

## Bayes Classifiers

Winter 2023

Clayton Scott

## 1 Probabilistic Setting for Classification

Probability theory provides a rigorous setting for machine learning. The reader is assumed to be familiar with the following concepts from probability: Random variables (discrete and continuous), jointly distributed random variables, (joint) pdf/ pmf, multivariate Gaussian distribution, independence, conditional distributions, Bayes rule, and the laws of total probability and expectation.

Let  $\mathbf{X} = [X_1, \dots, X_d]^T \in \mathbb{R}^d$  be a *feature vector* for a classification problem, and let  $Y \in \{1, \dots, K\}$  be the associated class label. For example, each  $X_j$  may be a different attribute of an individual applying to receive a loan, and  $Y \in \{1, 2\}$  where  $Y = 1$  if the individual will make all of their payments on time, and  $Y = 2$  otherwise. We use capital letters for the feature vector and label to emphasize that these quantities are viewed as random variables.

In these notes, and throughout much of this course, we make the following assumption:

$\mathbf{X}$  and  $Y$  are *jointly distributed* random variables.

There are two ways in which we will think about the joint distribution. The first is in terms of the marginal distribution of  $Y$ , together with the conditional distribution of  $\mathbf{X}|Y$ . The second is in terms of the marginal distribution of  $\mathbf{X}$ , together with the conditional distribution of  $Y|\mathbf{X}$ . These two perspectives are equivalent but one is typically easier to work with depending on the setting. We will use both perspectives.

**Example 1.** Figure 1 show a joint distribution for a binary ( $K = 2$ ) classification problem with a  $d = 2$  dimensional feature space. In this case, the marginal distribution of  $Y$  is uniform on the values 1 and 2. The conditional distribution of  $\mathbf{X}|Y$  is

$$\begin{aligned}\mathbf{X}|Y = 1 &\sim N(\boldsymbol{\mu}_1, \Sigma) \\ \mathbf{X}|Y = 2 &\sim N(\boldsymbol{\mu}_2, \Sigma)\end{aligned}$$

where

$$\boldsymbol{\mu}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \boldsymbol{\mu}_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \text{ and } \Sigma = \begin{bmatrix} .9 & .4 \\ .4 & .3 \end{bmatrix}.$$

The marginal distribution of  $Y$  and the conditional distribution of  $\mathbf{X}|Y$  provide a complete characterization of the joint distribution. These quantities allow us to answer any probabilistic question about the data. For example, suppose we wish to know the probability that  $X_2 \leq 0.5$ . Then

$$\begin{aligned}\Pr(X_2 \leq 0.5) &= \Pr(X_2 \leq 0.5|Y = 1) \Pr(Y = 1) + \Pr(X_2 \leq 0.5|Y = 2) \Pr(Y = 2) \\ &= 0.5\Phi(0.5; 1, \sqrt{.3}) + 0.5\Phi(0.5; -1, \sqrt{.3}) \\ &= 0.5888\end{aligned}$$

where  $\Phi(x; \mu, \sigma)$  is the CDF of a  $\mathcal{N}(\mu, \sigma^2)$  random variable.

It is also possible to represent this joint distribution using the marginal distribution of  $\mathbf{X}$  and the conditional distribution of  $Y|\mathbf{X}$ . By the law of total probability, the marginal density of  $\mathbf{X}$  is

$$\frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma) + \frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_2, \Sigma)$$

where  $\phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$  is the pdf of a  $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$  random variable. By Bayes rule, the conditional distribution of  $Y|\mathbf{X} = \mathbf{x}$  is

$$\Pr(Y = 1|\mathbf{X} = \mathbf{x}) = \frac{\frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma)}{\frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma) + \frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_2, \Sigma)}, \quad \Pr(Y = 2|\mathbf{X} = \mathbf{x}) = \frac{\frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_2, \Sigma)}{\frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma) + \frac{1}{2}\phi(\mathbf{x}; \boldsymbol{\mu}_2, \Sigma)}.$$

