

Experiment No. 3

Date: 23/01/2025

Title:- Implementation of Bayes Classifier for data Classification.

Aim:- To implement a Bayes classifier to classify data into different categories using probabilistic reasoning.

Theory:

The Bayes classifier is a Supervised Learning Algorithm based on Bayes Theorem, which provide a probabilistic approach to classify datapoint. It predict the class of a given data instance calculate the posterior probability of each class and selecting the one with the highest probability.

Bayes Theorem:
$$P(C_k|X) = \frac{P(X|C_k) \cdot P(C_k)}{P(X)}$$

$P(C_k|X)$: Posterior probability of class C_k given the data X .

$P(X|C_k)$: Likelihood of observing data X given the class C_k .

$P(C_k)$: Prior probability of class C_k (based on prior knowledge)

$P(X)$: Evidence or normalization factor, calculated as the sum of the probabilities of X over all classes.

The Naive Bayes classifier simplifies computation by assuming that the features are conditionally independent given the class level.

If the features are continuous, assuming that the data for each feature follows a Gaussian (Normal) Distribution. The probability density function (PDF):

$$P(X|c_k) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

x : feature value

μ : The mean (average) of the feature for class c_k

σ^2 : The variance of the feature for class c_k

$P(x|c_k)$: The probability of observing x given that the data belongs to class c_k .

for the classification,

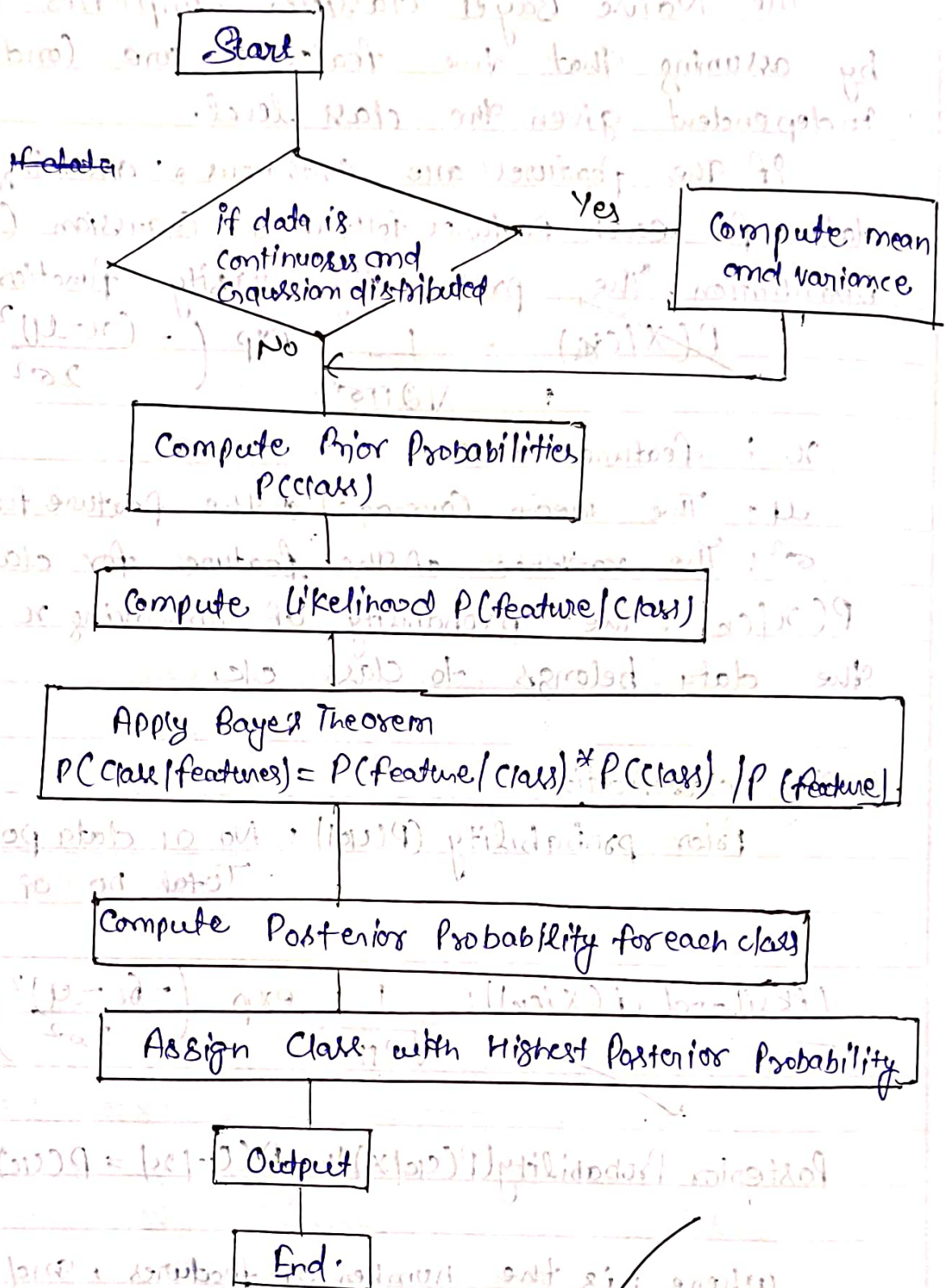
Prior probability $P(c_k)$: $\frac{\text{No of data points in class } c_k}{\text{Total no of data points.}}$

$$\text{Likelihood } P(X|c_k): \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\text{Posterior Probability } P(c_k|x): P(c_k|x) = P(c_k) \cdot \prod_{i=1}^n P(x_i|c_k)$$

where n is the number of features, and x_i is the value of the i th feature.

