

ECE523: Engineering Applications of Machine  
Learning and Data Analytics

Due 02/19/2020 @ 11:30PM (D2L)

Name: Zhaoyang Bu

Signature: Zhaoyang Bu

Date: Feb/21/2020

**Instructions:** Partial credit is given for answers that are partially correct. No credit is given for answers that are wrong or illegible. All work must be supported and code must be submitted for credit.

Theory: \_\_\_\_\_

Practice: \_\_\_\_\_

Total: \_\_\_\_\_

## Part A: Theory (15pts)

### (5pts) Linear Regression and Regularization

In class we derived and discussed linear regression in detail. Find the result of minimize the loss of sum of the squared errors; however, add in a penalty for an  $L_2$  penalty on the weights. More formally,

$$\arg \min_{\mathbf{w}} \left\{ \sum_i (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2 \right\} \quad (1)$$

How does this change the solution to the original linear regression solution? What is the impact of adding in this penalty?

### (5pts) Density Estimation

In  $k$ -nearest neighbors (KNN), the classification is achieved by majority vote in the vicinity of data. Suppose there are two classes of data each of  $n/2$  points overlapped to some extent in a 2-dimensional space. Describe what happens to the training error (using all available data) when the neighbor size  $k$  varies from  $n$  to 1.

### (5pts) Feature Selection & Preprocessing

A friend asks you for some help with a feature selection project. Your friend goes out and collects data,  $\mathcal{D}$ , for their project. Using  $\mathcal{D}$ , your friend tries many subsets  $\mathcal{F} \subset \mathcal{X}$  by adapting  $\mathcal{F}$  based on the error of a classifier. They return  $\mathcal{F}$  that corresponds to the smallest classification error. After they have the feature set  $\mathcal{F}$  they perform randomized trials to estimate the accuracy of the model using their feature selector.

This is the procedure they carry out to validate the impact of the feature selection routine. This procedure is repeated

- Make a new data set  $\mathcal{D}'$  with  $\mathcal{F}$  features using the feature selection routine.
- Repeat 50 times
  - Split  $\mathcal{D}'$  into randomized training & testing sets (80/20% splits)
  - Train a classifier and record its error
- Report the error averaged over 50 trials

Critique and respond to how your friend performed their analysis.

## Part B: Practice (25pts)

You are free to use functions already implemented in Python with the exception of problem 1. I recommend using Python's Scikit-learn (<http://scikit-learn.org/stable/>) as it implements most of the methods we will be discussing in this course... as well as problems in this homework!

## (10pts) Logistic Regression on Synthetic and Real-World Data

Write your own implementation of logistic regression and implement your model on either real-world (see Github data sets), or synthetic data. If you simply use Scikit-learn's implementation of the logistic regression classifier, then you'll receive zero points. A full 10/10 will be awarded to those that implement logistic regression using the optimization of cross-entropy using stochastic gradient descent.

## (5pts) Dimensionality Reduction

Choose ten data sets of your choice from the ECE523. Implement a comparison between either two classifiers of your choice, or a classifier with and without using a preprocessing step (i.e., feature selection, PCA, etc.), and report the accuracies of the two models in a table for (e.g., a  $10 \times 2$  table of classifier accuracies). Use a hypothesis testing procedure from homework #1 to determine if there is statistical significance (i.e., do both approaches perform equally well).

## (10pts) Density Estimation in Practice

The ECE523 Github page has code for generating data from a checkerboard data set. Generate checkerboard data from two classes and use any density estimate technique we discussed to classify new data using

$$\hat{p}_{Y|X}(y|x) = \frac{\hat{p}_{X|Y}(x|y)\hat{p}_Y(y)}{\hat{p}_X(x)}$$

where  $\hat{p}_{Y|X}(y|x)$  is your estimate of the posterior given you estimates of  $\hat{p}_{X|Y}(x|y)$  using a density estimator and  $\hat{p}_Y(y)$  using a maximum likelihood estimator. You should plot  $\hat{p}_{X|Y}(x|y)$  using a pseudo color plot (see <https://goo.gl/2SDJPL>). Note that you must model  $\hat{p}_X(x)$ ,  $\hat{p}_Y(y)$ , and  $\hat{p}_{X|Y}(x|y)$ . Note that  $\hat{p}_X(x)$  can be calculated using the Law of Total Probability.

Part A.

1) Original solution:  $W = (X^T X)^{-1} X^T Y$

$$L_2 \text{ Solution: } L(w) = \frac{1}{2} \sum_{i=1}^n (W^T x_i - y_i)^2 + \frac{\lambda}{2} \|W\|_2^2$$

$$= \frac{1}{2} (WX - Y)^T (WX - Y) + \frac{\lambda}{2} W^T W$$

$$\frac{dL}{dW} = 0 \Rightarrow X^T (WX - Y) + \lambda W = 0$$
$$X^T W X - X^T Y + \lambda W = 0$$

$$(X^T X + \lambda) W = X^T Y$$

$$W = \frac{X^T Y}{X^T X + \lambda}$$

$\Rightarrow L_2$  solution gives smaller values for  $w$ .

2). When  $k$  equals to  $n$ , we will use all training samples to estimate the classification of testing sample. As a result, the training error will become very large.

When  $k$  equals to 1, if the testing sample is locating outside the overlapping area, then the training error is zero, because your testing sample is surrounding by training samples that from one classification. However, if our testing sample is locating inside the overlapping area, the training error may still exist.

3) This method could help him to get the optimal features. However, this method may also very expensive. If our data set is very large, then he has to repeat 50 times of the procedure, which is a lot of work for him. Besides, he may get good features after 10 or 20 times, there is no need for him to do 50 times.

HW2

ECE523

Zhaoyang Bu