

# The Bayes Classifier

We have been starting to look at the supervised classification problem: we are given data  $(\mathbf{x}_i, y_i)$  for  $i = 1, \dots, n$ , where  $\mathbf{x}_i \in \mathbb{R}^d$ , and  $y_i \in \{1, \dots, K\}$ . In this section, we suppose that we know everything there is to know about the data (in a probabilistic sense): we assume that we know the *joint distribution* of  $(X, Y)$ . If we have full knowledge of the distribution, then we can design an optimal classifier without seeing any data at all.

We now make the mathematical setup completely concrete. The “feature vector”  $X$  is a random vector<sup>1</sup> in  $\mathbb{R}^d$ , and the “class label”  $Y$  is a discrete random scalar in  $\{1, \dots, K\}$ . When we say that we have a joint probability distribution for  $(X, Y)$ , it means that we have a rule that assigns probabilities to events that obeys the Kolmogorov axioms<sup>2</sup>. Given  $\mathcal{X} \subset \mathbb{R}^d$  and  $\mathcal{Y} \subset \{1, \dots, K\}$ , the joint distribution gives us a probability

$$\mathbb{P}[X \in \mathcal{X}, Y \in \mathcal{Y}] = \text{probability that } X \text{ is in } \mathcal{X} \text{ and } Y \text{ is in } \mathcal{Y}.$$

We will treat the entries in the feature vector as continuous-valued. Fixing the feature vector  $X$  at different points  $\mathbf{x}$  results in different conditional probability mass functions (pmfs) for the class label  $Y$ :

$$\eta_k(\mathbf{x}) := p_{Y|X}(k|\mathbf{x}) = \mathbb{P}[Y = k|X = \mathbf{x}]. \quad (1)$$

The pmf  $p_{Y|X}(k|\mathbf{x})$ , which is also called the **a posteriori distribution**, will play a central role in much of our discussion here and throughout the course. We encounter it so often that it is useful to give it the more compact notation  $\eta_k(\mathbf{x})$ .

---

<sup>1</sup>We use non-bold capital letters for all random variables in these notes, whether they are scalar-, vector-, matrix-, or whatever-valued.

<sup>2</sup>[https://en.wikipedia.org/wiki/Probability\\_axioms](https://en.wikipedia.org/wiki/Probability_axioms)

It is also useful to note that fixing the class label  $Y$  to different values  $y$  results in different conditional probability density functions (pdfs) for the feature vector  $X$ :  $f_{X|Y}(\mathbf{x}|y)$ , where

$$\mathrm{P}[X \in \mathcal{X} | Y = y] = \int_{\mathcal{X}} f_{X|Y}(\mathbf{x}|y) d\mathbf{x}.$$

$f_{X|Y}(\mathbf{x}|y)$  is the **class conditional distribution** of  $X$ , i.e., the distribution of  $X$  given that  $Y$  belongs to class  $y$ .

A classification rule or **classifier** is simply a function  $h : \mathbb{R}^d \rightarrow \{1, \dots, K\}$ ; that is, a function which takes a feature vector and returns a class label. We can specify this classification rule by **partitioning**  $\mathbb{R}^d$  into  $K$  regions  $\Gamma_1(h), \dots, \Gamma_K(h)$ , where  $\Gamma_k(h)$  is the set of point that  $h$  maps to  $k$ :

$$\Gamma_k(h) = \{\mathbf{x} \in \mathbb{R}^d : h(\mathbf{x}) = k\}.$$

We will judge the quality of a classifier by the probability that it makes a mistake:

$$R(h) = \mathrm{P}[h(X) \neq Y].$$

This is also called the **risk** of  $h$ , or the **probability of error**.

We can now ask a very well-defined question which has a clear-cut answer: What is the classifier that minimizes the probability of error? The answer is simple: given  $X = \mathbf{x}$ , choose the class label that maximizes the conditional probability in (1).

