



28/3/24

Week-2

Date \_\_\_\_\_  
Page \_\_\_\_\_Hands on Machine Learning with Scikit  
(Learn Keras and Tensor Flow)

## 1. Get the Data

Download Import as

the import to life

Data import urllib

```
def fetch_housing_data(housing_url = HOUSING_URL,
                        housing_path = HOUSING_PATH):
```

```
    urllib.request.urlretrieve(housing_url, housing_path)
```

```
    housing_tag = tempfile.open('name: tag_path')
```

```
    housing_tag.writeall(fetch_housing_data(housing_path))
```

```
    housing_tag.close()
```

```
    fetch_housing_data()
```

Load import pandas as pd.

def load\_housing\_data(housing\_path = HOUSING\_PATH)

data\_path = os.path.join

```
    return pd.read_csv(data_path)
```

```
    housing.head()
```

```
    housing.info()
```

```
    housing['ocean-proximity'].value_counts()
```

```
    housing.describe()
```

Create a def split\_train\_test(data, test\_ratio=0.2)

Test shuffled\_indices = np.random.permutation(len(data))

Let test\_set\_size = int(len(data) \* test\_ratio)

```
    test_indices = shuffled_indices[test_set_size:]
```

```
    return data[test_indices]
```

```
def test_set_check(indicator, test_ratio=0.2)
```

```
    total_size = 28 * 32
```

```
    indicator_max_val < (test_ratio * total_size)
```



def cantor\_pairing (n1, n2)  
 $n = ((n1 + n2) * (n1 + n2 + 1) / 2) + n2$   
 return n

def lat\_lon\_to\_index (lat, lon)  
 lat, lon = int (lat \* 100), int (lon \* 100)  
 lat, lon = 100 - 2 to N/lat, 100 - 2 to N/lon  
 index = cantor\_pairing (lat, lon)  
 return np.int64 (index)

def from\_2\_to\_N (z)  
 if z >= 0  
 $n = z * z$   
 else  
 $n = -z * z - 1$   
 return n

2. Discover and Visualize the data to gain insight.

strat\_train, test = train\_test\_split (data, target, test\_size=0.1, random\_state=42)  
 strat\_train, strat\_test = strat\_train, strat\_test  
 strat\_train, strat\_test = strat\_train, strat\_test

Visualizing geographical data

housing.plot (kind='scatter', x='longitude', y='latitude',  
 plot.show())

housing.plot (kind='scatter', x='longitude', y='latitude',  
 plot.show())

3 Prepare the Data for Machine Learning Algorithms

housing = data\_train.set\_drop(["median\_house\_value",  
axis=1])

housing\_labels = data\_train.set(["median\_house\_value"]  
copy())

housing.shape, housing\_labels.shape

→ Data Cleaning -

```
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='median')
imputer.fit(housing_train)
imputer.statistics_
X.shape
housing_train.median().values
housing_train.head()
```

→ Custom transformer.

```
from sklearn.base import TransformerMixin,
BaseEstimator
class CombinedAttributeAdder(BaseEstimator,
TransformerMixin):
    self.add_bedrooms_per_room =
    add_bedrooms_per_room
    def fit(self, X, y, y_ = None):
        rooms_per_household = X[:, 'rooms'] /
        X[:, 'household']
        housing_extra_attribs = attr_adder.transform
        (housing values)
```



## Transformation pipeline -

from sklearn.pipeline import Pipeline

from sklearn.preprocessing import StandardScaler

New = Pipeline([

    ('imputer', SimpleImputer(strategy='most\_frequent')),

    ('scaler', StandardScaler())

])

## 4. Select and Train a Model

from sklearn.linear\_model import LinearRegression

lin\_reg = LinearRegression()

some\_data = housing.data[0:50]

some\_data\_prepared = full\_pipeline.transform(some\_data)

lin\_reg.fit(some\_data\_prepared, housing.target[0:50])

lin\_reg.predict(some\_data\_prepared)

from sklearn.tree import DecisionTreeRegressor

dec\_reg = DecisionTreeRegressor()

dec\_reg.fit(some\_data\_prepared, housing.target[0:50])

dec\_reg.predict(some\_data\_prepared)

