Name : SHEHARYAR TAHIR
ID#      13484

**Q1**: **Distinguish between Classification and Regression with the help of relevant scenarios.**

**Answer:**

Classification and Regression are two major prediction problems which are usually dealt in Data mining. Predictive modeling is the technique of developing a model or function using the historic data to predict the new data. The significant difference between Classification and Regression is that classification maps the input data object to some discrete labels. On the other hand, regression maps the input data object to the continuous real values.

Let's take an **example in classification**, suppose we want to predict the possibility of the rain in some regions on the basis of some parameters. Then there would be two labels rain and no rain under which different regions can be classified.

Let's take the similar **example in regression** also, where we are finding the possibility of rain in some particular regions with the help of some parameters. In this case, there is a probability associated with the rain. Here we are not classifying the regions within rain and no rain labels instead we are classifying them with their associated probability.

## Comparison Chart

| COMPARING BASE | CLASSIFICATION | REGRESSION |
|---|---|---|
| Basic | The discovery of model or functions where the mapping of objects is done into predefined classes. | A devised model in which the mapping of objects is done into values. |
| Predicts | Discrete values | Continuous values |
| Algorithms used | Decision tree, logistic regression, etc. | Regression tree (Random forest), Linear regression, etc. |
| Nature of data | Unordered | Ordered |
| Calculation method | Measuring accuracy | Measurement of root mean square error |

**Q2: Perform Naïve Bayes or Decision tree classification for new instance where (SSN = 123-46-4455, Test1= 85, Test2= 31 and Final= 30) Find Grade.**

| "SSN" | "Test1" | "Test2" | "Final" | "Grade" |
|---|---|---|---|---|
| "123-45-6789" | 100 | 83 | 49 | "D" |
| "123-12-1234" | 96 | 97 | 48 | "D" |
| "567-89-0123" | 60 | 40 | 44 | "C" |

| | | | | |
|---|---|---|---|---|
| "087-65-4321" | 36 | 45 | 47 | "B-" |
| "456-78-9012" | 88 | 77 | 45 | "A-" |
| "234-56-7890" | 80 | 90 | 46 | "C-" |
| "345-67-8901" | -1 | 4 | 43 | "F" |
| "632-79-9939" | 30 | 40 | 50 | "B+" |

**Answer:**

| SSN | Test1 | Test2 | Final | Grade |
|---|---|---|---|---|
| X = 123-46-4455 | 85 | 83 | 30 | ? |

**\* Problem Statement :**
- Given Features $X_1, X_2, ---, X_n$
- Predict a Label Y

**\* Bayes Classifier :**
- A probablistic framework for solving classification problems.
- Conditional Probability

$$P(C|A) = \frac{P(A,C)}{P(A)}$$

$$P(A|C) = \frac{P(A,C)}{P(C)}$$

**\* Bayes Theorem**

$$P(C|A) = \frac{P(A|C) \cdot P(C)}{P(A)}$$

In general compute the posterian probability $P(C|A_1, A$
for all values of C using the Baye's Theorem. Choose th
value of C that maximizes $P(C|A_1, A_2 --- A_n)$

$$P(c) = \frac{Nom}{Total\ Sample}$$

* Class

| | | | | |
|---|---|---|---|---|
| $P(A-) = 1$ | | $P(A-) = 1/8$ | $\gg$ | $\infty$ |
| $P(B+) = 1$ | | $P(B+) = 1/8$ | $\gg$ | $\infty$ |
| $P(B-) = 1$ | | $P(B-) = 1/8$ | $\gg$ | $\infty$ |
| $P(C-) = 1$ | | $P(C-) = 1/8$ | $\gg$ | $\infty$ |
| $P(C) = 1$ | | $P(C) = 1/8$ | $\gg$ | $\infty$ |
| $P(D) = 2$ | | $P(D) = 2/8$ | $\gg$ | |
| $P(F) = 1$ | | $P(F) = 1/8$ | $\gg$ | $\infty$ |

