

# MLT: Week 8

GENERATIVE MODELS - NAIVE BAYES

Sherry Thomas

# Naive Bayes Algorithm

# Naive Bayes Algorithm

Given a dataset  $\{x_1, \dots, x_n\}$  where  $x_i \in \{0, 1\}^d$ ,  
let  $\{y_1, \dots, y_n\}$  be the labels, where  $y_i \in \{0, 1\}$ .

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

The Naive Bayes Algorithm is characterized by two main things:

1. Class Conditional Independence Assumption
2. Bayes Rule

The Naive Bayes Algorithm is characterized by two main things:

1. Class Conditional Independence Assumption
2. Bayes Rule

Based on the first characteristic, we get  $2d+1$  parameters where  $d$  is the number of dimensions of the features set.

Using MLE, the estimates of the parameters are given by,

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i$$

The Naive Bayes Algorithm is characterized by two main things:

1. Class Conditional Independence Assumption
2. Bayes Rule

Based on the first characteristic, we get  $2d+1$  parameters where  $d$  is the number of dimensions of the features set.

Using MLE, the estimates of the parameters are given by,

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i$$
$$\hat{p}_j^y = \frac{\sum_{i=1}^n \mathbb{I}(f_j^i = 1, y_i = y)}{\sum_{i=1}^n \mathbb{I}(y_i = y)} \quad \forall j \in \{1, 2, \dots, d\} \quad \forall y \in \{0, 1\}$$

Estimating the GMM parameters for our dataset would result in parameters whose values are zero. Therefore, we need to apply **Laplace Smoothing**.

After applying it on our dataset, we get,

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Estimating the GMM parameters for our dataset would result in parameters whose values are zero. Therefore, we need to apply **Laplace Smoothing**.

After applying it on our dataset, we get,

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

In this new dataset, parameters don't end up as zero probability estimates.



Estimating the GMM parameters for our dataset would result in parameters whose values are zero. Therefore, we need to apply **Laplace Smoothing**.

After applying it on our dataset, we get,

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

In this new dataset, parameters don't end up as zero probability estimates.

Calculating  $\hat{p}$  using MLE for our dataset, we get,

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i$$

Estimating the GMM parameters for our dataset would result in parameters whose values are zero. Therefore, we need to apply **Laplace Smoothing**.

After applying it on our dataset, we get,

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

In this new dataset, parameters don't end up as zero probability estimates.

Calculating  $\hat{p}$  using MLE for our dataset, we get,

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{8} (1 + 1 + 0 + 0 + 1 + 0 + 1 + 0) = \frac{4}{8} = 0.5$$

Estimating the GMM parameters for our dataset would result in parameters whose values are zero. Therefore, we need to apply **Laplace Smoothing**.

After applying it on our dataset, we get,

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

In this new dataset, parameters don't end up as zero probability estimates.

Calculating  $\hat{p}$  using MLE for our dataset, we get,

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{8} (1 + 1 + 0 + 0 + 1 + 0 + 1 + 0) = \frac{4}{8} = 0.5$$

$$\therefore \hat{p} = 0.5$$

Calculating  $\hat{p}_j^y$  using MLE for our dataset, we get,

$$\hat{p}_j^y = \frac{\sum_{i=1}^n \mathbb{I}(f_j^i = 1, y_i = y)}{\sum_{i=1}^n \mathbb{I}(y_i = y)}$$

Calculating  $\hat{p}_j^y$  using MLE for our dataset, we get,

$$\hat{p}_j^y = \frac{\sum_{i=1}^n \mathbb{I}(f_j^i = 1, y_i = y)}{\sum_{i=1}^n \mathbb{I}(y_i = y)}$$

$$\hat{p}_1^1 = \frac{1 + 1 + 1}{4} = 0.75 \qquad \hat{p}_1^0 = \frac{1}{4} = 0.25$$

Calculating  $\hat{p}_j^y$  using MLE for our dataset, we get,

$$\hat{p}_j^y = \frac{\sum_{i=1}^n \mathbb{I}(f_j^i = 1, y_i = y)}{\sum_{i=1}^n \mathbb{I}(y_i = y)}$$

$$\begin{aligned} \hat{p}_1^1 &= \frac{1+1+1}{4} = 0.75 & \hat{p}_1^0 &= \frac{1}{4} = 0.25 \\ \hat{p}_2^1 &= \frac{1+1}{4} = 0.50 & \hat{p}_2^0 &= \frac{1+1}{4} = 0.50 \end{aligned}$$

