# Linear Regression 3

Agenda
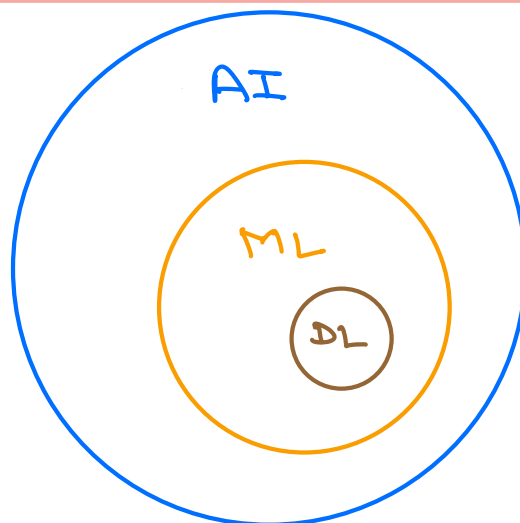
⇒ Recap

⇒ Adjusted r2_score

⇒ Need for scaling

⇒ Assumptions of Linear Regression

⇒ Sklearn and StatsModel implemention

# Types of ML



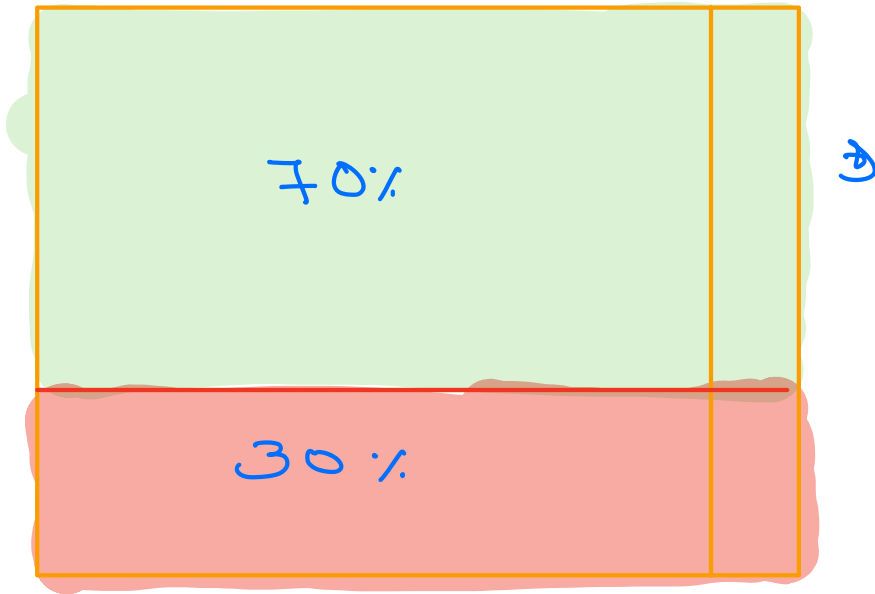Supervised
- Classification
- Regression

Semi

Unsupervised
- Clustering
- Association

## Regression

→ Deals with prediction of continous numeric values.

Ex:- Stock price prediction
Car Value prediction

| | $F_1$ | $F_2$ | $F_3$ | $F_4$ | $F_5 \ldots \ldots F_d$ | $y$ |
|---|---|---|---|---|---|---|
| $i = 1$ | | | | | | |
| 2 | | | | | | |
| 3 | | | | | | |
| $\ldots$ | | | | | $x_i$ | $y_i$ |
| 5 | | | | | | |

