Mehryar Mohri
Advanced Machine Learning 2018
Courant Institute of Mathematical Sciences
Homework assignment 2
April 30, 2018
Due: May 14, 2018


**A. Learning kernels**

In this problem, we will derive an alternative guarantee for learning kernels. We will use the notation adopted in class.

1. Using the results presented in class, prove the following equality:

$$\widehat{\mathfrak{R}}_S(H_1) = \frac{1}{m} \, \mathrm{E}\left[\max_{k \in [p]} \sqrt{\boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma}}\right].$$

   *Solution:* This follows immediately the equality at the bottom of slide 19. □

2. Compute $\mathrm{E}[\boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma}]$, for any $k \in [p]$.

   *Solution:* By definition, for any $k \in [p]$,

$$\begin{aligned}
\mathrm{E}[\boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma}] &= \mathrm{E}[\sum_{i,j=1}^{m} \sigma_i \sigma_j [\mathbf{K}_k]_{ij}] \\
&= \sum_{i,j=1}^{m} [\mathbf{K}_k]_{ij} \, \mathrm{E}[\sigma_i \sigma_j] \\
&= \sum_{i=1}^{m} [\mathbf{K}_k]_{ii} = \mathrm{Tr}[\mathbf{K}_k]. \qquad (\mathrm{E}[\sigma_i \sigma_j] = 0 \text{ for } i \neq j)
\end{aligned}$$

   □

3. Prove the following inequality for the empirical Rademacher complexity of $H_1$:

$$\widehat{\mathfrak{R}}_S(H_1) \leq \frac{1}{m} \sqrt{\max_k \mathrm{Tr}[\mathbf{K}_k] + m \lambda_{\max} \sqrt{\frac{\log p}{2}}},$$

where $\lambda_{\mathrm{max}}$ is the largest eigenvalue of a matrix $\mathbf{K}_k$, $k \in [p]$ (*hint*: you can use Jensen's inequality and the proof technique in Massart's Lemma).

*Solution:* By Jensen's inequality, in view of 1), we can write

$$\widehat{\mathfrak{R}}_S(H_1) = \frac{1}{m} \, \mathrm{E} \left[ \sqrt{\max_{k \in [p]} \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma}} \right] \leq \frac{1}{m} \sqrt{\mathrm{E} \left[ \max_{k \in [p]} \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \right]}.$$

Now, for any $t > 0$, we can write

$$\exp \left( t \, \mathrm{E} \left[ \max_{k \in [p]} \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \right] \right)$$

$$\leq \mathrm{E} \left[ \exp \left( t \max_{k \in [p]} \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \right) \right] \qquad \text{(Jensen's inequality)}$$

$$= \mathrm{E} \left[ \max_{k \in [p]} \exp \left( t \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \right) \right]$$

$$\leq \mathrm{E} \left[ \sum_{k \in [p]} \exp \left( t \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \right) \right]$$

$$= \sum_{k \in [p]} \mathrm{E} \left[ \exp \left( t \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \right) \right]$$

$$= \sum_{k \in [p]} e^{t \, \mathrm{Tr}[\mathbf{K}_k]} \, \mathrm{E} \left[ \exp \left( t \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} - t \, \mathrm{E}[\boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma}] \right) \right] \qquad \text{(in view of 2))}$$

$$\leq \sum_{k \in [p]} e^{t \, \mathrm{Tr}[\mathbf{K}_k]} e^{t^2 \lambda_{\mathrm{max}}^2 m^2 / 8} \qquad \text{(Hoeffding's inequality)}$$

$$\leq p \, e^{t \max_k \mathrm{Tr}[\mathbf{K}_k]} e^{t^2 \lambda_{\mathrm{max}}^2 m^2 / 8}.$$

When applying Hoeffding's inequality, we used the fact that

$$0 \leq \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \leq m \lambda_{\mathrm{max}}.$$

Taking the log of both sides gives

$$\mathrm{E} \left[ \max_{k \in [p]} \boldsymbol{\sigma}^\top \mathbf{K}_k \boldsymbol{\sigma} \right] \leq \max_k \mathrm{Tr}[\mathbf{K}_k] + \frac{\log p}{t} + \frac{t \lambda_{\mathrm{max}}^2 m^2}{8}.$$

Choosing $t$ to minimize the right-hand side ($t = \sqrt{8 \log p} / \lambda_{\mathrm{max}} m$) completes the proof. $\qquad \square$

## B. Model selection and convex surrogates

In class, we proved that the SRM technique benefits from very favorable learning guarantees. However, SRM requires solving multiple ERM problems, which in general are NP-hard problems. Here, we will discuss guarantees for using a convex surrogate loss instead of the original binary loss.

