

# STAT 471: Homework 2

James Kuang

Due: October 4, 2021 at 11:59pm

## Contents

<b>Instructions</b>	<b>2</b>
Setup . . . . .	2
Collaboration . . . . .	2
Writeup . . . . .	2
Programming . . . . .	2
Grading . . . . .	2
Submission . . . . .	3
<b>1 Case study: Bone mineral density (40 points for correctness; 10 points for presentation)</b>	<b>4</b>
1.1 Import (2 points) . . . . .	4
1.2 Explore (10 points) . . . . .	4
1.3 Model (15 points) . . . . .	6
1.4 Evaluate (6 points) . . . . .	8
1.5 Interpret (7 points) . . . . .	9
<b>2 KNN and bias-variance tradeoff (45 points for correctness; 5 points for presentation)</b>	<b>10</b>
Setup: Apple farming . . . . .	10
2.1 A simple rule to predict this season's yield (15 points) . . . . .	12
2.2 K-nearest neighbors regression (conceptual) (15 points) . . . . .	12
2.3 K-nearest neighbors regression (simulation) (15 points) . . . . .	12

# Instructions

## Setup

Pull the latest version of this assignment from Github and set your working directory to `stat-471-fall-2021/homework/homework-2`. Consult the [getting started guide](#) if you need to brush up on R or Git.

## Collaboration

The collaboration policy is as stated on the Syllabus:

“Students are permitted to work together on homework assignments, but solutions must be written up and submitted individually. Students must disclose any sources of assistance they received; furthermore, they are prohibited from verbatim copying from any source and from consulting solutions to problems that may be available online and/or from past iterations of the course.”

In accordance with this policy,

*Please list anyone you discussed this homework with:*

Alex Chen

*Please list what external references you consulted (e.g. articles, books, or websites):*

Stack Overflow, R Documentation

## Writeup

Use this document as a starting point for your writeup, adding your solutions after “**Solution**”. Add your R code using code chunks and add your text answers using **bold text**. Consult the [preparing reports guide](#) for guidance on compilation, creation of figures and tables, and presentation quality.

## Programming

The `tidyverse` paradigm for data wrangling, manipulation, and visualization is strongly encouraged, but points will not be deducted for using base R.

We’ll need to use the following R packages:

```
library(tidyverse) # tidyverse
library(kableExtra) # for printing tables
```

```
## Warning: package 'kableExtra' was built under R version 4.1.1
```

```
library(cowplot) # for side by side plots
```

```
## Warning: package 'cowplot' was built under R version 4.1.1
```

```
library(FNN) # for K-nearest-neighbors regression
```

```
## Warning: package 'FNN' was built under R version 4.1.1
```

We’ll also need the `cross_validate_spline` function from Unit 2 Lecture 3:

```
source("../..functions/cross_validate_spline.R")
```

## Grading

The point value for each problem sub-part is indicated. Additionally, the presentation quality of the solution for each problem (as exemplified by the guidelines in Section 3 of the [preparing reports guide](#) will be evaluated on a per-problem basis (e.g. in this homework, there are three problems). There are 100 points possible on this homework, 85 of which are for correctness and 15 of which are for presentation.

## Submission

Compile your writeup to PDF and submit to [Gradescope](#).

# 1 Case study: Bone mineral density (40 points for correctness; 10 points for presentation)

In this exercise, we will be looking at a data set (available [online](#)) on spinal bone mineral density, a physiological indicator that increases during puberty when a child grows.

Below is the [data description](#):

“Relative spinal bone mineral density measurements on 261 North American adolescents. Each value is the difference in spnbmd taken on two consecutive visits, divided by the average. The age is the average age over the two visits.”

Variables:

**idnum**: identifies the child, and hence the repeat measurements

**age**: average age of child when measurements were taken

**gender**: male or female

**spnbmd**: Relative Spinal bone mineral density measurement

The goal is to learn about the typical trends of bone mineral density during puberty for boys and girls.

