

Data Science

• Statistic

↳ sample

- ↳ mean
- ↳ std deviation
- ↳ variance

• parameter

↳ population

- ↳ mean
- ↳ std
- ↳ variance

• Descriptive statistics

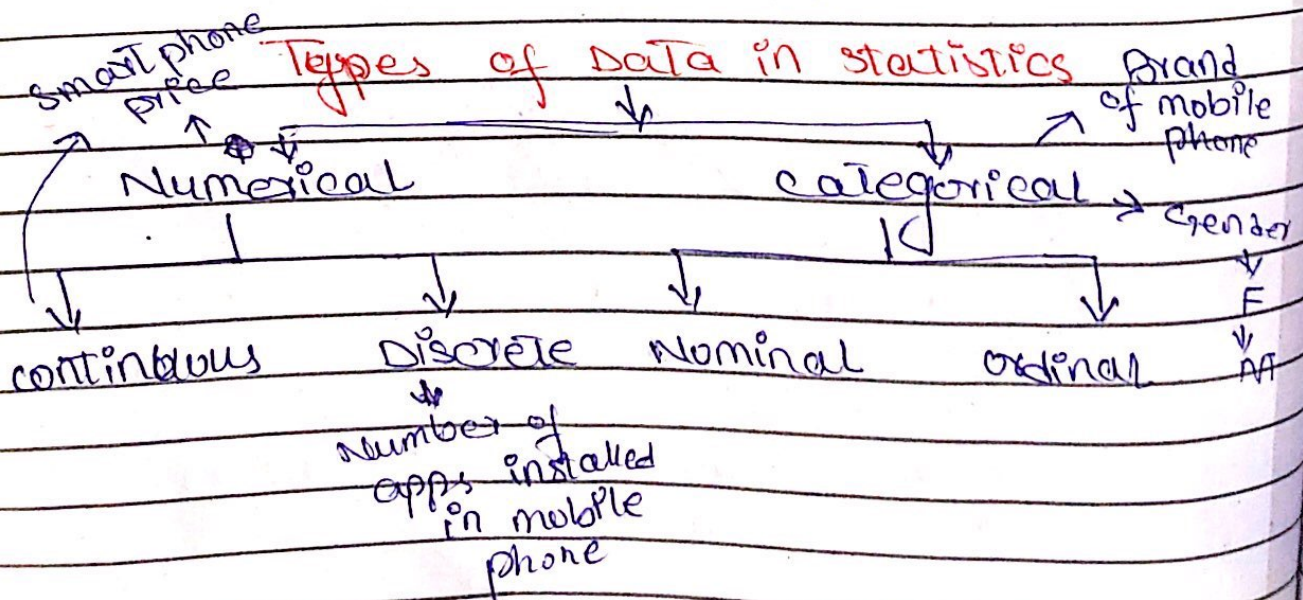
- mean
- median
- mode
- std deviation
- variance
- MAD

↳ sampling from population. Draw conclusion from sample

(two variables)
Bivariate
↳ scatterplot

• Inferential statistics

- confidence Interval
- Hypothesis Testing



→ order or relationship between values is ordinal.

• categorical (nominal) → pie chart
→ bar chart
(ordinal) → pie, bar

Numerical
→ histogram

→ no order or no relationship between values is nominal.

1st Visualization Technique:

Numerical → Histogram
Bivariate (two variables) → scatterplot
numerical → continuous → univariate

~~univariate~~

→ Boxplot
→ KDEplot

• Categorical (nominal, ordinal)

→ pie chart
→ bar chart

Bivariate

• Boxplot

→ x-axis → categorical
→ y-axis → numerical

→ one categorical, one numerical
→ Boxplot

• Both numerical values → scatter plot

• Both categorical → Bar chart
→ heatmap

2nd Measure of central Tendency:

→ where our values from data occurs the most

→ mean (when data have outlier, mean don't be used)

Boxplot is used to find outlier.

is too many
data to
draw it

→ When data have outliers, use median.

→ When range of data is high, median doesn't effect which is another good about median.

→ We need to find three of them (mean, median, mode), but use them properly and to get fair and proper idea of centrality.

3. Measure of spread

→ How much our data is spreaded

1. Range

2. Interquartile Range

3. Std deviation

4. Variance.

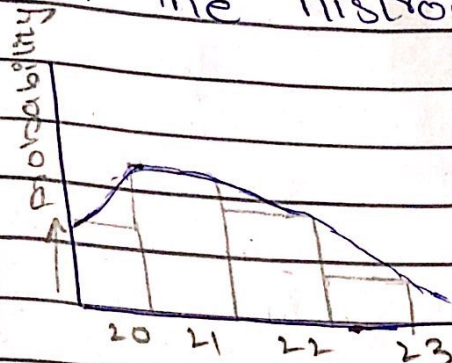
• Range → difference between smallest and largest values in data.

• Percentiles → indicating the below value percentage of given observations in a group of observations falls.

• Quartiles → divide data into quarters.

