

# Gaussian Naïve Bayes Classifier

*Vikrant Bhati*

## 1. Approach(es)

### Naïve Bayes

The Naïve Bayes classifier is a family of classifiers based on Bayes' Theorem. It makes an assumption that features are independent of each other given the class label, which is known as the naïve independence assumption. Despite this simplifying assumption, Naïve Bayes performs well in many real-world applications, particularly in text classification, spam detection, and medical diagnosis. Naïve Bayes can be represented as:

Bayes Theorem:

$$P(C | X) = \frac{P(X | C) \cdot P(C)}{P(X)}$$

Where:

- $P(C|X)$  = Posterior probability (the probability of class C given X)
- $P(X|C)$  = Likelihood (the probability of X given class C)
- $P(C)$  = Prior probability of class C
- $P(X)$  = Evidence (the total probability of X)

Where Naïve Bayes after naïve assumption:

$$P(X | C) = P(X_1 | C) \cdot P(X_2 | C) \cdot \dots \cdot P(X_n | C)$$

Where:

- $P(X|C)$  = Probability of X given class C
- $P(X_i|C)$  = Probability of  $X_i$  given class C

However, there are many scenarios where our input columns are continuous, and therefore for continuous features we use Gaussian Naïve Bayes.

### Gaussian Naïve Bayes

The Gaussian Naïve Bayes classifier is a probabilistic model based on Bayes' Theorem, assuming that the features are conditionally independent given the class label and follow a Gaussian (normal)

distribution. It is used for classification tasks, especially when the input features are continuous. The probability density function for a feature given a class is modeled as:

$$P(x | C) = (1 / \text{sqrt}(2\pi\sigma_C^2)) * \exp(-(x - \mu_C)^2 / (2\sigma_C^2))$$

Where:

- $\mu_C$  = Mean of the feature for class  $C$
- $\sigma_C$  = Standard deviation for class  $C$

### **Dataset**

For the experiments, I used the Breast Cancer Wisconsin dataset from scikit-learn where the output column is a binary classification based on 30 input columns.

### **Output Column:**

The target variable (output) is a binary classification:

- **0:** Benign (non-cancerous)
- **1:** Malignant (Cancerous)

### **Input Features:**

- The dataset consists of 30 real, positive features that are derived from the Breast Cancer Wisconsin dataset.

**Table 1: A summary of the breast cancer dataset**

Classes	2
Samples per class	212(M),357(B)
Samples total	569
Dimensionality	30
Features	real, positive

*Fig1: Summary of the breast cancer dataset*

Further, the data was split into the training test data using *train\_test\_split* with 80% data being used for training and the remaining 20% for testing.

## Adding Gaussian Noise

Zero mean Gaussian noise was added to see the impact on the accuracy score and confusion matrix. I used 6 noise variances of values: 50, 1000, 20000, 400000, 8000000 and 100000000.

## Evaluation

In this experiment, the performance of the Gaussian Naïve Bayes classifier was evaluated using two key metrics:

### **Accuracy Score:**

Accuracy is the proportion of correctly classified instances out of the total number of instances.

$$\text{Accuracy} = (\text{True Positives} + \text{True Negatives}) / \text{Total Number of Predictions}$$

### **Confusion Matrix:**

A confusion matrix summarizes the classification results by showing the counts of true positives, true negatives, false positives, and false negatives.

$$\begin{bmatrix} \text{True Positives} & \text{False Negatives} \end{bmatrix}$$
$$\begin{bmatrix} \text{False Positives} & \text{True Negatives} \end{bmatrix}$$

## **2. Experimental Results**

For the first experiment, I used the original data without any Gaussian noise and the output was:

- Accuracy Score: 0.9210526315789473
- Confusion Matrix:  $\begin{bmatrix} 33 & 6 \\ 3 & 72 \end{bmatrix}$

After adding the noise, a zero mean Gaussian noise with a variance of 50 was:

- Accuracy Score: 0.8859649122807017
- Confusion Matrix:  $\begin{bmatrix} 31 & 8 \\ 5 & 70 \end{bmatrix}$

After adding the noise, a zero mean Gaussian noise with a variance of 1000 was:

- Accuracy Score: 0.8596491228070176
- Confusion Matrix:  $\begin{bmatrix} 27 & 12 \\ 4 & 71 \end{bmatrix}$

After adding the noise, a zero mean Gaussian noise with a variance of 20000 was:

- Accuracy Score: 0.8859649122807017
- Confusion Matrix:  $\begin{bmatrix} 29 & 10 \\ 3 & 72 \end{bmatrix}$

After adding the noise, a zero mean Gaussian noise with a variance of 400000 was:

- Accuracy Score: 0.7807017543859649
- Confusion Matrix:  $\begin{bmatrix} 23 & 16 \\ 9 & 66 \end{bmatrix}$

After adding the noise, a zero mean Gaussian noise with a variance of 8000000 was:

- Accuracy Score: 0.5526315789473685
- Confusion Matrix:  $\begin{bmatrix} 12 & 27 \\ 24 & 51 \end{bmatrix}$

After adding the noise, a zero mean Gaussian noise with a variance of 100000000 was:

- Accuracy Score: 0.6403508771929824
- Confusion Matrix:  $\begin{bmatrix} 10 & 29 \\ 12 & 63 \end{bmatrix}$

### 3. Discussion

The Gaussian classifier performed well on the clean dataset, achieving an accuracy of 92.11%. The confusion matrix for this setup was  $\begin{bmatrix} 33 & 6 \\ 3 & 72 \end{bmatrix}$ . However, as noise variance increased, the accuracy steadily declined, reflecting the classifier's sensitivity to noise. Initially, the classifier maintained a high accuracy at lower noise levels, but the performance degraded significantly with higher noise variance.

