

Classification using Bayes Rule

Bayes rule allows us to say something very strong about a particular (and very important) problem in statistical inference: classification.

The problem is stated as follows. We are given (random) data points $X_1 = \mathbf{x}_1, X_2 = \mathbf{x}_2, \dots$ that we want to assign to one of K different *classes*. Underlying this is a probabilistic model, where the class (denoted by Y_n) controls which distribution the data point X_n comes from:

$$X_n \sim f_X(\mathbf{x}|Y = k), \quad \text{if } X_n \text{ comes from class } k.$$

The goal is to come up a *classification rule* (or **classifier** for short) that maps data points X to their respective classes Y as accurately as possible.

More formally, a **classifier** is a function $h : \mathbb{R}^D \rightarrow \{1, \dots, K\}$; this function takes a data point $X \in \mathbb{R}^D$ (in machine learning, this is sometime referred to as a “feature vector”) and returns a discrete class label. This means that the function h is *piecewise constant*: an equivalent way of specifying h is through a **partition** of \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 :

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

[PICTURE]

If we have a known probabilistic model (that is accurate), there is a clear-cut best answer for choosing the best classification rule. Here, we are treating the (data point, class label) pair (X, Y) as coupled random variables — we observe X , and want to infer Y . If we know the conditional densities

$$X|Y = k \sim f_X(\mathbf{x}|Y = k) \quad (1)$$

for all $k = 1, \dots, K$, and the marginal probabilities

$$p_Y(k) = P(Y = k) \quad (2)$$

again for all $k = 1, \dots, K$, then we have completely specified the joint distribution of the continuous-valued random vector X and the discrete-valued random scalar Y . For a given classification rule h , the probability of h making an error is

$$R(h) = P(h(X) \neq Y).$$

We might also write this as the probability that a sample arrives in the wrong partition specified by h :

$$R(h) = P(X \notin \Gamma_Y(h)).$$

The quantity $R(h)$ is sometimes referred to as the **risk** associated with the rule h (hence the R for notation).

Given the probabilistic model specified by (1) and (2), the best classifier is simply to choose the class that has maximum conditional probability. This is codified in the following theorem.

Theorem: Define the classifier

$$h^*(\mathbf{x}) = \arg \max_{k \in \{1, \dots, K\}} P(Y = k | X = \mathbf{x}). \quad (3)$$

Then every other classifier h has

$$R(h) \geq R(h^*). \quad (4)$$

The classifier given by (3) is call the **Bayes classifier**, and its probability of error in (4) is called the **Bayes risk**.

Proof: The optimality of h^* in (3) 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 $h = h^*$:

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

By Bayes rule,

$$\begin{aligned} P(Y = k | X = \mathbf{x}) &= \frac{P(Y = k) f_X(\mathbf{x} | Y = k)}{f_X(\mathbf{x})} \\ &= \frac{P(Y = k) f_X(\mathbf{x} | Y = k)}{\sum_{\ell=1}^K P(Y = \ell) f_X(\mathbf{x} | Y = \ell)}. \end{aligned}$$

Note that the denominator is a function of \mathbf{x} that is independent of k . Using this and the fact that the regions $\Gamma_k(h)$ are disjoint, we can

continue the string of equalities:

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

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 (3) chooses the regions so that the function inside the integral above is as large as possible; it is clear¹ that

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

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

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

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

¹For a fixed \mathbf{x} , only one term in each of the sums is non-zero, and by definition h^* chooses the term that will be the largest.

Examples:

1. $D = 1$, $K = 2$, $P(Y = 1) = P(Y = 2) = 1/2$, and $X|Y = k \sim \text{Normal}(k, 1)$.

2. $D = 2$, $K = 2$, $P(Y = 1) = P(Y = 2) = 1/2$, and $X|Y = k \sim \text{Normal}(\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})$,
where $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2$ are given vectors in \mathbb{R}^D .

Nearest-neighbor classification

The Bayes classifier is great ... if we have perfect knowledge of the probability model. But this is only the case in a few special applications² — in fact, it's not an exaggeration to say that the entire field of machine learning sprouted up around the fact that it is entirely unrealistic that we know these distributions at all, and that everything should be learned from examples.

There is a super-easy data-driven rule for classification that requires no probability model, is very intuitive, and whose (asymptotic) performance comes within a factor of 2 of the optimal Bayes risk. It is called the *nearest neighbor rule*.

Instead of a probability model, we are given examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$. Our analysis will still assume that these are random, have some joint distribution, each pair (\mathbf{x}_n, y_n) is independent from the other pairs. Our classification algorithm, however, will operate without knowledge of this joint distribution.

