

Statistical learning theory

Joris LIMONIER

February 25, 2022

Contents

1	Inclass exercise January 12, 2022	1
1.1	Exercise 1	1
1.2	Exercise 2	2
2	Inclass exercise January 21, 2022	2
2.1	Exercise 1	2
3	Inclass exercise January 28, 2022	4
3.1	Exercise 1	4
3.1.1	Question 1	4
3.1.2	Question 2	4
3.1.3	Question 3	5
4	In-class exercise February 4, 2022	6
4.1	Question 1	6
4.2	Question 2	6
5	In-class exercise February 22, 2022	6
5.1	Question 1	6
5.2	Question 2	8
6	In class exercises February 25, 2022	9
6.1	Question 1	9
6.2	Question 2	9
6.3	Question 3	10
6.4	Exercise	11

1 Inclass exercise January 12, 2022

1.1 Exercise 1

Show that

$$\mathbb{E} \left[\hat{\mathcal{R}}_S(h) \right] = \mathcal{R}_{D,f}(h) \quad (1)$$

$$\begin{aligned}
\mathbb{E} [\hat{\mathcal{R}}_S(h)] &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{h(x_i) \neq y_i} \right] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\mathbf{1}_{h(x_i) \neq y_i}] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{P}(h(x_i) \neq y_i) \\
&= \frac{1}{n} n \mathbb{P}(h(x_i) \neq y_i) \\
&= \mathbb{P}(h(x_i) \neq y_i) \\
&= \mathbb{P}(h(x_i) \neq f(x)) \\
&= \mathcal{R}_{D,f}(h)
\end{aligned}$$

1.2 Exercise 2

We must prove that the variance of $\hat{\mathcal{R}}_S(h) \rightarrow 0$

$$\begin{aligned}
\text{Var} [\hat{\mathcal{R}}_S(h)] &= \text{Var} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{h(x_i) \neq y_i} \right] \\
&= \text{Var} \frac{1}{n^2} \left[\sum_{i=1}^n \mathbf{1}_{h(x_i) \neq y_i} \right]
\end{aligned}$$

Let the Z_i be defined as follows:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{h(x_i) \neq f(x_i)} =: \frac{1}{n} \sum_{i=1}^n Z_i$$

(not finished, see lecture 1 slides)

2 Inclass exercise January 21, 2022

2.1 Exercise 1

Set $g(x) = \mathbb{P}(Y = 1 \mid X = x)$. We define the Bayes optimal predictor as:

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & g(x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

Question 1. Let $h : \mathcal{X} \rightarrow \{0, 1\}$ be a classifier. Show that

$$\begin{aligned}
&\mathbb{P}(h(X) \neq Y \mid X = x) \\
&= g(x) \cdot \mathbb{P}(h(X) = 0 \mid X = x) + (1 - g(x)) \cdot \mathbb{P}(h(X) = 1 \mid X = x)
\end{aligned}$$

$$\begin{aligned}
& g(x) \cdot \mathbb{P}(h(X) = 0 \mid X = x) + (1 - g(x)) \cdot \mathbb{P}(h(X) = 1 \mid X = x) \\
&= \mathbb{P}(Y = 1 \mid X = x) \cdot \mathbb{P}(h(X) = 0 \mid X = x) \\
&+ (1 - \mathbb{P}(Y = 1 \mid X = x)) \cdot \mathbb{P}(h(X) = 1 \mid X = x) \\
&= \mathbb{P}(Y = 1 \cap h(X) = 0 \mid X = x) \\
&+ \mathbb{P}(h(X) = 1 \mid X = x) - \mathbb{P}(Y = 1 \cap h(X) = 1 \mid X = x) \\
&= \mathbb{P}(Y = 1 \cap h(X) = 0 \mid X = x) + \mathbb{P}(Y = 0 \cap h(X) = 1 \mid X = x) \\
&= \mathbb{P}(h(X) \neq Y \mid X = x)
\end{aligned}$$

Question 2. Deduce that

$$\mathbb{P}(f_D(X) \neq Y \mid X = x) = \min(g(x), 1 - g(x))$$

$$\begin{aligned}
& \mathbb{P}(f_D(X) \neq Y \mid X = x) \\
&= \begin{cases} \mathbb{P}(1 \neq Y \mid X = x), & g(x) \geq 1/2 \\ \mathbb{P}(0 \neq Y \mid X = x), & g(x) < 1/2 \end{cases} \\
&= \begin{cases} 1 - g(x), & g(x) \geq 1 - g(x) \\ g(x), & g(x) < 1 - g(x) \end{cases} \\
&= \min(g(x), 1 - g(x))
\end{aligned}$$

