

# Linear Models

Rohan Saha

# About me

MSc student at UofA - Computing Science

Research : Machine Learning and Neuroscience

Volunteer at GDG Edmonton working as a technical speaker

I love bodybuilding and cooking (sometimes fail).



# Prediction Problem

Given some input  $\mathbf{X}$ , predict an output  $\mathbf{Y}$ .

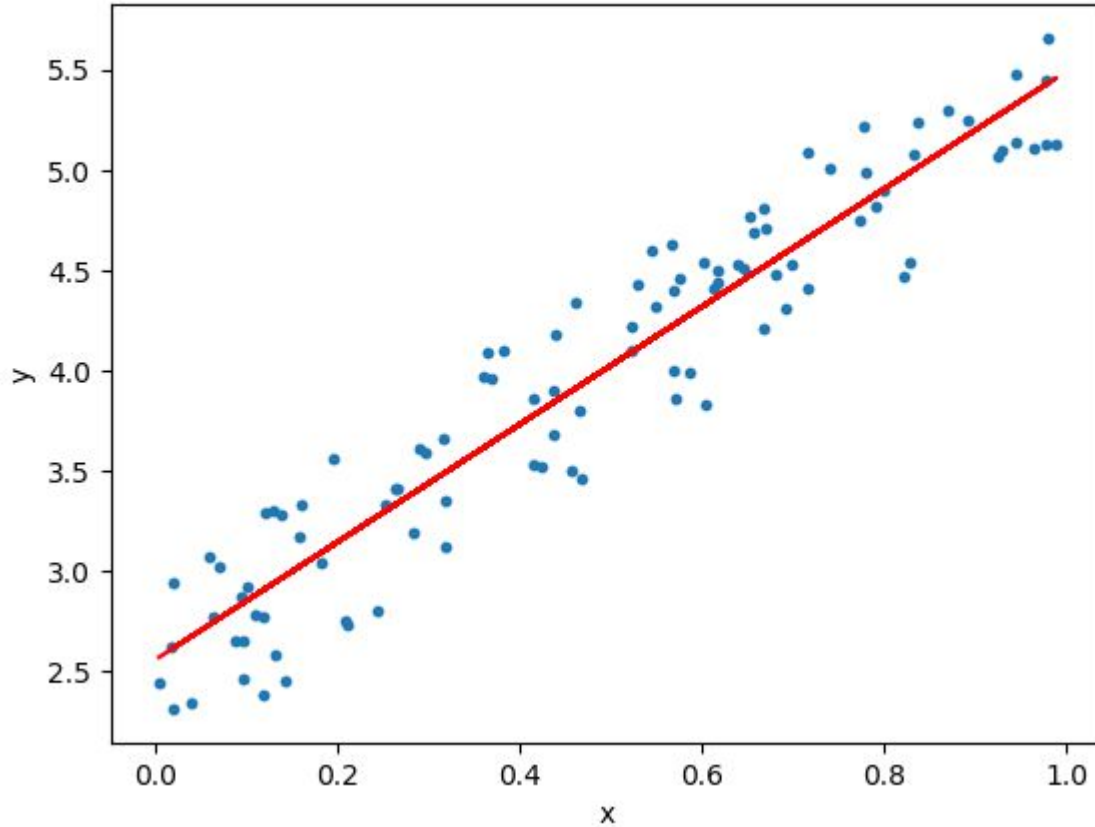
$$\mathbf{Y} = f(\mathbf{X})$$

$f$  can be linear or non-linear.

$\mathbf{X}$  is (usually) a matrix.

$\mathbf{Y}$  is (usually) a value of set of values.

We will focus on Linear functions for now.



$$Y = \Theta . X$$

Here  $\Theta$  is the transformation applied on X.

Task: Find (best)  $\Theta$

Goal : find  $\Theta$ .

$\Theta$  can take values like.

Hypothesis:  $Y = f(X) = \Theta.X$

$$\Theta = 0.5$$

For simplicity: Let  $X$  be one valued variable.

$$\Theta = 1.4$$

Therefore:  $Y = f(X) = \Theta_1.X$  is our hypothesis.

$$\Theta = 2$$

$\Theta$ 's are called parameters.