Calculating  $\hat{p}_j^y$  using MLE for our dataset, we get,

$$\hat{p}_j^y = \frac{\sum_{i=1}^n \mathbb{I}(f_j^i = 1, y_i = y)}{\sum_{i=1}^n \mathbb{I}(y_i = y)}$$

$$\begin{array}{ll} \hat{p}_1^1 = \frac{1+1+1}{4} = 0.75 & \hat{p}_1^0 = \frac{1}{4} = 0.25 \\ \hat{p}_2^1 = \frac{1+1}{4} = 0.50 & \hat{p}_2^0 = \frac{1+1}{4} = 0.50 \\ \hat{p}_3^1 = \frac{1}{4} = 0.25 & \hat{p}_3^0 = \frac{1+1}{4} = 0.50 \end{array}$$

## Prediction using Naive Bayes

If

$$\left( \prod_{i=1}^d \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^d \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p})$$

we predict  $y = 1$ , otherwise  $y = 0$ .



## Prediction using Naive Bayes

If

$$\left( \prod_{i=1}^d \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^d \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p})$$

we predict  $y = 1$ , otherwise  $y = 0$ .

Predicting  $y_1$  for our dataset,

$$y_1 = \mathbb{I} \left( \left( \prod_{i=1}^3 \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^3 \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p}) \right)$$

## Prediction using Naive Bayes

If

$$\left( \prod_{i=1}^d \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^d \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p})$$

we predict  $y = 1$ , otherwise  $y = 0$ .

Predicting  $y_1$  for our dataset,

$$\begin{aligned} y_1 &= \mathbb{I} \left( \left( \prod_{i=1}^3 \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^3 \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p}) \right) \\ &= \mathbb{I} \left( ((0.75)^1 (0.75)^0 (0.50)^0 (0.50)^1 (0.25)^0 (0.75)^1) 0.5 \geq ((0.25)^1 (0.75)^0 (0.50)^0 (0.50)^1 (0.50)^0 (0.50)^1) 0.5 \right) \end{aligned}$$

## Prediction using Naive Bayes

If

$$\left( \prod_{i=1}^d \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^d \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p})$$

we predict  $y = 1$ , otherwise  $y = 0$ .

Predicting  $y_1$  for our dataset,

$$\begin{aligned} y_1 &= \mathbb{I} \left( \left( \prod_{i=1}^3 \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^3 \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p}) \right) \\ &= \mathbb{I} \left( ((0.75)^1 (0.75)^0 (0.50)^0 (0.50)^1 (0.25)^0 (0.75)^1) 0.5 \geq ((0.25)^1 (0.75)^0 (0.50)^0 (0.50)^1 (0.50)^0 (0.50)^1) 0.5 \right) \\ &= \mathbb{I}(0.141 \geq 0.031) \end{aligned}$$

## Prediction using Naive Bayes

If

$$\left( \prod_{i=1}^d \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^d \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p})$$

we predict  $y = 1$ , otherwise  $y = 0$ .

Predicting  $y_1$  for our dataset,

$$\begin{aligned} y_1 &= \mathbb{I} \left( \left( \prod_{i=1}^3 \left( \hat{p}_i^1 \right)^{f_i} \left( 1 - \hat{p}_i^1 \right)^{1-f_i} \right) \hat{p} \geq \left( \prod_{i=1}^3 \left( \hat{p}_i^0 \right)^{f_i} \left( 1 - \hat{p}_i^0 \right)^{1-f_i} \right) (1 - \hat{p}) \right) \\ &= \mathbb{I} \left( ((0.75)^1 (0.75)^0 (0.50)^0 (0.50)^1 (0.25)^0 (0.75)^1) 0.5 \geq ((0.25)^1 (0.75)^0 (0.50)^0 (0.50)^1 (0.50)^0 (0.50)^1) 0.5 \right) \\ &= \mathbb{I}(0.141 \geq 0.031) \\ y_1 &= 1 \end{aligned}$$

Prediction table for all datapoints:

Prediction table for all datapoints:

Label	$P(\hat{y} = 1 \hat{x})$	$P(\hat{y} = 0 \hat{x})$	Prediction	Actual
$y_1$	0.141	0.031	1	1
$y_2$	0.141	0.031	1	1
$y_3$	0.047	0.094	0	0
$y_4$	0.016	0.094	0	0
$y_5$	0.047	0.031	1	1
$y_6$	0.047	0.031	1	0
$y_7$	0.047	0.094	0	1
$y_8$	0.047	0.094	0	0

# Gaussian Naïve Bayes Algorithm

# Gaussian Naive Bayes Algorithm

Given a dataset  $\{x_1, \dots, x_n\}$  where  $x_i \in \mathbb{R}^d$ ,  
let  $\{y_1, \dots, y_n\}$  be the labels, where  $y_i \in \{0, 1\}$ .

