



Write python code , consider filename as "housing.csv" in same folder

i) To load .csv file into the data frame

```
import pandas as pd
df = pd.read_csv("housing.csv")
```

ii) To display information of all columns

```
print(df)
```

iii) To display statistical information of all numerical

```
print(df.describe())
```

iv) To display the count of unique labels for "Ocean Proximity" column

```
print(df["Ocean Proximity"].value_counts())
```

v) To display which attributes in dataset have missing values count greater than zero

```
missing_values = df.isnull().sum()
```

```
missing_columns = missing_values[missing_values]
```

```
print(missing_columns)
```

For both datasets Diabetes and adult income

- which column in dataset had missing values
how did you handle them.

Dropped the missing values with the help of dropna since the number of rows with missing values were less compared

- what were the categorical columns
and how did you encode it?

diabetes.csv : ['Gender', 'CLASS']

adult.csv : ['workclass', 'education', 'marital status', 'occupation', 'relationship', 'race', 'gender', 'native-country', 'income']

and we used LabelEncoder() to encode it

- what is difference Min max scaling and standardization? when we use one over the other

~~Both Min-Max scaling and standardization are feature scaling techniques~~

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}} \quad X' = \frac{X - \mu}{\sigma}$$

we use min max when we do not have significant outliers
and we use standardization when we have outliers



Demonstrate steps to build ml model that predicts median house price

1. perform describe and info
describe the dataset
2. Plot histogram of each feature
median-income is right skewed
median-house-value is right skewed
3. Observe difference between test set random and stratified.
we take random record from the df
in stratified we target a particular feature
4. What does housing prices w.r.t to location show
areas closer to the coast have more prices
inland areas have lower prices
this is due to high population
5. Find most correlated column
Ocean proximity is correlated to numerical values using - get - columns then found correlation
6. List features that could be combined to improve correlation
rooms_per_household, bedrooms_per_room and population_per_household.

Lab - 3

①

x_i	y_i
1	2
2	4
3	5
4	9

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

$$y = 5$$

$$\bar{x} = 2.5, \text{slope} = \frac{45 + 0.5 + 0 + 6}{2 \cdot 2.5 + 0.25 + 0.25 + 2.25}$$

$$= \frac{11}{8} = 1.375$$

$$b = \bar{y} - m\bar{x}$$

$$= 5 - (2.2 \times 2.5)$$

$$= 5 - 5.5 = -0.5$$

$$y = 2.2x - 0.5$$

②

$$X = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 4 \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$XB = ((X^T X)^{-1} X^T) Y$$

$$= \begin{bmatrix} 1.0 & -0.5 \\ -0.5 & 0.2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$= \begin{bmatrix} 1.0 & 0 & -0.5 \\ 0.3 & -0.1 & 0.1 & 0.3 \end{bmatrix}$$

$$y = -0.5 + 2.2x$$

$$= \begin{bmatrix} -0.5 \\ 2.2 \end{bmatrix}$$

1 considering carda-per-capita-income.csv, salary.csv did you perform data preporcessing

Yes missing values were calculated by using simple imputer and label encoder

2 did you visualize regression line for the csv file

Yes, the regression line was plotted. The plot shows a strong linear relationship between year and per capita income meaning that as the year increases.

3 The predicted salary is dependent on the script depends on the fitted model's coefficient.

4 Did you encode categorical variables for 1000-companies.csv

Yes

Yes, it was done. Column was encoding using Label Encoder()

1. Binary classification with logistic regression.

a) $P(y=1|x)$ = $\frac{1}{1+e^{-(a_0+a_1x)}}$

b) $P(y=1|x=2) = \frac{1}{1+e^{(-0.5+0.8(2))}} = \frac{1}{1+e^{-0.6}} = 0.645$

c) predicted class : it is pass since it is above 0.5

2 applying softmax

$$\sigma(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

$$P_1 = \frac{e^2}{\sum e^i} = \frac{e^2}{11.107} = 0.665$$

$$P_2 = \frac{e^1}{\sum e^i} = \frac{e^1}{11.107} = 0.245 \quad P_3 = \frac{e^0}{\sum e^i} = \frac{e^0}{11.107} = 0.090$$

1. For iris.csv data set.

The accuracy was 98% indicating it classifies efficiently.

