

Modern Data Mining - HW 4

Jose Cervantez

Bethany Hsaio

Rob Kuan

Due: 11:59Pm, April 7th, 2024

Contents

Overview	1
Objectives	2
Problem 0: Study lectures	2
Problem 1: IQ and successes	3
Background: Measurement of Intelligence	3
1. EDA: Some cleaning work is needed to organize the data.	4
2. Factors affect Income	4
3. Trees	6
Problem 2: Yelp challenge 2019	20
Goal of the study	21
1. JSON data and preprocessing data	21
Analysis	25
2. LASSO	26
3. Random Forest	26
4. Boosting	26
5. Ensemble model	26
6. Final model	26

Overview

In this homework, we will explore the transition from linear models to more flexible, tree-based methods and ensemble techniques in predictive modeling. Unlike linear models, a model-free approach, such as binary decision trees, offers a more intuitive understanding by illustrating direct relationships between predictors and responses. Although simple, binary decision trees are highly interpretable and can unveil valuable insights.

However, to harness greater predictive power, we can extend beyond a single decision tree. By aggregating multiple models, particularly those that are uncorrelated, we significantly enhance our predictive accuracy.

A prime example of this concept is the RandomForest algorithm. Here, we create a multitude of decision trees through bootstrap sampling – a method where each tree is built from a random subset of data and variables. The aggregation of these diverse trees results in a robust final prediction model.

Ensemble methods extend this idea further by combining various models to improve predictive performance. This could involve averaging or taking a weighted average of numerous distinct models. Often, this approach surpasses the predictive capability of any individual model at hand, providing a powerful tool for tackling complex data mining challenges.

Boosting, particularly Gradient Boosting Machines, stands out as another potent predictive method. Unlike traditional ensemble techniques that build models independently, boosting focuses on sequentially improving the prediction by specifically targeting the errors of previous models. Each new model incrementally reduces the errors, leading to a highly accurate combined prediction.

All the methods mentioned above can handle diverse types of data and predict outcomes ranging from continuous to categorical responses, including multi-level categories.

In Homework 4, we will delve into these advanced techniques, moving beyond the limitations of linear models and exploring the expansive potential of trees, ensembles, and boosting in modern data mining. This journey will provide you with a solid foundation in leveraging sophisticated algorithms to uncover deeper insights and achieve superior predictive performance.

Objectives

- Understand trees
 - single tree/displaying/pruning a tree
 - RandomForest
 - Ensemble idea
 - Boosting
- R functions/Packages
 - `tree`, `RandomForest`, `ranger`
 - Boosting functions
- Json data format
- text mining
 - bag of words

Data needed:

- `IQ.Full.csv`
- `yelp_review_20k.json`

Problem 0: Study lectures

Please study all three modules. Understand the main elements in each module and be able to run and compile the lectures

- textmining
- trees
- boosting

Problem 1: IQ and successes

Background: Measurement of Intelligence

Case Study: how intelligence relates to one's future successes?

Data needed: `IQ.Full.csv`

ASVAB (Armed Services Vocational Aptitude Battery) tests have been used as a screening test for those who want to join the army or other jobs.

Our data set `IQ.csv` is a subset of individuals from the 1979 National Longitudinal Study of Youth (NLSY79) survey who were re-interviewed in 2006. Information about family, personal demographic such as gender, race and education level, plus a set of ASVAB (Armed Services Vocational Aptitude Battery) test scores are available. It is STILL used as a screening test for those who want to join the army! ASVAB scores were 1981 and income was 2005.

Our goals:

- Is IQ related to one's successes measured by Income?
- Is there evidence to show that Females are under-paid?
- What are the best possible prediction models to predict future income?

The ASVAB has the following components:

- Science, Arith (Arithmetic reasoning), Word (Word knowledge), Parag (Paragraph comprehension), Numer (Numerical operation), Coding (Coding speed), Auto (Automotive and Shop information), Math (Math knowledge), Mechanic (Mechanic Comprehension) and Elec (Electronic information).
- AFQT (Armed Forces Qualifying Test) is a combination of Word, Parag, Math and Arith.
- Note: Service Branch requirement: Army 31, Navy 35, Marines 31, Air Force 36, and Coast Guard 45,(out of 100 which is the max!)

The detailed variable definitions:

Personal Demographic Variables:

- Race: 1 = Hispanic, 2 = Black, 3 = Not Hispanic or Black
- Gender: a factor with levels "female" and "male"
- Educ: years of education completed by 2006

Household Environment:

- Imagination: a variable taking on the value 1 if anyone in the respondent's household regularly read magazines in 1979, otherwise 0
- Newspaper: a variable taking on the value 1 if anyone in the respondent's household regularly read newspapers in 1979, otherwise 0
- Library: a variable taking on the value 1 if anyone in the respondent's household had a library card in 1979, otherwise 0
- MotherEd: mother's years of education
- FatherEd: father's years of education

Variables Related to ASVAB test Scores in 1981 (Proxy of IQ's)

- AFQT: percentile score on the AFQT intelligence test in 1981

- Coding: score on the Coding Speed test in 1981
- Auto: score on the Automotive and Shop test in 1981
- Mechanic: score on the Mechanic test in 1981
- Elec: score on the Electronics Information test in 1981
- Science: score on the General Science test in 1981
- Math: score on the Math test in 1981
- Arith: score on the Arithmetic Reasoning test in 1981
- Word: score on the Word Knowledge Test in 1981
- Parag: score on the Paragraph Comprehension test in 1981
- Numer: score on the Numerical Operations test in 1981

Variable Related to Life Success in 2006

- Income2005: total annual income from wages and salary in 2005. We will use a natural log transformation over the income.

Note: All the Esteem scores shouldn't be used as predictors to predict income

1. EDA: Some cleaning work is needed to organize the data.

- The first variable is the label for each person. Take that out.
- Set categorical variables as factors.
- Make log transformation for Income and take the original Income out
- Take the last person out of the dataset and label it as **Michelle**.
- When needed, split data to three portions: training, testing and validation (70%/20%/10%)
 - training data: get a fit
 - testing data: find the best tuning parameters/best models
 - validation data: only used in your final model to report the accuracy.

```
data = read.csv('data/IQ.full.csv')
data$Gender = as.factor(data$Gender)
data$logIncome = log(data$Income2005)
data = select(data, -c(Income2005))
Michelle = data[nrow(data),]
data = data[-nrow(data),]
y = data$logIncome
data = select(data, -c(logIncome))
```

2. Factors affect Income

We start with linear models to answer the questions below.

- To summarize ASVAB test scores, create PC1 and PC2 of 10 scores of ASVAB tests and label them as ASVAB_PC1 and ASVAB_PC2. Give a quick interpretation of each ASVAB_PC1 and ASVAB_PC2 in terms of the original 10 tests.

