

20-09-2025

Agenda:

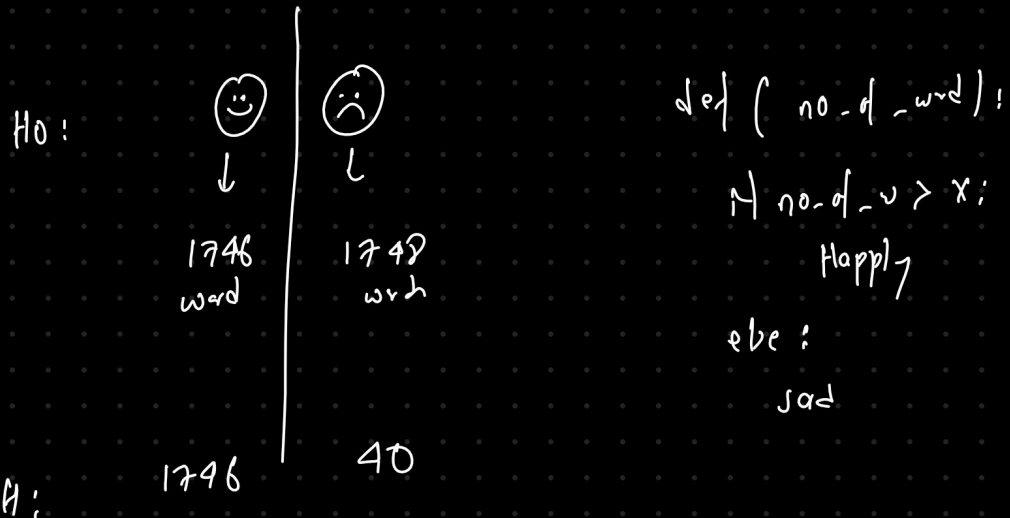
- EDA-1 → movie sentiment } Textual
- EDA-2 → movie sentiment - π
- wine quality } numerical
- commodity data } date, numerical, missing value



- can we claim this effect exist in the larger population, or is it just chance?

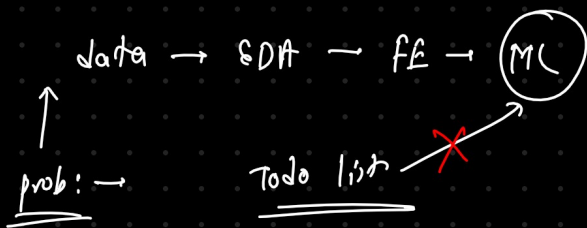
Tools:

- Hypothesis testing
- confidence interval



EDA

- group data / understand data
- data loader / visualization
- finding trends / patterns
- To clean data
- correlation
- skewness
- distribution

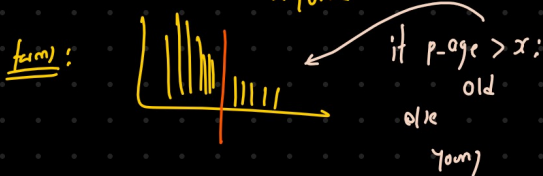


prob: separate young vs old people

data: age, income

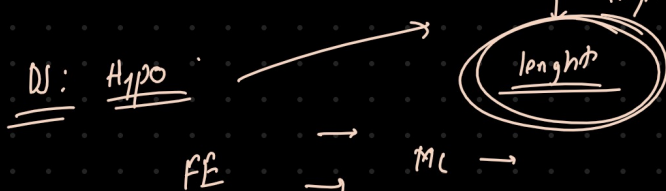
① $\text{data} \rightarrow \text{EDA} \rightarrow \text{FE} \rightarrow \text{ML}$

② data → EDA → think if ML required

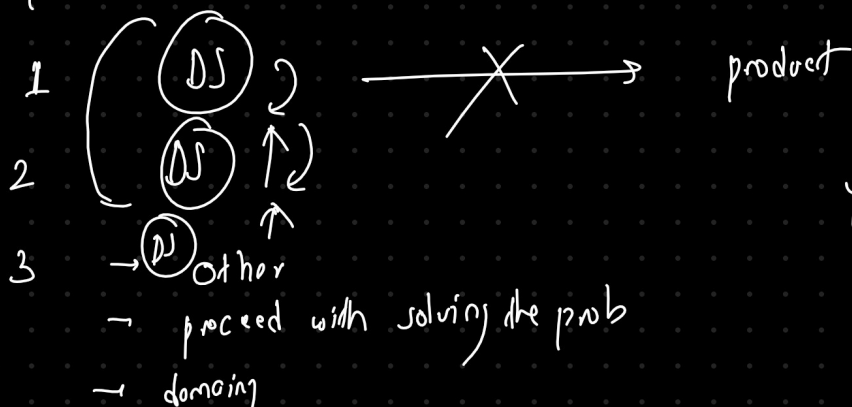


Hypothesis: mult:

diff. between -ve & +ve review length



Year



yt → 1 → 1161

→ 2 → MC

pphes: (DS)

