

INDEX

Name : Rachana R. Class : _____

Section : Roll No. : Subject : ML Lab.

Write a python program to import and export data using pandas library functions.

Code :

```
import pandas as pd
iris_data = pd.read_csv("iris.csv")
iris_data.head()
iris_data.columns = col_names.
```

O/P

	sl	sw	pl	pw	class
0	5.1	3.5	1.4	0.2	versicolor
1	4.9	3.0	1.3	0.2	versicolor
2	4.7	3.2	1.4	0.2	versicolor
3	4.6	3.1	1.5	0.2	versicolor
4	5.0	3.6	1.4	0.4	versicolor

Code :

```
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
col_names = ["sepal-length-in-cm",
             "sepal-width-in-cm",
             "petal-length-in-cm",
             "petal-width-in-cm"]
```

```
iris_data = pd.read_csv(url, names=col_names)
iris_data.head()
```

Output :

	sl	sw	pl	pw
0	5.1	3.5	1.4	0.2

```
iris_data.to_csv("cleaned-iris-data.csv")
```

Lab - 2

23 - Standard

- Get the Data
- Downloading other data sources from Kaggle.
- Import as
- import os
- import tarfile
- import urllib
- housing = download住房数据集
("https://raw.githubusercontent.com/mlpacker/housing-dataset/master/housing.tgz")
- Using these libraries, download and extract housing_data and load the data into "housing.csv".
- housing.head()
- housing.info()
- housing.describe()
- Using the above commands inspect the attributes of the loaded housing dataset.
- Using matplotlib and seaborn libraries, by plotting the histogram detect the outliers.
- Create the Test set
- Splitting the "dataset" on test:ratio = 0.2, i.e., training data is 80% of dataset and testing data is 20% of dataset.
- Stratified sampling is where each random chosen data are representative of a whole target population. Each homogeneous subgroup is called strata.

Lab - 03.

Discover and visualize the data to gain insights.

Visualize the data using matplotlib and seaborn libraries.
Calculating the standard correlation coefficient of every pair of columns.

Prepare the data for machine learning algorithms.

Data cleaning, handling text and categorical data, custom transformation, feature scaling, transformation pipelines etc. are done here.

Select and Train model

At first linear regression model is used to train but the model is overfitting the data.
To tackle this, Decision Tree Regression model is used as it is capable of finding non-linear relationships with data. But DTRM is also overfitting so badly that it performs worse than the linear RM.

Fine tune your model.

At last, fine tuning of model is carried out, evaluating on the test set and then adding more and maintaining the system.

(Part 1, 8.8.1) : (1d, 1e)

Lab-04

Python implementation of linear regression

```
import numpy as np
import matplotlib.pyplot as plt

def estimate_coef(x, y):
    n = np.size(x)
    m_x = np.mean(x)
    m_y = np.mean(y)

    ss_xy = np.sum(y*x) - n*m_y*m_x
    ss_xx = np.sum(x*x) - n*m_x*m_x

    b_1 = ss_xy / ss_xx
    b_0 = m_y - b_1 * m_x

    return (b_0, b_1)

def plot_regression_line(x, y, b):
    plt.scatter(x, y, color='m', marker='o')
    y_pred = b[0] + b[1]*x
    plt.plot(x, y_pred, color='g')

    plt.xlabel('x')
    plt.ylabel('y')

def main():
    x = np.array([0, 1, 2, 3, 4])
    y = np.array([1, 3, 2, 4, 5])
    b = estimate_coef(x, y)
    print(b)
    plot_regression_line(x, y, b)

    O/P
    (b0, b1) = (1.263, 1.169)
```

→ Multiple linear regression.
from sklearn import linear_model
.train-test-split "x_train" & "x_test"
import numpy as np
import datasets
from sklearn import datasets
linear_model, metrics

x_train = pd.read_csv("iris.csv")
raw_df = pd.read_csv("iris.csv")
x = np.vstack([x_train, raw_df])
raw_df["values"]
raw_df["values"]
y = raw_df["values"]

x_train, y_train, x_test, y_test =
train-test split (x, y, test_size=20%,
random_state=1)

reg = linear_model.LinearRegression()

reg.fit(x_train, y_train)

print("co-efficients:", reg.coef_)

print("variance score:", reg.score)