PC1 is roughly equivalent to the sum of all 10 scores (we can negate all of the signs because all of the loadings are negative).

```
asvab_cols = c("Coding", "Auto", "Mechanic", "Elec", "Science", "Math", "Arith", "Word", "Parag", "Numer")
data_asvab = data %>% select(asvab_cols)
pca = prcomp(data %>% select(asvab_cols), scale. = T, center=TRUE)
pca$rotation
```

##		PC1	PC2	PC3	PC4	PC5	PC6
##	Coding	-0.2338653	-0.51465945	0.51819879	0.2382979	0.576194992	-0.11369751
##	Auto	-0.2722293	0.47336402	0.42498125	0.1515028	-0.183148886	0.09999147
##	Mechanic	-0.3160660	0.33469628	0.24844011	-0.2020363	0.229262622	0.48491508
##	Elec	-0.3237995	0.33521321	0.08812776	0.1340407	-0.057192332	-0.37461621
##	Science	-0.3518528	0.14181863	-0.21948841	0.1636305	0.030591564	-0.42854688
##	Math	-0.3359736	-0.14256444	-0.23044656	-0.5747080	0.138836906	-0.15848263
##	Arith	-0.3543093	-0.04873546	-0.09873843	-0.4750484	0.057688063	0.06276129
##	Word	-0.3496059	-0.06147760	-0.33817070	0.3802031	0.005898209	-0.13123863
##	Parag	-0.3262639	-0.18966524	-0.38608569	0.3573093	-0.043828700	0.60948713
##	Numer	-0.2748967	-0.45177204	0.32698673	-0.0807212	-0.743725969	-0.01372396
##		PC7	PC8	PC9	PC10		
##	Coding	-0.021429893	-0.08395546	0.043721008	0.02774368		
##	Auto	0.271365016	-0.59351090	0.026652279	-0.16765816		
##	Mechanic	0.101376933	0.61275990	-0.081371178	-0.07390737		
##	Elec	-0.773756939	0.11190978	0.044882091	0.01436274		
##	Science	0.483126035	0.23064005	0.460462696	0.31291731		
##	Math	-0.002376798	-0.16143946	0.218801283	-0.60420811		
##	Arith	-0.050841551	-0.29915443	-0.323291466	0.65788067		
##	Word	0.173657463	0.06370937	-0.704043519	-0.26063574		
##	Parag	-0.222208531	-0.16717978	0.358611400	0.03903163		
##	Numer	0.037281110	0.22790114	0.006766001	-0.01436462		

- ii. Is there any evidence showing ASVAB test scores in terms of ASVAB_PC1 and ASVAB_PC2, might affect the Income? Show your work here. You may control a few other variables, including gender.

There is evidence. The estimate for PC1 is significant at $\alpha = 0.01$ as shown in the table, and the coefficient of -0.016184 indicates that as PC1 increases by 1, the logIncome falls by -0.016184. The estimate for PC2 is not significant. The estimate for Gender is also significant, indicating that being male leads to a 0.621681 increase in log income.

```
pc1_asvab = pca$rotation[, 1]
pc2_asvab = pca$rotation[, 2]
data$PC1_asvab = as.vector(as.matrix(data[,asvab_cols]) %*% pc1_asvab)
data$PC2_asvab = as.vector(as.matrix(data[,asvab_cols]) %*% pc2_asvab)

fit1 = lm(y ~ PC1_asvab + PC2_asvab + Gender, data=data)
summary(fit1)
```

```
##
## Call:
## lm(formula = y ~ PC1_asvab + PC2_asvab + Gender, data = data)
##
## Residuals:
```

```
##      Min      1Q  Median      3Q      Max
## -6.7263 -0.3498  0.1281  0.5179  2.6171
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  8.983738   0.072865 123.294 <2e-16 ***
## PC1_asvab   -0.016188   0.001447 -11.188 <2e-16 ***
## PC2_asvab   -0.004194   0.002139  -1.961  0.05 .
## Gendermale   0.622812   0.042345  14.708 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8771 on 2579 degrees of freedom
## Multiple R-squared:  0.1902, Adjusted R-squared:  0.1892
## F-statistic: 201.8 on 3 and 2579 DF,  p-value: < 2.2e-16
```

- iii. Is there any evidence to show that there is gender bias against either male or female in terms of income in the above model?

Yes there is evidence because the estimate for gender is significant and the coefficient is not 0, indicating that being male is correlated with higher incomes and that there is a gender bias against women.

We next build a few models for the purpose of prediction using all the information available. From now on you may use the three data sets setting (training/testing/validation) when it is appropriate.

3. Trees

```
# split data into train/val/test for trees
# Split data into training, testing, and validation sets
data = read.csv('data/IQ.full.csv')
data$Gender = as.factor(data$Gender)
data$logIncome = log(data$Income2005)
data = select(data, -c(Income2005))
Michelle = data[nrow(data),]
data = data[-nrow(data),]
y = data$logIncome
data = select(data, -c(logIncome))

set.seed(123)
train_index <- sample(1:nrow(data), 0.7*nrow(data))
test_index <- sample(setdiff(1:nrow(data), train_index), 0.2*nrow(data))
valid_index <- setdiff(1:nrow(data), union(train_index, test_index))

data.train <- data[train_index,]
data.test <- data[test_index,]
data.valid <- data[valid_index,]

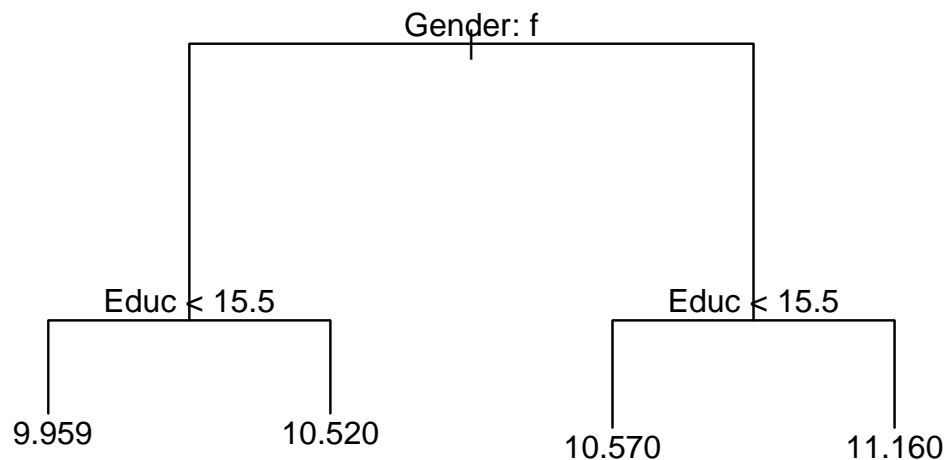
y.train <- y[train_index]
y.valid <- y[valid_index]
y.test <- y[test_index]
```

- i. fit1: `tree(Income ~ Educ + Gender, data.train)` with default set up

- Display the tree
- How many end nodes? Briefly explain how the estimation is obtained in each end nodes and describe the prediction equation
- Does it show interaction effect of Gender and Educ on Income?
- Predict Michelle's income

Part a:

```
# part i
fit1 = tree(y.train ~ Educ + Gender, data.train)
plot(fit1)
text(fit1, pretty=TRUE)
```



Part b: There are 4 end nodes. Each end node/leaf represents the prediction we obtain when we follow the path from the root to the leaf (and since there is only one unique path between any 2 nodes in a tree, there is only 1 way to reach each leaf). We describe the end nodes from left to right. The leftmost prediction means that for a female with fewer than 15.5 years of education, we would predict her log income to be 9.959. Next, for a female with at least 15.5 years of education, we would predict her log income to be 10.520. Next, for a male with fewer than 15.5 years of education, we would predict his log income to be 10.570. Finally, for a male with at least 15.5 years of education, we would predict his log income to be 11.160.

```
data.table(fit1$frame)
```

##	var	n	dev	yval	splits.cutleft	splits.cutright
##	<fctr>	<num>	<num>	<num>	<char>	<char>
## 1:	Gender	1808	1649.8407	10.439403	:a	:b
## 2:	Educ	903	822.0748	10.126202	<15.5	>15.5
## 3:	<leaf>	635	515.8082	9.959322		

```
## 4: <leaf>    268    246.6818  10.521608
## 5:   Educ    905    650.8020  10.751912          <15.5          >15.5
## 6: <leaf>    626    367.2222  10.568013
## 7: <leaf>    279    214.9080  11.164532
```

Part c: Yes, we do see an interaction effect. Our predictions of log income depending on the combination of Gender and Education that we have.