Given an unlabeled data point \mathbf{x} , the nearest-neighbor rule simply assigns the same label as the closest example point to \mathbf{x} . If $\text{NI}(\mathbf{x})$ is the index of the closest point to \mathbf{x} ,

$$\text{NI}(\mathbf{x}) = \arg \min_{n=1, \dots, N} \|\mathbf{x} - \mathbf{x}_n\|_2,$$

then the nearest-neighbor classifier is

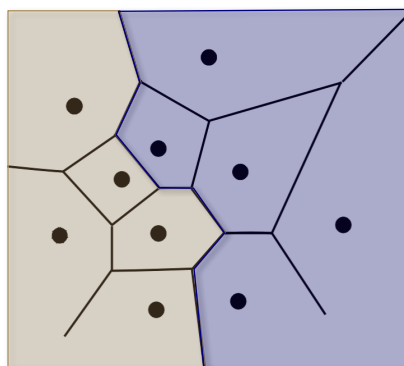
$$h_{\text{NN}}(\mathbf{x}) = y_{\text{NI}(\mathbf{x})}.$$

²Detectors for digital communications is one area where we are pretty confident that the class conditional densities are Gaussian, for example.

The example data point closest to \mathbf{x} is of course referred to as the “nearest neighbor” and is denoted

$$\text{NN}(\mathbf{x}) = \mathbf{x}_{\text{NI}(\mathbf{x})}.$$

The regions defined by the nearest neighbor classifier are a *Voronoi partition*.



As we can see, this rule is completely data-driven, and does not even make an attempt at estimating the conditional probability model $f_X(\mathbf{x}|Y = k)$. Nevertheless, its asymptotic performance is within a constant factor of the Bayes risk.

Theorem. As the number of examples N goes to infinity, the nearest-neighbor classifier $h_{\text{NN}}(\mathbf{x})$ has a probability of error that obeys

$$R(h_{\text{NN}}) \leq 2R(h^*),$$

where h^* is the Bayes classifier.

Proving this in detail is a bit technical, but we can get a good idea of how the argument works by just looking at the $K = 2$ case and

applying some basic conditional probability. First, we write the NN-risk as

$$P(h_{\text{NN}}(X) \neq Y) = \int_{\mathbf{x} \in \mathbb{R}^D} P(h_{\text{NN}}(X) \neq Y | X = \mathbf{x}) f_X(\mathbf{x}) \, d\mathbf{x}.$$

We now look at the probability we make a mistake at a fixed point $X = \mathbf{x}$. We have³

$$\begin{aligned} P(h_{\text{NN}}(X) \neq Y | X = \mathbf{x}) &= P(Y = 1 | X = \text{NN}(\mathbf{x})) P(Y = 2 | X = \mathbf{x}) \\ &\quad + P(Y = 2 | X = \text{NN}(\mathbf{x})) P(Y = 1 | X = \mathbf{x}) \end{aligned}$$

As $N \rightarrow \infty$, $\|\text{NN}(\mathbf{x}) - \mathbf{x}\|_2 \rightarrow 0$, so⁴

$$P(Y = k | X = \text{NN}(\mathbf{x})) \rightarrow P(Y = k | X = \mathbf{x}), \quad \text{for } k = 1, 2.$$

To ease the notation here, set

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

Thus we have

$$P(h_{\text{NN}}(X) \neq Y | X = \mathbf{x}) \rightarrow 2\eta_1(\mathbf{x})\eta_2(\mathbf{x}).$$

This relation holds pointwise at all $\mathbf{x} \in \mathbb{R}^D$.

Now we integrate. Let Γ_1^* and Γ_2^* be the regions in \mathbb{R}^D associated with where the Bayes classifier takes values 1 and 2 respectively:

$$\Gamma_k^* = \{\mathbf{x} : h^*(\mathbf{x}) = k\}, \quad k = 1, 2.$$

³Implicit here is the fact that the example data points $\{(X_n, Y_n)\}$ are independent from the unlabeled point (X, Y) we are analyzing.

⁴Obviously, this is where we are being mathematically imprecise. For the full argument, see Chapter 5 of *A Probabilistic Theory of Pattern Recognition* by Devroye, Györfi, and Lugosi (Springer, 1996).

The probability of error of the nearest-neighbor classifier is thus

$$\begin{aligned}
 R(h_{\text{NN}}) &= \mathbb{P}(h_{\text{NN}}(X) \neq Y) \\
 &\rightarrow 2 \int_{\mathbf{x} \in \Gamma_1^*} \eta_1(\mathbf{x})\eta_2(\mathbf{x}) \, d\mathbf{x} + 2 \int_{\mathbf{x} \in \Gamma_2^*} \eta_1(\mathbf{x})\eta_2(\mathbf{x}) \, d\mathbf{x} \\
 &\leq 2 \int_{\mathbf{x} \in \Gamma_1^*} \eta_2(\mathbf{x}) \, d\mathbf{x} + 2 \int_{\mathbf{x} \in \Gamma_2^*} \eta_1(\mathbf{x}) \, d\mathbf{x} \\
 &= 2R(h^*)
 \end{aligned}$$

where the first inequality comes from the fact that both $\eta_1(\mathbf{x}), \eta_2(\mathbf{x}) \leq 1$, and the second by the definition of the Bayes classifier.