

Reporting assignment  
 京都大学集中講義  
 「機械学習と深層学習の数理と応用」

2018-12-20

Taiji Suzuki  
 e-mail: taiji@mist.i.u-tokyo.ac.jp

Solve the following problems.

**Due date: 10 January/2019.**

1. Consider a linear model

$$Y = \mathbf{X}\beta^* + \epsilon$$

where  $Y \in \mathbb{R}^n$ ,  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , and  $\epsilon = (\epsilon_i)_{i=1}^n \in \mathbb{R}^n$ . Suppose that  $\epsilon_i$  is an i.i.d. noise such that  $\mathbb{E}[\epsilon_i] = 0$  and  $\mathbb{E}[\epsilon_i^2] = \sigma^2$  and  $\mathbf{X}^\top \mathbf{X}$  is full rank. In this setting, evaluate the in-sample predictive error

$$\mathbb{E}_{Y|X} \left[ \frac{1}{n} \|\mathbf{X}\beta^* - \mathbf{X}\hat{\beta}_{\text{LS}}\|^2 \right].$$

of the least squared estimator  $\hat{\beta}_{\text{LS}} \in \mathbb{R}^p$ .

2. (Stein's shrinkage estimator) Let  $X = [X_1, \dots, X_d]^\top \in \mathbb{R}^d$  be distributed from multivariate normal  $N(\boldsymbol{\mu}, I)$  (mean  $\boldsymbol{\mu}$  and variance-covariance  $I$ ). Assume  $d \geq 3$ . Then, show that

$$\boldsymbol{\delta} = \left( 1 - \frac{(d-2)}{\|X\|^2} \right) X$$

satisfies

$$\mathbb{E}_{X \sim N(\boldsymbol{\mu}, I)} [\|X - \boldsymbol{\mu}\|^2] > \mathbb{E}_{X \sim N(\boldsymbol{\mu}, I)} [\|\boldsymbol{\delta}(X) - \boldsymbol{\mu}\|^2] \quad (\forall \boldsymbol{\mu} \in \mathbb{R}^d).$$

You may use the following *Stein's identity* without proof: For a function  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  ( $X \mapsto \mathbf{f}(X) = [f_1(X), \dots, f_d(X)]^\top$ ) such that  $\mathbb{E}_X[f_i(X)]$  exists and  $f_i$  is differentiable almost everywhere for all  $i$ , it holds that

$$\mathbb{E}_X [-2\langle \boldsymbol{\mu} - X, \mathbf{f}(X) \rangle] = 2\sigma^2 \mathbb{E}_X \left[ \sum_{i=1}^d \frac{\partial f_i(X)}{\partial X_i} \right].$$

3. For  $1 \leq q < \infty$ , let  $\|w\|_q := (\sum_{i=1}^d |w_i|^q)^{1/q}$  for  $w \in \mathbb{R}^d$ , and  $\mathcal{H}_q := \{f(x) = w^\top x \mid \|w\|_q \leq 1, w \in \mathbb{R}^d\}$ . Given  $x_1, \dots, x_n \in \mathbb{R}^d$ , its empirical Rademacher complexity is denoted by

$$\hat{R}_n(\mathcal{H}_q) = \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{H}_q} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \mid x_1, \dots, x_n \right].$$

Now, suppose that  $1 \leq p, q < \infty$  satisfies  $q < p$ . Then, show that

$$\hat{R}_n(\mathcal{H}_p) \leq d^{1/p^* - 1/q^*} \hat{R}_n(\mathcal{H}_q)$$

where  $p^* = p/(p-1)$  and  $q^* = q/(q-1)$ .

4. Prove the Massart's theorem: For a finite set of functions  $\mathcal{F} = \{f_1, \dots, f_M\}$  where each  $f_i$  ( $i = 1, \dots, M$ ) is a function from  $\mathbb{R}^d$  to  $\mathbb{R}$  satisfying  $\sup_x |f_i(x)| \leq R$ , it holds that

$$\hat{R}_n(\mathcal{F}) \leq R \sqrt{\frac{2 \log M}{n}}.$$

Hint: You may use the following inequalities:

- Jensen's inequality:  $\exp(s\mathbb{E}_\sigma[g(\sigma_1, \dots, \sigma_n)]) \leq \mathbb{E}_\sigma[\exp(sg(\sigma_1, \dots, \sigma_n))]$  for  $s \in \mathbb{R}$  and  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ .
- $\exp(\max_{f \in \mathcal{F}} F(f)) \leq \sum_{f \in \mathcal{F}} \exp(F(f))$  for  $F : \mathcal{F} \rightarrow \mathbb{R}$ .
- Hoeffding's inequality:  $\mathbb{E}_{\sigma_i}[\exp(\sigma_i a)] \leq \exp(a^2/2)$  for  $a \in \mathbb{R}$ .

5. Suppose that  $\|x\|_\infty < 1$  (a.s.). Derive an upper bound of the Rademacher complexity of a neural network model:

$$\mathcal{F} = \left\{ \sum_{j=1}^M \alpha_j \eta(a_j^\top x) \mid \alpha_j \in \mathbb{R}, a_j \in \mathbb{R}^d, \max_{1 \leq j \leq M} \|a_j\|_p \leq C_1, \|\alpha\|_q \leq C_2 \right\}$$

where  $1 \leq p, q \leq \infty$  and  $\eta(u) = \max\{u, 0\}$ .