

Exploratory data Analysis Lab:

Basic Python packages:

Numpy:

```
import numpy as np
array = np.array([1, 2, 3, 4])
array_sum = np.sum(array)
array_mean = np.mean(array)
print("Array : ", array)
print("Sum : ", array_sum)
print("Mean : ", array_mean)
```

```
matrix = np.array([[1, 3], [3, 4]])
transpose = matrix.T
determinant = np.linalg.det(matrix)
print("Matrix : \n", matrix)
print("Transpose : \n", transpose)
print("Determinant : \n", determinant)
```

Pandas:

```
import pandas as pd
data = {'name': ['Alice', 'Bob', 'Charlie'], 'Age': [20, 32, 34]}
df = pd.DataFrame(data)
print(df)
print(df['name'])
print(df.describe())
```

④ data = { 'Name': ['Alice', 'Bob', 'Charlie'], 'Age': [24, 27, 22],
 'Score': [85, 83, 79] }
df = pd.DataFrame(data)
df['pass'] = df['Score'] > 80
print("DataFrame:\n", df)
print("Average Age:", df['Age'].mean())

Matplotlib:

⑤ import numpy as np
import matplotlib.pyplot as plt
y1 = [2, 1, 4.5]
y2 = [1, 1.5, 4]
plt.plot(y1)
plt.plot(y2)
plt.legend(["blue", "green"], loc="lower right")
plt.show()

⑥ import matplotlib.pyplot as plt
x = [0, 1, 2, 4, 7]
y = [0, 1, 9, 18]
plt.plot(x, y, marker='o', label='y=x^2')
plt.title("Simple plot")
plt.xlabel("x-axis")
plt.ylabel("y-axis")
plt.legend()
plt.grid()
plt.show()

Seaborn:

```
import seaborn as sns  
import matplotlib.pyplot as plt  
tips = sns.load_dataset("tips")  
sns.scatterplot(data=tips, x="total_bill", y="tip",  
                 hue="sex")  
plt.title("Tips dataset scatter plot")  
plt.show()
```

import seaborn as sns

```
import seaborn as sns  
import matplotlib.pyplot as plt  
tips = sns.load_dataset("tips")  
sns.boxplot(x="day", y="total_bill", data=tips)  
plt.title("Total bill distribution by day")  
plt.show()
```

Scikit-learn:

```
from sklearn.datasets import load_iris
```

```
iris = load_iris()
```

```
x = iris.data
```

```
y = iris.target
```

```
feature_names = iris.feature_names
```

```
target_name = iris.target_names
```

```
print("feature names:", feature_names)
```

```
print("target names:", target_names)
```

```
print("InType of x is:", type(x))
```

```
print("In first 5 rows of x:\n", x[:5])
```

```
⑥ from sklearn.datasets import load_iris  
from sklearn.model_selection import train_test_split  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.metrics import accuracy_score  
  
iris = load_iris()  
x,y = iris.data, iris.target  
x_train, x_test, y_train, y_test = train_test_split  
(x,y, test_size=0.3)  
  
clf = DecisionTreeClassifier()  
clf.fit(x_train, y_train)  
y_pred = clf.predict(x_test)  
accuracy = accuracy_score(y_test, y_pred)  
print("Accuracy of the Decision Tree classifier:",  
      accuracy)
```

```
① import numpy as np  
import matplotlib.pyplot as plt  
import pandas as pd  
dataset = pd.read_csv("C:/Users/android/Desktop/  
diabetes.csv")  
print(dataset)
```

④ dataset.sum()

`dataset.head()`

② dataset.tail()

0008395

卷之三

④ dataset.describe()

④ dataset.info()

④ # Preprocessing tech.

