

Data sets
in tutorials



Data sets in
the wild



A GENDA

- ① 2 - Way anova
- ② KS - Test
- ③ AB - Test
- ④ Parametric Vs Non-parametric tests
- ⑤ Feature Engineering → Heatmap
→ Skewness & Kurtosis
→ Outlier treatment
→ Normalizing Vs Standardizing

2-Way Anova

- * One-way anova ??

Cat Vs Numerical
↓
→ 2 Cat-variable

Aerofit Example → 3 Products Vs Num

10 Products vs Num - Value

① T-test → 2 Products → T-test
P-value
how many times ??

$10C_2 \Rightarrow 45$
tests

~~1000C₂~~ → we don't do this
→ too troublesome

2

Type-II / I error.

0.05 \rightarrow alpha (Level of significance)

F-value $<$ alpha \rightarrow H_0 is actually still true

Incorrectly reject H_0

0.05 - chance of error \rightarrow Z-test

us time

45×0.05

\rightarrow We don't use T-test for more than 2-cat value

→ anova → more than 2 cat values
1 cat - Variable

→ 2-way anova → 2 separate independent cat vari
at once

- (1) Flavour has an impact → sales
 - (2) location where it's sold → sales
 - (3) Impact of Flavour on sales
- depends on location or not
- Interaction effect - 2-way anova

II Diet / Exercise — Weight loss

- Does diet impact weight loss or not
- " Exercise " " " " " "

→ Does impact of diet on weight loss
↓
Impacted by exercise or not

III Impact of Fertilizer type
Watering Frequency,
growth of plant

(1) Impact Fert. → growth

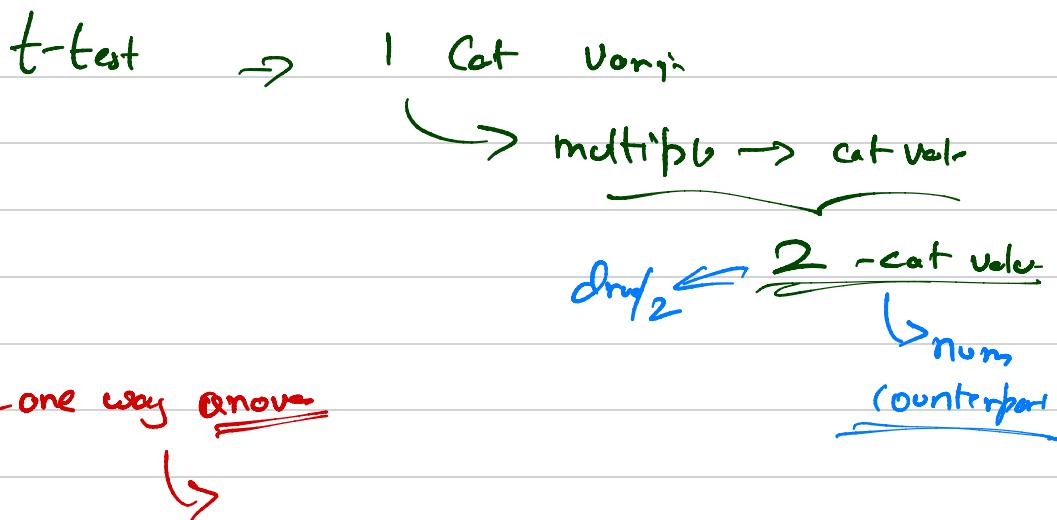
(2) " watering → ",

(3) Impact of Fert. on growth

↓
Impacted by watering freq or not.

t-test - sample (f

? , Value)
→ coming from sr



F -one way anova
 \downarrow

Null hy Po.

