

ECE523: Engineering Applications of Machine  
Learning and Data Analytics  
Due 02/09/2017 @ 11:59PM (D2L)

Name: David Allen

Signature: David A

Date: 2/10/18

**Instructions:** There are seven problems. 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

- ① In class we derived and discussed linear regression in detail. Find the results of the minimization of the loss of the sum of the squared errors. However add in the penalty for an L2 penalty on the weights. More formally:

$$\arg \min_w \left\{ \sum_i (\omega^T x_i - y_i)^2 + \lambda \|\omega\|_2^2 \right\}$$

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

- (a) Part 1: Find results of the minimization function "loss of the sum of the squared error"

Given:  $\arg \min_w \left\{ \sum_i (\omega^T x_i - y_i)^2 + \lambda \|\omega\|_2^2 \right\}$

- L2 regularization term:  $\lambda \|\omega\|_2^2$

- Loss function of squared errors:  $\sum_i (\omega^T x_i - y_i)^2$

Step 1: Minimize  $\arg \min_w \left\{ \sum_i (\omega^T x_i - y_i)^2 + \lambda \|\omega\|_2^2 \right\}$  — Definitions

(i)  $\omega$ : model we're trying to learn  
 $x$ : feature vector  
 $y$ : output

(ii)  $\omega^T x_i$ : prediction

$y_i$ : actual output

(iii)  $\lambda \|\omega\|_2^2$ : L2 regularization term

$\therefore$  we're trying to prevent overfitting with a regularization term (independent of the data and solely a function of the parameters) which can help minimize the error

-  $\arg \min_{\omega \in \mathbb{R}} \left\{ \sum E(\omega) + \lambda E_R(\omega) \right\} = \arg \min_w \left\{ \sum (\omega^T x_i - y_i)^2 + \lambda \|\omega\|_2^2 \right\}$   
 $\rightarrow$  regularization term where  $\lambda \geq 0$

Step 2: Take derivative of loss of sum of squared errors and set it to zero

-  $g(\omega) = \arg \min_{\omega \in \mathbb{R}} \left\{ E(\omega) + \lambda \|\omega\|_2^2 \right\}$  where  $E(\omega)$  is the cost function  $\sum_i (\omega^T x_i - y_i)^2$

$\frac{dg(\omega)}{d\omega} = 0$ ;  $-x^T(y - x\omega) + 2\lambda\omega = 0$

$$2\lambda\omega = x^T(y - x\omega)$$

$$2\lambda\omega = x^T y - x^T x \omega$$

$$\omega(2\lambda + x^T x) = x^T y$$

$$\boxed{\omega = \frac{x^T y}{2\lambda + x^T x}}$$

where  $\lambda \geq 0$



(b) Part 2: How does this change the solution to the original linear regression solution?

$\therefore$  The original minimization of the squared loss sum looks as follows:

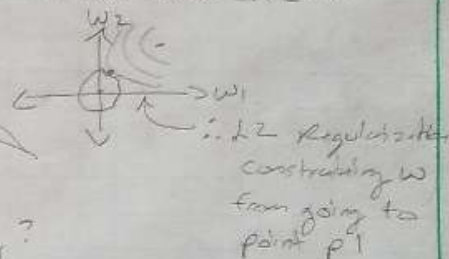
$$(i) \quad L(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 = (y - Xw)^T (y - Xw) \quad X \in \mathbb{R}^{n \times p}$$

$$(ii) \quad \text{Take } \frac{dL(w)}{dw} = 0; \quad \frac{dL(w)}{dw} = -X^T (y - Xw) = 0$$

$$(iii) \quad \text{Solve for } w; \quad w = (X^T X)^{-1} X^T y$$

— This is very similar to the newly captured  $w$  value of  $\frac{X^T y}{2\lambda + X^T X}$ .

The main difference is that we have  $2\lambda$  term in the denominator which will affect the overall value of  $w$ . Since  $\lambda \geq 0$ , we can assume that this will limit the capacity of the value  $w$  can take on (i.e. we're limiting  $w$  from getting too large).



(c) Part 3: What is the impact of adding this penalty?

$\therefore$  This is partly answered above in part 2: when we add in the penalty (i.e.  $L2$  regularization term) we are constraining how large  $w$  can get.

The  $\lambda$  term will ultimately dictate how much penalty will get enforced on  $w$  and it is given to be greater than or equal to 0.

(2) 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 2D space. Describe what happens to the training error (using all available data) when the neighbor size  $K$  varies from  $n$  to 1.

(a) Part 1: Describe what happens to the training error (using all available data) when the neighbor size  $K$  varies from  $n$  to 1

Step 1: Describe  $k$ -nearest neighbors procedure

— Procedure

(i) Select an initial volume around  $x$  that we want to estimate  $p(x)$

(ii) Grow the region until  $K$  samples fall inside the region

$$(iii) \quad p(x) = \frac{K/n}{V}$$

-(iii)  $P(x) = \frac{K/n}{V}$  describes the probability at  $x$  (i.e.  $P(x)$ ) is equal to the total number of points in the region ( $K$ ) divided by the total number of samples in data ( $n$ ) divided by the volume of the region ( $V$ ).

