

STATISTICAL MACHINE LEARNING FOR DATA SCIENCE

ASSIGNMENT 1

Due Date: January 10, 2021 - By 11:59pm

All your answers must be written on a separate sheet, properly typeset and submitted in the form of a report, in pdf format. No MS Word report will be accepted. Make sure your report in pdf format is uploaded to the designated folder before the deadline.

EXERCISE 1

Let $\mathcal{D}_n = \{(x_i, y_i) \stackrel{iid}{\sim} p_{xy}(x, y), x_i \in \mathbb{R}, y_i \in \mathbb{R}, i = 1, \dots, n\}$. Consider using the data \mathcal{D}_n , to build mappings $f: \mathcal{X} \rightarrow \mathcal{Y}$, such that $f \in \mathcal{H}$, where

$$\mathcal{H} := \left\{ x \mapsto f(x) = \theta x, \theta \in \mathbb{R}^* \right\} \quad (1)$$

Further suppose that $\forall i \in [n]$, we have

$$p(y_i | x_i, \theta, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - f(x_i))^2}, \quad (2)$$

where $f \in \mathcal{H}$. Finally, let $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top$, and $\mathbf{y} = (y_1, y_2, \dots, y_n)^\top$ and define

$$\text{SSE}_n(a) = \sum_{i=1}^n (y_i - ax_i)^2 = (\mathbf{y} - a\mathbf{x})^\top (\mathbf{y} - a\mathbf{x}) = \|\mathbf{y} - a\mathbf{x}\|_2^2. \quad (3)$$

1. Specify the input space in this problem, and clearly indicate its dimensionality.
2. Specify the output space in this problem, and clearly indicate its dimensionality.
3. Specify the dimensionality of \mathbf{x} .
4. Specify the dimensionality of \mathbf{y} .
5. Determine in this case the assumed conditional distribution of Y_i given x_i and deduce the distribution \mathbf{y} given \mathbf{x} .
6. Rewrite the model defined by Equation (2) in its additive form featuring the deterministic component (signal) and the stochastic component (noise or error term). Be sure to reflect the fact that $f \in \mathcal{H}$ as defined in Equation (1), but also the clear probability distribution used.
7. Indicate which type statistical machine learning task is being solved here. Justify your answer.

8. Use the appropriate tools to find $\frac{\partial \text{SSE}_n(a)}{\partial a}$

9. Show that

$$\hat{\theta} = \underset{\theta \in \mathbb{R}^*}{\operatorname{argmin}} \{ \text{SSE}_n(\theta) \} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}}.$$

Now, let $\hat{Y}_i = \hat{\theta} x_i$, for $i \in [n]$, and define $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top$.

10. Find $\text{mean}[\hat{\theta}] = \mathbb{E}[\hat{\theta}]$.

11. Find $\text{variance}[\hat{\theta}] = \mathbb{V}[\hat{\theta}]$.

12. Find $\text{mean}(Y_i|x_i) = \mathbb{E}[Y_i|x_i]$ for $i \in [n]$ and deduce $\text{mean}(\mathbf{y}|\mathbf{x}) = \mathbb{E}(\mathbf{y}|\mathbf{x})$.

13. Write down \mathbf{y} as a function of \mathbf{x} , θ and all other necessary parts of the assumed model in keeping with Equation (2).

14. Write down $\hat{\mathbf{y}}$ as function of \mathbf{x} and \mathbf{y} .

15. Find $\text{variance}[Y_i|x_i] = \mathbb{V}[Y_i|x_i]$ for $i \in [n]$ and deduce $\text{variance}(\mathbf{y}|\mathbf{x})$.

16. Find $\text{mean}(\hat{Y}_i|x_i) = \mathbb{E}[\hat{Y}_i|x_i]$ for $i \in [n]$ and deduce $\text{mean}(\hat{\mathbf{y}}|\mathbf{x}) = \mathbb{E}(\hat{\mathbf{y}}|\mathbf{x})$.

