Bayesian Decision Theory

Computer Science- Pattern Recognition Prof. Dr. Dhahir A. Abdullah



• Covariance matrices for all of the classes are identical,

• But covariance matrices are arbitrary.



$$g_i(x) = -\frac{1}{2}(x-\mu_i)^t \sum_{i}^{-1} (x-\mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

$$\frac{d}{2}\ln 2\pi$$
 and $\frac{1}{2}\ln |\Sigma_i|$ are independent of i.

So,

$$g_i(x) = -\frac{1}{2}(x-\mu_i)^t \sum_{i=1}^{-1} (x-\mu_i) + \ln P(\omega_i)$$



Discriminant functions are

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$
 (linear discriminant)

where

$$\mathbf{w}_{i} = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{i}$$
$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_{i}^{T} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{i} + \ln P(w_{i}).$$



Decision boundaries can be written as

 $\mathbf{w}^T(\mathbf{x} - \mathbf{x}_0) = 0$

where

$$\mathbf{w} = \mathbf{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$
$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\ln(P(w_i)/P(w_j))}{(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T \mathbf{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j).$$

Hyperplane passes through x_0 but is not necessarily orthogonal to the line between the means.





Figure 6: Probability densities with equal but asymmetric Gaussian distributions. The decision hyperplanes are not necessarily perpendicular to the line connecting the means.

• Covariance matrices are different for each category.

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \sum_{i}^{-1} (x - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

only
$$\frac{d}{2} \ln 2\pi$$
 is independent of i.

So,

$$g_i(x) = -\frac{1}{2}(x-\mu_i)^t \sum_{i=1}^{-1} (x-\mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

Discriminant functions are

 $g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$ (quadratic discriminant)

where

$$\mathbf{W}_{i} = -\frac{1}{2} \boldsymbol{\Sigma}_{i}^{-1}$$
$$\mathbf{w}_{i} = \boldsymbol{\Sigma}_{i}^{-1} \boldsymbol{\mu}_{i}$$
$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_{i}^{T} \boldsymbol{\Sigma}_{i}^{-1} \boldsymbol{\mu}_{i} - \frac{1}{2} \ln |\boldsymbol{\Sigma}_{i}| + \ln P(w_{i}).$$

Decision boundaries are hyperquadrics.



Figure 7: Arbitrary Gaussian distributions lead to Bayes decision boundaries that are general hyperquadrics.



Figure 8: Arbitrary Gaussian distributions lead to Bayes decision boundaries that are general hyperquadrics.

EXAMPLE:

Consider a two-category classification problem

with two-dimensional feature vector $X = (x_1, x_2)^t$.

The two categories are ω_1 and ω_2 .

$$p(X|\omega_1) \sim N\left(\begin{bmatrix} -1\\1 \end{bmatrix}, \Sigma_1\right),$$

$$p(X|\omega_2) \sim N\left(\begin{bmatrix} 1\\0 \end{bmatrix}, \Sigma_2\right), \qquad \Sigma_1 = \begin{bmatrix} 1 & 0\\0 & 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 2 & 2\\2 & 3 \end{bmatrix}.$$

$$P(\omega_1) = P(\omega_2) = \frac{1}{2}$$

Calculate the Bayes decision boundary.

Find the discriminant function for the first class.

$$g_{1}(x) = -\ln(2\pi) - \frac{1}{2}\ln|\Sigma_{1}| - \frac{1}{2}(\bar{x} - \bar{\mu}_{1})^{t}\Sigma_{1}^{-1}(\bar{x} - \bar{\mu}_{1})$$
$$= -\ln(2\pi) - \frac{1}{2}\ln|\Sigma_{1}| - \frac{1}{2}[x_{1} - \mu_{11} \quad x_{2} - \mu_{12}]\Sigma_{1}^{-1}\begin{bmatrix}x_{1} - \mu_{11}\\x_{2} - \mu_{12}\end{bmatrix}$$

Find the discriminant function for the first class.