The Confusion matrix shows no misclassification.

2. for petrol-consumption.csv dataset

Regression tree predicts based on features.

It predicts for continuous values.

$$MAE = 90.60$$

$$MSE = 16851.60$$

$$RMSE = 129.81$$

\rightarrow NO = 4 time no prediction output will be
Yes = 1 time.

$$HCS = - \left(\frac{4}{5} \log_2 \frac{1}{4} + \frac{1}{5} \log_2 \frac{1}{5} \right) = 0.72$$

$$H(Hot) = - \left(\frac{3}{4} \times (-0.42) + \frac{1}{4} \times (-200) \right)$$

$$H(Cool) = 0 \quad 0.81$$

it branches at a3.

Lab-8

Person	Age	Salary	Target	Distance
A	18	50	N	52.79
B	23	55	N	46.96
C	24	70	N	31.95 ②
D	41	60	Y	40.78 ③
E	43	70	Y	31.04 ①
F	38	40	Y	60.08

[Y, N, V] so final prediction is Y

① For this dataset : How to choose the value of KNN

We try different values of K and then based on the accuracy and the value of K which has the best accuracy is chosen. This helps avoiding overfitting and underfitting. Best K is the one which has accuracy rate is highest and accuracy error rate is lowest.

② For Diabetes dataset : What is the purpose of feature scaling ? How to perform ?

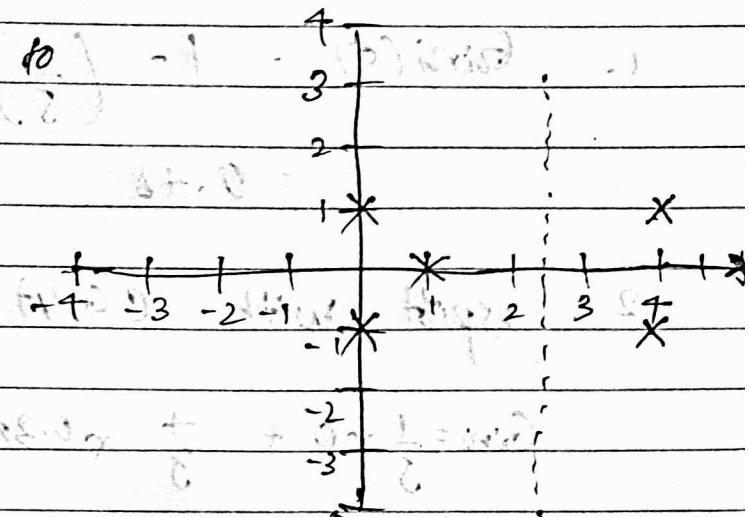
In the diabetes dataset, the features are measured in different units and scales. Some feature have values in hundred other decimals. If we do not scale the largest data entity will dominate. Hence scaling helps all features contribute equally. It is done by standardization and normalization.

Lab - 7

positive classes lie to the right $x=3$

negative class

points are left of $x=2$



$$\text{hyperplane} = 2.5$$

For Iris dataset

Linear Kernel

Accuracy : 1.0

Confusion matrix

$\begin{bmatrix} 10 & 0 & 0 \end{bmatrix}$

0.95 accuracy

$\begin{bmatrix} 0 & 0 & 11 \end{bmatrix}$

Accuracy : 1.0

Confusion Matrix :