data Cleaning

```
import pandas as pd
```

```
import pandas as pd  
from sklearn import import SimpleImputer
```

```
data = pd.DataFrame ({  
    'name': ['Vanya', 'Prithvi', 'Siri', 'Vanya', None],  
    'age': [21, 34, None, 21, 22],  
    'purchase amount': [100.5, None, None, 100.5, 50.5],  
    'Date of Purchase': ['2023/12/01', '2023/12/02', '2025/12/01',  
    '2023/12/01', '2023/12/03']})
```

importer = Simple Importer (strategy = 'mean')

```
imputer = SimpleImputer(strategy='ffill')
data[['age','purchase-amount']] = imputer.fit_transform
(data[['date-of-purchase']], error_if='raise')
```

~~Print (order)~~

```
# data integration
import pandas as pd
data1 = pd.DataFrame({
    'cust_id': [1, 2, 3],
    'name': ['Sam', 'Tom', 'Ravi'],
    'age': [28, 34, 29]
})
data2 = pd.DataFrame({
    'cust_id': [1, 3, 4],
    'purchase_amount': [100.5, 85.3, 45.0],
    'purchase_date': ['2023-01-01', '2023-02-02', '2023-12-03']
})
```

```
merged_data = pd.merge(data1, data2, on='cust_id', how='inner')
print(merged_data)
```

```
merged_data = pd.merge(data1, data2, on='cust_id', how='outer')
print(merged_data)
```

```
merged_data = pd.merge(data1, data2, on='cust_id', how='right')
print(merged_data)
```

```
merged_data = pd.merge(data1, data2, on='cust_id', how='left')
print(merged_data)
```

⑩ # data transformation

```
from sklearn.preprocessing import StandardScaler, OneHotEncoder
data = pd.DataFrame({
    'category': ['A', 'B', 'A', 'C', 'B'],
    'numerical_column': [10, 15, 10, 20, 15]
})
```

```
scaler = StandardScaler()
data['scaled_numeric_column'] = scaler.fit_transform(data[['category']])
```

```
columns = encoder.get_feature_names_out(['category']))
```

```
data = pd.concat([data, encoded_data], axis=1)
```

```
print(data)
```

```
#data reduction
from sklearn.decomposition import PCA
from sklearn.feature_selection import SelectKBest, chi2
data = pd.DataFrame({ 'feature1': [10, 20, 30, 40, 50],
                      'feature2': [1, 2, 3, 4, 5],
                      'feature3': [100, 200, 300, 400, 500],
                      'target' : [0, 1, 0, 1, 0] })
selector = SelectKBest(chi2, k=2)
selected_features = selector.fit_transform(data[ ['feature1',
                                                 'feature2', 'feature3'] ], data['target'])
print("Selected features (Select K Best):")
print(selected_features)
pca = PCA(n_components=2)
pca_data = pca.fit_transform(data[ ['feature1', 'feature2',
                                    'feature3'] ])
print("PCA reduced data:")
print(pca_data)
```

Data Visualization: Univariate Analysis

① Count plot:

```
import numpy as np  
import pandas as pd  
import seaborn as sns  
import matplotlib.pyplot as plt  
  
plt.figure(figsize=(15,8))  
sns.countplot(x='Department', data=df)  
plt.title("Count plot of department")  
plt.show()
```

② Pie chart:

```
plt.figure(figsize=(8,8))  
pdf['department'].value = count().plot.pie  
plt.title("Department Distribution")  
plt.show()
```

③ Histogram:

```
plt.figure(figsize=(10,5))  
sns.histplot(df['Salary'], bins=10, color='blue')  
plt.title("Histogram of salary")  
plt.show()
```

Histogram with distplot:

```
sns.histplot(df['Salary'], bins=10, kde=True, color='blue')
plt.title("Histogram of salary with KDE (distplot)")
plt.show()
```

Box Plot:

```
plt.figure(figsize(8,6))
sns.boxplot(y=df['Salary'])
plt.title("Salary Boxplot")
plt.show()
```

Bivariate & Multivariate Analysis:

Scatter plot

```
plt.figure(figsize(8,6))
sns.scatterplot(x="Experience", y='Salary', hue='Department',
                 data=df)
plt.title("Experience vs Salary")
plt.show()
```

