

WELCOMEAgenda:

- Summary - learnt
- Overview (look first overview)
- ML? vs SDE
- Tasks
- Learning
- Applications

① Python & DS Libraries3. ② Prob & Stats1. ③ Coordinate Geometry (Linear Algebra)2. ④ Calculus & Optimisation

Numpy, Pandas, Seaborn

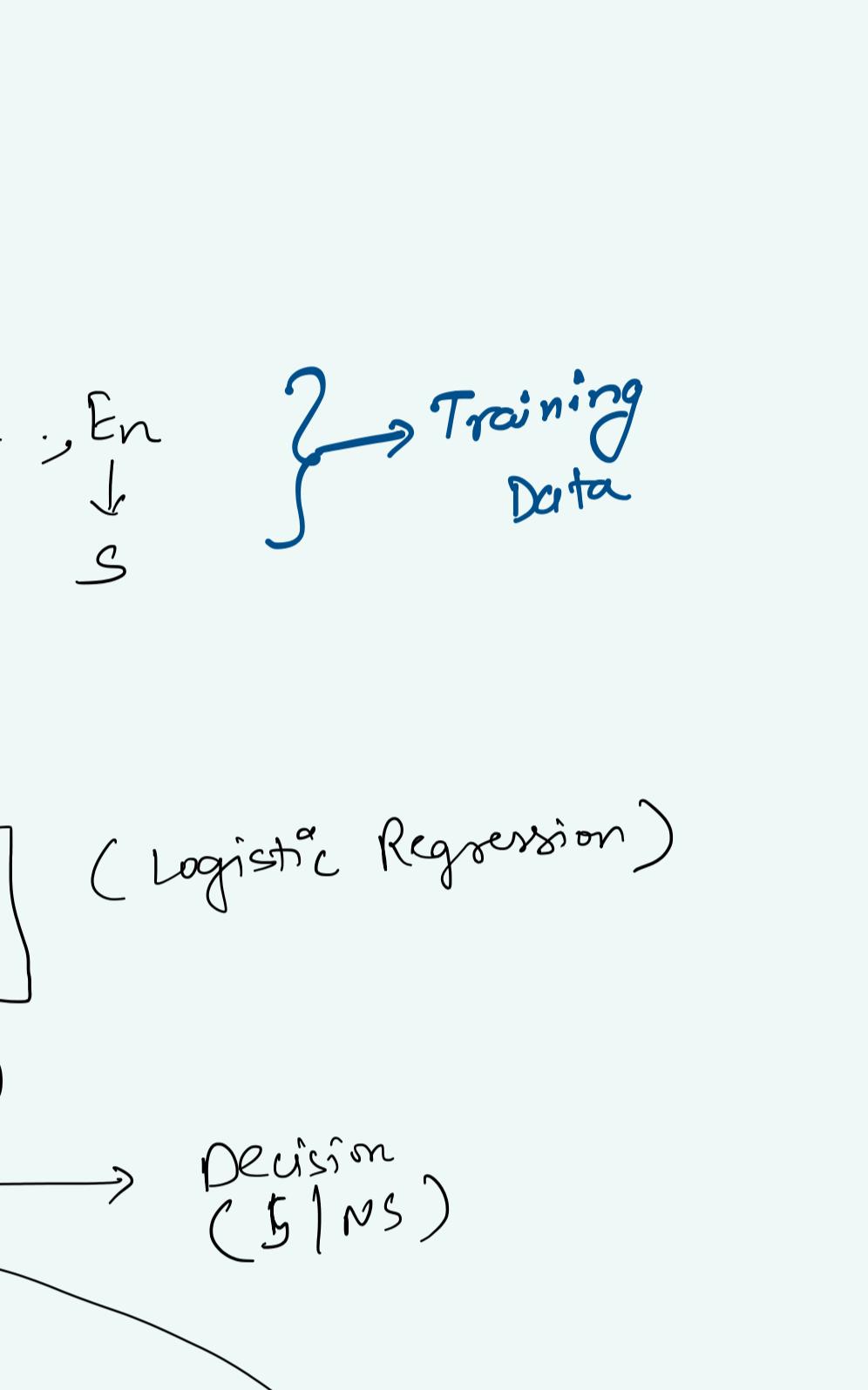
IV → Numerical : histogram, boxplot, kde, violinplot, distribution  
 Categorical : countplot, piechart

2V → N-N : scatterplot, hexbin, lineplot  
 N-C : boxplot/violin, boxplot,

C-C : heatmap, contingency table  
 Contingency table

Gender | Profession

→ M → N-N-C



Q: P(diabetic) ???

$$\frac{300}{1000} \Rightarrow 0.3$$

$$\left\{ P(\text{diabetic} \mid \text{Patient } w=120 \text{ kg} \& h=6') \Rightarrow 0.6 \right.$$

Patient<sub>i</sub> ⇒ w=120 kg & h=6'Case I: Patient<sub>i</sub> is diabeticCase II: Patient<sub>i</sub> has 90% probability of being diabetic

## # Linear Algebra:

Data →  $x^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}]$   
 (1-dim vector)  
 n-data points

ML Vs Classical Programming  
 Arthur Samuel → ML provides ability to learn w/o being explicitly programmed.

Example: Email → Spam or Not Spam

Data: Text-Email  
 Input: E<sub>1</sub>, E<sub>2</sub>, E<sub>3</sub>, ..., E<sub>n</sub>  
 Output: S, NS, NS, ..., S

Rules: lottery, 100 off / free, ASAP, immediately send money  
 if rule → Spam  
 else Non-Spam

Email → SDE Model (function) → Decision (S/NS)

(rules) (logic) → Decision (S/NS)

# Drawbacks:  
 → manual logic  
 → rigid rules → worst terrible

ML → New Email → ML Algo (function) → h (Hypothesis function) → Decision (S/NS) → Prediction

SDE: inputs → logic (pattern) → output  
 ML: Train data → h → output

"Pattern Recognition"

SDP: inputs → logic (pattern) → output  
 ML: Train data → h → output

SDP: inputs → logic (pattern) → output  
 ML: Train data → h → output

Tom M. Mitchell:

Task (T) → S/NS classification  
 Experience (E) → How well learned thing correctly?

ML Algo to train → ML Algo to train

Performance (P) (Metric):  
 Unseen Data / Test Data → Good Prediction.

100 New Email → 20 emails (S) → 10 emails (NS) →  $\frac{92}{100} = 92\%$

T: Predict Blood Glucose Level:

Continuous value

D:  $\{ (x^{(i)}, y^{(i)})_{i=1}^n ; x^{(i)} \in \mathbb{R}^d ; y^{(i)} \in \mathbb{R} \}$

Target output is real value.

ML:  $\{ (x^{(i)}, y^{(i)})_{i=1}^n ; x^{(i)} \in \mathbb{R}^d ; y^{(i)} \in \{0, 1\} \}$

Target is categorical

Customer Segmentation

f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>, ..., f<sub>d</sub>

Low Segment, Mid-Segment, High Segment

D:  $\{ (x^{(i)})_{i=1}^n \}$

f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>, ..., f<sub>d</sub>

clustering

Amazon (100M):  
 → Product?  
 → spent in 3 months  
 → location

clustering

Rec-Sys

Supervised / Unsupervised

Rec-Sys → item-item similarity

Collaborative filtering

Matrix factorization (SVD)

Rec-Sys → Matrix factorization (SVD)