$\begin{bmatrix} 10 & 0 & 0 \end{bmatrix}$

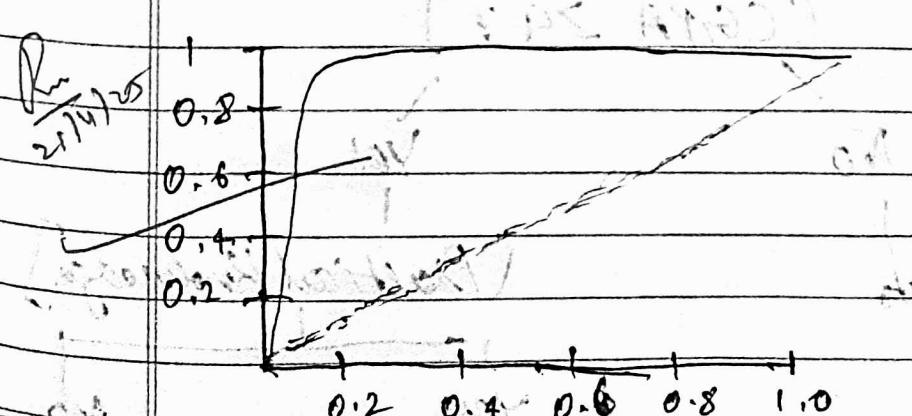
0.95 accuracy

$\begin{bmatrix} 0 & 0 & 11 \end{bmatrix}$

2. For letter dataset

accuracy : 0.95025

AUC : 1.00



Lab - 8

1. Gini(s) = $1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{8}\right)^2 = 1 - 0.36 - 0.125 = 0.48$

2. split with CGPA

$$\text{Gini} = \frac{1}{5} \times 0 + \frac{4}{8} \times 0.375 = 0.300$$

$$\Delta\text{Gini} = 0.48 - 0.300 = 0.18$$

3. growing CGPA ≥ 9 nodes based on interaction

$$\rightarrow \text{Gini} = \frac{2}{4} \cdot 0 + \frac{2}{4} \cdot 0.25 = 0.25$$

$$\Delta\text{Gini} = 0.375 - 0.25 = 0.125$$

\rightarrow based on practical knowledge

$$\text{Gini} = \frac{3}{4} \cdot 0 + \frac{1}{4} \cdot 0 = 0$$

$$\Delta\text{Gini} = 0.375 - 0 = 0.375 \text{ (base)}$$

final east

[CGPA ≥ 9 ?]

No

Yes

No

[Practical Knowledge]

Yes

No

1. Gini = $1 - \left(\frac{4}{8}\right)^2 - \left(\frac{1}{5}\right)^2 = 1 - 0.64 - 0.04$

$= 0.32$

2. Gini = $2 \times 0 + \frac{3}{5} \times 0.444 = 0.2664$

$\Delta \text{Gini} = 0.32 - 0.2664 = 0.0536$

growing the "no" branch with other features

→ CGPA:

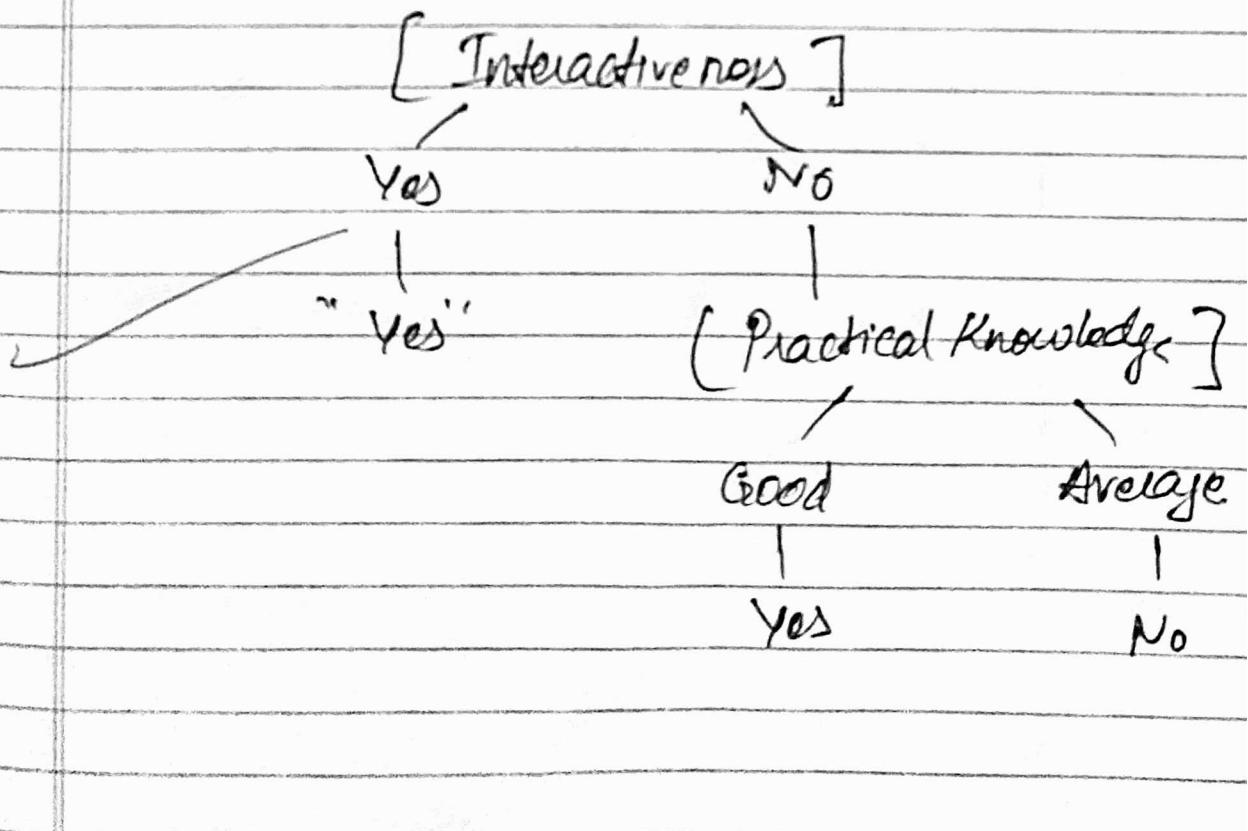
Weighted Gini = $\frac{1}{2} \times 0 + \frac{2}{3} \times 0.50 = 0.333$

