

UMC203: AI and ML

Naman Mishra

January 2024

Contents

0	The Course	1	
0.1	References	1	
1		1	
1.1	Probability Review	2	
1.2	Multivariate Gaussians	3	
1.3	How Good is the Bayes Classifier?	5	
1.4	Bayes' Decision Theory	6	Lecture
			01: Tue
			09 Jan
			'24
0	The Course		

Instructor: Prof. Chiranjib Bhattacharyya

Office: CSA, 254

Office hours: TBD

Lecture hours: TuTh 10:00–11:20

0.1 References

- (i) *Pattern Classification* by Duda, Hart, and Stork
- (ii) *Probabilistic Theory of Pattern Recognition* by Devroye, Györfi, and Lugosi
- (iii) *Pattern Recognition and Machine Learning* by Bishop

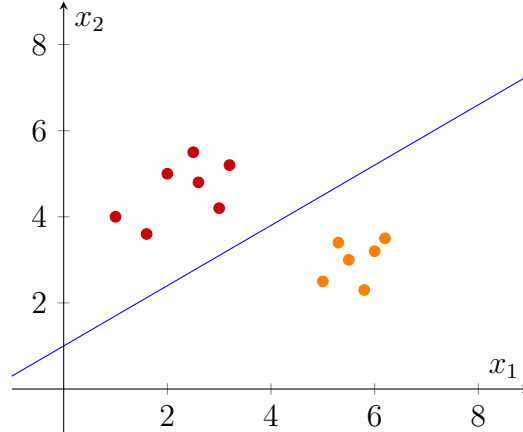
1

Consider a machine which can measure the diameter of any fruit placed on it. Can the machine distinguish between an apple and an orange? Now

suppose the machine also has the *capacity* to measure the weight of the fruit. Can it distinguish between an apple and an orange now?

$$\text{Fruit} \mapsto (x_1, x_2) \mapsto \{\text{Apple}, \text{Orange}\}$$

where x_1 is the diameter and x_2 is the weight. These are called *features*.



How do we measure how good a classifier is? This example has very few data points, so error is zero. Data is expensive, so accurate testing is expensive.

Let h be a classifier. We want to measure how good h is. We consider a random variable (of as yet unknown distribution) and compute the probability of error.

We consider this slightly more formally. Let the training data be

$$\mathcal{D} = \{(x^i, y^i)\}$$

where $x \in \mathbb{R}^2$ and $y \in \{1, -1\}$, where -1 and 1 represent apples and oranges respectively. Then h is a function from \mathbb{R}^2 to $\{-1, 1\}$. We wish to measure the probability $\Pr(h(X) \neq Y)$.

1.1 Probability Review

Suppose a coin is given with unknown probability of heads p . How do we estimate p ? We flip the coin n times and count the number of heads n_H . Then we estimate p as $\hat{p} = \frac{n_H}{n}$.

The rationale behind this is the weak/strong law of large numbers.

Theorem 1.1 (Weak Law of Large Numbers). Let X_1, X_2, \dots be i.i.d. random variables with mean μ . Then for any $\varepsilon > 0$,

$$\Pr \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| > \varepsilon \right) \rightarrow 0$$

as $n \rightarrow \infty$.

Theorem 1.2 (Strong Law of Large Numbers). Let X_1, X_2, \dots be i.i.d. random variables with mean μ . Then

$$\Pr \left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu \right) = 1$$

We have made several assumptions here. We know the structure of the problem at hand. For instance, we know that there is exactly one coin tossed each time. We have assumed that the coin tosses are independent. We have assumed that the probability of heads is the same for each toss.

Suppose we know the following in our earlier experiment:

- $\Pr(Y = 1)$
- $\Pr(Y = -1)$
- $\Pr(X = x \mid Y = 1)$
- $\Pr(X = x \mid Y = -1)$

Let $\eta(x) = \Pr(Y = 1 \mid X = x)$ given by Bayes' rule. Our rule for classification is

$$h(x) = \begin{cases} 1 & \text{if } \eta(x) > \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

This is called the *Bayes classifier*.

1.2 Multivariate Gaussians

$$f(x) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

where x is a d -dimensional column vector (so that the exponent is a scalar), μ is the mean, and Σ is the covariance matrix.

$$E[X] = \int_{x \in \mathbb{R}^d} x f(x) dx$$

$$\Sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)]$$

Reading assignment: Let A be a square symmetric real-valued matrix. What is known about the eigenvalues and positive definiteness of A ?

Definition 1.3 (Positive definite). A matrix $A_{n \times n}$ is *positive definite* if $u^\top A u > 0$ for all $u \in \mathbb{R}^n \setminus \{0\}$.