* For Discreate Attributes

$$P(A_i | C_k) = |A_{ik}| / N_{kc}$$

where, $|A_{ik}|$ is the numbers of instances having attributes $A$ and belongs to Class $C_k$

e.g ~~P(Status = married |No) = 4/7~~ } google example

~~P(Refund = Yes|Yes) = 0/3 = 0~~

① $P(X | Class = A-)$

$= P(Test1 = 85 | class = A-) *$

$P(Test2 = 83 | Class = A-) *$

$P(final = 30 | Class = A-)$

* For Non-Discreate values we have to use Normal Distribu

$$P(A_i | C_k) = \frac{1}{\sqrt{2\pi \sigma_{ij}^2}} e^{\left( \frac{A_i - \mu_{ij}}{2\sigma_{ij}^2} \right)}$$

where $\mu_{ij}$ = mean

$\sigma_{ij}^2$ = Varience

e.g $X = (5, 8, 7, 6, 9)$ mean

$$\sigma^2 = \frac{\sum (X_i - \bar{X})^2}{n - 1} \quad \text{no of sample}$$

where $\bar{X} = \dfrac{5 + 8 + 7 + 6 + 9}{5} = 7$

$$s^2 = \frac{(5-7)^2 + (8-7) + (7-7)^2 + (6-7)^2 + (9-7)^2}{5-1} \rightarrow P(T$$

$$s^2 = 3_{1}$$

① $\rightarrow P(\text{Test 1} = 85 \mid \text{Class} = A-)$

$$= \frac{1}{\sqrt{2\pi\sigma_{ij}}} \, e^{\left(\frac{A_i - \mu_{ij}}{2\sigma_{ij}^2}\right)} \quad \text{—— Ⓐ}$$

$A_1 = 85$ —①

$\mu_{ij} = \frac{85}{1} = 88$ —②

$\sigma_{ij}^2 = \frac{88-88}{1-1} = 0$ —③

Put in Ⓐ

$$= \frac{1}{\sqrt{2\pi(0)}} \times e^{\left(\frac{85-88}{2(0)}\right)}$$

$$= \frac{1}{0} \times e^{\left(\frac{-3}{0}\right)}$$

$$= \infty$$

$\rightarrow P(\text{Test 2} = 83 \mid \text{Class } A-)$

$A_i = 83$ —①

$\mu_{ij} = \frac{77}{1} = 77$ —②

$\sigma_{ij}^2 = \frac{77-77}{1-1} = 0$ —③

Put in Ⓐ $= \frac{1}{\sqrt{2\pi(0)^2}} \times e^{\left(\frac{83-77}{2(0)}\right)}$

$$= \frac{1}{0} \times e^{\left(\frac{6}{0}\right)}$$

$$= \infty$$

$\Sigma^2 + (\overline{q7})^2 \rightarrow P\left(\dfrac{Final}{Test4} = 30 \mid Class = A-\right)$

$A_f = 30$

$M_{ij} = \dfrac{45}{1} = 45$

$\partial_{ij}^2 = \dfrac{45 - 45}{L-1} = 0$

Put in (A)

$= \dfrac{1}{\sqrt{2\pi(0)}} \times e\left(\dfrac{30 - 45}{2(0)}\right)$

$= \infty$

② $P(Test1 = 85 \mid Class = B+)$

$= \dfrac{1}{2\pi\, c_{ij}^2} \times e\left(\dfrac{A_i - M_{ij}}{2\, c_{ij}^2}\right) \quad\longrightarrow \text{(A)}$

$A_i = 85$

$M_{ij} = 30/1 = 30$

$\partial_{ij}^2 = 30 - 30/_{1-1} = 0$

Put in A get $\infty$

$\Rightarrow P(Test\,2 = 88 \mid Class\,B+)$

$= \dfrac{1}{2\pi\, c_{ij}^2} \times e\left(\dfrac{A_i - M_{ij}}{2\, c_{ij}^2}\right)$

$A_i = 88$

$M_{ij} = 40$

$\partial_{ij}^2 = \quad = 0$

Put in A get $\infty$

$\Rightarrow P(Final = 80 \mid Class\,B+)$

$A_i = 80$

$M_{ij} = 80$

$\partial_{ij}^2 = \quad 0$

Put in A $\infty$

③ → $P(\text{Test 1} = 85 \mid \text{Class} = D)$

$$= \frac{1}{\sqrt{2\pi \, 6_{ij}^2}} \, e\left(\frac{A_i - M_{ij}}{2 6_{ij}^2}\right)$$

