

Write a Python program to import and export data using Pandas library functions.

Importing

```
import pandas as pd  
df = pd.read_csv("/content/laustinhousingdata.csv")  
df.head()
```

Exporting

```
url = "archive.ics.uci.edu/ml/machine-learning-databases"  
colnames = ("sepal-length", "sepal-width", "petal-length", "petal-width", "class")  
iris_data = pd.read_csv(url, names=colnames)  
iris_data.head()  
iris_data.to_csv("exported.csv")
```

## 2, Data preprocessing

- i) Import dataset using pandas
- ii) perform dataset shape & analysis
- iii) use `is-null()` function from pandas to analyse missing values.
- iv) Drop or fill missing values according to your use case example `dropna()` and `fill()`.
- v) we can generate dummy variables. (is a binary variable that indicates whether a separate categorical variable that indicates a specific values)



12/1/24 Lab 2:

## Decision Tree Algorithm

- create a Root node for the tree.
- If all examples are positive, return the single node tree Root, with label = +.
- If all examples are negative, return the single node tree Root, with label = -.
- If Attributes is empty, return the single node tree Root, with label = most common value of Target attribute in Examples.
- otherwise Begin
  - $A \leftarrow$  the attribute from attribute that ~~best~~<sup>best</sup> classifies examples
  - The decision attribute for Root  $\leftarrow A$
  - For each possible value,  $v_i$  of  $A$ .
    - Add a new branch below Root, corresponding to the test  $A = v_i$ .
    - let examples  $v_i$ , be the subset of Examples that have value  $v_i$  for  $A$ .
  - If examples  $v_i$  is empty
    - Then below this new branch add a leaf node with label = most common value of Target attribute in Examples.

- Else below this new branch add the subtree ID3 (Example  $v_i$ , Target attribute, Attributes  $\subseteq \{A\}$ ).

- End

- Return Root

\* The best attribute is the one with highest information gain.

output:

Entropy of the entire dataset = 0.94

① Highest information gain = outlook = 0.246.

② Highest information gain = rainy = 0.77

But attribute is windy.



lab 3:

## KNN Algorithm:

- 1) Input: unlabeled instance  $x_{test}$  to which the class label needs to be predicted.
- 2) calculate Euclidean distance.
- 3) Find nearest neighbour; select  $K$  instance with the smallest distance to  $x_{test}$ .
- 4) For classification:
  - Majority vote: If its classification test, count the occurrence of each class label among the  $K$ -nearest neighbour.
  - Assign label: Assign the class label with the highest count of the predicted label for  $x_{test}$ .

## OUTPUT:

```
model.predict([(7.7, 2.6, 8.9, 2.3)])  
- array(['Fris Virginia'], dtype=object)
```

---

# Linear regression:

## Algorithm:

### Training Phase:

1) Input: Training data with one predictor ( $x$ ) & target variable ( $y$ )

2) compute means

3) Estimate coefficients:

$$\beta_1 = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

### Prediction Phase:

1) Input: new  $x_{\text{test}}$

2) compute prediction: use  $\beta_0$  &  $\beta_1$  &  $x_{\text{test}}$  to predict  $y_{\text{pred}}$ .

### OUTPUT:

~~model.Predict = ((2))~~

~~- Array ((107.82727115))~~



# Logistic regression:

## Algorithm:

### Training Phase:

1) Initialize parameter: start with random or zero values for weights (a) & bias (b).

2) compute predictions: calculate predictions using the equation:

$$\hat{y} = a(\theta^T x + b), \text{ where } a \text{ is the Sigmoid function.}$$

3) compute loss: measure the error between predictions & actual labels using binary cross-entropy loss.

4) update parameter: adjust weights & bias using gradient descent to initialize.

5) repeat: Iterate step 2-4 until converge or a max no. of iterations.

Prediction Phase

880  
3/5/24

## Week-4

### Support Vector Machine:

#### Algorithm

1) Define Kernel functions

$$\text{Ex: } K(x_1, x_2) = x_1 \cdot x_2$$

2) solve the quadratic programming problem to find the  $\alpha$  value.

3) complete weight and bias.

4) Identify the support vectors.

5) Make predictions.

