ML | Stats | EDA & FE | DL    Brythd
4      1        2        3    5

Task 1 → PYTHON — DAILY at Ⓗ & W at Ⓢ & Ⓢ → 3, 3

Task 2 → STATS

Task 3 → EDA & Feature Engineering.

Task 4 → ML

## INTRO TO STATS:

1. Descriptive STATS

* Measure of control Tendency.
* Measure of Dispersion.
   ↓
Summarizing the data.

Probability, Permutation, Mean
Mode, variance, median, SD
standard Deviation.

1. Gaussian Distribution
2. Lognormal    "      } – PY
3. Binomial     "
4. Bernoullis   "
5. Pareto       " (Power law).
6. standard normal dist.
7. Transformation & standardization.
8. Q-Q Plot

2. Inferential stats.
→ Z-test — PY
→ T-test — PY
→ ANOVA
→ CHISQUARE
→ HYPOTHESIS testing.
      ↓
   NULL Hypothesis
      &
   Alternate Hypothesis

→ confidence Intervals.

→ Z table, t-table, Chisquare table.

# What is statistics?

* stats is a science of collecting, oraganizing and analysing data.
* It helps out how to used the data in perfect way.
* Better decision Making.

## Define data?

Pieces of information that can be ~~that~~ a measured.

Eg: The IQ of a class.

$\{98, 97, 60, 55, 75, 85\}$

Ages of students of a class.

$\{30, 25, 24, 23, 27, 28\}$ — DATA.

## Types of statistics:

① Descriptive stats.

It consists of organizing and summarizing data.

② Inferential stats.

It is technique where we used the data that we have measured to form conclusions.

1. Descriptive Stats.

Eg: * class room of Maths student (20)

   Marks of the 1$^{st}$ sem
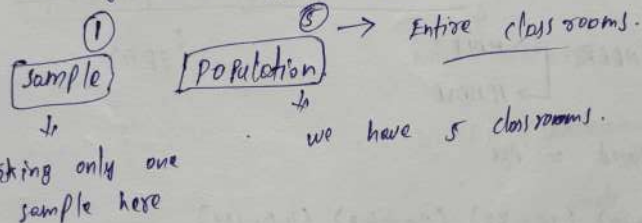
   84, 86, 78, 72, 75, 65, 80, 81, 92, 95, 96, 97, - - - -

2. Inferential Stats.
* What is the average ~~age~~ marks of the student in the class?
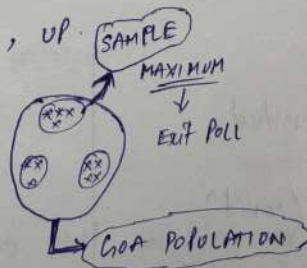
   Mean, Median, Mode (or) SD.

Eg: Are the ~~age~~ marks of the students of this classroom similar to the age of the maths class room in the college?



   ① Sample        ⑤ Population → Entire class rooms.

   ↓                    ↓
   taking only one      we have 5 class rooms.
   sample here

POPULATION AND SAMPLE                    POPULATION (N)

   Elections → Goa, UP  SAMPLE           SAMPLE (n)
                        MAXIMUM
                           ↓
   Exit Poll            Exit Poll

                    GOA POPULATION

1) Here Goa Population is big. | 2) So media will take samples from different locations.

SAMPLING TECHNIQUES:

1) Simple Random Sampling: Every member of the population (N) has an equal chance of being selected for your sample (n).

2) STRATIFIED SAMPLING: where the population (N) is split into non-overlapping groups (Strata).

Eg1: GENDER ⟶ MALE
              ⟶ FEMALE                              ↓ BP (1)

Eg2: Based on Age

(0-10) (10-20) (20-40) (40-100)

Eg3:

3) Systematic Sampling:

(N) → $n^{th}$ individual.

Eg: Mall ⟶ Survey (covid)
        ⟶ $8^{th}$ person → survey.
        ⟶ Every $1^{st}$ person → survey.

In whole population (N). Here we select sample (n) in a Systematic way.

4) convenience Sampling : Doing survey on specific topic. ⑤

    Survey will be based on domain knowledge (or) interest.

    Data Science
    ↓
we will survey only who has knowledge about data science.

    Stack overflow
        ↓
    we will survey only developers.