∴ Furthermore if  $p(x)$  is the convergence point then:

(iv) Assumption that the space will converge to  $p(x)$  the region will shrink uniformly:  $\lim_{n \rightarrow \infty} V_n = 0$

(v) Secondly,  $\lim_{n \rightarrow \infty} K_n = \infty$

(vi) Lastly,  $\lim_{n \rightarrow \infty} \frac{K_n}{n} = 0$

Step 2: Observe what happens to the training data as we vary  $k$  from 1 to  $n$

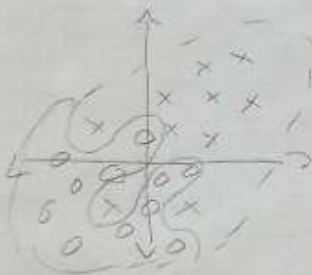
Example: Two classes,  $C1 = 0$  &  $C2 = x$



•  $C1 = 10$  training data points plotted ( $n/2$ )

•  $C2 = 10$  training data points plotted ( $n/2$ )

First case:  $K=1$



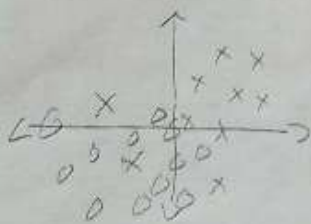
• When  $K=1$  the training error is 0, but as we can see the testing error and overfitting is quite large

Second case:  $K=10$



• When  $K=10$  we start to get training error. In the example we misclassify 4 data points out of 20 total points (10 that belong to  $C1$  and 10 that belong to  $C2$ )

Last case:  $K=20$  (or  $n$ )



• When  $K=n$  the KNN algorithm will select the all points under the class majority. For example if there are more  $C1$  points than  $C2$  points and  $K=n$  then all data will be classified to  $C1$  which will lead to more training data error. Since  $K=n$  and  $C1$  &  $C2$  are equally distributed (i.e.  $\frac{n}{2}$ ) KNN will not be able to choose a majority class which will lead to a deadlock - poor



∴ Overall, if  $k=1$  we will have zero training error but potentially large testing error. As  $k$  increases from 1 to  $n$  we begin to get training error (the training error count will increase as well as the percentage in this example since each class has  $n/2$  training data points). Even though the training error is increasing there will be a balance point (pareto optimality point) where overfitting & underfitting are normalized and the training error is acceptable. Once  $k$  reaches  $n$  then all data points will be classified to the majority class; in this example it's 50% class 1 & 50% class 2, so the algorithm will have to arbitrarily decide which classification to give all of the data points.

\* when  $k$  is too low  $\Rightarrow$  overfitting of data and small training error

\* when  $k$  is too large  $\Rightarrow$  larger training error and potentially underfitting

③ Feature selection & Preprocessing: A friend asks you for some help with a feature selection project. Your friend goes out and collects  $D$  data for their project. Using  $D$ , your friend tests many subsets  $F \subset X$  by adapting  $F$  based on the error of a classifier. They return  $F$  that corresponds to the smallest classification error. After they have the feature set  $F$  they perform randomized trials to estimate the accuracy of the model using their feature selector. This procedure they carry out to validate the impact of the feature selection routine. This procedure is repeated:

- Make a new data set  $D$  with  $F$  features using the feature selection routine
- Repeat 50 times
  - Split  $D$  into randomized training and 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

a) Part 1: Critique and respond to how your friend performed their analysis

Step 1: Describe Feature selection Process

∴ Feature selection is ultimately the task of reducing the feature set to a smaller one that contains only relevant & potentially non-redundant features.

- In an example:  $x \in \mathbb{R}^{n \times p}$  (collection of data)
  - $y \in \{ \pm 1 \}^p$  (two different states)
  - $\#p$  variables (inference for the state of the system)
    - $\rightarrow$  Relevance
    - $\rightarrow$  Irrelevance
    - $\rightarrow$  Redundancy

- In this example we want to eliminate redundant and irrelevant data and so there are various different feature selection methods we can employ

(i) Wrappers: Search through a sub-space of features using training/validation accuracy as a particular classifier

→ SUM-RFE (Recursive Feature Elimination)

essentially this algorithm can take a feature set and add or remove it based on its relevancy factor

→ Downside: it's computationally expensive

(ii) Embedded: Exploit the structure of a class of models to guide the feature selection process

→ Lasso:  $\arg \min \|y - X^T w\|_2^2 + \lambda \|w\|_1$  ;  $\|w\|_1 = \sum_{j=1}^r |w_j|$  ;  $y = w^T x$

o Applying L2 penalty to our margin of the output and predicted result  
owe can perform variable selection this way

(iii) Filters: Evaluate a classifier-independent source to evaluate feature relevance/redundancy

→ Filter assumption: ① Pick good features  
② Build the classifier

Step 2: Provide feedback based on step 1

- The critique: My friend does his due diligence by collecting many subsets

of  $F \times X$  based on the data captured  $D$  to assess the error of the classifier. He also adjusts  $F$  based on the error providing for a small feedback loop to minimize the error on the classifiers. He then returns the feature set with the smallest error. After this, the friend attempts to validate the accuracy of the model by performing the following steps:

① Make new data set  $D$  with  $F$  features

② Repeat 50 times

- Split  $D$  into random training and testing sets

- Train classifier and record its error

③ Report error averaged over 50 trials.

Overall this process is thorough but time consuming. Instead of repeatedly collecting new subsets to reduce the error (which is similar to step 1-(i) Wrapper feature selection method) my friend could have employed an embedded approach to minimize the data collection function by also applying an L2 regularization term to constrain the margin of the results. This would save him time selecting the appropriate feature set to validate the accuracy of the model. Additionally, this would reduce the necessity to run the validation step 50 times and report the error averaged over 50 times because constraining the feature selection within the embedded method we reduce training error as well as normalize potential overfitting or underfitting of data on the model.