

IBM21CS235



Name VAISHNAVI KAMATH Std _____ Sec _____

Roll No. _____ Subject ML LAB School/College _____

School/College Tel. No. _____ Parents Tel. No. _____

Sl. No.	Date	Title	Page No.	Teacher Sign / Remarks
1.	21/3/24	LAB-1	A	
2.	28/3/24	LAB-2	A	
3	4/4/24	LAB-3	A	
4	18/4/24	LAB-4	A	
5.	25/4/24	Insurance system		
6.	9/5/24	LAB-6	A	
7.	23/5/24	LAB-7		
8.	28/5/24	CAB-7	A	
9.	30/5/24	CAB-8 . PCA		
10.	30/5/24	CAB-8 (green)	A	

> import pandas as pd
 url = "https://...
 df = pd.read_csv(url), columns=[
 'sepal.length', 'sepal.width', 'petal.length',
 'petal.width', 'class'])

df.head()

> df.to_csv("cleaned_data.csv")

	sepal.length	sepal.width	petal.length	petal.width	class
0	5.1	3.5	1.4	0.2	Iris-setosa

PANDAS COURSE

1. Creating, Reading & Writing.

> import pandas as pd

> pd.DataFrame({'Yes': [1, 2], 'No': [3, 4]},
 index=['A', 'B'])

> df

> pd.Series([30, 40, 50], index=['2015 sales',
 '2016 sales', '2017 Sales'], name='ProductA')

2. Indexing, Selecting & Assigning.

> reviews.country (or) reviews['country']

> reviews.iloc[[0, 1, 27, 0]]

> reviews.iloc[-5:]
 ↳ index based sel.

> reviews.loc[:, ['a', 'b', 'c']]

- > reviews.set_index('NI')
- > reviews.loc[reviews.country == 'Italy' & (reviews.points >= 90)]
- > reviews.loc[reviews.country.isin(['Italy', 'France'])]
- > reviews.loc[reviews.price.notnull()]
- 3. summary fns & Maps
 - > df.describe()
 - > df.col.unique()
 - > df['A'].value_counts()
 - > df.points.map(lambda p: p - mean)

- #### 4. Grouping & sorting
- > reviews.groupby(['country', 'province']).apply(lambda df: df.loc[df.points.idmax()])
 - > reviews.groupby(['country']).price.agg(['min', 'med', 'max'])

- > countries - reviewed .reset_index()
- > df . sort_index()

5. Data Types and Missing Values

- > df . dtypes
- > df . points . astype ('float64')
- > reviews [pd . isnull (reviews . country)]
- > reviews ['A']. replace ('A', 'B')

6. Renaming & Combining

- > reviews . rename (columns = { 'points' : 'score'})
- > reviews . rename . axis ('titles', axis = 'row')
- > reviews . rename . axis ('fields', axis = 'column')
- > pd . concat ([df1, df2])
- > left . join (right, lsuffix = '_CAN', rsuffix = '_UK')

End to End Project

28/3/2023
Date
Page 3/2

- import os
- import tarfile
- import urllib
- ```
DOW_ROOT = "https://..."
HOUSING_PATH = os.path.join("data", "ol罕")
HOUSING_URL = DOW_ROOT + "dataset/housing.tgz"
```
- ```
def fetch(housing_url=HOUSING_URL,  
         housing_path=HOUSING_PATH):  
    os.makedirs(name=housing_path,  
                exist_ok=True)
```
- ```
tgz_path = os.path.join(housing_path,
 "housing.tgz")
```
- ```
urllib.request.urlretrieve(url=housing_url,  
                           filename=tgz_path)  
housing_tgz = tarfile.open(name=housing_tgz)
```
- ```
housing_tgz.extractall(path=housing_path)
```
- ```
housing_tgz.close()
```
- ```
fetch_housing_data()
```
- ```
import pandas as pd
```

def lhd(housing_path, housing, path):

data_path = os.path.join(housingpath,
housing.csv))

return pd.read_csv(data path)

→ housing = load('housing.csv')
housing.head()
housing['ocean_proximity'], value_counts()

→ import malplotlib.pyplot as plt
import seaborn as sns
sns.set(figure_dpi=50; figsize=(20, 15))
plt.show()

