

# Stats 315B: Homework 1

*Joe Higgins, Austin Wang, Jessica Wetstone*

*Due 5/20/2018*

## Question 1

Random forests predict with an ensemble of bagged trees each trained on a bootstrap sample randomly drawn from the original training data. Additional random variation among the trees is induced by choosing the variable for each split from a small randomly chosen subset of all of the predictor variables when building each tree. What are the advantages and disadvantages of this random variable selection strategy? How can one introduce additional tree variation in the forest without randomly selecting subsets of variables?

The advantages of choosing a random subset of predictor variables are:

- Because we are no longer doing an exhaustive search over all predictor variables, each split requires less processing power during the training process. Given modern computer architecture, this advantage matters less since each tree in the ensemble is independent and can be computed simultaneously in a different parallel process.
- The random variation introduced by choosing the variable for each split from a randomly chosen subset de-correlates the resulting trees in the forest, forcing them to be different from each other. This has the effect of decreasing the variance of the ensemble's estimate, since the variance of the mean of a collection random variables is directly related to the correlation of those variables. Decreasing the correlation of the trees therefore decreases the variance of the forest.

The disadvantage of this approach is that it can result in trees with very weak predictive power. Suppose that none of the best predictors were available for a given tree's first several splits – this would result in a highly sub-optimal tree.

Another way to inject randomness into the process and de-correlate the trees is to decrease the size of the bootstrapped sample of the training data. Each individual tree is now more likely to overfit the smaller sample, increasing tree variation within the forest.

## Question 2

**Why is it necessary to use regularization in linear regression when the number of predictor variables is greater than the number of observations in the training sample? Explain how regularization helps in this case. Are there other situations where regularization might help? What is the potential disadvantage of introducing regularization? Why is sparsity a reasonable assumption in the boosting context. Is it always? If not, why not?**

If the number of predictor variables  $n$  is greater than the number of observations in the training sample  $N$ , then there is an infinite number of solutions that will satisfy the linear regression problem – and there will be very high variance in the solution derived from that training sample. Regularization helps in this situation because it decreases the size of the function class, reducing variance. Regularization will help in any situation where we have a reason to favor sparse solutions by both encouraging sparsity and preventing high variance.

A potential disadvantage of regularization is that it could increase bias - especially if you have chosen a penalty that doesn't align with the best answer for the problem.

In the boosting context, each predictor is a possible tree that could be built on the training data from the set of all possible trees. It is reasonable to assume that most trees - out of the set of all possible trees (All combinations of variables and all possible splits) - will do very poorly on the overall prediction problem. So we would want most of our coefficients in the linear regression to be 0.

Sparsity is not always a necessary assumption in linear regression, especially if your feature space is small. For example, if we were predicting house prices from zip code and square footage, we would expect both variables to be important.

### Question 3

Show that the convex members of the power family of penalties, except for the lasso, have the property that solutions to  $\hat{a}(\lambda) = \operatorname{argmin}_a \hat{R}(a) + \lambda P_\gamma(a)$  have nonzero values for all coefficients at each path point indexed by  $\lambda$ . By contrast the convex members of the elastic net (except ridge) can produce solutions with many zero valued coefficients at various path points.

### Question 4

Show that the variable  $x_{j^*}$  that has the maximum absolute correlation with  $j^* = \operatorname{argmax}_{1 \leq j \leq J} |E(yx_j)|$  is the same as the one that best predicts  $y$  using squared—error loss  $j^* = \operatorname{argmin}_{1 \leq j \leq J} \min_\rho E[y - \rho x_j]^2$ . This shows that the base learner most correlated with the generalized residual is the one that best predicts it with squared—error loss.

Let  $f(\rho, x_j, y) = E[y - \rho x_j]^2$ . We first find  $\min_\rho f(\rho, x_j, y)$  by taking the partial derivative of  $f$  with respect to  $\rho$  and setting it equal to 0:

$$\begin{aligned} f(\rho, x_j, y) &= E[y - \rho x_j]^2 \\ &= E[y^2 - 2\rho yx_j + \rho^2 x_j^2] \\ &= E[y^2] - 2\rho E[yx_j] + \rho^2 E[x_j^2] \\ &= E[y^2] - 2\rho E[yx_j] + \rho^2 \quad \text{since } E[x_j^2] = 1 \end{aligned}$$

$$\begin{aligned} \implies \frac{\partial f}{\partial \rho} &= -2E[yx_j] + 2\rho \\ \implies 0 &= -2E[yx_j] + 2\rho^* \\ \implies \rho^* &= E[yx_j] \end{aligned}$$

$$\begin{aligned} \implies \min_\rho f(\rho, x_j, y) &= f(\rho^*, x_j, y) \\ &= E[y - \rho^* x_j]^2 \\ &= E[y^2] - 2\rho^* E[yx_j] + \rho^{*2} \\ &= E[y^2] - 2E[yx_j]E[yx_j] + (E[yx_j])^2 \\ &= E[y^2] - (E[yx_j])^2 \end{aligned}$$

Note that this is a unique minimum value because  $f$  is strictly convex in  $\rho$ . We know this because  $\frac{\partial^2 f}{\partial \rho^2} > 0$ :

$$\begin{aligned}
\frac{\partial f}{\partial \rho} &= -2E[yx_j] + 2\rho \\
\Rightarrow \frac{\partial^2 f}{\partial \rho^2} &= 2 \\
&> 0
\end{aligned}$$

Then:

$$\begin{aligned}
\operatorname{argmin}_{1 \leq j \leq J} \min_{\rho} E[y - \rho x_j]^2 &= \operatorname{argmin}_{1 \leq j \leq J} (E[y^2] - (E[yx_j])^2) \\
&= \operatorname{argmin}_{1 \leq j \leq J} -(E[yx_j])^2 \\
&= \operatorname{argmax}_{1 \leq j \leq J} (E[yx_j])^2 \\
&= \operatorname{argmax}_{1 \leq j \leq J} |E[yx_j]|^2 \\
&= \operatorname{argmax}_{1 \leq j \leq J} |E[yx_j]| \quad \text{since } f(x) = \sqrt{x} \text{ is monotonically increasing for } x > 0
\end{aligned}$$

Therefore,  $\operatorname{argmax}_{1 \leq j \leq J} |E[yx_j]| = \operatorname{argmin}_{1 \leq j \leq J} \min_{\rho} E[y - \rho x_j]^2$ , so the base learner most correlated with the generalized residual is the one that best predicts it with squared—error loss. ■

## Question 5

Let  $\mathbf{z}_l = \{z_1, \dots, z_l\}$  be a subset of the predictor variables  $\mathbf{x} = \{x_1, \dots, x_n\}$  and  $\mathbf{z}_{\setminus l}$  the complement subset, i.e.  $\mathbf{z}_l \cup \mathbf{z}_{\setminus l} = \mathbf{x}$ . Show that if a function  $F(\mathbf{x})$  is additive in  $\mathbf{z}_l$  and  $\mathbf{z}_{\setminus l}$ , i.e.

$$F(\mathbf{x}) = F_l(\mathbf{z}_l) + F_{\setminus l}(\mathbf{z}_{\setminus l})$$

then the partial dependence of  $F(\mathbf{x})$  on  $\mathbf{z}_l$  is  $F_l(\mathbf{z}_l)$  up to an additive constant. This is the dependence of  $F(\mathbf{x})$  on  $\mathbf{z}_l$  accounting for the effect of the other variables  $\mathbf{z}_{\setminus l}$ . Show that this need not be the case for  $E[F(\mathbf{x}) \mid \mathbf{z}_l]$  which is the dependence of  $F(\mathbf{x})$  on  $\mathbf{z}_l$  ignoring the other variables  $\mathbf{z}_{\setminus l}$ . Under what conditions would the two be the same?

## Question 6

Binary classification: Spam Email. The data set for this problem is `spam_stats315B.csv`, with documentation files `spam_stats315B_info.txt` and `spam_stats315B_names.txt`. The data set is a collection of 4601 emails of which 1813 were considered spam, i.e. unsolicited commercial email. The data set consists of 58 attributes of which 57 are continuous predictors and one is a class label that indicates whether the email was considered spam (1) or not (0). Among the 57 predictor attributes are: percentage of the word “free” in the email, percentage of exclamation marks in the email, etc. See file `spam_stats315B_names.txt` for the full list of attributes. The goal is, of course, to predict whether or not an email is “spam”. This data set is used for illustration in the tutorial *Boosting with R Programming*. The data set `spam_stats315B_train.csv` represents a subsample of these emails randomly selected from `spam_stats315B.csv` to be used for training. The file `spam_stats315B_test.csv` contains the remaining emails to be used for evaluating results.

