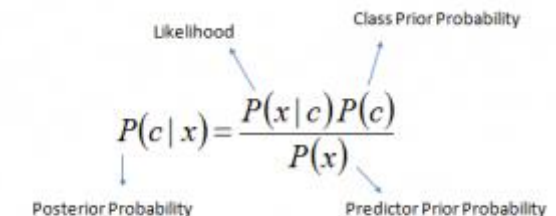


Lab 10

Naive Bayes:

It is a classification technique based on Bayes' Theorem with an independence assumption among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

Bayes theorem provides a way of computing posterior probability $P(c|x)$ from $P(c)$, $P(x)$ and $P(x|c)$. Look at the equation below:


$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$
$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Above,

- $P(c|x)$ is the posterior probability of class (c, target) given predictor (x, attributes).
- $P(c)$ is the prior probability of class.
- $P(x|c)$ is the likelihood which is the probability of the predictor given class.
- $P(x)$ is the prior probability of the predictor.

How Do Naive Bayes Algorithms Work?

Let's understand it using an example. Below we have a training data set of weather and corresponding target variable 'Play' (suggesting possibilities of playing). Now, we need to classify whether players will play or not based on weather condition. Let's follow the below steps to perform it.

- **Convert the data set into a frequency table**
In this first step data set is converted into a frequency table
- **Create Likelihood table by finding the probabilities**

Create Likelihood table by finding the probabilities like Overcast probability = 0.29 and probability of playing is 0.64.

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Likelihood table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
All	5	9
	=5/14	=9/14
	0.36	0.64

- Use Naive Bayesian equation to calculate the posterior probability

Now, use Naive Bayesian equation to calculate the posterior probability for each class.
The class with the highest posterior probability is the outcome of the prediction.

Problem: Players will play if the weather is sunny. Is this statement correct?

- We can solve it using the above-discussed method of posterior probability.
- $P(\text{Yes} \mid \text{Sunny}) = P(\text{Sunny} \mid \text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$
- Here $P(\text{Sunny} \mid \text{Yes}) * P(\text{Yes})$ is in the numerator, and $P(\text{Sunny})$ is in the denominator.
- Here we have $P(\text{Sunny} \mid \text{Yes}) = 3/9 = 0.33$, $P(\text{Sunny}) = 5/14 = 0.36$, $P(\text{Yes}) = 9/14 = 0.64$
- Now, $P(\text{Yes} \mid \text{Sunny}) = 0.33 * 0.64 / 0.36 = 0.60$, which has higher probability.

Question 1: Naïve bayse code from Scratch:

- These lines import necessary libraries. numpy is used for numerical operations, train_test_split is from scikit-learn and is used for splitting the dataset into training and testing sets, and accuracy_score is used to calculate the accuracy of the classifier.

```

7
8 import numpy as np
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import accuracy_score
11 # Self-created dataset
12 data = np.array([
13     [1, 'Sunny', 'Hot', 'High', 'No'],
14     [2, 'Sunny', 'Hot', 'High', 'No'],
15     [3, 'Overcast', 'Hot', 'High', 'Yes'],
16     [4, 'Rainy', 'Mild', 'High', 'Yes'],
17     [5, 'Rainy', 'Cool', 'Normal', 'Yes'],
18     [6, 'Rainy', 'Cool', 'Normal', 'No'],
19     [7, 'Overcast', 'Cool', 'Normal', 'Yes'],
20     [8, 'Sunny', 'Mild', 'High', 'No'],
21     [9, 'Sunny', 'Cool', 'Normal', 'Yes'],
22     [10, 'Rainy', 'Mild', 'Normal', 'Yes']
23 ])
24

```

- This separates the dataset into features (X) and labels (y). Features are columns 1 to 3 (indexing starts from 0), and labels are in column 4.
- **train_test_split**: This is a function from the scikit-learn library that splits a dataset into random train and test subsets.

```

25
26 # Split data into features and labels
27 X = data[:, 1:4] # Features
28 y = data[:, 4]   # Labels
29
30 # Split the data into training and testing sets
31 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
32

