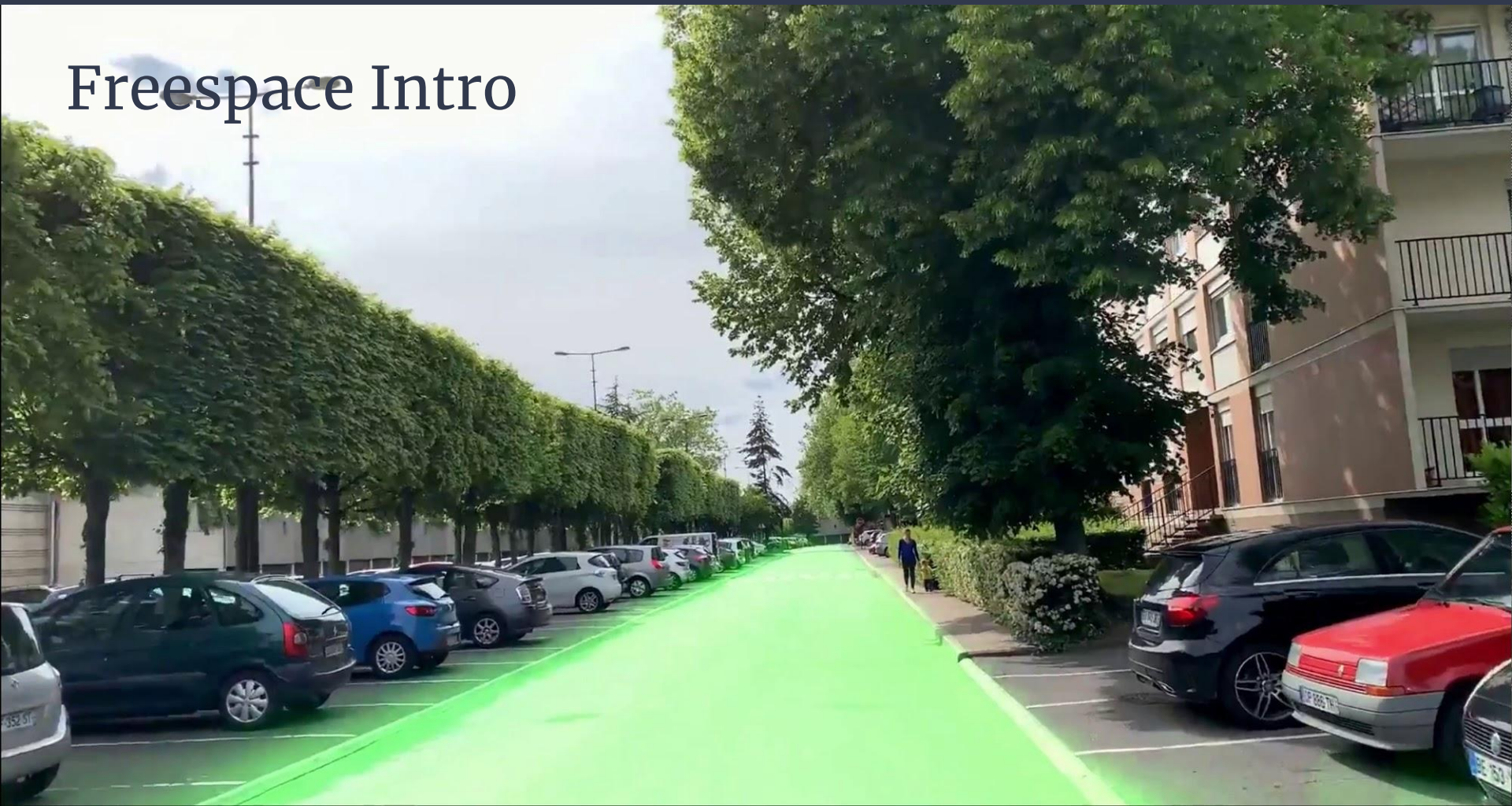
THE GOAL IS TO DETERMINE THE  
LINE OR CURVE THAT BEST FITS  
THE DATA

# Regression

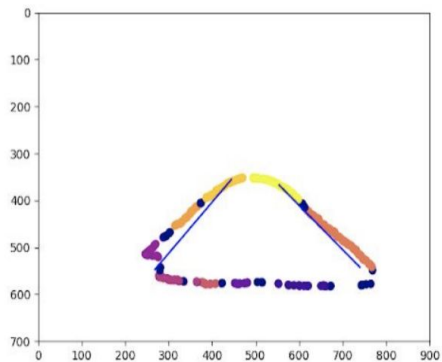


**THE GOAL IS TO DETERMINE THE  
LINE OR CURVE THAT BEST FITS  
THE DATA**

# Freespace Intro

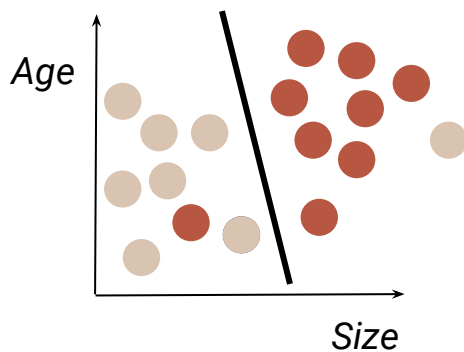


# Regression in Practice

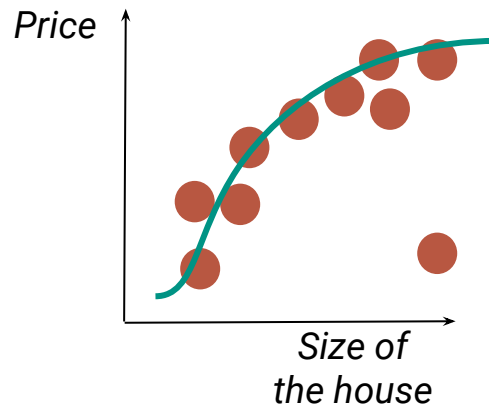


# Outliers

## CLASSIFICATION



## REGRESSION



**OUTLIERS ARE TO BE  
REMOVED WHEN TRAINING**



# Classification vs Regression



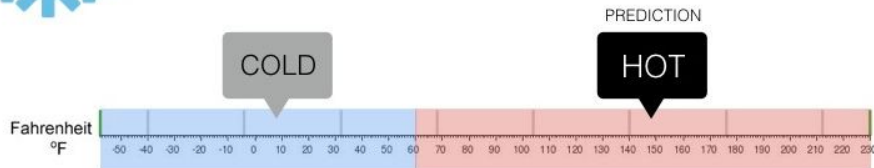
## Regression

What is the temperature going to be tomorrow?