$\Delta \text{Gini} = 0.111$

→ Practical Knowledge:

Weighted gini: 0.444

$\Delta \text{Gini} = 0.444$ (best)





Date :

Page No.: 12

RF - n-in-estimators fit mean-accuracy

1	10	0.9667
2	50	0.9667
3	100	0.9667
4	150	0.9667
5	200	0.9667

The best can be chosen with n-estimators

Lab - 9

1. \rightarrow CGPA

107. Actual Job

≥ 9

107.0 + 1

1/6

< 9

+ 1

1/6

≥ 9

- 1

1/6

< 9

- 1

1/6

≥ 9

+ 1

1/6

≥ 9

+ 1

1/6

≥ 9

+ 1

1/6

$$E_{CGPA} = \frac{2 \times 1}{6} = 0.333$$

$$\sigma_{CGPA} = 0.342 = \frac{1}{2} \ln [(-0.333)]$$

$$Z_{CGPA} = \frac{1}{6} \times 4 \times e^{-0.342} + \frac{1}{6} \times 2 \times e^{0.342}$$

$$\text{not } \frac{1/6 \times e^{-0.342}}{0.9428} \text{ not } \frac{1/6 \times e^{0.342}}{0.9428} = 0.1249 = 0.2501$$

\rightarrow 2 interactions

Actual off

102.0 Yes

Yes

Yes 0.1249

102.0 No

No

No 0.2501

102.0 No

No

No 0.2501

102.0 Yes

Yes

Yes 0.1249

102.0 Yes

Yes

Yes 0.1249

$$E_{\text{Interactions}} = 0.2501$$

$$\lambda = \frac{1}{2} \ln \left[\frac{1 - 0.2501}{0.2501} \right] = 0.8490$$

$$Z_{\text{Interactions}} = 0.1249 * 4 * e^{-0.549} + 0.2501 * 1 * e^{-0.549} + 0.2501 * 1 * e^{0.549} = 0.866$$

$$wt_{+ve} = \frac{0.1249 * e^{-0.549}}{0.866} = 0.0832$$

$$wt_{fre} = \frac{0.2501 * e^{-0.549}}{0.866} = 0.1667$$

$$wt_{-ve} = \frac{0.2501 * e^{0.549}}{0.866} = 0.5001$$

→ No instance is misclassified by practical knowledge.

