

# Introduction to Machine Learning - Final Exam

Landon Buell

Spring 2020

## Problem 1

Assume that you have a data set with a predictor (feature)  $X$  and the target  $Y$ . You run simple linear regression ( $Y \sim X$ ) and get the best fit with  $\text{RSS}=20$  and  $\text{TSS}=120$ .

### 1. What is the covariance between $X$ and $Y$ if the variances of $X$ and $Y$ are $\text{Var}(X)=15$ and $\text{Var}(Y)=20$ ?

Using the RSS and TSS metric, we can compute the  $R^2$  statistic for the linear regression fit given the two metrics.

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{20}{120} = \frac{5}{6} \quad (1)$$

Know that  $R^2 = r^2$  which is the correlation between  $X$  and  $Y$ , thus  $\text{Cor}(X, Y) = r = \pm\sqrt{R^2} = \pm\sqrt{\frac{5}{6}}$ . Finall, from ISL, 3.18, we know how the correlation coefficient relates to the Covariance, and the two variances:

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\sigma_X^2} \sqrt{\sigma_Y^2}} \quad (2)$$

Thus we solve for the covariance between  $X$  and  $Y$ :

$$\text{Cov}(X, Y) = \text{Cor}(X, Y) \sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)} = \pm \sqrt{\frac{5}{6}(15)(20)} = \pm\sqrt{250} \quad (3)$$

Thus for  $\text{RSS} = 20$ ,  $\text{TSS} = 120$ ,  $\text{Var}(X) = 15$  and  $\text{Var}(Y) = 20$ , the resulting covariance,  $\text{Cov}(X, Y) = \pm\sqrt{250}$ .

**2. Repeat if RSS=50 and TSS=40. Discuss the result.**

Repeat the same process as above, We can compute  $R^2 = -1/4$  which onto it's own creates a problem!  $R^2$  is bound on the inter val  $[0, 1]$  so a negative number certain raises some issues. Looking at the respective definitions of RSS and TSS, it seems impossible for this to occur in practice, thus I would attribute this to a numerical error, switching the two metrics or producing a fit so bad, we've broken statistics so much that we've pushed ourselves into the complex plane for the correlation coefficient. However, continuing the calculation would show that  $\text{Cov}(X, Y) = \frac{1}{2}i\sqrt{300}$ , which again, is physically unrealistic, uninterpretable and holds no physical meaning.

## Problem 2

Suppose that we estimate the regression coefficients in a LASSO model by solving:

$$\min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \frac{1}{\lambda} \sum_{j=1}^p |\beta_j| \quad (4)$$

for a particular value of  $\lambda$ . For parts (1) through (5), indicate which of i. through v. is correct. Justify your answer.

### 1. As we increase $\lambda$ from 0.1 to $\infty$ , the training RSS will typically:

As  $\lambda$  increases, the coefficient in front of the sum decreases. This decreasing value imposes a *smaller and smaller* constraint on the lasso regression, thus the constraint region (figure 6.7 of ISL) gets smaller and smaller,  $\sum |\beta_j| \leq \frac{1}{\lambda} \rightarrow 0$  as well. As the region shrinks in sizes, it is likely to get farther and farther from the  $\hat{\beta}$ , which is at RSS predicted coefficients. Thus the training RSS will (ii.) slowly decrease.

### 2. As we increase $\lambda$ from 0.1 to $\infty$ , the test MSE will typically:

Increasing  $\lambda$  causes  $1/\lambda$  to decrease to zero. Somewhere along the way, the appropriate value for  $\lambda$  will be chosen such that MSE is reduced as much as possible. (this would be an example of cross-validation) Thus, as the constrain decreases, the resulting MSE will (v.) decrease initially, and then start increasing in a 'U' shape, as indicated by figure 6.8 in ISL.

### 3. As we increase $\lambda$ from 0.1 to $\infty$ , the (squared bias) will typically:

Again, increasing  $\lambda$  causes  $1/\lambda$  to decrease to zero. Thus decreasing the volume of the constrain region, which will cause the squared bias to (ii.) steadily decrease, as also indicated by figure 6.8 in ISL. (In their cases,  $\lambda \rightarrow 0$ )

### 4. As we increase $\lambda$ from 0.1 to $\infty$ , the variance will typically:

Again, increasing  $\lambda$  causes  $1/\lambda$  to decrease to zero. Thus decreasing the volume of the constrain region, which will cause the variance to (iii.) steadily increase, as also indicated by figure 6.8 in ISL. (In their cases,  $\lambda \rightarrow 0$ )

### 5. As we increase $\lambda$ from 0.1 to $\infty$ , the irreducible error will typically:

As  $\lambda$  is adjusted, the irreducible error will be untouched, thus (i.) remain constant.

**6. How would I choose the best value for  $\lambda$ ?**

There is no single rule for what  $\lambda$  works best for any given problem, thus I would have to implement some sort of *cross-validation* method. I would essentially choose a subset of values to test, construct a LASSO model with each value, and then choose the  $\lambda$  that results in the lowest *testing* RSS error.