(reg.score(x_test, y_test))

plt.style.use('fivethirtyeight')

plt.scatter(reg.predict(x_train), reg.
predict(x_train) - y_train, color="green")
s = 10, label="Train data")

plt.scatter (reg.predict (x-test), y-test, color = "blue", s=10, label = "test data")
plt.lines (y=0, xmin=20, xmax=50, linewidth=2)

plt.legend (loc = 'upper right')
plt.title ("residual errors")
plt.show()

• design matrix test vs training
• test set vs judge test results
• (I = identity matrix)
• transpose: "the rows remain" (test vs train reference)
• (judge plant with new plant) (test vs train)
• (test vs train reference) where all
"rows" = rows (train - x) - (test - x) following
• (rows mean) = 1 std. dev. ± 0.1 ± 0.

TD3

```

import numpy as np
import pandas as pd.
eps = np.finfo(float).eps
Path = 'drive/My Drive/ML lab/datasets/plays-tennis.csv'
df = pd.read_csv(Path)
def find_entropy(df):
    target = df.keys()[-1]
    entropy = 0
    values = df[target].unique()
    for value in values:
        fraction = df[target].value_counts()[value] / len(df[target])
        entropy += - fraction * np.log2(fraction)
    return entropy
def average_information(df, attributes):
    target = df.keys()[-1]
    target_variables = df[target].unique()
    variables = df[attributes].unique()
    entropy_2 = 0.

```

```

for variable in variables:
    entropy = 0
    for target_variable in target_variables:
        num = len(df[df['attribute'] == target_variable][df['target'] == target_variable])
        iden = len(df[df['attribute'] == target_variable])
        entropy += (num / iden) * np.log((num / iden))
    print("Entropy for", variable, "is", entropy)

```

fraction = num / (den + eps)

entropy + = - fraction * log(fraction + eps)

$$\text{fraction 2} = \text{den} / \text{len}(\text{df})$$

~~entropy 2 + = -fraction 2 * entropy~~

return abs ~~(entropy 2)~~

~~Elicit questions - Intj. felt
about us - simple past tense~~

Chapman (Spartak J. fbs = 2nd day)

प्रत्येक विद्यार्थी का अपना अपना नाम है।

— and the Regent I hope — well.

(Eigenvectors) \vec{v}_1 , \vec{v}_2 , \vec{v}_3 = nontrivial
 \vec{v}_1 , \vec{v}_2 , \vec{v}_3 = linearly independent

(negative) \neg (positive) \vdash (negative) \neg (positive)

Citrus verna (L.) DC. var. *variegata*

• U-3C) was, I do -

Opinions. [Original] The = experts.

Umweltwiss. Studiengang = Wissenschafts- und Forschungsstudiengang

.0 = Superficial

R NB: Implementation

```
import numpy as np
import pandas as pd
from google.colab import drive
drive.mount('/content/drive')
```

```
dataset = pd.read_csv('/content/drive/My Drive/Iris.csv')
```

```
dataset.head()
```

```
dataset.groupby('species').size()
```

& species.

Iris-setosa 50

Iris-versicolor 50

Iris-virginica 50

feature columns = ['Sepal length',
'Sepal width', 'Petal length',
'Petal width']

```
x = dataset[feature columns].values
```

```
y = dataset['species'].values
```

from sklearn.preprocessing import
LabelEncoder

```
le = LabelEncoder()
```

```
y = le.fit_transform(y)
```

from sklearn.model_selection
import train_test_split

```
x_train, x_test, y_train, y_test
```

```
= train_test_split(x, y)
```

```
import matplotlib.pyplot as plt  
import seaborn as sns  
%matplotlib inline  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn import confusion_matrix  
from sklearn.model_selection import cross_val_score
```

```
classifier = KNeighborsClassifier(n_neighbors=3)  
classifier.fit(x_train, y_train)  
y_pred = classifier.predict(x_test)  
accuracy = accuracy_score(y_test,  
                           y_pred) * 100.
```