**Theorem:** Define the classifier

$$h^*(\mathbf{x}) = \arg \max_{k \in \{1, \dots, K\}} \eta_k(\mathbf{x}). \quad (2)$$

Then every other classifier  $h$  has

$$R(h) \geq R(h^*).$$

**Proof:** The optimality of  $h^*$  in (2) follows from carefully writing down the risk for an arbitrary classifier  $h$ , applying Bayes rule, and then showing that  $h^*$  optimizes the resulting expression. We start with an expression for  $1 - R(h)$ , which we will show is as *large* as possible when  $h = h^*$ :

$$\begin{aligned}
1 - R(h) &= \mathbb{P}[h(X) = Y] \\
&= \sum_{k=1}^K \mathbb{P}[Y = k] \cdot \mathbb{P}[h(X) = k | Y = k] \\
&= \sum_{k=1}^K \mathbb{P}[Y = k] \int_{\Gamma_k(h)} f_{X|Y}(\mathbf{x}|k) d\mathbf{x} \\
&= \sum_{k=1}^K \int_{\Gamma_k(h)} \mathbb{P}[Y = k] f_{X|Y}(\mathbf{x}|k) d\mathbf{x}.
\end{aligned}$$

By Bayes rule,

$$\eta_k(\mathbf{x}) = \frac{\mathbb{P}[Y = k] f_{X|Y}(\mathbf{x}|k)}{\sum_{\ell=1}^K \mathbb{P}[Y = \ell] f_{X|Y}(\mathbf{x}|\ell)}.$$

Note that the denominator is a function of  $\mathbf{x}$  that is independent of  $k$ ; it is in fact the marginal density  $f_X(\mathbf{x})$  for  $X$ . Using this and the fact that the regions  $\Gamma_k(h)$  are disjoint, we can continue the string of equalities:

$$1 - R(h) = \int_{\mathbb{R}^d} \left( \sum_{k=1}^K 1_{\Gamma_k(h)}(\mathbf{x}) f_X(\mathbf{x}) \eta_k(\mathbf{x}) \right) d\mathbf{x},$$

where  $1_{\mathcal{A}}(\mathbf{x})$  is the indicator function

$$1_{\mathcal{A}}(\mathbf{x}) = \begin{cases} 1, & \mathbf{x} \in \mathcal{A}, \\ 0, & \mathbf{x} \notin \mathcal{A}. \end{cases}$$

The way we choose  $h^*$  in (2) chooses the regions so that the function inside the integral above is as large as possible; it is clear that

$$\sum_{k=1}^K 1_{\Gamma_k(h)}(\mathbf{x}) f_X(\mathbf{x}) \eta_k(\mathbf{x}) \leq \sum_{k=1}^K 1_{\Gamma_k(h^*)}(\mathbf{x}) f_X(\mathbf{x}) \eta_k(\mathbf{x}),$$

for all  $\mathbf{x} \in \mathbb{R}^d$ . Thus

$$\begin{aligned} 1 - R(h) &\leq \int_{\mathbb{R}^d} \left( \sum_{k=1}^K 1_{\Gamma_k(h)}(\mathbf{x}) f_X(\mathbf{x}) \eta_k(\mathbf{x}) \right) d\mathbf{x} \\ &= 1 - R(h^*), \end{aligned}$$

and so  $R(h^*) \leq R(h)$ .

## The nearest neighbor classifier

We have just seen that the Bayes classifier is optimal. Unfortunately, it requires complete knowledge of the conditional probability mass function  $\eta_k(\mathbf{x})$ . In the context of machine learning, this is not a reasonable assumption. The **nearest neighbor classifier** is an *extremely* simple alternative. For any  $\mathbf{x}$ , we simply find the closest point  $\mathbf{x}_i$  in the training set and then assign  $\mathbf{x}$  the same label as its nearest neighbor.

This is an incredibly simple rule, but perhaps somewhat surprisingly we can show that as  $n \rightarrow \infty$ , i.e., as the size of our training data grows, this simple classifier is near-optimal. To see this, we will consider the risk of the nearest neighbor classifier  $h^{\text{NN}}$  conditioned on  $X = \mathbf{x}$  and compare this to the risk of the Bayes classifier  $h^*$ .

To make our discussion simpler, we will restrict our attention to the case of binary classification where  $y_i \in \{0, 1\}$ . We first note that the

risk of the Bayes classifier  $h^*$  conditioned on  $X = \mathbf{x}$  is given by

$$R^*(\mathbf{x}) := \mathbb{P}[Y \neq h^*(\mathbf{x}) | X = \mathbf{x}].$$

If  $h^*(\mathbf{x}) = 0$  then we have  $R^*(\mathbf{x}) = \mathbb{P}[Y = 1 | X = \mathbf{x}] = \eta_1(\mathbf{x})$ . Similarly, if  $h^*(\mathbf{x}) = 1$  we have  $R^*(\mathbf{x}) = \eta_0(\mathbf{x})$ . Since by definition  $h^*(\mathbf{x})$  selects the label that *maximizes*  $\eta_k(\mathbf{x})$ , we thus have that

$$R^*(\mathbf{x}) = \min\{\eta_0(\mathbf{x}), \eta_1(\mathbf{x})\}. \quad (3)$$

For the nearest neighbor classifier, note that

$$R^{\text{NN}}(\mathbf{x}) := \mathbb{P}[h^{\text{NN}}(\mathbf{x}) \neq Y | X = \mathbf{x}].$$