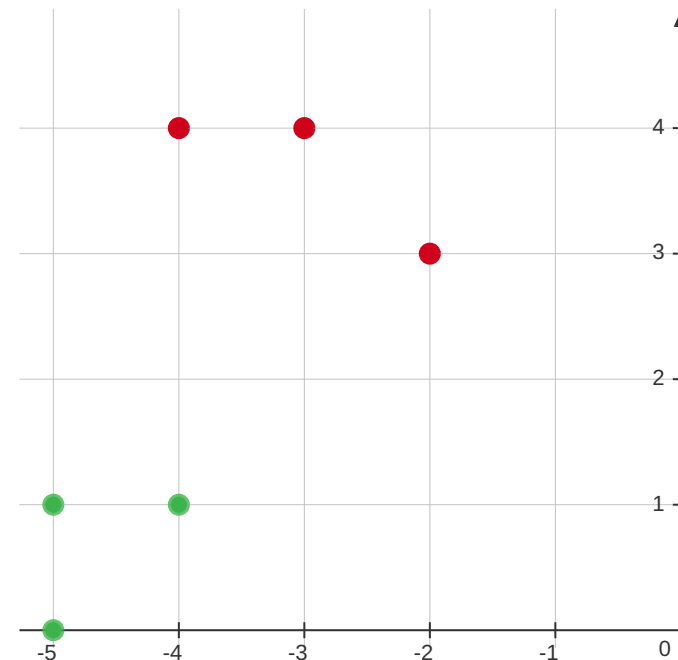
$$X = \begin{bmatrix} -3 & -4 & -2 & -4 & -5 & -5 \\ 4 & 4 & 3 & 1 & 1 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



# Gaussian Naive Bayes Algorithm

Given a dataset  $\{x_1, \dots, x_n\}$  where  $x_i \in \mathbb{R}^d$ ,  
let  $\{y_1, \dots, y_n\}$  be the labels, where  $y_i \in \{0, 1\}$ .

$$X = \begin{bmatrix} -3 & -4 & -2 & -4 & -5 & -5 \\ 4 & 4 & 3 & 1 & 1 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



Let  $P(x|y = 1) \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $P(x|y = 0) \sim \mathcal{N}(\mu_0, \Sigma_0)$ .

The parameters to be estimated are  $\hat{p}$ ,  $\hat{\mu}_0$ ,  $\hat{\mu}_1$ , and  $\hat{\Sigma}_0$  and  $\hat{\Sigma}_1$ .

Let  $P(x|y = 1) \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $P(x|y = 0) \sim \mathcal{N}(\mu_0, \Sigma_0)$ .

The parameters to be estimated are  $\hat{p}$ ,  $\hat{\mu}_0$ ,  $\hat{\mu}_1$ , and  $\hat{\Sigma}_0$  and  $\hat{\Sigma}_1$ .

Using Maximum Likelihood Estimation, we get the following results:

Let  $P(x|y = 1) \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $P(x|y = 0) \sim \mathcal{N}(\mu_0, \Sigma_0)$ .

The parameters to be estimated are  $\hat{p}$ ,  $\hat{\mu}_0$ ,  $\hat{\mu}_1$ , and  $\hat{\Sigma}_0$  and  $\hat{\Sigma}_1$ .

Using Maximum Likelihood Estimation, we get the following results:

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i$$

Let  $P(x|y = 1) \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $P(x|y = 0) \sim \mathcal{N}(\mu_0, \Sigma_0)$ .

The parameters to be estimated are  $\hat{p}$ ,  $\hat{\mu}_0$ ,  $\hat{\mu}_1$ , and  $\hat{\Sigma}_0$  and  $\hat{\Sigma}_1$ .