The hypotheses we consider are real-valued functions $h\colon X \to \mathbb{R}$. The sign of $h$ defines a binary classifier $f_h\colon X \to \{-1,+1\}$ defined for all $x \in X$ by $f_h(x) = 1_{h(x)\geq 0} - 1_{h(x)<0}$. The loss or error of $h$ at point $(x,y) \in X \times \{-1,+1\}$ is defined as the binary classification error of $f_h$:

$$1_{f_h(x)\neq y} = 1_{yh(x)<0} + 1_{h(x)=0\wedge y=-1} \leq 1_{yh(x)\leq 0}.$$

1. Show that, for any $h$, the generalization error of $h$ can be expressed as follows, where $\eta(x) = \mathbb{P}[y = +1|x]$ and where $D_X$ denote the marginal distribution over $X$:

   $$R(h) = \operatorname*{E}_{x\sim D_X} \left[\eta(x)1_{h(x)<0} + (1 - \eta(x))1_{h(x)\geq 0}\right].$$

   *Solution:*

   $$\begin{aligned}
   R(h) &= \operatorname*{E}_{(x,y)\sim D} \left[1_{f_h(x)\neq y}\right] \\
   &= \operatorname*{E}_{x\sim D_X} \left[\eta(x)1_{h(x)<0} + (1 - \eta(x))1_{h(x)>0} + (1 - \eta(x))1_{h(x)=0}\right] \\
   &= \operatorname*{E}_{x\sim D_X} \left[\eta(x)1_{h(x)<0} + (1 - \eta(x))1_{h(x)\geq 0}\right].
   \end{aligned}$$

   $\square$

2. Show that the Bayes classifier can be induced by $h^*$ defined for all $x \in X$ by $h^*(x) = \eta(x) - \frac{1}{2}$. We will denote by $R^*$ the Bayes error.

   *Solution:* In view of (1), the Bayes classifier can be defined as assigning label $+1$ to $x$ when $\eta(x) \geq \frac{1}{2}$, $-1$ and can therefore be induced by $h^*$.

   $\square$

3. Prove that the following equality holds for the excess error of any hypothesis $h\colon X \to \mathbb{R}$:

   $$R(h) - R^* = 2 \operatorname*{E}_{x\sim D_X} \left[|h^*(x)|\, 1_{h(x)h^*(x)\leq 0}\right].$$

*Solution:* For any $h$, we can write

$$R(h) = \mathop{\mathrm{E}}_{x \sim D_X} \left[ \eta(x) 1_{h(x)<0} + (1 - \eta(x)) 1_{h(x)\geq 0} \right]$$

$$= \mathop{\mathrm{E}}_{x \sim D_X} \left[ \eta(x) 1_{h(x)<0} + (1 - \eta(x))(1 - 1_{h(x)<0}) \right]$$

$$= \mathop{\mathrm{E}}_{x \sim D_X} \left[ [2\eta(x) - 1] 1_{h(x)<0} + (1 - \eta(x)) \right]$$

$$= \mathop{\mathrm{E}}_{x \sim D_X} \left[ 2h^*(x) 1_{h(x)<0} + (1 - \eta(x)) \right],$$

where we used for the last step the equation established earlier on. In view of that, for any $h$, the following holds:

$$R(h) - R(h^*) = \mathop{\mathrm{E}}_{x \sim D_X} \left[ 2[h^*(x)](1_{h(x)\leq 0} - 1_{h^*(x)\leq 0}) \right]$$

$$= \mathop{\mathrm{E}}_{x \sim D_X} \left[ 2[h^*(x)] \operatorname{sgn}(h^*(x)) 1_{(h(x)h^*(x)\leq 0)\wedge((h(x),h^*(x))\neq(0,0))} \right]$$

$$= 2 \mathop{\mathrm{E}}_{x \sim D_X} \left[ |h^*(x)| \, 1_{h(x)h^*(x)\leq 0} \right],$$

which completes the proof, since $R(h^*) = R^*$. $\qquad\square$

4. Let $\Phi \colon \mathbb{R} \to \mathbb{R}$ be a strictly convex and non-decreasing function so that for any $u \in \mathbb{R}$, $1_{u\leq 0} \leq \Phi(-u)$. For any $h$, define $\mathcal{L}_\Phi(h) = \mathrm{E}_{(x,y)\sim D} \left[ \Phi(-yh(x)) \right]$. Show that $\mathcal{L}_\Phi(h) = \mathrm{E}_{x\sim D_X}[L_\Phi(x, h(x))]$, where $L_\Phi(x, u) = \eta(x)\Phi(-u) + (1 - \eta(x))\Phi(u)$.