In our analysis, we will treat not only  $(X, Y)$  as random, but also the output  $h^{\text{NN}}(\mathbf{x})$  as random since it depends on the dataset, which is itself drawn at random from the same distribution as  $(X, Y)$ . This allows us to write

$$\begin{aligned} R^{\text{NN}}(\mathbf{x}) &= \mathbb{P}[Y = 0 | X = \mathbf{x}] \mathbb{P}[h^{\text{NN}}(\mathbf{x}) = 1 | X = \mathbf{x}] \\ &\quad + \mathbb{P}[Y = 1 | X = \mathbf{x}] \mathbb{P}[h^{\text{NN}}(\mathbf{x}) = 0 | X = \mathbf{x}]. \end{aligned} \quad (4)$$

If  $\mathbf{x}_{\text{NN}}$  denotes the nearest neighbor to  $\mathbf{x}$ , then we can write

$$\mathbb{P}[h^{\text{NN}}(\mathbf{x}) = k | X = \mathbf{x}] = \mathbb{P}[Y = k | X = \mathbf{x}_{\text{NN}}] = \eta_k(\mathbf{x}_{\text{NN}}).$$

As  $n \rightarrow \infty$ , we have that  $\|\mathbf{x}_{\text{NN}} - \mathbf{x}\| \rightarrow 0$ , and thus as  $n \rightarrow \infty$  we have

$$\eta_k(\mathbf{x}_{\text{NN}}) \rightarrow \eta_k(\mathbf{x}).$$

Plugging this back into (4) and simplifying, we obtain

$$\begin{aligned} R^{\text{NN}}(\mathbf{x}) &\rightarrow \eta_0(\mathbf{x})\eta_1(\mathbf{x}) + \eta_1(\mathbf{x})\eta_0(\mathbf{x}) \\ &= 2\eta_0(\mathbf{x})\eta_1(\mathbf{x}) \\ &\leq 2\min\{\eta_0(\mathbf{x}), \eta_1(\mathbf{x})\}, \end{aligned}$$

where the last inequality follows from the fact that both  $\eta_1(\mathbf{x})$  and  $\eta_2(\mathbf{x})$  are less than 1. Combining this with (3), this yields

$$\lim_{n \rightarrow \infty} R^{\text{NN}}(\mathbf{x}) \leq 2R^*(\mathbf{x}),$$

or in words, that asymptotically, the risk of the nearest neighbor classifier is at most twice the Bayes risk.

This can be strengthened by considering the more general  $k$ -nearest neighbors classifier. The idea here is to assign a label to  $\mathbf{x}$  by taking a majority vote over the  $k$  training points closest to  $\mathbf{x}$ . If  $R^{\text{kNN}}(\mathbf{x})$  denotes the risk of the  $k$ -nearest neighbor classifier, then one can show via a similar argument that

$$\lim_{n \rightarrow \infty} R^{\text{kNN}}(\mathbf{x}) \leq \left(1 + \sqrt{2/k}\right) R^*(\mathbf{x}).$$

Thus, by increasing  $k$  it is possible to drive this multiplicative constant arbitrarily close to 1. This results in a property known as **universal consistency**. Specifically, if  $R^*$  denotes the Bayes risk and  $R_n^{\text{kNN}}$  denotes the risk of the  $k$ -nearest neighbors classifier based on a dataset of size  $n$ , then one can show that as  $n \rightarrow \infty$ , if  $k \rightarrow \infty$  while  $k/n \rightarrow 0$ , then  $R_n^{\text{kNN}} \rightarrow R^*$ .

In words this is simply saying that for any possible distribution on the data, if we are given enough data eventually the risk of the  $k$ -nearest neighbor classifier will converge to the Bayes risk (i.e., to the optimal risk). Unfortunately (or fortunately, depending on your perspective), you might have to wait a very long time, so there is still a role for other machine learning algorithms to improve on this situation when we only have a finite amount of data.