

Census pay predictions

Rajesh Haridas

December 15, 2021

Preface

This document and the associated code is based on the Capstone HarvardX Professional Certificate in Data Science (PH125.9x) course work and additional reading material provided in the course.

The following files (3 types) are included in the upload:

- **CensusPay.rmd** - Markdown for the main summary report file
- **CensusPay.R** - Main R code for the project
- **DatasetProcessingCode.R** - The R code for downloading, cleaning and converting the dataset into tidy form to be used in the project
- **CensusPaySummaryReport.pdf** - Main summary report containing analysis
- **CensusPayExecutionReport.pdf** - Main execution report containing output of the runs

Contents

Preface	1
Introduction to Census income level predictions	3
Analysis	6
Data Exploration, cleaning, processing and, feature engineering	6
Process	23
Model Creation	24
Results	51
Conclusion	52
Appendix A - Complete code	52
Appendix B - Links	74
Citations	74

List of Figures

1	Income distribution	6
2	Education Frequency	8
3	Education Frequency	8
4	Education distribution	9
5	Education distribution	10
6	Occupation frequency	11
7	Occupation frequency	12
8	Occupation distribution	12
9	Sex feature characteristics	13
10	Sex feature characteristics	13
11	Sex feature characteristics	14
12	Race distribution	15
13	Race distribution	15
14	Race distribution	16
15	Marital Status frequency	17
16	Marital Status frequency	17
17	Marital Status distribution	18
18	Relationship distribution	19
19	Relationship distribution	19
20	Relationship distribution	20
21	Class distribution	21
22	Class distribution	21
23	Class distribution	22
24	GLM classification summary	29
25	naive bayes classification summary	32
26	KNN tuning	34
27	knn classification summary	36
28	knn tune classification summary	39
29	rpart classification summary	42
30	random forest classification summary	46
31	random forest tuned classification summary	50
32	Final Results	51

Introduction to Census income level predictions

The Current Population Survey (CPS), sponsored jointly by the U.S. Census Bureau and the U.S. Bureau of Labor Statistics (BLS), is the primary source of labor force statistics for the population of the United States. The adult census income data was extracted from the 1994 Census bureau database by Ronny Kohavi and Barry Becker (Data Mining and Visualization, Silicon Graphics).

This machine learning project uses the adult census income data to predict annual money incomes for adults, given a set of 13 employment and demographic attributes. Census money income is defined as income received on a regular basis before payments for personal income taxes, social security, union dues, Medicare deductions, etc. The income levels were categorized into two classes – more than \$50,000 and less than or equal to \$50,000.

More information can be found at <https://www.kaggle.com/uciml/adult-census-income>.

In this project for the Capstone course of the HarvardX Professional Certificate in Data Science (PH125.9x), we will explore the adult census income data set. The objective will be to analyze and model based on the classification of incomes and develop a machine-learning model by creating, training and test sets to predict income levels on a validation set with accuracy of more than 80% and a reasonable sensitivity and specificity

Here is a glimpse of adultpay dataset. The columns that are of interest are income, age, sex, education.num, education, race, marital.status, relationship, workclass, hours.per.week and occupation. There are 32,561 rows and 13 columns

```
# glimpse(adultpay)

# Rows: 32,561

# Columns: 13

# $ age <int> 90, 82, 66, 54, 41, 34, 38, 74, 68, 41, 45, 38, 52, 32, 51, 46,~

# $ workclass <chr> '?', 'Private', '?', 'Private', 'Private', 'Private',
# 'Private'~

# $ fnlwgt <int> 77053, 132870, 186061, 140359, 264663, 216864, 150601, 88638,
# 4~

# $ education <chr> 'HS-grad', 'HS-grad', 'Some-college', '7th-8th',
# 'Some-college'~

# $ education.num <int> 9, 9, 10, 4, 10, 9, 6, 16, 9, 10, 16, 15, 13, 14, 16,
# 15, 7, 14~

# $ marital.status <chr> 'Widowed', 'Widowed', 'Widowed', 'Divorced',
# 'Separated', 'Divo~

# $ occupation <chr> '?', 'Exec-managerial', '?', 'Machine-op-inspct',
# 'Prof-special~

# $ relationship <chr> 'Not-in-family', 'Not-in-family', 'Unmarried',
# 'Unmarried', 'Dw~

# $ race <chr> 'White', 'White', 'Black', 'White', 'White', 'White', 'White', ~

# $ sex <chr> 'Female', 'Female', 'Female', 'Female', 'Female', 'Female', 'Ma~
```

```
# $ hours.per.week <int> 40, 18, 40, 40, 40, 45, 40, 20, 40, 60, 35, 45, 20,
# 55, 40, 40, ~

# $ native.country <chr> 'United-States', 'United-States', 'United-States',
# 'United-Stat~

# $ income <chr> '<=50K', '<=50K', '<=50K', '<=50K', '<=50K', '<=50K', '<=50K',
# ~
```