Eg: EXIT POLL     RBI → survey with house hold.
                  ↓
                what sampling.
                  ↓
             WOMEN (convenience sampling
                          (or)
                     streets ).

Eg: DRUG TEST
      ↓
    what kind of sample?
      ↓

VARIABLES :

A variable is a property that can take any value.

   Eg: Height = {182, 172, 180, 190}
       weight = {78, 75, 85, 90}

  Two kinds of variables.

1) quantitative variables.

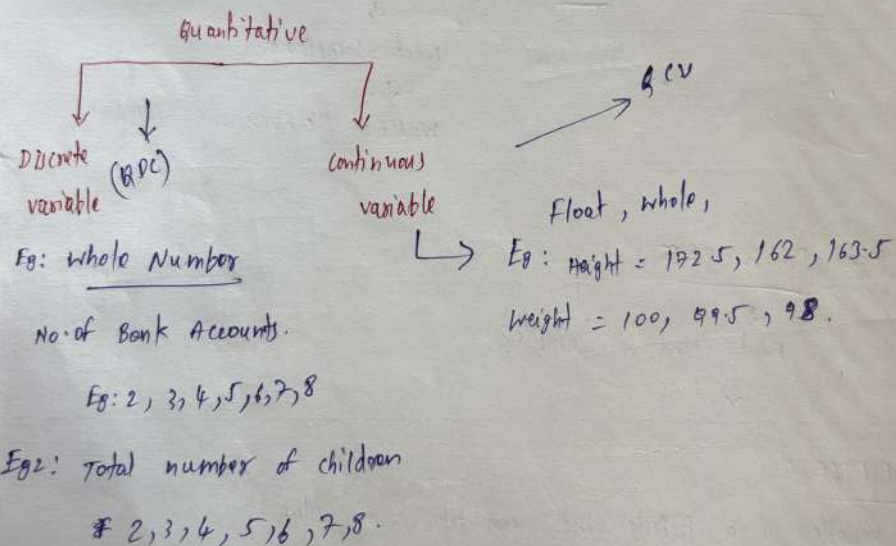2) qualitative variables (or) categorical variables.

1) Q V → Measured Numerically

    Eg: Age, weight, Height

2) C V → It should be converted to numeric data.

    Eg: Gender ⟶ M { Based on some characteristics we can
            ⟶ F { derive categorical variable y.

    Eg : T-shirt sizes.

        L, XL, M, S

           Quantitative

                                      → Q CV

Discrete (QDC)             Continuous
variable                   variable       Float, whole,

Eg: whole Number              ⟶ Eg: Height = 172·5, 162, 163·5

No. of Bank Accounts.             weight = 100, 99·5, 98.

    Eg: 2, 3, 4, 5, 6, 7, 8

Eg2: Total number of children

    ≢ 2, 3, 4, 5, 6, 7, 8.

1) What kind of variable  Gender is?  Q CV (Categorical)

                 Martial status ? CV    "

                 River length ? Q CV

                 Song length ! Q CV

Variable Measurement scales :

4 types of measured variable.

1) Nominal Data.

categorical data → colors, Gender, Type of flower.

2) ordinal data

order of the data matters. But value does not

Eg: Students (Marks)    RANK
                               → ordinal
       100            1         data.
        96            2
        85            3
        88            4

3) Interval data :

Both order of data & value matters.

Eg: Temperature.

*) Fahrenheit
   70-80    80-90    90-100    0
You will some range of values & order also matters.

*) 4) RATIO DATA:

# FREQUENCY DISTRIBUTION:

sample data set: Rose, lilly, sunflower, Rose, lilly, lilly, Rose

sunflower

frequency distribution table

| flower type | frequency |
|-------------|-----------|
| Rose | ·3 |
| lilly | 3 |
| Sunflower | 2 |

Total number of flowers

↓

cummulative frequency

3

6

8

## 1) BAR GRAPH



Discrete variables.

Rose    lilly    sunflower

## 2) HISTOGRAMS

→continuous variables.

Ages = {10, 12, 14, 18, 50, 80, 88}

bin = 10



kernel density estimator.

we do bins in histograms.

## BAR VS HISTOGRAM!

PDF: Probability Density function.

↓

smoothing of histograms.