$n \Rightarrow$ no of rows sample

$d \Rightarrow$ no of features

$i^{th}$ sample $\Rightarrow x_i$

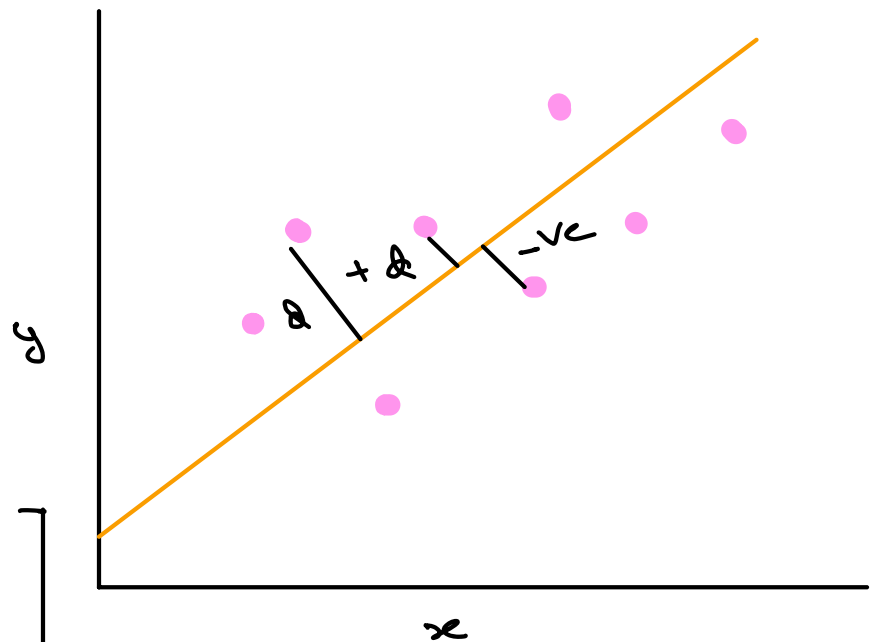$i^{th}$ label $\Rightarrow y_i$

70%

③

30%

# Linear Regression
## (ordinary Least Square)

Single Variable L.R.

$$\hat{y} = \omega_0 + \omega_1 (x)$$

Multi Variable L.R

$$\hat{y} = \omega_0 + \omega^T x$$

→ $\omega$ is a Vector → $\begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_d \end{bmatrix}$
d dims

→ $x$ is a Vector → $\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{bmatrix}$
d dims

$y$

$d$ | $+d$ | $-ve$

$x$

Goal: Find value of $\omega_0$ and $\omega_1$ for Best Fit line

| i | x | y | $\hat{y}$ | $y - \hat{y}$ |
|---|---|---|---|---|
| 0 | 1.5 | 3 | 2.5 | -0.5 |
| 1 | 3.5 | 4 | 3.5 | -0.5 |
| 2 | 6 | 2 | 5 | 3 |
| 3 | 5 | 5.5 | 5 | -0.5 |
| 4 | 7 | 5 | 5.5 | 0.5 |



$$RSS = \sum_{i=0}^{n} (y_i - \hat{y_i})^2$$

$$\hat{y_i} = w_1 x_i + w_0$$

**Loss/Cost Functions:**

$$MSE \Rightarrow \frac{1}{n} RSS = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y_i})^2$$

$$RMSE \Rightarrow \sqrt{MSE} \quad (\text{unit same as}) \; y$$

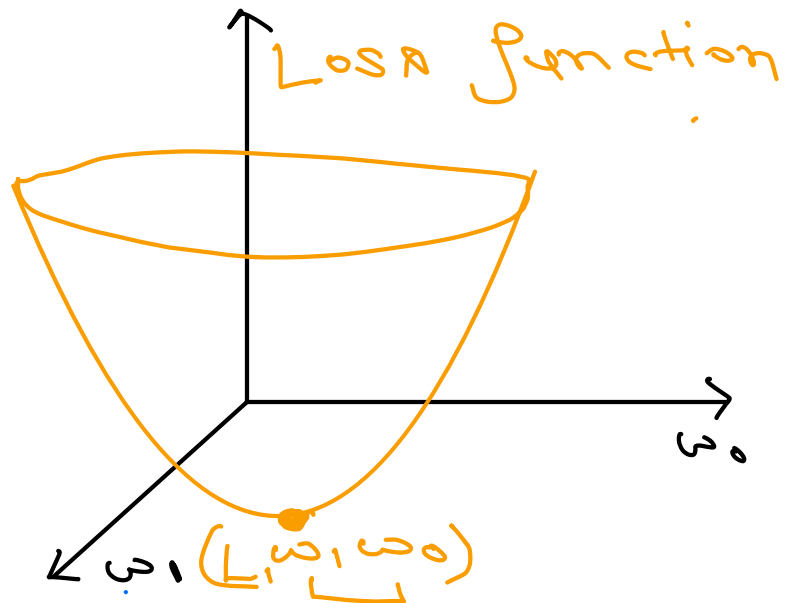$$MAE \Rightarrow \frac{1}{n} \sum_{i=0}^{n} |y_i - \hat{y_i}|$$

$$\mathcal{L} = \min_{w_0, w_1} \frac{1}{r} \sum_{i=0}^{r} (y_i - \hat{y}_i)^2$$



LoSA function

$(\mathcal{L}, w_1, w_0)$

$$\frac{\partial \mathcal{L}}{\partial w_1} \leftarrow \Delta w_1$$

$$\frac{\partial \mathcal{L}}{\partial w_0} \leftarrow \Delta w_0$$

$$w_1 = w_1 - \alpha \Delta w_1$$

$\eta \rightarrow$ learning rate
$\alpha$

How will the eq$^n$ change for MLR?

