

## Classification

## Supervised learning

It is the process of categorising / grouping the data based on predefined (labeled) data.

## Goal:

Is to develop a model that can effectively assign a label or category

## Phases:

1. First phase: A model is created based on the trained data
2. Second phase: applied the 1<sup>st</sup> phase to categories  
Ex: Training loan, input: income, age, loan, Output: can  
Method take 8 and 1 high & low

1. Specifying Boundaries [threshold value & Decision]
  - Executed by partitioning the potential decision tuple into distinct regions
  - Each region being linked to specific data.

2. Using probability Distribution



- Involves utilizing the probability density function (PDF)
- 3. Using posterior prob.  $P(c_2 | t_1)$   
 Denoted as  $P(t_1 | c_2)$   
 $t_1 \rightarrow$  likelihood of a specific point / attribute  
 $c_2 \rightarrow$  class  $c_2$

### ISSUES in classification

1. Overfitting
2. Underfitting
3. Noise in data
4. Feature Selection.
5. Scalability.

### Measuring Performance in classification

- Used to evaluate the accuracy of a classification algorithm in predicting the correct class for a given instance
- One of the approach is  
 \* Confusion matrix.

classified model: Accuracy, Recall, Precision, F1-score, AUC

Regression: RMSE,  $R^2$

Segmentation: IOU, Pixel Accuracy, Recall precision.

What data should be used to measure the performance

- As the model core tuned into 'training set' during the process, the output is predicted.
- Therefore, in order to estimate the 'general' error,  $\rightarrow$  the model required to eval through unseen data.
- $\therefore$  Use Test data.

### Confusion Matrix

- Used to describe the perform of a classified model.

Actual	Predicted Neg	Predicted Positive
Actual Negative	TN	FP
Actual Positive	FN	TP

1. True positive (TP): Instance that is correctly classified as belonging to the positive class
2. True Negative (TN): Instance that is correctly classified as not belonging to the positive class