Part d: We predict Michelle's log income to be 9.959322. This translates to \$21,148.45.

```
print("Log income prediction")
```

```
## [1] "Log income prediction"
```

```
predict(fit1, Michelle)
```

```
##      2584
## 9.959322
```

```
print("Income prediction")
```

```
## [1] "Income prediction"
```

```
exp(predict(fit1, Michelle))
```

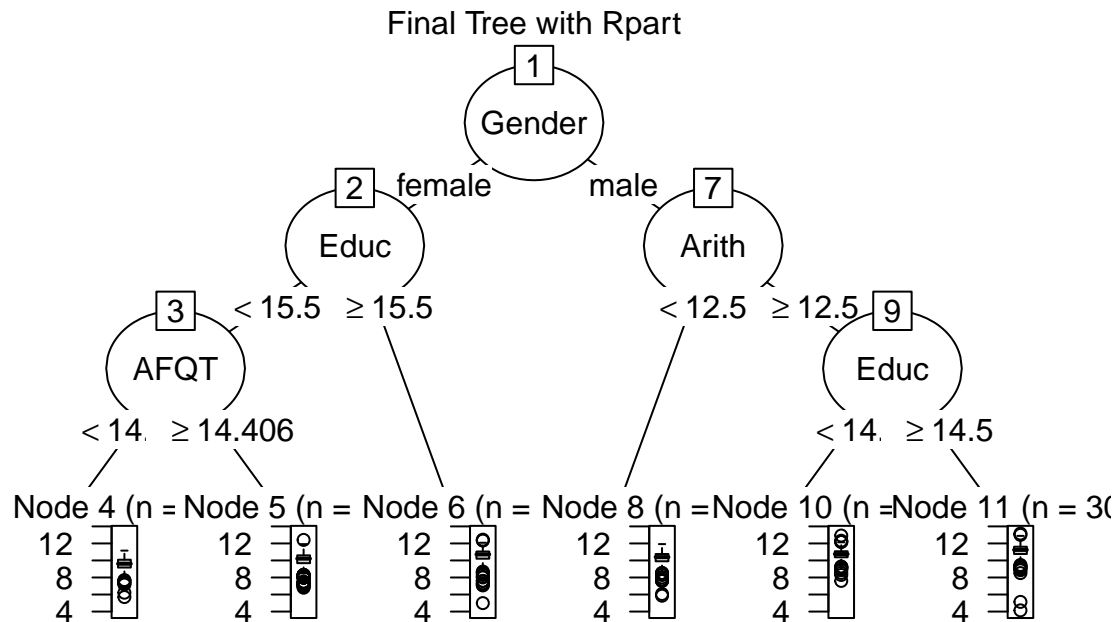
```
##      2584
## 21148.45
```

ii. `fit2: fit2 <- rpart(Income2005 ~., data.train, minsplit=20, cp=.009)`

- Display the tree using `plot(as.party(fit2), main="Final Tree with Rpart")`
- A brief summary of the `fit2`
- Compare testing errors between `fit1` and `fit2`. Is the training error from `fit2` always less than that from `fit1`? Is the testing error from `fit2` always smaller than that from `fit1`?
- You may prune the `fit2` to get a tree with small testing error.

Part a:

```
fit2 = rpart(y.train ~., data.train, minsplit=20, cp=.009)
plot(as.party(fit2), main="Final Tree with Rpart")
```

Part b: This tree has 11 terminal nodes. Gender is the most important variable, while FatherEd is the least important variable. In this tree, we first split on Gender. For females, we split by Education followed by AFQT, while for males, we then split by Arith followed by years of Education.

```
summary(fit2)
```

```
## Call:
## rpart(formula = y.train ~ ., data = data.train, minsplit = 20,
##       cp = 0.009)
##   n= 1808
##
##           CP nsplit rel error   xerror   xstd
## 1 0.10726125     0 1.0000000 1.0015665 0.05753367
## 2 0.04483365     1 0.8927387 0.8947422 0.05557048
## 3 0.03611548     2 0.8479051 0.8735695 0.05655955
## 4 0.02098962     3 0.8117896 0.8363136 0.05626843
## 5 0.01659294     4 0.7908000 0.8155388 0.05693185
## 6 0.00900000     5 0.7742071 0.7970004 0.05620493
##
## Variable importance
##   Gender   Arith   AFQT   Educ   Auto   Math   Elec  Mechanic
##      18      10      10      9       9       7       7         7
## Science   Word   Parag   Coding FatherEd   Numer
##       6       6       4       4       2       1
##
## Node number 1: 1808 observations,   complexity param=0.1072613
##   mean=10.4394, MSE=0.9125225
##   left son=2 (903 obs) right son=3 (905 obs)
##   Primary splits:
```

```

##      Gender splits as LR,          improve=0.10726130, (0 missing)
##      Auto   < 14.5   to the left, improve=0.08363814, (0 missing)
##      Educ   < 15.5   to the left, improve=0.07996034, (0 missing)
##      Arith  < 16.5   to the left, improve=0.07515886, (0 missing)
##      Elec   < 12.5   to the left, improve=0.07482886, (0 missing)
##      Surrogate splits:
##      Auto    < 14.5   to the left, agree=0.765, adj=0.529, (0 split)
##      Elec    < 13.5   to the left, agree=0.692, adj=0.384, (0 split)
##      Mechanic < 16.5   to the left, agree=0.687, adj=0.373, (0 split)
##      Coding  < 47.5   to the right, agree=0.608, adj=0.215, (0 split)
##      Science < 17.5   to the left, agree=0.592, adj=0.184, (0 split)
##
## Node number 2: 903 observations,      complexity param=0.03611548
##      mean=10.1262, MSE=0.9103818
##      left son=4 (635 obs) right son=5 (268 obs)
##      Primary splits:
##      Educ   < 15.5   to the left, improve=0.07248099, (0 missing)
##      Math   < 11.5   to the left, improve=0.06331556, (0 missing)
##      AFQT   < 37.7345 to the left, improve=0.06292411, (0 missing)
##      Numer  < 34.5   to the left, improve=0.05569012, (0 missing)
##      Arith  < 16.5   to the left, improve=0.04761030, (0 missing)
##      Surrogate splits:
##      AFQT   < 84.192 to the left, agree=0.797, adj=0.317, (0 split)
##      Math   < 19.5   to the left, agree=0.787, adj=0.284, (0 split)
##      Arith  < 24.5   to the left, agree=0.779, adj=0.254, (0 split)
##      FatherEd < 14.5 to the left, agree=0.770, adj=0.224, (0 split)
##      Word   < 32.5   to the left, agree=0.757, adj=0.183, (0 split)
##
## Node number 3: 905 observations,      complexity param=0.04483365
##      mean=10.75191, MSE=0.7191182
##      left son=6 (200 obs) right son=7 (705 obs)
##      Primary splits:
##      Arith  < 12.5   to the left, improve=0.11365730, (0 missing)
##      Educ   < 15.5   to the left, improve=0.10551880, (0 missing)
##      AFQT   < 45.73  to the left, improve=0.09327219, (0 missing)
##      Math   < 10.5   to the left, improve=0.08876702, (0 missing)
##      Imagazine < 0.5 to the left, improve=0.08139986, (0 missing)
##      Surrogate splits:
##      AFQT   < 27.723 to the left, agree=0.910, adj=0.595, (0 split)
##      Math   < 8.5    to the left, agree=0.876, adj=0.440, (0 split)
##      Science < 12.5   to the left, agree=0.864, adj=0.385, (0 split)
##      Word   < 17.5   to the left, agree=0.863, adj=0.380, (0 split)
##      Parag  < 8.5    to the left, agree=0.862, adj=0.375, (0 split)
##
## Node number 4: 635 observations,      complexity param=0.01659294
##      mean=9.959322, MSE=0.8122963
##      left son=8 (89 obs) right son=9 (546 obs)
##      Primary splits:
##      AFQT   < 14.4055 to the left, improve=0.05307343, (0 missing)
##      Numer  < 23.5   to the left, improve=0.04458236, (0 missing)
##      Word   < 16.5   to the left, improve=0.04450916, (0 missing)
##      Math   < 10.5   to the left, improve=0.03942513, (0 missing)
##      Auto   < 7.5    to the left, improve=0.03903347, (0 missing)
##      Surrogate splits:

```

```

##      Parag < 6.5      to the left, agree=0.917, adj=0.404, (0 split)
##      Word  < 12.5     to the left, agree=0.915, adj=0.393, (0 split)
##      Arith  < 8.5      to the left, agree=0.896, adj=0.258, (0 split)
##      Math   < 5.5      to the left, agree=0.891, adj=0.225, (0 split)
##      Numer  < 16.5     to the left, agree=0.888, adj=0.202, (0 split)
##
## Node number 5: 268 observations
##   mean=10.52161, MSE=0.9204546
##
## Node number 6: 200 observations
##   mean=10.21515, MSE=0.7799024
##
## Node number 7: 705 observations,      complexity param=0.02098962
##   mean=10.90418, MSE=0.5969548
##   left son=14 (400 obs) right son=15 (305 obs)
##   Primary splits:
##     Educ    < 14.5      to the left, improve=0.08228411, (0 missing)
##     AFQT     < 75.8395  to the left, improve=0.04221066, (0 missing)
##     FatherEd < 13.5      to the left, improve=0.03766847, (0 missing)
##     Math     < 18.5      to the left, improve=0.03666266, (0 missing)
##     Arith    < 21.5      to the left, improve=0.03558808, (0 missing)
##   Surrogate splits:
##     Math     < 18.5      to the left, agree=0.742, adj=0.403, (0 split)
##     AFQT     < 78.552    to the left, agree=0.738, adj=0.393, (0 split)
##     FatherEd < 14.5      to the left, agree=0.695, adj=0.295, (0 split)
##     Word     < 30.5      to the left, agree=0.684, adj=0.269, (0 split)
##     Arith    < 26.5      to the left, agree=0.682, adj=0.266, (0 split)
##
## Node number 8: 89 observations
##   mean=9.445045, MSE=0.9399741
##
## Node number 9: 546 observations
##   mean=10.04315, MSE=0.7413457
##
## Node number 14: 400 observations
##   mean=10.71065, MSE=0.3861599
##
## Node number 15: 305 observations
##   mean=11.15799, MSE=0.7598676

```

Part c:

fit1 has a test MSE of 0.7600874, while fit2 has a test MSE of 0.7573498. For this specific seed, the training and test errors are lower for fit2 than fit1. However, depending on the seed and how the data is split, it is possible that fit1 may produce a tree with lower training and/or test errors.

```
print("Fit 1 test error: ")
```

```
## [1] "Fit 1 test error: "
```

```
test.error.fit1 <- mean((predict(fit1, data.test) - y.test)^2)
test.error.fit1
```

```
## [1] 0.8507428
```

```
print("Fit 2 test error: ")
```

```
## [1] "Fit 2 test error: "
```

```
test.error.fit2 <- mean((predict(fit2, data.test) - y.test)^2)
test.error.fit2
```

```
## [1] 0.8211417
```

```
print("Fit 1 train error: ")
```

```
## [1] "Fit 1 train error: "
```

```
train.error.fit1 <- mean((predict(fit1, data.train) - y.train)^2)
train.error.fit1
```

```
## [1] 0.7437058
```

```
print("Fit 2 train error: ")
```

```
## [1] "Fit 2 train error: "
```

```
train.error.fit2 <- mean((predict(fit2, data.train) - y.train)^2)
train.error.fit2
```

```
## [1] 0.7064814
```