Using Maximum Likelihood Estimation, we get the following results:

$$\begin{aligned}\hat{p} &= \frac{1}{n} \sum_{i=1}^n y_i \\ \hat{\mu}_1 &= \frac{\sum_{i=1}^n \mathbb{I}(y_i = 1) x_i}{\sum_{i=1}^n \mathbb{I}(y_i = 1)} \\ \hat{\mu}_0 &= \frac{\sum_{i=1}^n \mathbb{I}(y_i = 0) x_i}{\sum_{i=1}^n \mathbb{I}(y_i = 0)}\end{aligned}$$

$$\widehat{\Sigma}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k) \left( (x_i - \widehat{\mu}_k)(x_i - \widehat{\mu}_k)^T \right)}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

Where  $\widehat{p}$  is the proportion of data points labeled 1,  $\widehat{\mu}_1$  is the sample mean of data points labeled 1,  $\widehat{\mu}_0$  is the sample mean of data points labeled 0, and  $\widehat{\Sigma}_1$  and  $\widehat{\Sigma}_0$  are the covariance matrices for classes 1 and 0 respectively.

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k) \left( (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T \right)}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

Where  $\hat{p}$  is the proportion of data points labeled 1,  $\hat{\mu}_1$  is the sample mean of data points labeled 1,  $\hat{\mu}_0$  is the sample mean of data points labeled 0, and  $\hat{\Sigma}_1$  and  $\hat{\Sigma}_0$  are the covariance matrices for classes 1 and 0 respectively.

Calculating  $\hat{p}$  using MLE for our dataset, we get,

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{6}(1 + 1 + 1 + 0 + 0 + 0) = \frac{3}{6} = 0.5$$

$$\therefore \hat{p} = 0.5$$

Calculating  $\hat{\mu}_k$  using MLE for our dataset, we get,

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k)x_i}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$



Calculating  $\hat{\mu}_k$  using MLE for our dataset, we get,

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k)x_i}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

$$\hat{\mu}_1 = \frac{\begin{bmatrix} -3 \\ 4 \end{bmatrix} + \begin{bmatrix} -4 \\ 4 \end{bmatrix} + \begin{bmatrix} -2 \\ 3 \end{bmatrix}}{3} = \begin{bmatrix} -3 \\ 3.666 \end{bmatrix}$$

$$\hat{\mu}_0 = \frac{\begin{bmatrix} -4 \\ 1 \end{bmatrix} + \begin{bmatrix} -5 \\ 1 \end{bmatrix} + \begin{bmatrix} -5 \\ 0 \end{bmatrix}}{3} = \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix}$$

Calculating  $\widehat{\Sigma}_k$  using MLE for our dataset, we get,

$$\widehat{\Sigma}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k) \left( (x_i - \widehat{\mu}_k)(x_i - \widehat{\mu}_k)^T \right)}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

Calculating  $\hat{\Sigma}_k$  using MLE for our dataset, we get,

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k) \left( (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T \right)}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

$$\hat{\Sigma}_1 = \frac{\left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T}{3}$$

Calculating  $\hat{\Sigma}_k$  using MLE for our dataset, we get,

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k) \left( (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T \right)}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

$$\hat{\Sigma}_1 = \frac{\left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T}{3}$$

$$\hat{\Sigma}_1 = \begin{bmatrix} 0.666 & 0.333 \\ -0.333 & 0.222 \end{bmatrix}$$

Calculating  $\hat{\Sigma}_k$  using MLE for our dataset, we get,

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k) \left( (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T \right)}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

$$\hat{\Sigma}_1 = \frac{\left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T}{3}$$

$$\hat{\Sigma}_1 = \begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}$$

$$\hat{\Sigma}_0 = \frac{\left( \begin{bmatrix} -4 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right) \left( \begin{bmatrix} -4 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -5 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right) \left( \begin{bmatrix} -5 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -5 \\ 0 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right) \left( \begin{bmatrix} -5 \\ 0 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right)^T}{3}$$

Calculating  $\hat{\Sigma}_k$  using MLE for our dataset, we get,

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^n \mathbb{I}(y_i = k) \left( (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T \right)}{\sum_{i=1}^n \mathbb{I}(y_i = k)}$$

$$\hat{\Sigma}_1 = \frac{\left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -3 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -4 \\ 4 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right) \left( \begin{bmatrix} -2 \\ 3 \end{bmatrix} - \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} \right)^T}{3}$$

$$\hat{\Sigma}_1 = \begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}$$

$$\hat{\Sigma}_0 = \frac{\left( \begin{bmatrix} -4 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right) \left( \begin{bmatrix} -4 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -5 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right) \left( \begin{bmatrix} -5 \\ 1 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -5 \\ 0 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right) \left( \begin{bmatrix} -5 \\ 0 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right)^T}{3}$$

$$\hat{\Sigma}_0 = \begin{bmatrix} 0.222 & 0.111 \\ 0.111 & 0.222 \end{bmatrix}$$