## 5. Fine Tune Your Model

GridSearch

from sklearn.model\_selection import GridSearchCV

param\_grid = {

    'criterion': ['mse', 'mae'],

    'max\_depth': [3, 4, 5],

    'min\_samples\_split': [2, 5, 10],

    'min\_samples\_leaf': [1, 2, 5]

grid search - best param

grid search best estimator

cross-validated search cv results

using eq of labels - interval confidence, for (approx-  
over - 1)

10 x - approx - error

6 Launch, Monitor and Maintain your systems.



# Linear Regression

input numpy as np  
input matplotlib.pyplot as plt

```
def calculate_coef(x,y):
```

```
    n = np.size(x)
```

```
    mx = np.mean(x)
```

```
    my = np.mean(y)
```

```
    Sxx = np.sum(x**2) - n*mx*mx
```

```
    Sxy = np.sum(x*y) - n*mx*my
```

```
    b1 = Sxy/Sxx
```

```
    b0 = my - b1*mx
```

```
    return (b0,b1)
```

```
def plot_regression_line(x,y,b):
```

```
    plt.scatter(x,y,color='r',marker='o',s=30)
```

```
    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([0,1,2,3,4,5,6,7,8,9])
```

```
    y = np.array([1,2,3,2,5,7,9,10,8,12])
```

```
    b = calculate_coef(x,y)
```

```
    print('estimated coefficients are: b0 =',b[0], 'b1 =',b[1])
```

```
    plt.show()
```

Output-

```
(b0,b1) = (1.2, 1.63, 1.6313)
```

### ③ Multiple Linear Regression

```
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import numpy as np
from sklearn import datasets, linear_model, metrics
```

```
data_url = 'url'
```

```
row_d1 = pd.read_csv(data_url, skiprows=23,
                      header=None)
```

```
X = row_d1.values[1::3, 1:]
y = row_d1.values[1::3, 2]
```

```
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.4, random_state=1)
```

```
reg = linear_model.LinearRegression()
reg.fit(X_train, y_train)
print("coefficients = ", reg.coef_)
print("variance score = %f" % reg.score(
    X_test, y_test))
```

```
plt.style.use('HovrThiz style')
plt.scatter(reg.predict(X_test), reg.predict(X_train),
            c='blue', s=10, label='fit data')
plt.scatter(y=0, x_min=0, x_max=50, width=1)
plt.legend(loc='upper right')
plt.title('Residuals')
plt.show()
```



## Week-4

①

Decision trees ID 3.

```
import numpy as np
import pandas as pd
df = pd.read_csv('v11')
df.head()
df.info()
df.describe()

def find_entropy(df):
    target = df.key[0][-1]
    entropy = 0
    value = df[target].unique()
    for value in value:
        fraction = df[target].value_counts()[value] /
            (len(df[target]))
        entropy += fraction * np.log2(fraction)
    return entropy
```

```
def buildTree(df, tree=None):
    target = df.key[0][-1]
    node = find_winner(df)
    att = np.unique(df[target])
    if tree is None:
        tree = {}
        tree[node] = {}
    for value in att:
        sub = get_subsets(df, node, value)
        distro_counts = np.unique(sub).count(target)
        return_counts = True
    if len(return_counts) == 1:
        tree[node][value] = value[0]
    else:
```

```
tree = buildTree(df)
input = print
print.print(tree)
```

## Decision Tree

(2)

Decision tree (sklearn):

```
import pandas as pd
import numpy as np
import sklearn.model_selection
import DecisionTree
from sklearn.tree import plotree
```

```
df = pd.read_csv('url')
```

```
df.head()
```

```
df.info()
```

```
df.isnull().sum()
```

```
cols = df.columns[0:-1]
```

```
for i in cols:
```

```
    plt.subplot(y=df[i])
```

```
    plt.show()
```

```
X = df.drop('ipower', axis=1)
```

```
y = df['ipower']
```

```
dt = DecisionTreeClassifier(max_depth=3)
```

```
dt.fit(X, y)
```

```
y_pred_train = dt.predict(X_train)
```

```
y_pred = dt.predict(X_test)
```