\rightarrow Flavour: No impact on Sale

Location: 1, 2, ...,

Interaction: There's no interaction

b/w Flavour & Location

RS - TEST

$M_1 \rightarrow \{5, 6, 8, \dots\}$ \rightarrow mean
8 to

$M_2 \rightarrow \{9, 10, \dots\}$ \rightarrow mean
- 8 to

1 talk sample

Use both t-test & z-test

↑
approximates z-test

For large samples

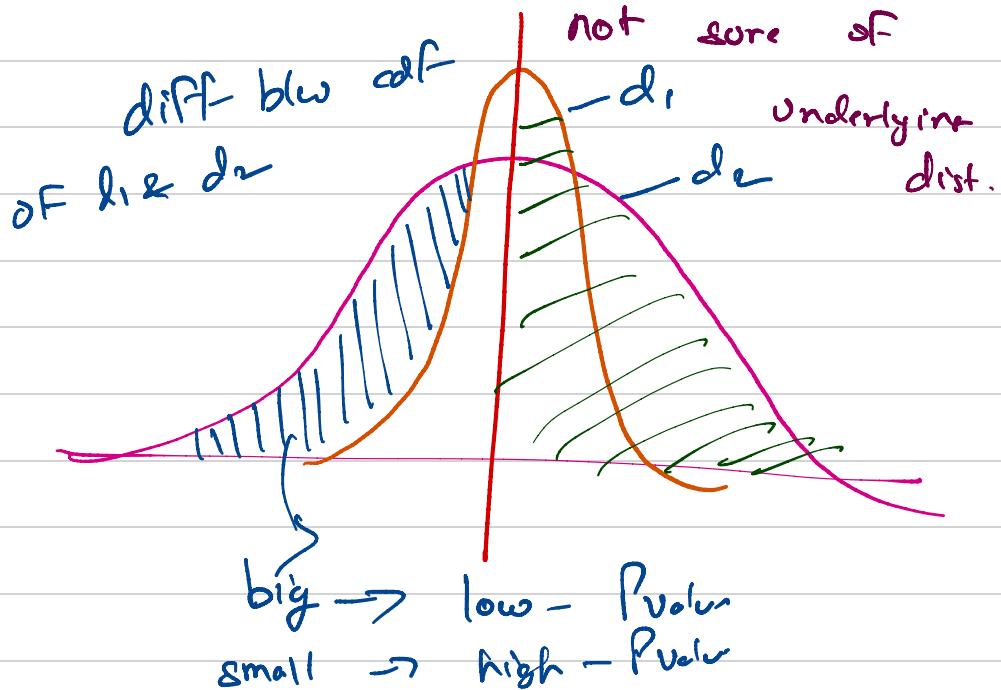
→ Here you make assumptions

①

Normality

K S - Test

→ distribution - Free



- ① Cat Vs Num →
 - < 2 cat vars → χ^2
 - ≥ 2 cat vars → one-way
 - 2 cat columns vs num → 2-way
- ② Cat Vs Cat → Chi-square test
- ③ Num Vs Num → Pearson corr
→ Spearman's

LS - test

H₀: The two dist are same
 $\mu_1 = \mu_2$

H_a: Not same

Reject this??

A B - TEST

→ Amazon → Π_1 model \rightarrow 80% / stable mod'

Π_2 model \rightarrow 90%

offline date

Concerns

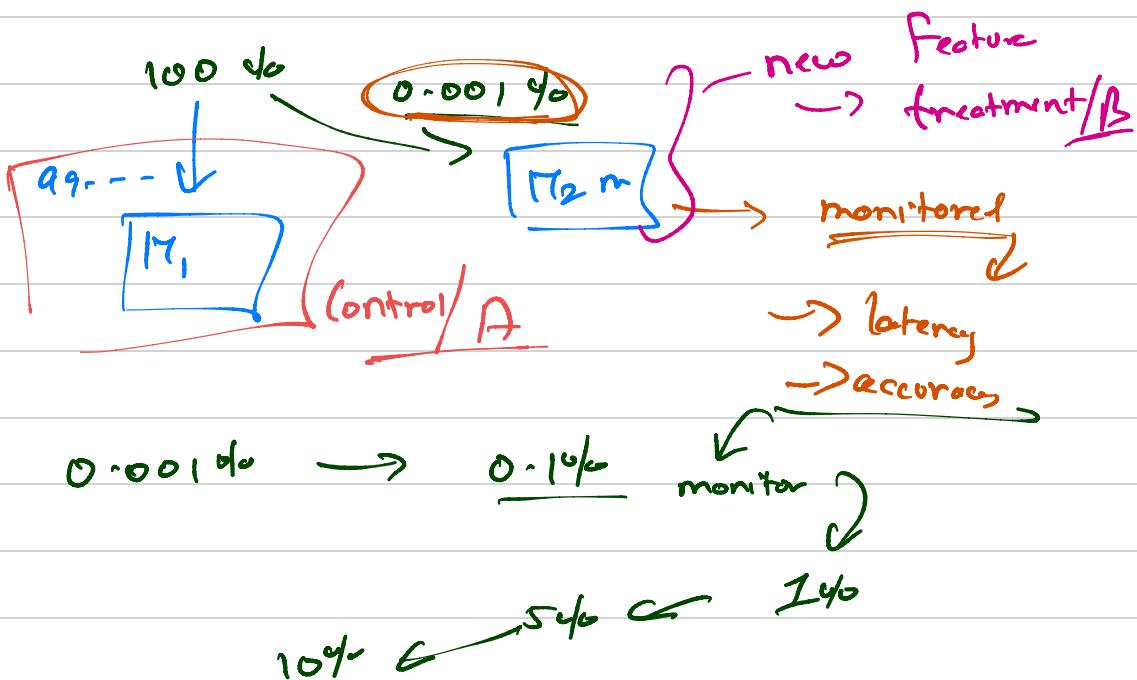
- ① Acc may not hold in real data
- ② May not be stable \rightarrow low throughput
 \rightarrow high latency.

