

28/10/21

Basic Statistics

What is statistics ?

It is measurable

It means summary of data or inference of data rather than entire data.

What is data ?

Structured data - excel , table , i.e tabular form

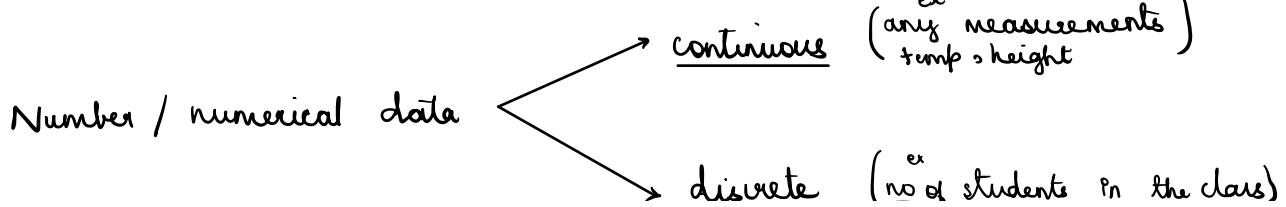
We know the maximum info we can get from the data .

Unstructured data - Video , audio , email

Semi structured data - Dictionary or JSON files } this can be both unstructured or structured
i.e key : value pairs

In machine learning we mainly work with structured data
With NLP we learn to convert unstructured text data to
structured data .

Different data types



Categorical data

Ordinal data
Nominal / categorical / Non ordinal data

→ when there is no sense of numbers. That data is called categorical data.

→ These are not numbers, some sort of category like yes or No , M or F

→ There are few categorical data that can be measured.

Ordinal data

→ These are categorical data that can be ordered.
ex: Movie review in stars , food experience ,
wealth of a family .

exercise

Volume of cornflakes - continuous

speed of car - continuous

population of a town - discrete

length of crocodile - continuous

No. of matches in a box = discrete

temp of oven - continuous

of t-shirts - discrete

No. of goals - discreet

Measures in Statistics

Descriptive Statistics : Those statistics that contain some description

of the data.

e.g.: Rate of employment is 75%.

Ratio of female is 1:2

→ It gives general description of data

- ① Min of data
- ② Max of data
- ③ Count of data - no. of data points
- ④ Median of data
- ⑤ Mean of data
- ⑥ Var of data
- ⑦ Standard deviation of the data
- ⑧ Q.M.

a) Mean \Rightarrow average

$$\frac{\text{sum of data}}{\text{count or number of data points}}$$

count or number of data points

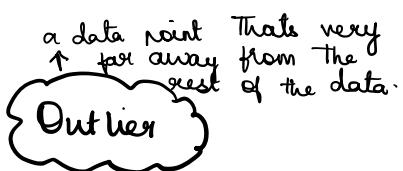
b) Median : middle value of ordered data (ascending / descending order)

* If mean and median are same then its uniformly distributed
otherwise data is not uniformly distributed.

For non uniformly distributed data, if we only look at the mean we get wrong data.

It is less trustworthy.

However Median is not affected by



Median is more trustworthy because it is not affected by outlier.

Median is more robust than mean

Mean is affected by outlier : not trustworthy

** If mean and median are nearly equal then there is no outlier in the data.

o Variance / dispersion / volatility

→ This is the variation in the data

→ When most of data is close to centre, there is less variance or when data is away from centre there is more variance.

Calculate Variance

$$\text{Var} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

→ Robust version of Variance

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

↓
we use median instead of mean.

ex: 1, 3, 5, 7, 9

$$\bar{x} = 5$$

$$\text{Var} = \frac{(1-5)^2 + (3-5)^2 + (5-5)^2 + (7-5)^2 + (9-5)^2}{5-1}$$

$$= \frac{4^2 + 2^2 + 0^2 + 2^2 + 4^2}{4}$$

$$\Rightarrow \frac{16+4+0+4+16}{4} = \frac{40}{4} = 10$$

$$\Rightarrow \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$$

$$\sum_{i=1}^n x_i = \bar{x}n$$

data : $x_1, x_2, x_3, \dots, x_n \rightarrow \text{mean is } m$

then mean of $x_1+c, x_2+c, x_3+c, \dots, x_n+c \rightarrow m+c$

Q1 : 1, 2, 3, 4, 5

$$\text{Mean} = \frac{15}{5} = 3$$

$$\text{Var} = \frac{2^2 + 1^2 + 0 + 1^2 + 2^2}{4} = \frac{10}{4} = 5/2 = 2.5$$

2, 4, 6, 8, 10

$$\text{Mean} = \frac{30}{5} = 6 \quad 3 \times 2 = 6$$

$$\text{Var} = \frac{(6-2)^2 + (4-6)^2 + (6-6)^2 + (8-6)^2 + (10-6)^2}{(5-1)=4}$$

$$= \frac{16 + 4 + 0 + 16}{4} = \frac{40}{4} = 10$$

When data is doubled, the new variance is 4-times of old variance.

Standard deviation

$$SD = \sqrt{\text{Variance}}$$

Mode : The data point that occurs maximum number of times in the data.
→ discrete or categorical data

Range : the value between min and max.

quantile :

Tool to detect outliers in the data using visualisation

90 percentile ⇒ means more than 90% and less than 10% of the people.

25% quantile ⇒ The number that is more than 25% of the data and less than 75% of the data.