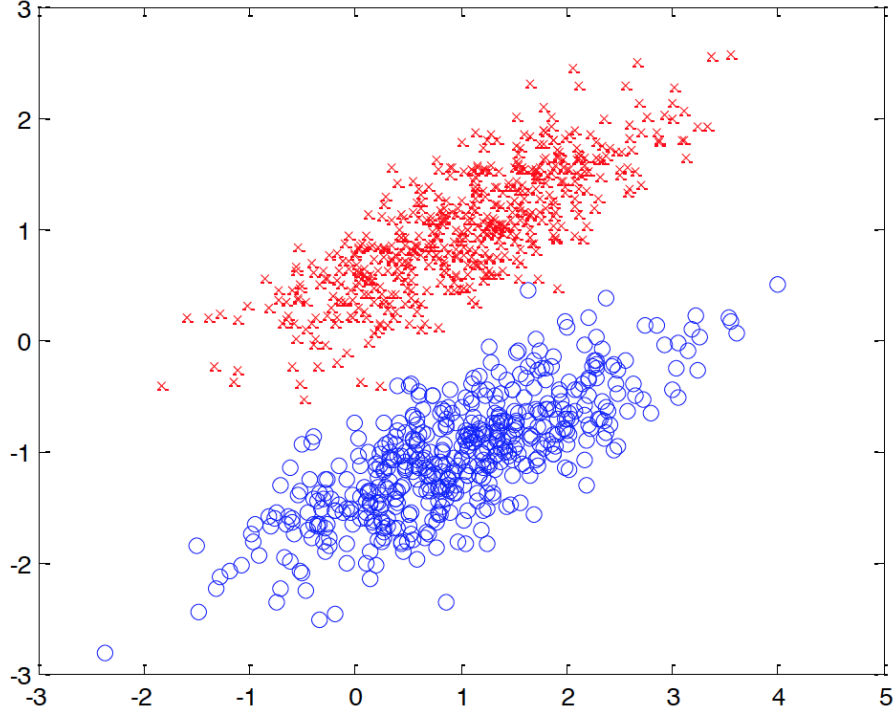


Figure 1: Two-dimensional data with Gaussian class-conditional distributions.

## 2 The Bayes Classifier

Given a joint distribution  $(\mathbf{X}, Y)$ , we may ask “what is the best possible classifier for that joint distribution?” To answer this question, we need to define what is meant by “best,” which means we need to define a performance measure.

A classifier is a function  $f : \mathbb{R}^d \rightarrow \{1, \dots, K\}$ . The most common performance measure in classification is the *probability of error*, defined by:

$$R(f) := \Pr(f(\mathbf{X}) \neq Y).$$

In words,  $R(f)$  is the probability that if we draw a random  $(\mathbf{X}, Y)$  from the joint distribution, that  $f$  predicts the wrong label for  $\mathbf{X}$ .

The probability of error is an example of a *risk*, a concept that will be defined later in the course. For now, we use the term “risk” to refer to the probability of error. The *Bayes risk* is the smallest risk of any classifier, and is denoted

$$R^* := \min_{\text{all } f} R(f),$$

where the minimization is over all possible classifiers. If  $R(f) = R^*$ ,  $f$  is called a *Bayes classifier*. It is important to note that  $R^*$  may not be zero. There are classification problems where the best classifier still makes errors. This occurs when there are feature vectors that could have arisen from multiple different classes.

We introduce the following terminology. The marginal distribution of  $Y$  is called the *prior distribution*, while the conditional distribution of  $Y|\mathbf{X} = \mathbf{x}$  is called the *posterior distribution*. The conditional distributions of  $\mathbf{X}|Y$  are called *class-conditional distributions*.

Let  $\pi_k = \Pr(Y = k)$  denote the prior class probabilities,  $g_k(\mathbf{x})$  the class-conditional pmfs/pdfs of  $\mathbf{X} | Y = k$ , and

$$\eta_k(\mathbf{x}) = \Pr(Y = k | \mathbf{X} = \mathbf{x})$$

the posterior class probabilities. Notice that  $\sum_i \pi_k = 1$  and that  $\forall \mathbf{x}, \sum_k \eta_k(\mathbf{x}) = 1$ .

**Theorem 1.** *The Bayes classifier is given by the equivalent formulas<sup>1</sup>*

$$f^*(\mathbf{x}) = \arg \max_{k=1, \dots, K} \eta_k(\mathbf{x}) \tag{1}$$

$$= \arg \max_{k=1, \dots, K} \pi_k g_k(\mathbf{x}). \tag{2}$$

where ties for the maximum may be broken arbitrarily.

