

Mean \rightarrow centre of gravity / mass

- The most common type of Average out there.
- Most influenced by the outliers.
- It changes when symmetry of the distribution/graph is affected.

$$\bar{x} = \frac{\sum_{i=1}^N x}{N}$$

$$\bar{x} = \frac{\sum f}{\sum} \rightarrow \text{frequency}$$

Mean \Rightarrow $\frac{\text{Sum of obs.}}{\text{No. of obs.}}$

$$[10, 10, 10, 10] \Rightarrow 10 \Rightarrow \frac{4 \times 10}{4} = 10$$

$$[10, 10, 10, 100] \Rightarrow 32.5$$

⇒ Mean is shifting towards higher value.

⇒ Mean is in love with outliers, so always shifts towards outliers.

$$\begin{array}{cccccc} & 0 & 1 & 2 & 3 & 4 & 5 \\ \hline [& 1 & , & 2 & , & 3 & , & 4 & , & 5 & , & 6 &] \\ \hline \end{array}$$

X

$$\text{Mean} = \frac{1 + 2 + 3 + 4 + 5 + 6}{6}$$

Mean \Rightarrow

$$\frac{\sum_{i=0}^5 X}{N}$$

6
summation

| x | f |
|---|---|
| 1 | 2 |
| 2 | 3 |
| 3 | 4 |

→ [1, 1, 2, 2, 2, 3, 3, 3, 3]

$$\text{Mean} \Rightarrow \frac{1+1+2+2+2+3+3+3+3}{9}$$

↓
Distribution table

$$\Rightarrow \frac{1 \times 2 + 2 \times 3 + 3 \times 4}{9}$$

df[col].count_values()
df[col].mean()
 ↓
 pandas

$$\Rightarrow \frac{\sum x f}{\sum f} \Rightarrow \text{Mean}$$

[1 2 3 4 5]

→ 3

[3 2 4 1 5]

Median

Sort ([1 2 3 4 5] ⇒ Median ⇒ 3

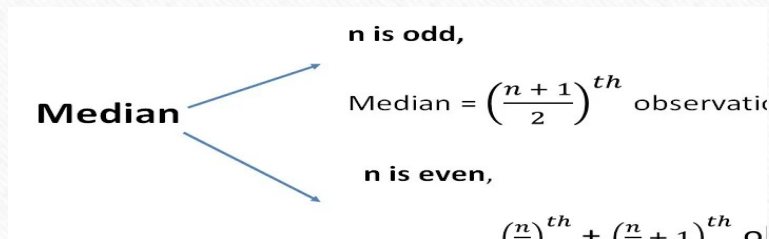


Median
(50-1.)

- It is actually the middle value, dividing the data into 2 equal halves.

- Steps:

- Sort the data in ascending order



~~Median = $l + \left(\frac{\frac{n}{2} - f}{f}\right) \times h$~~

- It is ~~not~~ influenced by outliers.

less

↘ special cases

| 1 st | 2 nd | 3 th | 4 th | 5 th |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| 11 | 12 | <u>13</u> | 14 | 15 |

Median \Rightarrow 13

getting
to right
position

$$\frac{5 + 1}{2} \Rightarrow 3^{\text{th}} \rightarrow 13$$

↑
Median

$\begin{matrix} 1^{st} & 2^{nd} & 3^{rd} & 3.5 & 4^{th} & 5^{th} & 6^{th} \\ \hline [& 11 & 12 & 13 & 14 & 15 & 16 &] \end{matrix}$

$len(list) = 6$

$$\left(\frac{3^{rd} + 4^{th}}{2} \right)$$

$$\left[\frac{\frac{6}{2} + \left(\frac{\frac{6}{2} + 1}{(3^{rd}) + 1} \right)}{2} \right]$$

positions

3.5

$$\Rightarrow \frac{13 + 14}{2} = 13.5$$

$[7, 9, 2, 17, 18, 28]$

Sort

| | | | | | | | |
|-------------|----|----|----|----|---------------|----|----|
| | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| $[2, 7, 9]$ | | | | | $17, 18, 28]$ | | |

$$\text{Median} = \frac{9 + 17}{2} = \textcircled{13}$$

①

[1, 2, 3, 4, 5]

↓ Median ⇒ (3)

②

[1 2 3 4 500]

↓

Median = 3

③

[1 2 3 4 500 600]

↓

Median = 3.5

④

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|-----|-----|-----|-----|
| 1 | 2 | 3 | 4 | 500 | 600 | 700 | 800 |

$$\frac{4^{\text{th}} + 5^{\text{th}}}{2} = \frac{4 + 500}{2} = 252$$

[B B B A C D E E]

Mode = B

Mode

[B B B A C D E E E]

Modes \Rightarrow B or E
Bimodal

- It works for both numerical and categorical data.
- Observation with highest frequency is the mode.
- A dataset can be unimodal or multimodal.
- For ungrouped data, we can count or directly observe in the table.
- For grouped data, we use

$$\text{Mode} = l_1 + \left(\frac{f_0 - f_{-1}}{2f_0 - f_{-1} - f_1} \right)$$

df.value_counts()
sns.count_plot()

Imputation of Missing values

- If outliers aren't present & feature is numerical, mean value is used for imputation.
- If outliers are present, median is used
- If col is categorical the mode is used!