Pair plot:

```
sns.pairplot(df, hue = 'Department')  
plt.show()
```

(8)

Correlation Heatmap:

```
plt.figure(figsize=(8,6))  
sns.heatmap(df, select_dtypes(include=['number']),  
            cbar=True, annot=True, cmap='coolwarm')  
plt.title("Correlation heatmap")  
plt.show()
```

Barplot:

```

plt.figure(figsize=(10,5))
sns.barplot(x='department', y='Salary', data=df, estimator=
             np.mean)
plt.title("Average salary by department")
plt.show()

```

Boxplot:

```

plt.figure(figsize=(10,6))
sns.boxplot(x='department', y='Salary')
plt.title('Salary distribution Across')
plt.show()

```

Time Series Analysis: Line Plot

```

df_sorted = df.sort_values('joining year')
plt.figure(figsize=(10,5))
sns.lineplot(x='joining year', y='Salary', data=df_sorted,
             marker='o')
plt.title("Salary trend over joining year")
plt.show()

```

Data Aggregation, Transformation, plot, table & charts

① import pandas as pd

```
data = {'Name': ['Tom', 'Nick', 'Krish', 'Jack'],  
       'Age': [20, 21, 19, 18]}
```

```
df = pd.DataFrame(data)
```

```
print(df)
```

② import pandas as pd

```
data = {'Name': ['Jai', 'Princi', 'Gaurav', 'Anuj'],  
       'Age': [27, 24, 28, 29],  
       'Qualification': ['MSc', 'MCA', 'MCA', 'Phd'],  
       'Address': ['Delhi', 'Kempur', 'Allahabad', 'Kannay']}
```

```
df = pd.DataFrame(data)
```

```
print(df[['Name', 'Qualification']])
```

③ df = pd.DataFrame({'Column1': ['A', 'B', 'C', 'A', 'C',
 'C', 'B', 'D', 'D', 'A'],
 'Column2': [5, 10, 15, 20, 25, 30, 35, 40, 45, 50]})

```
df = df.groupby('Column1')[['Column2']].apply(list)
```

```
df
```

④ df = pd.DataFrame({'animal': ['Falcon', 'Falcon', 'Parrot', 'Parrot',
 'maxspeed': [380, 370, 240, 260]})

```
df.groupby(['animal']).mean()
```

⑤

```
import numpy as np
import pandas as pd
df = pd.DataFrame({'key1': ['a', 'a', 'b', 'b', 'a'],
                   'key2': ['one', 'two', 'one', 'two', 'one'],
                   'date1': np.random.rand(5),
                   'date2': np.random.rand(5)})
df.groupby(['key1', 'key2']).size()
```

⑥ #grouping index levels

```
columns = pd.MultiIndex.from_arrays([['US', 'US', 'US', 'JP', 'JP'], [1, 3, 5, 1, 3]], names=['city', 'febor'])
hier_df = pd.DataFrame(np.random.rand(4, 5), columns=columns)
hier_df
```

hier-df

```
hier-df.groupby(level='city', axis=1).count()
```

⑦ #data Integration

```
df = pd.DataFrame({'key1': ['a', 'a', 'b', 'b', 'a'],
                   'key2': ['one', 'two', 'one', 'two', 'one'],
                   'date1': np.random.rand(5),
                   'date2': np.random.rand(5)})
```

```
df.describe()
```

Aggregation functions

df = pd.DataFrame([[1, 2, 3], [4, 5, 6], [7, 8, 9], [np.nan],
[np.nan, np.nan, np.nan]]).

columns = ['A', 'B', 'C']

df.agg(['sum', 'min'], 'B': {'min', 'max'})

aggregate over the columns

df = pd.DataFrame({'A': range(3), 'B': range(1, 4)})

df.

df = pd.DataFrame({'A': range(3), 'B': range(0, 4)})

df

22	0
201	1
1	2
121	3
23	A
201	B
02	C
121	F

S = pd.Series(range(3))

s = transform([np.sqrt, np.exp])