*Solution:* By definition,

$$\mathcal{L}_\Phi(h) = \mathop{\mathrm{E}}_{(x,y)\sim D} \left[ \Phi(-yh(x)) \right]$$

$$= \mathop{\mathrm{E}}_{x\sim D_X} \left[ \eta(x)\Phi(-h(x)) + (1 - \eta(x))\Phi(h(x)) \right].$$

$\qquad\square$

5. Let $h_\Phi^*$ be defined by $h_\Phi^*(x) = \operatorname{argmin}_{u\in[-\infty,+\infty]} L_\Phi(x, u)$. Prove that $f_{h_\Phi^*}$ is the Bayes classifier.

*Solution:* For a fixed $x \in X$, the minimum of $L_\Phi(x, u)$ is achieved for $-\eta(x)\Phi'(-u) + (1 - \eta(x))\Phi'(u) = 0$. Since $\Phi'$ is non-decreasing, we have $u > 0$ iff $\eta(x) > \frac{1}{2}$. $\qquad\square$

4

6. Assume that there exists $s \geq 1$ and $c > 0$ such that the following holds for all $x \in X$: $|h^*(x)|^s = |\eta(x) - \frac{1}{2}|^s \leq c^s [L_\Phi(x, 0) - L_\Phi(x, h_\Phi^*(x))]$. Use Jensen's inequality to prove the following

$$R(h) - R(h^*) \leq 2c \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ [\Phi(0) - L_\Phi(x, h_\Phi^*(x))] \, 1_{h(x)h^*(x) \leq 0} \right]^{\frac{1}{s}}.$$

*Solution:* By definition, we can write

$$R(h) - R(h^*)$$
$$= \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ |2\eta(x) - 1| \, 1_{h(x)h^*(x) \leq 0} \right]$$
$$\leq \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ |2\eta(x) - 1|^s \, 1_{h(x)h^*(x) \leq 0} \right]^{\frac{1}{s}} \qquad \text{(Jensen's ineq.)}$$
$$\leq 2c \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ [\Phi(0) - L_\Phi(x, h_\Phi^*(x))] \, 1_{h(x)h^*(x) \leq 0} \right]^{\frac{1}{s}} \qquad \text{(assumption)}.$$

$\square$

7. Prove that $\Phi(-2h^*(x)h(x)) \leq L_\Phi(x, h(x))$ (*hint:* use convexity).

*Solution:* Using the convexity of $\Phi$, we can write

$$\Phi(-2h^*(x)h(x)) = \Phi((1 - 2\eta(x))h(x))$$
$$= \Phi(\eta(x)(-h(x)) + (1 - \eta(x))h(x))$$
$$\leq \eta(x)\Phi((-h(x))) + (1 - \eta(x))\Phi(h(x)) = L_\Phi(x, h(x)).$$

$\square$

8. Use the previous inequalities to prove the following bound on the excess error in terms of the excess surrogate loss:

$$R(h) - R^* \leq 2c \left[ \mathcal{L}_\Phi(h) - \mathcal{L}_\Phi^* \right]^{\frac{1}{s}}.$$

*Solution:* Starting with the inequality previously proven, we can write

$$R(h) - R(h^*)$$

$$\leq 2c \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ \left[ \Phi(0) - L_\Phi(x, h_\Phi^*(x)) \right] 1_{h(x)h^*(x) \leq 0} \right]^{\frac{1}{s}}$$

$$\leq 2c \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ \left[ \Phi\big(-2h^*(x)h(x)\big) - L_\Phi(x, h_\Phi^*(x)) \right] 1_{h(x)h^*(x) \leq 0} \right]^{\frac{1}{s}} \quad (\Phi \text{ non-decreasing})$$

$$\leq 2c \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ \left[ L_\Phi(x, h(x)) - L_\Phi(x, h_\Phi^*(x)) \right] 1_{h(x)h^*(x) \leq 0} \right]^{\frac{1}{s}} \quad (\text{convexity ineq.})$$

$$\leq 2c \underset{x \sim \mathcal{D}_X}{\mathrm{E}} \left[ L_\Phi(x, h(x)) - L_\Phi(x, h_\Phi^*(x)) \right]^{\frac{1}{s}}.$$

$$\square$$

9. [Bonus point] Show that the assumption holds for the logistic loss with
   where $(\Phi(u) = \log_2(1 + e^u))$, with $s = 2$ and $c = \frac{1}{\sqrt{2}}$. Show the same
   for the exponential loss.