

91

# ST 563 601 – SPRING 2025 – POST Final Exam Tablet

Student's Name:

Dev Kewlani

Date of Exam: Monday, April 28, 2025 - Wednesday, April 30, 2025

Time Limit: 90 minutes

Allowed Materials: None (closed book & closed notes)

## Student – NC State University Pack Pledge

I, Dev Kewlani, have neither given nor received unauthorized aid on this exam or assignment. I have read the instructions and acknowledge that this is the correct exam.

STUDENT'S PRINTED NAME

Dev Kewlani

STUDENT'S SIGNATURE

DATE

**Exam must be turned in by:**

3:14 pm

EXAM END TIME

DKL

STUDENT'S  
INITIAL  
AGREEMENT

**NOTE: Failure to turn in exam  
on time may result in penalties  
at the instructor's discretion.**

## Final Exam

Please write your answers below each question. You should not have access nor use any materials during this exam.

A reminder that, by taking this exam, you are required to uphold the NC State honor pledge:  
"I have neither given nor received unauthorized aid on this test or assignment."

1. We know that a multiple linear regression model fits a (hyper) plane as the response surface (or a curved hyperplane with higher order polynomial or interaction terms). How does a standard regression tree model the response surface?

The standard regression tree models the responses in rectangles in a plane ✓ where each rectangle represents a region of a split we built the tree on. We can have many such rectangles & regions based on the splits we do.

2. For a standard regression tree that uses recursive binary splitting, suppose we have two predictors  $X_1$  and  $X_2$ . What criterion is used to determine the first split? Describe how this first split is decided upon. Be specific on both of these!

we use a greedy algorithm, so basically for full range of values for each  $X_1$  &  $X_2$  we decide on a value and consider that value as split point, so for eg. for  $X_1$ : split at  $c$ , for  $X_2$ : split at  $d$ . then we have regions of the form -  $X_1 < c \quad X_1 \geq c$  or  $X_2 < d \quad X_2 \geq d$  we send all values to whichever split they belong to and take the mean of all

values in each split, that's our prediction for that split.

so we have  $\hat{Y}(x_1 < c)$ ,  $\hat{Y}(x_1 \geq c)$   
or

$\hat{Y}(x_2 < d)$ ,  $\hat{Y}(x_2 \geq d)$

then we find residuals comparing each  $y_i$  with  $\hat{Y}$ , and sum up across both  $(y_i - \hat{Y}_{(x_1 < c)})^2 + (y_i - \hat{Y}_{(x_1 \geq c)})^2$

$$y_i - \hat{Y}_{(x_2 < d)}^2 + (y_i - \hat{Y}_{(x_2 \geq d)})^2$$

whichever from  $x_1$  or  $x_2$  yields a lower RSS is selected as the split and we do this for a range of values for both  $x_1$  and  $x_2$ . Then we decide on our first split.

3. Suppose we have a large data set where we want to perform a regression task. We want to determine the best overall model between a kNN model and a ridge regression model. We want to use a train test split and compare the best kNN and ridge regression model on the test set. We wish to determine the appropriate tuning parameters on the training set only using the bootstrap. Fully outline the process for splitting the data, tuning, comparing, and fitting a final overall best model.

- 1) Split the data between ~~train and test~~ (80%) ~~(20%)~~
- 2) Since we are using bootstrap; we will draw 10 bootstrap samples of training data [which ~~is~~ sampling training data with replacement and sample size same as training data], so some observations will naturally be left out. That is our OOB samples for validation.
- 3) we define a grid of both  $K$  &  $\lambda$  values.
- 4) for each  $K$  and each  $\lambda$ , we train KNN & ridge regression respectively on the bootstrap sample and ~~test~~ the performance on OOB samples, we keep track of RSS across all 10 OOB samples.
- 5) we repeat this for all  $K$  and all  $\lambda$ .
- 6) So let's say we have a grid of 10  $K$  and 10  $\lambda$  values, we get a ~~total of 20 RSS~~ values averaged across ~~O~~

10 bootstrapped OOB observations.

- 7) we select the best  $\lambda$  and best  $K$   
whichever has the lowest RSS among  
all ridge and all KNN
- 8) we then train the best  $K$  KNN  
model and best  $\lambda$  ridge model on ~~the~~  
the full ~~data~~ training data (80%)  
and evaluate its performance on  
test set (20% remaining)
- 9) whichever yields better results  
we select that model.
- 10) We train that model on the  
full dataset. [This is the best  
model].

4. We discussed two ways to do 'early stopping' in a regression or classification tree.  
What are those two methods?
- 1) Either we stop when ~~one~~ <sup>our</sup> terminal nodes exceed a threshold.
  - 2) we can stop when number of observations in leaf node is less than a certain threshold.
5. In a standard multilayer feed-forward neural network, what are two common activation functions?

