

## PC2. FOUNDATIONS OF MACHINE LEARNING. (MDC\_51006\_EP - 2025-2026)

**EXERCISE 1 (LINK ESTIMATION/PREDICTION)** Assume we observe a sample  $((X_i, Y_i))_{i=1}^n$  where  $Y_i \in \{-1, 1\}$  and  $X_i \in \mathbb{R}^d$  are independent random variables following the same law than a generic pair  $(X, Y)$ . We want to predict  $Y$  from  $X$  using a predictor  $f : \mathbb{R}^d \rightarrow \{-1, 1\}$  and we measure the quality of our predictor using the 0/1 loss:

$$\ell^{0/1}(y, y') = \mathbf{1}_{y \neq y'}.$$

The optimal Bayes classifier is given by

$$f^*(X) = \text{sign}(2p_1(X) - 1)$$

where  $p_1(X) = \mathbb{P}(Y = 1|X)$ .

Assume we have an estimate of the conditional law of  $Y|X$  and denote  $\hat{p}_1(X) = \mathbb{P}(\hat{Y} = 1|X)$ . We define the plug-in classifier as

$$\hat{f} = \text{sign}(2\hat{p}_1 - 1).$$

We want to prove that

$$\begin{aligned} \mathbb{E}[\ell^{0/1}(Y, \hat{f}(X))] - \mathbb{E}[\ell^{0/1}(Y, f^*(X))] \\ \leq \mathbb{E}[\|\hat{Y}|X - Y|X\|_1] \\ \leq \left( \mathbb{E}[2\text{KL}(Y|X, \hat{Y}|X)] \right)^{1/2} \end{aligned}$$

1. Prove that for any predictor  $f$ ,

$$\mathbb{E}[\ell^{0/1}(Y, f(X))] = \mathbb{E}_X[(1 - p_1(X)) + (2p_1(X) - 1)\mathbf{1}_{f(X)=-1}]$$

2. Deduce that

$$\begin{aligned} \mathbb{E}[\ell^{0/1}(Y, \hat{f}(X))] - \mathbb{E}[\ell^{0/1}(Y, f^*(X))] \\ \leq 2\mathbb{E}_X[|p_1(X) - \hat{p}_1(X)|] = \mathbb{E}_X[\|p(X) - \hat{p}(X)\|_1] \end{aligned}$$

3. Finish the proof using  $\|P - Q\|_1 \leq \sqrt{2\text{KL}(P, Q)}$

**EXERCISE 2 (LINEAR REGRESSION AND DESIGN)** Assume we observe some points  $((X_i, Y_i))_{i=1}^n$  with  $X_i \in \mathbb{R}^p$  and  $Y_i \in \mathbb{R}$  are assumed to be independent copies of a generic  $(X, Y)$ , regressing those points means estimating  $\mathbb{E}[Y|X]$  at least on the  $X_i$ . Indeed, there are two classical problems that can be studied

- **Fixed Design** We assume that the  $X_i$  are fixed and want to estimate  $f(X) = \mathbb{E}[Y|X]$  by a predictor  $\hat{f}$  only at the observed points  $X_i$ . We measure the quality of  $\hat{f}$  as the average error on a replication set  $((X_i, Y'_i))_{i=1}^n$ .

$$\text{Err}_F = \mathbb{E}_{(Y_i)_{i=1}^n} \left[ \mathbb{E}_{(Y'_i)_{i=1}^n} \left[ \frac{1}{n} \sum_{i=1}^n \ell(Y'_i, \hat{f}(X_i)) \right] \right]$$

where the  $X_i$  are the same and  $Y'_i|X_i$  has the same law as  $Y_i|X_i$ .