Question 3. Show that

$$\mathbb{P}(h(X) \neq Y \mid X = x) \geq \mathbb{P}(f_D(x) \neq Y \mid X = x)$$

$$\begin{aligned}
\mathbb{P}(f_D(x) \neq Y \mid X = x) &= \min(g(x), 1 - g(x)) \\
&= \min(g(x), 1 - g(x)) \\
&\cdot (\mathbb{P}(h(X) = 0 \mid X = x) + \mathbb{P}(h(X) = 1 \mid X = x)) \\
&\leq g(x) \cdot (\mathbb{P}(h(X) = 0 \mid X = x) \\
&+ (1 - g(x)) \cdot \mathbb{P}(h(X) = 1 \mid X = x)) \\
&= \mathbb{P}(h(X) \neq Y \mid X = x)
\end{aligned}$$

Question 4. Prove that

$$\mathcal{R}_D(f_D) \leq \mathcal{R}_D(h)$$

$$\begin{aligned}
& \mathbb{P}(f_D(x) \neq Y \mid X = x) \leq \mathbb{P}(h(X) \neq Y \mid X = x) \\
\implies \mathbb{E}[\mathbb{P}(f_D(x) \neq Y \mid X = x)] &\leq \mathbb{E}[\mathbb{P}(h(X) \neq Y \mid X = x)] \\
\implies \mathcal{R}_D(f_D) &\leq \mathcal{R}_D(h)
\end{aligned}$$

3 Inclass exercise January 28, 2022

3.1 Exercise 1

Let Z be a random variable with a second moment such that $\mathbb{E}[Z] = \mu$ and $\text{Var}(Z) = \sigma^2$.

3.1.1 Question 1

Let $g : t \mapsto \mathbb{E}[(Z - t)^2]$. Show that g is minimum at $t = \mu$.

$$\begin{aligned} g(t) &= \mathbb{E}[(Z - t)^2] \\ &= \mathbb{E}[Z^2 + t^2 - 2tZ] \\ &= \mathbb{E}[Z^2] + \mathbb{E}[t^2] - \mathbb{E}[2tZ] \\ &= \mathbb{E}[Z^2] + t^2 - 2t\mathbb{E}[Z] \\ &= \sigma^2 - \mu^2 + t^2 - 2t\mu \\ &= \sigma^2 - \mu^2 + t^2 - 2t\mu \end{aligned}$$

We differentiate with respect to t :

$$\begin{aligned} \frac{\partial}{\partial t} g(t) &= 0 \\ \implies \frac{\partial}{\partial t} [\sigma^2 - \mu^2 + t^2 - 2t\mu] &= 0 \\ \implies 2t - 2\mu &= 0 \\ \implies t &= \mu \end{aligned}$$

3.1.2 Question 2

Assume $Z \in [a, b]$ almost surely. Use the previous question to show that

$$\text{Var}(Z) \leq \frac{(b - a)^2}{4}$$

$$\begin{aligned}
& g(\mu) \leq g(t) \\
\implies \text{Var}(Z) & \leq \mathbb{E}[(Z - t)^2] \\
\implies \text{Var}(Z) & \leq \frac{1}{4} \mathbb{E}[(2Z - a - b)^2] \\
\implies \text{Var}(Z) & \leq \frac{1}{4} \mathbb{E}[(Z - a) + (Z - b)]^2 \\
\implies \text{Var}(Z) & \leq \frac{1}{4} \mathbb{E}[(Z - a) - (b - Z)]^2 \\
\implies \text{Var}(Z) & \leq \frac{1}{4} \mathbb{E}[|Z - a| - |b - Z|]^2 \\
\implies \text{Var}(Z) & \leq \frac{1}{4} \mathbb{E}[|(Z - a) - (Z - b)|^2] \\
\implies \text{Var}(Z) & \leq \frac{1}{4} \mathbb{E}[|b - a|^2] \\
\implies \text{Var}(Z) & \leq \frac{(b - a)^2}{4}
\end{aligned}$$

3.1.3 Question 3

Let $Z_1, \dots, Z_n \sim Z$ be i.i.d. Use Chebyshev inequality to obtain a concentration inequality for

$$Z := \frac{1}{n} \sum_{i=1}^n Z_i$$

Chebyshev inequality:

$$\mathbb{P}(|Z - \mathbb{E}[Z]| \geq a) \leq \frac{\text{Var } Z}{a^2} \quad (2)$$

$$\begin{aligned}
\text{Var}(Z) &= \text{Var}\left(\frac{1}{n^2} \sum_{i=1}^n Z_i\right) \\
&= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(Z_i) \quad (Z_i \text{'s independent}) \\
&\leq \frac{1}{n^2} \sum_{i=1}^n \frac{(b - a)^2}{4} \\
&\leq \frac{(b - a)^2}{4n}
\end{aligned}$$