- 1) ReLU ✓
- 2) Tanh

6. What task is a Recurrent neural network well-suited for?

- ~~Time Series Tasks~~ Next token prediction  
Stock market prediction
- ~~when the output we want to model is sequential.~~

7. True or False questions (write True or false next to each letter):

- a. Random forest and bagged tree models generally require you to standardize your predictors ~~False~~ True
- b. kNN models generally require you to standardize your predictors ~~False~~ True
- c. The number of trees we use in a random forest model is important because we can overfit with too many trees. ~~False~~ True
- d. When using BART we need to remove the first few prediction models. ~~True~~ True
- e. SVM models can only be used in classification tasks. ~~False~~ True
- f. KMeans clustering does not necessarily create the same clusters in each run of the algorithm. ~~True~~ True
- g. Hierarchical clustering requires you to know the 'true' underlying groupings to use it effectively. ~~False~~ True
- h. In a standard multilayer neural network, all inputs are 'connected to' all first level activations. ~~False~~ True
- i. KNN provides a discriminant for classifying our observations ~~False~~ True
- j. The Naive Bayes provides a discriminant for classifying our observations ~~True~~ True

8. Consider the piecewise polynomial regression model. Here we define our knots to be  $c_1, \dots, c_M$  and use the indicator functions

$h_1(X) = I(c_1 \leq X < c_2), \dots, h_{M-1} = I(c_{M-1} \leq X < c_M), h_M(X) = I(X > c_M)$   
in our regression equation given by

$$Y_i = \beta_0 + h_1(X_i)\beta_1 + \dots + h_M(X_i)\beta_M + \epsilon_i$$

Suppose we have  $n$  observations and we fit the model.

- a. What is the estimate of  $\beta_0$  in this model?

$$\beta_0 = \underbrace{\sum_{i=1}^n I(X_i \leq c_1)}_{\checkmark}$$

- b. What is the estimate of  $\beta_1$  in the model?

$$\text{For } \beta_1 = \beta_2 - \beta_M = 0 \Rightarrow c_1 \leq X_i < c_2$$

$$\beta_1 = \frac{\sum_{i=1}^n I(c_1 \leq X_i < c_2)}{n} - \frac{\sum_{i=1}^n I(X_i \geq c_M)}{n}$$

9. What are the three most common tuning parameters associated with a boosted tree model?

Tree Depth : how much should the tree be extended vertically based on how many splits. [more depth more overfit]

Learning Rate : how fast or slow we change the movement of our current prediction.

No. of trees to fit : more trees, more overfitting

Subsampling ratio : how many rows or columns to select for building the current tree [Idea of stochastic gradient Boosting]

10. Why do random forests for a regression task generally improve prediction over the basic bagged tree model?

By considering only a subset of columns at each split (usually  $m = \sqrt{p}$ ) it decorrelates individual trees, which lead to more improved stable estimates when we take mean across all trees.

11. Describe the algorithm for fitting a basic boosted regression tree model.

- 1) Start with  $\hat{f}^0(x) = 0 \rightarrow r_i = \text{residual} = y_i$
- 2) for tree  $t = 1 \dots T$   
fit a tree [usually a small tree as weak learner]  
get prediction  $\hat{f}^t(x)$ ,  
update  $\hat{f}(x) = \hat{f}^0(x) + \lambda \hat{f}^t(x)$   
update  $r_i = r_{i-1} - \lambda \hat{f}^t(x)$   
continue this process for all iteration [Trees T]
- 3) for a new sample, our prediction is,  
$$\hat{f}(x) = \sum_{t=1}^T \lambda \hat{f}^t(x)$$

12. When fitting a support vector machine model for classification, what are support vectors?

Support vectors are those data points which lie inside the margin or on the wrong side of hyperplane. Our margin is influenced only by these data points. Based on our cost, we can have a wider or

a narrower margin. More wider margin leads to more points violating & more support vectors, less wide & narrow margin & less support vectors.

13. When we wish to apply the SVM model to a classification task with more than two levels, we discussed the one-versus-one approach. Describe how this SVM model works.

let's say we have 3 classes, {1, 2, 3}. Then we train 3 models, 1 vs (2 & 3), 2 vs (1 & 3), 3 vs (1 & 2). where (n vs) means that they are assigned some class if they are assigned some class. By doing this we get probability for label  $\hat{y}$ , doing this we get probability for each class based on our model fit. Which was the highest - we select that class.

14. Why do we often run the kmeans clustering algorithm multiple times?

Because k-means is heavily influenced by the initial configuration that we start with when we assign each data point in one of the  $k$  clusters, hence it is important to run multiple times to select the one with minimum loss to get confidence in our algo.

15. When doing hierarchical clustering, how does the 'single' linkage create a dissimilarity measure?

- we find all pairwise dissimilarities between observations from cluster A and cluster B.  
- we then select the minimum of all these pairwise dissimilarities to get our dissimilarity

16. What is a biplot and how can it be useful?

A biplot is a powerful tool to visualize the (direction) loadings of principal components and data points of original features that existed in our original space. If it is an effective

Visualization tool that helps us understand  
and see which principal component have  
how much coefficients or to put it loosely,  
how much content in which direction  
and  
from the original feature space.

e.g.: PC1 has  $0.2x_1 - 0.2x_2 + 0.7x_3$  then  
we can see the directions of ~~the~~  
these original features and PC1 w.r.t  
data points through the biplot.