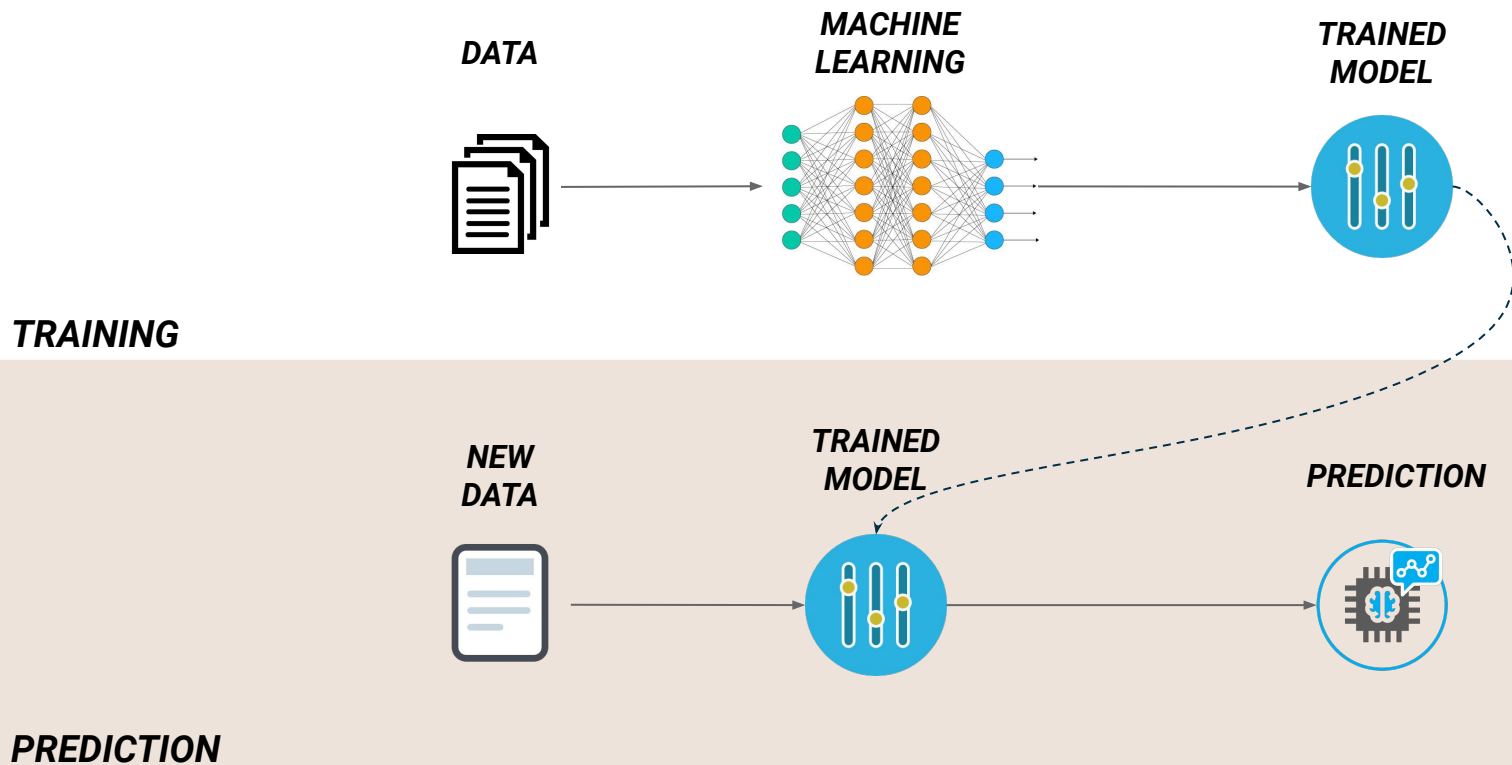


## Classification

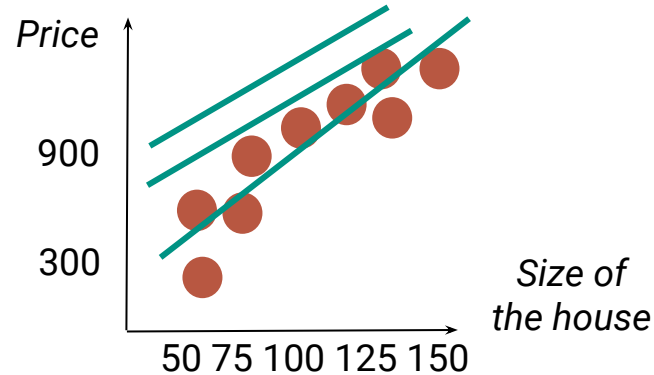
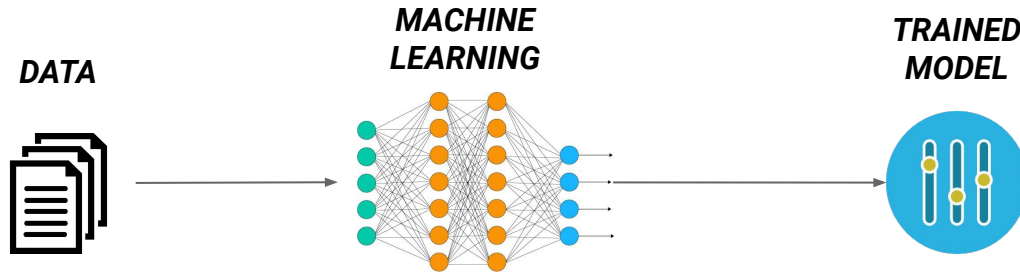
Will it be Cold or Hot tomorrow?



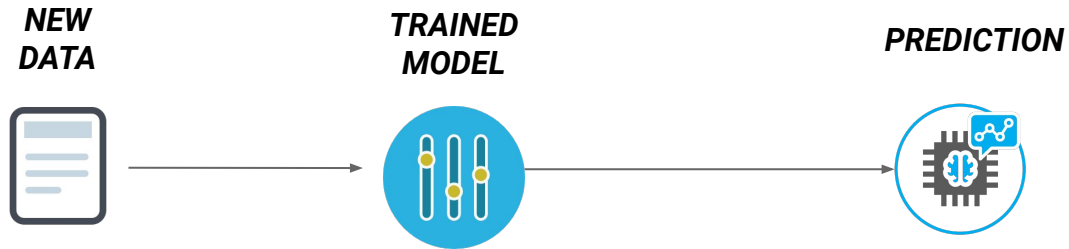
# Machine Learning Process



# Training Linear Regression

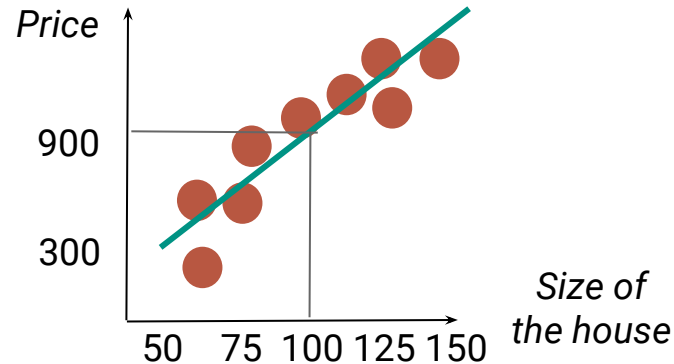


# Testing Linear Regression



I want to buy a house of  
100 square meters.  
How much will it cost?