Then we apply (2):

$$\begin{aligned} \mathbb{P}(|Z - \mathbb{E}[Z]| \geq \varepsilon) &\leq \frac{\text{Var } Z}{\varepsilon^2} \\ \Rightarrow \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n Z_i - \mu\right| \geq \varepsilon\right) &\leq \frac{(b-a)^2}{4n\varepsilon^2} \end{aligned}$$

4 In-class exercise February 4, 2022

4.1 Question 1

We define our loss as:

$$\ell(y, y') = |y - y'|$$

Show:

$$\forall c \in \mathbb{R}, \begin{cases} |c| = \min_{a \geq 0} a \\ s.t. \quad a \geq c \\ \quad a \geq -c \end{cases}$$

A function study of $x \mapsto |x|$ gives the result.

4.2 Question 2

ERM consists in finding the following quantity:

$$\min_{w \in \mathbb{R}} \mathcal{R}_S(w) = \min_{w \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n |\langle w, x_i \rangle - y_i|$$

5 In-class exercise February 22, 2022

5.1 Question 1

Show that ERM with the logistic loss is equivalent to minimizing

$$F(w) = \sum_{i=1}^n \log(1 + \exp(-\tilde{y}_i \langle w, x_i \rangle))$$

with $\tilde{y}_i = \text{sign}(y_i - 0.5)$. Deduce that $\hat{\mathcal{R}}$ is a convex function of w .
We have:

$$\begin{aligned}
\ell(y, y_i) &= \begin{cases} -\log(1 - \hat{y}) & y = 0 \\ -\log(\hat{y}) & y = 1 \end{cases} \\
&= \begin{cases} -\log\left(1 - \frac{1}{1+e^{-w^T x_i}}\right) & y = 0 \\ -\log\left(\frac{1}{1+e^{-w^T x_i}}\right) & y = 1 \end{cases} \\
&= \begin{cases} -\log\left(\frac{e^{-w^T x_i}}{1+e^{-w^T x_i}}\right) & y = 0 \\ -\log\left(\frac{1}{1+e^{-w^T x_i}}\right) & y = 1 \end{cases} \\
&= \begin{cases} \log\left(\frac{1+e^{-w^T x_i}}{e^{-w^T x_i}}\right) & y = 0 \\ \log\left(1 + e^{-w^T x_i}\right) & y = 1 \end{cases} \\
&= \begin{cases} \log\left(1 + \frac{1}{e^{-w^T x_i}}\right) & y = 0 \\ \log\left(1 + e^{-w^T x_i}\right) & y = 1 \end{cases} \\
&= \begin{cases} \log\left(1 + e^{w^T x_i}\right) & y = 0 \\ \log\left(1 + e^{-w^T x_i}\right) & y = 1 \end{cases} \\
&= \log\left(1 + e^{-\tilde{y}_i w^T x_i}\right) \\
&= \log\left(1 + e^{(-1)^{y_i} w^T x_i}\right)
\end{aligned}$$

$$\begin{aligned}
&\frac{\partial^2}{\partial w^2} \log\left(1 + e^{(-1)^{y_i} w^T x_i}\right) \\
&= \frac{\partial}{\partial w} \frac{-y_i x_i e^{-\tilde{y}_i w^T x_i}}{1 + e^{-\tilde{y}_i w^T x_i}} \\
&= \frac{(y_i x_i)^2 e^{-\tilde{y}_i w^T x_i} (1 + e^{-\tilde{y}_i w^T x_i}) - (-y_i x_i e^{-\tilde{y}_i w^T x_i})^2}{(1 + e^{-\tilde{y}_i w^T x_i})^2} \\
&= \frac{x_i^2 e^{-\tilde{y}_i w^T x_i} (1 + e^{-\tilde{y}_i w^T x_i}) - x_i^2 e^{-2\tilde{y}_i w^T x_i}}{(1 + e^{-\tilde{y}_i w^T x_i})^2} \\
&= \frac{x_i^2 e^{-\tilde{y}_i w^T x_i}}{(1 + e^{-\tilde{y}_i w^T x_i})^2} \\
&\geq 0
\end{aligned}$$

5.2 Question 2

Compute the gradient of $\hat{\mathcal{R}}$ with respect to w .