We will start with data exploration and cleaning unwanted data and processing. Then we will mimic the ultimate evaluation process by splitting the data into two parts - training and validation and act as if we don't know the outcome for one of these. We will split the dataset once we cleaned and processed the dataset. We want the test set to be large enough so that we obtain a stable prediction without fitting an impractical number of models. We will then progressively apply various algorithms on the training set and test the predictions on the test set and improve on the overall accuracy.

We will choose 13 attributes, income being be the outcome and the rest 12 attributes being the predictors or features. Looking at the dimensions we will randomly choose 10% of the dataset to be the test set and remaining 90% will be the training set. 90% is sizable enough to run various algorithms like lm, glm, knn, random forest with some tuning.

```
# > dim(adultpayclean)

# [1] 29170 13

# > dim(adultpayclean_train)

# [1] 26252 13

# > dim(adultpayclean_validation)

# [1] 2918 13

# > colSums(is.na(adultpayclean))

# age fnlwgt education education.num marital.status occupation relationship

# 0 0 0 0 0 0 0

# race sex hours.per.week native.country income class

# 0 0 0 0 0 0
```

This is a classification problem and the outcome is binary (income above-50K or at-or-below-50K). On reviewing the summary of the dataset we have a choice of predictors. We will keep only USA specific data. Almost all of the predictors have data spread out. There are no NAs in the dataset. There is sufficient degree of variability.

```
# summary(adultpayclean)

# age fnlwgt education eduyears maritalstatus

# Min. :17.00 Min. : 12285 HSgrad :9702 Min. : 1.00 Divorced : 4162
```

```

# 1st Qu.:28.00 1st Qu.: 115895 Somecollege:6740 1st Qu.: 9.00 MarriedAFspouse
# : 23

# Median :37.00 Median : 176730 Bachelors :4766 Median :10.00 Marriedciuspouse
# :13368

# Mean :38.66 Mean : 187069 Masters :1527 Mean :10.17 Marriedspouseabsent: 253

# 3rd Qu.:48.00 3rd Qu.: 234139 Assocvoc :1289 3rd Qu.:12.00 Nevermarried :
# 9579

# Max. :90.00 Max. :1484705 11th :1067 Max. :16.00 Separated : 883

# (Other) :4079 Widowed : 902

# occupation relationship race sex hoursperweek

# Execmanagerial:3735 Husband :11861 Amer-Indian-Eskimo: 296 Female: 9682 Min.
# : 1.00

# Profspecialty :3693 Notinfamily : 7528 Asian-Pac-Islander: 292 Male :19488
# 1st Qu.:40.00

# Craftrepair :3685 Otherrelative: 696 Black : 2832 Median :40.00

# Admclerical :3449 Dunchild : 4691 Other : 129 Mean :40.45

# Sales :3364 Unmarried : 3033 White :25621 3rd Qu.:45.00

# Otherservice :2777 Wife : 1361 Max. :99.00

# (Other) :8467

# native income class

# Length:29170 Above50K : 7171 Private :20135

# Class :character AtBelow50K:21999 Selfempnotinc: 2313

# Mode :character Localgov : 1956

# Unknown : 1659

# Stategov : 1210

# Selfempinc : 991

# (Other) : 906

```

We can also see the number or above-50K income is significantly less than (30%) the at-or-below-50K income. This implies that we will encounter some prevalence towards the at-or-below 50K group