(a) Based on the training data, fit a gbm model for predicting whether or not an email is “spam”, following the example in the tutorial. What is your estimate of the misclassification rate? Of all the spam emails of the test set what percentage was misclassified, and of all the non-spam emails in the test set what percentage was misclassified?

```
rm(list = ls())
data_path <- paste(getwd(), '/data', sep='')
setwd(data_path)

#data labels
rflabs<-c("make", "address", "all", "3d", "our", "over", "remove",
  "internet", "order", "mail", "receive", "will",
  "people", "report", "addresses", "free", "business",
  "email", "you", "credit", "your", "font", "000", "money",
  "hp", "hpl", "george", "650", "lab", "labs",
  "telnet", "857", "data", "415", "85", "technology", "1999",
  "parts", "pm", "direct", "cs", "meeting", "original", "project",
  "re", "edu", "table", "conference", ";", "(", "[", "!", "$", "#",
  "CAPAVE", "CAPMAX", "CAPTOT", "type")

#load the data
train <- read.csv(file="spam_stats315B_train.csv", header=FALSE, sep=",")
test <- read.csv(file="spam_stats315B_test.csv", header=FALSE, sep=",")
colnames(train)<-rflabs
colnames(test)<-rflabs

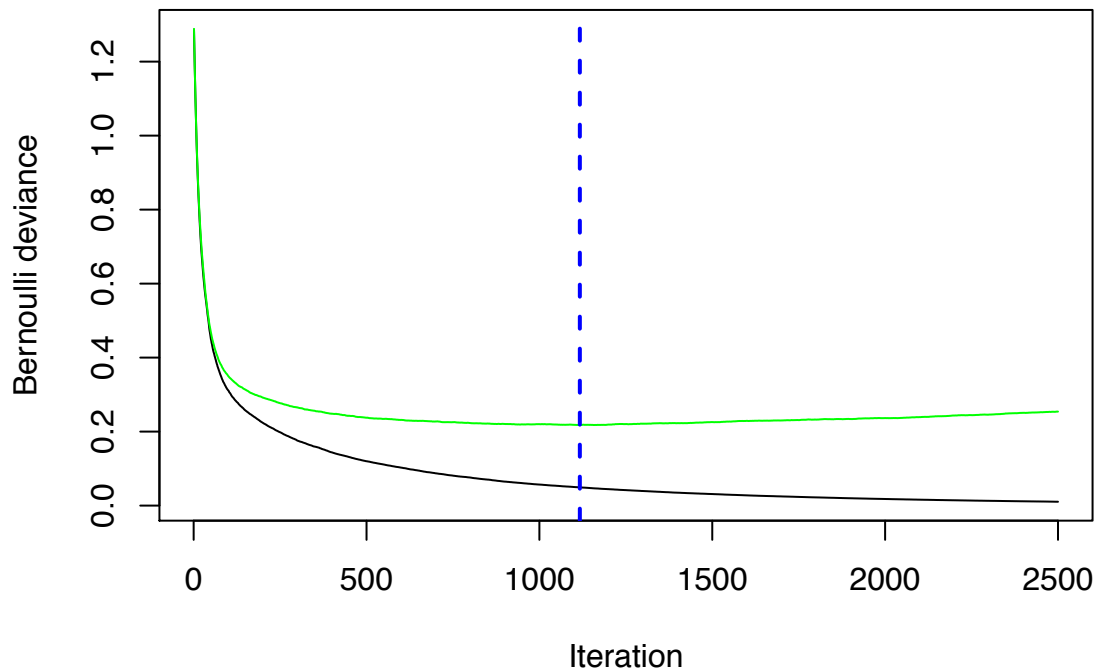
#set labels to avoid confusion
spam <- 1
nospam <- 0
threshold = 0.5

#randomize order of data
set.seed(131)
x <- train[sample(nrow(train)),]

#fit model on training data
set.seed(444) # random for bag.fraction
gbm0 <- gbm(type~., data=x, train.fraction=1,
```

```
interaction.depth=4, shrinkage=.05,
n.trees=2500, bag.fraction=0.5, cv.folds=5,
distribution="bernoulli", verbose=F)
```

```
bestTreeForPrediction = gbm.perf(gbm0)
```



```
#Training misclassification rate
y_hat_train <- predict(gbm0, x, type="response", n.trees=bestTreeForPrediction)
y_hat_train[y_hat_train > threshold] <- spam
y_hat_train[y_hat_train <= threshold] <- nonspam
y_train <- x$type
correct_train <- y_train == y_hat_train
pct_correct_train <- sum(correct_train)/length(correct_train)

#What is your estimate of the misclassification rate?
print("Training overall misclassification rate:")
1-pct_correct_train

#misclassification rate test
y_hat_test <- predict(gbm0, test, type="response", n.trees=bestTreeForPrediction)
y_hat_test[y_hat_test > threshold] <- spam
y_hat_test[y_hat_test <= threshold] <- nonspam
y_test <- test$type
correct <- y_test == y_hat_test
pct_correct_test <- sum(correct)/length(correct)
```

```

#Overall
print("Estimate of overall misclassification rate on test set:")
1-pct_correct_test

#Of all the spam emails of the test set what percentage was misclassified?
spam_correct_test <- correct[y_test==spam]
pct_spam_correct_test <- sum(spam_correct_test)/length(spam_correct_test)
print("Test spam misclassification rate:")
1-pct_spam_correct_test

#Of all the non-spam emails in the test set what percentage was misclassified?
nonspam_correct_test <- correct[y_test==nonspam]
pct_nonspam_correct_test <- sum(nonspam_correct_test)/length(nonspam_correct_test)
print("Test nonspam misclassification rate:")
1-pct_nonspam_correct_test

## Using cv method...
## [1] "Training overall misclassification rate:"
## [1] 0.001956309
## [1] "Estimate of overall misclassification rate on test set:"
## [1] 0.02346806
## [1] "Test spam misclassification rate:"
## [1] 0.02588997
## [1] "Test nonspam misclassification rate:"
## [1] 0.02183406

```

(b) Your classifier in part (a) can be used as a spam filter. One of the possible disadvantages of such a spam filter is that it might filter out too many good (non-spam) emails. Therefore, a better spam filter might be the one that penalizes misclassifying non-spam emails more heavily than the spam ones. Suppose that you want to build a spam filter that “throws out” no more than 0.3% of the good (non-spam) emails. You have to find and use a cost matrix that penalizes misclassifying “good” emails as “spam” more than misclassifying “spam” emails as “good” by the method of trial and error. Once you have constructed your final spam filter with the property described above, answer the following questions:

```

w <- x$type
w[x$type == spam] <- 1
w[x$type == nonspam] <- 2 #100

set.seed(444) # random for bag.fraction
gbm1 <- gbm(type~., data=x , train.fraction=1,
            interaction.depth=4, shrinkage=.05,
            n.trees=2500, bag.fraction=0.5, cv.folds=5,
            distribution="bernoulli", verbose=F,
            weights = w)

bestTreeForPrediction <- gbm.perf(gbm1, plot.it=FALSE,method="cv") #method={"OOB","test","cv"}

