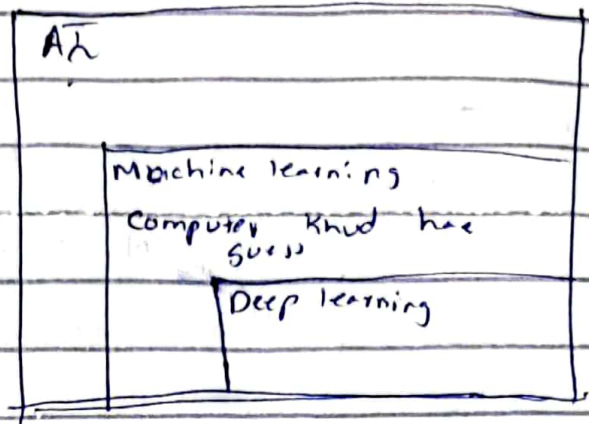


Session 2

Machine Learning

"Field of study that gives computer the ability to learn without being explicitly programmed"

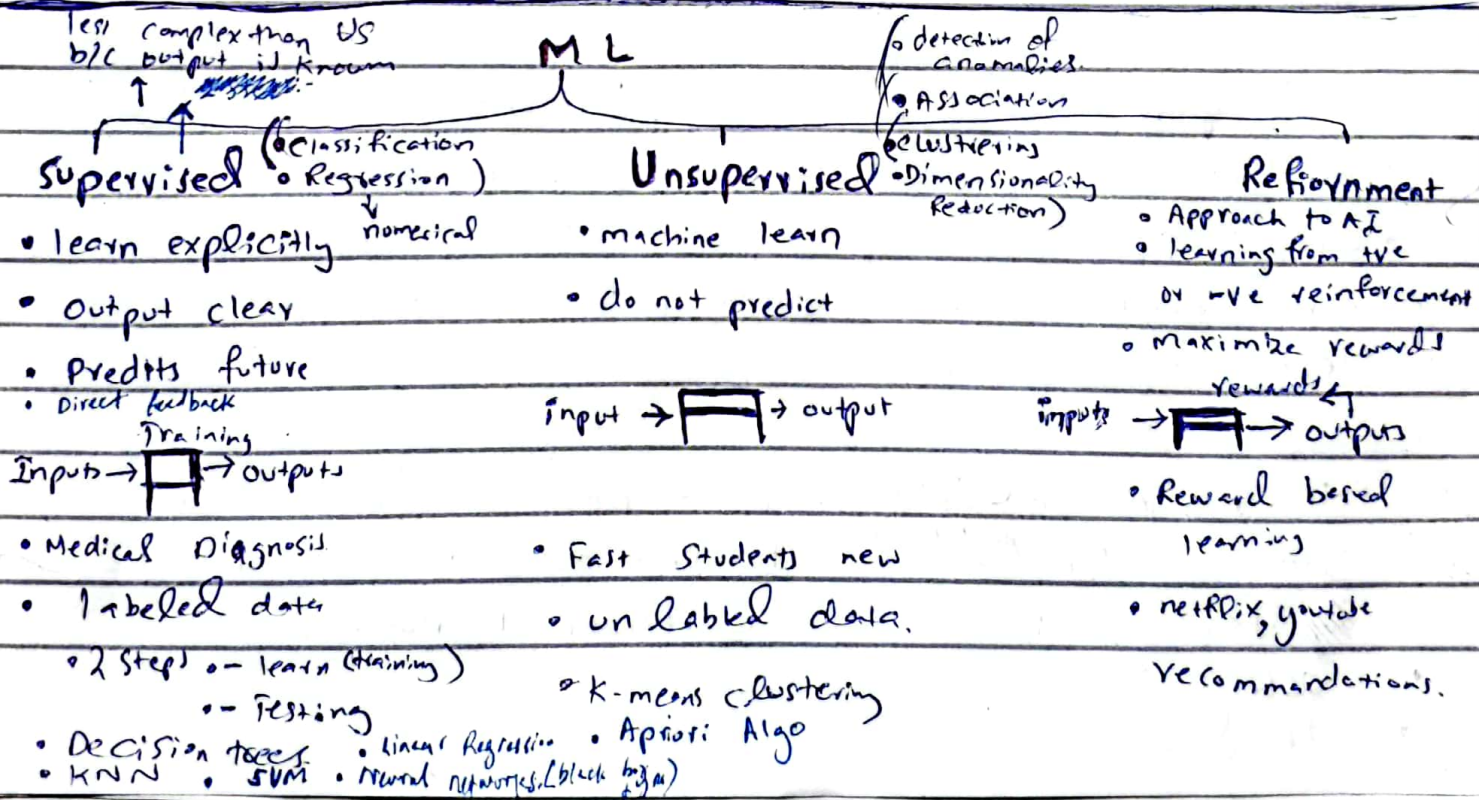
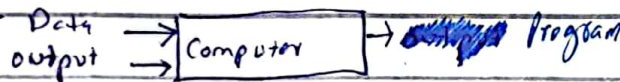


- Start with data $\rightarrow E$ (experience)
- perform some sort of task $\rightarrow T$ (Task) (Task is called supervised learning, classification or inductive learning)
- Validation test \rightarrow Performance measure P