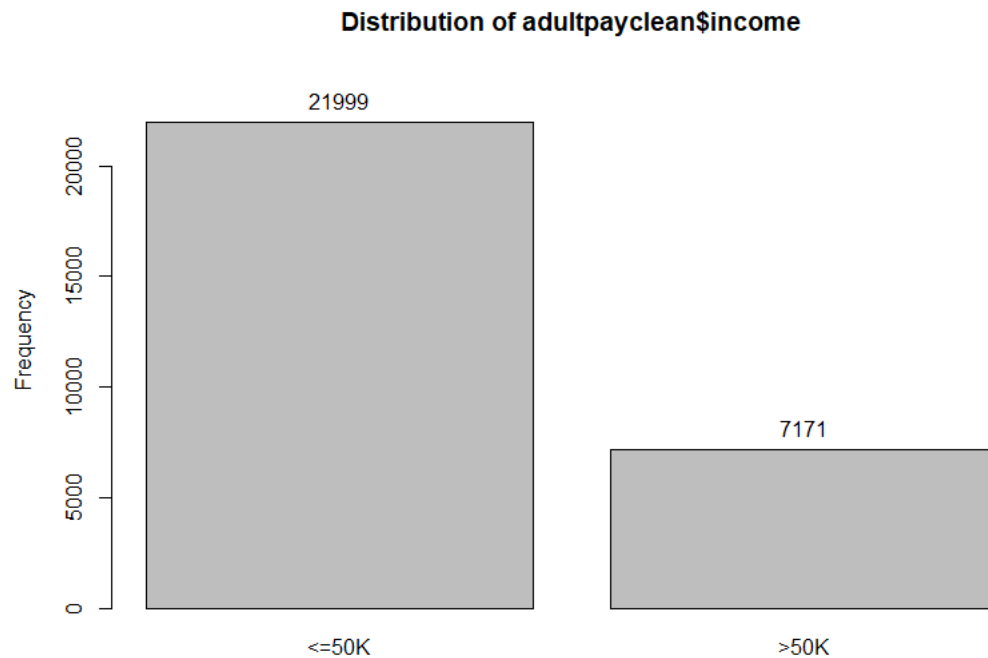


Figure 1: Income distribution

So we begin...

Analysis

Data Exploration, cleaning, processing and, feature engineering

The adult pay dataset from census bureau is a ML ready database. For convenience, I have downloaded it from kaggle and stored in the same github repository as this code and markdown. The dataset is in zip format and is downloaded from <https://github.com/rajeshharidas/havardxwork2/blob/main/adult.csv.zip>. Then it is read into a data frame. The data is ML ready for the most part. However, there are some additional cleaning tasks needed before using the dataset. This processing is done in DatasetProcessingCode.R file. Code fragment is shown below. We perform the following cleaning tasks

- 1) Filter to keep only USA data
- 2) Remove all ? marks data from all the columns
- 3) Remove intermediate columns and columns that are not used in the project
- 4) Rename columns to exclude non-alphanumeric characters
- 5) convert character labels to factors

```
# adultpayclean <-
# adultpay %>% filter (native.country == 'United-States') %>%
# mutate (class = ifelse(workclass == '?', 'Unknown',
```

```

# str_replace_all(workclass, '-', '')) %>%

# select(-workclass, -capital.gain, -capital.loss) %>%

# rename( c( eduyears = education.num, maritalstatus = marital.status,
# hoursperweek = hours.per.week, native = native.country ) ) %>%

# mutate (maritalstatus = ifelse( maritalstatus == '?', 'Unknown',
# str_replace_all(maritalstatus, '-', '' ) ) %>%

# mutate (occupation = ifelse( occupation == '?', 'Unknown',
# str_replace_all(occupation, '-', '' ) ) %>%

# mutate (education = ifelse(education == '?', 'Unknown',
# str_replace_all(education, '-', '')) %>%

# mutate (relationship = ifelse( relationship == '?', 'Unknown',
# str_replace_all(relationship, '-', '' ) ) %>%

# mutate (native = ifelse(native == '?', 'Unknown', str_replace_all(native,
# '-', ''))) %>%

# mutate (income = ifelse( income == '?', 'Unknown', str_replace_all(income,
# '<=50K', 'AtBelow50K') ) ) %>%

# mutate (income = ifelse( income == '?', 'Unknown', str_replace_all(income,
# '>50K', 'Above50K') ) )

```

After performing this cleaning we see 29170 rows and 13 columns. We further analyze the makeup of the categorical columns