3 False Negative:- Instance is incorrectly classified as not belonging to the positive class

4 False Positive:- Instance is incorrectly classified as belonging to positive class

1. Accuracy:- Probability of correct prediction.

$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN)$$

2. Precision:- Proportion of true positive among the instances that the model is predicted as positive.

$$\text{Precision} = \frac{TP}{(TP + FP)}$$

3. Recall:- 
$$\frac{TP}{(TP + FN)}$$

4. F1 Score:- 
$$2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$$

# Statistical Based algorithms

classmate

Date

Page

## classification models

1. Regression Analysis
2. Bayesian Classification

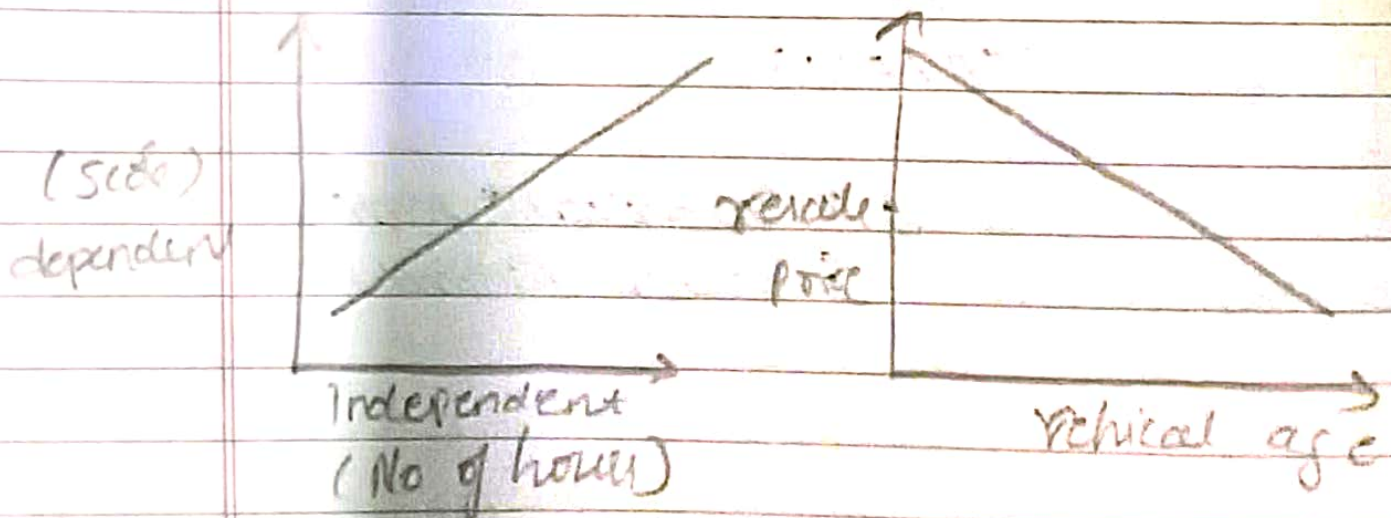
## Regression Analysis

Regression is a statistical method that helps us to understand & predict the relationship b/w variables.

variable  $\rightarrow$  terms  $\leftarrow$  Independent  
Dependent

- Describe how one variables varies if another changes.

Ex: Score based on hours, resale price based on vehical age



+ve linear regression

-ve linear regression



## 1. Data Mining Metrics

1. Return on Investment (ROI)
2. Business or Usefulness perspective
3. Space and Time complexity
4. Accuracy in classification
5. Specific Metrics for DM task
6. Effectiveness of the implementation

## 2. Statistical Perspective on DM

1. Point Estimation
2. Bias of an estimator (EBD)
3. Squared error
4. Mean square error (MSE)
5. Root mean square
6. Unbiased of estimator
- 7.
8. Confidence interval
9. Jackknife Estimate
10. Maximum likelihood Estimate (MLE)

3.

Metrics for set similarity in DM

classmate

Date  
Page

$$1. \text{ Cosine Similarity } \cos \theta = \frac{|A \cap B|}{|A| |B|}$$

$$2. \text{ Jaccard Sim} \\ J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{\sqrt{|A|} \sqrt{|B|}} \\ = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

3. Pearson correlation coefficient.

$$P_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

cov  $\Rightarrow$  covariance $\sigma \Rightarrow$  SD

d. Sorensen - Dice coefficient.

$$S(A, B) = \frac{2|A \cap B|}{|A| + |B|}$$

 $|A| \ \& \ |B| \longrightarrow$  Sizes of sets



## 4. Dissimilarity measures

↓ ↑ Dissimilarity Distance ↑ ↓

- 1 → complete dissimilarity
- 0 → "Similar"

### 1. Euclidean Distance

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$

### 2. Manhattan Distance

$$d(x, y) = \sum |x_i - y_i|$$

## 5. Measuring performance in classification

1. Accuracy =  $\frac{TP + TN}{(TP + TN + FP + FN)}$

2. Precision =  $\frac{TP}{(TP + FP)}$

3. Recall =  $\frac{TP}{(TP + FN)}$

4. F1 score =  $2 * \frac{(Pre \times Recall)}{(Pre + Recall)}$



## 6. Attribute Selection Measure (PSM)

### 1. Information Gain:

a. (al Entropy

b. (al average entropies

c. Gain

$$\text{Gain}(A) = \text{Entropy}(D) - \text{Entropy}(D_A)$$

### 2. Gain Ratio:

$$\text{Gain Ratio} = \frac{\text{Gain}(A)}{\text{Split Info}(D)}$$

$$\text{Split Info}(D) = - \sum \frac{|D_i|}{|D|} \log \frac{|D_i|}{|D|}$$

### 3. Gini Index:

$$\text{Gini}(D) = 1 - \sum p_i^2$$

$$\text{Gini}_A(D) = \frac{|D_1|}{|D|} \text{Gini}(D_1) + \frac{|D_2|}{|D|} \text{Gini}(D_2)$$

$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D)$$