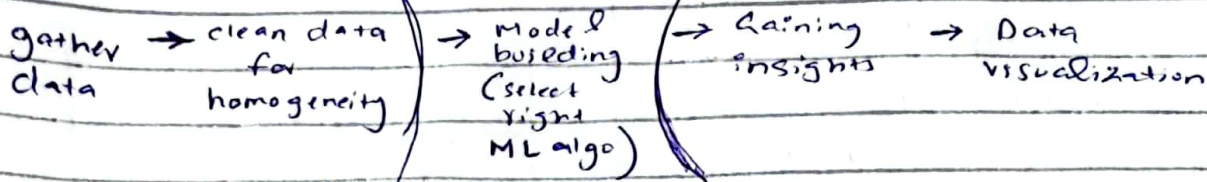
Traditional



ML



Machine learning Process



Accuracy: (classification problem)

$$\text{Accuracy} = \frac{\text{No. of correct classification}}{\text{Total no of test cases}}$$

Classification: (2 things)

- We predict y labels (classes) for input x.
- Yes or No (fraud or not fraud)

$$y = f(x) \rightarrow \text{input}$$

↓
output

↓
classification function

Regression: Numerical data (1 thing)

Supervised

- feedback
- used for prediction
- known number of classes
- less complex than →

Examples

- Text recognition. (training data)
- Spam detection
- ~~Hand~~ face detection
- OCR
- Fraud detection.

Unsupervised

- no feedback
- used for analysis
- unknown number of classes.

Examples

- Noise removal from datasets
- Recommendation engines
- Customer Behaviour prediction.

• - Different Metrics :-

Confusion Matrix

Accuracy

Precision

AUC (Area under the curve)

MAPE (Mean absolute % error)

MAE (Mean absolute error)

MSE (Mean squared error)

F1 Score

We can not apply unbalance data to accuracy.

Predicted \ Actual	Actual	
	TP	FP
Predicted	FN	TN

• - Accuracy :- $\frac{TP + TN}{Total}$

$\frac{TP}{TP + TN}$

• accuracy score

• - Precision :-

sb kay sb predicted cheezon ma

~~TP + FP~~

TP

TP + FP

$\frac{TP}{TP + FP}$

(How many did we watch)

• - Recall or Sensitivity :-

↓

If every person has cancer, recall is 100%.

$\frac{TP}{TP + FN}$

$\frac{TP}{TP + FN}$

(How many we miss)

• - If focus on minimize FN, Recall $\rightarrow 100\%$ without precision too bad

• - If on " FP, Precision $\rightarrow 100\%$.

-- Specificity :- (Opposite to recall)

$$\frac{TN}{TN+FP}$$

$$\frac{TN}{TN+FP}$$

-- F1 Score:-

Single score that represent Precision and recall.

F1 Score = Harmonic Mean (Precision, recall)

$$F1 \text{ Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

-- AUC:- (Scale - invariant) (0-1 range) (If 100% correct AUC is 1 no. of samples)

a) True-Positive rate (TPR)

$$TPR = \frac{TP}{TP+FN}$$

b) False-Positive rate (FPR)

$$FPR = \frac{FP}{FP+TN}$$

-- MAE:-

$$= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)$$

(Avg diff b/w original and predicted values)

-- MSE:-

$$= \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

-- MAPE:-

$$= \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$

Decision Tree

(if one attribute 0, answer is zero)

- Calculate entropy and gain of all the columns
- From these columns whose gain is more it means that impurity is less and we select that as a root node.
- Step 1:- Calculate entropy of whole dataset.

$$- \frac{C(P)}{C(T)} \times \log_2 \frac{C(P)}{C(T)} - \frac{C(N)}{C(T)} \times \log_2 \frac{C(N)}{C(T)}$$

Step 2:- Entropy of all attributes of that column.
→ Same formula.

Step 3:-
$$\text{gain} = \underset{\substack{\uparrow \\ \text{Entropy of attribute}}}{\text{Entropy(whole data)}} - \frac{\text{1st attribute total values} \times \text{Total}}{\text{Total}}$$

Step 4:- Maximum gain is Ka and Kb root node.

Root node \Rightarrow splitting
attribute.

Attribute Selection Measures (ASM) / splitting rules.

Popular. are \Rightarrow • Information gain (decrease in entropy)

• Gain ratio

• Gini index

Data imbalance \Rightarrow Red block] To man na samjhega the

Performance metrics

① Data Preparation \rightarrow Feature engineering (Features nikalta hai) \rightarrow data modeling \rightarrow Performance measure.

② Confusion matrix (used in classification) (like ek banday ko Cancer machine Yes. No, Pakein agr uske ulta ho jae, Yes ko No aur No ko Yes. Ab aisi table ko confusion matrix)

		Actual (Columns)		
		Positive	Negative	
Predicted (Rows)	Truly Positive	TP 2	FP 3	$\Rightarrow 5$
	Negative	FN 1	TN 3	$\Rightarrow 4$
		3	6	

(We have to overcome FP and FN according to the situations)
 Small cases minimize (email spam)
 Big cases minimize (Cancer)
 Falsly hesitive (Agr false lekha hwa hai tou agr word galat ho ga) \rightarrow galat N hai (e.g. FN)

Accuracy

TP	FP
FN	TN

Accuracy = $\frac{TP+TN}{TP+FP+FN+TN}$

(Accuracy achi nahi hai, data imbalance)

Precision

$$\frac{TP}{TP+FP}$$