② df = pd.DataFrame({
 'Date': ["2015-05-08", "2017-05-07", "2015-05-06",
 "2015-05-05", "2015-05-08", "2017-05-07",
 "2017-05-06", "2015-05-05"],
 'Data': [5, 8, 6, 1, 50, 100, 60, 120]})
df

③ df.groupby('Date')[['Data']].transform(sum)

import seaborn as sns
tips = sns.load_dataset('tips')
tips

import seaborn as sns

tips = sns.load_dataset("tips")

tips = pd.pivot_table(tips, index="sex", values=[1 + tips])

pd.pivot_table(tips, index="sex", values=['tip'], aggfunc='sum')

pd.pivot_table(tips, index="Sex", values=['tip'], columns=['time'])

pd.pivot_table(tips, index=["Sex"], values=['total_bill', "tip"], columns=['time', "smoker"])

pd.crosstab(tips["sex"], tips["time"], normalize='columns')

pd.crosstab(grip['Sex'], tips['time'], normalize='columns')

```
pd.crosstab(tips["sex"], tips["time"], normalize="all")
```

```
pd.crosstab(tips["Sex"], tips["Time"], value="tips",  
            aggfunc="mean")
```

```
pd.crosstab(tips[["sex"]], tips[["time"]], values=tips[["tip"]],  
aggfunc="mean", rownames=[["Gender"]],  
columns=[["Time"]])
```

pd. crosstab(tips[["sex"]], tips[["time"]], tips[["smoker"]])

1920 February 2000 lot of 3M Skins and 9100
02 27 226 38 PP skins 0
18 27 226 44 PP neck 1

⑧ pd.crosstab(tabs["sex"], tips[["day"]], [tips[["time"]],
tips[["smoker"]]])

④ pd.crosstab(tips["day"], tips["sex"], normalize="index")
).plot.bar(stacked=False)

managed room, I didn't

⑤ import pandas as pd
student = ({'Name': ['Vaish', 'Sham', 'Anj', 'Anu'],
'math': [99, 88, 76, 63],
'ML': [88, 77, 89, 89],
'total_score': [155, 233, 266, 224],
'grading': ['A', 'B', 'A', 'B'],
'Age': [20, 21, 20, 20]})

df = pd.DataFrame(student)

df

print(df[['name', 'grading']])

~~print(df.groupby('name')[['total_score']].apply(list))~~

Linear Regression:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge, Lasso, ElasticNet
from sklearn.preprocessing import PolynomialFeatures
from sklearn.metrics import mean_squared_error
np.random.seed(42)
x = 2 * np.random.rand(100, 1)
y = 4 + 3 * x + np.random.randn(100, 1)

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)

lin_reg = LinearRegression()
lin_reg.fit(x_train, y_train)
y_pred = lin_reg.predict(x_test)

print("Single Linear Regression MSE:", mean_squared_error(y_test, y_pred))

plt.figure(figsize=(8, 6))
plt.scatter(x, y, color='blue', label='Actual data')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.show()
```

```
x-mult = np.c_[x, x**2]
x-train-m, x-test-m, y-train-m, y-test-m = train-
test -> split(x-mult, y, test_size=0.2,
random_state=42)

lin-reg-multi = LinearRegression()
lin-reg-multi.fit(x-train-m, y-train-m)
y-pred-multi = lin-reg-multi.predict(x-test-m)
print("Multiple Linear Regression MSE", mean-
squared_error(y-test-m, y-pred-multi))

plt.figure(figsize=(8, 6))
plt.scatter(x, y, color='blue', label='Actual Data')
plt.plot(x-test-m, y-pred-multi, color='green', label=
'Multiple Linear Regression')