*Proof.* Let  $\mathbf{1}_{\{\text{pred}\}}$  be 1 if the predicate **pred** is true and 0 otherwise. Then

$$\begin{aligned} R(f) &= \mathbb{E}[\mathbf{1}_{\{f(\mathbf{X}) \neq Y\}}] \\ &= \mathbb{E}_{\mathbf{X}} \mathbb{E}_{Y|\mathbf{X}}[\mathbf{1}_{\{f(\mathbf{X}) \neq Y\}}] \\ &= \mathbb{E}_{\mathbf{X}} \left[ \sum_{k=1}^K \eta_k(\mathbf{X}) \mathbf{1}_{\{f(\mathbf{X}) \neq k\}} \right] \end{aligned}$$

where the second step uses the law of total expectation. In the last expression, for a given  $\mathbf{X}$ , note that all of the indicator functions are 1 except for the one with  $k = f(\mathbf{X})$ . Therefore, to minimize the expression in brackets,  $f$  should be such that for each  $\mathbf{x}$ ,  $f(\mathbf{x}) = \arg \max_k \eta_k(\mathbf{x})$ . This proves (1). Furthermore, this argument works for any  $f$  such that  $f(\mathbf{x})$  outputs a  $k$  maximizing  $\eta_k(\mathbf{x})$ . Therefore ties for the maximum can be broken arbitrarily.

To prove (2), note that by Bayes rule,

$$\eta_k(\mathbf{x}) = \frac{\pi_k g_k(\mathbf{x})}{\sum_{\ell=1}^K \pi_{\ell} g_{\ell}(\mathbf{x})},$$

and that the denominator does not depend on  $k$ . □

### 3 Plug-in Classifiers

In machine learning, we don't know the joint distribution of  $(\mathbf{X}, Y)$ , and therefore we cannot determine the Bayes classifier. Instead, we observe training data

$$(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n) \stackrel{iid}{\sim} P_{\mathbf{X}Y},$$

where iid stands for *independent and identically distributed*. One approach to classification is to use the training data to estimate either the posterior probability function  $\eta_y(\mathbf{x})$ , or the class-conditional distributions and prior probabilities, and “plug” these estimates in to the corresponding formula for the Bayes classifier.

---

<sup>1</sup>Technically, the  $\arg \max$  operator returns a set, since multiple values might attain the maximum. In this case, any maximizer can be assigned and still yield an optimal classifier.

We will look at three examples of plug-in classifiers. Each is based on a different assumption about the joint distribution:

- Linear discriminant analysis assumes the class-conditional distributions are multivariate Gaussian with common covariance matrix.
- Naïve Bayes assumes the features  $X_i$  are independent when conditioned on the label.
- Logistic regression assumes  $\eta_y(\mathbf{x})$  is described by a logistic probability model.

## Exercises

1. (★) Let  $f$  be the constant classifier that always predicts a fixed class  $k$  regardless of  $\mathbf{x}$ . Given a simple formula for  $R(f)$ .
2. (★★) Consider multiclass classification with  $K$  classes. As a function of  $K$ , what is the largest possible value of the Bayes Risk  $R^*$ ? Justify your answer.
3. This problem concerns binary classification. For this problem, let the labels be 0 or 1, and define  $\eta(\mathbf{x}) := \Pr(Y = 1 \mid \mathbf{X} = \mathbf{x})$ .

- (a) (★★) Let the labels be  $\pm 1$ . Prove that for any classifier  $f$ ,

$$R(f) - R^* = \mathbb{E}_{\mathbf{X}}[|2\eta(\mathbf{X}) - 1| \mathbf{1}_{\{f(\mathbf{X}) \neq \text{sign}(2\eta(\mathbf{X}) - 1)\}}]. \quad (3)$$

The results says that the excess risk depends on how much (on average)  $\eta(\mathbf{X})$  deviates from  $\frac{1}{2}$  at points where  $h$  disagrees with the Bayes classifier. *Hint*: Refer to the proof for the Bayes classifier in the lecture notes, and for the binary case rewrite the proof in terms of  $\eta$ .

- (b) (★★) Prove that

$$R^* = \mathbb{E}_{\mathbf{X}}[\min\{\eta(\mathbf{X}), 1 - \eta(\mathbf{X})\}].$$

This result says that the closer  $\eta(X)$  is to  $\frac{1}{2}$  on average, the larger the Bayes risk.

- (c) (★★★) Let  $\alpha \in (0, 1)$ , and define the  $\alpha$ -weighted misclassification cost to be

$$R_{\alpha}(f) := \mathbb{E}[(1 - \alpha)\mathbf{1}_{\{f(\mathbf{X})=0, Y=1\}} + \alpha\mathbf{1}_{\{f(\mathbf{X})=1, Y=0\}}].$$

This assigns different weights to “false positives” and “false negatives.” Note that when  $\alpha = \frac{1}{2}$ ,  $2R_{\alpha}(f) = R(f)$ . Determine a formula for the Bayes classifier analogous to (1), and also prove a formula analogous to (3).