## 1.1 Import (2 points)

- Using `readr`, import the data from the above URL into a tibble called `bmd`. Specify the column types using the `col_types` argument.
- Print the imported tibble (no need to use `kable`).

```
library(readr) # Read in the data
bmd <- read_tsv("https://web.stanford.edu/~hastie/ElemStatLearn/datasets/bone.data")
bmd
```

```
## # A tibble: 485 x 4
##   idnum   age gender  spnbmd
##   <dbl> <dbl> <chr>    <dbl>
## 1     1   11.7 male    0.0181
## 2     1   12.7 male    0.0601
## 3     1   13.8 male    0.00586
## 4     2   13.2 male    0.0103
## 5     2   14.3 male    0.211
## 6     2   15.3 male    0.0408
## 7     3   11.4 male   -0.0296
## 8     3   12.4 male   -0.00643
## 9     3   13.4 male    0.0566
## 10    4   10.6 female  0.108
## # ... with 475 more rows
```

## 1.2 Explore (10 points)

- To keep things simple, let's ignore the fact that we have repeated measurements on children. To this end, remove the `idnum` column from `bmd`.
- What is the number of boys and girls in this dataset (ignoring the fact that there are repeated measurements)? What are the median ages of these boys and girls?

- Produce boxplots to compare the distributions of `spnbmd` and `age` between boys and girls (display these as two plots side by side, one for `spnbmd` and one for `age`). Are there apparent differences in either `spnbmd` or `age` between these two groups?
- Create a scatter plot of `spnbmd` (y axis) versus `age` (x axis), faceting by `gender`. What trends do you see in this data?

```
bmd['idnum'] <- NULL # Removing the idnum column
summarized_bmd <- bmd %>% # Summarizing the data
  group_by(gender) %>%
  summarise(
    count = n(), # Counting rows
    median_age = median(age) # Finding Median
  )
```

The number of girls and boys are 259 and 226, respectively. The median age is 15.3 and 15.6 for girls and boys, respectively.

```
#Creating a boxplot graph on gender and age
box_age <- ggplot(bmd, aes(x = gender, y = age)) +
  geom_boxplot() + # Creating the plot
  theme_bw() + # Formatting Changes
  scale_x_discrete(labels = c("Female", "Male")) +
  labs(y= "Age", x = "Gender")
#Creating a boxplot graph on gender and age
box_spnbmd <- ggplot(bmd, aes(x = gender, y = spnbmd)) +
  geom_boxplot() + # Creating the plot
  theme_bw() + # Formatting fixes
  scale_x_discrete(labels = c("Female", "Male")) +
  labs(y= "Relative Spinal Bone Mineral Density", x = "Gender")
plot_grid(box_age, box_spnbmd) # Putting both plots together
```

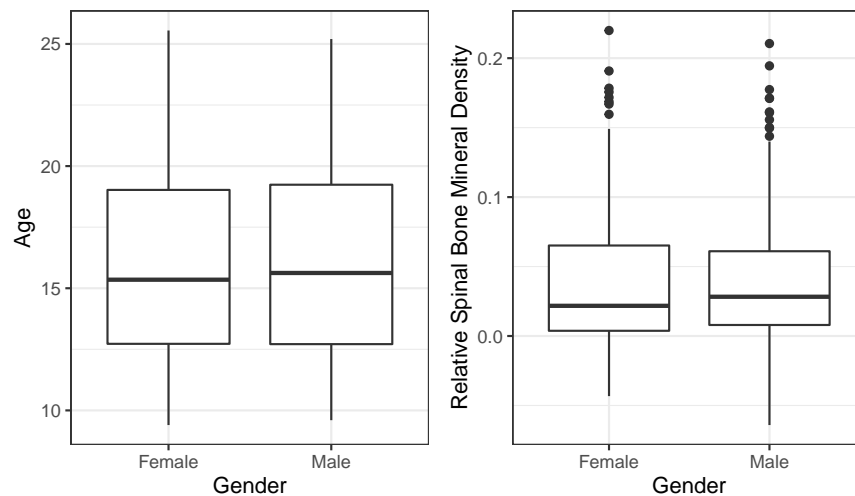


Figure 1: Boxplots for Age and Relative Spinal Bone Mineral Density by Gender

There are a few differences between the two genders. In Figure 1, the median for age and relative spinal bone mineral density appear a bit lower for females than men. The distribution of female relative spinal bone mineral density also appears a bit higher than male's.

```
# Creating a scatter plot based on age and spinal bone density
ggplot(bmd, aes(x=age, y=spnbmd)) +
  geom_point() +
  labs(x = "Age", y = "Relative Spinal Bone Mineral Density") +
  theme_bw() +
  facet_wrap(~gender)
```