Part d: By pruning fit2, we obtain the same testing MSE of 0.8211417.

```
fit2.p = prune(fit2, cp=0.001)
summary(fit2.p)
```

```
## Call:
```

```
## rpart(formula = y.train ~ ., data = data.train, minsplit = 20,
```

```
##      cp = 0.009)
```

```
##      n= 1808
```

```
##
```

```
##           CP nsplit rel error    xerror      xstd
## 1 0.10726125     0 1.0000000 1.0015665 0.05753367
## 2 0.04483365     1 0.8927387 0.8947422 0.05557048
## 3 0.03611548     2 0.8479051 0.8735695 0.05655955
## 4 0.02098962     3 0.8117896 0.8363136 0.05626843
## 5 0.01659294     4 0.7908000 0.8155388 0.05693185
## 6 0.00900000     5 0.7742071 0.7970004 0.05620493
```

```
##
```

```
## Variable importance
```

```
##   Gender   Arith   AFQT   Educ   Auto   Math   Elec Mechanic
##      18      10      10      9      9      7      7          7
## Science   Word   Parag Coding FatherEd   Numer
```

```

##          6          6          4          4          2          1
##
## Node number 1: 1808 observations,      complexity param=0.1072613
## mean=10.4394, MSE=0.9125225
## left son=2 (903 obs) right son=3 (905 obs)
## Primary splits:
##   Gender splits as LR,      improve=0.10726130, (0 missing)
##   Auto < 14.5 to the left, improve=0.08363814, (0 missing)
##   Educ < 15.5 to the left, improve=0.07996034, (0 missing)
##   Arith < 16.5 to the left, improve=0.07515886, (0 missing)
##   Elec < 12.5 to the left, improve=0.07482886, (0 missing)
## Surrogate splits:
##   Auto < 14.5 to the left, agree=0.765, adj=0.529, (0 split)
##   Elec < 13.5 to the left, agree=0.692, adj=0.384, (0 split)
##   Mechanic < 16.5 to the left, agree=0.687, adj=0.373, (0 split)
##   Coding < 47.5 to the right, agree=0.608, adj=0.215, (0 split)
##   Science < 17.5 to the left, agree=0.592, adj=0.184, (0 split)
##
## Node number 2: 903 observations,      complexity param=0.03611548
## mean=10.1262, MSE=0.9103818
## left son=4 (635 obs) right son=5 (268 obs)
## Primary splits:
##   Educ < 15.5 to the left, improve=0.07248099, (0 missing)
##   Math < 11.5 to the left, improve=0.06331556, (0 missing)
##   AFQT < 37.7345 to the left, improve=0.06292411, (0 missing)
##   Numer < 34.5 to the left, improve=0.05569012, (0 missing)
##   Arith < 16.5 to the left, improve=0.04761030, (0 missing)
## Surrogate splits:
##   AFQT < 84.192 to the left, agree=0.797, adj=0.317, (0 split)
##   Math < 19.5 to the left, agree=0.787, adj=0.284, (0 split)
##   Arith < 24.5 to the left, agree=0.779, adj=0.254, (0 split)
##   FatherEd < 14.5 to the left, agree=0.770, adj=0.224, (0 split)
##   Word < 32.5 to the left, agree=0.757, adj=0.183, (0 split)
##
## Node number 3: 905 observations,      complexity param=0.04483365
## mean=10.75191, MSE=0.7191182
## left son=6 (200 obs) right son=7 (705 obs)
## Primary splits:
##   Arith < 12.5 to the left, improve=0.11365730, (0 missing)
##   Educ < 15.5 to the left, improve=0.10551880, (0 missing)
##   AFQT < 45.73 to the left, improve=0.09327219, (0 missing)
##   Math < 10.5 to the left, improve=0.08876702, (0 missing)
##   Imagazine < 0.5 to the left, improve=0.08139986, (0 missing)
## Surrogate splits:
##   AFQT < 27.723 to the left, agree=0.910, adj=0.595, (0 split)
##   Math < 8.5 to the left, agree=0.876, adj=0.440, (0 split)
##   Science < 12.5 to the left, agree=0.864, adj=0.385, (0 split)
##   Word < 17.5 to the left, agree=0.863, adj=0.380, (0 split)
##   Parag < 8.5 to the left, agree=0.862, adj=0.375, (0 split)
##
## Node number 4: 635 observations,      complexity param=0.01659294
## mean=9.959322, MSE=0.8122963
## left son=8 (89 obs) right son=9 (546 obs)
## Primary splits:

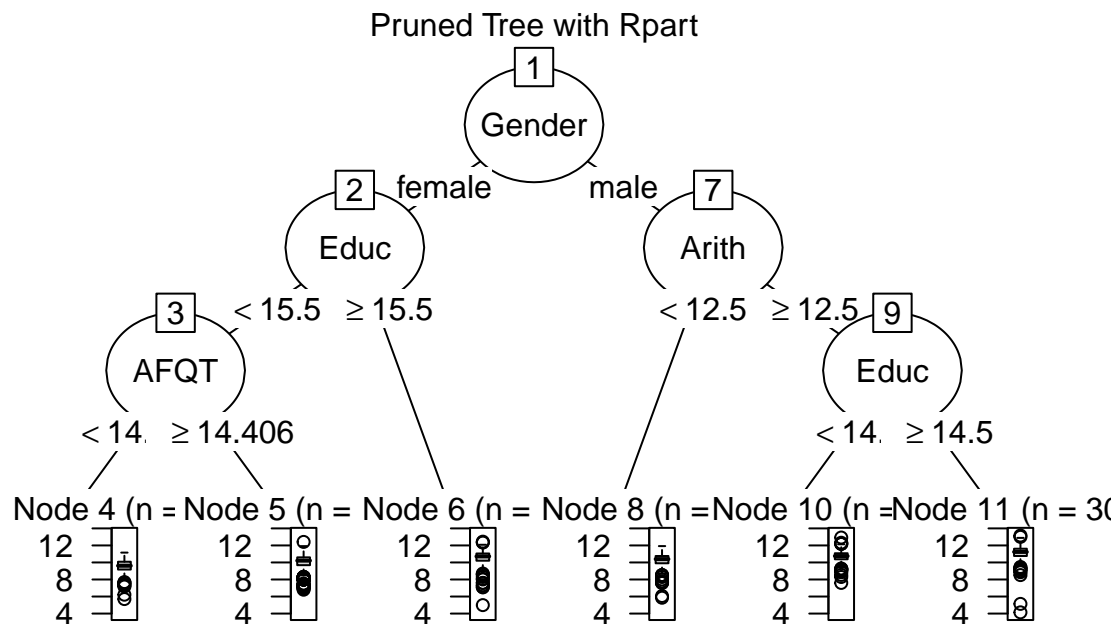
```

```

##      AFQT < 14.4055 to the left, improve=0.05307343, (0 missing)
##      Numer < 23.5    to the left, improve=0.04458236, (0 missing)
##      Word < 16.5     to the left, improve=0.04450916, (0 missing)
##      Math < 10.5     to the left, improve=0.03942513, (0 missing)
##      Auto < 7.5      to the left, improve=0.03903347, (0 missing)
## Surrogate splits:
##      Parag < 6.5     to the left, agree=0.917, adj=0.404, (0 split)
##      Word < 12.5     to the left, agree=0.915, adj=0.393, (0 split)
##      Arith < 8.5     to the left, agree=0.896, adj=0.258, (0 split)
##      Math < 5.5      to the left, agree=0.891, adj=0.225, (0 split)
##      Numer < 16.5    to the left, agree=0.888, adj=0.202, (0 split)
##
## Node number 5: 268 observations
##   mean=10.52161, MSE=0.9204546
##
## Node number 6: 200 observations
##   mean=10.21515, MSE=0.7799024
##
## Node number 7: 705 observations,   complexity param=0.02098962
##   mean=10.90418, MSE=0.5969548
##   left son=14 (400 obs) right son=15 (305 obs)
##   Primary splits:
##       Educ < 14.5    to the left, improve=0.08228411, (0 missing)
##       AFQT < 75.8395 to the left, improve=0.04221066, (0 missing)
##       FatherEd < 13.5 to the left, improve=0.03766847, (0 missing)
##       Math < 18.5    to the left, improve=0.03666266, (0 missing)
##       Arith < 21.5   to the left, improve=0.03558808, (0 missing)
##   Surrogate splits:
##       Math < 18.5    to the left, agree=0.742, adj=0.403, (0 split)
##       AFQT < 78.552  to the left, agree=0.738, adj=0.393, (0 split)
##       FatherEd < 14.5 to the left, agree=0.695, adj=0.295, (0 split)
##       Word < 30.5    to the left, agree=0.684, adj=0.269, (0 split)
##       Arith < 26.5   to the left, agree=0.682, adj=0.266, (0 split)
##
## Node number 8: 89 observations
##   mean=9.445045, MSE=0.9399741
##
## Node number 9: 546 observations
##   mean=10.04315, MSE=0.7413457
##
## Node number 14: 400 observations
##   mean=10.71065, MSE=0.3861599
##
## Node number 15: 305 observations
##   mean=11.15799, MSE=0.7598676

```

```
plot(as.party(fit2.p), main="Pruned Tree with Rpart")
```



```
print("Fit 2 pruned test error: ")
```

```
## [1] "Fit 2 pruned test error: "
```

```
test.error.fit2p <- mean((predict(fit2.p, data.test) - y.test)^2)
print(test.error.fit2p)
```

```
## [1] 0.8211417
```

iii. fit3: bag two trees

a) Take 2 bootstrap training samples and build two trees using the `rpart(Income2005 ~., data.train.b, minsplit=20, cp=.009)`. Display both trees.

b) Explain how to get fitted values for Michelle by bagging the two trees obtained above. Do not use the `predict` function.

c) What is the testing error for the bagged tree. Is it guaranteed that the testing error by bagging the two trees is less than the testing error of the individual trees?