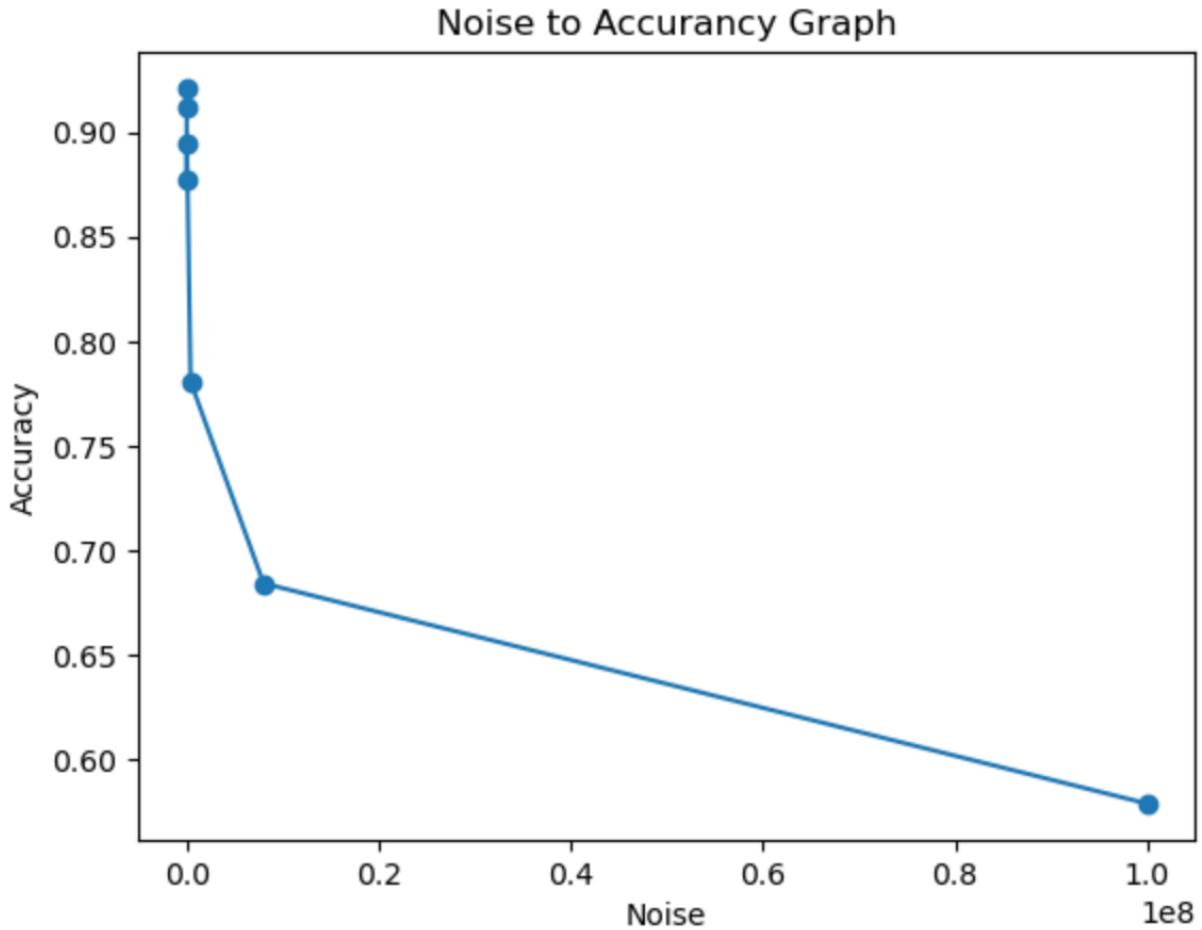
#### Trend Observations

**Noise Variance: 50 to 20,000** – Accuracy fluctuated between 85.96% and 88.60%, indicating the classifier's robustness to moderate noise.

**Noise Variance: 400,000** – Accuracy dropped to 78.07%, showing the beginning of a significant decline.

**Noise Variance: 8,000,000 and 100,000,000** – The accuracy dropped to 55.26% and 64.04%, respectively. The high misclassification rate at these noise levels indicates that the noise overwhelmed the signal in the data.

This trend is evident in blow graph that showcase the accuracy vs noise.



*Fig2. Graph that showcases the accuracy vs noise*

The confusion matrices at higher noise levels showed a marked increase in false positives and false negatives, indicating that the classifier struggled to maintain decision boundaries.

The results demonstrate a clear negative correlation between noise variance and accuracy. While the Gaussian Naïve Bayes classifier is reasonably robust to small amounts of noise, its performance deteriorates as noise becomes more dominant in the feature space.