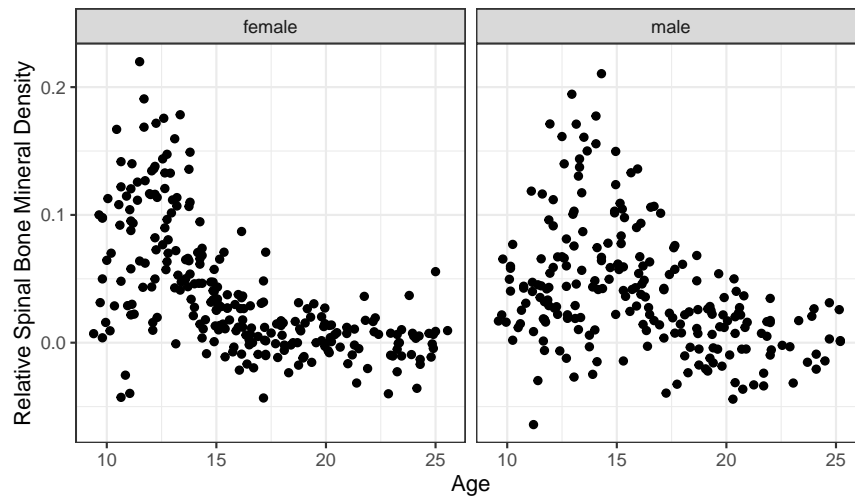


Figure 2: Relative Spinal Bone Mineral Density by Age for Men and Women

In Figure 2, for both genders, it appears that relative spinal bone mineral density decreased as the subjects got older. This may be bone density growth leveling off upon reaching some number.

### 1.3 Model (15 points)

There are clearly some trends in this data, but they are somewhat hard to see given the substantial amount of variability. This is where splines come in handy.

#### 1.3.1 Split

To ensure unbiased assessment of predictive models, let's split the data before we start modeling it.

- Split `bmd` into training (80%) and test (20%) sets, using the rows in `train_samples` below for training. Store these in tibbles called `bmd_train` and `bmd_test`, respectively.

```
set.seed(5) # seed set for reproducibility (DO NOT CHANGE)
n = nrow(bmd)
train_samples = sample(1:n, round(0.8*n))
```

```
bmd_train <- bmd[train_samples,] # Creating a train tibble
bmd_test <- bmd[-train_samples,] # Creating a test tibble
```

#### 1.3.2 Tune

- Since the trends in `spnbmd` look somewhat different for boys than for girls, we might want to fit separate splines to these two groups. Separate `bmd_train` into `bmd_train_male` and `bmd_train_female`, and likewise for `bmd_test`.

- Using `cross_validate_spline` from Lecture 3, perform 10-fold cross-validation on `bmd_train_male` and `bmd_train_female`, trying degrees of freedom 1,2,...,15. Display the two resulting CV plots side by side.
- What are the degrees of freedom values minimizing the CV curve for boys and girls, and what are the values obtained from the one standard error rule?
- For the sake of simplicity, let's use the same degrees of freedom for males as well as females. Define `df.min` to be the maximum of the two `df.min` values for males and females, and define `df.1se` likewise. Add these two spline fits to the scatter plot of `spnbmd` (y axis) versus `age` (x axis), faceting by `gender`.
- Given our intuition for what growth curves look like, which of these two values of the degrees of freedom makes more sense?

```
# Splitting everything into male and female
bmd_train_male <- bmd_train %>% filter(gender == "male")
bmd_train_female <- bmd_train %>% filter(gender == "female")
bmd_test_male <- bmd_test %>% filter(gender == "male")
bmd_test_female <- bmd_test %>% filter(gender == "female")

# Using the function to find the best splines for male and females
male_fit <- cross_validate_spline(bmd_train_male$age,
                                bmd_train_male$spnbmd, 10, 1:15)
female_fit <- cross_validate_spline(bmd_train_female$age,
                                   bmd_train_female$spnbmd, 10, 1:15)

# Making the plots better
male_plot <- male_fit$cv_plot +
  labs(x = "Degrees of Freedom (Male Fit)")
female_plot <- female_fit$cv_plot +
  labs(x = "Degrees of Freedom (Female Fit)")
#Actually plotting the graphs
plot_grid(male_plot, female_plot)
```