Exercise 1.4. Compute the Bayes classifier under the assumption that X under class 1 and class 2 has multivariate Gaussian distribution with means μ_1 and μ_2 with same covariance matrix Σ .

Lecture
02: Thu
11 Jan
'24

We come back to apples and oranges.

$$\begin{aligned}\mathcal{X} &\subseteq \mathbb{R}^d && \text{instance space} \\ \mathcal{Y} &= \{-1, 1\} && \text{label space} \\ \mathcal{D} &= \{(x_i, y_i) \mid x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}\end{aligned}$$

We also know *priors*

$$\Pr(Y = 1) = p_1 \qquad \Pr(Y = -1) = 1 - p_1 =: p_2$$

and *class condition distributions*

$$\Pr(X = x \mid Y = 1) = f_1(x) \qquad \Pr(X = x \mid Y = -1) = f_2(x)$$

Remarks. We will always write probabilities like this, but understand them to be densities whenever appropriate.

From Bayes' rule, we have the *posterior* $\eta: \mathcal{X} \rightarrow [0, 1]$ defined by

$$\begin{aligned}\eta(x) := \Pr(Y = 1 \mid X = x) &= \frac{\Pr(X = x \mid Y = 1) \Pr(Y = 1)}{\Pr(X = x)} \\ &= \frac{f_1(x)p_1}{f_1(x)p_1 + f_2(x)p_2}\end{aligned}$$

We can then define the *Bayes classifier* as

$$\begin{aligned}h^*(x) &:= \text{sgn}(2\eta(x) - 1) \\ &= \begin{cases} 1 & \text{if } \eta(x) > \frac{1}{2}, \\ -1 & \text{otherwise} \end{cases} \\ &= \begin{cases} 1 & \text{if } f_1(x)p_1 > f_2(x)p_2, \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

For the specific case of multivariate Gaussians, *i.e.*, f_1 and f_2 of the form

$$N(x \mid \mu, C) = \frac{1}{\sqrt{(2\pi)^d \det(C)}} \exp\left(-\frac{1}{2}(x - \mu)^\top C^{-1}(x - \mu)\right),$$

with same covariance C but different means μ_1 and μ_2 , we write

$$h^*(x) = \begin{cases} 1 & \text{if } \log \frac{\eta(x)}{1-\eta(x)} > 0 \\ -1 & \text{otherwise} \end{cases}$$

Now

$$\begin{aligned} \log \frac{\eta(x)}{1-\eta(x)} &= \log \frac{f_1(x)p_1}{f_2(x)p_2} \\ &= \log \frac{p_1}{p_2} - \frac{1}{2} ((x - \mu_1)^\top C^{-1} (x - \mu_1)) \\ &\quad + \frac{1}{2} ((x - \mu_2)^\top C^{-1} (x - \mu_2)) \\ &= \log \frac{p_1}{p_2} + (\mu_1 - \mu_2)^\top C^{-1} x - \frac{1}{2} (\mu_1^\top C^{-1} \mu_1 - \mu_2^\top C^{-1} \mu_2) \\ &= w^\top x - b \end{aligned}$$

where $w = C^{-1}(\mu_1 - \mu_2)$ (since C is symmetric) and b is something. Thus $h^*(x) = \text{sgn}(w^\top x - b)$.

Remarks. $w^\top x = b$ is a hyperplane in \mathbb{R}^d , dividing the space into two half-spaces: $w^\top x < b$ and $w^\top x > b$. So a line is a very good guess for a classifier!

Exercise 1.5. Examine the special case of $C_1 = \sigma_1^2 I$ and $C_2 = \sigma_2^2 I$.

Solution. We have

$$\log \frac{f_1(x)}{f_2(x)} = d \log \frac{\sigma_2}{\sigma_1} - \frac{1}{2\sigma_1^2} \|x - \mu_1\|^2 + \frac{1}{2\sigma_2^2} \|x - \mu_2\|^2$$

Thus we choose class 1 when

$$d \log \sigma_1 + \frac{1}{2\sigma_1^2} < d \log \sigma_2 + \frac{1}{2\sigma_2^2}$$

and class 2 otherwise.

