

14-09-2025

Agenda:

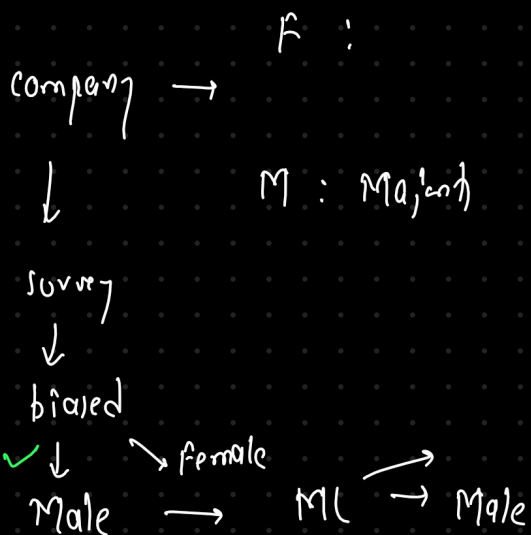
→ FE - 6

→ EDA - 1

→ Data : house price, FE (1 to 5)  
: diabetes dataset (FE - 6)

Imbalanced Data?

→ classification prob, one class has more sample than other.



conv prediction

1000 people → 90 → ML  
→ 910

→ Solution : resampling

910      90  
↓      ↓  
non-C    C

### Techniques to handle Imbalance:

- Random Under-sample (RUS)
- Random Over-sample (ROS)
- SMOTE (Synthetic Minority oversampling Technique)
- ADASYN (Adaptive Synthetic Sampling)

→ 2 US

Male : 100

Female : 30

- Reduce Majority class by randomly removing sample.

Majority  $\rightarrow$  male  $\rightarrow$  30  $\searrow$  Dataset  
female  $\rightarrow$  30  $\swarrow$

- Disadvantage → we are throwing useful data.

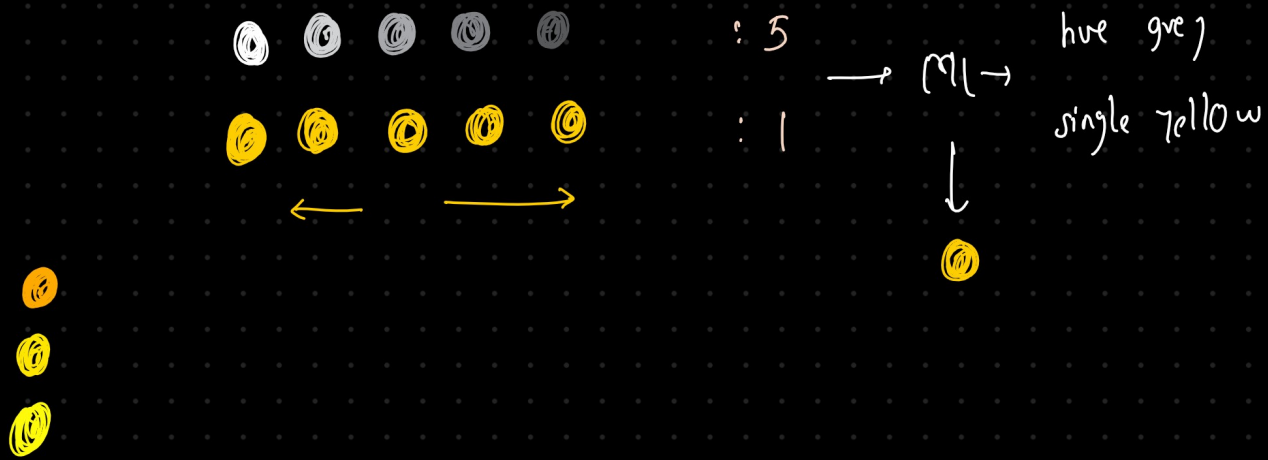
→ ROS

- Duplicate samples from minority class.

Male : 100

Female: 30  $\rightarrow$  ROS  $\rightarrow$  Female: 30 + 70  
 $\nearrow$   
by duplication

- can cause overfitting



→ SMOTE

→ instead of duplicating data, create synthetic data by interpolating between minority sample.

	Age	Sugar	
A →	50	200	
B →	60	220	
C →	55	210	→ within A & B

Better than ROS

→ works only with numerical data

→ ADASYN

→ similar to SMOTE, but smarter.

→ generate more synthetic samples near boundaries

→

low ← 70 - 100 - 140 → High  
 ← normal →

100 people  $\rightarrow$  70%  $\rightarrow$  70 - 140

15 → > 140

15  $\rightarrow$   $< 70$

→ help generate data of under represented categories

M : 1056

$$F = 30$$

```
def fun ( )
```

50 → smoke

50 - ADA 57N

# EDA - I

Dataset  $\rightarrow$  set of analysis  $\rightarrow$  pattern/relationship

```

graph LR
    data --> info
    info --> count_n[how many colms are numerical?]
    info --> count_c[categorical?]
    info --> dt_head[dt:head()]
    dt_head --> count_cp[categorical -> count plot -> t]
    data --> read[read about columns]

```

→ Analysis → visualization

↓

stability method

Data  $\rightarrow$  EDA  $\rightarrow$  Feature-Engineering  $\rightarrow$  ML  $\rightarrow$  Deployment

Data → NLP → ML  
↓

A → i don't like this topic → negative  
→ i like this topic → positive  
→ Hello → neutral  
→ topic is okay → neutral



Summarize and visualize what the data looks like

- Central tendency (mean, median, mode)
- Variability (variance, std, IQR)
- Correlation & Associations
- Frequency counts, histogram, boxplot, plots...
- handling missing values, outliers

ASK : What is happening in the dataset?

Make inferences or generalizations about the population from the sample

- Hypothesis Test
- Confidence interval
- Effect size
- Regression based inference

ASK : Can we claim this effect exist in the larger population or is it just by chance?

NLP → language

nltk → nlp task → english  
french

corpus - collection of text/data

stop words: a, an, the, she, it, on, at  
(common words  
which doesn't add any significant meaning)

n-grams:

sequence of n words that appear together in a text:

$n=1$ , Unigram: movie, good, bad, mona, Hello

$n=2$ , Bigram: great movie, not good, Hello Alok

$n=3$ , Trigram: waste of time, best movie ever

Unigram → not, good → not good  
↓  
+ve

Bigram: not good

↓  
-ve

Trigram: not good direction

useful