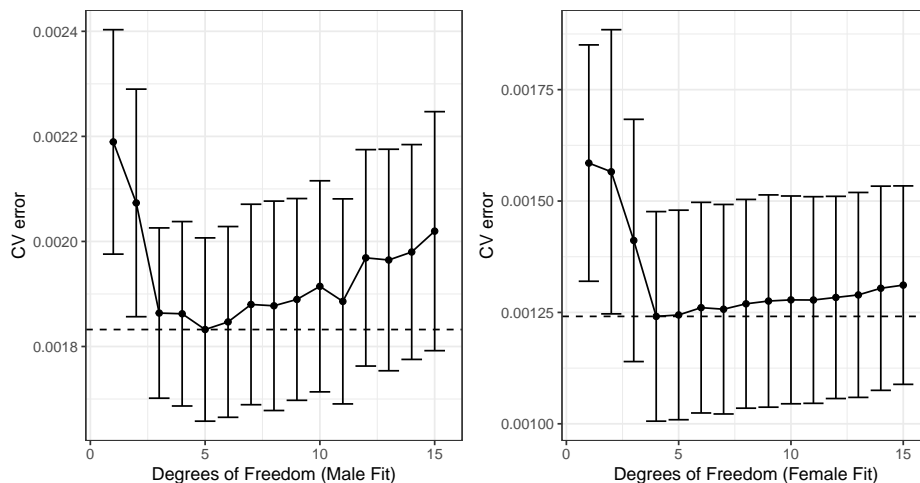


Figure 3: CV Error for Different Degrees of Freedom by Gender

From Figure 3, it is possible to see that the degrees of freedom minimizing the CV error are 5 and 4 for boys and girls, respectively. The values obtained from the one standard error rule, however, are 3 for both boys and girls.

```
# Creating a scatter plot based on age and spinal bone density
ggplot(bmd, aes(x=age, y=spnbmd)) +
  geom_point() + # The scatter plot
  geom_smooth(method = "lm", # Plotting the df=5 fit as a red line
    formula = "y ~ splines::ns(x, df = 5)",
    se = FALSE, aes(color = "5")) +
  geom_smooth(method = "lm", # Plotting the df=3 fit as a blue line
    formula = "y ~ splines::ns(x, df = 3)",
    se = FALSE, aes(color = "3")) +
  facet_wrap(~gender) + # Facetting by gender
  labs(x = "Age", y = "Relative Spinal Bone Mineral Density") +
  theme_bw() +
  scale_colour_manual(name="Degrees of Freedom", values=c("blue", "red"))
```

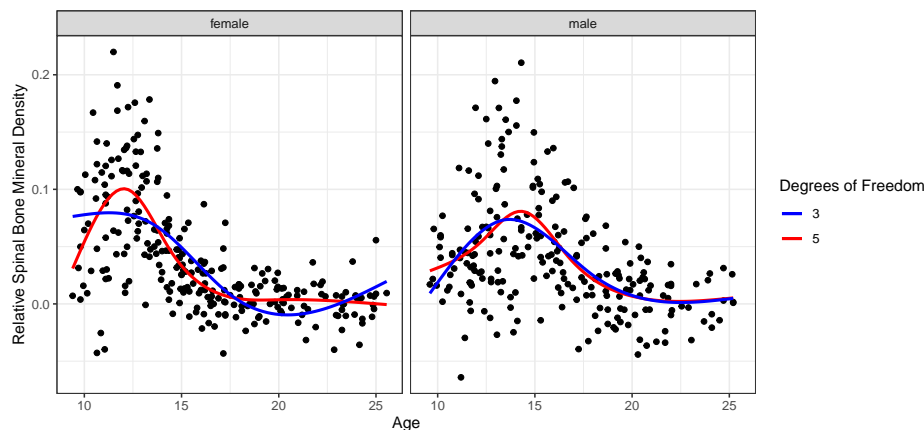


Figure 4: The Splines Fitted to the Data by Gender

In Figure 4, the fit with 5 degrees of freedom appears to be the best. Both splines do roughly match our intuition. People should be growing through puberty, and stop once they reach their twenties. However, with 3, there is an uptick for women once they reach past 20 years old. That seems strange given that most people should have stopped growing at that point, so 5 degrees of freedom fits better with my intuition.