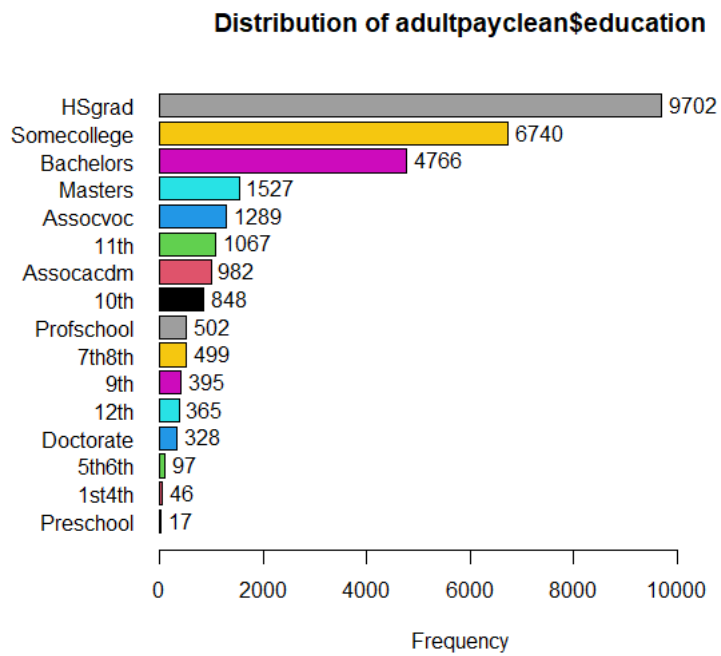


Figure 2: Education Frequency

	Frequency	Percent	Cum. percent
HSgrad	9702	33.3	33.3
Somecollege	6740	23.1	56.4
Bachelors	4766	16.3	72.7
Masters	1527	5.2	77.9
Assocvoc	1289	4.4	82.4
11th	1067	3.7	86.0
Assocacdm	982	3.4	89.4
10th	848	2.9	92.3
Profschool	502	1.7	94.0
7th8th	499	1.7	95.7
9th	395	1.4	97.1
12th	365	1.3	98.3
Doctorate	328	1.1	99.5
5th6th	97	0.3	99.8
1st4th	46	0.2	99.9
Preschool	17	0.1	100.0
Total	29170	100.0	100.0