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

poly-features = PolynomialFeatures(degree=2,
include_bias=False)
x-poly = poly-features.fit_transform(x)

x-train-p, x-test-p, y-train-p, y-test-p =
train-test-split(x-poly, y-test-size=0.2, random_
state=42)

poly-reg = LinearRegression()
poly-reg.fit(x-train-p, y-train-p)
y-pred-poly = poly-reg.predict(x-test-p)
print("Polynomial Regression MSE", mean-squared-
error(y-test-p, y-pred-poly))
```

```
plt.figure(figsize=(8,6))
```

```
plt.scatter(x,y, color='blue', label='Actual Data')  
plt.plot(x-test-p[:,0], y-pred-poly, color='orange',  
label='Polynomial Regression')
```

```
plt.xlabel('x')
```

```
plt.ylabel('y')
```

```
plt.legend()
```

```
plt.show()
```

```
Ridge-reg = Ridge(alpha=1.0)
```

```
ridge-reg-fit(x-train-p, y-train-p)
```

```
y-pred-ridge = ridge-reg.predict(x-test-p)
```

```
print("Ridge Regression MSE", mean_squared_error  
      (y-test-p, y-pred-ridge))
```

```
plt.figure(figsize=(8,6))
```

```
plt.scatter(x,y, color='blue', label='Actual Data')
```

```
plt.plot(x-test-p[:,0], y-pred-ridge, color='purple',  
label='Ridge-Regression')
```

```
plt.xlabel('x')
```

```
plt.ylabel('y')
```

```
plt.legend()
```

```
plt.show()
```

```
Lasso-reg = Lasso(alpha=0.1)
```

```
lasso-reg-fit = (x-train-p, y-train-p)
```

```
p-pred-lasso = lasso-reg.predict(x-test-p)
```

```
print("Lasso Regression MSE", mean_squared_error  
      (y-test-p, y-pred-lasso))
```

```
plt.figure(figsize=(8,6))
```

```
plt.scatter(x,y, color='blue', label='Actual Data')
plt.plot(x-testp[0], y-pred-lans0, color='brown',
         label='Lasso Regression')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.show()
```

```
Elastic-net = ElasticNet( alpha=0.1, l1_ratio=0.5)
elastic-net.fit(x-train-p, y-train-p)
y-pred-elastic = elastic-net.predict(x-test-p)
print("Elastic Net Regression MSE:", mean_squared_error(y-testp, y-predelastic))
```

```
plt.figure(figsize=(8,8))
plt.scatter(x,y, color='blue', label='Actual Data')
plt.plot(x-testp[0], y-pred-elastic, color='pink',
         label='Elastic Net Regression')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.show()
```

Logistic Regression:

27/2/25

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.datasets import load_iris
data = load_iris()
X, y = data.data, data.target
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

Binary Logistic

```
binary_y_train = (y_train == 0).astype(int)
binary_y_test = (y_test == 0).astype(int)
log_reg_binary = LogisticRegression()
log_reg_binary.fit(X_train, binary_y_train)
preds_binary = log_reg_binary.predict(X_test)
print("Binary Logistic Regression accuracy: ", accuracy_score(binary_y_test, preds_binary))
```

Multiclass logistic

```
log_reg_over = LogisticRegression(multi_class='ovr',
                                    solver='lbfgs', max_iter=200)
log_reg_over.fit(X_train, y_train)
preds_over = log_reg_over.predict(X_test)
print("One-vs-Rest Logistic Regression accuracy: ", accuracy_score(y_test, preds_over))
```

```
# Logistic regression w/ softmax  
log-reg-softmax = LogisticRegression(multi_class='ovr', solver='lbfgs', max_iter=200)  
log-reg-softmax.fit(x-train, y-train)  
preds-softmax = log-reg-softmax.predict(x-test)  
print("Softmax Logistic Regression accuracy:",  
      accuracy_score(y-test, preds-softmax))
```

```
# L1, L2, Elastic Net.  
log-reg-l1 = LogisticRegression(penalty='l1', solver='liblinear', max_iter=200)  
log-reg-l1.fit(x-train, y-train)  
preds-l1 = log-reg-l1.predict(x-test)          preds-l1  
print("L1 logistic accuracy:", accuracy_score(y-test, ^  
      preds-l1))
```