**915 k€**



# 3 Datasets

## Training Set



~70% of the dataset

Used to **train the model**

## Validation Set



~20% of the dataset

Used to **test** the model  
and **generalize better to  
new data**

## Test Set

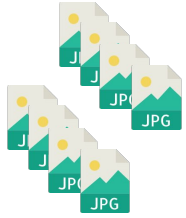


~10% of the dataset

Used **only once** when  
both accuracies are  
good and ready for  
real-world

# Train Test Process

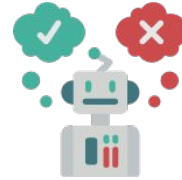
Training Set



Machine Learning  
Algorithm



TRAINING



Trained Model

TESTING



Test Set



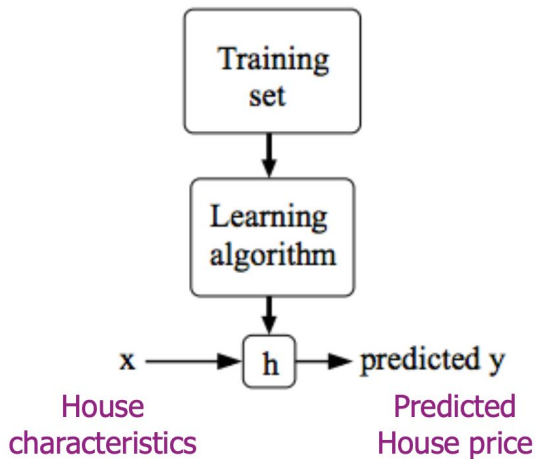
Trained Model



Prediction

# Model

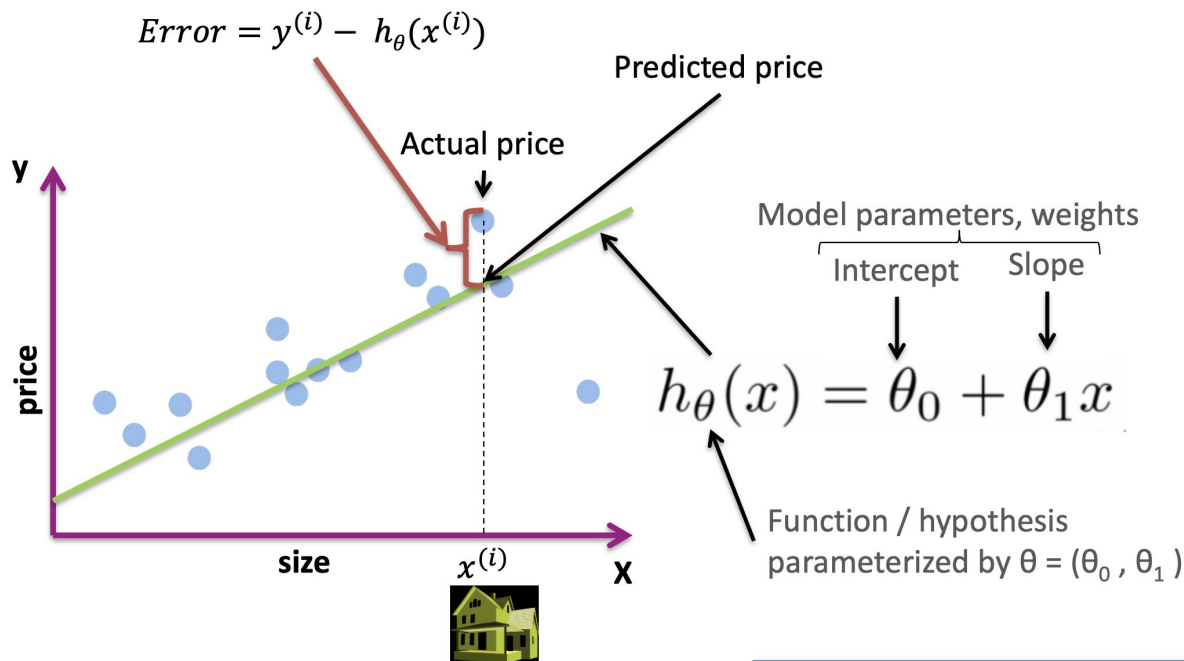
House characteristics (size, neighborhood, ...): Feature **X**  
House Price : Label **Y**



**$h: X \rightarrow Y$**

Hypothesis or function that takes as input the house's characteristics to estimate its price.

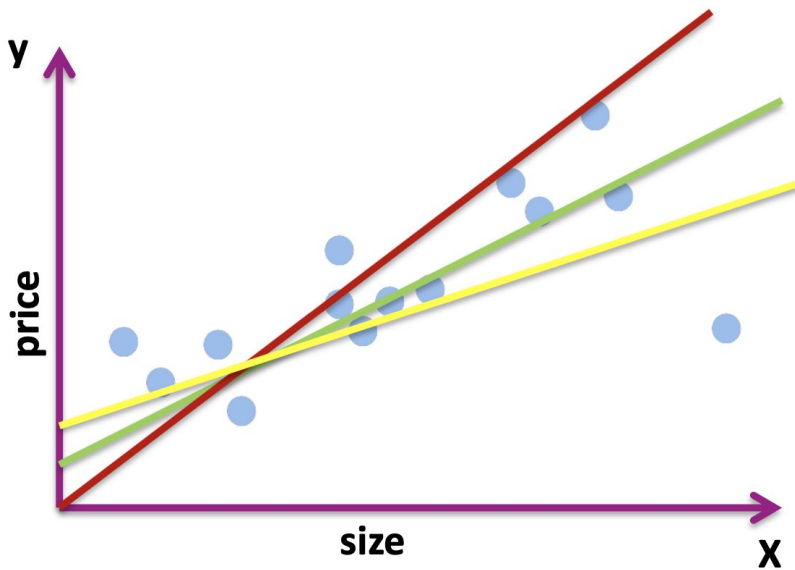
# Simple Linear Regression



$$y = \theta_0 + \theta_1 x + \varepsilon$$

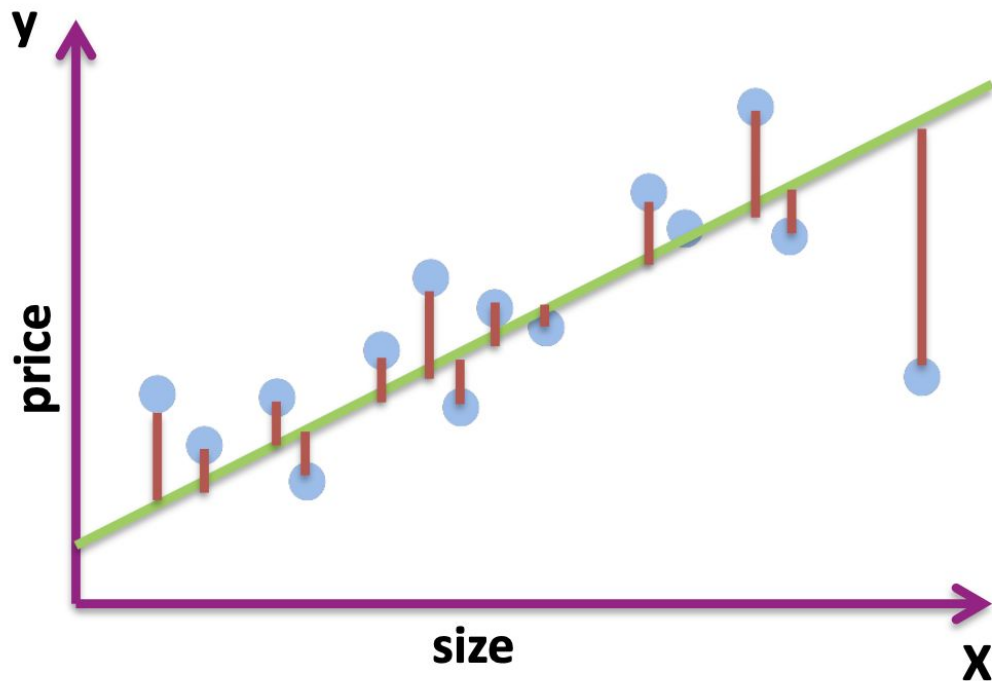


# The Best Fit



Each line has different  
parameters  $\theta = (\theta_0, \theta_1)$   
→ different errors