#### output:

→ Model = sum(x)

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

predictions = model.predict(x\_test)

accuracy(y\_test, predictions)

0.962300

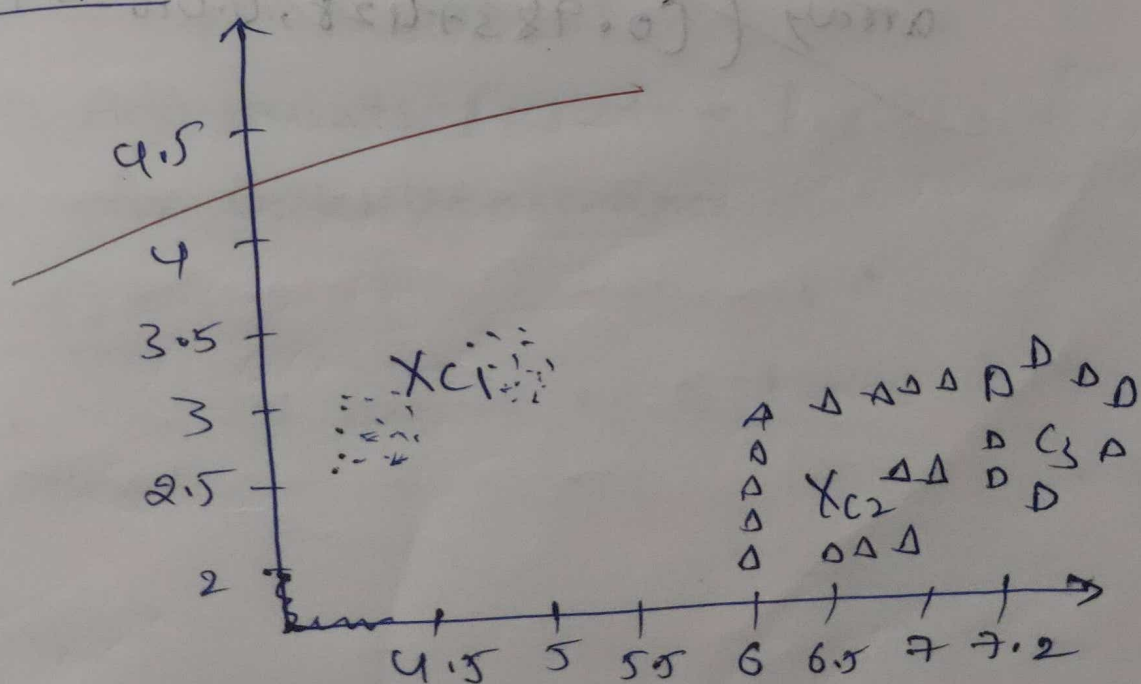
→ model.predict([-0.47069, -0.1604...  
0.17695]) array(0)



## K Means clustering Algorithm:

1. Select the number  $K$  to decide the no. of clusters.
2. Select random  $K$  points or centroids.
3. Assign each point to the closest centroid, which will form the predefined  $K$  cluster.
4. Calculate the variance and place a new centroid of each cluster.
5. Repeat step 3, reassign the centroid.
6. If any reassignment go to step-4 else go to finish.
7. Model is Ready.

Output:



### 3) Principal Component Analysis

#### Algorithm

- 1) calculate Mean.
- 2) calculation of covariance matrix.
- 3) Eigen values of covariance matrix.
- 4) computation of the Eigen Vector - with eigen vectors.
- 5) computation of first principal components.
- 6) Geometric meaning of first principal components.

#### Output:

~~pca~~ -

pca.explained\_variance\_ratio.

array (0.9837428, 0.01620478)



Lab 5:

## ANN with Backpropagation:

### Algorithm:

- 1) Initialization: Randomly set weights and biases.
- 2) Forward pass: compute output using current weights and biases.
- 3) Calculate Error: Find the difference between predicted and actual output.
- 4) Backward pass: update weights and biases to minimize error using gradient descent.
- 5) Repeat: Iterate steps 2-4 for multiple epochs.
- 6) Evaluation: Assess performance on a separate dataset.
- 7) Adjustments: Fine-tune hyperparameters for better performance.

### Output:

100 correct answers while testing : 8/9  
(Accuracy = 0.88)

# Random forest:

## Algorithm:

### 1) Initialization:

- choose number of trees and random feature subsets.

### 2) Tree Construction:

#### For each tree:

- Randomly select data with replacement and features.
- Grow tree recursively with random feature selection at each split.

### 3) Prediction:

- Aggregate predictions from all trees.

### 4) Optional:

- Calculate feature importance.
- Fine-tune hyperparameters.

### 5) Evaluation:

- Assess performance on validation data.



output:

accuracy = 0.93

### Adaboost Algorithm:

#### 1) Initialization:

- Initialize weights for each training sample uniformly
- choose the base learner (ex: decisionTree)

#### 2) for each iteration $t$ :

- Train the base learner on the training data with current sample weights.
- compute the weighted error of the base learner.
- compute the weight of the base learner in the final ensemble.

#### 3) update sample weights:

- Increase weights of misclassified samples.
- Decrease weights of correctly classified samples.

4) Normalize sample weights.

5) Repeat 2-4 for a fixed number of iterations or until a stopping criterion is met.

6) Final prediction:

- combine predictions of all weak learners using their respective weights.

7) Evaluation:

- Assess the performance of the AdaBoost ensemble on a separate validation dataset.

Output:

Accuracy = 0.944

21/5/24