

Statistical Learning-Classification

STAT 441 / 841, CM 763

Assignment 3

Department of Statistics and Actuarial Science
University of Waterloo

Due: Friday Nov 9 at 1pm

Policy on Lateness: No assignment are accepted after the due date.

1. a) Write a program to fit an RBF network. In implementing RBF, you need to cluster the data and find the center and spread of each cluster. You don't need to implement a clustering algorithm yourself. ¹.
- b) Use the Ionosphere dataset (Ion.mat) . Use the standard cross validation, standard Leave one out cross validation, and Leave one out cross validation as expressed in (1) (The method explained in Question 6 shows how LOO can be performed without iteration.) and find the optimum number of basis function for each model. Compute the test error in each case and complete the following table.

In this table

CV is standard cross validation

LOO is standard leave one out cross validation

CLOO is Leave one out cross validation as expressed in (1).

Method	Target function	
	<i>TrainingError</i>	<i>TestError</i>
CV		
LOO		
CLOO		

Note: Attach your code to your assignment as an appendix and submit the code to the assignment drop box in Learn as well.

2. Support Vector Machine

¹You can use any clustering algorithm and any clustering routine in any programming language based on your preference

a) Write a function $[b, b_0] = \text{HardMarg}(X, y)$ which takes a $d \times n$ matrix X and $n \times 1$ vector of target labels y and returns: a $d \times 1$ vector of weights b and a scalar offset b_0 , corresponding to the maximum margin linear discriminant classifier.

b) Write a function $[b, b_0] = \text{SoftMarg}(X, y, \gamma)$ which takes an additional scalar argument γ and returns b and b_0 corresponding to the maximum soft margin linear discriminant classifier.

c) Write a function $[yhat] = \text{classify}(Xtest, b, b_0)$ which takes a $d \times m$ matrix $Xtest$, a $d \times 1$ vector of weights b , and a scalar b_0 , and returns a $m \times 1$ vector of classifications $yhat$ on the test patterns.

d) For each of the datasets linear, noisylinear, and quadratic on Piazza solve for each kind of discriminant function:

$$[bh, b_0h] = \text{HardMarg}(X, y), [bs, b_0s] = \text{SoftMarg}(X, y, 0.5)$$

produce a 2D plot of the training data and the two hypotheses corresponding to bh, b_0h and bs, b_0s and report the mean misclassification error (i.e., the sum of misclassification errors divided by the number of data points) that each of the two hypotheses obtained on the training data and on the test data.

Hand in a plot and two tables for each dataset.

Note 1: Submit your code to Learn drop box.

Note 2: Your function must be able to handle arbitrary d , n , γ , and m .

- Let \hat{f} be an estimator of the quantity f , show that its mean-squared error can be decomposed as follows:

$$\begin{aligned} E(\hat{f} - f)^2 &= E[\hat{f} - E(\hat{f})]^2 + [E(\hat{f}) - f]^2 \\ &= \text{Var}(\hat{f}) + \text{Bias}^2(\hat{f}) \end{aligned}$$

- Consider a neural network that consists of a single neuron with d inputs. The neuron has d weights, $w \in \mathbb{R}^d$. The output of the neuron for an input pattern $x \in \mathbb{R}^d$ is given by $\hat{y} = \Phi(x \cdot w)$, where Φ is an activation function.

For any differentiable activation function Φ , there exists a *matching loss*, denoted by $\text{err}_\Phi(y, \hat{y})$, such that when using Φ and its matching loss $\text{err}_\Phi(y, \hat{y})$, the error function of a single neuron is convex and thus has only one minimum. The *matching loss* can be computed as:

$$\text{err}_\Phi(y, \hat{y}) = \int_{\Phi^{-1}(y)}^{\Phi^{-1}(\hat{y})} (\Phi(z) - y) dz$$

- Find the matching loss for the activation function $\Phi_1(z) = z$.

- b) Find the matching loss for the activation function $\Phi_2(z) = \frac{1}{1+e^{-z}}$.
- c) Suppose you want to use this simple network for a classification task. Which of these two loss functions (Φ_1 or Φ_2) is more appropriate? Briefly explain why.

Only for Grad Students

5. Given a set of data points $\{\mathbf{x}_i\}$, we can define the *convex hull* to be the set of all points \mathbf{x} given by

$$\mathbf{x} = \sum_i \alpha_i \mathbf{x}_i$$

where $\alpha_i \geq 0$ and $\sum_i \alpha_i = 1$. Consider a second set of points $\{\mathbf{y}_i\}$ together with their corresponding convex hull. By definition, the two sets of points will be linearly separable if there exist a vector $\hat{\mathbf{w}}$ and a scalar w_0 such that $\hat{\mathbf{w}}^T \mathbf{x}_i + w_0 > 0$ for all \mathbf{x}_i , and $\hat{\mathbf{w}}^T \mathbf{y}_i + w_0 < 0$ for all \mathbf{y}_i .

Show that if their convex hulls intersect, the two sets of points cannot be linearly separable.

6. **Leave-one-out cross validation.** Consider the model $y_i = f(\mathbf{x}_i) + \epsilon_i$. When $f(\mathbf{x}_i) = \beta_0 + \beta^T \mathbf{x}_i$, this is known as an ordinary least square (OLS). Let H be the hat matrix associated with OLS. Show that

$$y_i - \hat{f}^{(-i)}(\mathbf{x}_i) = \frac{y_i - \hat{f}(\mathbf{x}_i)}{1 - H_{ii}} \quad (1)$$

where H_{ii} denote the i -th diagonal element of H ; and $\hat{f}^{(-i)}(\mathbf{x}_i)$ denotes estimating $\hat{f}(\mathbf{x}_i)$ using an \hat{f} that is obtained without using the i -th observation. Thus show that the leave-one-out cross validation can be computed without iteration.