	Comm Skill	Actual Job offer	weight
Good	Yes	Yes	0.0832
Moderate	Yes	No	0.5001
Moderate	No	No	0.1667
Good	No	Yes	0.0832
Moderate	Yes	Yes	0.0832
Moderate	Yes	Yes	0.0832



$$\varepsilon_{CS} = 0.7497$$

$$d = \frac{1}{2} \ln \left[\frac{1 - 0.7497}{0.7497} \right] \\ = -0.5485$$

$$Z_{CS} = 0.0832 \times 1 \times e^{0.5485} \\ + 0.1667 \times 1 \times e^{-0.5485} \\ + 0.5001 \times 1 \times e^{-0.5485} \\ + 0.0832 \times 3 \times e^{-0.5485} \\ = 0.868$$

$\alpha_{CF,PA}$	discreetion	$\alpha_{Committed}$	Avg	Pred
Yes	Yes	Yes	0.3485	Y
No	No	No	0	N
Yes	No	No	0.372	Y
No	No	Yes	-0.5485	N
Yes	Yes	No	0.896	Y
Yes	Yes	No	0.896	Y

2. Adaboost	n-estimators	mean-accuracy
1	10	0.82037
2	50	0.830023
3	100	0.832050
7	180	0.832296
Best	5	0.832624

The best accuracy is from
200 estimators

Lab-10

Record	C1	C2	dist.	C1	C2	dist.	Assignment
1,1	3.04	3.0	0	7.21	7.21	0	C1
1,5,2	7.12	7.12	0	6.12	6.12	0	C1
3,4	3.61	3.61	0	3.61	3.61	0	C1
5,7	7.21	7.21	0	0	0	0	C2
3,5,5	7.12	7.12	0	2.5	2.5	0	C2
4,5,5	7.12	7.12	0	2.06	2.06	0	C2
3,5,4,5	7.12	7.12	0	2.92	2.92	0	C2

$$C_1 = [1.83, 2.33]$$

$$C_2 = [7.12, 5.37]$$

Record	C1	C2	dist.	Assignment
1,1	1.52	5.32	3.81	C1
1,5,2	0.47	4.27	4.27	C1
3,4	2.04	1.27	0.77	C2
5,7	5.64	1.85	4.79	C2
3,5,5	3.15	0.72	2.43	C2
4,5,5	3.28	0.53	2.75	C2
3,5,4,4	2.87	1.07	1.80	C2
3,6,3	0.1	0.1	0.1	C1

$$C_1 = [1.25, 1.5]$$

$$C_2 = [3.9, 5.1]$$

The elbow plot shows a sharp elbow at $K=3$, indicating that three clusters is the optimal choice for the petal length of iris dataset.

Lab-11

$$\begin{array}{c|c|c|c|c|c} x_1 & 4 & 8 & 13 & 7 & \bar{x}_1 = 8 \\ \hline x_2 & 11 & 4 & 5 & 14 & \bar{x}_2 = 8.5 \end{array}$$

$$B = \begin{bmatrix} -4 & 0 & 5 & -1 \\ 2.5 & -4.5 & -3.5 & 5.5 \end{bmatrix}$$

$$S = \frac{1}{N-1} [B B^T]$$

$$= \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

$$\lambda^2 - \text{Tr}(S)\lambda + \det(S) = 0$$

$$\lambda^2 - 37\lambda + 201 = 0$$

$$\lambda = 30.3849, 6.6151$$

$$\begin{bmatrix} -16.3849 & -11 \\ -11 & -7.3849 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} -11 / -16.3849 \\ 1 \end{bmatrix}$$

$$\hat{u}^T = \hat{u}^T = \begin{bmatrix} -0.5574 & 0.8303 \end{bmatrix}$$

$$\hat{u}^T \begin{bmatrix} -0.5574 & 0.8303 \end{bmatrix} \begin{bmatrix} -4.5 & 0 & 5 & -1 \\ 2.5 & -4.5 & -3.5 & 5.5 \end{bmatrix}$$

$$= \begin{bmatrix} 9.305 & -3.736 & -5.693 & +5.124 \end{bmatrix}$$

for heart.csv dataset

logistic regression: 0.901+

SVM = 0.8526

Random forest = 0.8545

After PCA

logistic regression : 0.8689

SVM : 0.8689

Random forest : 0.8852

This shows that with a slight tradeoff
with accuracy and lesser computation
the model yields similar results.