$$\ln(2\pi) - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} [x_1 - \mu_{11} \quad x_2 - \mu_{12}] \Sigma_1^{-1} \begin{bmatrix} x_1 - \mu_{11} \\ x_2 - \mu_{12} \end{bmatrix} =$$

$$= -\ln(2\pi) - \frac{1}{2} \ln 1 - \frac{1}{2} [x_1 + 1 \quad x_2 - 1] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 + 1 \\ x_2 - 1 \end{bmatrix}$$

$$= -\ln(2\pi) - \frac{1}{2} ((x_1 + 1)^2 + (x_2 - 1)^2)$$

$$= -\ln(2\pi) - \frac{1}{2} (x_1^2 + 2x_1 + 1 + x_2^2 - 2x_2 + 1)$$

Similarly, find the discriminant function for the second class.

$$g_{2}(x) = -\ln(2\pi) - \frac{1}{2}\ln|\Sigma_{2}| - \frac{1}{2}(\bar{x} - \bar{\mu}_{2})^{t}\Sigma_{2}^{-1}(\bar{x} - \bar{\mu}_{2})$$

$$= -\ln(2\pi) - \frac{1}{2}\ln|\Sigma_{2}| - \frac{1}{2}[x_{1} - \mu_{21} \quad x_{2} - \mu_{22}]\Sigma_{2}^{-1}\begin{bmatrix}x_{1} - \mu_{21} \\ x_{2} - \mu_{22}\end{bmatrix}$$

$$= -\ln(2\pi) - \frac{1}{2}\ln 2 - \frac{1}{2}[x_{1} - 1 \quad x_{2}]\frac{1}{2}\begin{bmatrix}3 & -2 \\ -2 & 2\end{bmatrix}\begin{bmatrix}x_{1} - 1 \\ x_{2}\end{bmatrix}$$

$$= -\ln(2\pi) - \frac{1}{2}\ln 2 - \frac{1}{4}[3(x_{1} - 1) - 2x_{2} \quad -2(x_{1} - 1) + 2x_{2}]\begin{bmatrix}x_{1} - 1 \\ x_{2}\end{bmatrix}$$

$$= -\ln(2\pi) - \frac{1}{2}\ln 2 - \frac{1}{4}(3(x_{1} - 1)^{2} - 2(x_{1} - 1)x_{2} - 2(x_{1} - 1)x_{2} + 2x_{2}^{2})$$

$$= -\ln(2\pi) - \frac{1}{2}\ln 2 - \frac{1}{4}(3x_{1}^{2} - 6x_{1} + 3 - 4x_{1}x_{2} + 4x_{2} + 2x_{2}^{2})$$

The decision boundary:

$$g(x) = g_{1}(x) - g_{2}(x)$$

$$= -\ln(2\pi) - \frac{1}{2} (x_{1}^{2} + 2x_{1} + 1 + x_{2}^{2} - 2x_{2} + 1) + \ln(2\pi) + \frac{1}{2} \ln 2 - \frac{1}{4} (3x_{1}^{2} - 6x_{1} + 3 - 4x_{1}x_{2} + 4x_{2} + 2x_{2}^{2})$$

$$= -\frac{1}{2} (x_{1}^{2} + 2x_{1} + 1 + x_{2}^{2} - 2x_{2} + 1) + \frac{1}{2} \ln 2 + \frac{1}{4} (3x_{1}^{2} - 6x_{1} + 3 - 4x_{1}x_{2} + 4x_{2} + 2x_{2}^{2})$$

$$= -\frac{1}{2} x_{1}^{2} - x_{1} - \frac{1}{2} - \frac{1}{2} x_{2}^{2} + x_{2} - \frac{1}{2} + \frac{1}{2} \ln 2 + \frac{3}{4} x_{1}^{2} - \frac{6}{4} x_{1} + \frac{3}{4} - x_{1}x_{2} + x_{2} + \frac{1}{2} x_{2}^{2}$$