```

(i) What is the overall misclassification error of your final filter and what is the percentage of good emails and spam emails that were misclassified respectively?

```
#misclassification rate on training set
y_hat_train <- predict(gbm1, x, type="response", n.trees=bestTreeForPrediction)
y_hat_train[y_hat_train > threshold] <- spam
y_hat_train[y_hat_train <= threshold] <- nonspam
y_train <- x$type
correct_train <- y_train == y_hat_train
spam_correct_train <- correct_train[y_train==spam]
nonspam_correct_train <- correct_train[y_train==nonspam]

pct_correct_train <- sum(correct_train)/length(correct_train)
pct_spam_correct_train <- sum(spam_correct_train)/length(spam_correct_train)
pct_nonspam_correct_train <- sum(nonspam_correct_train)/length(nonspam_correct_train)

#misclassification rate on test set
y_hat_test <- predict(gbm1, test, type="response", n.trees=bestTreeForPrediction)
y_hat_test[y_hat_test > threshold] <- spam
y_hat_test[y_hat_test <= threshold] <- nonspam
y_test <- test$type
correct_test <- y_test == y_hat_test
spam_correct_test <- correct_test[y_test==spam]
nonspam_correct_test <- correct_test[y_test==nonspam]

pct_correct_test <- sum(correct_test)/length(correct_test)
pct_spam_correct_test <- sum(spam_correct_test)/length(spam_correct_test)
pct_nonspam_correct_test <- sum(nonspam_correct_test)/length(nonspam_correct_test)

#Overall misclassification:
print("Train/test overall misclassification rate:")
1-pct_correct_train
1-pct_correct_test

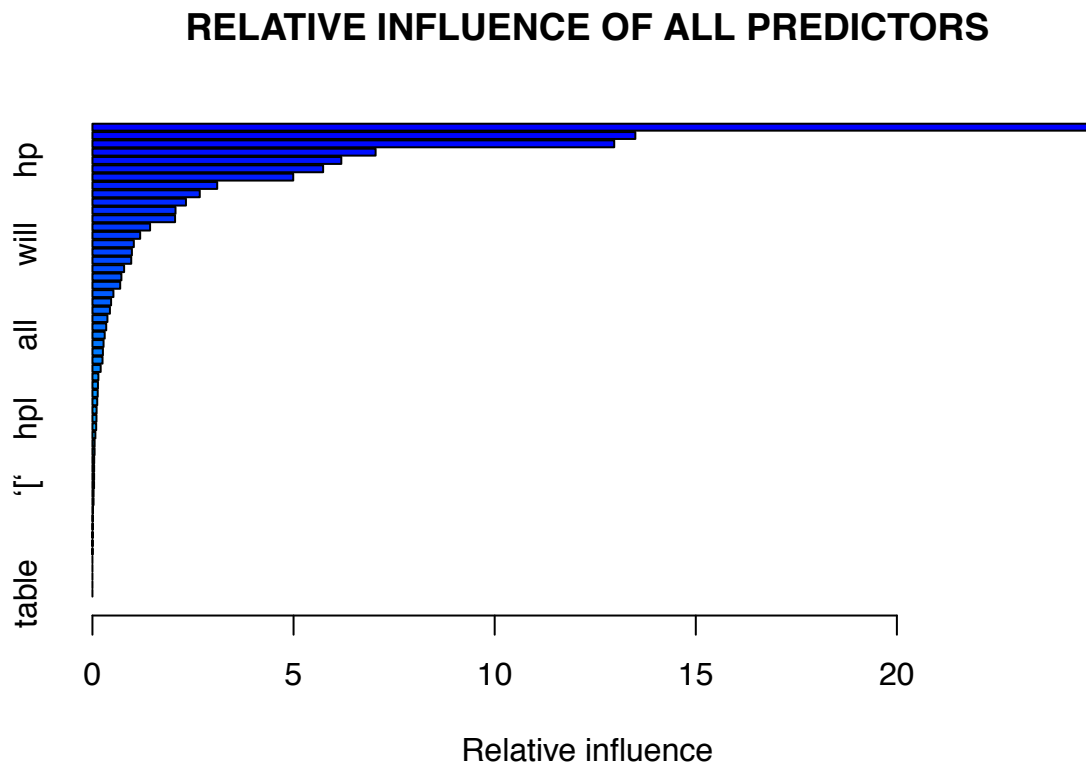
#Spam misclassification:
print("Train/test spam misclassification rate:")
1-pct_spam_correct_train
1-pct_spam_correct_test

#Nonspam misclassification:
print("Train/test nonspam misclassification rate:")
1-pct_nonspam_correct_train
1-pct_nonspam_correct_test

## [1] "Train/test overall misclassification rate:"
## [1] 0.003586567
## [1] 0.02346806
## [1] "Train/test spam misclassification rate:"
## [1] 0.009031199
## [1] 0.03398058
## [1] "Train/test nonspam misclassification rate:"
## [1] 0
## [1] 0.01637555
```

(ii) What are the important variables in discriminating good emails from spam for your spam filter?

```
top5 <-summary(gbm1,main="RELATIVE INFLUENCE OF ALL PREDICTORS")$var[1:5]
```



```
top5 <- gsub("`", "", top5)
print(top5)
```

```
## [1] "remove" "!" "$" "free" "hp"
```

(iii) Using the interpreting tools provided by gbm, describe the dependence of the response on the most important attributes.

Partial dependence plots show what the average change in the value of our model’s output is when varying a subset of variables. In this case, we are looking at single variables. Looking at the partial dependence, we can understand what the effect of increasing or decreasing that variable’s value is on the model’s output (holding all other variables constant). This gives us a better understanding of the behavior of otherwise “black box” models.

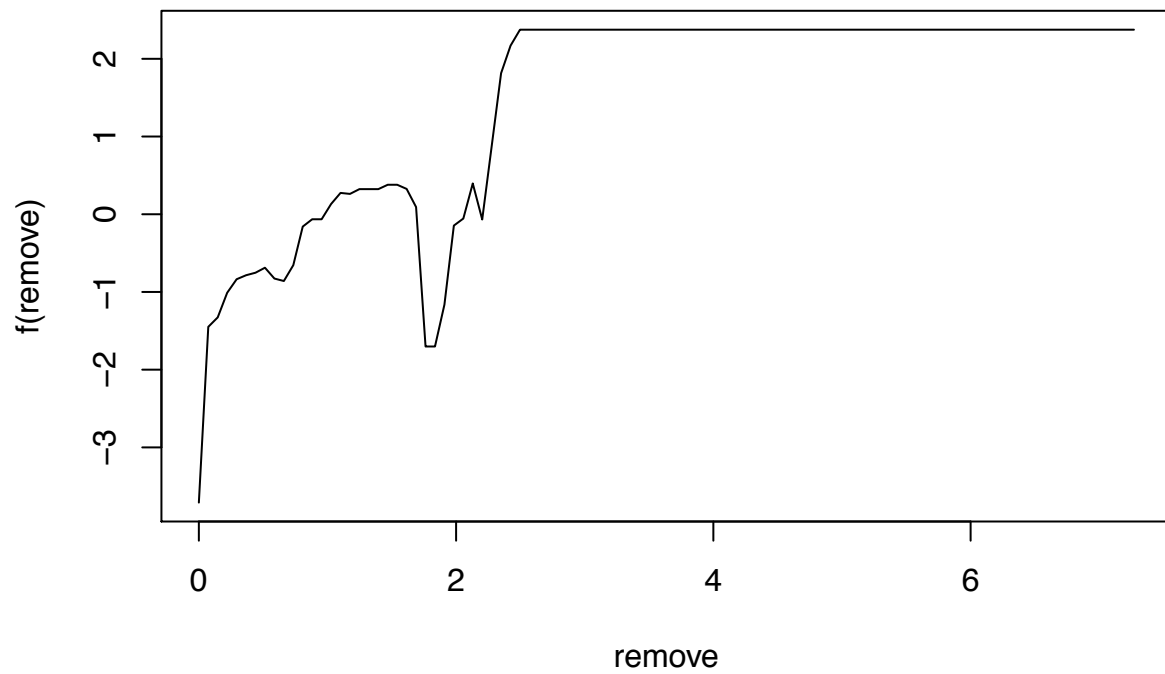
Looking at the variable with the most relative influence, “remove”, we can see that as we increase the frequency of the word remove as a proportion of an email’s length, it increases the average score of the model’s output. This means that emails with a high frequency of “remove” relative to their length are indicative of spam since spam is the positive class and nonspam is the zero class.

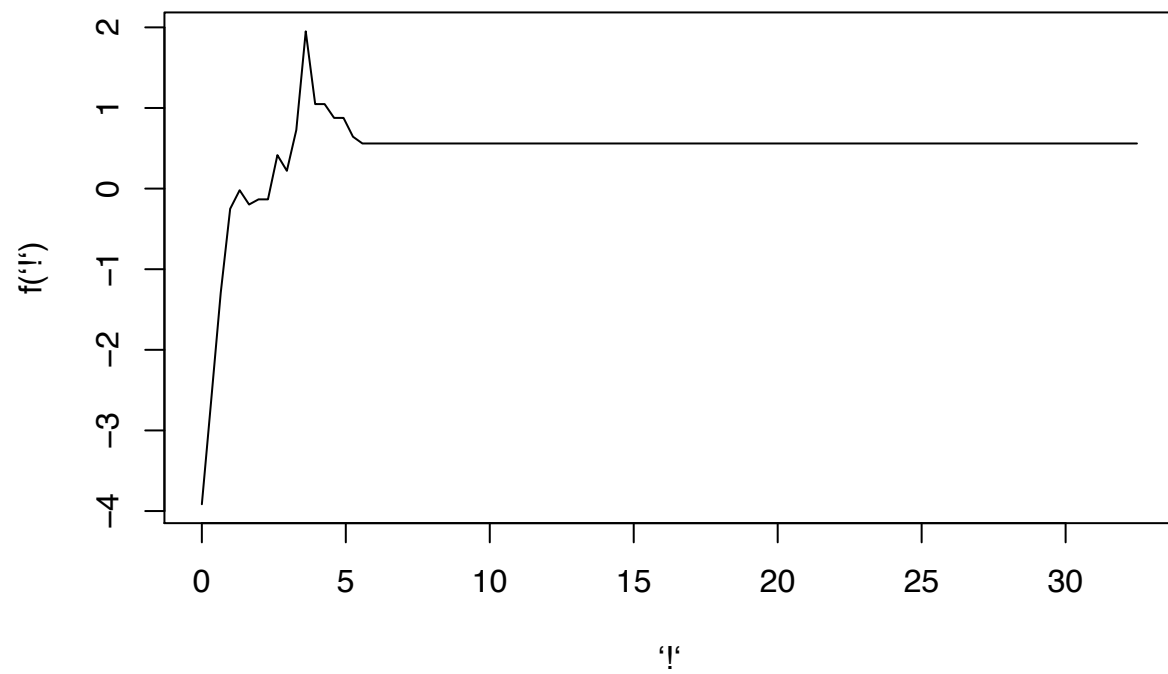
We see opposite behavior in “hp”. Since George and his colleagues worked at Hewlett-Packard, emails with “hp” display the opposite effect and the higher we move “hp” in frequency the more likely it is that the email

