

# Univariate Linear Regression

ARYAN KUMAR

# Machine learning

ARYAN KUMAR

- Given an experience  $E$ , performance measure  $P$ , task  $T$ .
- A machine is said to be learning if its  $P$  wrt a task  $T$  improves with experience,  $E$ .

## Types of Machine learning

- Supervised:

Data is labelled or we know the correct answers as it's provided.

Types:

Regression:

→ Response is quantitative

Classification:

↓  
Response is categorical

- Unsupervised:

Data is not labelled. We usually find

Structure  
Patterns  
Groupings (Clusters)  
Outliers (Anomalies)

Types:

Clustering  
Anomaly Detection  
Dimensionality Reduction

# Reinforcement learning ?

To be Updated

Here is a basic example :

Eg. Teaching a dog to play fetch  
using treats as a reward

# Univariate linear Regression

How to do Univariate linear Regression:

We could use python and do it one line using :

Sklearn, Statsmodel, pytorch, tensorflow, numpy.

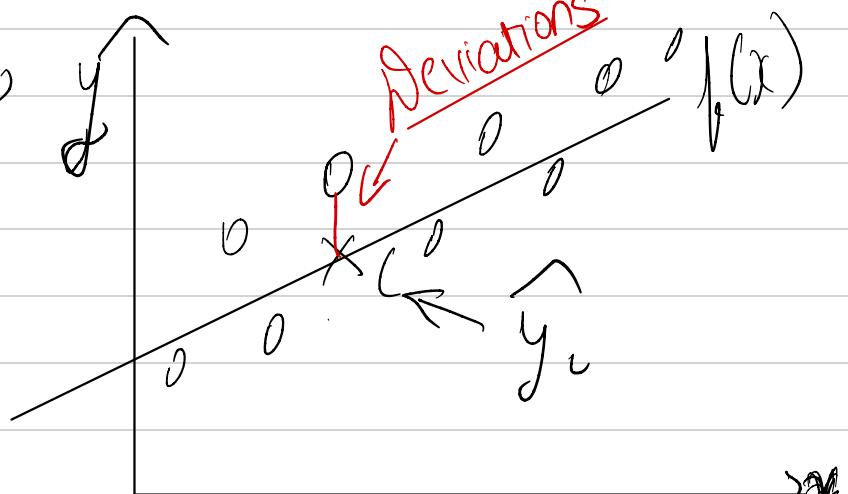
1) Pick a **form** for our model  $f$ .

2) Choose a **cost func.  $J$** .  $J$  will help us in getting the best model parameters.

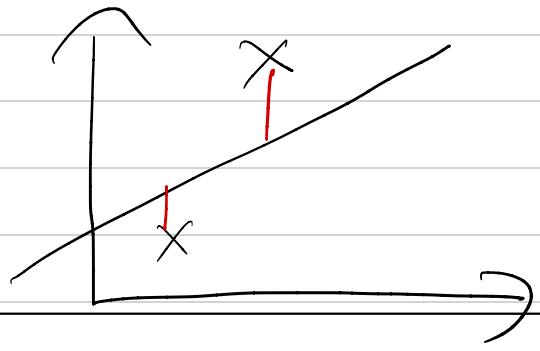
3) Train the model. Apply **Gradient Descent** to find optimal parameters.

Pick a general form for our model,  $f$ :

- $f_{w,b}(x_i) = w x_i + b + (\epsilon)$   $\epsilon$  - Irreducible error
- $w$  &  $b$  are model parameters.
- $w$  (weight or  $\beta_1$  in stats)  $y = \beta_0 + x\beta_1 + \epsilon$
- $w$  is slope, for  $\theta$ ,
- $b$  is  $y$  intercept



- We have observed value  $y_i$  & predicted  $\hat{y}_i$



- We want a model that minimizes these deviations

- Define a cost function: We choose squared error func.