### 1.3.3 Final fit

- Using the degrees of freedom chosen above, fit final spline models to `bmd_train_male` and `bmd_train_female`.

```
# Fit the splines for both men and women
spline_fit_male = lm(spnbnmd ~ splines::ns(age, df = 5), data = bmd_train_male)
spline_fit_female = lm(spnbnmd ~ splines::ns(age, df = 5), data = bmd_train_female)
```

## 1.4 Evaluate (6 points)

- Using the final models above, answer the following questions for boys and girls separately: What percent of the variation in `spnbmd` is explained by the spline fit in the training data? What is the training RMSE? What is the test RMSE? Print these three metrics in a nice table.
- How do the training and test errors compare? What does this suggest about the extent of overfitting that has occurred?



```

library(Metrics) # Using a package for this
# Predicting the values for the train error for men
y_hat_train_male <- predict(spline_fit_male, newdata = bmd_train_male)
# Getting the error using the library
train_rmse_male <- rmse(bmd_train_male$spnbmd, y_hat_train_male)
# Predicting the values for the test error
y_hat_test_male <- predict(spline_fit_male, newdata = bmd_test_male)
# Getting the error using the library
test_rmse_male <- rmse(bmd_test_male$spnbmd, y_hat_test_male)

# Predicting the values for the train error for women
y_hat_train_female <- predict(spline_fit_female, newdata = bmd_train_female)
# Getting the error using a formula
train_rmse_female <- rmse(bmd_train_female$spnbmd, y_hat_train_female)
# Predicting the values of the test error
y_hat_test_female <- predict(spline_fit_female, newdata = bmd_test_female)
# Getting the error using the library
test_rmse_female <- rmse(bmd_test_female$spnbmd, y_hat_test_female)

# Printing the test metrics in a table
tribble(
  ~Statistic, ~Male, ~Female,
  #-----/-----/-----
  "% of Variation Explained", summary(spline_fit_male)$r.squared,
    summary(spline_fit_female)$r.squared,
  "Training RMSE", train_rmse_male, train_rmse_female,
  "Testing RMSE", test_rmse_male, test_rmse_female
) %>%
  kable(format = "latex", row.names = NA,
    booktabs = TRUE, digits = 2,
    caption = "Metrics for the Final Fit on Test Data") %>%
  kable_styling(position = "center", latex_options = "HOLD_position")

```

Table 1: Metrics for the Final Fit on Test Data

Statistic	Male	Female
% of Variation Explained	0.30	0.51
Training RMSE	0.04	0.03
Testing RMSE	0.03	0.04

In Table 1, the test and training errors are quite close together for both male and female fits. This suggests that there is very little overfitting.

## 1.5 Interpret (7 points)

- Using the degrees of freedom chosen above, redo the scatter plot with the overlaid spline fits, this time without faceting in order to directly compare the spline fits for boys and girls. Instead of faceting, distinguish the genders by color.
- The splines help us see the trend in the data much more clearly. Eyeballing these fitted curves, answer the following questions. At what ages (approximately) do boys and girls reach the peaks of their growth spurts? At what ages does growth largely level off for boys and girls? Do these seem in the right

ballpark?

```
# Plotting the data with the final spline
ggplot(bmd, aes(x=age, y=spnbmd, color = gender)) +
  geom_point() +
  geom_smooth(method = "lm",
    formula = "y ~ splines::ns(x, df = 5)", se = FALSE) +
  theme_bw() +
  labs(x = "Age", y = "Relative Spinal Bone Mineral Density", color = "Gender")
```