$V \cdot V \cdot V$. accurate \rightarrow entire day??

Compromises b/w

Speed & Accuracy.

M_2 model \rightarrow live traffic



A/B TEST

Drug / Medicine

(200)

(Random)

(100)

Control or A

(100)

Treatment or B

A/B Test

randomized control-
Treatment test

(Paracetamol)

or

Placebo

(New medicine)

days to
recovery

old model

recommendation

algo

new model

U2 - recommendation
algo

Engagement metric

Z-Proportion
test

$$\frac{t + o_F}{t + o_P}$$

engaged user

user

throughput \rightarrow no of requests you can handle in one minute

latency \rightarrow amount of time

\rightarrow respond to one req/^{10 ms}

$$\rightarrow \frac{1000}{60} \rightarrow \frac{60}{1000}$$

throughput \rightarrow 1000 \rightarrow 1 min

$$\rightarrow 1 \rightarrow \frac{1 \text{ min}}{1000}$$

A D Sense \rightarrow one of the highest throughput

9 month \rightarrow women \rightarrow birth

2 women \rightarrow 4.5 month

1 - - - - 2 num data

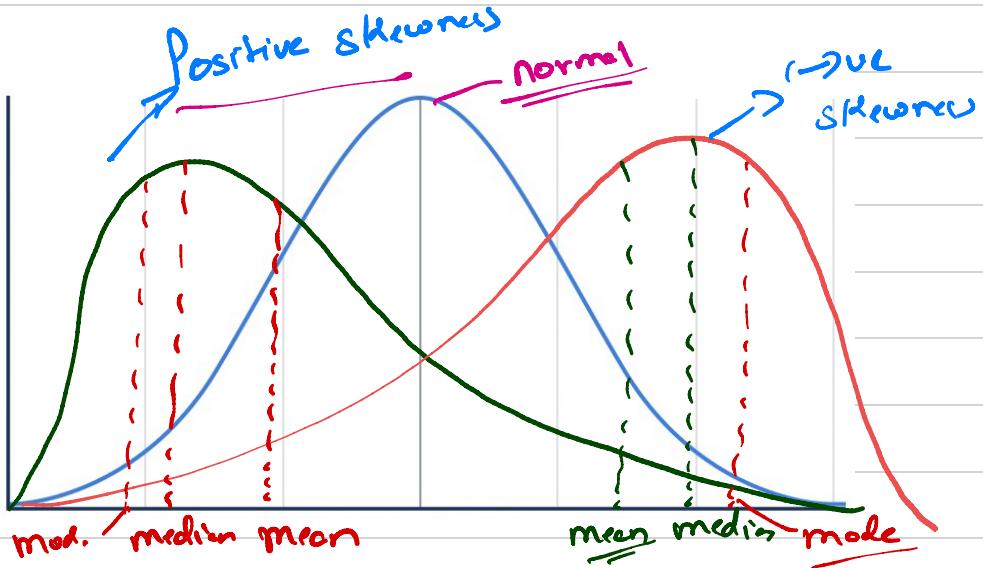
1 - - - } num data

Parametric Us Non-Parametric

make assumptions
data dist

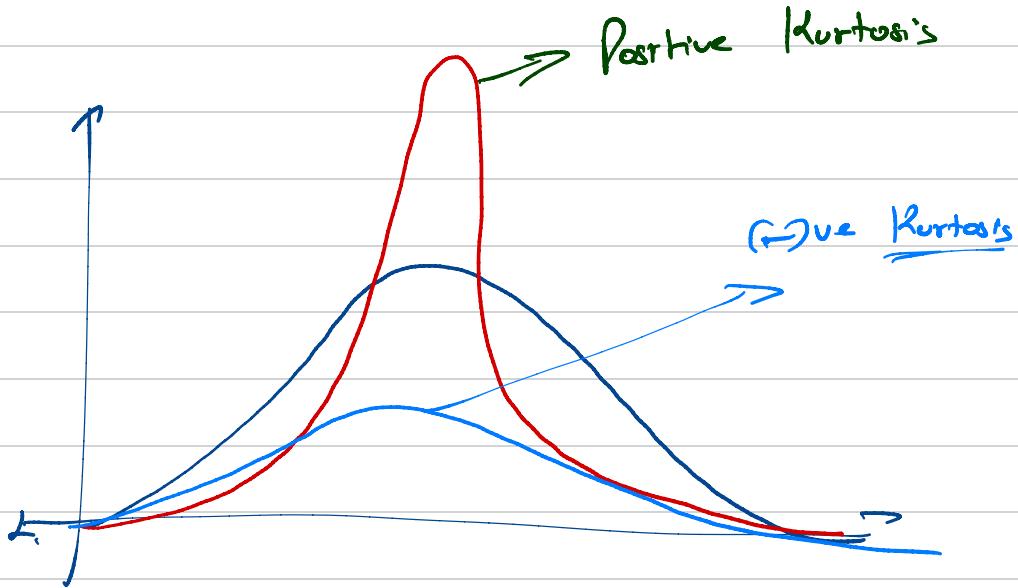
They don't

SKEWNESS



→ (+)ve Skewness median < mean
(-)ve i. median > mean

KURTOSIS



Leptokurtic → Pointy → (+)ve Kurtosis

Mesokurtic → Normal → normal

Platykurtic → Flat one → (-)ve Kurtosis

→ Remove missing Value

- (1) Mean
- (2) Median
- (3) Mode
- (4) Interpolate

