# Lecture 10 Data Preparation

جب ہم اپنے data کو کسی مشین لرننگ ماڈل کے ساتھ استعمال کرتے ہیں
اس سے پہلے ہم Data Preprocessing کرتے ہیں۔
data set میں جو data ہوتا ہے اسے صاف ستھرا کرتے ہیں۔

## Steps : Preprocessing :

**Imp** 1) **Remove the outliers :** ہم ایک error threshold define کرتے ہیں M2 Huber SSE

جو error ہوتا ہے اگر outlier زیادہ ہو تو اسے remove کر دیتے ہیں۔

max ▭—— min — **Method 2** box Plot سے بھی data سے outlier remove کرتے ہیں
min اور max value کو remove کرتے ہیں۔

## 2) Missing Data :

**Inf. loss** — 1) جہاں missing ہو اسے remove کر دیتے ہیں۔

2) جہاں پر جس colum میں missing ہو اس کا average لے کر جہاں missing ہے وہاں average ڈال دیتے ہیں۔

3) اس data کو دیکھتے ہیں جہاں value سب سے زیادہ ہے اسے missing میں ڈال دیتے ہیں۔

## 3) Non → Numric Value handling : (one hot Ecoding)
*to numric*

1) one hot Ecoding میں ہم unique value find کرتے ہیں، پھر ان کو matric میں تبدیل کرتے ہیں۔

Matric Inter undergr grad

$$
\begin{bmatrix}
1 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 \\
0 & 0 & 1 & 0 \\
0 & 0 & 0 & 1
\end{bmatrix}
$$

Matric / replace
[1 0 0 0]

Inter
[0 1 0 0]

non-numric value کو replace / vector سے بدل دیتے ہیں۔

سب سے پہلے first vector کو لے کر one سمجھ کر، اس

belong کر کہ وہ کس میٹرک میں مترک سے جاتے ہیں۔

⇒ One More Method:

اس میں ہم اپنی مرضی سے unique value کے لیے assign کریں گے

Red ⟶ 1

Green ⟶ 2

Yellow ⟶ 3

again      Red ⟶ 1

اس میں ماڈل کو آسانی سے sequence learn کرنے میں مدد ملتی ہے۔

4) Feature Scaling:       اس کو ہم بعد میں پڑھیں گے

5) 5 number Summary: Box Plot

Steps:
1) Sort dataset    Min
2) Select Min & Max
3) IQR = $Q_2 - M_1$
4) Lower limit = $Q_1 - 1.5 \times IQR$
5) Upper limit = $Q_3 + 1.5 \times IQR$

⇒ Inter Qurtial Range: IQR.

اس کی مدد سے ہم آسانی سے کسی بھی data کو سمجھ سکتے ہیں

$M_1$ Median1        Median        $M_2$ Median2        Max

6) Nois Injection:

جب بھی ہمارا data سائز میں کم ہوتا ہے تو اس کو Increase کرنے کے لیے dataset

ہم Nois Injection method لگاتے ہیں۔ اس میں

ہم original data میں تھوڑا سا شور ڈال کر نیا data بنا لیتے ہیں

**=) Scaling:**

0 اور 1 کے data / column / attribute کو رکھتے ہیں

column / attribute کی اگر کوئی ایک range کی

instance کو select کریں high value

دوسرے کو divide کریں maximum

Training time

Exp. ← Training کی iteration کو کم کرتے ہیں

| Size (house | Scaling Size |
|---|---|
| 2000 | $2000/6000 =$ |
| 4000 | $4000/6000 =$ |
| max 6000 | $6000/6000 =$ |

$$\text{Formula : } Scaling = \frac{Size\ (x)}{max\ (size)}$$

**Feature Scaling Method:**

**1) Mean Normalization:**

1) اس میں سب سے پہلے column کا mean

calculate                         vise

کرتے ہیں، value سے mean منفی کرتے ہیں

3) پھر اس کو column کے maximum سے divide کرتے ہیں۔

Size                Size Column

$$\text{Normalization} = \frac{value - mean}{max}$$

**2) Min - Max Scaling:**

$$x_{new} = \frac{x - min}{x_{max} - x_{min}}$$

Scaled
feature

**3) Standard deviation Scaling:**

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{n}}$$

## Box Plot & Five Number Summary

Dataset = $[3, 10, 14, 19, 22, 29, 32, 36, 49, 70]$

Step1 : Sort the dataset

Dataset Sorted = $[3, 10, 14, 19, 22, 29, 32, 36, 49, 70]$

outlier

Step2 : Five number Summary

1 - Minimum Value

2) First Quatial $Q_1$ (25%)

3) Median

4) Thrid Quatial $Q_3$ (75%)

5) Maximam.

Step3 :

1) Upper limit = $Q_3 + 1.5 \, (IQR)$

$$= \overset{36}{36} + 1.5(\overset{22}{22}) = \cancel{108} \boxed{69}$$

2) Lower limit = $Q_1 - 1.5 (IQRS)$

$$= \overset{14}{14} - 1.5(22) = \cancel{} \boxed{-19}$$

Step 4 : Calculte $Q_1$, $Q_3$, IQR

1) $Q_1 = \dfrac{25}{100}$ (Total Instance + 1)

$$= \dfrac{25}{100}(10+1) = \dfrac{275}{100} \boxed{= 2.75}$$

$\boxed{Q_1 = 2.75}$ → replace $\boxed{14}$

2) $Q_3 = \dfrac{75}{100}$ (total instance + 1)

$$Q_3 = \dfrac{75}{100}(10+1) = \dfrac{825}{100} = 8.25$$

$\boxed{Q_3 = 8.25}$ → $\overset{8\text{here}}{\text{replace}}$ value $\boxed{Q_3 = 36}$

3) IQR = $Q_3 - Q_1$

IQR = $\cancel{} \, 36 - 14$

$\boxed{IQR = \cancel{} \, 22}$

→ Draw box Plot :

Dataset = [3, 10, 14, 19, 22, 29, 32, 36, 49, 70] ↙ outlier

Five number Summary :

1) Minimum = 3
2) First Quential $Q_1$ = 14
3) Median = 22
4) Third Quential $Q_3$ = 36
5) Maximum = 49

Draw box Plot