

## "Machine Learning and Computational Statistics"

### 5<sup>th</sup> Homework

#### Exercise 1:

Consider the **regression problem**  $y=g(x)+\eta$

It is known that  $E[y|x]$  is the **minimum MSE estimate** of  $y$  given  $x$ . Consider the estimator  $f(x;D)$ .

- (a) Under what conditions (theoretically) the quantity  $E_D[(f(x;D) - E[y|x])^2]$  becomes zero?
- (b) Why this cannot be achieved in practice?

#### Exercise 2:

Consider a regression task modelled by the relation  $y = g(x) + \eta$ . Let  $f_{\theta}$  be an estimator of  $g$ , parameterized by the vector  $\theta$ . Let  $Tr = \{(x_n, y_n), n = 1, \dots, N_1\}$  be the training set (the set that will be used for the estimation of  $\theta$ ) and  $Te = \{(x_n, y_n), n = 1, \dots, N_2\}$  be the test (which will be used for testing the performance of  $f_{\hat{\theta}}$  (where  $\hat{\theta}$  is the estimate of  $\theta$  based on  $Tr$ )).

- (a) What indicates a **large** error value on the training set  $Tr$ ?
- (b) What may indicate a **large** error value on the test set  $Te$ ?
- (c) What indicates a **small** error value on the training set  $Tr$ ?
- (d) What may indicate a **small** error value on the test set  $Te$ ?

#### Exercise 3:

Consider a regression task  $y = g(x) + \eta$ , where  $y$  and  $x$  are modeled by the random variables  $y$  and  $x$ . The joint pdf of  $y$  and  $x$  is:

$$p(x, y) = \frac{3}{2}, \text{ for } x \in (0, 1), y \in (x^2, 1).$$

Determine the optimum MSE estimate  $E[y|x]$ , for a given  $x$ , by performing the following steps:

- (a) Verify that  $p(x, y)$  is a pdf (prove that it integrates to 1).
- (b) Compute the marginal pdf of  $x$ ,  $p_x(x)$ .
- (c) Compute the conditional pdf of  $y$ , given  $x$ .
- (d) Compute and plot  $E[y|x]$ .

Hint: It is  $\int_a^b x^n dx = \left[ \frac{1}{n+1} x^{n+1} \right]_a^b = \frac{1}{n+1} b^{n+1} - \frac{1}{n+1} a^{n+1}$

**Exercise 4 (python code + text):**

Consider the **regression problem** (1-dep., 1-indep. variables)

$$y = g(x) + \eta$$

where **y** and **x** are **jointly distributed** according to the **normal distribution**  $p(y, x) = N(\mu, \Sigma)$

with  $\mu = \begin{bmatrix} \mu_y \\ \mu_x \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$   $\Sigma = \begin{bmatrix} \sigma_y^2 & \sigma_{yx} \\ \sigma_{yx} & \sigma_x^2 \end{bmatrix} = \begin{bmatrix} 3 & 2 \\ 2 & 4 \end{bmatrix}$

- Determine  $E[y|x]$  and plot the corresponding curve (recall the relevant theory concerning the normal distribution case).
- Generate **100** data sets  $D_i$ ,  $i=1, \dots, 100$ , each one consisting of  **$N=50$**  randomly selected pairs  $(y_n, x_n)$ ,  $n=1, \dots, N$ , from  $p(y, x)$ .
- Adopt a linear estimator  $f(x; D)$  and determine its instances  $f(x; D_1), \dots, f(x; D_{100})$ , utilizing the LS criterion.
- Plot in a single figure **(i)** the lines corresponding to the above 100 estimates (**blue color**) and **(ii)** the line corresponding to the optimal MSE estimate (**green color**).
- Repeat steps (b)-(d) where now each data set consists of  **$N=5000$**  points.
- Discuss the results (in your discussion, take into account the decomposition of the MSE to a variance and a bias term).

**Exercise 5 (python code + text):**

Consider the set up of exercise 3 and recall the  $E[y|x]$  determined there.

- Generate a single data set  **$D$**  of 100 pairs  $(y_n, x_n)$ ,  $n=1, \dots, 100$  from  $p(y, x)$ .
- Determine the linear estimate  $f(x; D)$  that minimizes the MSE criterion, based on  **$D$** .
- Generate randomly a set  **$D'$**  of additional 50 points  $(y'_n, x'_n)$ ,  $n=1, \dots, 50$ . For each  $x'_n$  determine the estimate  $y_n' = f(x_n; D')$  (50 numbers (estimates) should be finally computed).
- Again, for the 50  $x'_n$ 's determine the associated estimates  $\hat{y} = E[y|x]$ .
- Based on the previous derived estimates for the 50 points from both  $f(x_n; D)$  and  $E[y|x]$ , propose and use a (practical) way for quantifying the performance of the two estimators  $f(x_n; D')$  and  $E[y|x]$ .

**Exercise 6 (python code + text):** Consider the setup of exercise 3. Generate a set  **$D$**  of  $N = 100$  data pairs  **$z_n = (y_n, x_n)$** .

- For each  $x_n$  compute the optimal MSE estimate (use the results of exercise 3).

- (b) Compute  $\boldsymbol{\mu} = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N x_n \\ \frac{1}{N} \sum_{n=1}^N y_n \end{bmatrix}$  and  $\Sigma = \frac{1}{N} \sum_{n=1}^N (\boldsymbol{\mu} - \mathbf{z}_n)(\boldsymbol{\mu} - \mathbf{z}_n)^T$ .
- (c) Pretend that you do not know the true distribution that generates the data and you (erroneously) assume that the joint pdf of  $x$  and  $y$  is a normal one with mean and covariance matrix those computed in (b). Derive the optimum MSE estimate for this case and compute the MSE estimate for each one of the 100  $x_n$ 's.
- (d) Discuss the results obtained from (a) and (c).

**NOTE:** Please give **brief explanations** in all **exercises**.