```

- Function to calculate the probabilities of different labels in the dataset.

```

28 # Helper function to calculate probabilities
29 def calculate_probabilities(data, label_column):
30     probabilities = {}
31     labels, counts = np.unique(data[:, label_column], return_counts=True)
32     total_samples = len(data)
33
34     for label, count in zip(labels, counts):
35         probabilities[label] = count / total_samples
36
37     return probabilities
38

```

- Function to train the Naive Bayes classifier by calculating probabilities for each label and feature value.

```

39 # Helper function to train Naive Bayes classifier
40 def train_naive_bayes(X, y):
41     num_features = X.shape[1]
42     unique_labels = np.unique(y)
43     probabilities = {}
44
45     for label in unique_labels:
46         label_indices = np.where(y == label)[0]
47         label_data = X[label_indices]
48
49         probabilities[label] = []
50         for i in range(num_features):
51             feature_values, counts = np.unique(label_data[:, i], return_counts=True)
52             feature_probabilities = dict(zip(feature_values, counts / len(label_data)))
53             probabilities[label].append(feature_probabilities)
54
55     return probabilities
56

```

- Function predicts the label for a given instance using the trained Naive Bayes classifier.

```

57 # function to predict using Naive Bayes classifier
58 def predict_naive_bayes(instance, probabilities):
59     predicted_label = None
60     max_probability = -1
61
62     for label, label_probabilities in probabilities.items():
63         instance_probability = 1.0
64         for i, value in enumerate(instance):
65             if value in label_probabilities[i]:
66                 instance_probability *= label_probabilities[i][value]
67             else:
68                 instance_probability = 0 # If the value is unseen in training data, probability is 0
69
70         if instance_probability > max_probability:
71             max_probability = instance_probability
72             predicted_label = label
73
74     return predicted_label
75

```

- Prediction based on test instance

```

76 # Train the Naive Bayes classifier
77 probabilities = train_naive_bayes(X, y)
78
79 # Example instance for prediction
80 test_instance = ['Sunny', 'Hot', 'High']
81
82 # Make a prediction
83 prediction = predict_naive_bayes(test_instance, probabilities)
84
85 print(f'The predicted label for the instance {test_instance} is: {prediction}')
86

```

In case your dataset is a CSV file:

```

import pandas as pd
import numpy as np

# Read data from CSV file
filename = 'play_data.csv'
df = pd.read_csv(filename)

# Display the loaded data
print("Loaded Data:")
print(df)

# Split data into features and labels
X = df.iloc[:, :-1].values # Features (all columns except the last one)
y = df.iloc[:, -1].values # Labels (last column)

# ... code...(rest of the Naive Bayes code remains the same)

```

Question 2: Naïve bayse code using built-in library:

Sklearn

Scikit-learn (Sklearn) is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistence interface in Python. This library, which is largely written in Python.

```
7
8 import pandas as pd
9 from sklearn.preprocessing import LabelEncoder
10 from sklearn.model_selection import train_test_split
11 from sklearn.naive_bayes import GaussianNB
12
13 # Reading CSV files
14 df = pd.read_csv('E:\\playsheet_dataset.csv')
15 print(df)
16
17 # Encoding the strings to Numericals
18 Numerics = LabelEncoder()
19
20 inputs = df.drop('Play', axis='columns')
21 target = df['Play']
22 print(target)
23
24 # Creating the new dataframe
25 inputs['outlook_n'] = Numerics.fit_transform(inputs['Outlook'])
26 inputs['Temp_n'] = Numerics.fit_transform(inputs['Temp'])
27 inputs['Humidity_n'] = Numerics.fit_transform(inputs['Humidity'])
28 inputs['windy_n'] = Numerics.fit_transform(inputs['Windy'])
29 print(inputs)
30
31
32 # Dropping the string values
33 inputs_n = inputs.drop(['Outlook', 'Temp', 'Humidity', 'Windy'], axis='columns')
34 print(inputs_n)
35
36 # Splitting the data into training and testing sets
37 X_train, X_test, y_train, y_test = train_test_split(inputs_n, target, test_size=0.2, random_state=42)
38
39 # Applying the Gaussian Naive Bayes to the training set
40 Classifier = GaussianNB()
41 Classifier.fit(X_train, y_train)
42
43 # Checking the accuracy on the training set
44 accuracy_train = Classifier.score(X_train, y_train)
45 print(f'Training Set Accuracy: {accuracy_train * 100:.2f}%')
46
47 # Checking the accuracy on the testing set
48 accuracy_test = Classifier.score(X_test, y_test)
49 print(f'Testing Set Accuracy: {accuracy_test * 100:.2f}%')
50
51 # Making a prediction
52 prediction = Classifier.predict([[0, 0, 0, 1]])
53 print(f'Prediction: {prediction}')
```