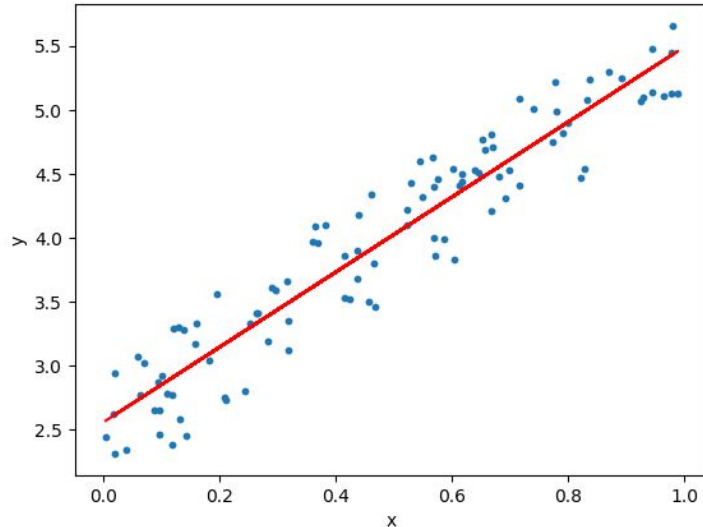
$\Theta$  = any real value.

Different  $\Theta$  -> Different  
hypothesis

# Cost Function

In the hypothesis, we also add a bias term for the intercept.

So hypothesis is now:  $Y = f(X) = \Theta_0 + \Theta_1 \cdot x$ , where  $\Theta_0$  is the bias.



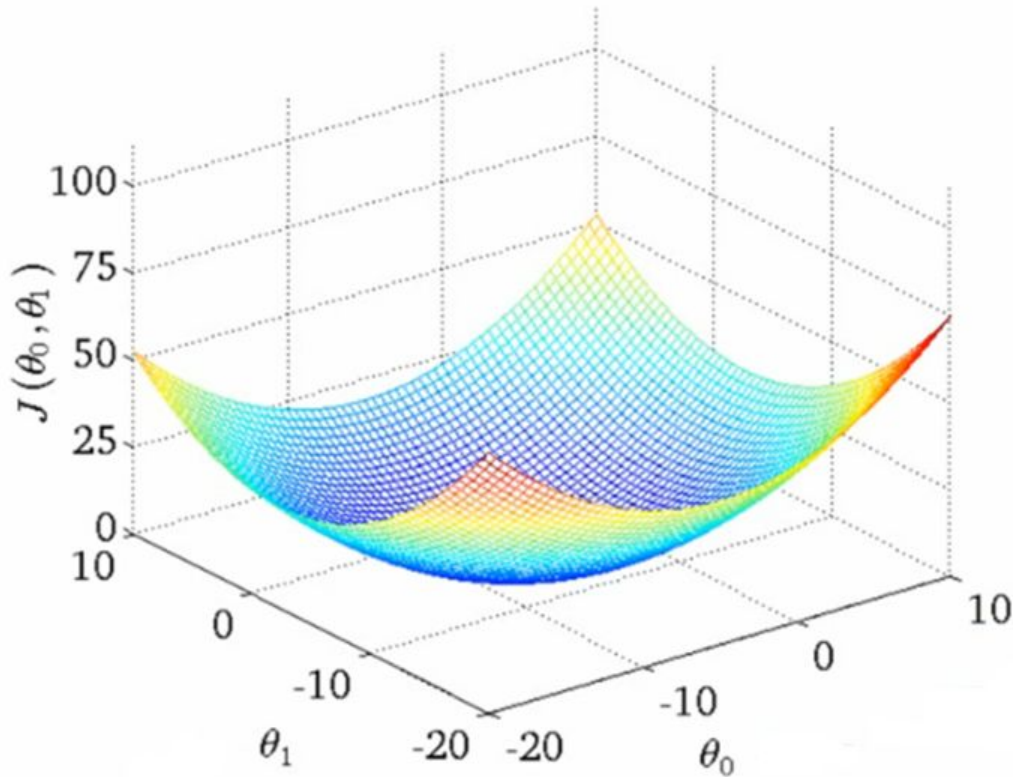
$$MSE = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

MSE = Mean Squared Error cost function.

$$MSE = J(\Theta_0, \Theta_1)$$

Goal: Minimize cost function.