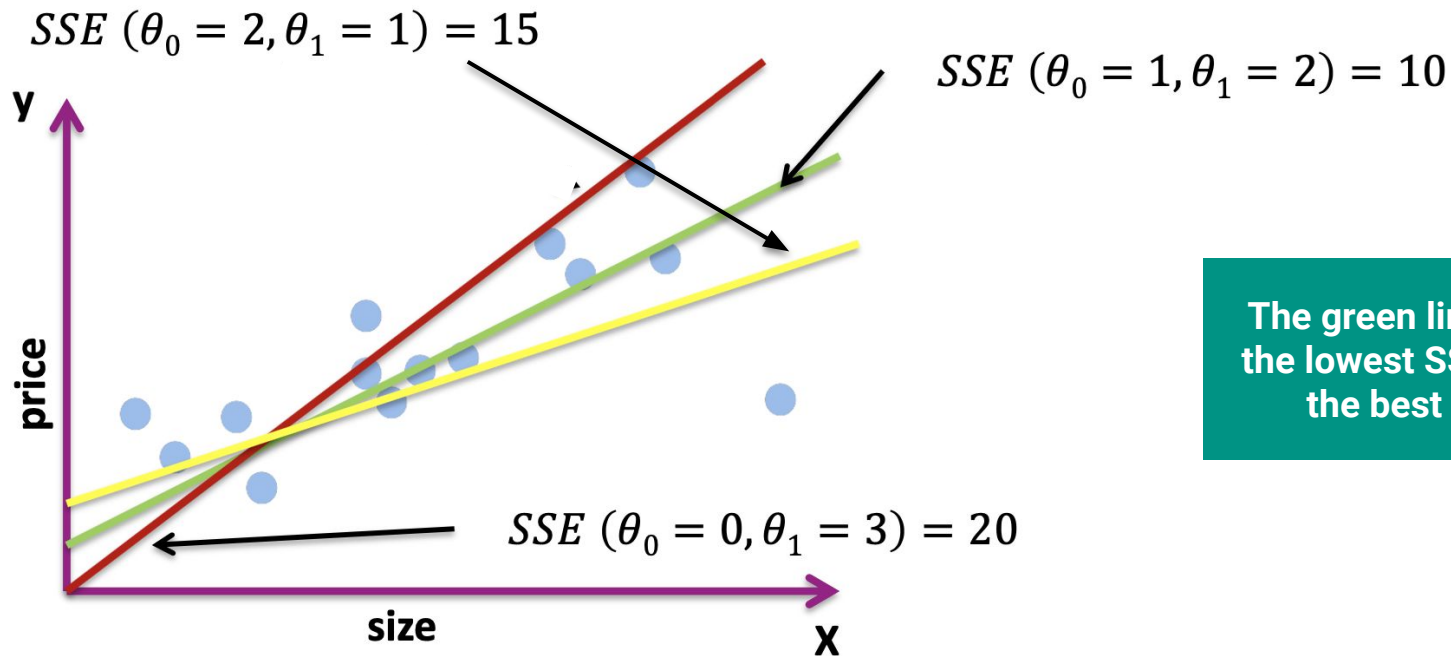
# Sum of Squared Errors (SSE)



$$SSE(\theta_0, \theta_1) =$$

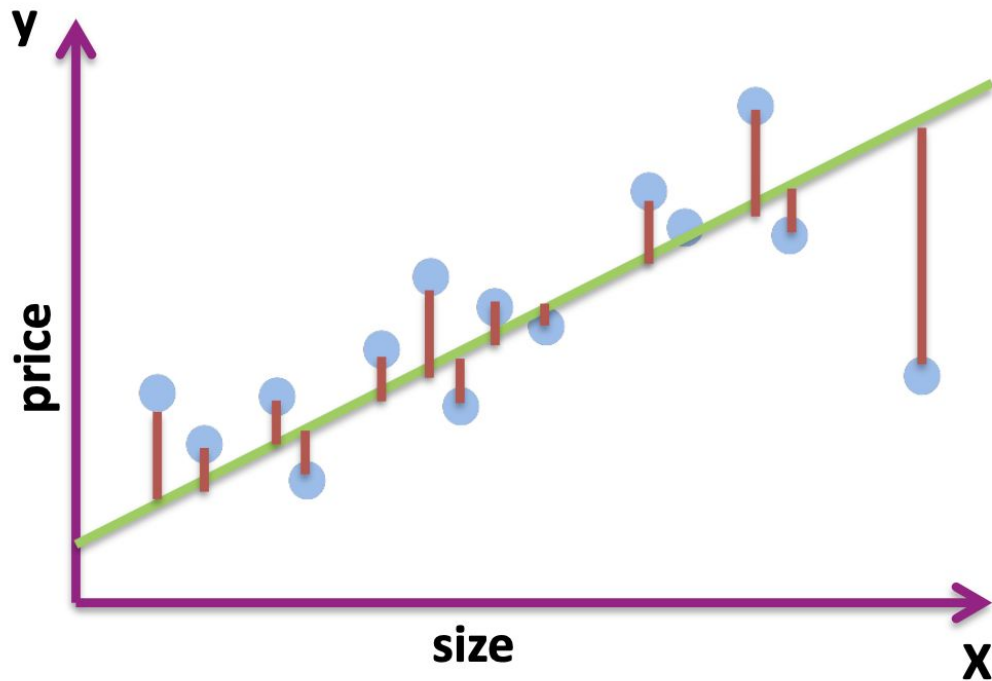
$$\begin{aligned} & (y^{(1)} - (\theta_0 + \theta_1 * x^{(1)}))^2 + \\ & (y^{(2)} - (\theta_0 + \theta_1 * x^{(2)}))^2 + \\ & \dots\dots\dots + \\ & (y^{(m)} - (\theta_0 + \theta_1 * x^{(m)}))^2 \end{aligned}$$

# The Best Fit



The green line has the lowest SSE and the best fit!

# Mean Squared Error



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

# The Best Fit

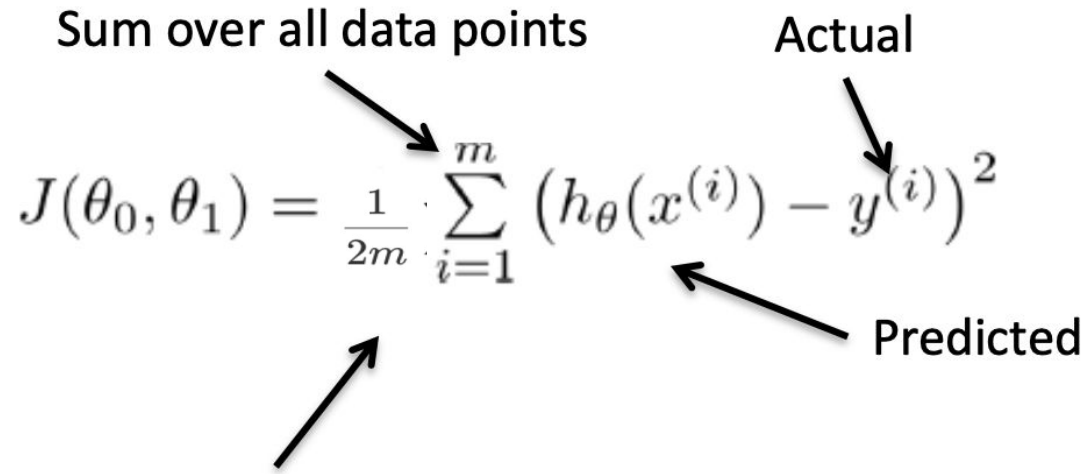
The best hypothesis  $h_{\theta}(x)$  is the one that minimizes the cost function

Sum over all data points

Actual

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

Predicted



$1/m$  - means we determine the average

$1/2m$  - simplifies the maths

# The Best Fit

The learning algorithm should find  $\theta^* = (\theta_0^*, \theta_1^*)$  that minimizes this cost

Several algorithms:

- Gradient descent
- Ordinary least square (OLS): Used in the linear regression in python (sklearn)

# Gradient Descent

# Notation – Simple Linear Regression

<b>X</b>	<b>Size (m<sup>2</sup>)</b>	<b>Y</b>	<b>Price (1000\$)</b>
	20		300
	37		540
	88		986
	...		...