- **Random Design** We assume that the  $X_i$  are i.i.d. and we want to estimate  $f(X) = \mathbb{E}[Y|X]$  by a predictor  $\hat{f}$  everywhere. We measure the quality as the average error on new observation  $(X', Y')$  having the same law as  $(X_i, Y_i)$ :

$$\text{Err}_R = \mathbb{E}_{((X_i, Y_i))_{i=1}^n} \left[ \mathbb{E}_{(X', Y')} [\ell(Y', \hat{f}(X'))] \right].$$

Linear least square regression is classically studied in the fixed design setting, while the random design setting corresponds to the statistical learning setting. In this exercise, we are going to look at the performance of the linear least square on those two settings.

Following Rosset and Tibshirani, we introduce an intermediate setting:

- **Repeated Fixed Design** We assume that the  $X_i$  are random, possibly of different law and want to estimate  $f(X) = \mathbb{E}[Y|X]$  by a predictor  $\hat{f}$  only at the observed points  $X_i$ . We measure the quality by the Fixed Design error averaged over all possible designs:

$$\text{Err}_{RF} = \mathbb{E}_{((X_i, Y_i))_{i=1}^n} \left[ \mathbb{E}_{(Y'_i)_{i=1}^n} \left[ \frac{1}{n} \sum_{i=1}^n \ell(Y'_i, \hat{f}(X_i)) \right] \right].$$

as well as the average empirical estimate of the quality:

$$\text{Err}_{Emp} = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \hat{f}(X_i)).$$

We assume that all the samples are independent conditionally to  $X$ , that the  $X_i$  are fixed in the Fixed design setting, independent in the Repeated Fixed design setting and i.i.d. in the Random design one. Except, for the warm-up, we will assume a homoscedastic setting in which  $Y = f(X) + \epsilon$  with  $\epsilon$  a centered random noise of variance  $\sigma^2$  and that the loss is the square one.

#### 1. Warm-up

- (a) Verify that  $\mathbb{E}_{(X_i, Y_i)_{i=1}^n} [\text{Err}_F] = \text{Err}_{RF}$  and, if the  $(X_i, Y_i)$  are i.i.d.,

$$\text{Err}_R = \mathbb{E}_{((X_i, Y_i))_{i=1}^n} \left[ \mathbb{E}_{((X'_i, Y'_i))_{i=1}^n} \left[ \frac{1}{n} \sum_{i=1}^n \ell(Y'_i, \hat{f}(X'_i)) \right] \right]$$

where  $(X'_i, Y'_i)$  are independent copies of  $(X_i, Y_i)$ .

- (b) Verify that if  $\hat{f} = \text{argmin}_g \sum_{i=1}^n \ell(Y_i, g(X_i))$  then  $\mathbb{E}_{(Y_i)_{i=1}^n} [\text{Err}_{Emp}] \leq \text{Err}_F$  and  $\mathbb{E}_{((X_i, Y_i))_{i=1}^n} [\text{Err}_{emp}] \leq \text{Err}_{RF}$ .

#### 2. We will now assume that $\ell$ is the $\ell^2$ loss: $\ell(Y, f(X)) = |Y - f(X)|^2$ .

- (a) Prove that in the fixed design setting:

$$\text{extErr}_F = \sigma^2 + \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}_{(Y_i)_{i=1}^n} [\hat{f}(X_i)])^2 + \frac{1}{n} \sum_{i=1}^n \text{Var}_{(Y_i)_{i=1}^n} [\hat{f}(X_i)]$$

- (b) Show that if we assume that the  $X_i$  are i.i.d. then for the repeated fixed design setting:

$$\text{Err}_{RF} = \sigma^2 + \mathbb{E}_{(X_i)_{i=1}^n} \left[ (f(X_1) - \mathbb{E}_{(Y_i)_{i=1}^n} [\hat{f}(X_1) | (X_i)_{i=1}^n])^2 \right] + \mathbb{E}_{(X_i)_{i=1}^n} \left[ \text{Var}_{(Y_i)_{i=1}^n} [\hat{f}(X_1) | (X_i)_{i=1}^n] \right]$$