Data Science

- showing
- probability's ~~func~~ density function (PDF) is KDE
 - smoothen the histogram is simply PDF



• PDF → tell information about only the certain point

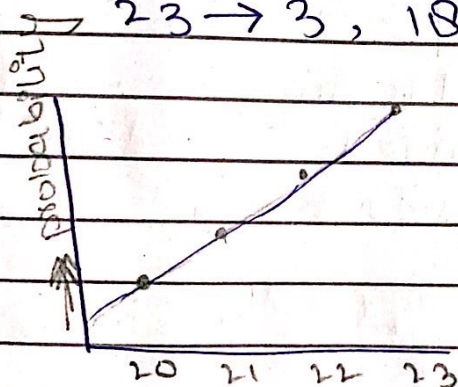
CDF :

20 → 5

21 → 7, $7+5=12$

22 → 6, $12+6=18$

23 → 3, $18+3=21$



• CDF → tells about all values before the certain point.

• Pdf = $\text{count} / \text{sum}(\text{count})$

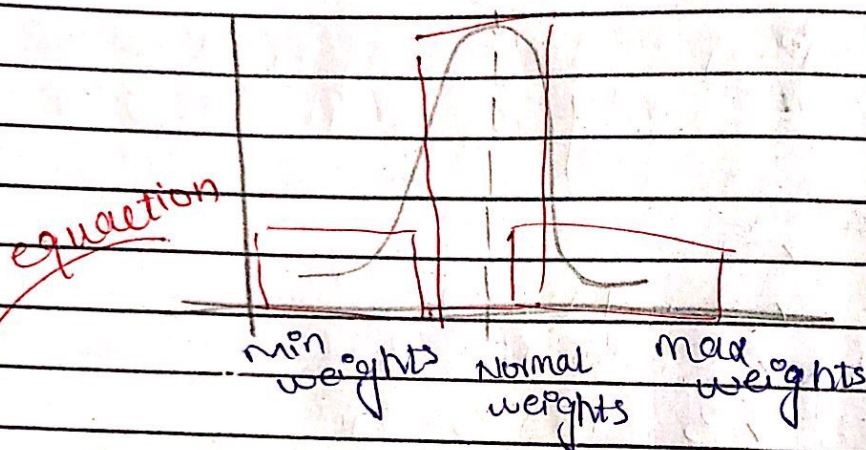
• Cdf = $\text{cumsum}(\text{pdf})$

→ Cdf tells us how much ↑ probability of our result is be correct and how much there probability of of our result be wrong.

→ CDF supports PDF by addition of sum results.

◦ Normal Distribution

- Poisson
- Binomial
- Bernoulli
- Normal Distribution



→ Central Limits Theorem:

We can convert any distribution into normal distribution

$$p(x=\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$\mu \rightarrow$ mean

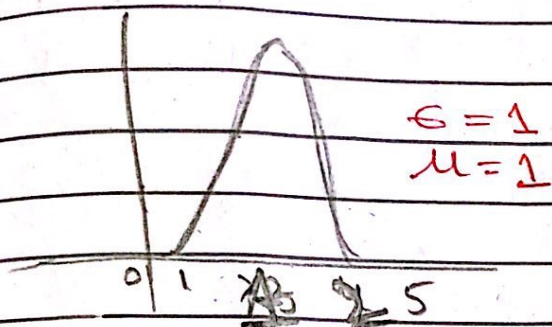
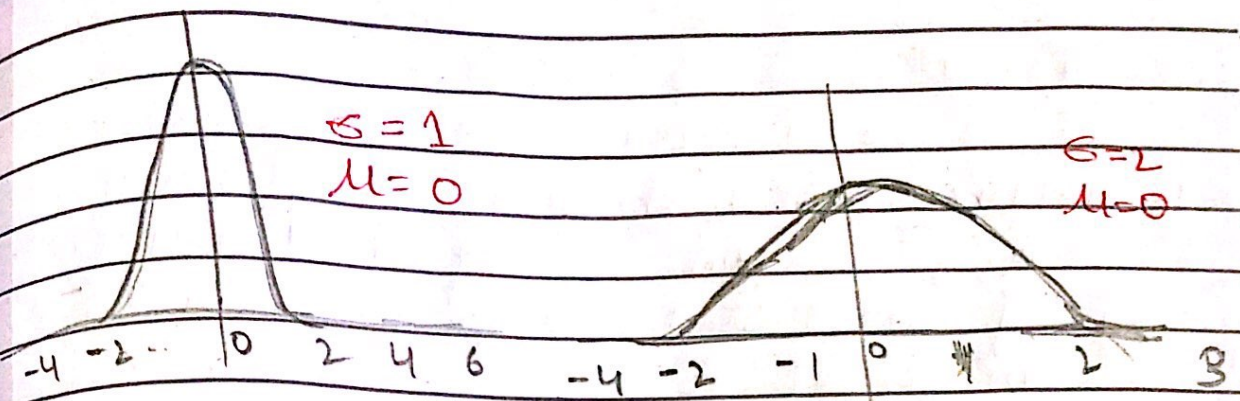
$\sigma \rightarrow$ std deviation

$\sigma^2 \rightarrow$ variance

when $\mu=0$, $\sigma=1$ (gaussian distribution)

$$PDF = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$$

$\frac{1}{\sqrt{2\pi}} \rightarrow$ constant



→ calculating z-score

→
$$z = \frac{\text{total score} - \text{mean}}{\text{std}}$$

→ z value in z-table

→ then multiply that z-table value with 100.

→ Mean, mode, Median is always same for normal distribution.