50% quantile ⇒ Median i.e. 50% of data is less than this value and 50% of data is more than this value.

doubt

Ex: 101, 102, 103, 104, ..., 200

$$25\% = 1/4$$

$$25^{\text{th}} \text{ data} = 125$$

$$26^{\text{th}} \text{ data} = 126$$

$$\therefore 25 \text{ quantile} = 125 + 126 = 125.5$$

$$50\% = 1/2 = 150 + 151 = 150.5$$

$$75\% = 175.5$$

ex: 101, 102, 103, ..., 199

$$25\% \text{ quantile: } 99 \times 1/4 = 24.75 \approx 25$$

$$= 125$$

$$50\% \text{ quantile: } 99 \times 1/2 = 49.5 \approx 50 = 150$$

$$75\% \text{ quantile: } 99 \times 3/4 = 74.25 \approx 75 = 175$$

09/10/21

Boxplot: It's a tool to detect an outlier

→ data

→ 50% is the mid quartile

→ Calculate 25% and 75% quartile
1st quartile 3rd quartile

ex: 101, 102, ..., 200

$$25\% \text{ quantile} | 1^{\text{st}} \text{ quartile} = 100 \times 1/4 = 25$$

$$= \frac{125 + 126}{2} = 125.5$$

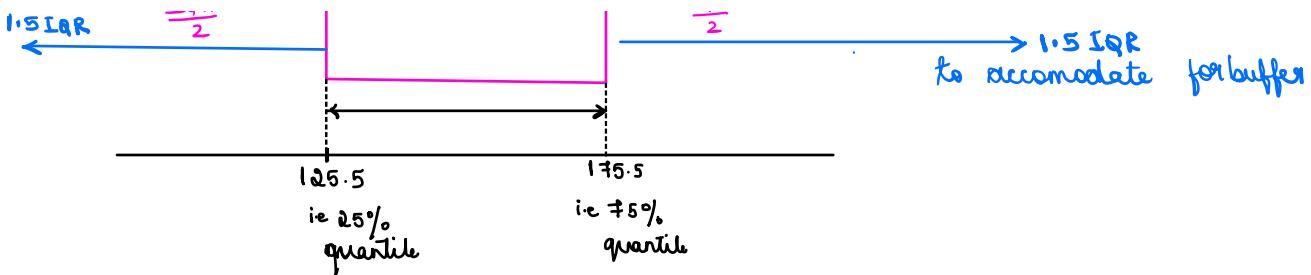
$$75\% \text{ quantile} | 3^{\text{rd}} \text{ quartile} = 100 \times 3/4 = 75$$

$$= \frac{175 + 176}{2} = 175.5$$

IQR → 50% of the data

25%
= IQR





- IQR → Inter quartile range. It's the distance between 1st and 3rd quartile.
- It has 50% of the data

On left & right we should have allocated only $\frac{1}{2}$ IQR but to give buffer space to accomodate non-uniform data we give $\frac{3}{2}$ IQR

outlier

Any data outside of this buffer space ie $\frac{3}{2}$ IQR is considered to be an outlier.

ex: 0, 101, 102, 103, , 198, 199, 300

Count = 101

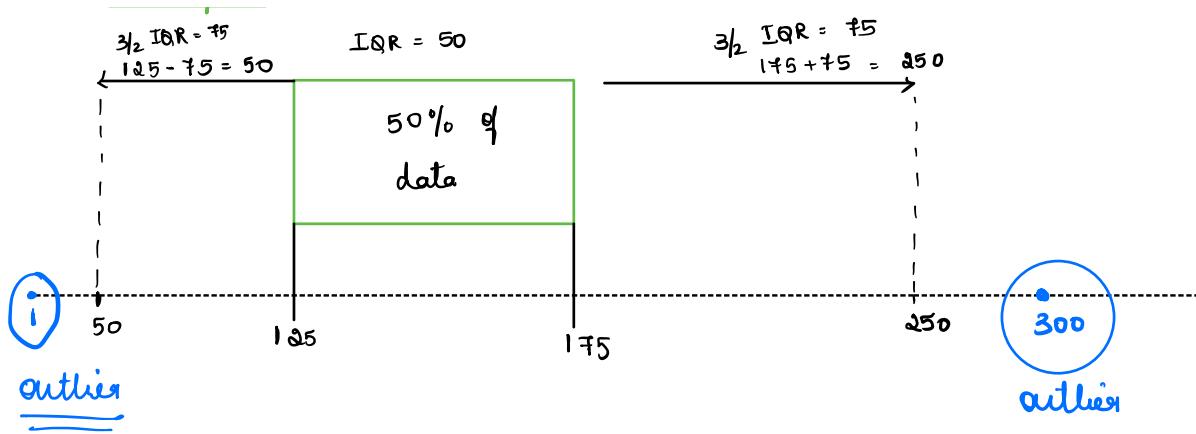
$$25\% \text{ quantile} / 1^{\text{st}} \text{ quartile} = 101 \times 1/4 = 25.25 \approx 26 \\ = 125$$

$$75\% \text{ quantile} / 3^{\text{rd}} \text{ quartile} = 101 \times 3/4 = 75.75 \approx 76 \\ = 175$$

$$\text{IQR} = 175 - 125 = 50$$

$$\frac{3}{2} \text{ IQR} = 50 \times \frac{3}{2} = 75$$

Box plot



Mython