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

*Handwritten signature and date*  
07.05.21



Logistic Regression -

```
import numpy as np
import pandas as pd
import sklearn.datasets import load_diabetes
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score,
classification_report, confusion_matrix
```

```
diabetes = load_diabetes()
```

```
X, y = diabetes.data, diabetes.target
```

```
y_binary = (y > np.median(y)), astype(int)
```

```
X_train, X_test, y_train, y_test = train_test_split(
    X, y_binary, test_size=0.2, random_state=42)
```

```
scaler = StandardScaler()
```

```
X_train = scaler.fit_transform(X_train)
```

```
X_test = scaler.transform(X_test)
```

```
model = LogisticRegression()
```

```
model.fit(X_train, y_train)
```

```
print("Accuracy: %.2f" % format(accuracy_score(y_test, y_pred)))
```

```
print("Confusion matrix:", confusion_matrix(y_test, y_pred))
```

```
print("Classification method:", classification_report(y_test, y_pred))
```

Confusion matrix

```
[[ 35 13]]
```

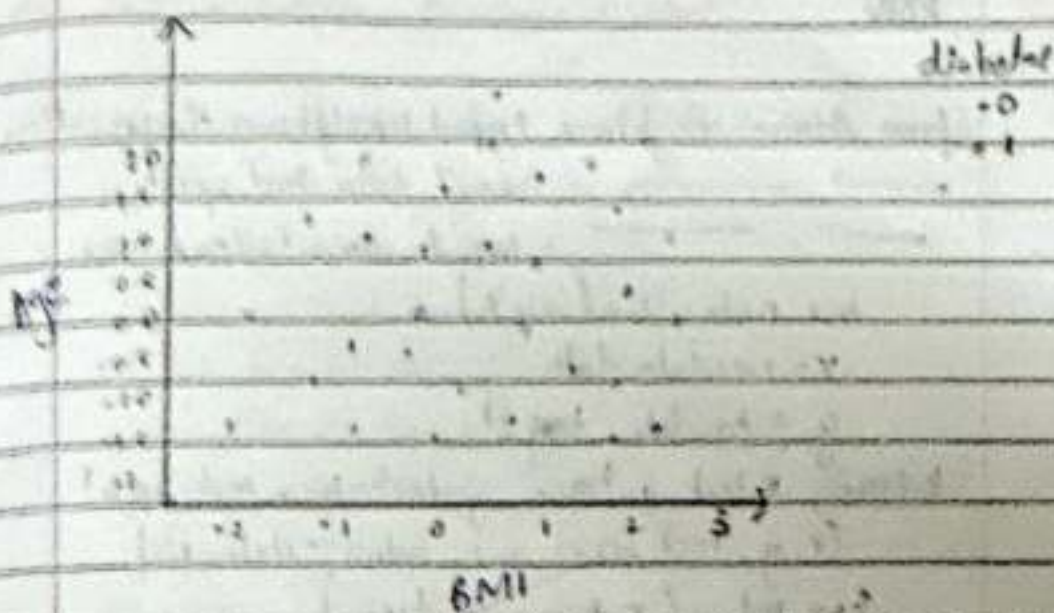
```
[[ 11 29]]
```

```
plt.scatter(X_test[2], y=X_test[8], label=y_test)
```

```
palette = {0: 'blue', 1: 'red', marker='o'}
```

```
plt.legend(title="Diabetes", loc='upper left')
```

```
plt.show()
```