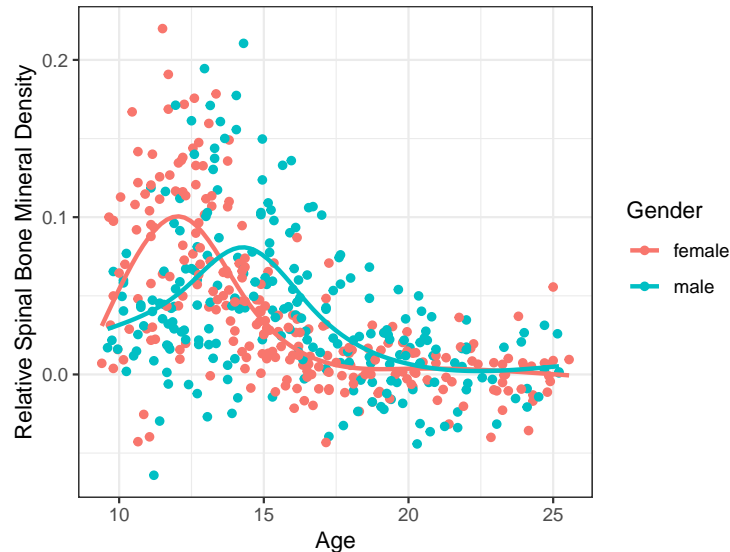


Figure 5: Plotted Data with Splines by Gender

In Figure 5, the peak of the growth curve is around 12 for women, and their growth curve levels off at around 17. The peak is around 14 for men, while their curve levels off at around 20. This makes sense as it's around when people enter and exit puberty. I also remember from high school biology that women tend to enter puberty earlier than men, so the fact that the female curve peaks earlier than the male one also makes sense.

## 2 KNN and bias-variance tradeoff (45 points for correctness; 5 points for presentation)

### Setup: Apple farming

You own a square apple orchard, measuring 200 meters on each side. You have planted trees in a grid ten meters apart from each other. Last apple season, you measured the yield of each tree in your orchard (in average apples per week). You noticed that the yield of the different trees seems to be higher in some places of the orchard and lower in others, perhaps due to differences in sunlight and soil fertility across the orchard.

Unbeknownst to you, the yield  $Y$  of the tree planted  $X_1$  meters to the right and  $X_2$  meters up from the bottom left-hand corner of the orchard has distribution

$$Y = 50 + 0.001X_1^2 + 0.001X_2^2 + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad \sigma = 4.$$

The data you collected are as in Figure 6.

The underlying trend is depicted in Figure 7, with the top right-hand corner of the orchard being more fruitful.

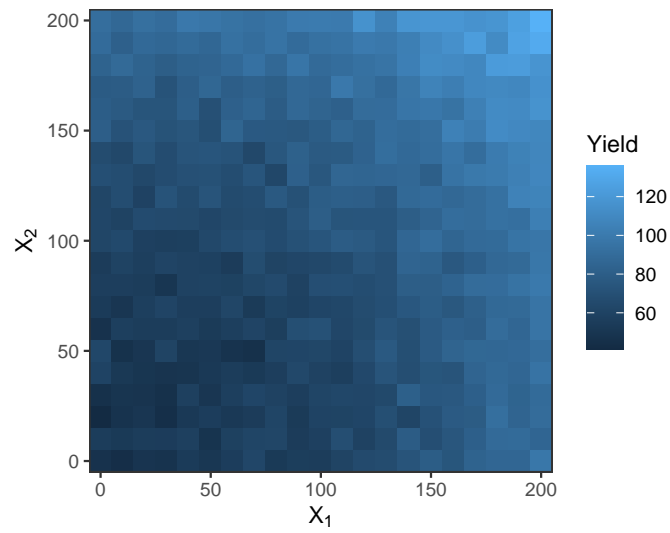


Figure 6: Apple tree yield for each 10m by 10m block of the orchard in a given year.

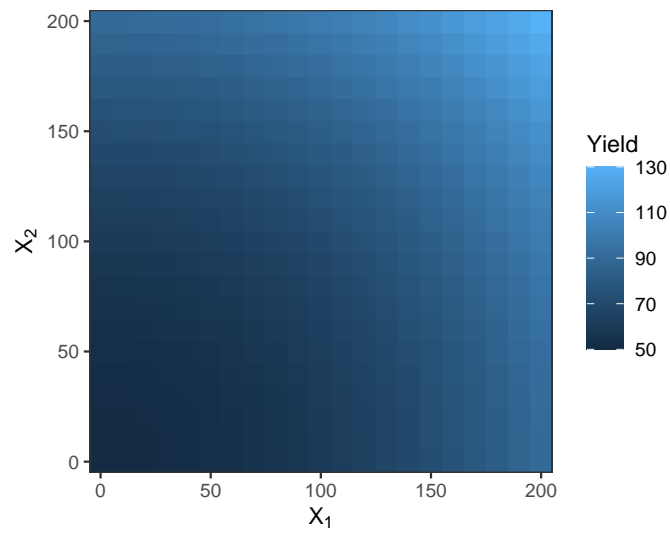


Figure 7: Underlying trend in apple yield for each 10m by 10m block of the orchard.