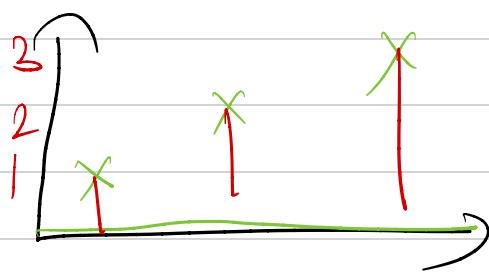
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 f(x)$$

$$= \frac{1}{2m} \sum_{i=1}^m (w x_i + b - y_i)^2$$

$$= \underset{w, b}{\operatorname{argmin}} J(w, b)$$

Eg  $f(x) = w x$

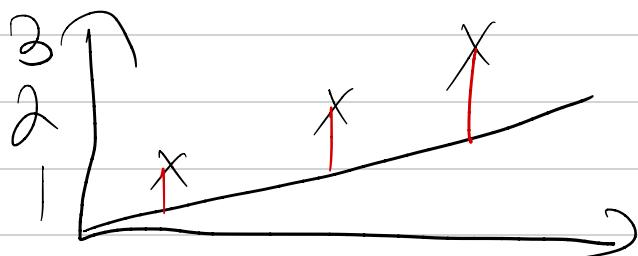
$$w=0 \Rightarrow \boxed{y=0} \text{ for all } x$$



$$\begin{aligned} J &= \frac{1}{2m} \sum_{i=1}^3 (w x_i - y_i)^2 \\ &= \frac{1}{2m} ((-1)^2 + (1)^2 + (3)^2) \end{aligned}$$