## KNN

```
from sklearn.neighbors import KNeighborsClassification
import train_test_split
import load_iris

iris_data = load_iris()
X = iris_data.data
y = iris_data.target

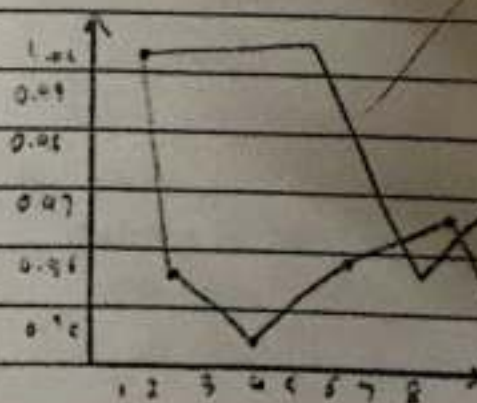
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42)

knn = KNeighborsClassifier()
knn.fit(X_train, y_train)
print(knn.predict(X_test))

neighbors = up_range(1, 9)
for i, k in enumerate(neighbors):
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    train_accuracy[i] = knn.score(X_train, y_train)
    test_accuracy[i] = knn.score(X_test, y_test)

plt.plot(neighbors, test_accuracy, label='Testing set Accuracy')
plt.plot(neighbors, train_accuracy, label='Training set Accuracy')

plt.legend()
plt.xlabel('n_neighbors')
plt.ylabel('Accuracy')
plt.show()
```



23/5/24

## \* Support Vector Machine (svm)

```

from sklearn.datasets import load_breast_cancer
import matplotlib.pyplot as plt
from sklearn.svm import SVC
cancer = load_breast_cancer()
x = cancer.data[:, :2]
y = cancer.target
svm = SVC(kernel="rbf", gamma=0.5, C=10)
svm.fit(x, y)
plt.scatter(x[:, 0], x[:, 1], c=y, s=20)
plt.show()

```

## \* PCA using sklearn -

```

from sklearn.decomposition import PCA

```

```

pca = PCA(n_components=2)

```

```

pca.fit(2)

```

```

x_pca = pca.transform(2)

```

```

df_pca = pd.DataFrame(x_pca, columns=['PC1',
                                     'PC2'])

```

```

print(df_pca)

```

output:

	PC1	PC2
0	9.1861	1.9468

## \* K-Means Clustering -

```

from sklearn.cluster import KMeans

```

```

data = list(zip(x, y))

```

```

init = []

```

```

for i in range(1, 11):

```

```

    kmeans = KMeans(n_clusters=2)

```

```

    kmeans.fit(data)

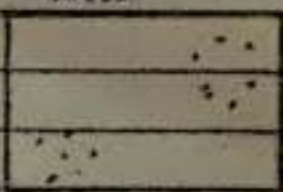
```

```

plt.show()

```

Result



x y 0 10 12



4 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 = load_breast_cancer()
data.keys()
```

```
print(data['feature_names'])
dpl = pd.DataFrame(data['data'], columns=data
```

```
scaling = StandardScaler()
```

```
scaling.fit(data)
```

```
principal = PCA(n_components=3)
```

```
X = principal.transform(scaled_data)
```

```
plt.figure(figsize=(12,10))
```

```
plt.scatter(X[:,0], X[:,1], c=data['target'])
```

```
plt.xlabel('PC1')
```

```
plt.ylabel('PC2')
```

## Random forest -

```
import pandas as np
import pandas as pd
from sklearn.ensemble import Random Forest
Classification
from sklearn import datasets
iris = datasets.load_iris
type(iris)
x = iris.data[:, 1:3].ravel()
y = iris.data[:, 4].ravel()
model = Random Forest Classifier()
model.fit(x_train, y_train)
y_pred = model.predict(x_test)
accuracy = score(y_test, y_pred)
```

## Ensemble boosting -

```
from sklearn import datasets, ensemble
X_train = 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 = AdaBoostClassifier(n_estimators=
    30, learning_rate=1)
adaBoost.fit(x_train, y_train)
y_pred = adaBoost.predict(x_test)
accuracy = score(y_test, y_pred)
```

*Don*  
30.05.21