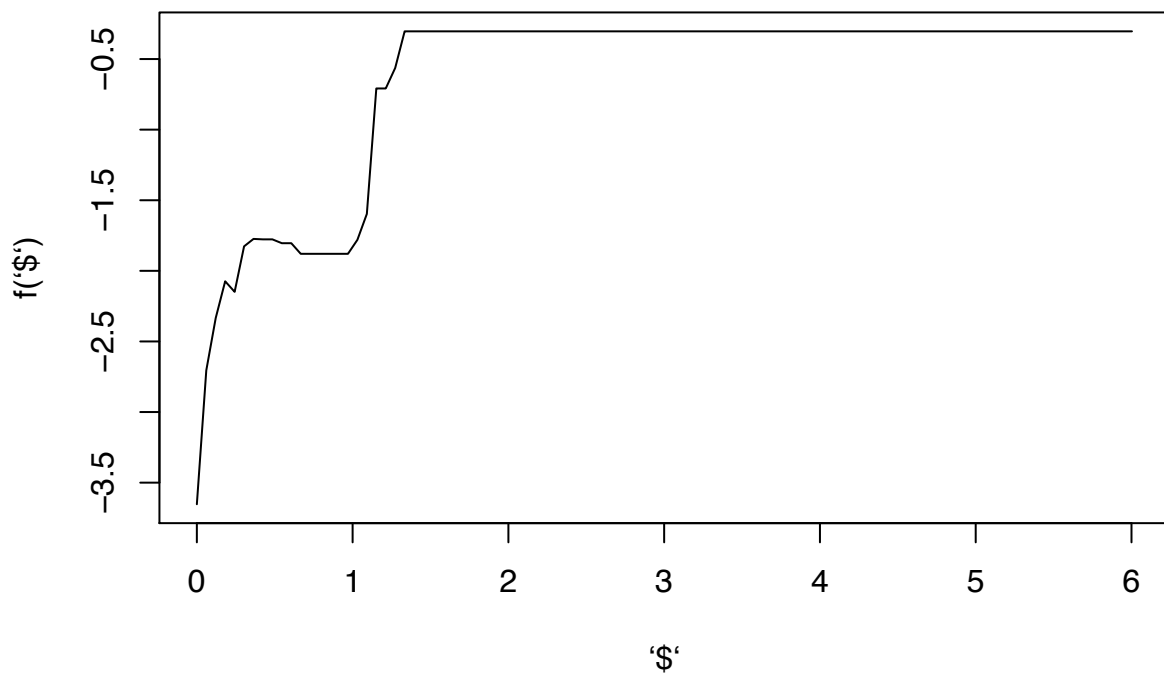


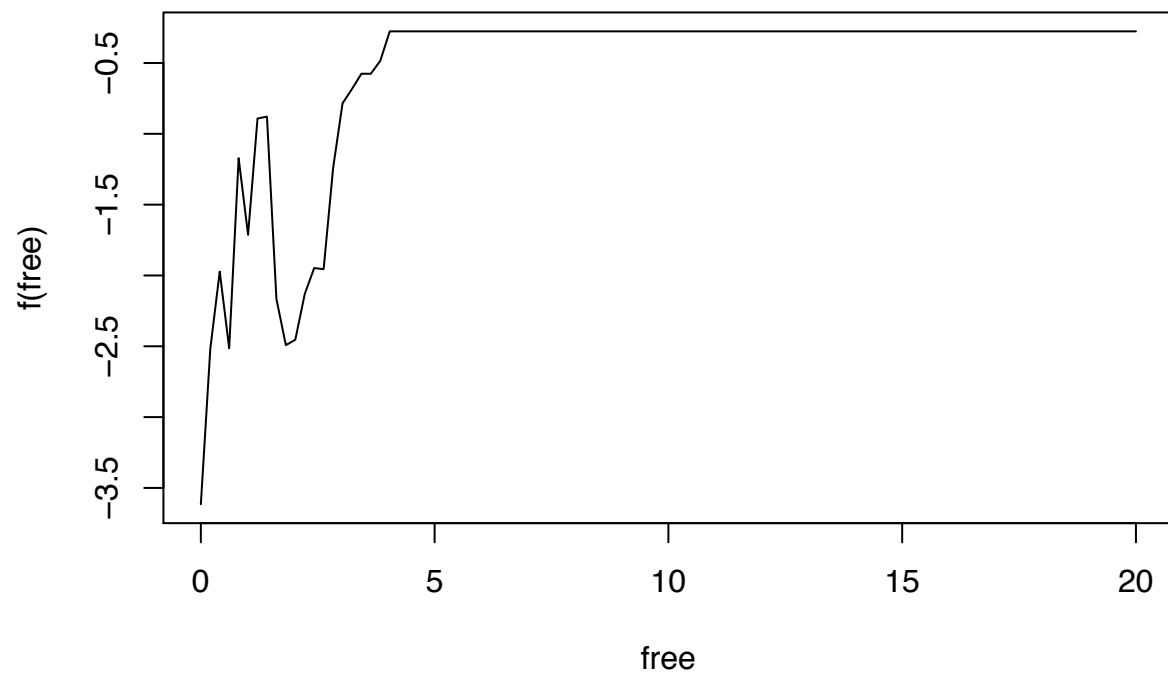
is legitimate and non-spam (according to our model).

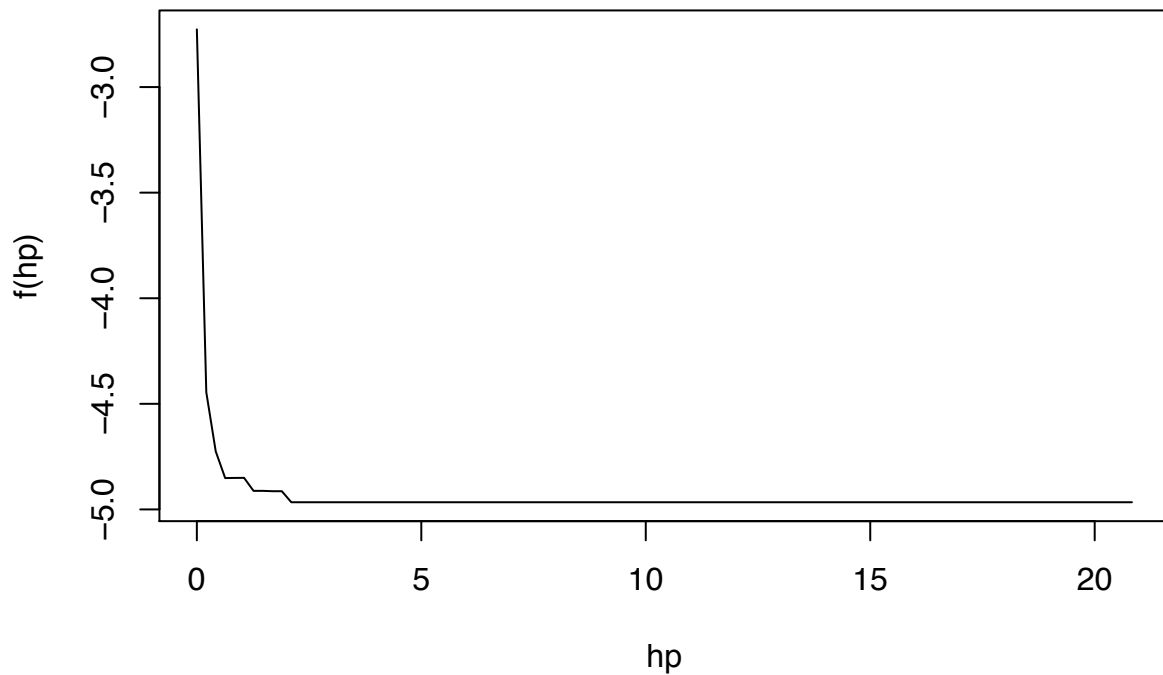
```
top5_indexes <- match(top5, rflabs)
for(i in top5_indexes){
  plot(gbm1,i,bestTreeForPrediction)
}
```