O/P : Accuracy of our Model is
equal 96.67%
X-axis. [Sepal length] feature = X
X-axis. [Sepal width] feature = Y
target values which must be same
as our target label

(Sepal length) = 1
(Sepal width) = 2
target values which must be same
as our target label

test = np.array([5, 3])
(5, 3) file 2nd row =

Lab - 7. Python

K-Means Implementation.

```
import pandas as pd
data = pd.read_csv('Iris.csv')
data.head()
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import datasets
iris = datasets.load_iris()
X = iris.data
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=3, random_state=2)
kmeans.fit(X)
```

~~pred = kmeans.fit_predict(X)~~

~~pred~~ (Int32) \rightarrow np.array

$2 \cdot 0 \rightarrow$ np.array

$[0] \rightarrow$ int32 \rightarrow np.array

$[0] \rightarrow$ int32 \rightarrow np.array

$0 \in \{0, 1, 2\} \times \{0, 1, 2\}$ where $0, 1, 2 \in \{0, 1, 2\}$

Lab - 9.

SVM.

from sklearn.datasets import
load_breast_cancer

from sklearn.inspection import
DecisionBoundaryDisplay

from sklearn.svm import SVC
cancer = load_breast_cancer()

X = cancer.data[:100]

y = cancer.target[:100]

svm = SVC(kernel='rbf', gamma=0.5)
(svm, fit) = (X, y)

Decision Boundary Display. from sklearn
estimator (

svm,

responsible - method = 'predict'

map = plt.cm.spectral,

'alpha' = 0.8,

xlabel = cancer.feature_names[0],

ylabel = cancer.feature_names[1],

)

plt.scatter(X[:, 0], X[:, 1], c=y, s=200,
plt.show().

for

Lab -7.

```
PCA.  
import pandas as pd.  
import numpy as np.  
import matplotlib.pyplot as plt  
%matplotlib inline  
from sklearn.decomposition  
import PCA.  
from sklearn.preprocessing import  
StandardScaler  
from sklearn.datasets import  
load_breast_cancer()  
data = data[feature_names]  
perm = data[feature_names]  
df1 = pd.DataFrame(data[feature_names])  
columns = data[feature_names]  
plt.figure(figsize=(10, 10))  
plt.scatter(x[:, 0], x[:, 1], c=claster[  
    'target'], cmap='plasma')  
plt.xlabel('PCA')  
plt.ylabel('PCA')  
(x1, x2) = columns  
(y1, y2) = columns  
(x1-x2, y1-y2) = df1  
(x1-x2) = df1  
(y1-y2) = df1
```

Lab 8

Random Boosting Forest

Import numpy as np
 Import pandas as pd
 Import matplotlib.pyplot as plt
 Import seaborn as sns
 Import sklearn
 from sklearn import datasets
 from sklearn.model_selection import train_test_split
 from sklearn.ensemble import RandomForestClassifier
 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
 This feature names data = pd.read_csv('iris.csv')
 This data.head()
~~X = data.loc[:, :-1].values~~
~~y = data.loc[:, -1].values~~
~~X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)~~
 model = RandomForestClassifier()
 model.fit(X_train, y_train)
~~y_pred = model.predict(X_test)~~
~~accuracy_score(y_test, y_pred)~~

```
model1 = RandomForestClassifier(n_estimators=10)
```

```
model1.fit(x_train, y_train)
```

Ada Boosting

```
from sklearn.datasets import load_iris  
from sklearn.model_selection import train_test_split
```

```
from sklearn.ensemble import AdaBoostClassifier  
from sklearn.metrics import accuracy_score
```

```
iris = load_iris()
```

```
X = iris.data
```

```
y = iris.target
```

```
x_train, x_test, y_train, y_test = train-  
test-split (X, y, test-size=0.4).
```

```
adaboost_clt = AdaBoostClassifier(n_estimators=30,  
learning_rate=1, random_state=2, 2)
```

```
adaboost_clt.fit(x_train, y_train)
```

```
y_pred = adaboost_clt.predict(x_test)
```

```
accuracy = accuracy_score(y_test, y_pred)
```

```
print("Accuracy: ")
```

O/P accuracy = 0.9666

~~Adaboost~~
~~30.0%~~