```
log-reg-l2 = LogisticRegression(penalty='l2', max_iter=200)  
log-reg-l2.fit(x-train, y-train)  
preds-l2 = log-reg-l2.predict(x-test)  
print("L2 logistic accuracy:", accuracy_score(y-test, preds-l2))  
log-reg-elastic = LogisticRegression(penalty='elasticnet',  
                                     solver='saga', l1_ratio=0.5, max_iter=200)
```

```
log-reg-elastic.fit(x-train, y-train)  
preds-elastic = log-reg-elastic.predict(x-test)
```

~~Elastic Net~~
~~accuracy score~~
print("Elastic Net logistic Regression Accuracy",
 accuracy_score(y-test, preds_elastic))

~~bar plot to compare~~
regression-types = ['Multiclass', 'multiclass with softmax',
 'regularization', 'L1', 'L2', 'elasticNet']
accuracy = [0.75, 0.88, 0.91, 0.94, 0.96]
plt.bar(regression-types, accuracy, color='skyblue')
plt.title('Regression types vs accuracies')
plt.show()

K-Nearest Neighbor Classifier:

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, classification_report
from sklearn.datasets import make_classification.
```

```
X, y = make_classification(n_samples=300, n_features=5,
                           n_classes=2, random_state=42)
x_train, x_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.2, random_state=42)
```

```
scaler = StandardScaler()
```

```
x_train = scaler.fit_transform(x_train)
```

```
x_test = scaler.transform(x_test)
```

```
k = 7
```

```
knn = KNeighborsClassifier(n_neighbors=k)
```

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

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

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

```
print(f"accuracy: {accuracy:.2f}%")
```

```
print("Classification Report")
```

```
print(classification_report(y_test, y_pred))
```

Decision Tree:

```
from sklearn.datasets import load_iris  
from sklearn.tree import DecisionTreeClassifier,  
                         plot_tree  
from sklearn.model_selection import train_test_split  
import matplotlib.pyplot as plt  
  
X,y = iris.data,iris.target  
  
X-train,X-test,y-train,y-test = train_test_split(  
    X,y,test_size=0.2,random_state=42)  
  
clf = DecisionClassifier(criterion='gini',max_depth=3,  
                         random_state=42)  
  
clf.fit(X-train,y-train)  
  
y-pred = clf.predict(X-test)  
  
accuracy = accuracy_score(y-test,y-pred)  
print(f"Accuracy: {accuracy:.2f}")  
  
plt.figure(figsize=(12,8))  
  
plot_tree(clf, feature_names=iris.feature_names,  
          class_names=iris.target_names, filled=True)  
plt.title("Decision tree Visualization")  
plt.show()
```

SVM :

```
from sklearn import svm, dataset
import sklearn.model_selection as model_selection
from sklearn.metrics import accuracy_score
from sklearn.metrics import f1_score

iris = load_iris()
X = iris.data[:, :2]
y = iris.target

X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, train_size=0.80, test_size=0.2, random_state=101)

rbf = svm.SVC(kernel='rbf', gamma=0.5, C=0.1).fit(X_train, y_train)

poly = svm.SVC(kernel='poly', degree=3, C=1).fit(X_train, y_train)

poly_pred = poly.predict(X_test)
rbf_pred = rbf.predict(X_test)

poly_accuracy = accuracy_score(y_test, poly_pred)
poly_f1 = f1_score(y_test, poly_pred, average='weighted')
print("Accuracy (Polynomial kernel):", "%.2f" % (poly_accuracy * 100))

rbf_accuracy = accuracy_score(y_test, rbf_pred)
rbf_f1 = f1_score(y_test, rbf_pred, average='weighted')
print("Accuracy (RBF kernel):", "%.2f" % (rbf_accuracy * 100))

print("F1 (RBF kernel):", "%.2f" % (rbf_f1 * 100))
```