## Question 7

Regression: California Housing. The data set `calif_stats315B.csv` consists of aggregated data from 20,640 California census blocks (from the 1990 census). The goal is to predict the median house value in each neighborhood from the others described in `calif_stats315B.txt`. Fit a gbm model to the data and write a short report that should include at least

```
rm(list = ls())
data_path <- paste(getwd(), '/data', sep='')
setwd(data_path)

#data labs
labels <-c(
  "house_value",
  "median_income",
  "housing_median_age",
  "average_no_rooms",
  "average_no_bedrooms",
  "population",
  "average_occupancy",
  "latitude",
  "longitude"
)
```

```

#load the data
data <- read.csv(file="calif_stats315B.csv", header=FALSE, sep=",")
colnames(data) <- labels

#fit model on training data
set.seed(444) # random for bag.fraction
model <- gbm(house_value~., data=data, train.fraction=1,
             interaction.depth=4, shrinkage=.05,
             n.trees=2500, bag.fraction=0.5, cv.folds=5,
             distribution="gaussian", verbose=F)

bestTreeForPrediction <- gbm.perf(model, plot.it=FALSE,method="cv") #method={"OOB","test","cv"}

```

**(a) The prediction accuracy of gbm on the data set.**

The MSE of the response variable house\_value is  $\sim 0.15$ . Since the response is in units of \$100,000, this means on average the squared error of the median house\_value predicted by our model is  $\sim 0.154 \times \$100,000$ , or  $\sim \$15,400$ .

```

#MSE for accuracy
y <- data$house_value
y_hat <- predict(model, data, type="response", n.trees=bestTreeForPrediction)
MSE <- sum((y_hat - y)^2)/length(y)
MSE

## [1] 0.1543481

```

**(b) Identification of the most important variables.**

The most important variable for predicting house\_value for a particular neighborhood in California is median\_income, which intuitively makes sense.

The second-tier of variables in terms of importance had similar levels of importance: average\_occupancy, latitude and longitude.

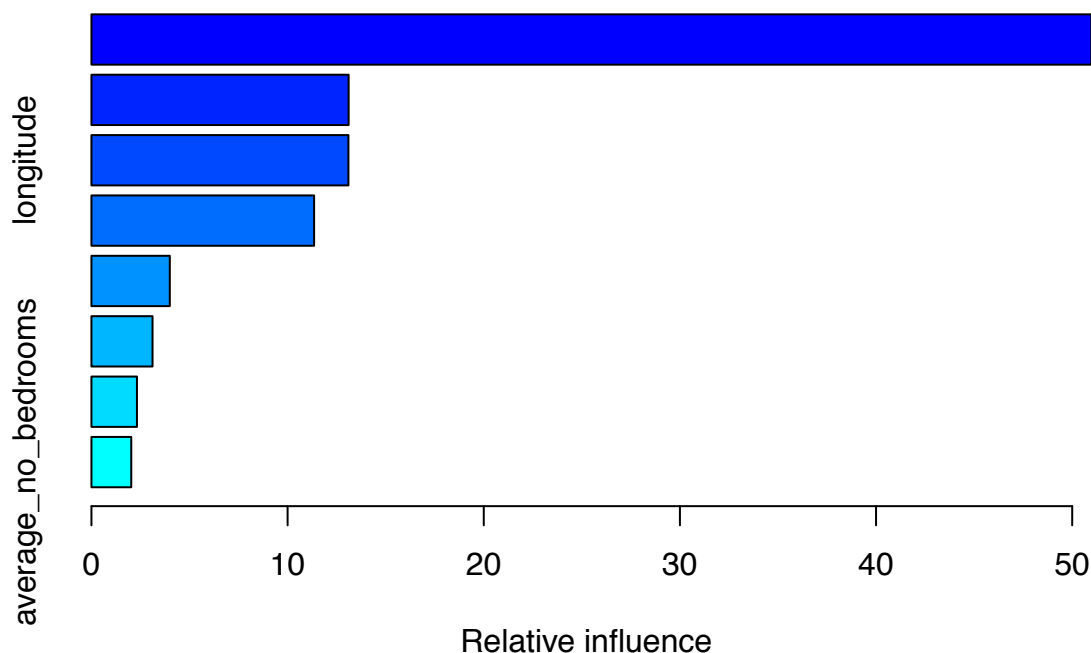
The 4 variables that all had the similar levels of insignificance in predicting house\_value for a given neighborhood: housing\_median\_age, average\_no\_rooms, population, average\_no\_bedrooms.

```

top4 <-summary(model,main="RELATIVE INFLUENCE OF ALL PREDICTORS")$var[1:4]

```

## RELATIVE INFLUENCE OF ALL PREDICTORS



```
top4 <- gsub("`", "", top4)
print(top4)
```

```
## [1] "median_income"      "average_occupancy" "longitude"
## [4] "latitude"
```

(c) Comments on the dependence of the response on the most important variables (you may want to consider partial dependence plots (plot) on single and pairs of variables, etc.).

We can see from the individual dependence plots that as median\_income increases holding other variables constant, our model's output for home\_value also increases.

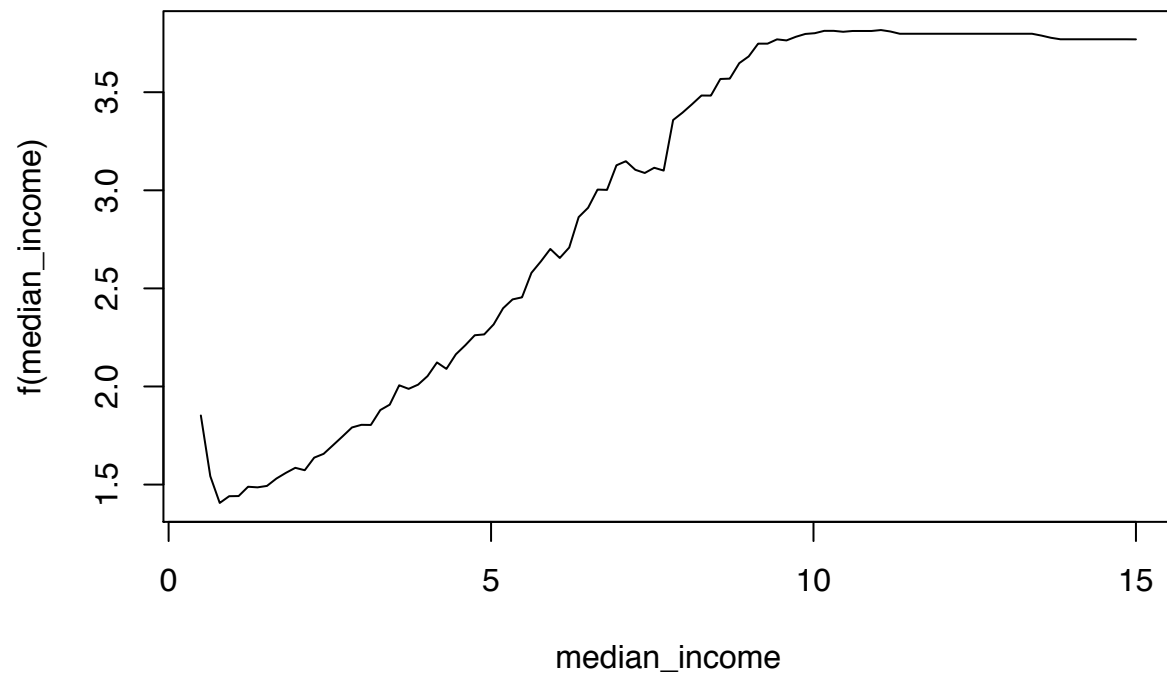
The interaction effects between latitude and longitude are enlightening. The border between the pink and blue zones have interior corner points at the Bay Area ( $latlng = (38.0, -122.0)$ ) and Los Angeles ( $latlng = (34.0, -118.0)$ ). It is unsurprising that these major metropolitan areas are correlated with home\_value.

