

# Statistical Machine Learning (SML)

Winter 2021

## Assignment 3

Maximum Marks - 100

Due Date: 23.59 hrs., 31<sup>st</sup> March, 21

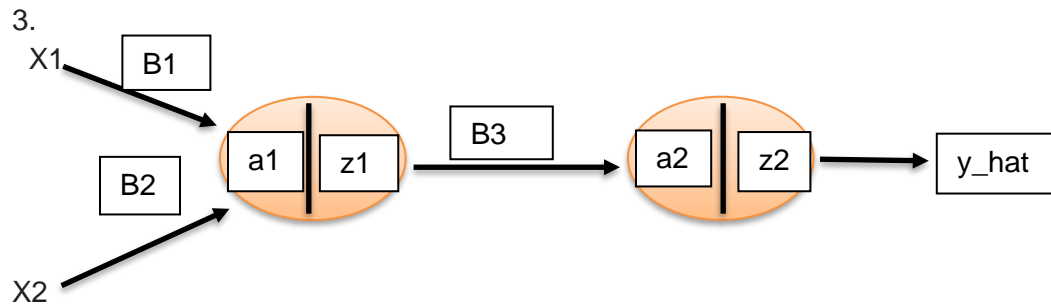
---

### Instructions:

1. You are free to use either python or MATLAB for this assignment.
  2. You can use inbuilt libraries for Math, plotting, and handling the data (eg. NumPy, Pandas, Matplotlib).
  3. Usage instructions for other libraries can be found in the question.
  4. Only (\*.py) and (\*.m) files should be submitted for code.
  5. Create a (\*.pdf) report explaining your assumptions, approach, results, and any further detail asked in the question.
  6. You should be able to replicate your results if required.
- 

1. Use [MNIST](#) data for this question, and perform the following tasks.
  - a. [5] Visualize 5 samples from each class in the form of images.
  - b. [10] Implement FDA for multiple classes from scratch, and find the coefficient vector  $W$ .  
Note: computation of  $W$  will use training samples only.
  - c. [3] Project the training data ( $X$ ) using  $W$ , and call the projection  $Y$ .
  - d. [10] Use the projected data  $Y$  to classify the testing samples using QDA (Quadratic Discriminant Analysis).  
Note: You can reuse the implementation of QDA from assignment 2.
  - e. [2] Report the accuracy (the ratio of correctly classified samples to the total number of samples tested).
2. In this problem, you will explore **Gaussian process regression (GPR)**.
  - a. [5] Generate 5 random samples from a uniform distribution in  $[0, 10]$ , call it  $X_{\text{train}}$ .  
Generate  $Y_{\text{train}}$  using  $Y_{\text{train}} = X_{\text{train}} \cdot \exp(X_{\text{train}})$ .
  - b. [10] Compute the matrices  $K$ ,  $K^*$ ,  $K^{**}$  and use cross-validation for obtaining  $\sigma$  and  $l \rightarrow$  the parameters of RBF kernel. Consider a range of values for  $\sigma$  and  $l$ .  
Perform cross-validation as follows:
    - b.1 For a particular combination of  $\sigma$  and  $l$ , take 4 samples to train the GPR and call the remaining sample test point, compute prediction for the test point, and run this 5 times, each time take a different set of samples as training and testing points. Find the error for each run and compute their mean.
    - b.2 Repeat b.1 for each combination of  $\sigma$  and  $l$  and choose the values which result in a minimum mean error.

- c. [5] Generate 50 random samples from a uniform distribution in  $[0,10]$  and call them  $X_{\text{test}}$ . Generate  $Y_{\text{test}}$  using  $Y_{\text{test}} = Y_{\text{test}} * \exp(X_{\text{test}})$ .
- d. [5] Compute the prediction  $Y_{\text{pred}}$  for test samples.
- e. [5] Plot the actual values( $Y_{\text{test}}$ ) and predicted values( $Y_{\text{pred}}$ ) of test samples.



$$\begin{aligned}
 a1 &= B1.X1 + B2.X2 + B0 & z1 &= \sigma(a1) \\
 a2 &= B3.Z1 + B^* & z2 &= \sigma(a2) \\
 y_{\text{hat}} &= z2 & E &= (y - y_{\text{hat}})^2
 \end{aligned}$$

Note: 1)  $B0$  and  $B^*$  are bias.  
 2)  $\sigma$  denotes sigmoid functions.

Refer to the given network to perform the following tasks:

- a. [10] Determine an expression for all the weights using backpropagation. (Pen-Paper problem)
- b. [5] Generate  $X$ : Sample 100 points from  $N(\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix})$

Generate  $Y$ : Sample 100 random points from Gaussian distribution (1 dimension), which acts as a label for  $X$ .

Use 50 of those samples for training and remaining for testing.

- c. [10] Implement the expression obtained in part a.
- d. [10] Cycle through each point and make an update for the complete training set. Call this as epoch, and do 5 such epoch.
- e. [5] Compute MSE for the test set.