K-means:

K-means clustering is one of the simplest and popular unsupervised machine learning algorithms. The objective of K-means is simple: group similar data points together and discover underlying patterns. To achieve this objective, K-means looks for a fixed number (k) of clusters in a dataset. A cluster refers to a collection of data points aggregated together because of certain similarities.

Question 3: K-means from scratch:

```
7
8 import numpy as np
9 import matplotlib.pyplot as plt
10
11 def kmeans(X, k, max_iterations=100):
12     # Randomly initialize centroids
13     centroids = X[np.random.choice(X.shape[0], k, replace=False)]
14
15     for _ in range(max_iterations):
16         # Assign each data point to the nearest centroid
17         distances = np.linalg.norm(X[:, np.newaxis] - centroids, axis=2)
18         labels = np.argmin(distances, axis=1)
19
20         # Update centroids based on the mean of data points in each cluster
21         for i in range(k):
22             centroids[i] = np.mean(X[labels == i], axis=0)
23
24     return labels, centroids
25
26 # Create a self-created dataset
27 np.random.seed(42)
28 data = np.concatenate([np.random.normal(loc=5, scale=1, size=(50, 2)),
29                        np.random.normal(loc=10, scale=1, size=(50, 2))])
30
31 # Run K-means clustering
32 k = 2
33 labels, centroids = kmeans(data, k)
34
35 # Plot the data points and centroids
36 plt.scatter(data[:, 0], data[:, 1], c=labels, cmap='viridis', marker='o')
37 plt.scatter(centroids[:, 0], centroids[:, 1], c='red', marker='X', label='Centroids')
38 plt.title('K-means Clustering')
39 plt.xlabel('Feature 1')
40 plt.ylabel('Feature 2')
41 plt.legend()
42 plt.show()
43
```

In case your dataset is a CSV file:

```
//... kmean code...
```

```
# Load data from CSV file
```

```
df = pd.read_csv('data.csv')
```

```
# Extract features from the DataFrame
```

```
X = df[['Feature1', 'Feature2']].values
```

```
//... plot graph...
```

Question 4: Kmeans code using built-in library:

```
34 import numpy as np
35 import matplotlib.pyplot as plt
36 from sklearn.cluster import KMeans
37
38 # Generate random data
39 np.random.seed(42)
40 X = np.concatenate([np.random.normal(loc=5, scale=1, size=(50, 2)),
41                    np.random.normal(loc=10, scale=1, size=(50, 2))])
42
43 # Run K-means clustering using scikit-learn
44 k = 2
45 kmeans = KMeans(n_clusters=k, random_state=42)
46 labels = kmeans.fit_predict(X)
47 centroids = kmeans.cluster_centers_
48
49 # Plot the data points and centroids
50 plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis', marker='o')
51 plt.scatter(centroids[:, 0], centroids[:, 1], c='red', marker='X', label='Centroids')
52 plt.title('K-means Clustering')
53 plt.xlabel('Feature 1')
54 plt.ylabel('Feature 2')
55 plt.legend()
56 plt.show()
57 |
```