Mean A = 0

Mean B = 0

Median A = 0

Median B = 0

A

-5

0

+5

B

-10

0

+10

Measures of Spread

Measures

Range

Variance

Standard
Deviation

Quartile
Range(IQR)

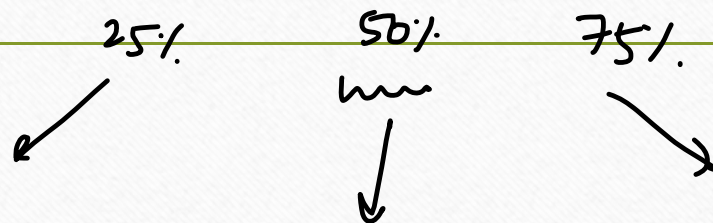
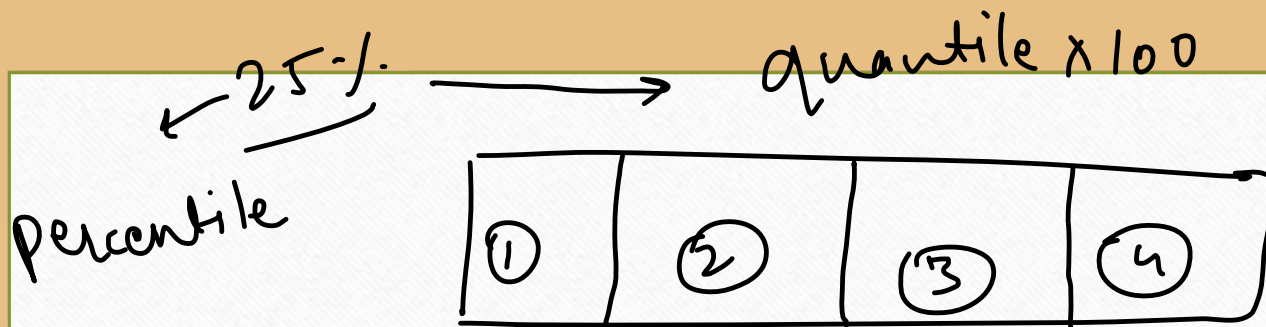
$$[1 \ 2 \ 3 \ 4 \ 5] \quad \text{Range} = 5 - 1 = 4$$

$$[1 \ 3 \ 5 \ 10 \ 200] \quad \text{Range} = 200 - 1 = 199$$

\times

- The range is a way of measuring how spread out a set of values are.
- The range only describes the width of the data, not how it's dispersed between the bounds.
- Range is very sensitive to outliers.

$$\text{Range} = X_{\max} - X_{\min}$$



Q1
↓
Quantile 1

Median

↓
Q2
↑
Quantile 2

Q3
↓
Quantile 3

quantile \neq quantile
↑
encapsulates

Quantile \Rightarrow $\frac{\text{Percentile}}{100}$
 $\Rightarrow 0.25$

Inter Quartile Range

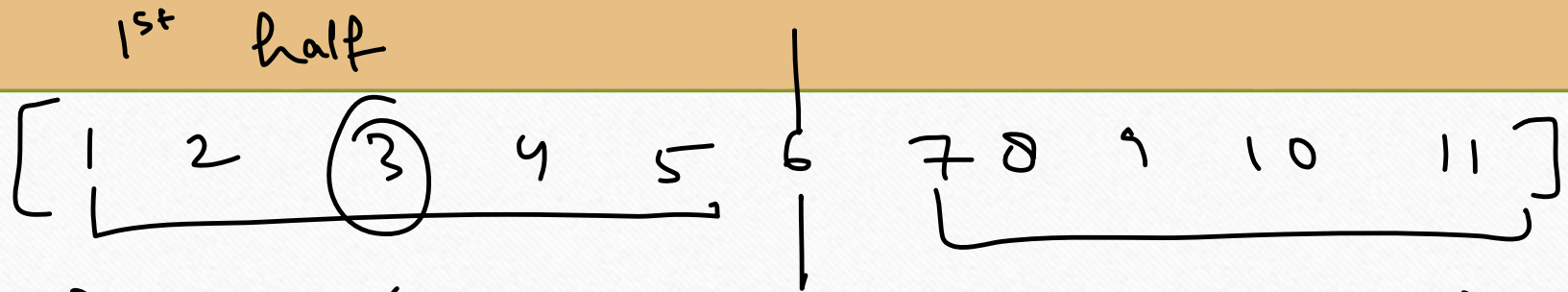
Quartiles (Q_1, Q_2, Q_3)
(between the quartiles)
max-min
 $IQR = Q_3 - Q_1$

- Quartiles: The numbers that separate data into 4 equal parts.