Hint: show that $\phi'(z) = \phi(z)(1 - \phi(z))$

$$\begin{aligned}
 \frac{d}{dz}\phi(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\
 &= -\frac{-e^{-z}}{(1 + e^{-z})^2} \\
 &= \frac{1 - 1 + e^{-z}}{(1 + e^{-z})^2} \\
 &= \frac{1}{1 + e^{-z}} \frac{-1 + e^{-z}}{1 + e^{-z}} \\
 &= \frac{1}{1 + e^{-z}} \frac{e^{-z}}{1 + e^{-z}} \\
 &= \frac{1}{1 + e^{-z}} \frac{1 + e^{-z} - 1}{1 + e^{-z}} \\
 &= \phi(z)(1 - \phi(z))
 \end{aligned}$$

$$\begin{aligned}
 \hat{\mathcal{R}}(w) &= \sum_{i=1}^n \ell(y_i, \hat{y}_i) \\
 &= \sum_{i=1}^n -(1 - y_i) \log(1 - \hat{y}_i) - y_i \log \hat{y}_i \\
 &= \sum_{i=1}^n -(1 - y_i) \log(1 - \phi(w^T x_i)) - y_i \log \phi(w^T x_i)
 \end{aligned}$$

For some $1 \leq j \leq n$:

$$\frac{\partial}{\partial w_j} \ell(y, \hat{y}) = \frac{\partial}{\partial w_j} -(1 - y) \log(1 - \phi(w^T x_i)) - y_i \log \phi(w^T x_i)$$

Final result:

$$\hat{\mathcal{R}} = \sum_{i=1}^n (\phi(w^T x_i) - y_i) x_{ij}$$

6 In class exercises February 25, 2022

6.1 Question 1

Prove that the following function is a kernel:

$$k(x, y) = 2^{x+y}$$

Symmetry Trivial

Positive definiteness

$$\begin{aligned}\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j 2^{x_i+x_j} \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j 2^{x_i} 2^{x_j} \\ &= \sum_{i=1}^n c_i 2^{x_i} \sum_{j=1}^n c_j 2^{x_j} \\ &= \left[\sum_{i=1}^n c_i 2^{x_i} \right]^2 \\ &\geq 0\end{aligned}$$

6.2 Question 2

Prove that the following function is a kernel:

$$k(x, y) = (x^T y)^2$$

Symmetry Trivial

Positive definiteness

$$\begin{aligned}\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j (x_i^T x_j)^2 \\&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j (x_i^T x_j)(x_i^T x_j) \\&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j (x_j^T x_i)(x_i^T x_j) \\&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \operatorname{tr}(x_j^T x_i x_i^T x_j) \\&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \operatorname{tr}(x_i x_i^T x_j x_j^T) \\&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \langle x_i x_i^T, x_j x_j^T \rangle \\&= \sum_{i=1}^n c_i x_i x_i^T \sum_{j=1}^n c_j x_j x_j^T \\&= \left[\sum_{i=1}^n c_i x_i x_i^T \right]^2 \\&\geq 0\end{aligned}$$

6.3 Question 3

Prove that the following function is a kernel:

$$k(x, y) = \cos(x - y)$$

Symmetry Trivial

Positive definiteness

$$\begin{aligned}
& \sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \\
&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \cos(x_i - x_j) \\
&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j [\cos(x_i) \cos(x_j) + \sin(x_i) \sin(x_j)] \\
&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \cos(x_i) \cos(x_j) + \sum_{i=1}^n \sum_{j=1}^n c_i c_j \sin(x_i) \sin(x_j) \\
&= \left[\sum_{i=1}^n c_i \cos(x_i) \right]^2 + \left[\sum_{i=1}^n c_i \sin(x_i) \right]^2 \\
&\geq 0
\end{aligned}$$

6.4 Exercise

Show that the Gaussian kernel, which is given by:

$$k(x, y) := \exp \left(-\frac{\|x - y\|^2}{2\nu^2} \right)$$

is actually a kernel.

We have that:

$$\exp \left(-\frac{\|x - y\|^2}{2\nu^2} \right) = \lim_{n \rightarrow \infty} \sum_{p=1}^n \frac{1}{p!} \left(-\frac{\|x - y\|^2}{2\nu^2} \right)^p$$

$$\begin{aligned}
& \|x - y\|^2 = (x - y)^T (x - y) \\
\Rightarrow & \|x - y\|^2 = x^T x - x^T y - y^T x + y^T y \\
\Rightarrow & \|x - y\|^2 = \|x\|^2 - 2x^T y + \|y\|^2 \\
\Rightarrow & -\|x - y\|^2 = \underbrace{2x^T y}_{\text{kernel}} - \|x\|^2 - \|y\|^2
\end{aligned}$$

Let us show that $\exp(-\|x\|^2 - \|y\|^2)$ is a kernel:

$$\begin{aligned}
\exp(-\|x\|^2 - \|y\|^2) &= \exp(-\|x\|^2) \exp(-\|y\|^2) \\
&= \langle \exp(-\|x\|^2), \exp(-\|y\|^2) \rangle
\end{aligned}$$