1.3 How Good is the Bayes Classifier?

We wish to compute the error $\Pr(h(X) \neq Y)$ for some rule h .

$$\begin{aligned} \Pr(h(X) \neq Y) &= E_{XY} \mathbf{1}_{h^*(X) \neq Y} \\ &= E_X E_{Y|X} \mathbf{1}_{h^*(X) \neq Y} \end{aligned}$$

but

$$\begin{aligned} E_{Y|X=x} \mathbf{1}_{h^*(X) \neq Y} &= \begin{cases} 1 - \eta(x) & \text{if } h(x) = 1 \\ \eta(x) & \text{if } h(x) = -1 \end{cases} \quad (*) \\ &= \eta(x) \mathbf{1}_{h^*(x)=-1} + (1 - \eta(x)) \mathbf{1}_{h^*(x)=1} \end{aligned}$$

It is clear from (*) that whenever $\eta(x) > 1 - \eta(x)$, setting $h(x) = 1$ minimizes the error, and whenever $\eta(x) < 1 - \eta(x)$, setting $h(x) = -1$ minimizes the error.

More rigorously, upon comparing with h^* ,

$$\begin{aligned} E_{Y|X}(\mathbf{1}_{h(X) \neq Y} - \mathbf{1}_{h^*(X) \neq Y}) &= E_{Y|X}(\mathbf{1}_{h^*(X)=Y} - \mathbf{1}_{h^*(X) \neq Y}) \\ &= \eta(x)(\mathbf{1}_{h^*(X)=1} - \mathbf{1}_{h^*(X)=-1}) \\ &\quad + (1 - \eta(x))(\mathbf{1}_{h^*(X)=-1} - \mathbf{1}_{h^*(X)=1}) \\ &= (2\eta(x) - 1)(\mathbf{1}_{h^*(X)=1} - \mathbf{1}_{h^*(X)=-1}) \\ &= (2\eta(x) - 1)(2\mathbf{1}_{h^*(X)=1} - 1) \end{aligned}$$

The second term is 1 when the first term is positive, and -1 when it is negative.

Thus

$$\begin{aligned} E_{XY}(\mathbf{1}_{h(X) \neq Y} - \mathbf{1}_{h^*(X) \neq Y}) &= E_X E_{Y|X}(\mathbf{1}_{h(X) \neq Y} - \mathbf{1}_{h^*(X) \neq Y}) \\ &= E_X |2\eta(x) - 1| \geq 0. \end{aligned}$$

This proves that the Bayes classifier is the classifier with the lowest probability of error.

(This is theorem 2.1 in DGL.)

Lecture
03: Tue
16 Jan
'24

1.4 Bayes' Decision Theory

We have a $x \in \mathbb{R}^d$ with label $y \in \{-1, 1\}$. We predict $\hat{y} \in \{-1, 1\}$. We have a loss function $\ell: \{-1, 1\} \times \{-1, 1\} \rightarrow \mathbb{R}_{\geq 0}$. We wish to minimize the expected loss

$$R(h) = E_{XY} \ell(h(X), Y).$$

This is minimized by

$$\tilde{h}(x) = \arg \min_h E_{Y|X=x} \ell(h(x), Y).$$

For a given instance, this rule chooses the label which yields the minimum loss.

Now

$$\begin{aligned} E_{Y|X=x} \ell(1, Y) &= \ell(1, 1) \Pr(Y = 1 | X = x) + \ell(1, -1) \Pr(Y = -1 | X = x) \\ &= \ell(1, 1) \eta(x) + \ell(1, -1) (1 - \eta(x)) \end{aligned}$$

Similarly

$$E_{Y|X=x}\ell(-1, Y) = \ell(-1, 1)\eta(x) + \ell(-1, -1)(1 - \eta(x))$$

\tilde{h} minimises the loss if whenever $E_{Y|X=x}\ell(1, Y) < E_{Y|X=x}\ell(-1, Y)$, we choose $\tilde{h}(x) = 1$, and $\tilde{h}(x) = -1$ otherwise.

If $\ell(1, 1) = \ell(-1, -1) = 0$ and $\ell(1, -1) = \ell(-1, 1) = 1$, this reduces to the Bayes classifier.

Let us now extend the Bayes classifier to multiple classes. We have

$$\begin{aligned}\mathcal{X} &\subseteq \mathbb{R}^d && \text{instance space} \\ \mathcal{Y} &= \{1, \dots, M\} && \text{label space} \\ \mathcal{D} &= \{(x_i, y_i) \mid x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}\end{aligned}$$

We also know priors

$$p_i = \Pr(Y = i)$$

and class condition distributions

$$f_i(x) = \Pr(X = x \mid Y = i)$$

We define the posteriors $\eta_i: \mathcal{X} \rightarrow [0, 1]$ by

$$\eta_i(x) = \Pr(Y = i \mid X = x) = \frac{f_i(x)p_i}{\sum_{j=1}^M f_j(x)p_j}$$

and the Bayes classifier $h^*: \mathcal{X} \rightarrow \mathcal{Y}$ by

$$h^*(x) = \arg \max_{i \in \mathcal{Y}} \eta_i(x).$$