## 2.1 A simple rule to predict this season's yield (15 points)

This apple season is right around the corner, and you'd like to predict the yield of each tree. You come up with perhaps the simplest possible prediction rule: predict this year's yield for any given tree based on last year's yield from that same tree. Without doing any programming, answer the following questions:

- What is the expected training error of such a rule?
- Averaged across all trees, what is the squared bias, variance, and ETE of this prediction rule?
- Why is this not the best possible prediction rule?

**The ETE of such a rule would be 0. The squared bias would be 0. The variance would be 16. The ETE then would be 32. This isn't the best possible prediction rule as you are overfitting on the data. You are only guessing by 1 tree, which is subject to a large variance.**

## 2.2 K-nearest neighbors regression (conceptual) (15 points)

As a second attempt to predict a yield for each tree, you average together last year's yields of the  $K$  trees closest to it (including itself, and breaking ties randomly if necessary). So if you choose  $K = 1$ , you get back the simple rule from the previous section. This more general rule is called *K-nearest neighbors (KNN) regression* (see ISLR p. 105).

KNN is not a parametric model like linear or logistic regression, so it is a little harder to pin down its degrees of freedom.

- What happens to the model complexity as  $K$  increases? Why?
- The degrees of freedom for KNN is sometimes considered  $n/K$ , where  $n$  is the training set size. Why might this be the case? [Hint: consider a situation where the data are clumped in groups of  $K$ .]
- Conceptually, why might increasing  $K$  tend to improve the prediction rule? What does this have to do with the bias-variance tradeoff?
- Conceptually, why might increasing  $K$  tend to worsen the prediction rule? What does this have to do with the bias-variance tradeoff?

**The model complexity decreases as  $K$  increases. If degrees of freedom are how many parameters you are fitting into a model, then you could consider all groups of  $k$  close together points to each be a single parameter. Thus, the degrees of freedom is  $N/k$ . Increasing  $K$  could improve the prediction rule by reducing the amount of overfitting that is being done. This means decreasing variance at the cost of some bias (i.e. trading off variance for some bias). Increasing  $K$  could also worsen the prediction rule by leading to underfitting. This would still be trading off variance for some bias, but, in this case, addition of bias would be the more important factor.**

## 2.3 K-nearest neighbors regression (simulation) (15 points)

Now, we try KNN for several values of  $K$ . For each, we compute the bias, variance, and ETE for each value based on 50 resamples. The code for this simulation, provided for you below (see Rmd file; code omitted from PDF for brevity), results in Figure 8.

- Based on Figure 8, what is the optimal value of  $K$ ?
- We are used to the bias decreasing and the variance increasing when going from left to right in the plot. Here, the trend seems to be reversed. Why is this the case?
- The squared bias has a strange bump between  $K = 1$  and  $K = 5$ , increasing from  $K = 1$  to  $K = 2$  but then decreasing from  $K = 2$  to  $K = 5$ . Why does this bump occur? [Hint: Think about the rectangular grid configuration of the trees. So for a given tree, the closest tree is itself, and then the next closest four trees are the ones that are one tree up, down, left, and right from it.]