# Visualizing Cost Function

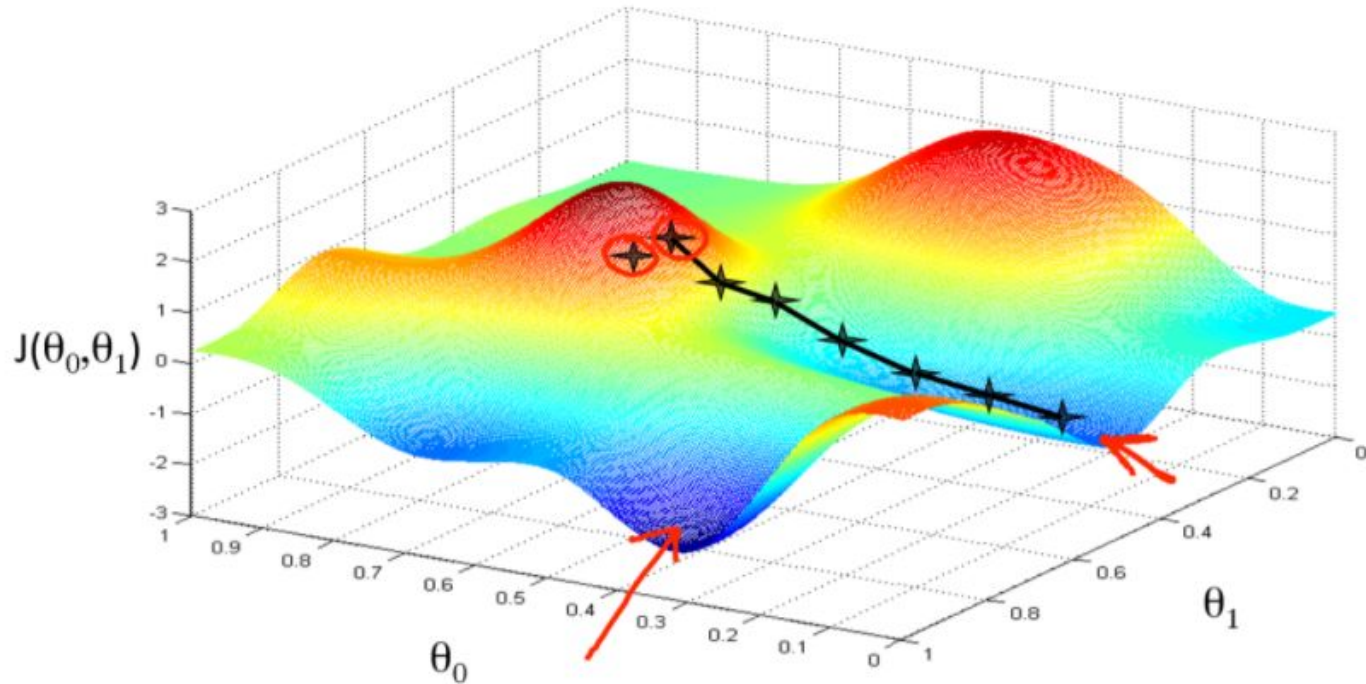


Shows how the cost function varies with different parameter values.

Linear Regression has a convex cost function.

We want to find values of parameters to obtain the minimal  $J$  value.

# Gradient Descent





Repeat until convergence {

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

$\alpha$  is the learning rate.

$J(\Theta)$  is the cost function.

## Types of Gradient Descent:

- Stochastic Gradient Descent
  - Update after each sample
- Batch Gradient Descent
  - Update after going through all the samples
- Mini-Batch Gradient Descent
  - Update after going through a subset of samples

# Linear Regression with Multiple Variables

- Input data has more than one variable (feature)

- $\mathbf{X} = (x_1, x_2, x_3, x_4, x_5, \dots, x_n)$

- More parameters

- $\Theta = (\Theta_1, \Theta_2, \Theta_3, \Theta_4, \dots, \Theta_n)$

- Helps to Vectorize using Linear Algebra!

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$\Rightarrow X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}, \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\Rightarrow h_{\theta}(x) = \theta^T x,$$

Q: What is the value of  $x_0$ ?

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:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$$

}

(simultaneously update for every  $j = 0, \dots, n$ )

# Closed-form solution

Usually used when we have many features (for the input data)

No iterations required.

No step size selection

Can be computationally intensive if #features is large.

$$\theta = (X^T X)^{-1} X^T y$$