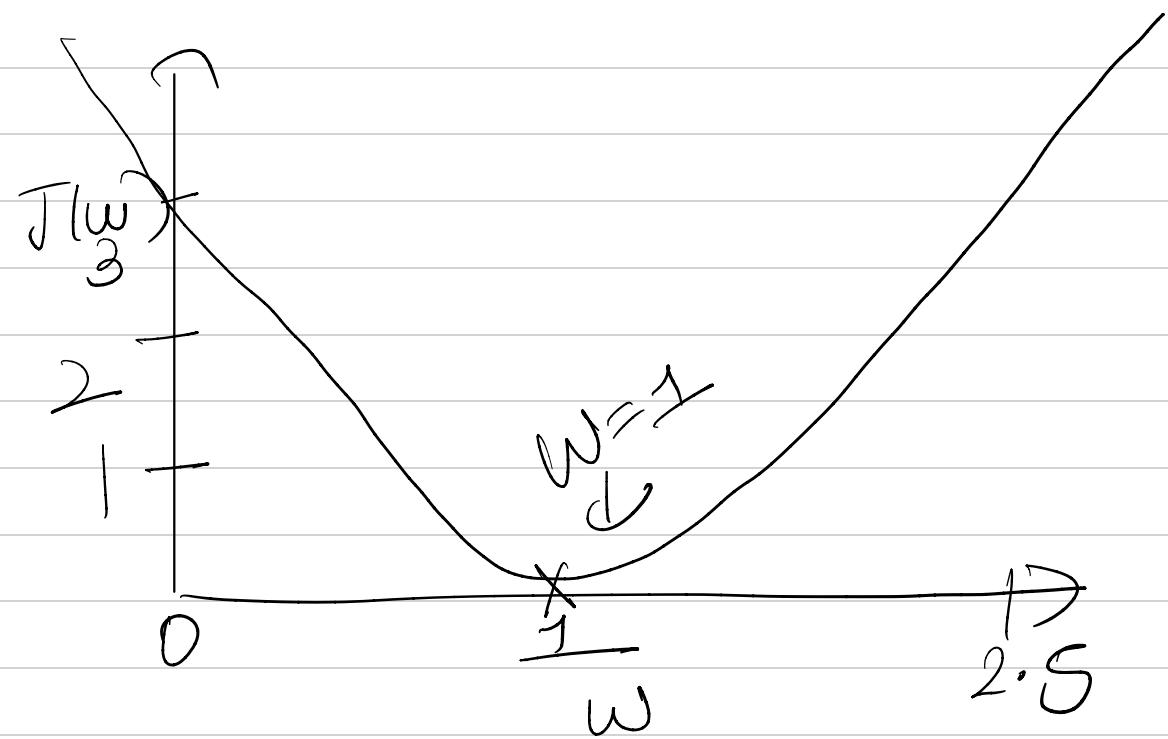
$$J(0) = 14/6$$

$$w = 1/2, \boxed{y = x/2}$$



$$\begin{aligned} J &= \frac{1}{2m} \sum_{i=1}^3 (w x_i - y_i)^2 \\ &= \frac{1}{2m} \left( \left(\frac{1}{2} - 1\right)^2 + \left(\frac{1}{2} - 1\right)^2 + \left(\frac{3}{2} - 3\right)^2 \right) \\ &= \frac{1}{2m} \left( \frac{1}{4} + 1 + \frac{9}{4} \right) = 7/12 \\ &= \boxed{7/12} \end{aligned}$$

$$w=1, J(1)=0$$



We see  $J(w)/w$  is a convex parabola

- We manually calculated minimum value of  $J(w)$ .
- Not feasible so we use gradient descent.

### Gradient Descent

1) Initialize random  $w$

2) Repeat till convergence.

$$w = w - \alpha \frac{\partial J}{\partial w}$$

$\alpha \in \mathbb{R} > 0$ , learning rate.

Finding  $\frac{\partial J(w)}{\partial w}$ : This is just slope.

$$\begin{aligned} J &= \frac{1}{2m} \sum_{i=1}^m (f(x) - y_i)^2 \\ &= \frac{1}{2m} \sum (wx_i - y_i)^2 \\ &= \frac{1}{2m} ((wx_1 - y_1)^2 + (wx_2 - y_2)^2 + \dots + (wx_m - y_m)^2) \end{aligned}$$

Take one at a time

$$\begin{aligned} \frac{\partial J}{\partial w} &= \frac{1}{2m} \left( \frac{\partial (wx - y)^2}{\partial w} \right) \\ \Rightarrow \Delta &= wx - y \Rightarrow \frac{\partial K}{\partial w} = \boxed{x} \end{aligned}$$

$$J \Rightarrow \frac{1}{2m} \sum K^2$$

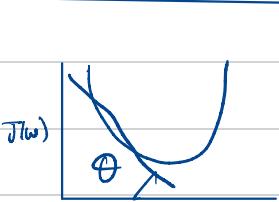
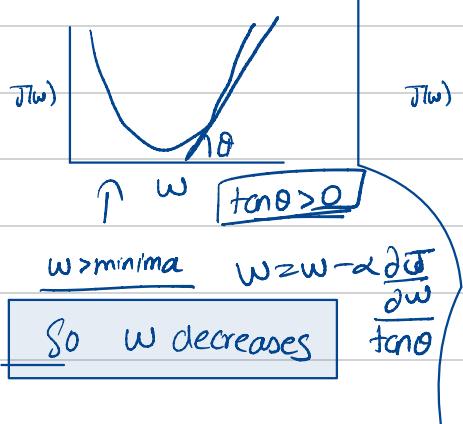
$$\underline{\frac{\partial J}{\partial K}} \Rightarrow K \quad (\Rightarrow \frac{\partial J}{\partial w} = \underline{\frac{\partial J}{\partial K}} \underline{\frac{\partial K}{\partial w}} \Rightarrow K(x))$$

$$\frac{\partial J}{\partial w} \Rightarrow (wx - y)x$$

$$\begin{aligned} \stackrel{0}{\circ} \quad \stackrel{0}{\circ} \quad \frac{\partial J}{\partial w} &= \frac{1}{m} \sum_{i=1}^m (wx_i - y_i) x_i \\ \Rightarrow \frac{\partial J}{\partial w} &= \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) x_i \quad \curvearrowleft a_i \end{aligned}$$