The univariate dependence plots show that home\_value tends to decrease as homes get further from the coast. It also shows that home\_value tends to decrease as homes get further north (which makes sense since the Bay area is at the ~half way mark of California in terms of latitude). However, there is important nuance in the interaction effects between these two variables as discussed in the paragraph above.

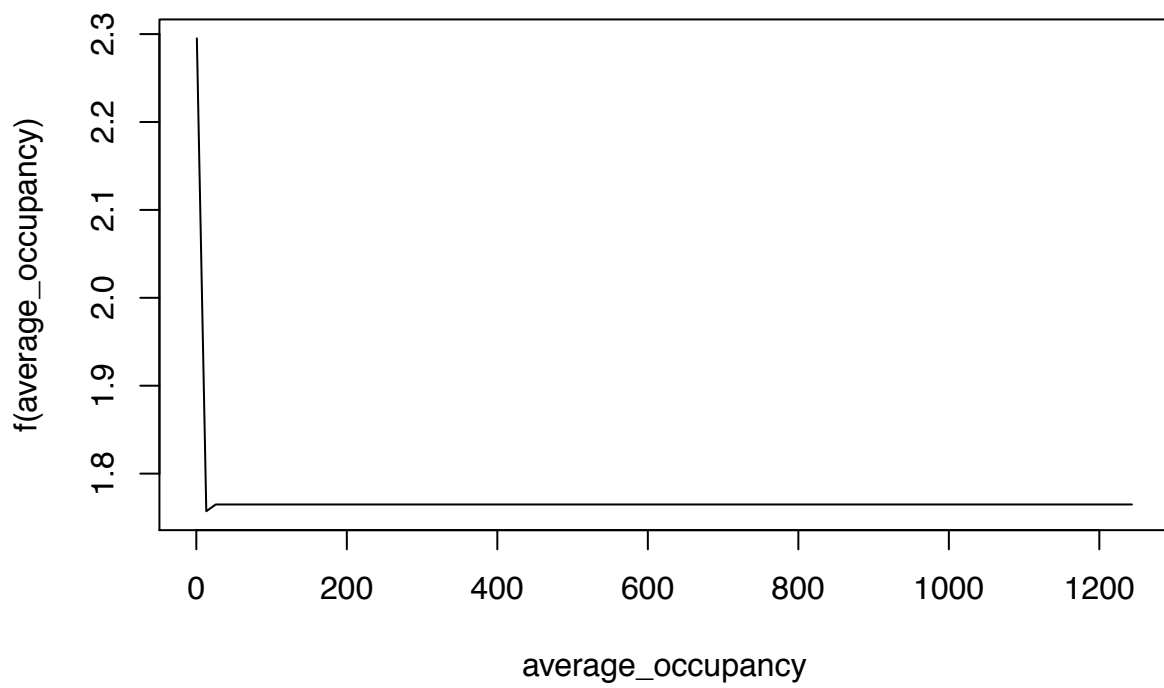
```
best.iter <- gbm.perf(model, plot.it=FALSE, method='cv')
top4_indexes <- match(top4, labels) - 1 # -1 because first label is response var

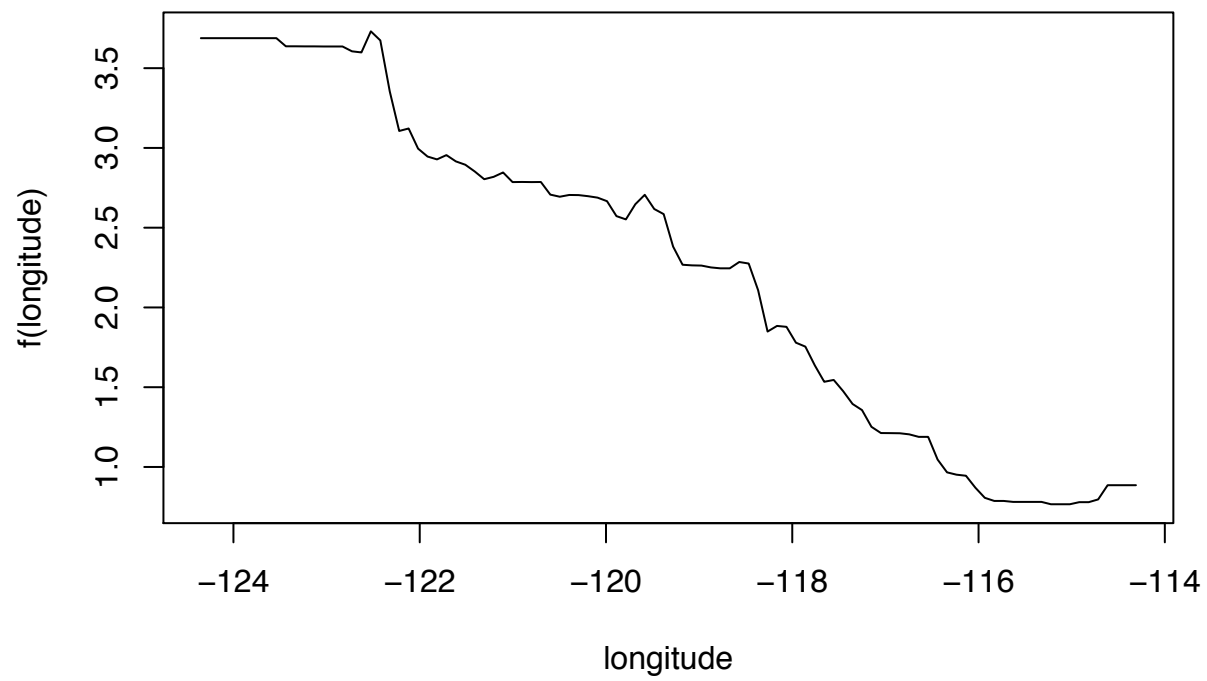
#main effects
for(i in top4_indexes){
```

