

1 Warm-up: Bayes classifier for scalar Gaussian mixtures

Let $(X_i, Y_i)_{1 \leq i \leq n}$ be independent variables in $\mathbb{R} \times \{0, 1\}$. Assume that $\mathbb{P}(Y_1 = 0) = 1/2$. Assume also that the distribution of X_1 given $\{Y_1 = 0\}$ (resp. $\{Y_1 = 1\}$) is Gaussian with mean μ_0 (resp. μ_1) and variance 1. The probability density function of X_1 is written g . Write

$$g_0 : x \mapsto (2\pi)^{-1/2} \exp(-(x - \mu_0)^2/2) \quad \text{and} \quad g_1 : x \mapsto (2\pi)^{-1/2} \exp(-(x - \mu_1)^2/2).$$

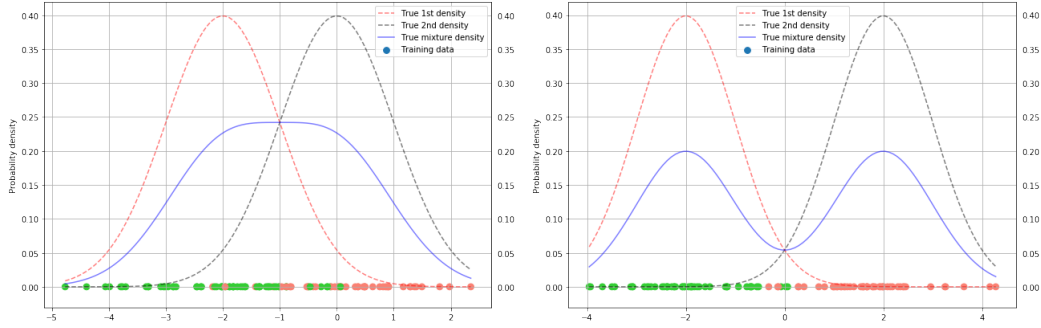


Figure 1: Samples and density when $\mu_0 = -2$ et $\mu_1 = 0$ (left) and $\mu_0 = -2$ and $\mu_1 = 2$ (right).

1. Provide an expression of a classifier h_* minimizing $h \mapsto \mathbb{P}(h(X) \neq Y)$.

The classifier h_ such that $h_*(X) = 1$ if and only if $\mathbb{P}(Y = 1|X) > \mathbb{P}(Y = 0|X)$ minimizes the missclassification error:*

$$h_* \in \operatorname{Argmin}_{h: \mathbb{R} \rightarrow \{0,1\}} \{\mathbb{P}(h(X) \neq Y)\}.$$

2. Using Bayes rule, show that h_* depends only on g_1/g_0 .

By Bayes formula, $\mathbb{P}(Y = 1|X) = \mathbb{P}(Y = 1)g_1(X)/g(X)$, which yields

$$\frac{\mathbb{P}(Y = 1|X)}{\mathbb{P}(Y = 0|X)} = \frac{g_1(X)}{g_0(X)}.$$

Then, $h_(X) = 1$ if and only if $g_1(X)/g_0(X) > 1$.*

3. Show that the Bayes classifier uses the mean between μ_0 and μ_1 to classify samples.

$h_(X) = 1$ if and only if $\log g_1(X) - \log g_0(X) > 0$, so that, assuming without loss of generality that $\mu_1 > \mu_0$:*

$$\begin{aligned} h_*(X) = 1 &\Leftrightarrow (X - \mu_0)^2 - (X - \mu_1)^2 > 0, \\ &\Leftrightarrow 2(\mu_1 - \mu_0)X + \mu_0^2 - \mu_1^2 > 0, \\ &\Leftrightarrow X > \frac{\mu_1^2 - \mu_0^2}{2(\mu_1 - \mu_0)}, \\ &\Leftrightarrow X > \frac{\mu_1 + \mu_0}{2}. \end{aligned}$$

This criterion can lead to very poor performance if means are close (see Figure 1).

2 Bayes classifier

2.1 Uniform distributions

Assume that $(X, Y) \in \mathbb{R} \times \{0, 1\}$ is defined on $(\Omega, \mathcal{F}, \mathbb{P})$ with $\mathbb{P}(Y = 1) = \pi \in (0, 1)$. Assume that conditionally on $\{Y = 0\}$ (resp. $\{Y = 1\}$) X has a uniform distribution on $[0, \theta]$ with $\theta \in (0, 1)$ (resp. on $[0, 1]$). Compute $\eta(X) = \mathbb{P}(Y = 1|X)$.