Figure 3: Education Frequency

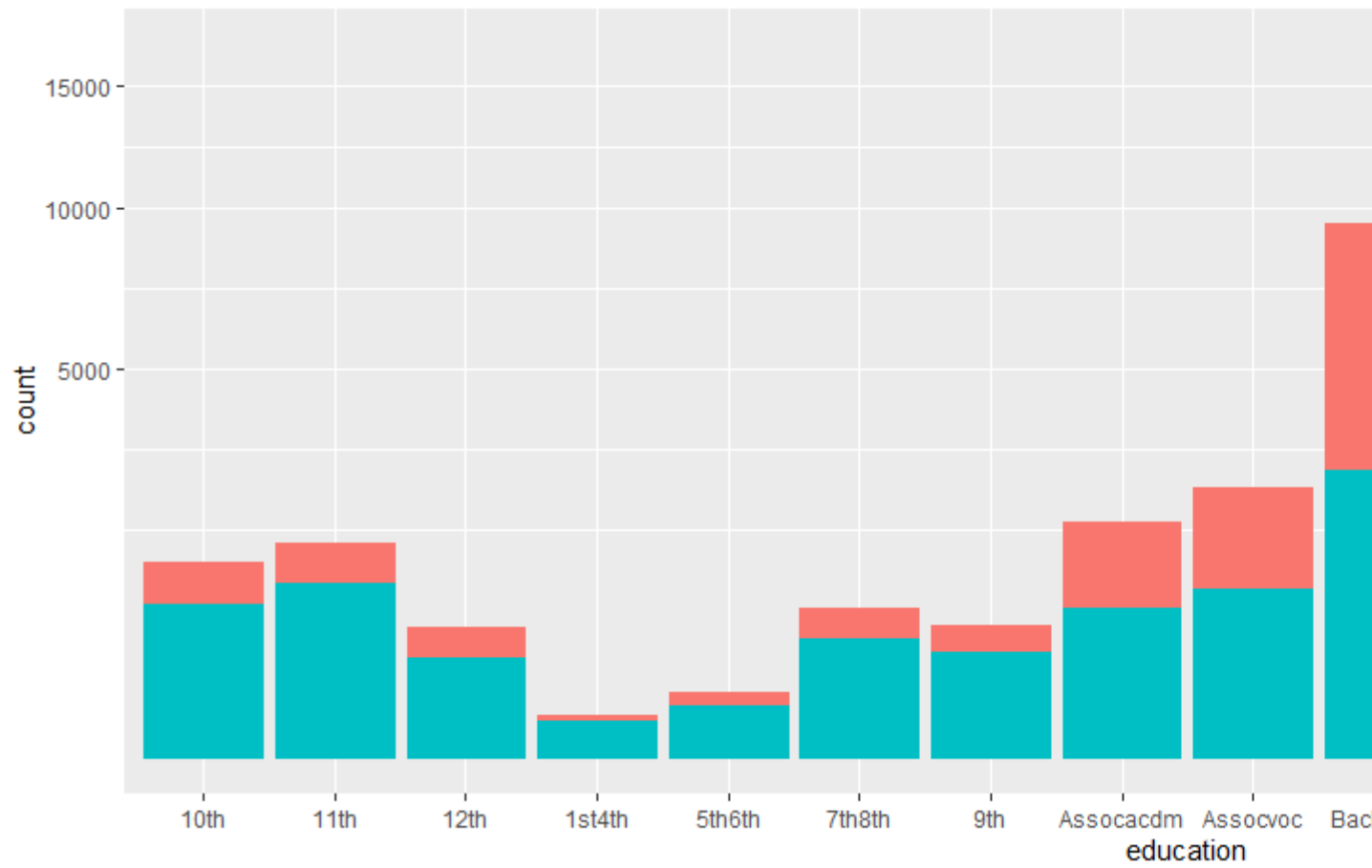


Figure 4: Education distribution

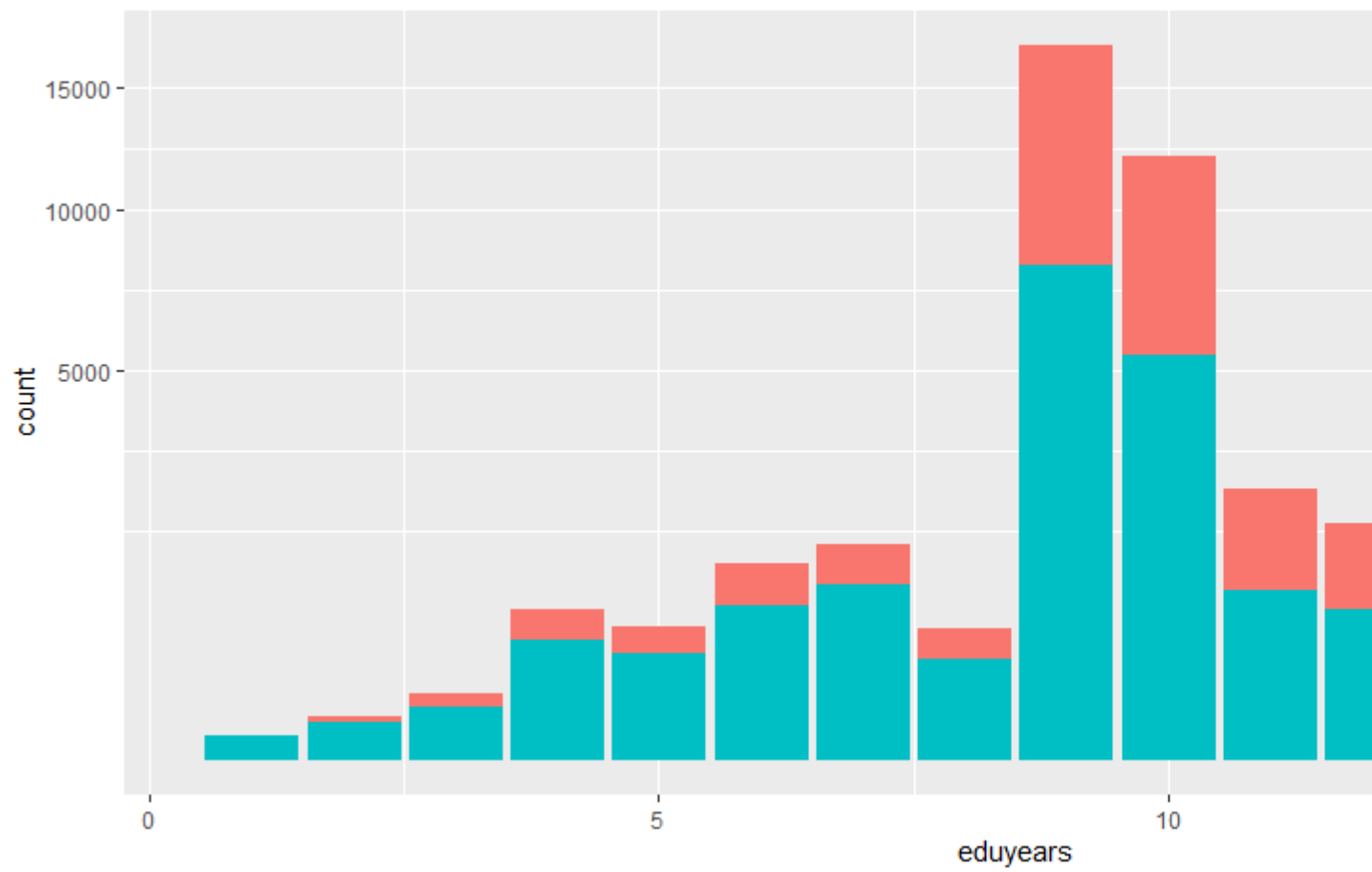


Figure 5: Education distribution

The number of high-school graduates are most frequent than other groups. There are 2 income distributions for education. One by the level of education and one by the number of years spent in education. Sure enough, most of the Above-50K earners have at least high-school education. Bachelors and some college have more above-50K earners than other categories. Interestingly enough, in comparison to Bachelors and some college masters have lesser Above-50K. Additionally, the group that has at least 12 years of education have more Above-50K than others. There are more at-or-below 50K earners