**m = 47**

m: Number of examples

**x(2) = [37]**

x(i): Input of the i-th training example



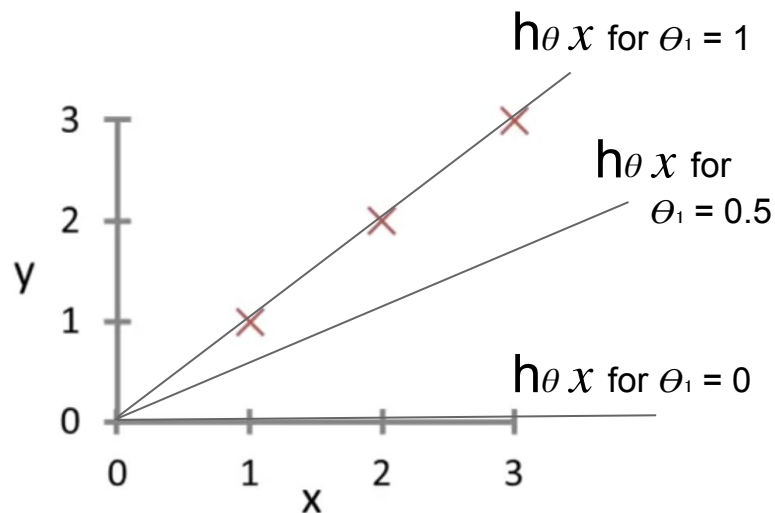
# Parameter influence

Hypothesis :  $\theta_0 = 0$

$$h_{\theta}(x) = \theta_1 x$$

# Parameter influence

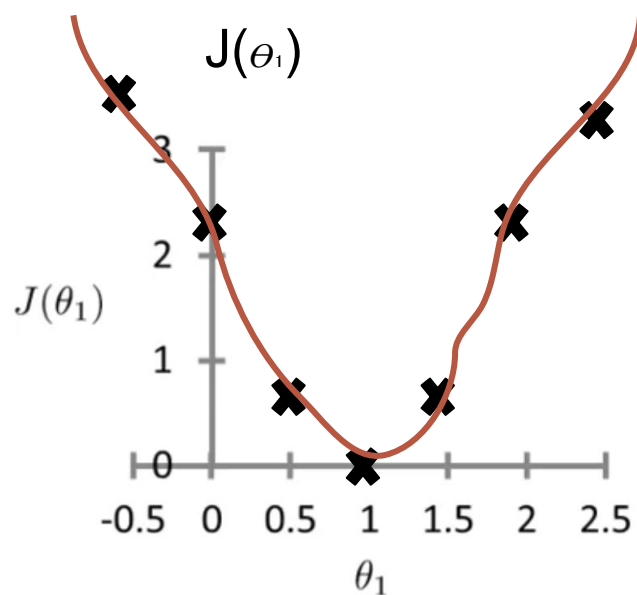
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$



$$J(1) = 1/2m (0^2 + 0^2 + 0^2) = 0$$

$$J(0.5) = 1/2m ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2) \sim 0.6$$

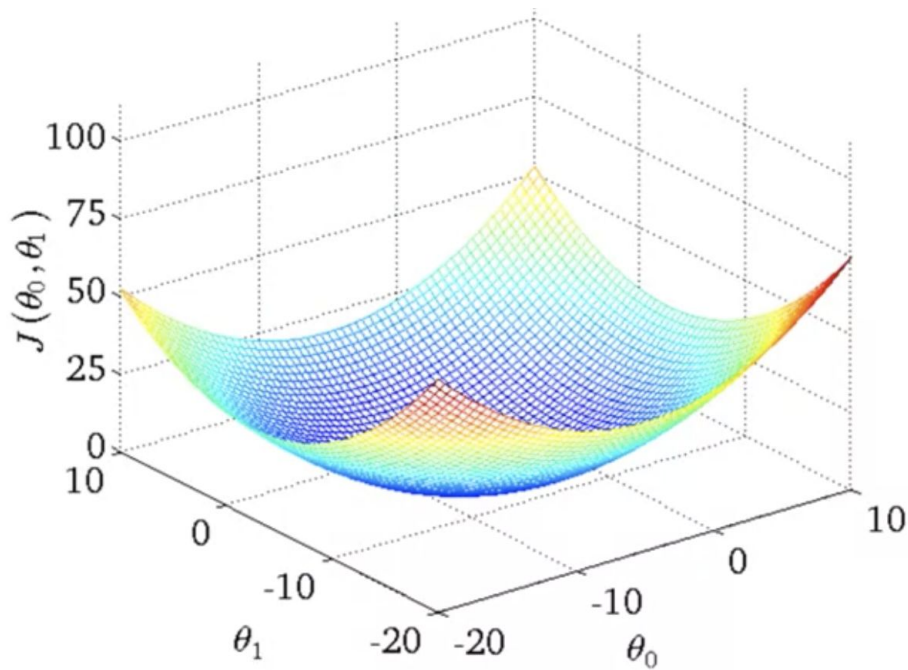
$$J(0) = 14/6 \sim 2.3$$



We want to minimize  $J(\theta)$ .