$$= -\frac{1}{2} x_{1}^{2} + \frac{3}{4} x_{1}^{2} - x_{1} - \frac{6}{4} x_{1} - x_{1}x_{2} + x_{2} + x_{2} + \frac{1}{2} x_{2}^{2} - \frac{1}{2} - \frac{1}{2} - \frac{1}{2} + \frac{3}{4} + \frac{1}{2} \ln 2$$

$$= \frac{1}{4} x_{1}^{2} - \frac{10}{4} x_{1} - x_{1}x_{2} + 2x_{2} - \frac{1}{4} + \frac{1}{2} \ln 2$$

The decision boundary:

$$g(x) = g_1(x) - g_2(x) = \frac{1}{4}x_1^2 - \frac{10}{4}x_1 - x_1x_2 + 2x_2 - \frac{1}{4} + \frac{1}{2}\ln 2$$

$$= x_1^2 - 10x_1 - 4x_1x_2 + 8x_2 - 1 + 2\ln 2$$

Using MATLAB we can draw the decision boundary:

 $g(x) = x_1^2 - 10x_1 - 4x_1x_2 + 8x_2 - 1 + 2\ln 2$

(to draw the decision boundary in MATLAB)
>> s = 'x^2-10*x-4*x*y+8*y-1+2*log(2)';
>> ezplot(s)

Using MATLAB we can draw the decision boundary:





Error Probabilities and Integrals

For the two-category case

$$P(error) = P(\mathbf{x} \in \mathcal{R}_2, w_1) + P(\mathbf{x} \in \mathcal{R}_1, w_2)$$

= $P(\mathbf{x} \in \mathcal{R}_2 | w_1) P(w_1) + P(\mathbf{x} \in \mathcal{R}_1 | w_2) P(w_2)$
= $\int_{\mathcal{R}_2} p(\mathbf{x} | w_1) P(w_1) d\mathbf{x} + \int_{\mathcal{R}_1} p(\mathbf{x} | w_2) P(w_2) d\mathbf{x}$.

Error Probabilities and Integrals

For the multicategory case

$$P(error) = 1 - P(correct)$$

= $1 - \sum_{i=1}^{c} P(\mathbf{x} \in \mathcal{R}_i, w_i)$
= $1 - \sum_{i=1}^{c} P(\mathbf{x} \in \mathcal{R}_i | w_i) P(w_i)$
= $1 - \sum_{i=1}^{c} \int_{\mathcal{R}_i} p(\mathbf{x} | w_i) P(w_i) d\mathbf{x}.$

Error Probabilities and Integrals



Figure 9: Components of the probability of error for equal priors and the non-optimal decision point x^* . The optimal point x_B minimizes the total shaded area and gives the Bayes error rate.

• Another measure of distance between two Gaussian distributions.

• found a great use in medicine, radar detection and other fields.

- Consider the two-category case and define
 - w₁: target is present,
 - w₂: target is not present.

Table 1: Confusion matrix.

		Assigned			
		w_1	w_2		
True	w_1	correct detection	mis-detection		
	w_2	false alarm	correct rejection		

- Mis-detection is also called false negative or Type I error.
- False alarm is also called false positive or Type II error.

If we use a parameter (e.g., a threshold) in our decision, the plot of these rates for different values of the parameter is called the *receiver operating characteristic* (ROC) curve.



Figure 10: Example receiver operating characteristic (ROC) curves for different settings of the system.

Let p_i be the probability that patient *i* will get a positive diagnosis (i.e., the patient is ill) and q_i be patient *i*'s probability of a positive test. The *prevalence*, *P*, of the positive diagnosis in the population² is theoretically $P = \text{mean}(p_i)$. The *level* of the test, *Q*, is $Q = \text{mean}(q_i)$. We also define P' = 1 - P and Q' = 1 - Q.