## Prediction using Gaussian Naïve Bayes

## Prediction using Gaussian Naïve Bayes

Predict  $y_i = 1$  if:

$$f(x_i; \hat{\mu}_1, \hat{\Sigma}_1) \hat{p} \geq f(x_i; \hat{\mu}_0, \hat{\Sigma}_0)(1 - \hat{p})$$



## Prediction using Gaussian Naïve Bayes

Predict  $y_i = 1$  if:

$$f(x_i; \hat{\mu}_1, \hat{\Sigma}_1) \hat{p} \geq f(x_i; \hat{\mu}_0, \hat{\Sigma}_0) (1 - \hat{p})$$
$$e^{-(x_i - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1} (x_i - \hat{\mu}_1)} \hat{p} \geq e^{-(x_i - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (x_i - \hat{\mu}_0)} (1 - \hat{p})$$

## Prediction using Gaussian Naïve Bayes

Predict  $y_i = 1$  if:

$$f(x_i; \hat{\mu}_1, \hat{\Sigma}_1) \hat{p} \geq f(x_i; \hat{\mu}_0, \hat{\Sigma}_0) (1 - \hat{p})$$

$$e^{-(x_i - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1} (x_i - \hat{\mu}_1)} \hat{p} \geq e^{-(x_i - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (x_i - \hat{\mu}_0)} (1 - \hat{p})$$

$$-(x_i - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1} (x_i - \hat{\mu}_1) + \log(\hat{p}) \geq -(x_i - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (x_i - \hat{\mu}_0) + \log(1 - \hat{p})$$

## Prediction using Gaussian Naïve Bayes

Predict  $y_i = 1$  if:

$$f(x_i; \hat{\mu}_1, \hat{\Sigma}_1) \hat{p} \geq f(x_i; \hat{\mu}_0, \hat{\Sigma}_0) (1 - \hat{p})$$

$$e^{-(x_i - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1} (x_i - \hat{\mu}_1)} \hat{p} \geq e^{-(x_i - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (x_i - \hat{\mu}_0)} (1 - \hat{p})$$

$$-(x_i - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1} (x_i - \hat{\mu}_1) + \log(\hat{p}) \geq -(x_i - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (x_i - \hat{\mu}_0) + \log(1 - \hat{p})$$

$$x_i^T \left( \hat{\Sigma}_1^{-1} - \hat{\Sigma}_0^{-1} \right) x_i - 2 \left( \hat{\mu}_1^T \hat{\Sigma}_1^{-1} - \hat{\mu}_0^T \hat{\Sigma}_0^{-1} \right) x_i + \left( \hat{\mu}_0^T \hat{\Sigma}_0^{-1} \hat{\mu}_0 - \hat{\mu}_1^T \hat{\Sigma}_1^{-1} \hat{\mu}_1 \right) + \log \left( \frac{1 - \hat{p}}{\hat{p}} \right) \geq 0$$

## Prediction using Gaussian Naïve Bayes

Predict  $y_i = 1$  if:

$$\begin{aligned} f(x_i; \hat{\mu}_1, \hat{\Sigma}_1) \hat{p} &\geq f(x_i; \hat{\mu}_0, \hat{\Sigma}_0) (1 - \hat{p}) \\ e^{-(x_i - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1} (x_i - \hat{\mu}_1)} \hat{p} &\geq e^{-(x_i - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (x_i - \hat{\mu}_0)} (1 - \hat{p}) \\ -(x_i - \hat{\mu}_1)^T \hat{\Sigma}_1^{-1} (x_i - \hat{\mu}_1) + \log(\hat{p}) &\geq -(x_i - \hat{\mu}_0)^T \hat{\Sigma}_0^{-1} (x_i - \hat{\mu}_0) + \log(1 - \hat{p}) \\ x_i^T \left( \hat{\Sigma}_1^{-1} - \hat{\Sigma}_0^{-1} \right) x_i - 2 \left( \hat{\mu}_1^T \hat{\Sigma}_1^{-1} - \hat{\mu}_0^T \hat{\Sigma}_0^{-1} \right) x_i &+ \left( \hat{\mu}_0^T \hat{\Sigma}_0^{-1} \hat{\mu}_0 - \hat{\mu}_1^T \hat{\Sigma}_1^{-1} \hat{\mu}_1 \right) + \log \left( \frac{1 - \hat{p}}{\hat{p}} \right) \geq 0 \end{aligned}$$