## Problem 3

- 1. Do you expect random forests to achieve a smaller or larger training error than bagged trees with the same number of trees? Justify your answer.**

Bagging trees is a method used to create  $N$  different decision trees based on  $N$  different subsets of a training data set. Each subset is then used to produce  $N$  different boot strapped models,  $\hat{f}^{*1}, \hat{f}^{*2}, \dots, \hat{f}^{*N}$ . The *bagged tree* model is then a uniform sum of each of the bootstrapped models. On the other hand, a random forest only generate new leaf splits based on a subset of all possible features, as opposed to even bagged tree use all possible features for splits. This subset of predictors lets forest have wildly different trees to account for different feature strengths, whereas boosted trees will all have roughly similar leaf structures. With this, I would expect random forests to have a *smaller* training error than bagged trees, with the same number of trees.

- 2. How would you expect the training error of bagged trees to compare with boosted trees? Justify your answer.**

Bagged trees as described above, contrast with Boosted trees. Boosted trees involve taking a subset of the design matrix, and using residuals, rather than the start target. In each iteration, some small decision tree is produced based on that feature or sample subset and target and then that small tree is added to the larger model (along with a shrinkage parameters). This process repeats for  $B$  trees as outlined in ISL Algorithm 8.2. By modifying these parameters as needed, the trees can grow into all different shapes and make splits and leafs based on different parameters. In general, this would cause of the training error of boosted trees to be lower than that of bagged trees.

- 3. How would you decide how many trees to use for any particular data set in a random forest? It is best to choose the number that minimizes the training error? Why or why not?**

As always, there is no set rule to determine the number of trees to use in any problem, it varies between data sets. It is not wise to choose a number of trees that reduces the *training error* because this indicates a model that will over fit that particular training set and thus not do well on a testing set. Thus, to choose the appropriate number of trees for a given data set or problem, I would again apply some sort of *cross-validation* procedure to find what number of decision trees in a random forest reduces the *testing error* when trained on a similar data set. This would allow for the determining of the correct number of trees, along with a roughly minimized testing error.

## Problem 4

This problem examines the differences between SVC (linear SVM), SVM with a polynomial kernel, and other linear classifiers.

**1. For an arbitrary training set, would you expect for SVC (linear) or SVM (polynomial kernel) to work better on the training set? Why?**

For any given data set, I would predict that a polynomial kernel SVM would work better on a training set. By design, a linear Support Vector Classifier, as the name implies, uses linear support vectors to create it's decision boundary. This model will only have a low training and testing error if the decision boundary between features is linear or can be approximated as linear. In most real-world data sets do not have linear decision boundaries.

On the other hand, a polynomial kernel Support Vector Machine, is built off of polynomial decision boundaries. Adding more terms and features is guaranteed to decrease the training error, but not necessarily the test error. Furthermore, In some cases, features with almost-linear decision boundary can be approximated with a polynomial where higher degrees have extremely small coefficients (Almost like Taylor-Series). In this cases, a polynomial kernel would still perform well on a linearly separable data set.

**2. If the Bayes decision boundary between the two classes is linear, would you expect SVC or SVM to work better on the training set?**

If we know the decision boundary is linear, then a SVC, which assumes a linear decision boundary, would be the optimal choice. However, an SVM, which often produces a non-linear boundary might perform better on a *training* set because it has the advantage of a greater number of features, which will always reduce the training error. Thus a linear SVC would do better on a test set and a polynomial SVM would still do better on a training set due to its flexibility to fit that set.

**3. True or False: There is no need to use slack variables in SVMs with polynomial kernels because the decision boundary can be non-linear. Justify your answer.**

**4. LDA, Logistic regression, and SVC (linear) all fit a linear decision boundary. Will their fits be the same? What would be some reasons for you to prefer LDA over SVC?**

The fits between LDA, Log Reg and SVC are all based on linear decision boundaries, but each one is constructed differently and comes with different purposes - thus the fits are not guaranteed to be the same, and in most cases, will not be the same. For example, Logistic Regression a probability concept based on the logistic function. The value of the function is then used as a percentage probability that a sample existent within a certain class given a feature space, split along the 50/50 line, which creates the decision boundary. Support Vector Classifiers also build decision boundaries but also included maximum margin ranges, which are the support vectors. This does not really include a formal probability metrics, but does add a tolerance band for class features.

In the case of LDA, it also attempts to build a class prediction based on a set of predictors, but allows for a few other ideas that may make it favorable to SVC's in a few contexts: As stated in ISL, pg. 138, LDA performs well with well-defined decision boundaries, and offers a probability for it's sample to be in a class, given it's predictors, whereas a SVC mostly gives a decision boundary and a maximum margin tolerance. Additionally, if the predictors in a data set are vaguely Gaussian, LDA will perform well because this distribution is one of it's core assumptions.

## Problem 5

1. Show that the method of least absolute deviations is the maximum likelihood solution to fitting a linear function to a target when assuming that the prediction error is distributed according to the Laplace distribution.
2. Is the method of least absolute deviations the same as or different from median regression (quantile regression with  $\tau = 1/2$ )? prove it?