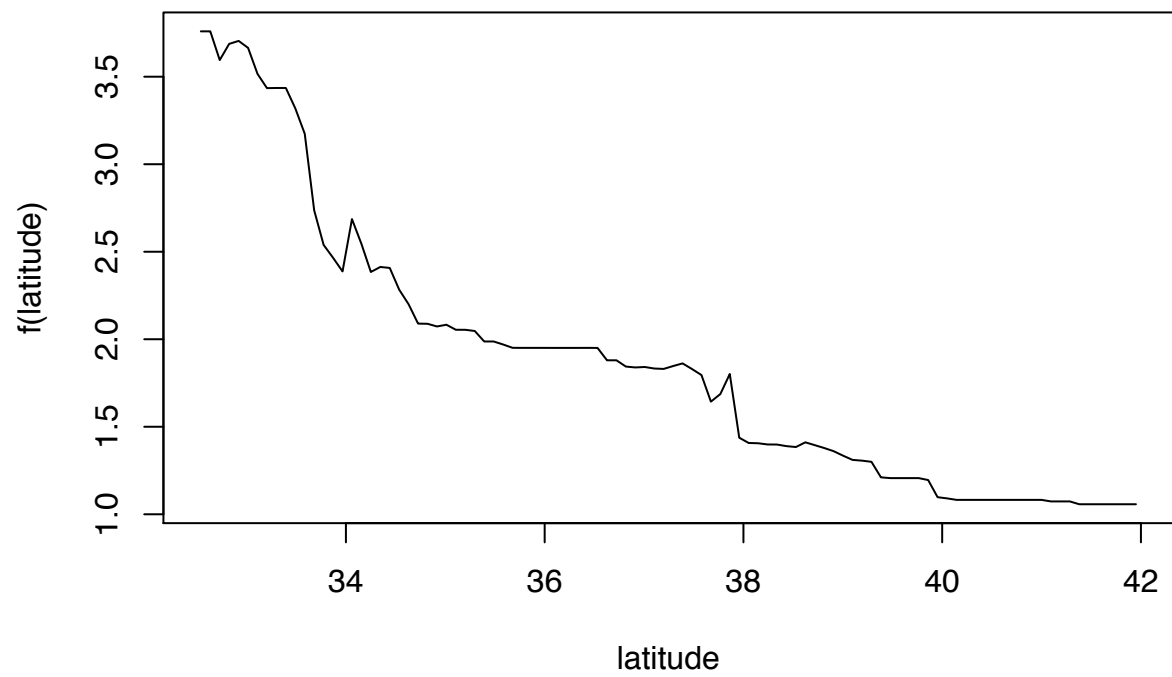
```
plot(model,i,best.iter)
}
```



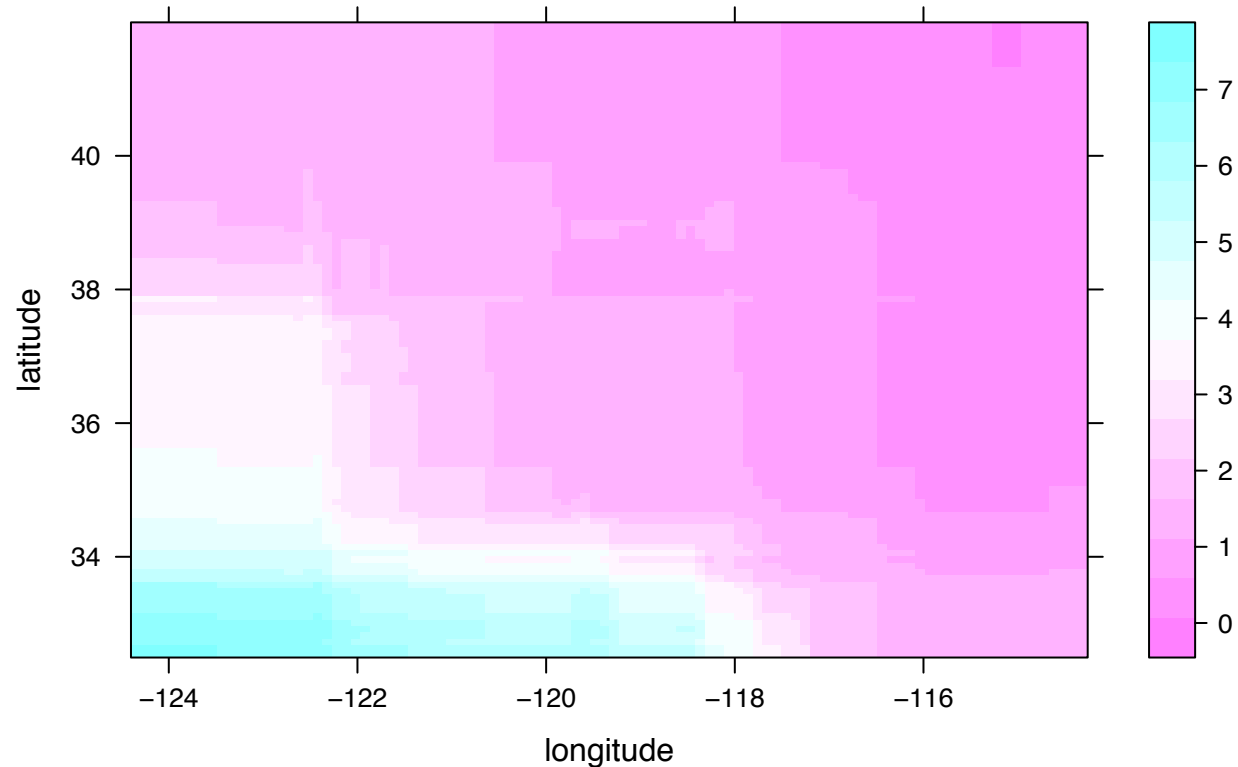








```
#longitude, latitude interaction effects  
plot(model,c(8,7),best.iter)
```



## Question 8

Regression: Marketing data. The data set `age_stats315B.csv` was already used in Homework 1. Review `age_stats315B.txt` for the information about order of attributes etc.

(a) Fit a gbm model for predicting age form the other demographic attributes and compare the accuracy with the accuracy of your best single tree from Homework 1.

```
rm(list = ls())
data_path <- paste(getwd(), '/data', sep='')
setwd(data_path)

#Read and type data
age_data <- read.csv('age_stats315B.csv')
factor_columns <- c(
  'Occup',
  'TypeHome',
  'sex',
  'MarStat',
  'DualInc',
  'HouseStat',
  'Ethnic',
  'Lang'
```

```

)
age_data[factor_columns] <- lapply(age_data[factor_columns], as.factor)

#fit a GBM model
set.seed(444) # random for bag.fraction
model_gbm <- gbm(age~., data=age_data, train.fraction=1,
  interaction.depth=4, shrinkage=.05,
  n.trees=2500, bag.fraction=0.5, cv.folds=5,
  distribution="gaussian", verbose=F)

#fit a tree
model_tree <- rpart(age ~ ., data = age_data, method = "anova",
  control=rpart.control(minbucket = 10,
    xval = 10,
    maxsurrogate = 5,
    usesurrogate = 2,
    cp=0.0001))

# Find the minimum cross-validation error + one SD
min_error_window <- min(model_tree$cptable[, "xerror"] + model_tree$cptable[, "xstd"])

# Find the simplest model with xerror within the min_error_window
best_cp <- first(model_tree$cptable[which(model_tree$cptable[, "xerror"] < min_error_window), "CP"])
best_single_tree <- prune(model_tree, cp = best_cp)

bestTreeForPrediction <- gbm.perf(model_gbm, method='cv', plot.it=FALSE) #method={"OOB", "test", "cv"}

#Make predictions
y_hat_gbm <- predict(model_gbm, age_data, type="response", n.trees=bestTreeForPrediction)
y_hat_tree <- predict(best_single_tree, age_data)

#Compare errors
y <- age_data$age
MSE_gbm <- sum((y_hat_gbm - y)^2)/length(y)
MSE_tree <- sum((y_hat_tree - y)^2)/length(y)

MSE_gbm

## [1] 0.5648001
MSE_tree

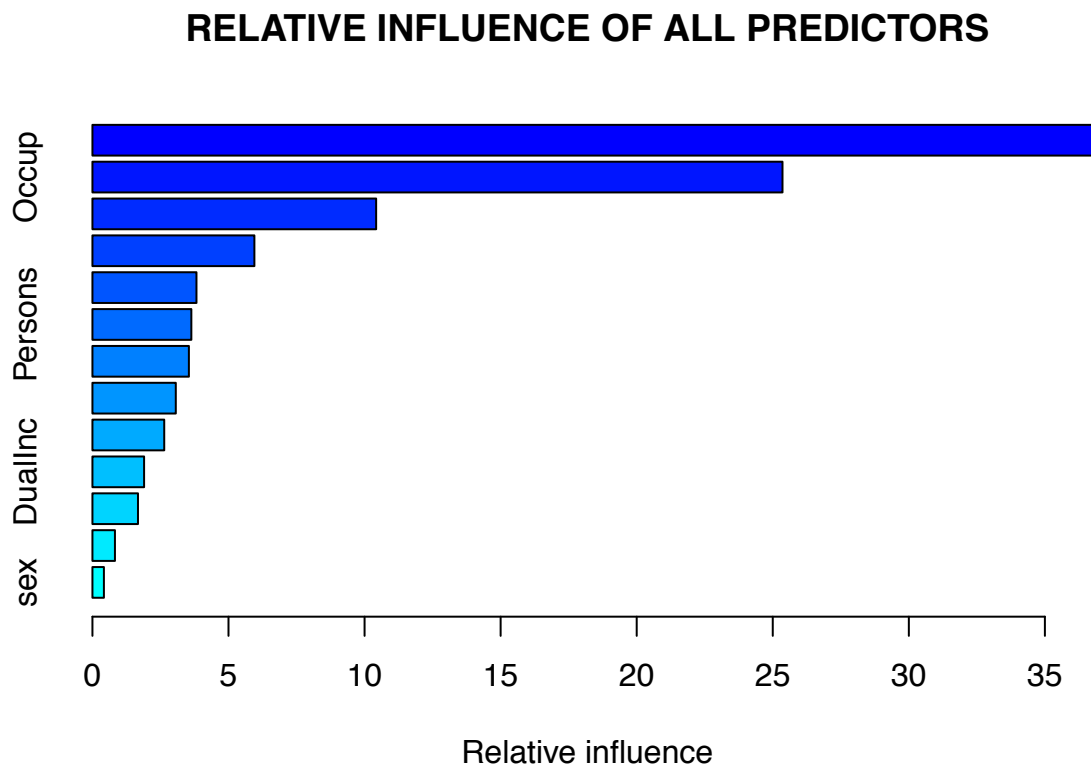
## [1] 0.6992926

```

**(b) Identify the most important variables.**

The most important variables were also identified as the most important in our tree from homework 1: MarStat, Occup, HouseStat, and Edu

```
top4 <-summary(model_gbm,main="RELATIVE INFLUENCE OF ALL PREDICTORS")$var[1:4]
```



```
top4 <- gsub("`", "", top4)
print(top4)
```

```
## [1] "MarStat" "Occup" "HouseStat" "Edu"
```