$A_i = 85$

$M_{ij} = \dfrac{100 + 96}{2} = \dfrac{196}{2} = 98$

$6_{ij}^2 = \dfrac{(100-98)^2 + (96-98)^2}{2-1} = \dfrac{(2)^2 + (-2)^2}{1} = \dfrac{4+4}{1} = 8$

Put in Ⓐ

$$= \frac{1}{\sqrt{2\pi(8)}} \, e\left(\frac{85-98}{2(8)}\right)$$

$$= \frac{1}{\sqrt{50}} \, e\left(\frac{-13}{16}\right)$$

$$= \frac{1}{7} \, e(-0.8)$$

$$= 0.1 \times 0.4$$

$$= 0.4$$

$P(\text{Test 2} = 83 \mid \text{Class N})$

$$= \frac{1}{\sqrt{2\pi \, 6_{ij}^2}} \, e\left(\frac{A_i - M_{ij}}{2 \, 6_{ij}^2}\right)$$

$A_i = 83$

$M_{ij} = \dfrac{83 + 97}{2} = \dfrac{180}{2} = 90$

$6_{ij}^2 = \dfrac{(83-90)^2 + (97-90)^2}{2-1} = \dfrac{49 + 49}{1} = 98$

Put in Ⓐ $= \dfrac{1}{2\pi(98)} \, e\left(\dfrac{83-90}{2(98)}\right)$

$$= \frac{1}{615} \, e\left(\frac{-7}{196}\right)$$

$$= 0.0016 \times e(-0.03) = 0.0043$$

$$P(\text{final} = 30 \mid \text{Class} = D)$$

$$\mu_i = 30$$

$$\mu_{ij} = \frac{49 + 48}{2} = \frac{97}{2} = 48.5$$

$$\partial_{ij}^2 = \frac{(49 - 48.5)^2 + (48 - 48.5)^2}{2 - 1} = \frac{(0.5)^2 + (1.05)^2}{1} = 0.5$$

Put in Ⓐ

$$= \frac{1}{\sqrt{2\pi(0.5)}} \times e\left(\frac{30 - 48.5}{2(0.5)}\right)$$

$$= \frac{1}{\sqrt{3.14}} \times e\left(\frac{-18.5}{1}\right) = \frac{1}{1.77} \times -18.7$$

$$= -8.87$$

---

✗ $$P(X \mid \text{class} = D)$$

$$= P(\text{Test1} = 85 \mid \text{class} ) D \times$$
$$P(\text{Test 2} = 83 \mid \text{class}) D \times$$
$$P(\text{Final} = 30 \mid \text{class}) P$$

$$= 0.4 \times 0.0043 \times -8.87$$

$$= \pm 0.01$$

So Sample X have Grades D

---

**Q3: Find a Dataset related to any field and perform several classification techniques (Naïve Bayes, Decision tree, SVM, or any) to predict a class of a new**

**instance using WEKA. Compare the results (Accuracy, Precision, Recall, MARE, MMRE) of classification algorithms in a Table.**

**Take snapshots of the all the steps you perform for the classification.**

**Answer:**

**Dataset used:** Weka default dataset named diabetes

**Classification techniques applied:** Simple Logistic, Naïve Bayesian, Random Forest and OneR



- OneR (71%) is significantly worse than Random Forest (76%)
- OneR (71%) is significantly worse than Naïve Bayesian (75%)
- OneR (71%) is significantly worse than Simple Logistic (77%)
- Shows that Simple Logistic (77%) performs significantly better for the particular data set.

1. Below is the snap shot from WEKA explorer taking the diabetes data set showing its class with 2 attributes **1)** Tested Positive **2)** Tested Negative

2. Below is the snap shot from Weka explorer by applying classifier Simple Logistic:

3. Below is the snap shot from Weka explorer by applying classifier Naïve
   Bayesian:

4. Below is the snap shot from Weka explorer by applying classifier Random Forest

5. Below is the snap shot from Weka explorer by applying classifier OneR

## COMPARISON TABLE

| Applied algorithms | Precision | Recall | F measure | MARE |
|---|---|---|---|---|
| Naive Baysian | 0.759 | 0.763 | 0.760 | 62.5028% |
| Random Forest | 0.754 | 0.758 | 0.755 | 68.3406% |
| Simple Logistic | 0.770 | 0.775 | 0.766 | 69.84% |
| OneR | 0.703 | 0.715 | 0.699 | 62.7298% |