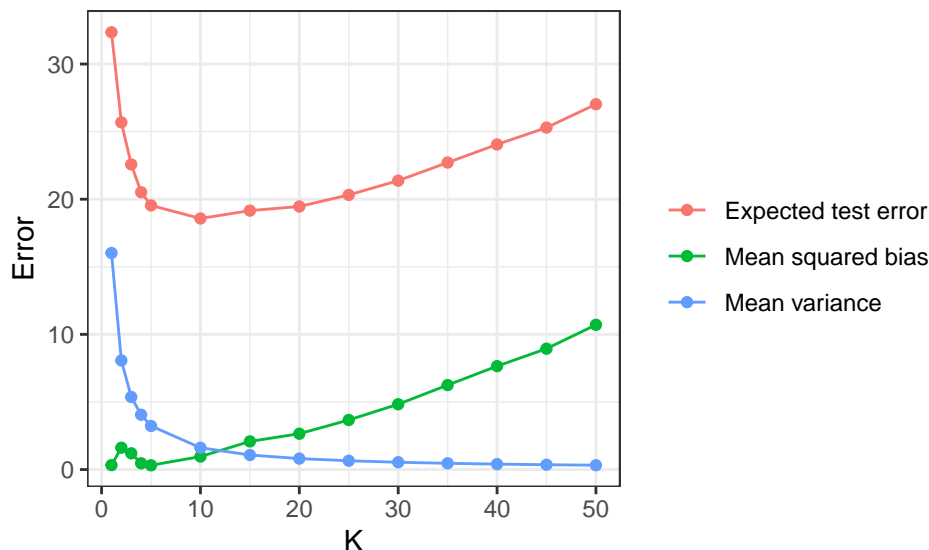


Figure 8: Bias-variance trade-off for KNN regression.

- The data frame `training_results_summary` contains the bias and variance for every tree in the orchard, for every value of  $K$ . Which tree and which value of  $K$  gives the overall highest absolute bias? Does the sign of the bias make sense? Why do this particular tree and this particular value of  $K$  give us the largest absolute bias?
- Redo the bias-variance plot above, this time putting  $df = n/K$  on the x-axis. What do we notice about the variance as a function of  $df$ ? Derive a formula for the KNN variance and superimpose this formula onto the plot as a dashed curve. Do these two variance curves match? [Hint: To derive the KNN variance, focus first on the prediction of a single tree. Recall the facts that the variance of the sum of independent random variables is the sum of their variances, and that the variance of a constant multiple of a random variable is the square of that constant times its variance.]

Base on the Figure 8, the optimal value seems to be  $k=10$ . Normally, the plot has the degrees of freedom on the x-axis. However, as  $K$  increases the degrees of freedom actually decreases, so the x-axis is reversed. Further, there is a little bump in figure 8, because when going from 1 to 2 means that, for each, the algorithm choose only additional tree. What that tree adds to the average will always bias it away from the true value. Then, as you go from 2 to 5, you add more trees so the bias of the 2nd tree is cancelled out. The top most right tree at  $k=50$  gives the highest bias of -20.57. The sign of this number makes sense as the top right tree has the highest yield, and so trying to estimate it by averaging all the other trees will naturally involve a negative bias. It also make sense that the highest bias is in the top right as the underlying trend is quadratic. So, not only are the values to the top, right drastically higher than the bottom, left, they are also further from the average.

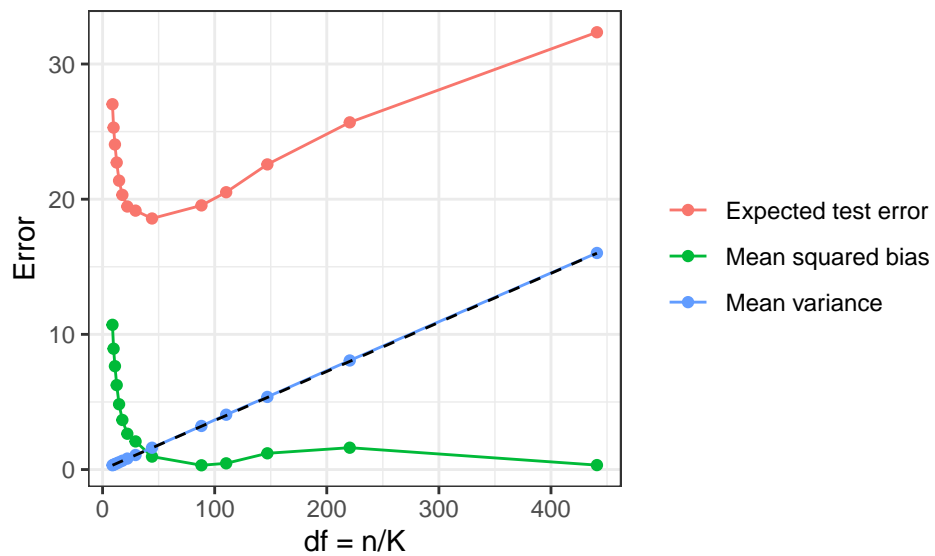


Figure 9: Bias-variance trade-off for KNN regression.

In figure 9, the variance is a straight line. My formula is  $\text{Var} = \sigma^2/k$  (where  $\sigma^2$  is 16).