Hence, we can say that the decision function is of the form  $x^T Q x - 2b^T x + c \geq 0$

where  $Q = \hat{\Sigma}_1^{-1} - \hat{\Sigma}_0^{-1}$ ,  $b = \hat{\mu}_1^T \hat{\Sigma}_1^{-1} - \hat{\mu}_0^T \hat{\Sigma}_0^{-1}$ , and  $c = (\hat{\mu}_0^T \hat{\Sigma}_0^{-1} \hat{\mu}_0 - \hat{\mu}_1^T \hat{\Sigma}_1^{-1} \hat{\mu}_1) + \log(\frac{1 - \hat{p}}{\hat{p}})$ .

Predicting  $y_1$  for our dataset,

$$y_1 = \mathbb{I} \left( x_1^T \left( \widehat{\Sigma}_1^{-1} - \widehat{\Sigma}_0^{-1} \right) x_1 - 2 \left( \widehat{\mu}_1^T \widehat{\Sigma}_1^{-1} - \widehat{\mu}_0^T \widehat{\Sigma}_0^{-1} \right) x_1 + \left( \widehat{\mu}_0^T \widehat{\Sigma}_0^{-1} \widehat{\mu}_0 - \widehat{\mu}_1^T \widehat{\Sigma}_1^{-1} \widehat{\mu}_1 \right) + \log \left( \frac{1 - \widehat{p}}{\widehat{p}} \right) \geq 0 \right)$$

Predicting  $y_1$  for our dataset,

$$\begin{aligned}
y_1 &= \mathbb{I}\left(x_1^T(\widehat{\Sigma}_1^{-1} - \widehat{\Sigma}_0^{-1})x_1 - 2(\widehat{\mu}_1^T \widehat{\Sigma}_1^{-1} - \widehat{\mu}_0^T \widehat{\Sigma}_0^{-1})x_1 + (\widehat{\mu}_0^T \widehat{\Sigma}_0^{-1} \widehat{\mu}_0 - \widehat{\mu}_1^T \widehat{\Sigma}_1^{-1} \widehat{\mu}_1) + \log\left(\frac{1 - \widehat{p}}{\widehat{p}}\right) \geq 0\right) \\
&= \mathbb{I}\left(\begin{bmatrix} -3 \\ 4 \end{bmatrix}^T \left( \begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}^{-1} - \begin{bmatrix} 0.222 & 0.111 \\ 0.111 & 0.222 \end{bmatrix}^{-1} \right) \begin{bmatrix} -3 \\ 4 \end{bmatrix} - 2 \left( \begin{bmatrix} -3 \\ 3.666 \end{bmatrix}^T \begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}^{-1} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix}^T \begin{bmatrix} 0.222 & 0.111 \\ 0.111 & 0.222 \end{bmatrix}^{-1} \right) \begin{bmatrix} -3 \\ 4 \end{bmatrix} \\
&\quad + \left( \begin{bmatrix} -3 \\ 3.666 \end{bmatrix}^T \begin{bmatrix} 0.222 & 0.111 \\ 0.111 & 0.222 \end{bmatrix}^{-1} \begin{bmatrix} -3 \\ 3.666 \end{bmatrix} - \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix}^T \begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}^{-1} \begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix} \right) + \log\left(\frac{1 - 0.5}{0.5}\right) \geq 0 \right)
\end{aligned}$$