DS → stat

→ CV ← master → lib

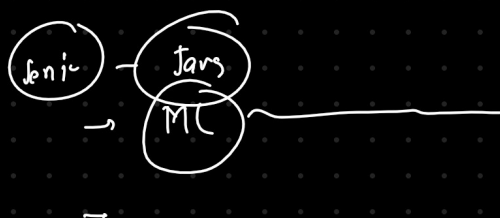
agent



→ CV ← research → bad

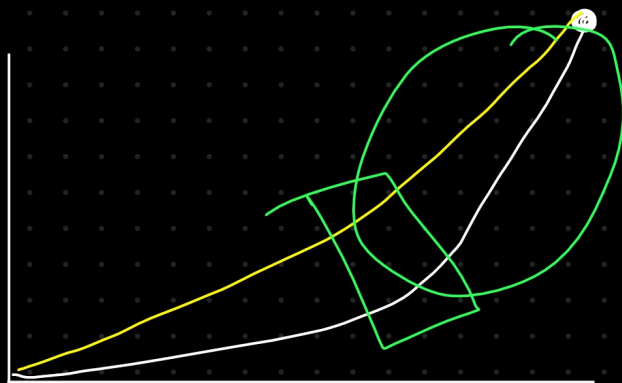
→ CV ← true

→ CV ←



Tech : CS → 5% job

Non-Tech : spec, ↑ → 30%



(X) (Y)

$x_1 m_1$ $x_2 m_2$ $x_3 m_3$

↓
ml
↓

Time-series data:

data that contains date time

Gold price:

date	price	-	-	-
Jan-01-2000				
Jan-02-2000				
⋮				
Jan-01-2025				

time-series → analysis

↳ regression → predict
↳ forecast

Mon Tue Wed Thu Friday
— — — — —

Summary of fill

fill → forward fill

pat9:

10
NaN
15
NaN
20

$$\frac{10 + 15 + 20 + 40 + 50}{5}$$

⇒ 27

Time series

10

15

20

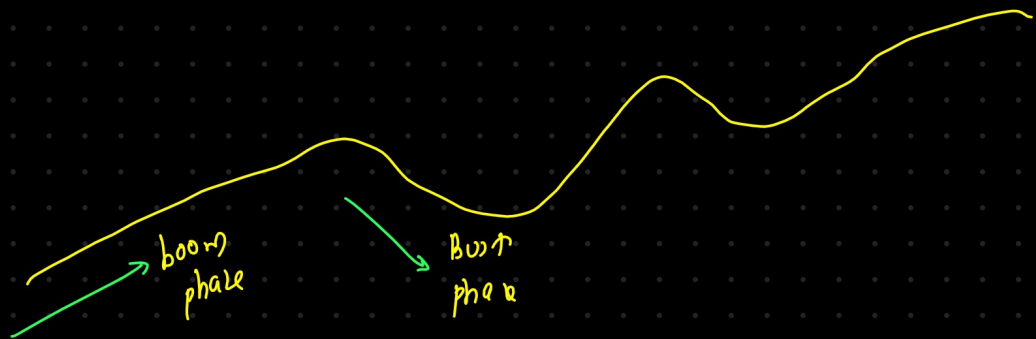
40

50

ffill → missing values would be replaced with the value row above it.

bfill → missing values would be replaced with the value below it.

NaN		NaN		NaN
10	→ ffill	10	→ bfill	10
20		20		20
30		30		30
40		40		40



pct-change

2025-01-01	100	}	→	NaN	←	$(105 - 110) / 110$
2025-01-02	110		0.10			
2025-01-03	105		-0.045			
2025-01-04	126		0.20			

Gold → [.10, -1.5, -0.045, .20]

↓
std

std → mean → spread

A - 30 40 50 60 70

- 10

B - 0 20 50 60 100



std - ↑