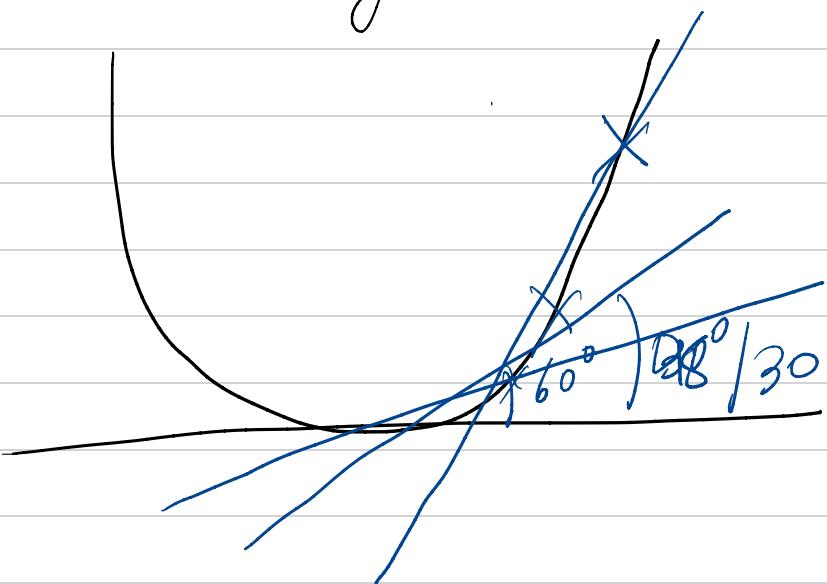
$$w = w - \alpha \frac{\partial J}{\partial w} \Rightarrow w = w - \alpha \left( \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) x_i \right)$$

Repeating till convergence



- Learning rate  $\alpha$
- $\bullet$  If  $\alpha$  is too small it converges slowly
- $\bullet$  If  $\alpha$  is too large it might not converge.

# Convergence



$$\tan 60^\circ = \frac{\sqrt{3}}{2}$$

$$w = w - \alpha \tan 60^\circ$$

$$w = w - \alpha \tan 30^\circ$$

$$w = w - \alpha \tan 30^\circ$$

\* We see as  $\theta \downarrow$ ,  $\tan \theta \downarrow$

so slope or  $\frac{\partial J}{\partial w} \downarrow$

: We do not need to change  $\alpha$  and the rate at which it converges reduces as we approach minima.

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

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f(x_i) - y_i)^2$$

$$= \frac{1}{2m} \sum_{i=1}^m (wx_i + b - y_i)^2$$

Gradient Descent:

- Initialize random  $w, b$
- Repeat till convergence :  $w = w - \alpha \frac{\partial J}{\partial w}$   
 $b = b - \alpha \frac{\partial J}{\partial b}$

$$\partial J / \partial w$$

$$J \Rightarrow \frac{1}{2m} ((wx_1 + b - y_1)^2 + (wx_2 + b - y_2)^2 + \dots + (wx_m + b - y_m)^2)$$

$$\partial J / \partial w \Rightarrow \frac{1}{m} (x_1 (wx_1 + b - y_1)) \quad | \quad \partial J / \partial w = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)$$

$$J_1 = \frac{1}{2m} (wx_1 + b - y_1)^2 \Rightarrow J_1 = \frac{1}{2m} \times u^2 \quad | \quad u = wx_1 + b - y_1$$

$$\frac{du}{db} = 1 \quad | \quad \frac{\partial J}{\partial u} = \frac{\partial u}{\partial b}$$

$$\frac{\partial J_1}{\partial b} = \frac{1}{m} (wx_1 + b - y_1)$$

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (wx_i + b - y_i)$$

Convergence

$$w = w - \alpha \frac{\partial J}{\partial w} \Rightarrow w - \alpha \left( \frac{1}{m} \sum (\hat{y}_i - y_i) x_i \right)$$

$$b = b - \alpha \frac{\partial J}{\partial b} \Rightarrow b - \alpha \left( \frac{1}{m} \sum (\hat{y}_i - y) \right)$$

They need to be calculated at the same iteration