Flowchart:



Class prediction: Select the class with the highest posterior probability = $\max_{\substack{C_k \\ k=1 \dots K}} (P(C_k|X))$.

The Algorithm:

- Dataset with n features and m data points
- Target classes C_1, C_2, \dots, C_K .
- for each class C_k , compute the mean μ and variance (σ^2) of each feature using the training data if the data is continuous gaussian distributed.
- Compute the prior probability for each class.

$$P(C_k) = \frac{\text{No of data point in class } C_k}{\text{total no of data points}}$$

→ Likelihood Calculation:

for a given test point $X = \{x_1, x_2, \dots, x_n\}$ calculate the likelihood of each feature using

$$P(x_i|C_k) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right)$$

→ Posterior Calculation:

Combine the prior and likelihood for all features to calculate the posterior probability for each class

$$P(C_k|X) = P(C_k) \cdot \prod_{i=1}^n P(x_i|C_k)$$

- Class C_k with the highest posterior probability.

Result:

Dataset:

Person	height (feet)	weight (lbs)	foot size (inches)
male	6	180	12
male	5.92	190	11
male	5.58	170	12
male	5.92	165	10
female	5	100	6
female	5.5	150	8
female	5.42	130	7
female	5.75	150	9

Sample data. (height: 6, weight: 180, foot size: 8)

fruit	Long	Sweet	Yellow	total
Banana	400	350	450	500
Orange	0	150	300	300
Other	100	150	50	200
Total	500	650	800	1000

--- Mean and Variance for Male Data ---

	Mean	Variance
Height	5.855	0.035033
Weight	176.25	122.92
Foot Size	11.25	0.91667

--- Mean and Variance for Female Data ---

	Mean	Variance
Height	5.4175	0.097225
Weight	132.5	558.33
Foot Size	7.5	1.6667

--- Probabilities for Test Data ---

Male probability for Height: 1.578883

Female probability for Height: 0.223459

Male probability for Weight: 0.000006

Female probability for Weight: 0.016789

Male probability for Foot Size: 0.001311

Female probability for Foot Size: 0.286691

--- Prior Probabilities ---

Prior probability for Male: 0.50

Prior probability for Female: 0.50

--- Posterior Probabilities ---

Posterior probability for Male: 0.00011)

Posterior probability for Female: 0.99998 }

Result: The test data is classified as Female

Discussion:

- In the first dataset included height, weight and foot size measurement for male and female. These features were used to calculate probabilities and classify the given test data.
- For each feature the mean and variance were computed separately for male and female data. Using the gaussian probability density function, the likelihood $P(\text{feature}|\text{male})$ and $P(\text{feature}|\text{female})$ were calculated for the test data. [6, 130, 8]
- The prior probabilities $P(\text{male})$ & $P(\text{female})$ were calculated based on the size of male and female data.
- The posterior probabilities $P(\text{male}|\text{data})$ and $P(\text{female}|\text{data})$ were computed by multiplying the likelihoods of all features with their respective prior.
- The Based on the posterior probabilities the test data was classified as female. The test data better matched the male data distribution than the female distribution, resulting in a female classification.
- The bayesian classifier successfully combined prior probabilities and feature likelihood to reach an interpretable and probabilistic decision.
- The bayesian classifier is effective in gender classification using continuous features. However increasing the dataset size and incorporating dependencies b/w features could improve performance.

--- Prior Probabilities ---

Banana: 0.4048

Orange: 0.1786

Other: 0.4167

--- Likelihood Probabilities ---

Banana ($P(\text{Features}|\text{Fruit})$): 0.1212

Orange ($P(\text{Features}|\text{Fruit})$): 0.0000

Other ($P(\text{Features}|\text{Fruit})$): 0.0400

--- Posterior Probabilities ---

Banana ($P(\text{Fruit}|\text{Features})$): 0.7463

Orange ($P(\text{Fruit}|\text{Features})$): 0.0000

Other ($P(\text{Fruit}|\text{Features})$): 0.2537

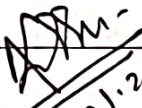
The given data (Long, Sweet, Yellow)
corresponds to: Banana

- For the second dataset with three fruit (Banana, Orange, Apple) and with feature long, sweet, and yellow. Each fruit was representing by its association with these features in terms of count.
- The prior probability of each fruit was calculated based on the total no of features count for that fruit relative to the dataset.
- The likelihood $P(\text{features}/\text{fruit})$ was computed for each fruit by analyzing how strongly the features matched the characteristics of each fruits. Banana showed the highest one, reflecting its strong alignment with the features.
- Using the Bayesian theorem the posterior probabilities were calculated by combining the likelihoods and prior. Banana emerged with the highest posterior probability, indicating that the features most likely correspond to a banana.
- The experiment classified clear and interpretable result, with probabilities calculated at each step. This transparency helps understand why a particular fruit was chosen.
- The dataset assumed all features contribute equally to Classification. This could limit accuracy if certain features are more important than others.
- The small dataset constrained the generalizability of result. Larger dataset with more diverse example would improve the model's reliability.

Conclusion:

The Bayesian classifier effectively classified the test data as female, leveraging height, weight, and foot size. By combining prior and likelihood probabilities, the model provided an interpretable and accurate decision. However, a larger dataset and consideration of feature dependencies could enhance performance.

Using this, the test data was accurately classified as female, based on features. The method's probabilistic approach effectively matched features with prior knowledge, demonstrating its robustness for categorical classification tasks.


30.01.25