$\frac{\text{Pre} + \text{Next}}{2}$
Current
Pre 1 Next

Simple Imputer

OUTLIERS

① IQR $\rightarrow Q_3 - Q_1$ Percentile $\downarrow 25\%$

IQR $(Q_3 - Q_1) \times 1.5$ Percentile $\downarrow 20\%$

$$\left. \begin{array}{l} \text{upper } Q_1 - IQR \\ \text{lower } Q_3 + IQR \end{array} \right\} \quad \begin{array}{l} \text{Points outside} \\ \text{upper \& lower} \\ \underline{\text{outliers}} \end{array}$$

② Z-Score \rightarrow Compute z-score
for all the values

$$\frac{(\bar{x} - \mu)}{\sigma} \quad \begin{cases} < 10 \\ > 10 \end{cases} \quad \underline{\text{z-score}}$$

$50 < 50 = 50$

$$\left. \begin{array}{l} \text{Upper} - \text{mean} + 3\sigma \\ \text{lower} - \text{mean} - 3\sigma \end{array} \right\} \quad \begin{array}{l} \text{Remove them} \end{array}$$



Standardize Vs Normalize

Normalize → Numerical Column

$$\rightarrow [0, 1]$$

$$a = \{1, 2, 3, 4, 5\}$$

Scale of 0
to 1

$$\text{aif} = \frac{\text{a} - \min(a)}{\max(a) - \min(a)}$$

Formula
for
normalizing

ML model → house price prediction

Feature → no of rooms
Ex-price → 120,000, 80,000
directly → ML model

Standardize \rightarrow Reduce data

$$a = [1, 2, 3, 4, 5]$$

\rightarrow 0 mean

1 Variance

$$\frac{x - \text{mean}}{\text{std}}$$

When do we use which?? don't have gaussian dist

If you know Range \rightarrow go for normalization
.. " don't " .. \rightarrow Standardization

have gaussian dist

Stock market