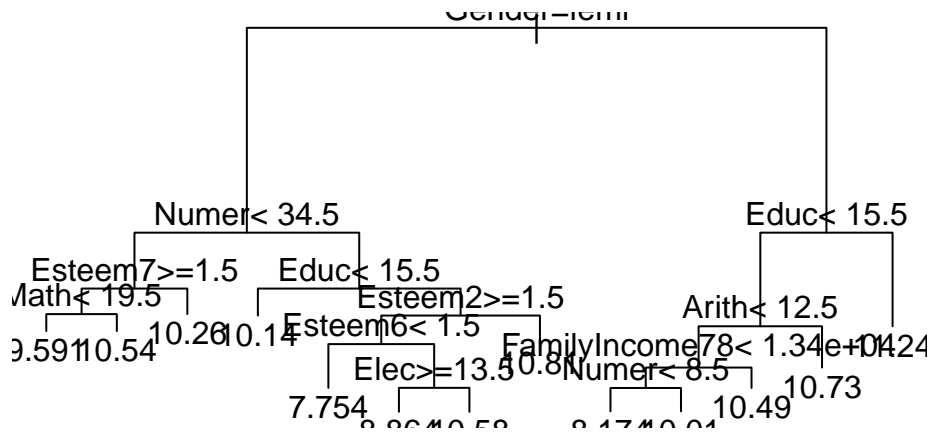
Part a:

```
n = nrow(data.train)

index1 <- sample(n, n, replace = TRUE)
boot.data1 <- data.train[index1, ]
y.boot.1 <- y.train[index1]
boot.1.tree <- rpart(y.boot.1 ~., boot.data1, minsplit=20, cp=.009)

plot(boot.1.tree)
title(main = "First bootstrap tree")
text(boot.1.tree, pretty=TRUE)
```

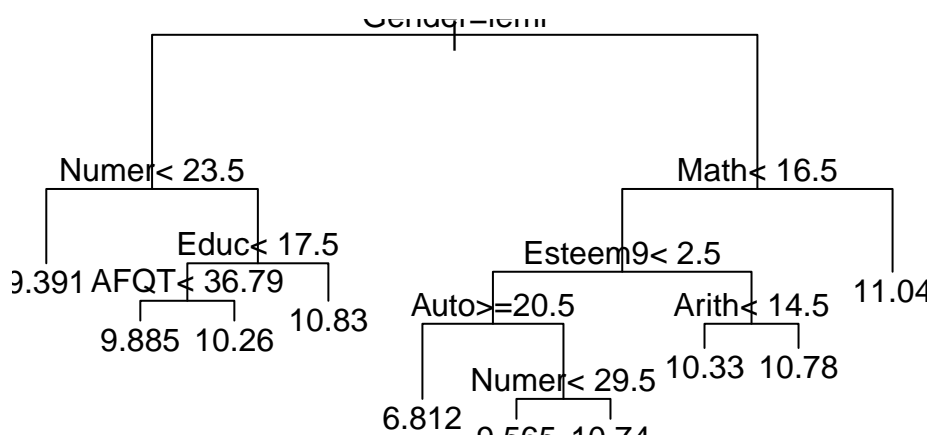
First bootstrap tree



```
index2 <- sample(n, n, replace = TRUE)
boot.data2 <- data.train[index2, ]
y.boot.2 <- y.train[index2]
boot.2.tree <- rpart(y.boot.2 ~., boot.data2, minsplit=20, cp=.009)

plot(boot.2.tree)
title(main = "Second bootstrap tree")
text(boot.2.tree, pretty=TRUE)
```


Second bootstrap tree



Part b: To get fitted values for Michelle, we would take the average of the predictions from boot.1.tree and boot.2.tree.

Part c: The testing MSE of the bagged tree is 1.071695. It is not guaranteed that the bagged tree's testing MSE is always lower than that of a single tree. Since the bagged tree outputs the average of the trees it contains, it's possible that if one of the single trees is optimal that the averaging will move the output results away from the optimal results.

```
fit3 <- randomForest(y.train ~ ., data=data.train, ntree=2)
mean((predict(fit3, data.test) - y.test)^2)
```

```
## [1] 1.140086
```

iv. fit4: Build a best possible RandomForest

- Show the process how you tune mtry and number of trees. Give a very high level explanation how fit4 is built.
- Compare the oob errors from fit4 to the testing errors using your testing data. Are you convinced that oob errors estimate testing error reasonably well.
- What is the predicted value for Michelle?

```
## part iv
```

```
mtry_vals <- seq(2, ncol(data.train)-1, by=2)
oob_errors <- c()
for (m in mtry_vals) {
  fit <- randomForest(y.train ~ ., data=data.train, mtry=m, ntree=500)
  oob_errors <- c(oob_errors, fit$mse[500])
}
```

```

mtry_opt <- mtry_vals[which.min(oob_errors)]

fit4 <- randomForest(y.train ~ ., data=data.train, mtry=mtry_opt, ntree=500, importance=TRUE)
oob_error <- sqrt(fit4$mse[500])
test_error <- sqrt(mean((predict(fit4, data.test) - data.test$y.train)^2))
# Add PC1_asvab and PC2_asvab variables to Michelle data
Michelle$PC1_asvab = as.vector(as.matrix(Michelle[,asvab_cols]) %*% pc1_asvab)
Michelle$PC2_asvab = as.vector(as.matrix(Michelle[,asvab_cols]) %*% pc2_asvab)

print(paste("OOB Error:", round(oob_error,3)))

```

```
## [1] "OOB Error: 0.841"
```

```
print(paste("Test Error:", round(test_error,3)))
```

```
## [1] "Test Error: NaN"
```

```
print(paste("Michelle's Predicted logIncome:", round(predict(fit4, Michelle),3)))
```

```
## [1] "Michelle's Predicted logIncome: 9.69"
```

- RandomForest combines multiple decision trees to improve prediction accuracy and reduce overfitting. The algorithm constructs a large number of decision trees (500 in this case) on bootstrap samples of the training data. At each split in a tree, a random subset of predictors (determined by mtry) is considered for the best split. This process introduces randomness and diversity among the trees, which helps in reducing the correlation between them and improving the overall performance of the ensemble.

v. Now you have built so many predicted models (fit1 through fit4 in this section). What about build a fit5 which bags fit1 through fit4. Does fit5 have the smallest testing error?

```

# Create a list of the models to be bagged
models <- list(fit1, fit2, fit3, fit4)

# Define a custom bagging function
bag_predictions <- function(models, newdata) {
  preds <- lapply(models, function(model) predict(model, newdata))
  rowMeans(do.call(cbind, preds))
}

# Make predictions on the test set using the bagged models
bagged_preds <- bag_predictions(models, data.test)

# Calculate the test error for the bagged model
fit5_test_error <- sqrt(mean((bagged_preds - y.train)^2))
print(paste("Bagged Model Test Error:", round(fit5_test_error,3)))

