ASSIGNMENT-II          V.Raja Rajeshwari

DSR :

1. What is Predictive modeling? Discuss about evaluation of Predictive models?

Ans Predictive modelling is the process of using known results to create, process and validate a model that can be used to forecast future outcome using statistics

=) It is a tool used in predictive Analytics, a data mining technique that is concerned with forecasting probabilities and kernels

=) examples of specific types of forecasting that benefit business and demand forecasting

headcount planning
clean Analysis
competition Analysis
finacial risks
Types of predictive modeling
* Regression
* clustering
* Neural networks

* classification
* Time series model
* forecasting

2. What is linear regression & list out the critical assumptions of linear regression?

Ans Linear regression is the supervised machine learning model in which the model finds the best fit linear line between the independent and dependent variables

Critical Assumptions

We can use to understand the relationship b/w two variables x and y

$$y = ax + b \quad ①$$

① Linear relationship

There exists a linear relationship b/w the independent variable x and the dependent variable y.

② Independence

③ Homoscedasticity

④ Normality.

③ Why logistic regression is used for the classification
Explain model building strategies for logistic regression

Ans logistic regression is used as a classification
technique. It uses logistic function to model
the dependent variable

⟹ The dependent variable is dischotomours
in nature

⟶ As a result, this technique is used while dealing
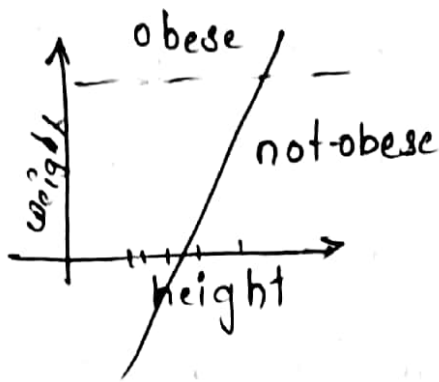with binary data.

Eg:- customers choosen
spam email / website
⟹ purposeful solution of variables includes the following

Steps:
1) univariable analysis
2) multivariable model comparison
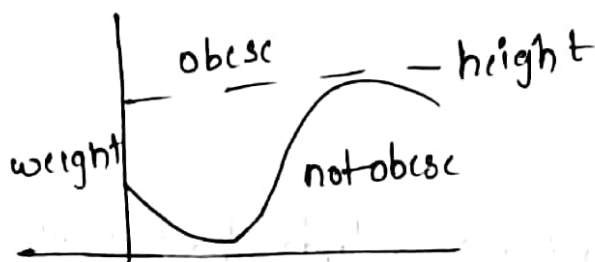3) linearity Assumption
4) Interactions Among covariate.

A) ⑤ Write in detail about linear and
logistic regression?

Ans linear regression is used to predict the
continuous dependent variable using a given
set of Independent variables

Logistic Regression is used to predict the categorical dependent variables using a given set of independent variables

$$\log\left(\frac{y}{1-y}\right) = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + \cdots + b_n X_n$$



5) List out the various control structures supported by R programming language?

Ans — control statements are expressions used to control the execution and flow of the program based on the condition provided in the statement

In R programming
* If-condition
* If-else
* for-loop
* nested-loop

6. Breifly describe the data structures in R-programming language?

Ans A datastructure is a particular way of organizing data in a computer so that it can be used effectividely
=) the most essential datastructure used in 'R' includes:

o Vectors
• Lists
• Dataframes
• Matrics
. Arrays
° Factors

\* Vectors

A vector is ordered collection of basic type of given length.

X = c(1,3,5,7,8)
print(n)
olp: [1] 1 3 5 7 8

Lists
A list is a generic object consisting of an ordered collection of objects
cid = c(1,2,3,4)

* while loop
* repeat and break statement
* Return statement
* next statement
* If-condition

Syntax
  if (expression)
  {
  statements;
  }

* for loop

Syntax
  for (value in vectors)
  {
  statements
  }

* nested loop
  m ← matrix(2:15,2)
      for(a in seq (nrow(m)))
          {
          for(c in seq (ncol(m)))
              {
              print (m[a,c])
              }}

[1] 2        (1) 8
(1) 4        [1] 10
[1] 6        [1] 12

7. Define object. List the methods for measuring Distance b/w objects?

Ans **objects**

object are the instance of the class. Also everything and to known more look at datatypes in R.