Predicting  $y_1$  for our dataset,

$$\begin{aligned}
y_1 &= \mathbb{I}\left(x_1^T\left(\widehat{\Sigma}_1^{-1}-\widehat{\Sigma}_0^{-1}\right)x_1-2\left(\widehat{\mu}_1^T\widehat{\Sigma}_1^{-1}-\widehat{\mu}_0^T\widehat{\Sigma}_0^{-1}\right)x_1+\left(\widehat{\mu}_0^T\widehat{\Sigma}_0^{-1}\widehat{\mu}_0-\widehat{\mu}_1^T\widehat{\Sigma}_1^{-1}\widehat{\mu}_1\right)+\log\left(\frac{1-\widehat{p}}{\widehat{p}}\right)\geq 0\right) \\
&= \mathbb{I}\left(\begin{bmatrix} -3 \\ 4 \end{bmatrix}^T\left(\begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}^{-1}-\begin{bmatrix} 0.222 & 0.111 \\ 0.111 & 0.222 \end{bmatrix}^{-1}\right)\begin{bmatrix} -3 \\ 4 \end{bmatrix}-2\left(\begin{bmatrix} -3 \\ 3.666 \end{bmatrix}^T\begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}^{-1}-\begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix}^T\begin{bmatrix} 0.222 & 0.111 \\ 0.111 & 0.222 \end{bmatrix}^{-1}\right)\begin{bmatrix} -3 \\ 4 \end{bmatrix} \\
&\quad +\left(\begin{bmatrix} -3 \\ 3.666 \end{bmatrix}^T\begin{bmatrix} 0.222 & 0.111 \\ 0.111 & 0.222 \end{bmatrix}^{-1}\begin{bmatrix} -3 \\ 3.666 \end{bmatrix}-\begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix}^T\begin{bmatrix} 0.666 & -0.333 \\ -0.333 & 0.222 \end{bmatrix}^{-1}\begin{bmatrix} -4.66 \\ 0.666 \end{bmatrix}\right)+\log\left(\frac{1-0.5}{0.5}\right)\geq 0) \\
&= \mathbb{I}(1.0 \geq 0)
\end{aligned}$$

Predicting  $y_1$  for our dataset,

$$\begin{aligned}
y_1 &= \mathbb{I}\left(x_1^T\left(\widehat{\Sigma}_1^{-1}-\widehat{\Sigma}_0^{-1}\right)x_1-2\left(\widehat{\mu}_1^T\widehat{\Sigma}_1^{-1}-\widehat{\mu}_0^T\widehat{\Sigma}_0^{-1}\right)x_1+\left(\widehat{\mu}_0^T\widehat{\Sigma}_0^{-1}\widehat{\mu}_0-\widehat{\mu}_1^T\widehat{\Sigma}_1^{-1}\widehat{\mu}_1\right)+\log\left(\frac{1-\widehat{p}}{\widehat{p}}\right)\geq 0\right) \\
&= \mathbb{I}\left(\begin{bmatrix}-3 \\ 4\end{bmatrix}^T\left(\begin{bmatrix}0.666 & -0.333 \\ -0.333 & 0.222\end{bmatrix}^{-1}-\begin{bmatrix}0.222 & 0.111 \\ 0.111 & 0.222\end{bmatrix}^{-1}\right)\begin{bmatrix}-3 \\ 4\end{bmatrix}-2\left(\begin{bmatrix}-3 \\ 3.666\end{bmatrix}^T\begin{bmatrix}0.666 & -0.333 \\ -0.333 & 0.222\end{bmatrix}^{-1}-\begin{bmatrix}-4.66 \\ 0.666\end{bmatrix}^T\begin{bmatrix}0.222 & 0.111 \\ 0.111 & 0.222\end{bmatrix}^{-1}\right)\begin{bmatrix}-3 \\ 4\end{bmatrix} \\
&\quad +\left(\begin{bmatrix}-3 \\ 3.666\end{bmatrix}^T\begin{bmatrix}0.222 & 0.111 \\ 0.111 & 0.222\end{bmatrix}^{-1}\begin{bmatrix}-3 \\ 3.666\end{bmatrix}-\begin{bmatrix}-4.66 \\ 0.666\end{bmatrix}^T\begin{bmatrix}0.666 & -0.333 \\ -0.333 & 0.222\end{bmatrix}^{-1}\begin{bmatrix}-4.66 \\ 0.666\end{bmatrix}\right)+\log\left(\frac{1-0.5}{0.5}\right)\geq 0\right) \\
&= \mathbb{I}(1.0 \geq 0) \\
y_1 &= 1
\end{aligned}$$



Prediction table for all datapoints:

Prediction table for all datapoints:

Label	$P(\hat{y} = 1 \hat{x})$	$P(\hat{y} = 0 \hat{x})$	Prediction	Actual
$y_1$	1.0	0	1	1
$y_2$	1.0	0	1	1
$y_3$	1.0	0	1	1
$y_4$	0	1.0	0	0
$y_5$	0	1.0	0	0
$y_6$	0.00019	0.99981	0	0

# Gaussian Naive Bayes Algorithm