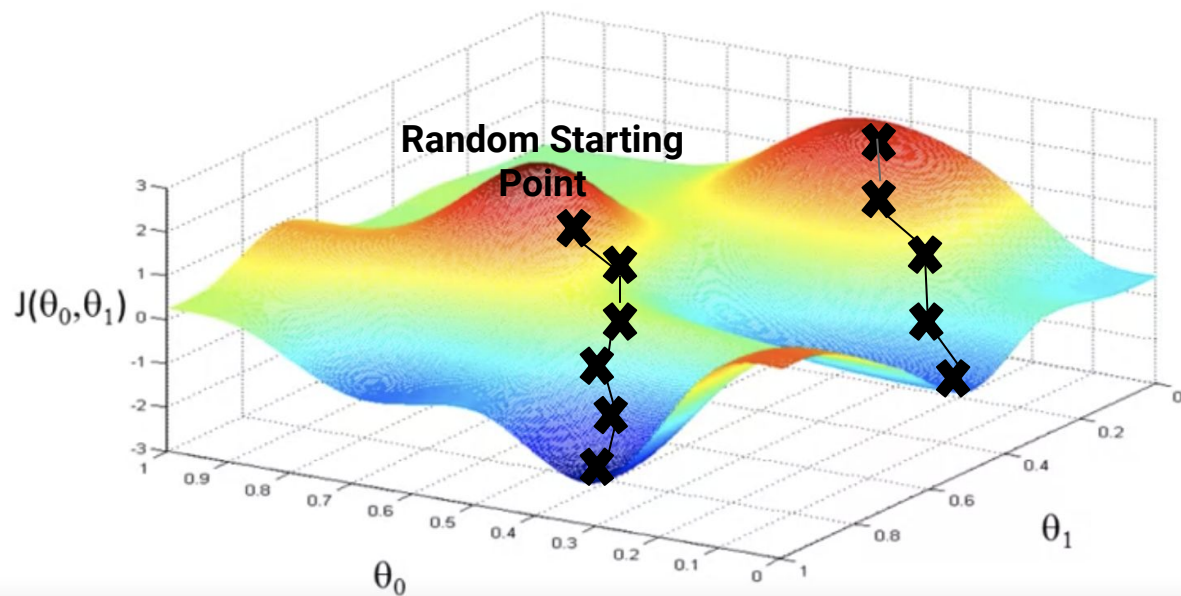
# Parameter influence

$\theta_0 \neq 0$



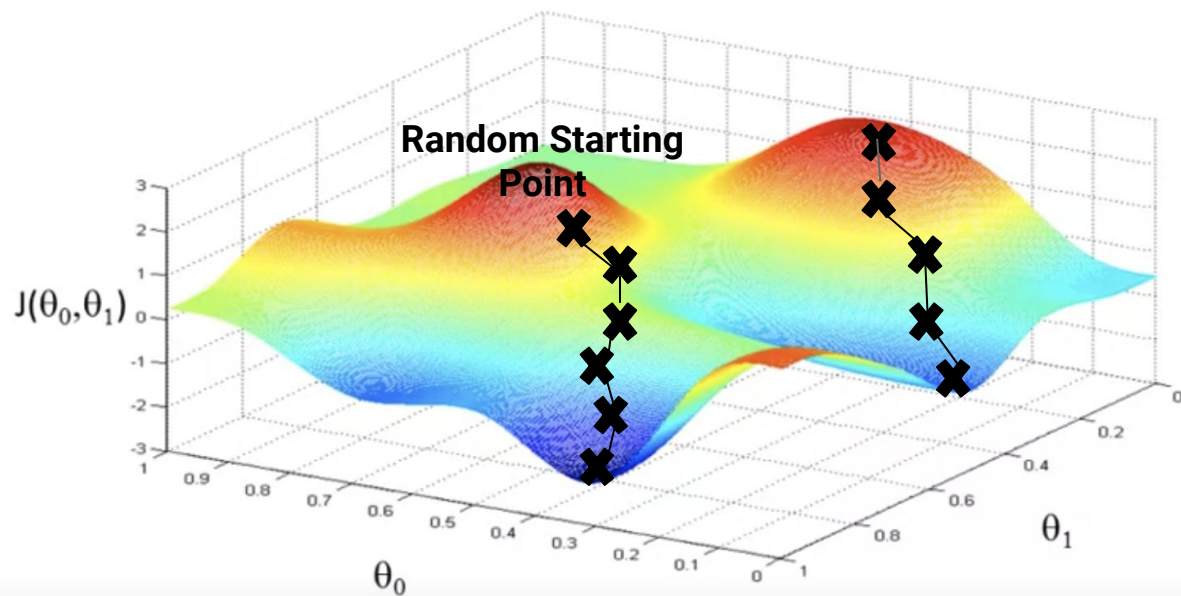
# Gradient Descent

- Start with random  $\theta_0$  and  $\theta_1$
- Change  $\theta_0$  and  $\theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we find a minimum

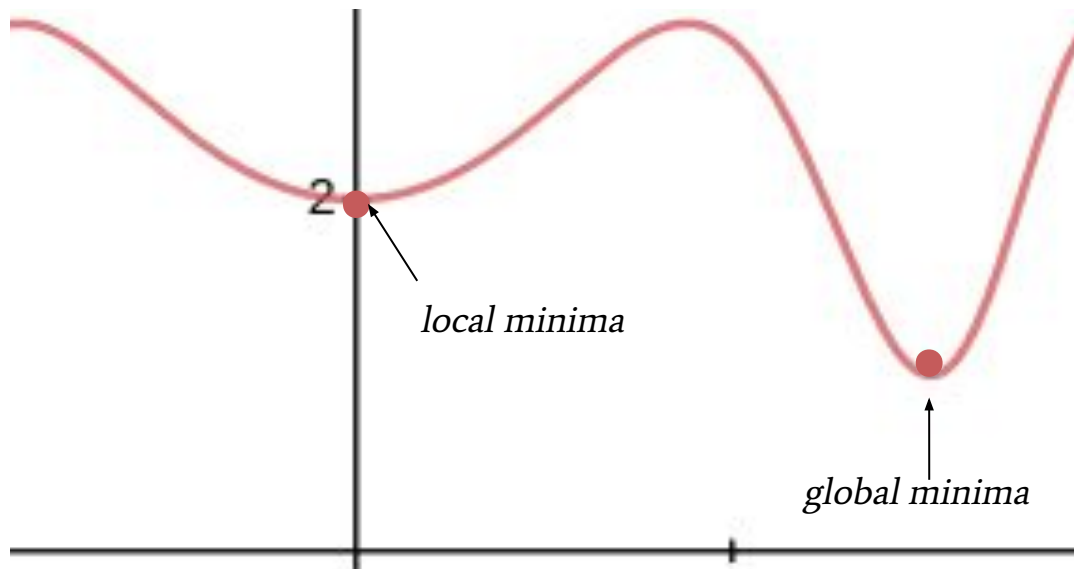


# Gradient Descent

For different initializations, we might have different results



# Local-Global minima



# Gradient Descent

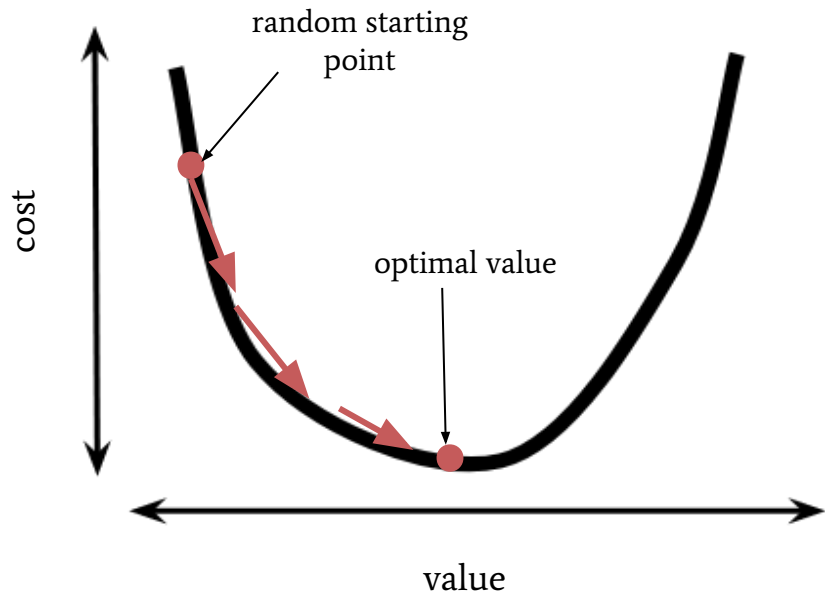
Repeat until convergence :