→ method for measuring distance b/w objects

**Euclideandistance**

The most common distance is euclidean distance

$$Edist(x,y) \leftarrow sqrt((x[1]-y[1])^2 + (x[2]-y[2])^2 + \cdots).$$

* **Manhattandistance**

Manhattan distance measure distance in the no. of horizontal and vertical units it takes to get from one point to the other

$$mdist \leftarrow sum(abs(x[1]-y[1]) + abs(x[2]-y[2]) + \cdots)$$

* **Cosine similarity**

cosine similarity is a common similarity metric in text analysis

$$dot(x,y) \leftarrow sum(x[1]*y[1] + x[2]*y[2] + \cdots)$$
$$cosine(x,y) \leftarrow dot(x,y)/sqr(dot(x,x) \text{ or } dot(y,y))$$
$$cos(\theta) = \frac{A \cdot B}{|A||B|}$$

## Dataframes

Dataframes are generic data objects of R which are used to store tabular form

=) Dataframes are the foremost popular data objects in R programming

=) Each item in a single column must be of same datatype

```
Name = c("Sweety", "Madhu", "Ramani")
lang = c("R", "python", "C++)
df = data.frame(Name, lang)

print(df)
    Name      lang
  Sweety      R
  Madhu       python
  Ramani      C++
```

## Matrices

A matrix as a rectangular arrangement of number in rows and columns

```
A = matrix(c(1,2,3,4,5,6,7,8,9),
        nrow =3, ncol= 3, byrow=TRUE)
print(A)
       [,1]  [,2]  [,3]
  [1,]    1    2    3
  [2,]    4    5    6
  [3,]    7    8    9
```

8  Define list and dataframe in R and explain various operations on lists and dataframes with suitable example.

Ans  A list is generated using list() function.

## Creating lists

Geek-list←list("Geek", TRUE, 27)
print (Geek-list)

* Naming the elements of a list
Accessing elements of a list!
Adding, deleting and updating elements of a list
Merging elements of a list
converting a list to vector

## Dataframes

Dataframes are generic data objects of R which are used to store the tabular data

creating a dataframe using vector
Accessing rows and columns
selecting of the subset of the dataframe
Editing dataframes

5. Explain k-Nearest Neighbours Algorithm and its implementation in R programming languages

Ans. K-Nearest Neighbors (KNN)

KNN can be defined as a supervised machine learning algorithm that can classify new data point into target class based on neighbouring data points features

For example: consider a KNN algorithm for a machine that differentiates between apples and mangoes.

To perform this a dataset of apples and mangoes are must be provided as input and then model must be trained in such a way that it detects that fruit based on certain characteristics.

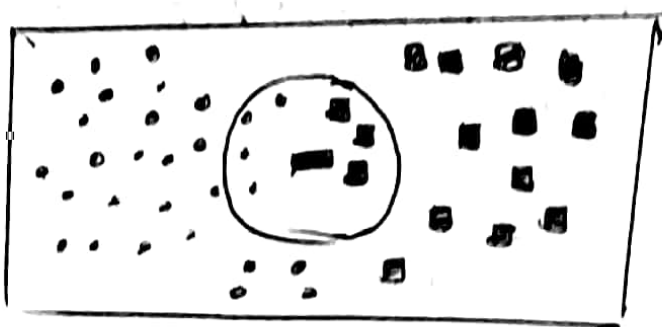For examples, features of an apple would be red in colour, round in shape etc

similarly features of mango would be yellow in colour, oval in shape etc

An image is provided to the model as input

measures such as eucledian and mahantann distance are used to classify the closeness b/w data points

When k=3 the nearest neighbours are 2 squares and 1 one dot. When k=3 it will assigned to class B

But ∩ k=7 the neighbours would be A dots and 8 squares



When it is classified the new data point will be assigned to class A.

## Implementation of KNN in R programming language.

Consider a problem to analyze the bank data and then build a machine learning model that can predict whether loan for applicant can be approved or rejected depending on its socio-economic profile

The bank credit dataset consists includes the information about 1000's of applicant's such as name, age, account balance, loan records etc.

Based upon this data, prediction can be made that whether loan can be approved or rejected. This

Problem can be solved by KNN algorithm by classifying loan apprequest into two classes

1. Approved
2. Rejected.

steps to solve this problem as follows -

Step1: Importing the dataset -

Import DATA the dataset of applicants.

loan ← read csv(" c:/usas/desktop/DATASET/knn/credits-csv")

The structure of data would be

Str (loan)

'dataframe': 1000 obs. of 21 variables:-

$ Credibility: int 1 1 1 1 1 1 1 1 1--:

$ Account Balance: int 1 2 11 11 4 2

$ Duration of credits month: int 18 9 12 12 10 8 6 18 24

$ purpose int 2 0 9 0 0 0 1 2 3

$ credit amount int 1049 2799 841 2122 2171

& values savings stocks : int 1 1 2 1 1 1 1 3---

& length of current employment: int 2 3 4 3 3 2 4 2 1 1

$ instalment percent : int 4 2 2 3 4 1 2 4 1

$ Marial status : int 4 2 3 2 3 2 2 1 1 1

$ Guarantors int 1 1 1 1 1 1 1 1 --

$ Duration in current Address: int 4 2 4 2 4 3 4 4 4--

$ Age years: int 21 36 23 39 38 40

$ No of dependices : int 1 2 1 1 2 1 2 1 1

$ Telephone : int 1 1 1 1 1 1 - - -

## Step: Data cleaning

The structure of dataset consists of predictable variables that are used to decide whether loan for applicant can be approved or rejected.

Some of the variables might not be useful in predicting the loan. For example - credit concurrent credits telephones etc. can be removed because such type of variable leads to complexity of machine learning model

loan-subset ← loan [c(' credibility', 'Age-years', 'Account balance')]]

The dataset now would be

str (loan-subset)

'dataframe' 1000 obs of 8 variables

```
$ credibility            int 1 1 1 1 1 1 ...
$ age - years            int 29 30 49 ...
                         int 2 3 71 11
$ Marital status         int 3929 4232 - -
$ occupation
$ Account balane         int 1 2 2 1 7 1 4 21
$ length of current employment   int 23 4 3 3 2 44 1
$ purpose                int 20 9 0 0 0 0 33 - -
```

# Steps: Data Normalization.

Normalization of data set is mandatory to make the output unbiased consider below observations

head (loan subset)

| | credibility | ag.yeans | Maritalstatus | occupation | Account Balance |
|---|---|---|---|---|---|
| | 1 | | 4 | 1 | 1 |
| 2 | 1 | 49 | 3 | 1 | 2 |
| 3 | 1 | 63 | 2 | 2 | 4 |
| 4 | 1 | 78 | 1 | 1 | 3 |
| 5 | 1 | 58 | 1 | 1 | 2 |
| 6 | 1 | 15 | | | |

| credit amount | length of employments | Purpose |
|---|---|---|
| | | 2 |
| 2943 | 2 | 0 |
| 5678 | 3 | 9 |
| 9629 | 4 | 0 |
| 1423 | 3 | 6 |
| 2428 | 3 | 0 |
| 1928 | 2 | 6 |

Here credit amount variable has value dcale in 1000s and the remaining are in dingle of two digitits. If the data is not normalized outcome is in baised form

normalize ← function(x)

$$return \ (x - min(x)) / (max(x) - min(x))$$

Now store the normalized data set in loan subset. and then remove crediability variable because it is

# Explain K-Means Algorithm and its implementation in R programming

## K-means

The k-means is an iterative clustering algorithm in which objects are moved among set of clusters until desired set is achieved.

It is most popular and commonly used method. the algorithm is based on the concept of most specified input parameter $k$). A set of $n$ objects are divided into $k$ clusters. A high degree of similarity among objects. elements in clusters is obtained.

## Working of k-means Algorithm

### Algorithm's input
the number of desired clusters are denoted by $k$. a dataset containing $n$ objects denoted by $D$.

### Algorithm's output
$k$: A set consisting of $k$ clusters

### Procedure:

Step1: Initially select $k$ objects randomly from $D$, as initial cluster centers

Step2: Depending upon the distance between the object and the cluster mean, each remaining object assigned to cluster to which it is most similar.

euclidean distance b/w mean and objects to
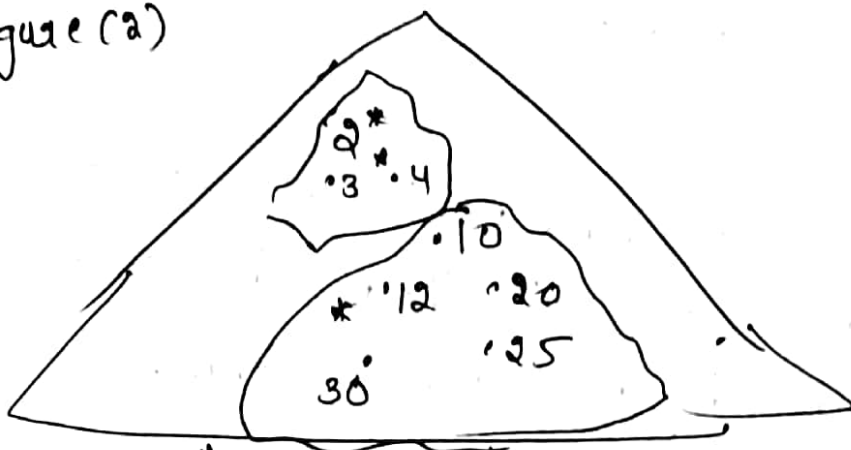classify objects into two clusters