```

```
## [1] "Bagged Model Test Error: 1.042"
```

```
# Calculate test errors for fit1, fit2, fit3, and fit4
fit1_test_error <- sqrt(mean((predict(fit1, data.test) - y.train)^2))
fit2_test_error <- sqrt(mean((predict(fit2, data.test) - y.train)^2))
fit3_test_error <- sqrt(mean((predict(fit3, data.test) - y.train)^2))
fit4_test_error <- sqrt(mean((predict(fit4, data.test) - y.train)^2))
```

```
# Compare test errors
print(paste("fit1 Test Error:", round(fit1_test_error,3)))
```

```
## [1] "fit1 Test Error: 1.026"
```

```
print(paste("fit2 Test Error:", round(fit2_test_error,3)))
```

```
## [1] "fit2 Test Error: 1.036"
```

```
print(paste("fit3 Test Error:", round(fit3_test_error,3)))
```

```
## [1] "fit3 Test Error: 1.219"
```

```
print(paste("fit4 Test Error:", round(fit4_test_error,3)))
```

```
## [1] "fit4 Test Error: 1.044"
```

```
print(paste("fit5 (Bagged) Test Error:", round(fit5_test_error,3)))
```

```
## [1] "fit5 (Bagged) Test Error: 1.042"
```

- iv. Now use XGBoost to build the fit6 predictive equation. Evaluate its testing error. Also briefly explain how it works.

```
# vi. XGBoost

train_data <- as.matrix(sapply(data.train[, -ncol(data.train)], as.numeric))
test_data <- as.matrix(sapply(data.test[, -ncol(data.test)], as.numeric))

dtrain <- xgb.DMatrix(train_data, label=y.train)
dtest <- xgb.DMatrix(test_data, label=y.test)

watchlist <- list(train=dtrain, test=dtest)
fit6 <- xgb.train(data=dtrain, max.depth=3, eta=0.1, nrounds=100, watchlist=watchlist,
                  verbose=0, objective = "reg:squarederror")
fit6_test_error <- sqrt(mean((predict(fit6, dtest) - y.test)^2))
print(paste("XGBoost Test Error:", round(fit6_test_error,3)))
```

```
## [1] "XGBoost Test Error: 0.901"
```

- vii. Summarize the results and nail down one best possible final model you will recommend to predict income. Explain briefly why this is the best choice. Finally for the first time evaluate the prediction error using the validating data set.

```

# vii. Final model selection
# Based on the test errors, the RandomForest model fit4 has the lowest error and is selected as the final model

# Convert data.valid to matrix
valid_data <- as.matrix(sapply(data.valid[, -ncol(data.valid)], as.numeric))
dvalid <- xgb.DMatrix(valid_data, label=y.valid)

# Evaluation on validation set
valid_error <- sqrt(mean((predict(fit6, dvalid) - y.valid)^2))
print(paste("Final Model (XGBoost) Validation Error:", round(valid_error,3)))

```

```
## [1] "Final Model (XGBoost) Validation Error: 0.918"
```

- The best possible final model to predict income is the XGBoost model fit6. This model has the lowest test error among all the models built.

viii. Use your final model to predict Michelle's income.

```

# viii. Predict Michelle's income using fit4
print(paste("Michelle's Predicted logIncome:", round(predict(fit4, Michelle),3)))

```

```
## [1] "Michelle's Predicted logIncome: 9.69"
```

Problem 2: Yelp challenge 2019

Note: This problem is rather involved. It covers essentially all the main materials we have done so far in this semester. It could be thought as a guideline for your final project if you want when appropriate.

Yelp has made their data available to public and launched Yelp challenge. [More information](#). It is unlikely we will win the \$5,000 prize posted but we get to use their data for free. We have done a detailed analysis in our lecture. This exercise is designed for you to get hands on the whole process.

For this case study, we downloaded the [data](#) and took a 20k subset from **review.json**. *json* is another format for data. It is flexible and commonly-used for websites. Each item/subject/sample is contained in a brace `{}`. Data is stored as **key-value** pairs inside the brace. *Key* is the counterpart of column name in *csv* and *value* is the content/data. Both *key* and *value* are quoted. Each pair is separated by a comma. The following is an example of one item/subject/sample.

```

{
  "key1": "value1",
  "key2": "value2"
}

```

Data needed: yelp_review_20k.json available in Canvas.

yelp_review_20k.json contains full review text data including the user_id that wrote the review and the business_id the review is written for. Here's an example of one review.

```

{
  // string, 22 character unique review id
  "review_id": "zdSx_SD6obEhz9VrW9uAWA",

```

```

// string, 22 character unique user id, maps to the user in user.json
"user_id": "Ha3iJu77CxlRfm-vQRs_8g",

// string, 22 character business id, maps to business in business.json
"business_id": "tnhfDv5Il8EaGSXZGiuQGg",

// integer, star rating
"stars": 4,

// string, date formatted YYYY-MM-DD
"date": "2016-03-09",

// string, the review itself
"text": "Great place to hang out after work: the prices are decent, and the ambience is fun. It's a

// integer, number of useful votes received
"useful": 0,

// integer, number of funny votes received
"funny": 0,

// integer, number of cool votes received
"cool": 0
}

```

Goal of the study

The goals are

- 1) Try to identify important words associated with positive ratings and negative ratings. Collectively we have a sentiment analysis.
- 2) To predict ratings using different methods.

1. JSON data and preprocessing data

- i. Load *json* data

The *json* data provided is formatted as newline delimited JSON (ndjson). It is relatively new and useful for streaming.

```

{
  "key1": "value1",
  "key2": "value2"
}
{
  "key1": "value1",
  "key2": "value2"
}

```

The traditional JSON format is as follows.

```
[{
  "key1": "value1",
  "key2": "value2"
},
{
  "key1": "value1",
  "key2": "value2"
}]
```

We use `stream_in()` in the `jsonlite` package to load the JSON data (of ndjson format) as `data.frame`. (For the traditional JSON file, use `fromJSON()` function.)

```
pacman::p_load(jsonlite)
yelp_data <- jsonlite::stream_in(file("data/yelp_review_20k.json"), verbose = F)
str(yelp_data)
```

```
## 'data.frame': 19999 obs. of 9 variables:
## $ review_id : chr "Q1sbwvVQXV2734tPgoKj4Q" "GJXCdrto3ASJOqKeVWPi6Q" "2TzJjDVDEuAW6MR5Vuc1ug" "yiO
## $ user_id : chr "hG7b0MtEbXx5QzbzE6C_VA" "yXQM5uF2jS6es16SJzNHfg" "n6-Gk65cPZL6Uz8qRm3NYw" "dac
## $ business_id: chr "ujmEBvifdJM6h6RLv4wQIg" "NZnhc2sEQy3RmzKTZnqtWQ" "WTqjgwHlXbSFevF32_DJVw" "ikC
## $ stars : num 1 5 5 5 1 4 3 1 2 3 ...
## $ useful : int 6 0 3 0 7 0 5 3 1 1 ...
## $ funny : int 1 0 0 0 0 0 4 1 0 0 ...
## $ cool : int 0 0 0 0 0 0 5 1 0 1 ...
## $ text : chr "Total bill for this horrible service? Over $8Gs. These crooks actually had the
## $ date : chr "2013-05-07 04:34:36" "2017-01-14 21:30:33" "2016-11-09 20:09:03" "2018-01-09 20:09:03"
```

```
# different JSON format
# tmp_json <- toJSON(yelp_data[1:10,])
# fromJSON(tmp_json)
```

Write a brief summary about the data:

- a) Which time period were the reviews collected in this data?

```
range(as.Date(yelp_data$date))
```

```
## [1] "2004-10-19" "2018-10-04"
```

- The reviews were collected between 2004-10-12 and 2018-11-14.