CREATING TEST SET

```
→ import numpy as np
def split(data, test_ratio=0.2):
    si = np.random.permutation(len(data))
    test = int(len(data) * test_ratio)
    test_ind = si[:test]
    train_ind = si[test:]
    return data.iloc[train_ind], data.iloc[test_ind]
```

$\rightarrow \text{train}, \text{test} = \text{split}(\text{train-test}(\text{data} = \text{housing}))$

\rightarrow (entrain set), (test set)

HASH STRATEGY

```
from glib import crc32
def test_set(identifier, test_ratio=0.2):
    total_size = 2 * 32
    hex_rep = crc32(np.int64(identifier)) & 0xffffffff
    in_test = hex_rep < test_ratio * total_size
    return in_test
def split_train_by_id(data, test_ratio, id_column):
    ids = data[id_column]
    in_ids = ids.apply(lambda id: in_test_set(id, test_ratio))
    return [data.loc[in_ids], data.loc[~in_ids]]
h_w_i = housing.reset_index()
train, test = split_train_by_id(h_w_i, test_ratio=0.2, id_column="index")
def Z_ton(z):
    if (z >= 0):
        n = 2 * z
    else:
```

$$n = 2 \times 3$$

else:

$$n = -2 \times (3 - 1)$$

3. Discover & Visualise Data

\rightarrow strat_train.set.shape, strat_test.set.shape.

\rightarrow strat_train.set.reset_index().to_feather(path='')

\rightarrow housing = strat_train.set.copy()

\rightarrow housing.plot(kind='scatter', x='longitude', y='latitude', alpha=0.1)

\rightarrow corr = housing.corr()

\rightarrow corr_matrix['median_hvh'].sort_values(ascending=False)

ATTRIBUTE COMBINATIONS

\rightarrow housing[['rooms_per_hh']] = housing[['bedroom']]
housing[['households']]

4. PREPARE DATA FOR ML ALgos

h = strat_train.set.drop(['med_house_value'], axis=1)

h_labels = strat_train.set['median_house_value'].copy()

h.shape, h_labels.shape

Data cleaning

Date _____
Page _____

from sklearn.impute import
SimpleImputer

imputer = SimpleImputer(strategy
median)

housing_num = housing.drop(["ocean_proximity"], axis=1)

imputer.fit(housing_num)

imputer.statistics()

housing_num.median().values

X = imputer.transform(housing_num)
X.shape

housing_tr = pd.DataFrame

(data=X, index=housing_num.index, columns=housing_num.columns)

housing_tr.read()

housing_cat = housing[["ocean_proximity"]]

housing_cat.head(10)

housing_cat["ocean_proximity"].value_counts()



from sklearn.preprocessing import
OrdinalEncoder

ordinal_encoder = OrdinalEncoder()

housing_cat_encoded[:10]

ordinal_encoder.categories_

→ from sklearn.preprocessing import
OneHotEncoder

one_hot_encoder = OneHotEncoder()

housing_cat_1hot = one_hot_encoder.
fit_transform(housing_cat.values)

housing_cat_1hot.toarray()

one_hot_encoder.categories_

↳ values returns memory representation
of a given array.

of Bob Dylan's books provided

1922-50 - 1922-50 - 1922-50 - 1922-50 - 1922-50 -

~~Wesleyen~~ • Johnson told me

Chlorine had to be added

Chagall's Soldier

A HISTORY OF THE AMERICAN PEOPLE

A HISTORY OF THE AMERICAN PEOPLE

ANSWER

Sample sentence
of = fed, read as
(by head)

→ sis. desplot (df) [
plt. show()]

→ plt.scatter(df)
plt.show()

SPLIT DATA

$\rightarrow X = \text{train}$, $X = \text{test}$

$\rightarrow y = \text{train}$, $y = \text{test}$

\rightarrow reg-linear Regress.

of pred test =
of pred train =

point (seq. ref.) point (seq. ref.)