Fig(1) initial partitioning with $m_1 = 2$ and $m_2 = 4$

Now! evaluate new mean for the resulting clusters
and again partition is done base of euclidean
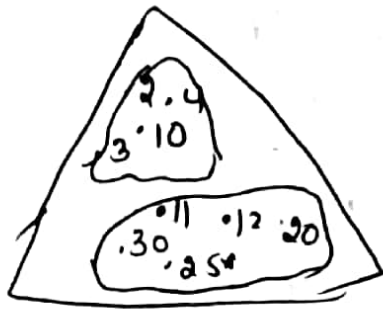distance as follows

$$m_1 = \frac{2+3}{2} = 2.5$$

$$m_2 = \frac{4+10+11+12+20+25+30}{7} = 16$$

$m_1 = 2.5$ and $m_2 = 16$ and resulting cluster is shown
in figure (2)



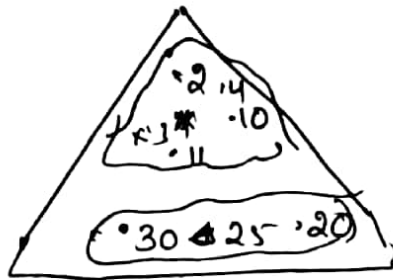Fig(2) New clusters formed with $m_1 = 2.5$ and $m_2 = 16$.

The process iterates for the successive values of the mean
to enhance the partition of clusters

clusters = $m_1 = 3$  $m_2 = 18$



clusters with $m_1 = 4.75$ and $m_2 = 19.5$



clusters with $m_1 = 7$ and $m_2 = 25$

Thus finally the correct partition of the numered objects into two clusters is achieved.

## K-Means Implementation in R

The K-means algorithm can be implemented by predefining. Here K represents number of clusters to be defined.

steps to implement K-means are as follows.

1. In the first step K centres are defined and every cluster is assigned to the cluster that has closest centre to it.

2. In the second step the centres are redefined by using the observation of each cluster. The column means are used for defining the centroid

The above two steps are repeated until the centres are converge therefore K-means algorithm said to be iterative.

List out the various performance metrics for classification?

Performance Measures

Performance of a classification model is evaluated once the predictive model is build. It determines whether the model can predict the outcome of new observation test data. that is not used in the training the model. The performance of a model can be accessed by comparing the predicted outcome values against known outcome values.

The commonly used metrics and methods to access the performance of predictive classification models are as follows.

1. Average classification Accuracy

It represents the proportion of correctly classified observation. The complete classification rate corresponds to the part of observations that are classified correctly. The first step of accessing the performance of a model is to determine the raw classification accuracy.

The raw classification error rate can be inversly defined as proportion of observation that are misclassified.

$$error\ rate = 1 - accuracy.$$

The raw classification accuracy and error can be determined by comparing the observed classes in test data and predicted classes by model

7accuracy ← mean(observed class == predicted class)

7accuracy

[1] 0·08

7error ← mean(observed class! == predicted class)
7error

[1] 0·192

The binary classifier in this example makes two types of error

=7 It can wrongly class assign invidiudal who has diabetes positive to diabetes negative

=7 It can wrongly assign individual who has diabetics negative to diabetics positive

2. Confusion matrix

The main factor required for performance evaluation of meta classification models includes correct or incorrect predictions of no. of test cases by the model. The actual test records are given in a tabular form known as confusion matrix.

| Actual Class | Predicted class | |
|---|---|---|
| | classY | classX |
| classy | fyy | fyx |
| classX | fxy | fxx |

The above specified table outlines the binary classification problem and each individual entry is given as $fxy$ Here $x$ represents the no of records belongs to class $x$ and $y$ represents no of records belongs to class $y$

$fxx$
This entry specifies the no of records belonging to class 'x' There is no mis-prediction in this entry, since the records of class 'x' are predicted as records from class 'x'

$fxy$
this entry specifies the no of records belonging to class 'x' There is mis prediction in this entry since the records of class are incorrectly predict as records from class 'y'.

$fyy$
This entry specifies the no of records belonging to class 'y'. There is no mis prediction in this entry

## A) ROC curve

The graphical representation that exhibits the performance of predicted classification model is known as receiver operating characteristics curve.

## Applications of ROC curve

ROC curves depicts the tradeoff blus porportion of correctly determined positive tuples known as true positive rate and the porportion of negative tuples incorrectly tuples as positive tuples known as negative rate for given model