in the high-school grad and less than 10 years of education categories. The rest of the groups have sparse data. This could be because of missing or bad data that was removed during initial cleaning before it was loaded into kaggle.

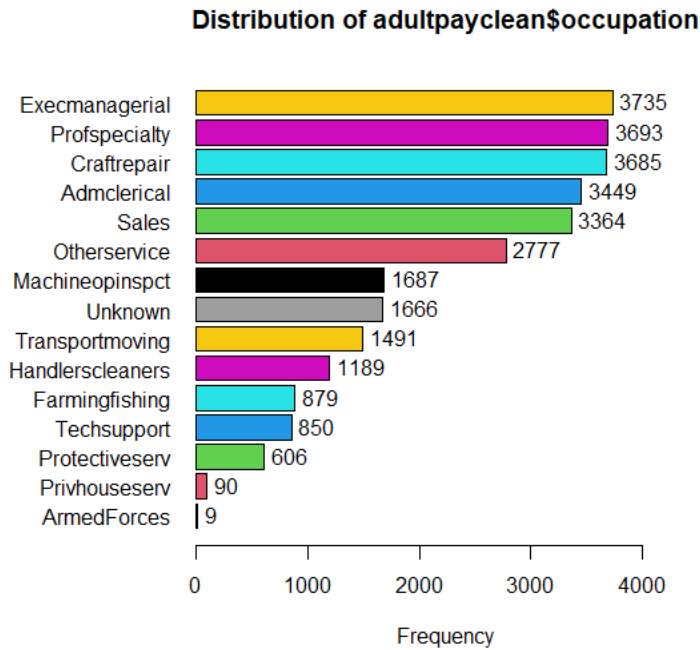


Figure 6: Occupation frequency

	Frequency	Percent	Cum. percent
Execmanagerial	3735	12.8	12.8
Profspecialty	3693	12.7	25.5
Craftrepair	3685	12.6	38.1
Admclerical	3449	11.8	49.9
Sales	3364	11.5	61.5
Otherservice	2777	9.5	71.0
Machineopinspct	1687	5.8	76.8
Unknown	1666	5.7	82.5
Transportmoving	1491	5.1	87.6
Handlerscleaners	1189	4.1	91.7
Farmingfishing	879	3.0	94.7
Techsupport	850	2.9	97.6
Protectiveserv	606	2.1	99.7
Privhouseserv	90	0.3	100.0
ArmedForces	9	0.0	100.0
Total	29170	100.0	100.0

