

## Lab-01

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

→ Pandas read csv()

```
import pandas as pd.
```

```
airbnb_data = pd.read_csv('austin_flowing_data.csv')
```

```
airbnb_data.head()
```

# Exporting dataframe to a csv file.

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

# importing dataset using url

```
url = "https://andvine.com/ml/machine-learning-dataset/iris_data"
```

N/D  
3/5/24

## Lab 02

Demonstrate various data preprocessing techniques

### Algorithm

- 1) Downloading the dataset.
- 2) Reading csv file  
`pd.read_csv(data-path)`
- 3) Observing dataset using `head()`, `info()` & `describe()`
- 4) Visualization of dataset using matplotlib lib.
- 5) Create test set & train set `split(housing, test_size=0.2, random_state=42)`
- 6) Visualise the data to gain insights.
- 7) Find the correlation bet<sup>n</sup> the categories
- 8) Data cleaning by dropping ~~NA~~ values  
`housing.dropna()`
- 9) Imputation of missing values
- 10) Encode categorical values with `one-hot`
- 11) Scaling data using standardization or min-max strategy
- 12) Training the linear regression model
- 13) Calculating root mean square error
- 14) Training decision tree
- 15) Cross validation using mean & standard deviation
- 16) Fitting test data & calculating accuracy.



- Q. Implement linear regression algorithm using appropriate dataset.

### Algorithm

1. Import necessary libraries.
2. Import dataset
3. Visualization of dataset using diff plots like heatmap, distribution plot, scatterplot etc.
4. preprocess the data, convert or encode categorical data
5. Split the dataset into training set & test set.

~~from sklearn.model\_selection import train\_test\_split~~

~~X\_train, X\_test, y\_train, y\_test = train\_test\_split~~

~~X, y, test\_size=0.3, random\_state=23)~~

6. Build model

from sklearn.linear\_model import LinearRegression

lin\_reg = LinearRegression()

7. Fit the dataset to model & train it  
lin\_reg.fit(X\_train, y\_train)

8. Calculate the accuracy using mean square error

Q. Implement multilinear regression.

1) Import necessary libraries like LinearRegression

2) Import dataset

3) Visualize the dataset using matplotlib.  
pyplot.

4) Encode categorical data

```
ct = ColumnTransformer(transformers =  
    [['encode', OneHotEncoder [33]],  
    remainder = 'pass through']
```

5) Split dataset to training & test dataset.

6) We can see multiple independent variables.

7) Create regressor model.

```
regressor = LinearRegression()
```

8) Fit the train set.

9) Test the model using test set.

10) Compare the actual value & predicted value.

21/5/24



Q) Use appropriate dataset for building the decision tree ID3 & apply the knowledge to classify new sample.

- 1) Import the libraries like DecisionTreeClassifier.
- 2) Import diabetes dataset.
- 3) Define calculate-entropy function to calculate entropy for each value.

$v\_count = \text{no. of values}$

$proportion = v\_count / \text{total rows}$

$entropy = proportion + \text{math.log2}(proportion)$

- 4) Define information-gain function

$information\_gain = entropy\_outcome \rightarrow$

~~$weighted\_entropy$~~

- 5) Access the feature with highest information gain.

- 6) Create a decision tree using DecisionTree Classifier function.

- 7) Plot the tree.

- 8) ~~or~~ Define ID3 function & pass new sample as parameter & obtain classification.

Q. Build logistic regression model for given dataset.

- 1) Import all required libraries.
- 2) Import the given dataset
- 3) Preprocess the dataset to standard state.
- 4) Split the dataset into test & train.
- 5) Code the logistic Regression model.

LR = LogisticRegression(C=0.01, solver='liblinear')  
fit(X\_train, y\_train)

- 6) Predict the test set using the model.

yhat = LR.predict(X\_test)

yhat\_probab = LR.predict\_proba(X\_test)

- 7) Calculate the performance or accuracy of the models.

Ans  
solution

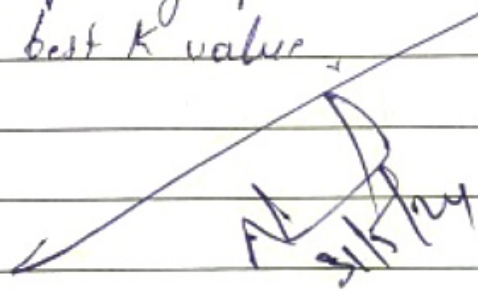


Q. Build KNN classification model for a given dataset

- 1) Import the necessary libraries
- 2) Import diabetes dataset.
- 3) Preprocess the dataset.
- 4) Split the dataset as test & train set
- 5) Build the KNN model.

$KNN = K\text{neighbors Classifier (n-neighbours = K)}$

- 6) Fit the model & test the accuracy.
- 7) We can plot graph of  $K$ 's accuracy to choose the best  $K$  value.



31-5-24

PAGE NO.

DATE / /

## 7) Build SVM model for a given dataset

from sklearn import datasets

cancer = datasets.load\_breast\_cancer()

print("Features:", cancer.feature\_names)

print("Labels:", cancer.target\_names)

cancer.data.shape

print(cancer.data[0:5])

print(cancer.target)

from sklearn.model\_selection import train\_test\_split

x\_train, x\_test, y\_train, y\_test = train\_test\_split(  
cancer.data, cancer.target, test\_size=0.3,  
random\_state=109)

from sklearn import svm

clf = svm.SVC(kernel='linear')

clf.fit(x\_train, y\_train)

y\_pred = clf.predict(x\_test)

from sklearn import metrics

print("Accuracy:", metrics.accuracy\_score(  
y\_test, y\_pred))

print("Precision", metrics.precision\_score(  
x\_test, y\_pred))

print("Recall", metrics.recall\_score(y\_test, y\_pred))



## 8) Random Forest algorithm

- Data pre-processing step
- fitting the random forest algorithm to training set
- predicting the test result
- test accuracy of the result (creation of confusion matrix)
- Visualizing the test set result

Step 1:- Select random  $k$  data points from the training set

Step 2: Build the decision trees associated with the selected data points.

Step 3: Choose the no.  $N$  for decision trees that you want to build

Step 4: Repeat step 1 & 2

Step 5: For new datapoints, find the predictions of each decision tree & assign new datapoints to category that wins the majority votes

AL  
3/16/24

1) Implement boosting ensemble method on given dataset

1) → Initialize weights :- assign equal weights to all training samples

2) → For each base learner (weak learner):  
Train a weak learner on training data with current weights.

3) → Update sample weights :- Increase weights of the incorrectly classified samples so that they are more likely to be selected in next iteration

→ Repeat steps 2 & 3 for a fixed no. of iterations or until a stopping criterion is met

→ Aggregate the predictions

→ output the ensemble model.



- 10) Build K-means algorithm to cluster a set of data stored in a CSV file.

function kMeans(X, K, max\_iterations):

- 1) Initialize centroids randomly
- 2) Repeat until convergence or max\_iterations:
  - a. assign each data point to nearest centroid
  - b. Update centroids based on the mean of data points in each cluster
  - c. check convergence
- 3) Return the final centroids.

- 11) Dimensionality reduction using PCA

1) Standardize data :- compute mean  $\mu$  & standard deviation  $\sigma$  for each feature.

- 2) Compute the covariance matrix

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

- 3) Compute eigenvectors & eigenvalues:

- 4) Select the top k eigenvectors:

select the top k eigenvectors corresponding to largest eigenvalues to form projection matrix W

- 5) Project the data onto the new feature space

$$X_{pca} = X \cdot W$$

- 6) output :- return the transformed dataset  $X_{pca}$ .