$$w_0 = w_0 - \alpha \Delta w_0$$
$$w_1 \Rrightarrow w_1 - \alpha w_1$$
$$w_2 \downarrow w_2 - \alpha w_2$$
$$w_3 \downarrow w_3 - \alpha w_3$$
$$w_d \downarrow w_d - \alpha w_d$$

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ -- \\ w_d \end{bmatrix}$$

Global minima = Best Value of $\leftarrow \vec{w}$ and $w_0$

## How do we find Global Minima?

### Gradient Descent

$$\omega_0 = \omega_0 - \alpha \frac{\partial L}{\partial \omega}$$

$$\omega_1 = \omega_1 - \alpha \frac{\partial L}{\partial \omega_1}$$

$$\vdots$$

$$\omega_d = \omega_d - \alpha \frac{\partial L}{\partial \omega_d}$$

$$L = (y_i - \hat{y}_i)^2$$

$$z^2$$

$$\frac{dL}{d\omega} \Rightarrow \boxed{\frac{dz^2}{dz} \times \frac{dz}{d\omega}}$$

$$2z \times \frac{dz}{d\omega}$$

Let's calculate partial derivatives manually:
for simplicity Let's assume only 2 features
i.e.  $\hat{y} = \omega_0 + \omega_1 x_1 + \omega_2 x_2$

$$L = \left[ y - (\omega_0 + \omega_1 x_1 + \omega_2 x_2) \right]^2$$

$$\downarrow$$

$$\frac{\partial L}{\partial \omega_0} \Rightarrow \frac{\partial z^2}{\partial \omega_0} \Rightarrow \frac{\partial z^2}{\partial z} \times \frac{\partial z}{\partial \omega_0}$$

$$\omega_0 + \omega_1 x_1 + \omega_2 x_2$$

$$z \Rightarrow y - \hat{y}$$

$$\frac{\partial (y - \hat{y})^2}{\partial (y - \hat{y})} \times \frac{\partial (y - \hat{y})}{\partial \omega_0}$$

$$\Rightarrow \quad 2 \times (y - \hat{y}) \times \left( \frac{\partial y}{\partial w_0} - \frac{\partial w_0}{\partial w_0} - \frac{\partial w_1 x_1}{\partial w_0} - \frac{\partial w_2 x_2}{\partial w_0} \right)$$

$$\Rightarrow \quad -2 \times (y - \hat{y})$$

$$\frac{\partial (y - w_0 - w_1 x_1 - w_2 x_2)}{\partial w_0}$$

$$0 \qquad -1 \qquad 0 \qquad -x_1$$

**Similarly**

$$\frac{\partial L}{\partial w_1} = -2(y - \hat{y}) \times x_1$$

**and**

$$\frac{\partial L}{\partial w_2} = -2(y - \hat{y}) \times x_2$$

$$\frac{\partial L}{\partial w_d} = -2(y - \hat{y}) \times x_d$$

$$\boxed{\begin{aligned}
\frac{\partial L}{\partial w_0} &= \frac{1}{r} \sum_{i=1}^{r} -2(y_i - \hat{y}_i) \\[2em]
\frac{\partial L}{\partial w_d} &= \frac{1}{r} \sum_{i=0}^{r} -2(y_i - \hat{y}_i) \times x_{id}
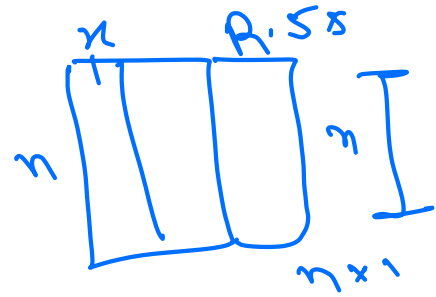\end{aligned}}$$

$y$

## Tricks to Speedup Training

### 1) Vectorization

in code we can Vectorize $\frac{\partial L}{\partial \omega_0}$ using $\cdot$ product

$$\frac{\partial}{\partial} \sum_{i=1}^{n} -2 (y_i - \hat{y}_i) \times x_{id} \qquad \text{if } d = 1$$