- IQR: difference between the 3rd quartile and 1st quartile.
- Less sensitive to outliers

$$IQR = Q_3 - Q_1$$



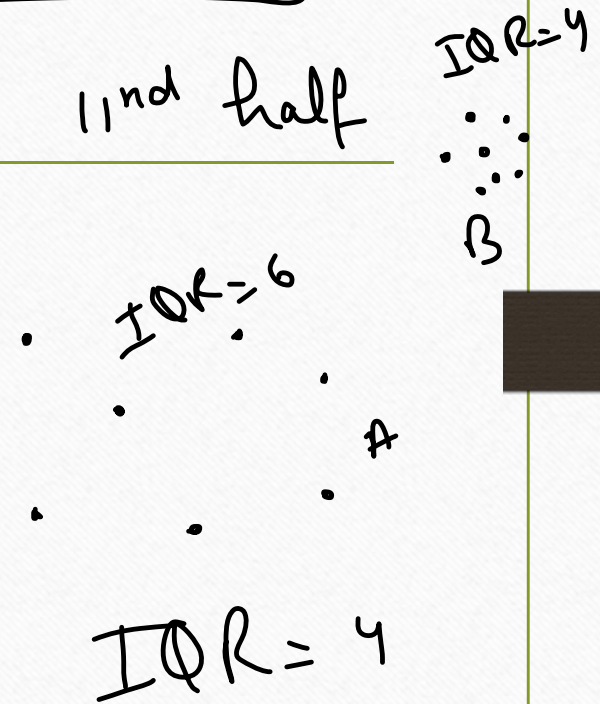
$Q1 = 3$ (Median of first half)

2nd half

$Q2 = \text{Median} \Rightarrow 6$

$Q3 = 9$ (Median of second half)

$$IQR = 9 - 3 = 6$$



"Outlier Detection"

$$UL = 9 + 1.5 \times 6 = 18$$

$$LL = 3 - 1.5 \times 6 = -6$$

[1 2 3 4 5 6 7 8 9 10 11]

$$Q1 = 3$$

$$Q2 = 6$$

$$Q3 = 9$$

$$IQR = 6$$

$$UL \Rightarrow$$

$$Q3 + 1.5 IQR \rightarrow \text{beyond it everything is an outlier}$$

$$LL \Rightarrow Q1 - 1.5 IQR \rightarrow \text{below it everything is an outlier}$$

$Q = [17, 17, 18, 19, 20, 22, 23, 25, 33, 64]$ outlier

$$Q_1 = 18$$

$$UL = Q_3 + 1.5 IQR = 25 + 1.5 \times 7 = 35.5$$

$$Q_2 = 21$$

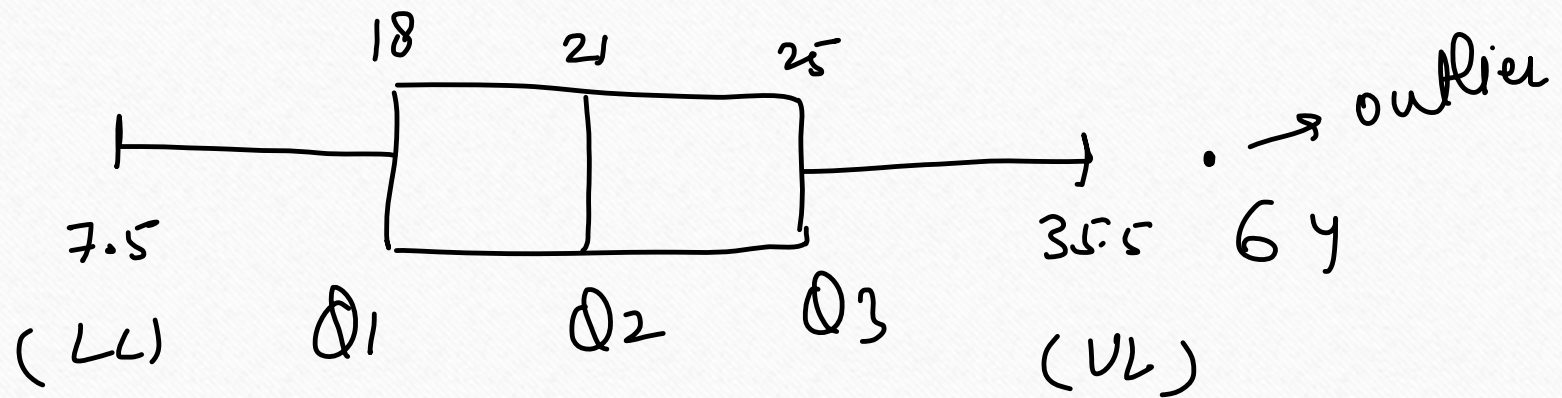
$$LL = Q_1 - 1.5 IQR = 18 - 1.5 \times 7$$