## 1) ARITHMETIC MEAN for POPULATION & sample.

Mean (Average)

Population (N)          Sample (n)

1, 1, 2, 2, 3, 3, 4, 5, 5, 6.

$$M = \sum_{i=1}^{N} \frac{x_i}{N}$$

$x = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$

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

$$M = \frac{1+1+2+2+3+3+4+5+5+6}{10} \xrightarrow{\text{(n)}} \text{(N)}$$

$= 3.2$

$M = \frac{32}{10} = \boxed{3.2} \rightarrow \text{Average}$

Mean = is a part of central tendency.

## central Tendency:

→ It refers to the measure used to determine the centre of the distribution of data.

→ centre of the data.

1) Mean    2) Median    3) Mode.

## 2) Median.

New element

Data = $\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$, $\widehat{100}$ ← we consider those as $\boxed{\text{outliers}}$

$\text{Mean} = \frac{32 + 100}{11} \Rightarrow \frac{132}{11} \Rightarrow 12 \quad \boxed{M = 12}$

① $M = 3.2$

② $M = 12$ huge difference

* we must be careful with outliers.

Median          outliers

MEDIAN:   $\{1, 1, 2, 2, 3, \widehat{3, 4}, 5, 5, 6, \boxed{100}, 112$     odd = 12

to do

1st thing ↑ 1) Sort the numbers in median.

$\text{Mode} = \frac{3+4}{2} = \boxed{3.5}$

Median

odd number = 11

Median = 3

when you got ② outliers Median works better than mean.

MEDIAN WORKS WELL WITH $\boxed{\text{OUTLIERS}}$

MODE: $\{1, 2, 2, 3, 4, 5, 6, 6, 6, 7, 8, 100, 200\}$

Measure of central tendency.

MODE = { MOST FREQUENT ELEMENT }

MODE = 6

| TYPE of flower | Petal length | Petal width | DATASET |
|---|---|---|---|
| ROSE | | | |
| LILLY | | | |
| SUNFLOWER | | | |
| — | | | |
| — | | | |
| — | | | |

→ Missing value → can be Replaced with Most frequent occuring element.
⇓
work well with categorical Variable.

Ages of Students.

Age

25 → we use Mean

26

—
—
—

32

34

38

MEASURE of DISPERSION:                     Dispersion

    └→ SPREAD

1) VARIANCE                        How Good your data is spread.

2) STANDARD DEVIATION

1) VARIANCE

POPULATION VARIANCE                    SAMPLE VARIANCE

$$\sigma^2 = \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{N}$$

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

| $x$ | $\mu$ | $x - \mu$ | $(x_i - \mu)^2$ |
|---|---|---|---|
| 1 | 2.83 | -1.83 | 3.34 |
| 2 | 2.83 | -0.83 | 0.68 |
| 2 | 2.83 | +0.83 | 0.68 |
| 3 | 2.83 | 0.17 | 0.03 |
| 4 | 2.83 | 1.17 | 1.37 |
| 5 | 2.83 | 2.17 | 4.71 |

$\mu = \frac{17}{6}$

$\mu = 2.83$

$\Rightarrow \frac{10.84}{6}$

$\Rightarrow 1.81$

BELL CURVE



2.83  4.12

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

      ① 2.83
      1.34

$= \sqrt{1.81}$    4.17

$= 1.345$

1)        2)



Which will have more
variance (Spread) ②

Percentiles and Quartlies { 1st step to find outliers }. (12)

Percentage : 1, 2, ③, 4, 5.

   % of the numbers that are odd?

$$\Rightarrow \% = \frac{\text{\# of numbers that are odd}}{\text{Total numbers}}.$$

$$\% = \frac{3}{5} = 0.6 = 60\%.$$

Percentiles :

   A percentile is a value below which a certain percentage of observation lie.

Data set : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8,8, 8, 8,8, 9, 9, 10, 11, 11, 12.

   what is the percentile ranking of 10?

Percentile Rank of 10 = x

$$x = 10 \Rightarrow \frac{\text{\# of values of below } x}{n} \times 100$$

$$= \frac{16}{20} \times 100$$

$$= 80\%.$$

$$\Rightarrow \quad x = 11 \Rightarrow \frac{17}{20} \times 100$$

$$\Rightarrow 85\%.$$

# Five number Summary:

1) MINIMUM    2) First Quartile ($Q_1$)

2) Median    4) Third Quartile ($Q_3$)

5) MAXIMUM

* This is used to [remove outliers].