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

Diagram illustrating the gradient descent update formula with annotations:

- $\theta_j :=$  assignment
- $\theta_j$  current value
- $-$  learning rate
- $\alpha$  learning rate
- $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$  derivative

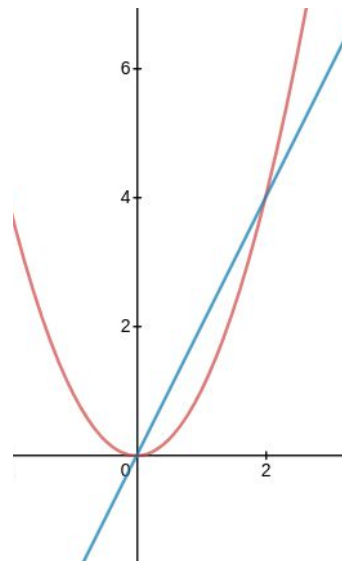
# Gradient Descent



gradient / rate of change / slope / derivative.

exemple :  $f'(x^2) = 2x$

gradient for  $x = 2$  is 4

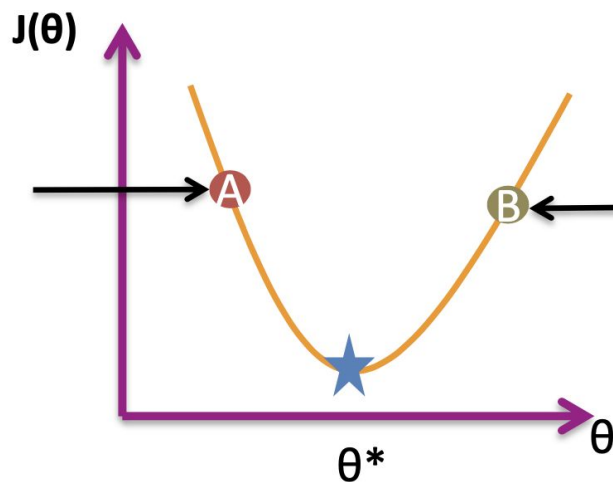




# Gradient Descent

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

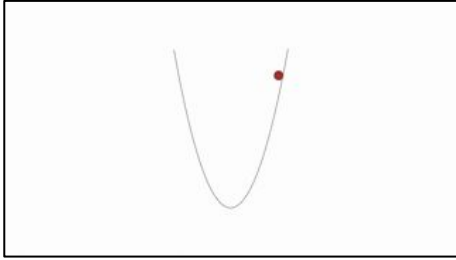
In this case, the derivative  
(gradient)  $\partial J(\theta) / \partial \theta < 0$  .  
 $\theta^A - \alpha * \partial J(\theta) / \partial \theta > \theta^A$   
 $\theta^A$  is moving to the right  
 $\theta$  is increasing



In this case, the derivative  
(gradient)  $\partial J(\theta) / \partial \theta > 0$  .  
 $\theta^B - \alpha * \partial J(\theta) / \partial \theta < \theta^B$   
 $\theta^B$  is moving to the left  
 $\theta$  is decreasing

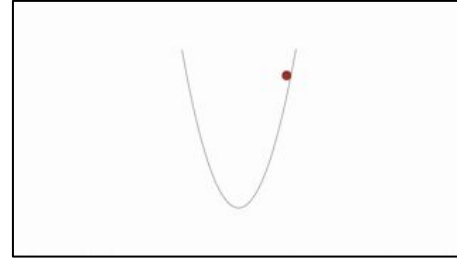
# Learning Rate

GOOD LEARNING RATE



- Not too low
- Allows the network to find the right parameters

BAD LEARNING RATE



- Too high
- Learns fast but quickly overshoots and ends up increasing the error

# Gradient Descent Recap

## Gradient Descent

(n=1):

Repeat {

$$\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_0} J(\theta)}$$

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

(simultaneously update  $\theta_0, \theta_1$ )

}

# Multivariate Linear Regression

# Notation – Multivariate Linear Regression

$n=4$

Size $X_1$	Number of Bedrooms $X_2$	Floor $X_3$	Age (years) $X_4$	Price (1000\$) $Y$
20	0	3	30	300
37	1	0	45	540
88	3	4	60	986
...				...

$m = 47$

$m$ : Number of examples

$n$ : Number of features

$x(i)$ : Input (features) of the  $i$ -th training example

$x_j(i)$ : Value of feature  $j$  for the  $i$ -th training example

$$x(2) = \begin{bmatrix} 37 \\ 1 \\ 0 \\ 45 \end{bmatrix}$$

$$x_3(2) = 0$$

# Multivariate Linear Regression

In simple linear regression,  $h_{\theta}(x) = \theta_0 + \theta_1 x$

In multiple linear regression,  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$

For convenience of notation, define  $x_0 = 1$  ( $x_0^{(i)} = 1$ )

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

# Multivariate Linear Regression

Rewrite in matrix notation

$$\begin{matrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \dots \\ \theta_n \end{matrix} \quad \theta, x \in R^{n+1} \quad \begin{matrix} x_0 \\ x_1 \\ x_2 \\ \dots \\ x_n \end{matrix}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$h_{\theta}(x) = \sum_{i=1}^n \theta_i x_i = \theta^T x$$

For correct multiplication, we need  $\theta^T$

$$\theta^T = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n]$$

$$X = \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{pmatrix} \Rightarrow 1$$

# Multivariate Linear Regression

Cost function:

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



# Gradient Descent

## Simple Linear Regression

$$\begin{aligned} &\text{Repeat } \{ \\ &\theta_0 := \theta_0 - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_0} J(\theta)} \\ &\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \\ &\quad \text{(simultaneously update } \theta_0, \theta_1 \text{)} \} \end{aligned}$$

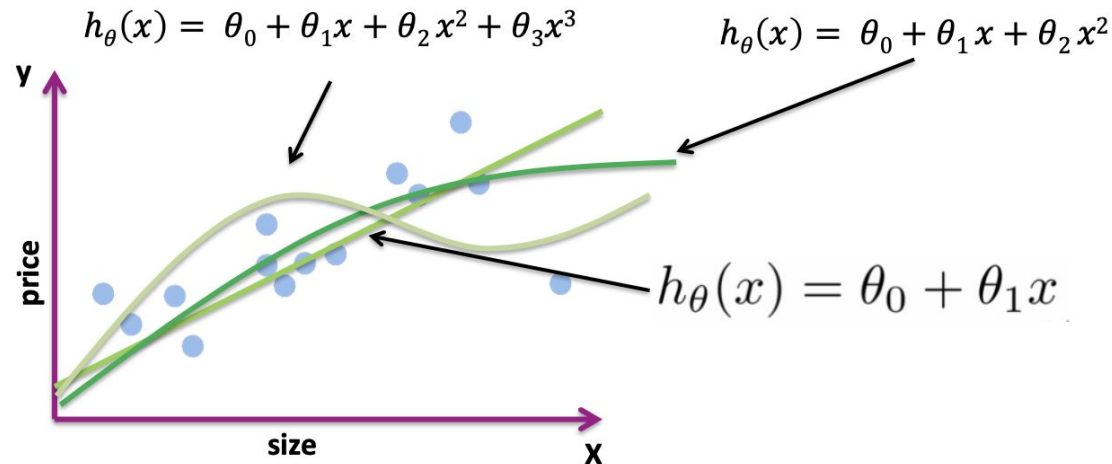
## Multivariate Linear Regression

$$\begin{aligned} &\text{Repeat } \{ \\ &\quad \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ &\quad \quad \text{(simultaneously update } \theta_j \text{ for } j = 0, \dots, n) \\ &\} \\ &\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ &\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ &\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} \\ &\dots \end{aligned}$$