Figure 7: Occupation frequency

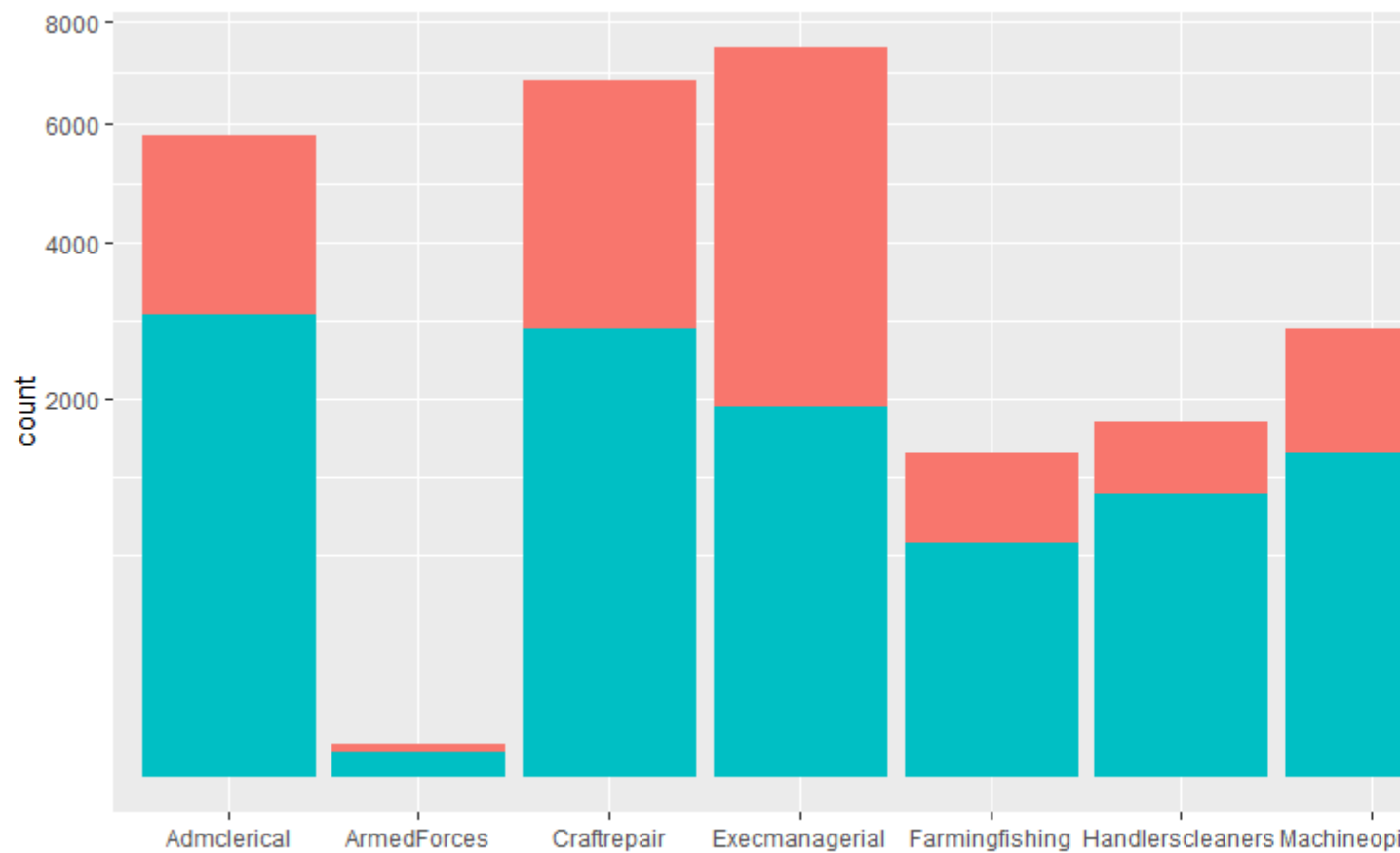


Figure 8: Occupation distribution

There are more executive and managerial occupations followed by professionals. The same groups earn more and are in the Above-50K category. There are fairly significant number of at-or-below-50K earners in the other categories. There are 1600 or so missing data. This makes up 5% of the cleaned dataset.