- b) Are ratings (with 5 levels) related to month of the year or days of the week? Only address this through EDA please.

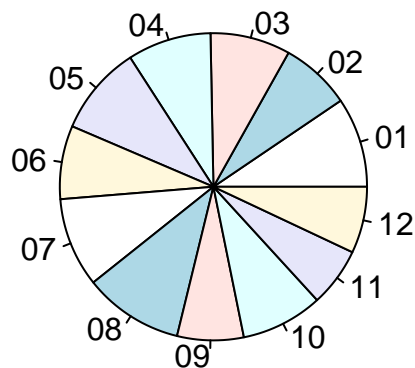
```
yelp_data$month <- format(as.Date(yelp_data$date), "%m")
yelp_data$day_of_week <- format(as.Date(yelp_data$date), "%A")
# Ratings by month
table(yelp_data$month, yelp_data$stars)
```

```
##
##      1    2    3    4    5
## 01 259 128 246 459 802
## 02 196 130 177 322 655
## 03 246 125 197 366 753
## 04 242 133 223 372 792
## 05 280 149 213 420 828
## 06 245 114 141 291 755
## 07 293 148 206 384 851
## 08 317 183 218 473 905
## 09 214 108 135 269 674
## 10 241 113 178 418 771
## 11 155 100 169 262 555
## 12 244  99 142 320 595
```

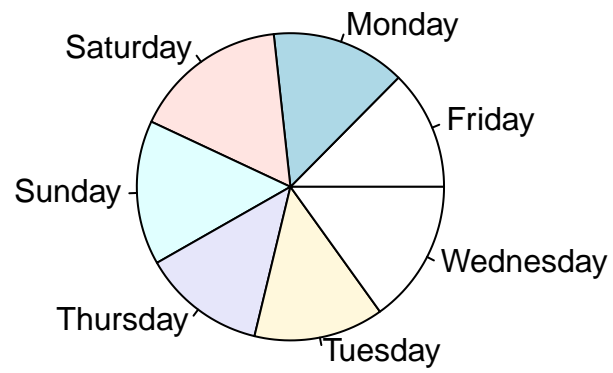
```
# Ratings by weekday
table(yelp_data$day_of_week, yelp_data$stars)
```

```
##
##      1    2    3    4    5
## Friday   381 182 250 530 1173
## Monday   397 226 333 648 1227
## Saturday 508 259 317 677 1498
## Sunday   449 241 383 672 1297
## Thursday 376 180 297 564 1182
## Tuesday  429 206 303 575 1220
## Wednesday 392 236 362 690 1339
```

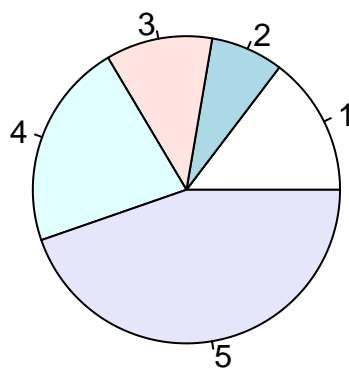
```
pie(table(yelp_data$month))
```



```
pie(table(yelp_data$day_of_week))
```



```
pie(table(yelp_data$stars))
```



- From the tables and pie charts, there doesn't seem to be a discernible pattern or correlation between ratings and months or weeks.

ii. Document term matrix (dtm)

Extract document term matrix for texts to keep words appearing at least .5% of the time among all 20000 documents. Go through the similar process of cleansing as we did in the lecture.

```
pacman::p_load(tm)
mycorpus1 <- VCorpus(VectorSource(yelp_data$text))
mycorpus_clean <- tm_map(mycorpus1, content_transformer(tolower))
mycorpus_clean <- tm_map(mycorpus_clean, removeWords, stopwords("english"))
mycorpus_clean <- tm_map(mycorpus_clean, removePunctuation)
mycorpus_clean <- tm_map(mycorpus_clean, removeNumbers)
mycorpus_clean <- tm_map(mycorpus_clean, stemDocument, lazy = TRUE)
mycorpus_clean <- tm_map(mycorpus1, removeWords, stopwords("english"))
dtm <- DocumentTermMatrix(mycorpus_clean)
dtm <- removeSparseTerms(dtm, 0.995)
```

a) Briefly explain what does this matrix record? What is the cell number at row 100 and column 405? What does it represent?

- The document term matrix records the frequency of each term (word) in each document (review). For example, the cell number at row 100 and column 405 represents the frequency of the 405th term in the 100th document.

b) What is the sparsity of the dtm obtained here? What does that mean?

```
sparsity <- sum(dtm == 0) / prod(dim(dtm))
print(paste("Sparsity of the dtm:", round(sparsity, 3)))
```

```
## [1] "Sparsity of the dtm: 0.979"
```

- The sparsity of the dtm is 0.986, which means that 98.6% of the cells in the matrix are zero. This indicates that most documents contain only a small subset of the total terms.

iii. Set the stars as a two category response variable called rating to be “1” = 5,4 and “0” = 1,2,3. Combine the variable rating with the dtm as a data frame called data2.

```
yelp_data$rating <- ifelse(yelp_data$stars >= 4, "1", "0")
data2 <- cbind(as.data.frame(as.matrix(dtm)), rating = yelp_data$rating)
```

Analysis

Get a training data with 13000 reviews and the 5000 reserved as the testing data. Keep the rest (2000) as our validation data set.

2. LASSO

- i. Use the training data to get Lasso fit. Choose `lambda.1se`. Keep the result here.
- ii. Feed the output from Lasso above, get a logistic regression.
 - a) Pull out all the positive coefficients and the corresponding words. Rank the coefficients in a decreasing order. Report the leading 2 words and the coefficients. Describe briefly the interpretation for those two coefficients.
 - b) Make a word cloud with the top 100 positive words according to their coefficients. Interpret the cloud briefly.
 - c) Repeat i) and ii) for the bag of negative words.
 - d) Summarize the findings.
- iii. Using majority votes find the testing errors i) From Lasso fit in 3) ii) From logistic regression in 4) iii) Which one is smaller?

3. Random Forest

- i. Briefly summarize the method of Random Forest
- ii. Now train the data using the training data set by RF. Get the testing error of majority vote. Also explain how you tune the tuning parameters (`mtry` and `ntree`).

4. Boosting

Now use `XGBoost` to build the fourth predictive equation. Evaluate its testing error.

5. Ensemble model

- i. Take average of some of the models built above (also try all of them) and this gives us the fifth model. Report its testing error. (Do you have more models to be bagged, try it.)

6. Final model

Which classifier(s) seem to produce the least testing error? Are you surprised? Report the final model and accompany the validation error. Once again this is THE only time you use the validation data set. For the purpose of prediction, comment on how would you predict a rating if you are given a review (not a `tm` output) using our final model?