

# SPPU-TE-COMP-CONTENT – KSKA Git

Q1) Explain the concept of Naive Bayes Theorem Classifier.

- ANS.
- The Naive Bayes Classifier Theorem is a fundamental concept in the field of Machine learning and statistical Inference.
  - It provides a Framework for making predictions or classifying data based on the principles of Bayes Theorem, which describes the relationship between conditional and Marginal probabilities.
  - The Theorem (Naive Bayes) is particularly well-suited for classifying emails as either spam or not spam.
  - It states that for a given Feature or attribute, the probability of a particular class can be calculated by multiplying the prior probability of the class with the likelihood of the Feature given that class, and then dividing the result by the probability of the Feature.

Mathematically,  $P(\text{class} | \text{Feature}) = \frac{P(\text{class}) \cdot P(\text{Feature} | \text{class})}{P(\text{Feature})}$

The Naive Bayes Classifier Theorem can be expressed as:-

$$P(\text{class} | \text{Feature}) = \frac{P(\text{class}) \cdot P(\text{Feature} | \text{class})}{P(\text{Feature})}$$

- $P(\text{class} | \text{Feature})$  is the posterior probability of the class given the Feature
- $P(\text{class})$  is the prior probability of the class.
- $P(\text{Feature} | \text{class})$  is the likelihood of the Feature given the class.
- $P(\text{Feature})$  is the prior probability of the Feature.

Q2) Consider the Observation For the car Theft Scenario having 3 attributes, color, type and Origin. Find the probability of car Theft having scenarios Red SUV and Domestic.

ANS. (1) Calculate the Prior Probabilities.

$$\Rightarrow \text{i.) } P(\text{stolen} = \text{yes}) = \frac{\text{Count of yes}}{\text{Total Observation}} = \frac{5}{10} = 0.5$$

$$\text{ii.) } P(\text{stolen} = \text{No}) = \frac{\text{Count of No}}{\text{Total Observation}} = \frac{5}{10} = 0.5$$



# SPPU-TE-COMP-CONTENT - KSKA Git

(2) Calculate Conditional Probabilities for  $P(\text{stolen} = \text{yes})$

$$- P(\text{Color} = \text{Red} \mid \text{stolen} = \text{yes}) = \frac{(\text{Red, yes count})}{(\text{Total stolen yes})}$$

$$= \frac{3}{5} \approx \underline{0.6}$$

$$- P(\text{Type} = \text{SUV} \mid \text{stolen} = \text{yes}) = \frac{(\text{SUV, yes count})}{(\text{Total stolen yes})} = \frac{1}{5}$$

$$= \frac{1}{5} \approx \underline{0.2}$$

$$- P(\text{Origin} = \text{Domestic} \mid \text{stolen} = \text{yes}) = \frac{(\text{Domestic, yes count})}{(\text{Total stolen yes})}$$

$$= \frac{2}{5} \approx \underline{0.4}$$

(3) Combine Probabilities :-

$$P(\text{stolen} = \text{yes} \mid \text{RED, SUV, DOMESTIC}) = 0.5 \times 0.6 \times 0.2 = \underline{0.024}$$

Answers The Probability of Car Theft having scenario Red, SUV and Domestic (car) is 0.024

Q3) Write a Python code for the Data pre-processing mentioned in step 4 and Explain every step in detail.

Ans. Python Code :-

```
import pandas as pd
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.impute import SimpleImputer
```

```
df = pd.read_csv("https://archive.iris.csv")
```

```
names = ["sepal-length", "sepal-width", "petal-length", "petal-width", "class"]
```

```
df = df[names]
```

# SPPU-TE-COMP-CONTENT – KSKA Git

```
imputer = SimpleImputer(strategy = 'mean')
df.iloc[:, :-1] = Imputer.fit_transform(df.iloc[:, :-1])

encoder = LabelEncoder()
df['class'] = encoder.fit_transform(df['class'])

scaler = StandardScaler()
df.iloc[:, :-1] = scaler.fit_transform(df.iloc[:, :-1])

print("Preprocessed Data Sample:")
print(df.head())
```

## # EXPLANATION: -

### → (1) Load the Dataset.

- The Iris Dataset (iris.csv) is read from the .csv file

### (2) Handle Missing Values.

- We use SimpleImputer to replace missing values in numerical columns with their mean.

### (3) Encode the Categorical Data.

- Convert Categorical to Numerical Labels.

### (4) Feature Scaling.

- We use StandardScaler() to normalize numerical features, ensuring that they have a mean of 0 (zero) and standard deviation of 1.

### (5) Display the Processed Data.

- The first few rows of pre-processed DF are printed to verify changes.