{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6,6, 6, 6, 7, 8, 8, 9, (27)}
-50

When ever you want to remove outlier ?                         outlier

[Lower fence ⟷ Higher fence]

Lower fence = $Q_1 - 1.5 (IQR)$

upper fence = $Q_3 + 1.5 (IQR)$   Interquartile Range.
                                      ↳

$(IQR) = Q_3 - Q_1$  ⟹  $Q_3 = (75\%)$

                          $Q_1 = (25\%)$

After applying z-score.

Data = $\{1, 2, 7, 4, 5, 6, 7, 8\}$

$$z(1) = \frac{1-4}{1} = -3, \quad z(2) = \frac{2-4}{1} = -2, \quad z(3) = \frac{3-4}{1} = -1$$

It is easy to calculate S.D. using z-score.

$$SD = \{-3, -2, -1, 0, 1, 2, 3, 4\}$$

$\{1, 2, 3, 4, 5, 6, 7, 8\}$ After Applying z-score $\{-2, -3, -4, 0, 2, 3, 4\}$

↳ This is normal distribution.    ↳ standard normal distribution

## PRACTICAL APPLICATION:

our mean should be $M = 0$, & $SD = 1$

DATA SET:

$$\boxed{\mu = 0, \ \sigma = 1}$$

| AGE | SALARY | WEIGHT |
|-----|--------|--------|
| Years | RS | Kg |
| 24 | 40k | 70 |
| 25 | 80k | 80 |
| 26 | 60k | 85 |
| 27 | 80k | 99 |

↳ Convert this to Standard normal distribution using z-score.

↓

This Process is called as

$$\boxed{\text{STANDARDIZATION}}$$

↓

Inside z-score will be applied

## NORMALIZATION:

↳ MINMAX SCALER → (0 to 1)

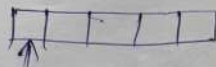Normalization gives you the Process where you can define lower and upper bond.

CNN → Image classification

0 - 225 is converted to ①

B - 1

```
┌──┬──┬──┬──┬──┐
│↑ │  │  │  │  │
└──┴──┴──┴──┴──┘
```
Normalization.

Practical eg:

In 2021

$M = 250, \ x_i = 240, \ \sigma = 10$

In 2020

$M = 260, \ x_i = 245, \ \sigma = 12$.

DAY-3

1) Distributions

↳ Normal dist

↳ standard

{ India v/s SA }

1) ODI Series        2021

Score Average $2021 = 250 \Rightarrow M$

Standard deviation $= 10 \Rightarrow \sigma$

Rishbpant ~~final~~ score $= 240$
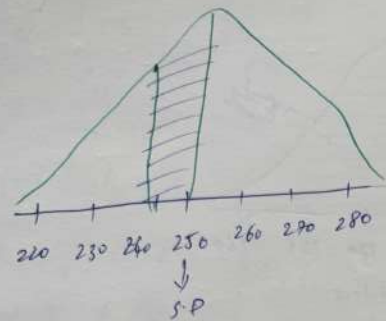          Avg

2020

Score Average $2020 = 260$

$SD = 12$.

Rishbpant ~~final~~ Score $= ~~tb~~ \ 245$
          Avg

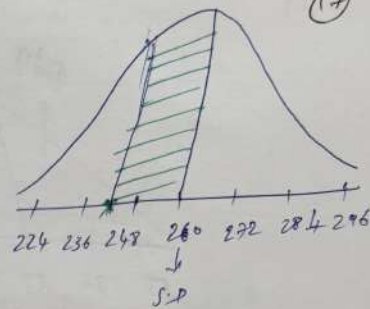compare both the series in which year Rishab Pant final Score was better?

→ Here we apply z-score for 2021

$z \ score = \dfrac{x_i - M}{\sigma} = \dfrac{240 - 250}{10} \Rightarrow \dfrac{240}{10} \not{18}$    $\dfrac{-10}{10} = -1$

$z \ score = \dfrac{x_i - M}{\sigma} = \dfrac{245 - 260}{12} \Rightarrow \dfrac{192}{12} \Rightarrow = \dfrac{-15}{12}$

$= -1.25$

220  230  240  250  260  270  280

$\downarrow$
S.P

z score = -1

224  236  248  260  272  284  296

$\downarrow$
S.P

z-score = -1.25.

## Example for z score :

stats interview question.