PART-B

Date 11/1/24 Page 4

→ Import data

```
df = pd.read_csv("glass.csv")  
df.head()
```

→ df.describe()

```
sns.distplot(df['Salary'])  
plt.show()
```

```
→ plt.scatter(df['Years of'], df['Sal'])  
plt.show()
```

SPLIT DATA

```
→ X = df.iloc[:, 1:]  
y = df.iloc[:, 0]
```

```
→ X_train, X_test, y_train, y_test  
= train_test_split(X, y, test_size=0.2,  
random_state=0)
```

```
→ reg = LinearRegression()  
reg.fit(X_train, y_train)
```

```
→ y_pred = reg.predict(X_test)  
y_train = reg.predict(X_train)
```

```
point(reg, y_of=1)  
point(reg, index_of=1)
```

MULTIPLE LINEAR REG.

→ df = pd.read_csv('... .csv')
df.head()

→ df.describe()

→ sns.distplot(df['profit'])

→ plt.scatter(df[['R&D Spend']],
df[['Profit']])

→ plt.show()

→ X = df.iloc[:, :-1].values
y = df.iloc[:, -1].values

One hot encoding of categorical data:

→ ct = ColumnTransformer([('encoder', OneHotEncoder(), [3]),
[('remander', 'passthrough', random_state=0)])

→ X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=0)

→ reg = LinearRegression()
reg.fit(X_train, y_train)

→ y_pred = reg.predict(X_test)

Wetland soils are often saturated with water and may have
high mineral content. They are usually acidic and
have low infiltration rates.

Wednesday, X - November 16
Went to the library - took out
books - read them - wrote
notes - did some work -
and then went home.

Consequently, many countries

Decision Tree

```

→ url = " - - - "
df = pd.read_csv(url, header=0)
names = [c.replace(' ', '_') for c in df.columns]
names[0] = 'id'
df = df[names]
df['species'].value_counts()
df['species'].value_counts().plot(kind='bar')
df['species'].value_counts().plot(kind='bar').set_title('Species Distribution')

X = df.drop(['species'], axis=1)
y = df[['species']]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

```

\Rightarrow oft = DecisionTreeClassifier(max_depth)
min_samples_leaf = 10
random_state = 1

- dot-fit (train, if train)
- from IPython, display report
 - Image
 - from sklearn.tree import
 - export graphviz
 - ! pip install graphviz
 - import dotgraph

- features = X, columns
- set data = ~~export graphing~~
C:\tfl\set file = ~~one~~, feature names = features
- graph = ~~syndicates~~, graph from set data (dot, data)



hyperparameter tuning:

→ dt = DecisionTreeClassifier(
random_state=1,
params = { 'max_depth': [2, 3, 4, 5],
'min_samples_split': [2, 3, 4, 5],
'min_samples_leaf': [1, 2, 3, 4, 5]})

qsearch = GridSearchCV(dt,
param_grid = params, cv=3)
qsearch.fit(X, y)
qsearch.best_params_

→ dt2 = DecisionTreeClassifier(
(# search. best_params)
random_state=1)
dt2.fit(X_train, y_train)
y_pred_train = dt2.predict(X_train)
y_pred_train = dt2.predict_proba(X_train)[:, 1]

yf_pred = dt2.predict(X_test)
yf_prob = dt2.predict_proba(X_test)
[:, 1]

→ point(accuracy_score(y_train, yf_train))
→ point(accuracy_score(yf_test, yf_prob))

petal width $\lambda = 0.8$
petal length $\lambda = 1.3882$
 $y_{true} = 0.667$
samples = 150
values = [50, 50, 50]

True

False

$y_{true} = 0.0$
Samples = 50
Value = [50, 90]

petal width
cm = 1.75
petal length
cm = 0.5
samples = 150
values = [50, 50, 50]

No classification error.
using $\lambda = 0.8$.
using $\lambda = 1.3882$.

Classification error.
using $\lambda = 0.8$.
using $\lambda = 1.3882$.

Classification error.
using $\lambda = 0.8$.
using $\lambda = 1.3882$.

Classification error.
using $\lambda = 0.8$.
using $\lambda = 1.3882$.