# Polynomial Regression

# Polynomial Regression

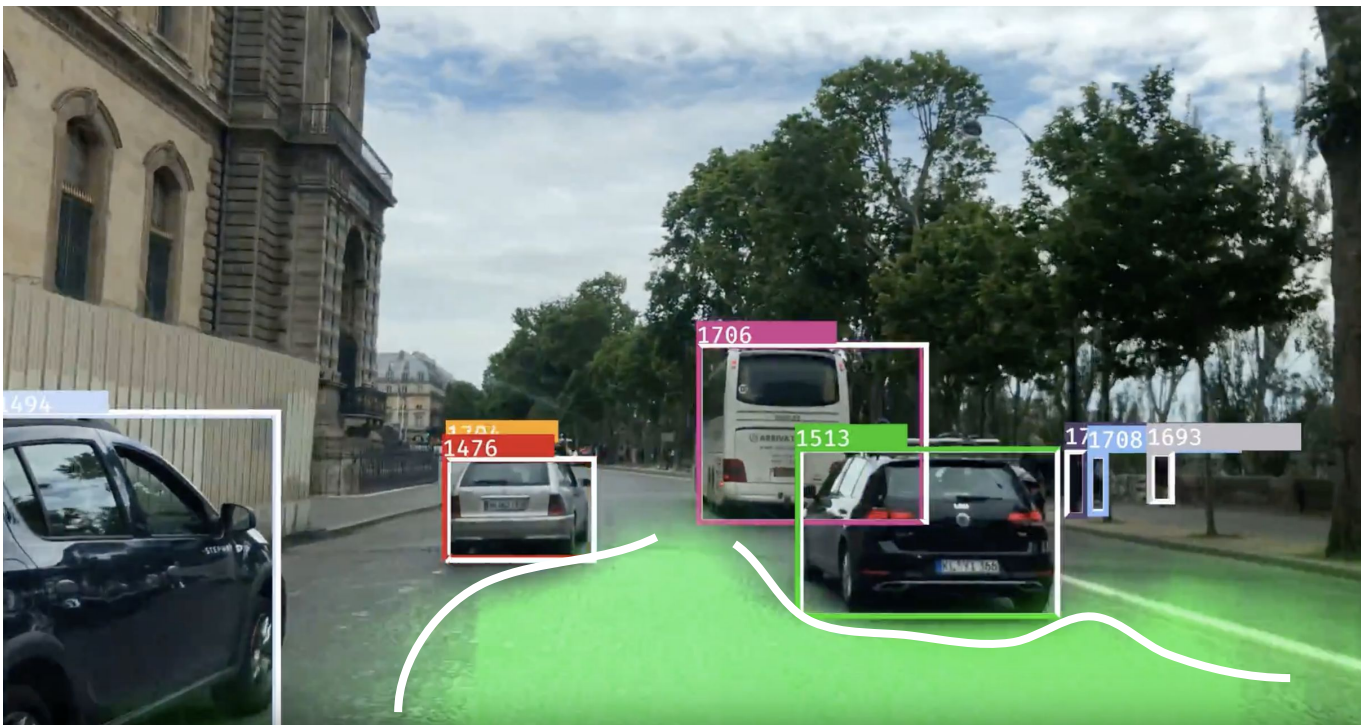
Polynomial regression is a particular case of multiple regression where the features are powers of one single feature  $x$



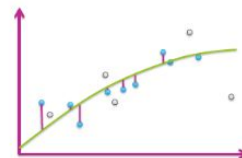
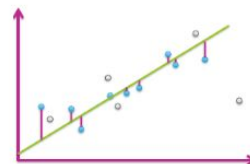
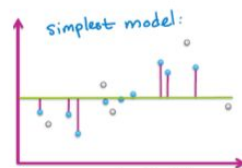
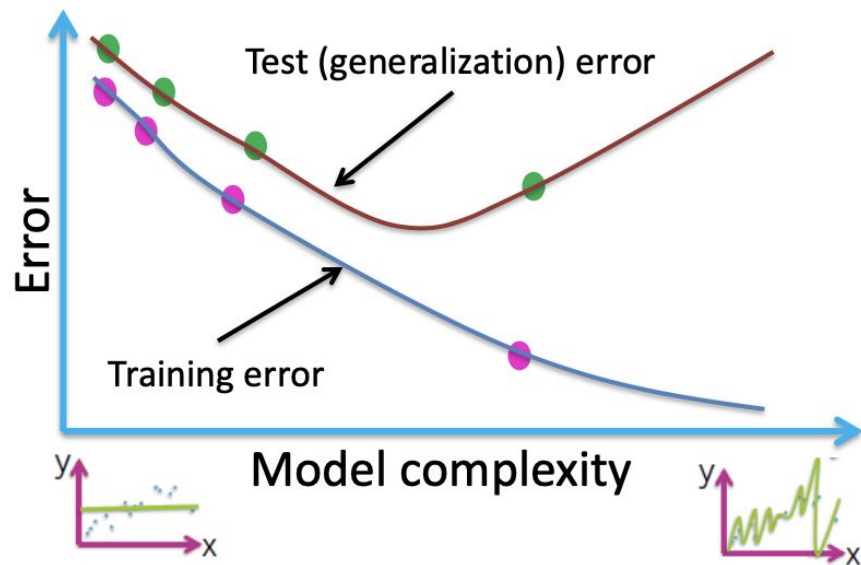
General model:

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots + \theta_p x^p + \varepsilon$$

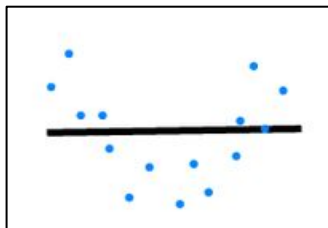
# Polynomial Regression in Practice



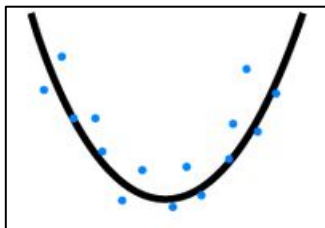
# Performance



# Performance



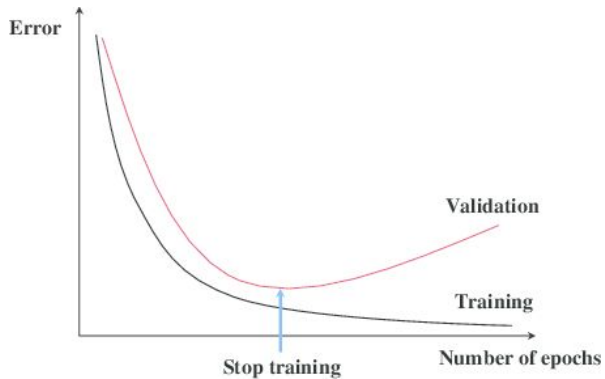
**UNDERFITTING**



**JUST RIGHT**



**OVERFITTING**



# Thank

# You

[jeremycohen.podia.com](https://jeremycohen.podia.com)

<https://www.linkedin.com/in/jeremycohen2626/>