- (c) Finally prove that for the random design setting

$$\text{Err}_R = \sigma^2 + \mathbb{E}_{(X_i)_{i=1}^n, X'} \left[ (f(X') - \mathbb{E}_{(Y_i)_{i=1}^n} [\hat{f}(X') | (X_i)_{i=1}^n])^2 \right] + \mathbb{E}_{(X_i)_{i=1}^n, X'} \left[ \text{Var}_{(Y_i)_{i=1}^n} [\hat{f}(X') | (X_i)_{i=1}^n] \right]$$

#### 3. We define the optimism $\text{Opt}_F$ (respectively $\text{Opt}_{RF}$ ) as the difference between $\text{Err}_F$ (respectively $\text{Err}_{RF}$ ) and the corresponding expectation of $\text{Err}_{emp}$ .

- (a) Justify the name *optimism* when the estimator is the minimizer of the empirical risk.  
(b) Verify that in the fixed design setting:

$$\text{Opt}_F = \text{Err}_F - \mathbb{E}_{(Y_i)_{i=1}^n} [\text{Err}_{emp}] = \mathbb{E}_{(Y_i)_{i=1}^n, (Y'_i)_{i=1}^n} \left[ \frac{1}{n} \sum_{i=1}^n (Y'_i - \hat{f}(X_i))^2 - \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}(X_i))^2 \right]$$

(c) Prove that

$$\text{Opt}_F = \frac{2}{n} \sum_{i=1}^n \text{Cov} \left[ Y_i, \hat{f}(X_i) \right]$$

(d) Deduce that in the repeated fixed design setting

$$\text{Opt}_{RF} = \frac{2}{n} \sum_{i=1}^n \mathbb{E}_{(X_i)_{i=1}^n} \left[ \text{Cov} \left[ Y_i, \hat{f}(X_i) \right] \middle| (X_i)_{i=1}^n \right]$$

4. We consider now the random design and define  $\text{Opt}_R$  in the same way.

(a) Verify that

$$\begin{aligned} \text{Err}_R = \text{Err}_F &+ \underbrace{\mathbb{E}_{(X_i)_{i=1}^n, X'} \left[ (f(X') - \mathbb{E}_{(Y_i)} [\hat{f}(X') | (X_i)_{i=1}^n])^2 \right] - \mathbb{E}_{(X_i)} \left[ (f(X_1) - \mathbb{E}_{(Y_i)} [\hat{f}(X_1) | (X_i)])^2 \right]}_{\Delta_B} \\ &+ \underbrace{\mathbb{E}_{(X_i)_{i=1}^n, X'} \left[ \text{Var}_{(Y_i)} [\hat{f}(X') | (X_i)_{i=1}^n] \right] - \mathbb{E}_{(X_i)_{i=1}^n} \left[ \text{Var}_{(Y_i)} [\hat{f}(X_1) | (X_i)_{i=1}^n] \right]}_{\Delta_V} \end{aligned}$$

(b) Deduce that

$$\text{Opt}_{RF} = \text{Opt}_F + \Delta_B + \Delta_V$$

(c) Why is this reasonable to expect that both  $\Delta_B$  and  $\Delta_V$  are non negative?

5. We consider now the linear model  $f_\beta(X') = \langle X', \beta \rangle$  and pick  $\beta$  by minimizing the empirical loss.

(a) Prove that, assume the design matrix  $\mathbb{X}$  is such that  $\mathbb{X}^\top \mathbb{X}$  is an invertible  $d \times d$  matrix, the estimate  $\hat{\beta}$  is given by

$$\hat{\beta} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y}$$

so that  $\hat{f}(X_i) = \left( \text{Proj}_{\text{span}(X_i)} \mathbf{Y} \right)_i$ .

(b) Deduce that  $\text{Opt}_R = \text{Opt}_{RF} = 2\sigma^2 \frac{d}{n}$

(c) Prove that  $\Delta_B \geq 0$ .

(d) Prove that  $\Delta_V \geq 0$  (one can use  $\mathbb{E}[(\mathbb{X}^\top \mathbb{X})^{-1}] - (\mathbb{E}[\mathbb{X}^\top \mathbb{X}])^{-1}$  is s.d.p.)