~~Objectives~~

Insurance System

import pandas as pd
from sklearn import preprocessing
from sklearn import model_selection
import train_test_split
from sklearn.linear_model import LogisticRegression
for LogisticRegression

```
df = pd.read_csv('c id.csv')
```

```
X_train, X_test, y_train, y_test = train_test_split(df['age'],  
df['height'],  
test_size=0.2)
```

```
# model = LogisticRegression()  
model = fit(X_train, y_train)  
y_predicted = model.predict(X_test)  
print(y_predicted)
```

```
# print(model.predict_proba(X_test))  
# print(model.score(X_test, y_test))
```

Accuracy:

0.5

```
def sigmoid(z):  
    return 1/(1+math.exp(-z))
```

```
def predict():
```

$$z = 0.042 * \text{age} - 1.53$$

$$y = \text{sigmoid}(z)$$

```
return y
```

```
print(predict(15))  
print(predict(43))
```

```
Predictions : 0.485  
0.5685
```

SVM - Iris Dataset

\rightarrow df = pd.read_csv('IRIS.csv')
 \rightarrow bnpairplot (gova = df, hue = 'species',
palette = 'Set2')

\rightarrow from sklearn.model_selection import train_test_split.

```
x = df.iloc[:, :-1]
y = df.iloc[:, -1]
x_train, x_test, y_train, y_test
= train_test_split(x, y, test_size = 0.3)
```

\rightarrow from sklearn.svm import SVC
model = SVC(C = 5, kernel = 'rbf')
model.fit(x_train, y_train)
pred = model.predict(x_test)
print(classification_report(y_test, pred))
print(classification_report(y_test, pred))

O/P:

Accuracy : 0.98
Precision : 0.98
Recall : 0.98

KNN

Date _____
Page _____

from sklearn.neighbors import
KNeighborsClassifier
 \rightarrow from sklearn.metrics import
confusion_matrix, accuracy_score

cls = KNeighborsClassifier(n_neighbors=3)

cls.fit(x_train, y_train)
 $y_{pred} = \text{cls.predict}(x_{test})$

\rightarrow print(confusion_matrix(y_test, y_pred))
 \rightarrow print(accuracy_score(y_test, y_pred))

CROSS VALIDATIONS

~~MSE = $\frac{1}{k} \sum \text{for } k \in \text{list}$~~

$k \text{ list} = \text{list}(\text{range}(1, 50, 2))$
 $\text{on scores} = []$

$\text{for } k \text{ in } k \text{ list}:$
 $\text{nn} = \text{KNeighborsClassifier(n_neighbors}=k\text{)}$
 $\text{scores} = \text{cross_val_score}(\text{nn},$
 $x_{train}, y_{train}, cv=10)$
 $\text{scoring} = \text{'accuracy'})$

$\text{mse} = \sum \text{for } k \text{ in scores}$
 plt.figure(1)

$\text{best_k} = \text{k} \text{ list MSE. index (min(mse))}$
 print(best_k)



Q.D:

Optimal no. of neighbors = 7
Accuracy = 97.78%.

((1B) 5 of 5))

الله يحيى، سعيد

MATERIALS AND METHODS

Longfellow - 1857

Dolby Digital 5.1 Surround Sound

卷之三

2007-12-15 10:00:00

卷之三

100

Amidst the clouds of life

and 2 sides 1000

2009-2010

the first time I ever saw him

no) a) 3 persons (not).

1. [S] good boy

23/3/2021

3. Back propagation

A 23-5-21

Date _____
Page _____

→ `db = np.loadtxt("duke-breast-cancer.txt")`
`print(np.shape(db))`

→ `np.random.shuffle(db)`
~~x_train, x_test, y_train, y_test
= train_test_split(x, y, test_size=0.1)~~

→ `hidden_layer = np.zeros(72)`
`weights = np.random.random((len(x[0]), 72))`

`output_layer = np.zeros(2)`
`hidden_weights = np.random.random((72, 2))`

`def back_propagation(hidden_layer,
output_layer, one_hot_encoding,
learning_rate, x):`

`op_derivative = np.zeros(2)`
`op_gradient = np.zeros(2)`

`for i in range(0, len(output_layer)):`
`output_derivative[i] = (1.0 - output_layer[i]) *`
`output_layer[i]`

for i in range(0, len(db)):
op_grad[i] = op_derivative[hidden_layer[i]] * output

hidden_derivative = np.zeros(3)
hidden_gradient = np.zeros(3)
for i in range(0, len(db)):
hidden_derivative[i] = (1.0 - hidden_layer[i]) * (hidden_layer[i]) * (1.0 + hidden_layer[i])