17. Find $\text{variance}(\hat{Y}_i|x_i) = \mathbb{V}[\hat{Y}_i|x_i]$ for $i \in [n]$ and deduce $\text{variance}(\hat{\mathbf{y}}|\mathbf{x}) = \mathbb{V}(\hat{\mathbf{y}}|\mathbf{x})$.

18. Based on all the above, determine the distribution $\hat{\theta}$.

19. Based on all the above, determine the distribution \hat{Y}_i given x_i and deduce the distribution of $\hat{\mathbf{y}}$ given \mathbf{x} .

20. Find the estimator $\hat{\sigma}^2$ of σ^2 .

EXERCISE 2

```
prostate <- read.csv('prostate-cancer-1.csv') # DNA MicroArray Gene Expression
```

1. Find by all means possible the history and description of this data and comment on it.
2. Plot the distribution of the response for this dataset and comment.
3. Comment on the shape of this dataset in terms of the sample size and the dimensionality of the input space
4. Comment succinctly from the statistical perspective on the type of data in the input space. It is absolutely crucial here to provide as many details as possible including distributional aspects via boxplots and others on randomly selected subsets of variables.

EXERCISE 3

Consider a linear regression analysis in \mathbb{R}^2 with $Y_i = \theta_1 x_{i1} + \theta_2 x_{i2} + \sigma Z_i$ where $Z_i \sim \mathcal{N}(0, 1)$ for $i = 1, \dots, n$, and $\sigma \in \mathbb{R}_+^*$. You are given the data matrix \mathbf{X} of Equation (4).

$$\mathbf{X} = \begin{bmatrix} 1 & -2 \\ -2 & 1 \\ 4 & 1 \end{bmatrix} \quad (4)$$

You are also given $\mathbf{Y} = (-5, 4, -3)^\top$ and the vector of observed response values.

1. Compute the important $\mathbf{X}^\top \mathbf{X}$
2. Comment on the shape of $\mathbf{X}^\top \mathbf{X}$
3. Find $(\mathbf{X}^\top \mathbf{X})^{-1}$ in the most straightforward way
4. Compute $\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
5. Compute the vector $\hat{\mathbf{Y}}$ of estimated responses
6. Compute the vector $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}}$ of residual values
7. Find the value of $\text{SSE}(\hat{\theta})$
8. Find the estimate $\hat{\sigma}^2 = \frac{\text{SSE}(\hat{\theta})}{n-2}$
9. Find and write down $\text{Variance}(\hat{\theta}) = \hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})^{-1}$
10. Write the matrix $\mathbf{X}^\top \mathbf{X}$ as a function of the identity matrix when the data matrix is

$$\mathbf{X} = \begin{bmatrix} 1 & -2 \\ 2 & 1 \end{bmatrix}. \quad (5)$$

EXERCISE 4

Consider the `gifted children` data set. For this dataset, the response variable is `score`, which contains the score of the child on a battery of tests. The explanatory variables in this case are self explanatory by their names. The goal here is to attempt to build the best linear model that captures the relationship between the response `score` and the other variables

1. Generate an upper triangular pairwise scatterplot for this data, and comment based on the scatterplots regarding which of the explanatory variables are more strongly related to the response. Can you tell from the plot the strongest of all the predictor variables?
2. Generate the correlation matrix for this data (please do not include the p-values in the matrix for now). Which variable does the correlation matrix appear to indicate as the strongest?
3. Plot a histogram of the response variable and also perform a test of normality on it
4. Perform a simple linear regression (SLR) model fitting featuring the response and the variable you singled out as the most important.
5. Generate the 4 residual analysis plots and comment on the suitability of your built model. Are there observations that might have badly influenced the estimation of your parameters?
6. Let's assume for a little while that you are to use the above SLR model. Then give an interpretation of your estimated slope in layman's term.
7. Generate both confidence bands and the prediction bands for this model and provide intelligent comments on what you get.
8. Build a multiple linear regression (MLR) model for this data comprising of all the provided explanatory variables