Let g be the probability density function of X . For any measurable set A ,

$$\begin{aligned} \mathbb{P}(X \in A) &= \mathbb{P}(Y = 0)\mathbb{P}(X \in A|Y = 0) + \mathbb{P}(Y = 1)\mathbb{P}(X \in A|Y = 1), \\ &= (1 - \pi)\theta^{-1} \int \mathbb{1}_A(x) \mathbb{1}_{[0, \theta]}(x) dx + \pi \int \mathbb{1}_A(x) \mathbb{1}_{[0, 1]}(x) dx, \\ &= \int \mathbb{1}_A(x) \{ (1 - \pi)\theta^{-1} \mathbb{1}_{[0, \theta]}(x) + \pi \mathbb{1}_{[0, 1]}(x) \} dx. \end{aligned}$$

Therefore, $g : x \mapsto (1 - \pi)\theta^{-1} \mathbb{1}_{[0, \theta]}(x) + \pi \mathbb{1}_{[0, 1]}(x)$. Then, using Bayes rules and writing g_1 the probability density of the distribution of X given $\{Y = 1\}$,

$$\eta(X) = \mathbb{P}(Y = 1|X) = \frac{\mathbb{P}(Y = 1)g_1(X)}{g(X)} = \frac{\pi \mathbb{1}_{[0, 1]}(X)}{(1 - \pi)\theta^{-1} \mathbb{1}_{[0, \theta]}(X) + \pi \mathbb{1}_{[0, 1]}(X)}.$$

2.2 Weighted risk

Assume that $(X, Y) \in \mathbb{R} \times \{0, 1\}$ is defined on $(\Omega, \mathcal{F}, \mathbb{P})$. Using $\omega_0, \omega_1 > 0$, with $\omega_0 + \omega_1 = 1$, we consider the weighted risk:

$$R(h) = \mathbb{E}[2\omega_Y \mathbb{1}_{Y \neq h(X)}].$$

Compute a classifier h_* minimizing $h \mapsto R(h)$ and $R(h_*)$.

For all classifiers h , writing $\eta(X) = \mathbb{P}(Y = 1|X)$,

$$\begin{aligned} R(h) &= \mathbb{E}[2\omega_Y \mathbb{1}_{Y \neq h(X)}] = \mathbb{E}[2\omega_Y \mathbb{1}_{Y=1} \mathbb{1}_{h(X)=0} + 2\omega_Y \mathbb{1}_{Y=0} \mathbb{1}_{h(X)=1}], \\ &= \mathbb{E}[2\omega_1 \mathbb{1}_{Y=1} \mathbb{1}_{h(X)=0} + 2\omega_0 \mathbb{1}_{Y=0} \mathbb{1}_{h(X)=1}], \\ &= \mathbb{E}[2\omega_1 \eta(X) \mathbb{1}_{h(X)=0} + 2\omega_0 (1 - \eta(X)) \mathbb{1}_{h(X)=1}], \end{aligned}$$

Therefore, choosing $h_* : x \mapsto \mathbb{1}_{\omega_1 \eta(x) \geq \omega_0 (1 - \eta(x))}$ yields,

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

Then, by definition, for all $x \in \mathbb{R}^d$,

$$h_*(x) = 1 \Leftrightarrow \omega_1 \eta(x) \geq \omega_0 (1 - \eta(x))$$

and

$$2\omega_1 \eta(x) \mathbb{1}_{h_*(x)=0} + 2\omega_0 (1 - \eta(x)) \mathbb{1}_{h_*(x)=1} = 2(\omega_1 \eta(x)) \wedge (\omega_0 (1 - \eta(x))).$$

This yields

$$R(h_*) = 2\mathbb{E}[(\omega_1 \eta(X)) \wedge (\omega_0 (1 - \eta(X)))].$$

3 Additional exercises

3.1 Bayes classifier: excess risk

Let $(X, Y) \in \mathbb{R}^d \times \{0, 1\}$ be random variables defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P})$. For any classifier $h : \mathcal{X} \rightarrow \{0, 1\}$, define its classification error by

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

The classifier h_* defined by:

$$h_*(x) = \mathbb{1}_{\eta(x) \geq \frac{1}{2}}$$

where

$$\eta(X) = \mathbb{P}(Y = 1|X) ,$$

minimizes $h \mapsto R(h)$.

1. Prove that

$$R(h_*) = \mathbb{E}[\eta(X) \wedge (1 - \eta(X))] \leq \frac{1}{2} .$$

For all classifiers h , as h and Y take values in $\{0, 1\}$,

$$R(h) = \mathbb{E}[\mathbb{1}_{h(X) \neq Y}] = \mathbb{E}[h(X)(1 - Y) + (1 - h(X))Y] .$$

As $\mathbb{E}[Y|X] = \eta(X)$ this yields,

$$R(h) = \mathbb{E}[h(X)(1 - \eta(X)) + (1 - h(X))\eta(X)]$$

and

$$R(h_*) = \mathbb{E}[h_*(X)(1 - \eta(X)) + (1 - h_*(X))\eta(X)] = \mathbb{E}[\eta(X) \wedge (1 - \eta(X))] .$$

2. Prove that for all classifiers h , the excess risk is given by

$$R(h) - R(h_*) = \mathbb{E}[|1 - 2\eta(X)| |h(X) - h_*(X)|] .$$

By the previous question, for all classifiers h ,

$$\begin{aligned} R(h) - R(h_*) &= \mathbb{E}[(h(X) - h_*(X))(1 - \eta(X)) + (h_*(X) - h(X))\eta(X)] , \\ &= \mathbb{E}[(h(X) - h_*(X))(1 - 2\eta(X))] . \end{aligned}$$