for i in range(0, len(db)):
if db[i][0] == 0:
hidden_derivative[i] = 0

for j in range(0, len(db)):
sum += output[j] * hidden_weights[j]

recalculate weights (using
hidden_weights, output, hidden_layer)

recompute bias?
refine
(valenceal BP - valenceal BP)
refine
valenceal BP, (grad = 0.0)
iterations in loop

for i in range(0, len(outputlayer)):
 op_weight[i] = op_derivative[i] * x
 (Concate one [i] * output layer[i])

hidden_derivative = input_layer[0]
hidden_grohent = np.zeros([72])

for i in range(0, len(hiddenlayer)):
 hidden_derivative[i] = (1.0 - hidden_layer[i]) * (1.0 - hidden_layer[i]) * (1.0 - hidden_layer[i])

for i in range(0, len(hiddenlayer)):
 bias[i] = 0
 for j in range(0, len(outputlayer)):
 hidden_weight[i][j] = 0

bias[i] = output_weight[i]

recalculate weights (learning rate)
hidden_weight[i][j] = learning_rate * hidden_derivative[i][j] * output_weight[j][i] * hidden_layer[i] * output_layer[j]

update weights
new_weight[i][j] = old_weight[i][j] + learning_rate * hidden_derivative[i][j] * output_weight[j][i] * hidden_layer[i] * output_layer[j]

new_weight[i][j] = new_weight[i][j] + learning_rate * hidden_derivative[i][j] * output_weight[j][i] * hidden_layer[i] * output_layer[j]

elements in the learning rate (learning)

g3\5(B)
q(1) & b) Random forest q

Adaboost

- + dataset = pd.read_csv('faults.csv')
- + con_matrix = import_xr.load('con_matrix.csv')
- + upper = con_matrix.where(np.triu(np.ones(con_matrix.shape), k=1).astype(np.bool))

to

to_drop = [col for col in upper.columns if any(upper[column] > 0.95)]

print(to_drop)

+ X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.1)

+ n_estimators = 7
n_estimators = 100
RANDOM_SEED = 13.

names = ['RF classifier', 'AdaBoost classifier']
models = [RandomForestClassifier(n_estimators=n_estimators),
AdaBoostClassifier(decision_function_shape='ovo',
learning_depth=None), n_estimators
= n_estimators], n_estimators
for counter, model in enumerate(models):



model: fit(X_{train} , y_{train})

$y_{\text{pred}} = \text{model}.\text{predict}(X_{\text{test}})$

point([nans, [sumter] + 1, +
metrice, accuracy, score(y - pred)])

$$RF = 0.7684$$

$$\text{Adaboost loss function} = 0.7341$$

WPS 2021 30/5/24

K-Means Clustering (line)

```
#> iris = pd.read_csv('iris.csv')  
#> iris.info()  
#> iris[0:10]  
#> iris + outcome = pd.concat([iris, outcome], axis=1)  
#> outcome = np.array(outcome)  
#> outcome = np.insert(outcome, 0, 1, axis=1)  
#> outcome = np.delete(outcome, 0, axis=1)
```

```
#> # Finding optimum no. of clusters  
#> # for kmeans.  
#> from sklearn.cluster import KMeans  
#> kmeans = [ ]
```

```
for i in range(1, 11):  
    km = KMeans(n_clusters=i)  
    init = km.fit([[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]] * 300)  
    n_init = 10; random_state=0)
```

```
kmeans = fit(x)  
wcss.append(kmeans.inertia_)
```

```
#> #> #> kmeans = KMeans(n_clusters=3)  
#> #> init = ('kmeans++', max_iter=300,  
#> #> n_init=10, random_state=0)
```

y_kmeans = kmeans.fit_predict(x)



30/5/24

PCA (Breast Cancer Dataset)

- cancer = load_breast_cancer()
- scaler = StandardScaler()
- scaled_data = scaler.transform(df)
- from sklearn.decomposition import PCA
- pca = PCA(n_components=2)
- pca.fit(scaled_data)
- α pca = pca.transform(scaled_data)
- scaled_data.shape → 2 pca.components_
- df['comp'] = pd.DataFrame(pca.components_, columns=['feature_0', 'feature_1'])
- plt.figure(figsize=(12, 6))
sns.heatmap(df['comp'].T, cmap='plasma')