	Test result			
	Positive	Negative		
Diagnosis Positive	TP	FN	Р	
Negative	FP	TN	P'	
	Q	Q'	1	

• If both diagnosis and test are positive, it is called a true positive. The probability of a TP to occur is estimated by counting the true positives in the sample and divide by the sample size.

• If the diagnosis is positive and the test is negative it is called a false negative (FN).

• False positive (FP) and true negative (TN) are defined similarly.

The values described are used to calculate different measurements of the quality of the test.
The first one is sensitivity, SE, which is the probability of having a positive test among the patients who have a positive diagnosis.

SE = TP/(TP + FN) = TP/P.

• Specificity, SP, is the probability of having a negative test among the patients who have a negative diagnosis.

SP = TN/(FP + TN) = TN/P'.



Figure 1.1: A sample population of N = 95 patients. The minus signs denote patients with a negative diagnosis and the plus signs denotes patients with a positive diagnosis. The colour shows the result of a test. White is a negative test and black is a positive test.

	#Test result		
	Positive	Negative	
• Example (cont.): #Diagnosis Positive	30	3	33
Negative	20	42	62
	50	45	95

TP	Ξ	30/95		$0.316(\pm 0.048)$	P	=	33/95	=	$0.347(\pm 0.049)$
FN	=	3/95	=	$0.032(\pm 0.018)$	Q	=	50/95	=	$0.526(\pm 0.051)$
FP	=	20/95	=	$0.211(\pm 0.042)$	SE	_	30/33	_	$0.909(\pm 0.050)$
TN	=	42/95	=	$0.442(\pm 0.051)$	SP	=	42/62	=	$0.677(\pm 0.059)$



Bayes Decision Theory – Discrete Features

- Assume x = (x1.....xd) takes only m discrete values
- Components of x are binary or integer valued, x can take only one of m **discrete values**

*V*₁, *V*₂, ..., *V*_m

Bayes Decision Theory – Discrete Features (Decision Strategy)

Maximize class posterior using bayes
 decision rule

$$P(\omega_1 \mid \underline{x}) = \underline{P(x \mid \omega_j) \cdot P(\omega_j)} \\ \Sigma P(\underline{x} \mid \omega_j) \cdot P(\omega_j)$$

Decide ω_i if $P(\omega_i | \underline{x}) > P(\omega_i | \underline{x})$ for all $i \neq j$ or

Minimize the overall risk by selecting the action α_i = α_{*} : deciding ω₂
 α_{*} = arg min R(α_i / x)

Here, we consider a 2-category problem in which the components of the feature vector are binary-valued and conditionally independent (which yields a simplified decision rule):

$$x = (x_1, \dots, x_d)^t, x_i \in \{0, 1\}$$

We also assign the following probabilities (p and q) to each x_i in X:

$$p_i = \Pr[x_i = 1 | \omega_1]_{\text{and}} q_i = \Pr[x_i = 1 | \omega_2]$$

If $p_i > q_i$, we expect to x_i to be 1 more frequently when the state of nature is w_1 than when it is w_2 . If we assume *conditional independence*, we can write $P(X|w_i)$ as the product of probabilities for the components of **X**. The class conditional probabilities are then:

$$P(\mathbf{x}|\omega_1) = \prod_{i=1}^d p_i^{x_i} (1-p_i)^{1-x_i}$$
$$P(\mathbf{x}|\omega_2) = \prod_{i=1}^d q_i^{x_i} (1-q_i)^{1-x_i}$$

Since this is a two class problem the discriminant function $g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$ where:

$$g_1(x) = \log \left[p(\mathbf{x}|\omega_1) p(\omega_1) \right]_{\text{and}} g_2(x) = \log \left[p(\mathbf{x}|\omega_2) p(\omega_2) \right]$$