By definition of h_* , $h(X) - h_*(X)$ and $1 - 2\eta(X)$ have the same sign so that

$$R(h) - R(h_*) = \mathbb{E}[|1 - 2\eta(X)| |h(X) - h_*(X)|] .$$

3.2 Plug-in classifier

Let $(X, Y) \in \mathbb{R}^d \times \{-1, 1\}$ be random variables defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P})$. For any classifier $h : \mathcal{X} \rightarrow \{-1, 1\}$, define its classification error by

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

The classifier h_* defined by:

$$h_*(x) = \text{sign}(\eta(x) - 1/2) ,$$

where

$$\eta(X) = \mathbb{P}(Y = 1|X) ,$$

minimizes $h \mapsto R(h)$. Given n independent couples $\{(X_i, Y_i)\}_{1 \leq i \leq n}$ with the same distribution as (X, Y) , an empirical surrogate for h_* is obtained from a possibly nonparametric estimator $\hat{\eta}_n$ of η :

$$\hat{h}_n : x \mapsto \text{sign}(\hat{\eta}_n(x) - 1/2) .$$

1. Prove that for any classifier $h : \mathcal{X} \rightarrow \{-1, 1\}$,

$$\mathbb{P}(Y \neq h(X)|X) = (2\eta(X) - 1)\mathbb{1}_{h(X)=-1} + 1 - \eta(X)$$

and

$$R(h) - R(h_*) = 2\mathbb{E} \left[\left| \eta(X) - \frac{1}{2} \right| \mathbb{1}_{h(X) \neq h_*(X)} \right].$$

For all classifiers h ,

$$\begin{aligned} \mathbb{P}(Y \neq h(X)|X) &= \mathbb{P}(Y = -1, h(X) = 1|X) + \mathbb{P}(Y = 1, h(X) = -1|X), \\ &= \mathbb{1}_{h(X)=1}\mathbb{P}(Y = -1|X) + \mathbb{1}_{h(X)=-1}\mathbb{P}(Y = 1|X), \\ &= \mathbb{1}_{h(X)=-1}(2\eta(X) - 1) + 1 - \eta(X). \end{aligned}$$

Then,

$$R(h) - R(h_*) = \mathbb{E} \left[(\mathbb{1}_{h(X)=-1} - \mathbb{1}_{h_*(X)=-1}) (2\eta(X) - 1) \right] = 2\mathbb{E} \left[\left| \eta(X) - \frac{1}{2} \right| \mathbb{1}_{h(X) \neq h_*(X)} \right].$$

2. Prove that

$$|\eta(x) - 1/2| \mathbb{1}_{\hat{h}_n(x) \neq h_*(x)} \leq |\eta(x) - \hat{\eta}_n(x)| \mathbb{1}_{\hat{h}_n(x) \neq h_*(x)},$$

where

$$\hat{h}_n : x \mapsto \text{sign}(\hat{\eta}_n(x) - 1/2).$$

Deduce that

$$R(\hat{h}_n) - R(h_*) \leq 2\mathbb{E}[|\eta(X) - \hat{\eta}_n(X)|^2]^{1/2}.$$

Note that, for all $x \in \mathbb{R}^d$, $\hat{h}_n(x) \neq h_*(x)$ if and only if i) $\eta(x) > 1/2$ and $\hat{\eta}_n(x) \leq 1/2$ or ii) $\eta(x) \leq 1/2$ and $\hat{\eta}_n(x) > 1/2$. If $\eta(x) > 1/2$ and $\hat{\eta}_n(x) \leq 1/2$, then $|\eta(x) - \hat{\eta}_n(x)| = \eta(x) - \hat{\eta}_n(x) \geq \eta(x) - 1/2$. On the other hand, if $\eta(x) \leq 1/2$ and $\hat{\eta}_n(x) > 1/2$, $|\eta(x) - \hat{\eta}_n(x)| = \hat{\eta}_n(x) - \eta(x) \geq 1/2 - \eta(x)$. Therefore, for all $x \in \mathbb{R}^d$,

$$|\eta(x) - 1/2| \mathbb{1}_{\hat{h}_n(x) \neq h_*(x)} \leq |\eta(x) - \hat{\eta}_n(x)| \mathbb{1}_{\hat{h}_n(x) \neq h_*(x)}.$$

By the first question and Cauchy-Schwarz inequality,

$$\begin{aligned} R(\hat{h}_n) - R(h_*) &= 2\mathbb{E} \left[|\eta(X) - 1/2| \mathbb{1}_{h_*(X) = \hat{h}_n(X)} \right], \\ &\leq 2\mathbb{E} \left[|\eta(X) - \hat{\eta}_n(X)| \mathbb{1}_{\hat{h}_n(X) \neq h_*(X)} \right], \\ &\leq 2\mathbb{E}[|\eta(X) - \hat{\eta}_n(X)|^2]^{1/2}. \end{aligned}$$