$$Q_3 = 25$$

$$= 7.5$$

$$IQR = 7$$

Box plot → Documentation of
Seaborn / Matplotlib
What is UL & LL?



Checking for Outliers

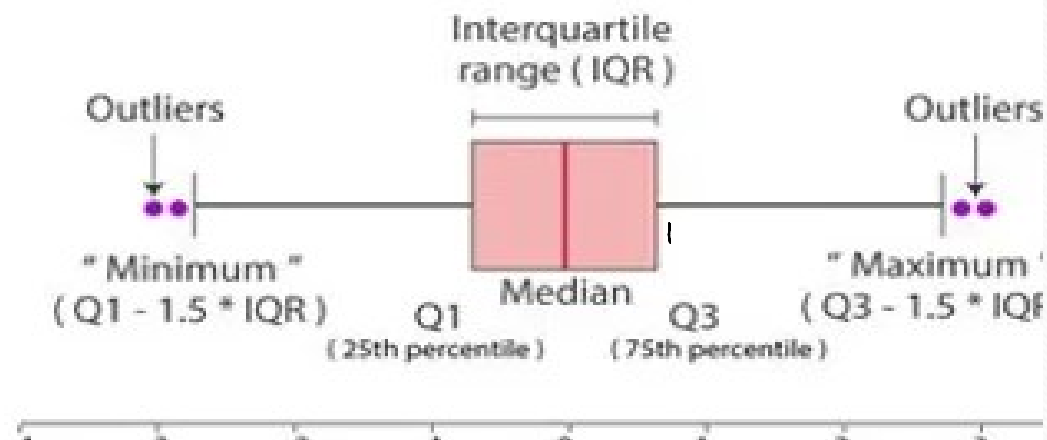
- With Range, Outlier detection isn't possible but with IQR it is possible!!

$$L = Q1 - (1.5 * IQR)$$

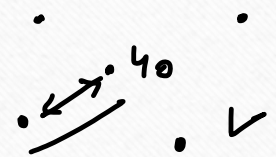
$$H = Q3 + (1.5 * IQR)$$

- Where L is the lower outlier
- H is the higher outlier
- Q1 and Q3 are the average values of those quartiles

Box Plot & 5 Number Summary



$[10, 10, 10, 10]$ $VK: [80, 0, 1, 40, 25]$ ✓
 $[10, 0, 5, -10, 0]$ $MSD: [40, 41, 40, 39, 40]$ ✓
 Variance 40 \therefore ✓



- The variance is a way of measuring spread, and it's the average of the distance of values from the mean squared.

(1)
$$VAR = \frac{1}{n} \sum_{i=0}^n (x - \mu)^2$$
 (2)
$$VAR = \frac{1}{n-1} \sum_{i=0}^n (x - \bar{x})^2$$

$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$

$$v = \frac{\sum_{i=1}^n x_i^2}{n} - \mu^2$$

Handwritten notes: "length" with an arrow pointing to the summation index i in formula (1); "population mean" with an arrow pointing to μ in formula (1); "sample mean" with an arrow pointing to \bar{x} in formula (2).

- The problem with the variance is that it can be quite difficult to think about spread in terms of distances squared.

(2) \rightarrow sample variance \rightarrow estimating pop. variance from sample

[1 2 3 4 5] \rightarrow Mean \Rightarrow 3

variance \Rightarrow $\frac{1}{n} \sum (x - \mu)^2$
avg

$$= \frac{1}{5} \left[(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2 \right]$$

$$= \frac{1}{5} [4 + 1 + 0 + 1 + 4] \Rightarrow \frac{10}{5} = 2$$

$$[-5 \quad 0 \quad 5] \Rightarrow \text{python} \rightarrow \text{variance}$$

$$\frac{(-5-0)^2 + (0-0)^2 + (5-0)^2}{3} = \frac{50}{3}$$

↓

$$\frac{50}{3-1} = \frac{50}{2} = 25$$

$$\text{python} \Rightarrow 25$$