The likelihood ratio is therefore given by:

$$\frac{P(\mathbf{x}|\omega_1)}{P(\mathbf{x}|\omega_2)} = \prod_{i=1}^d \left(\frac{p_i}{q_i}\right)^{x_i} \left(\frac{1-p_i}{1-q_i}\right)^{1-x_i}$$

which yields the discriminant function as follows:

$$g(x) = \ln \frac{P(\mathbf{x}|\omega_1)}{P(\mathbf{x}|\omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)} = \sum_{i=1}^{d} \left[x_i \ln \frac{p_i}{q_i} + (1 - x_i) \ln \frac{1 - p_i}{1 - q_i} \right] + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

If we notice that this function is linear in x_i , we can rewrite it as a linear function of x_i

$$g(x) = \sum_{i=1}^d w_i x_i + w_0$$

where

$$w_i = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)}, i = 1,...,d$$

and

$$w_{0} = \sum_{i=1}^{d} \ln \frac{1 - p_{i}}{1 - q_{i}} + \ln \frac{P(\omega_{1})}{P(\omega_{2})}$$

The discriminant function g(x) will therefore indicate whether the current feature vector belongs to class 1 and class 2. It is important to note that w_0 and w_i are weights calculated for the linear discriminant. A decision boundary lies wherever g(x) = 0. This decision boundary can be a line, or hyper-plane depending upon the dimension of the feature space.



Let's consider a three dimensional binary feature vector $X = (x_1, x_2, x_3) = (0, 1, 1)$ that we will attempt to classify with one of the following classes:



and lets say that the prior probability for class 1 is $P(\omega_1) = 0.6$ while for class 2 is $P(\omega_2) = 0.4$. Hence, it is already evident that there is a bias towards class 1.

Additionally, we know that likelihoods of each independent feature is given by p and q where:

 $p_i = P(x_i=1|\omega_1)$ and $q_i = P(x_i=1|\omega_2)$

meaning that we know the probability (or likelihood) of each independent feature given each class - these values are known and given:

 $p = \{0.8, 0.2, 0.5\}$ and $q = \{0.2, 0.5, 0.9\}$

therefore, the discriminant function is $g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$ or by taking the log of both sides:

$$g(x) = \log \frac{p(\mathbf{x}|\omega_2)}{p(\mathbf{x}|\omega_2)} + \log \frac{p(\omega_1)}{p(\omega_2)}$$

however, since the problem definition assumes that X is independent, the discriminant function can be calculated by:

$$g(x) = \sum_{i=1}^{d} w_i x_i + w_c$$

with

$$w_i(x) = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)}$$
 $i = 1, \dots, d$

$$w_i(x) = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)}$$
 $i = 1, \dots, d$

$$w_1 = \ln \frac{0.8(1-0.2)}{0.2(1-0.8)} = 2.77; w_2 = \ln \frac{0.2(1-0.5)}{0.5(1-0.2)} = -1.39; w_3 = \ln \frac{0.5(1-0.9)}{0.9(1-0.5)} = -2.19$$

$$w_0 = \ln\left(\frac{0.6}{0.4}\right) + \ln\left(\frac{1-0.8}{1-0.2}\right) + \ln\left(\frac{1-0.2}{1-0.5}\right) + \ln\left(\frac{1-0.5}{1-0.9}\right) = 1.0986$$

 $g(x) = 2.77x_1 - 1.39x_2 - 2.19x_3 + 1.0986$

After inputting the x_i values into the discriminant function, the answer $g(\mathbf{x}) = -2.4849$. Therefore this belongs to class 2. Below is a plot of the decision boundary surface.



All points above the plane belong to class ω_2 since if X = (0, 1, 1), $g(\mathbf{x}) = -2.4849 < 0$.

APPLICATION EXAMPLE

Bit – matrix for machine – printed characters



 p_i is the probability that $X_i = 1$ for letter A,B,...

Summary

To minimize the overall risk, choose the action that minimizes the conditional risk R (α |x).
To minimize the probability of error, choose the class that maximizes the posterior probability P (wj |x).
If there are different penalties for misclassifying patterns from different classes, the posteriors must be weighted according to such penalties before taking action.
Do not forget that these decisions are the optimal ones under the assumption that the "true" values of the probabilities are known.