## Question 9

**Multiclass classification: marketing data.** The data set `occup_stats315B.csv` comes from the same marketing database used in Homework 1. The description of the attributes can be found in `occup_stats315B.txt`. The goal in this problem is to fit a gbm model to predict the type of occupation from the 13 other demographic variables.

```
rm(list = ls())
data_path <- paste(getwd(), '/data', sep='')
setwd(data_path)

#Read and type data
labels <-c(
  "Occup",
  "TypeHome",
  "sex",
  "MarStat",
```

```

"age",
"Edu",
"Income",
"LiveBA",
"DualInc",
"Persons",
"Under18",
"HouseStat",
"Ethnic",
"Lang"
)
factor_columns <- c(
  'TypeHome',
  'sex',
  'MarStat',
  'Occup',
  'LiveBA',
  'DualInc',
  'HouseStat',
  'Ethnic',
  'Lang'
)

occup_labels <- c(
  "Professional/Managerial",
  "Sales Worker",
  "Factory Worker/Laborer/Driver",
  "Clerical/Service Worker",
  "Homemaker",
  "Student, HS or College",
  "Military",
  "Retired",
  "Unemployed"
)

#load the data
house_data <- read.csv('occup_stats315B.csv', header=FALSE, sep=",")
colnames(house_data) <- labels
house_data[factor_columns] <- lapply(house_data[factor_columns], as.factor)

#fit GBM model for occupation
set.seed(444) # random for bag.fraction
model <- gbm(Occup ~ ., data=house_data, train.fraction=1,
  interaction.depth=4, shrinkage=.05,
  n.trees=2500, bag.fraction=0.5, cv.folds=5,
  distribution="multinomial", verbose=F)

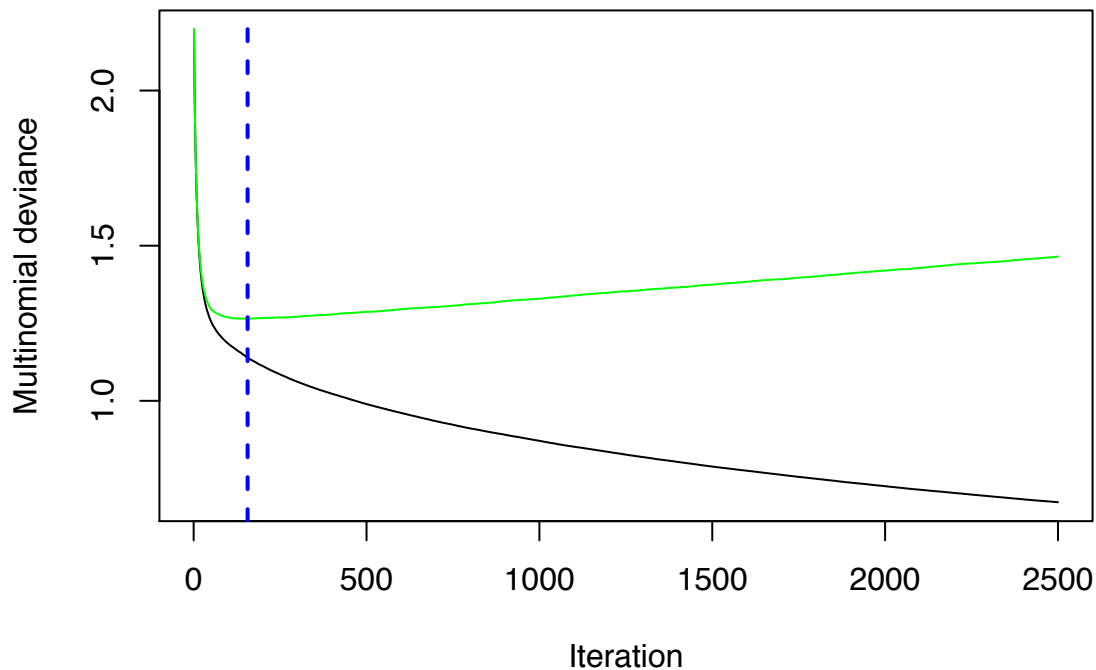
```

(a) Report the test set misclassification error for gbm on the data set, and also the misclassification error for each class.

```

bestTreeForPrediction <- gbm.perf(model, method='cv') #method={"OOB", "test", "cv"}

```



```
#determine misclassification rate
y_hat <- predict(model, house_data, type="response", n.trees=bestTreeForPrediction)
y_hat_max_p <- apply(y_hat,1,which.max)
y <- house_data$Occup
correct <- y == y_hat_max_p
```

```
#overall
pct_correct <- sum(correct)/length(correct)
print("Overall misclassification rate:")
```

```
## [1] "Overall misclassification rate:"
```

```
1-pct_correct
```

```
## [1] 0.3880546
```

```
for(occupation in levels(house_data$Occup)){
  occupation_i <- strtoi(occupation)
  print(occup_labels[occupation_i])
  correct_in_class <- correct[y==occupation_i]
  pct_correct_in_class <- sum(correct_in_class)/length(correct_in_class)
  print(1-pct_correct_in_class)
}
```

```
## [1] "Professional/Managerial"
```

```
## [1] 0.1829878
```

```
## [1] "Sales Worker"
```

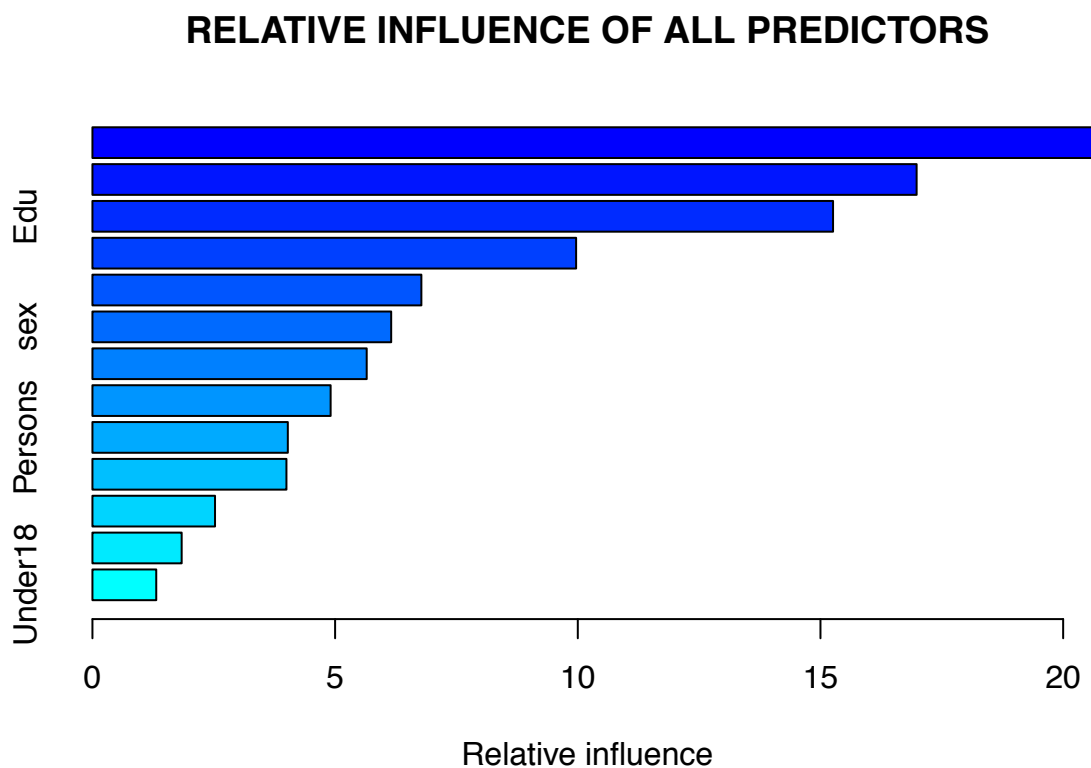
```
## [1] 0.8994911
```



```
## [1] "Factory Worker/Laborer/Driver"
## [1] 0.610596
## [1] "Clerical/Service Worker"
## [1] 0.6758818
## [1] "Homemaker"
## [1] 0.3636364
## [1] "Student, HS or College"
## [1] 0.1681356
## [1] "Military"
## [1] 0.6141732
## [1] "Retired"
## [1] 0.1776504
## [1] "Unemployed"
## [1] 0.7439353
```

(b) Identify the most important variables.

```
top4 <-summary(model,main="RELATIVE INFLUENCE OF ALL PREDICTORS")$var[1:4]
```



```
top4 <- gsub("`", "", top4)
print(top4)
```

```
## [1] "age"      "Income"   "Edu"      "HouseStat"
```