$$-2(y - \hat{y}) \cdot X$$

$$\downarrow \qquad\qquad \downarrow$$

$$n \times 1 \qquad\qquad n \times d$$
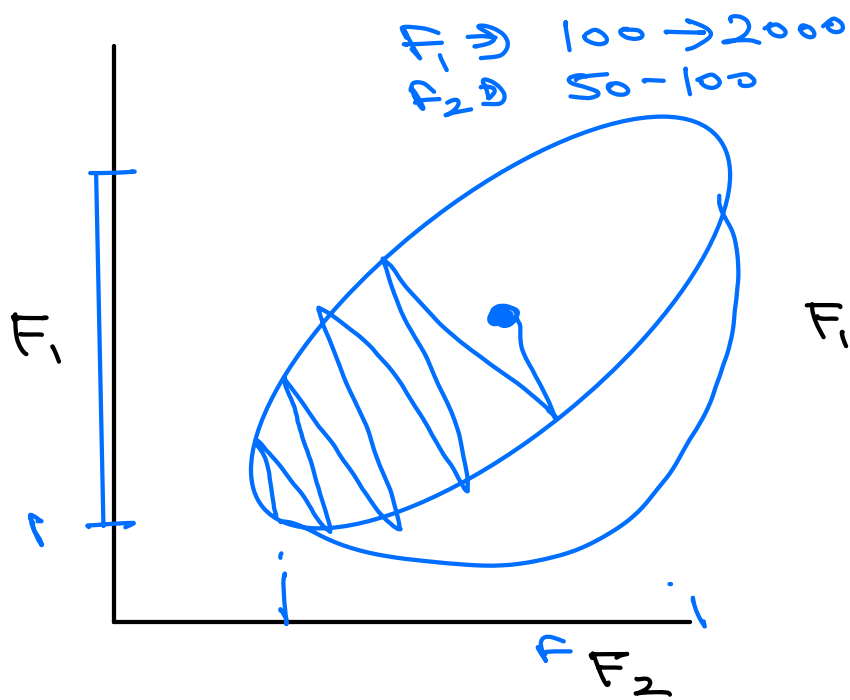


$n \times 1$

d derivatives over full dataset in single ops

$$x^T \cdot (y_i - \hat{y}) \Rightarrow \text{Replacement For Loop}$$

$d \times 1$

$$d \times n \cdot n \times 1 \Rightarrow d \times 1$$

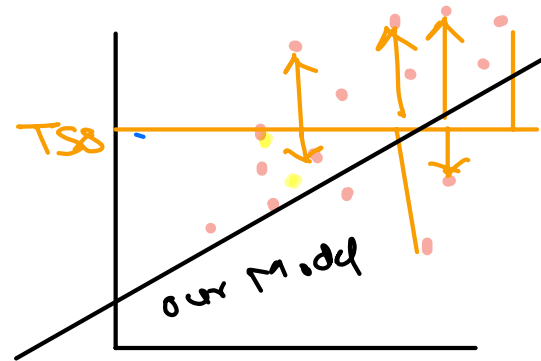$$\boxed{\frac{\partial L}{\partial \omega} \Rightarrow \frac{-2}{n}\left(x^T \cdot (y - \hat{y})\right)}$$

② Scaling    Why?

$F_1 \Rightarrow 100 \rightarrow 2000$
$F_2 \Rightarrow 50 - 100$

Standars
0
Std = 1

$F_1 \Rightarrow 0 - 1$
$F_2 \Rightarrow 0 - 1$

$F_1$

$F_2$

0.5

0.25

$F_1$

$F_2$

→ Scaling Ensures No feature dominates

**=Q: What does RSS or MSE of 10000 mean?**

$$R2\text{-score} = 1 - \frac{RSS}{TSS}$$

$$\sum_{i=0}^{n} \left( y_{mean} - y_i \right)^2$$

→ How good is LR model from mean Model



**Q Range of R-2 Score ?**

$(-\infty, 1]$

**Q Problem with R-2 Score**

Case 1: d-feature + 1 more relevant feature

r2-score ⇒ ↑ increase

Case 2: d - feature + 1 irrelevant feature

r2-score ⇒ Same
↓ x≠d
wd=0 ↑

⦿: How to mitigate this

$$adj\text{-}r2\text{-}score \Rightarrow \left[ 1 - \frac{(1-r2\_score) \times (n-1)}{n-d-1} \right]$$
↑

n ⇒ no of row
d ⇒ no of features

Case 1 ⇒ adj_r2_score ↑

Case 2 ⇒ adj_r2_score ↓

RSS

$$1 - \dfrac{RSS}{TSS}$$

$\neq > 1$

SS

r2_score $\Rightarrow$ -ve

$$r^2 = 1 - \frac{RSS}{TSS} \Rightarrow 1$$

$$-4000$$
$$-2000$$

RSS2

RSS = 0

TSS

⟶ Transformer

Generative AI

⟶ Diffusion Models

Text     audio     Video     Image