$M = 4$
$\sigma = 1$

What percentage of scores falls above 4.25?

x.



Body          Tail

1  2  3  (4)  5  6  7  8

$$z = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1}$$

$$\boxed{z = 0.25}$$

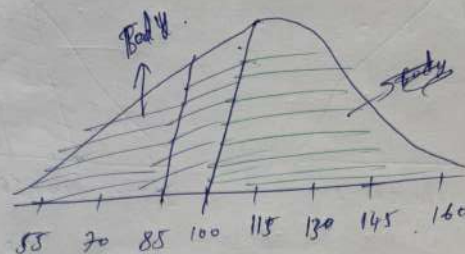z score helps you to find area of the body curve.

Look At  [ z -table ]

$1 - 0.5987 \Rightarrow 0.4013$

* In india the avg IQ is 100, with a S.D of 15. What % of population would you expect to have an IQ lower than 85 ?

$M = 100$
$\sigma = 15$

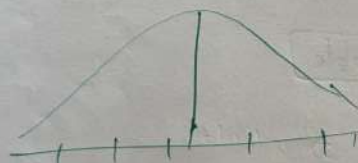$$z \text{ score} = \frac{85 - 100}{15} = \frac{-15}{15} = -1$$

$$= -1 - 0.84884$$

$$= \underline{-1.1586} = .$$
$$4$$

Iq b/w 90 to 120.

$$z = \frac{90 - 100}{15} \Rightarrow$$

$$z = 120.$$

DAY - 4



SD 1 $\Rightarrow$ 68%.

SD 2 $\Rightarrow$ 95%.

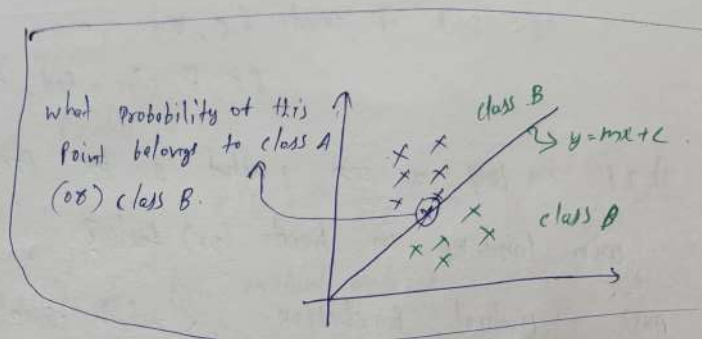SD 3 $\Rightarrow$ 99.7%.

After (SD3) all are outliers.

It can be treated as outlier after SD3.

$$Z\text{-score} = \frac{x_i - \mu}{\sigma} \quad (\text{or}) \quad \frac{x_i - \mu}{\dfrac{\sigma}{\sqrt{n}}}$$

## * Probability:

Probability is a measure of the likelihood of an event.

what probability of this point belongs to class A (or) class B.



Eg: Roll a dice $\{1,2,3,4,5,6\}$.

What is $Prob(6) = \dfrac{\#\ of\ way\ an\ event\ can\ occur}{\#\ of\ Possible\ outcomes}$

$$= \frac{1}{6}$$

Eg: Toss a coin. $(H, T)$

$$P(H) = 1/2$$

2) Addition Rule (Probability, "or")

Mutual Exclusive Events.

Two events are mutually exclusive if they cannot occur at same time.

Eg: Rolling a dice $\{1,2,3,4,5,6\}$

## Non -Mutual Exclusive.

Multiple events can occur at the same time.

Eg: Deck of cards $\{Q, \heartsuit\}$

$\{K, \heartsuit, color, Red, Black\}$.

1) If i toss a coin, what is the Probability of the coin landing on heads (or) tails?

Ans) Mutual Exclusive. → which is also called $\omega$ mutual exclusive.

$Pr ( H \text{ (or) } T )$

$Pr (A \text{ or } B) = P(A) + P(B)$

$= \frac{1}{2} + \frac{1}{2}$

$\boxed{Pr (A \text{ or } B) = 1}$

Roll a dice

$Pr(1 \text{ or } 3 \text{ or } 6) = P(A) + P(B) + P(c)$

$= \frac{1}{6} + \frac{3}{6} + \frac{6}{6}$

$= 1 + 2 + 1 \qquad = \frac{3}{6} = \frac{3}{2}$

$= 7$

$= 0.5$