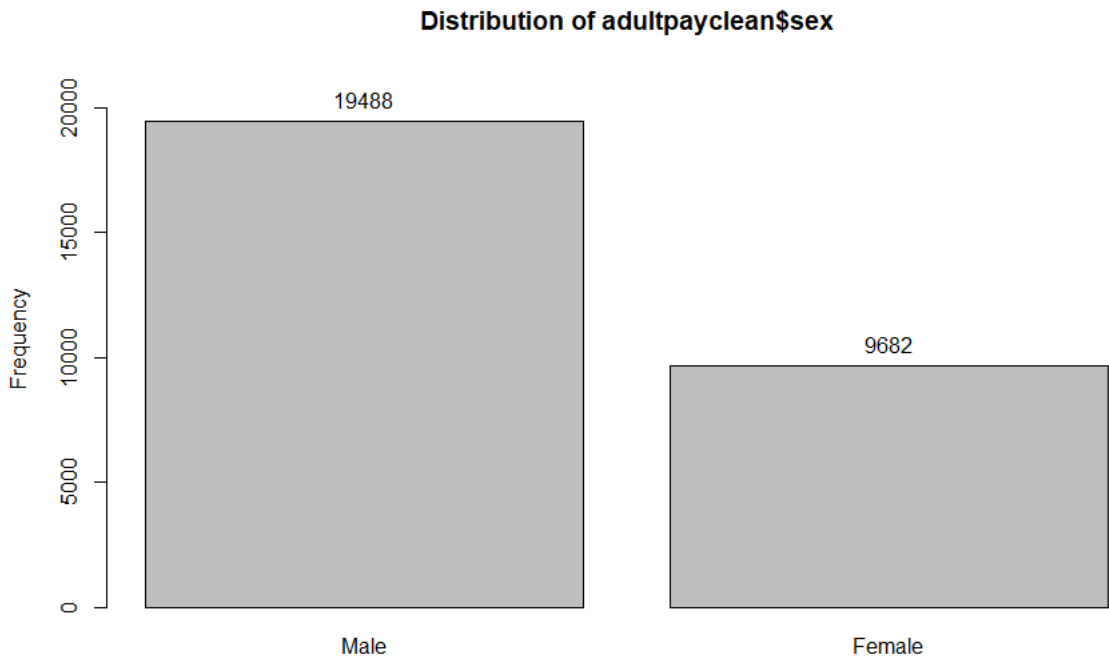


Figure 9: Sex feature characteristics

	Frequency	Percent	Cum. percent
Male	19488	66.8	66.8
Female	9682	33.2	100.0
Total	29170	100.0	100.0

Figure 10: Sex feature characteristics

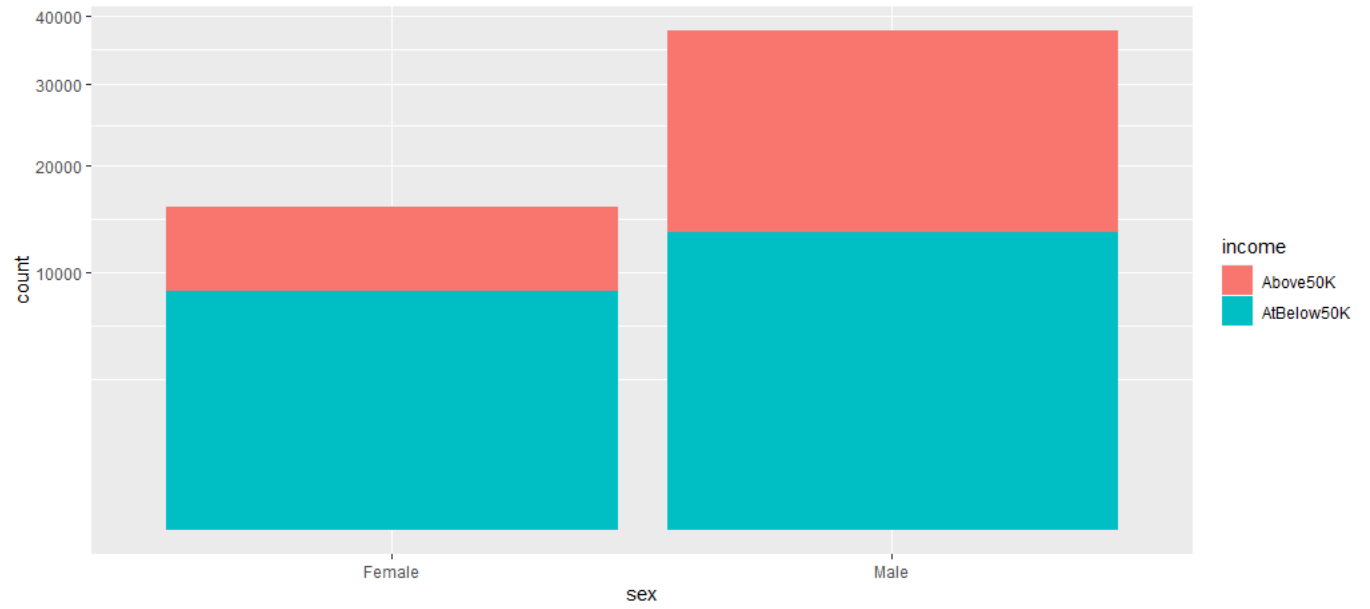


Figure 11: Sex feature characteristics

There are more males in the workforce than females. The distribution also indicates that the number of males at Above-50K are more than the number of females. Additionally, there are more males at-or-below-50K than females as well.

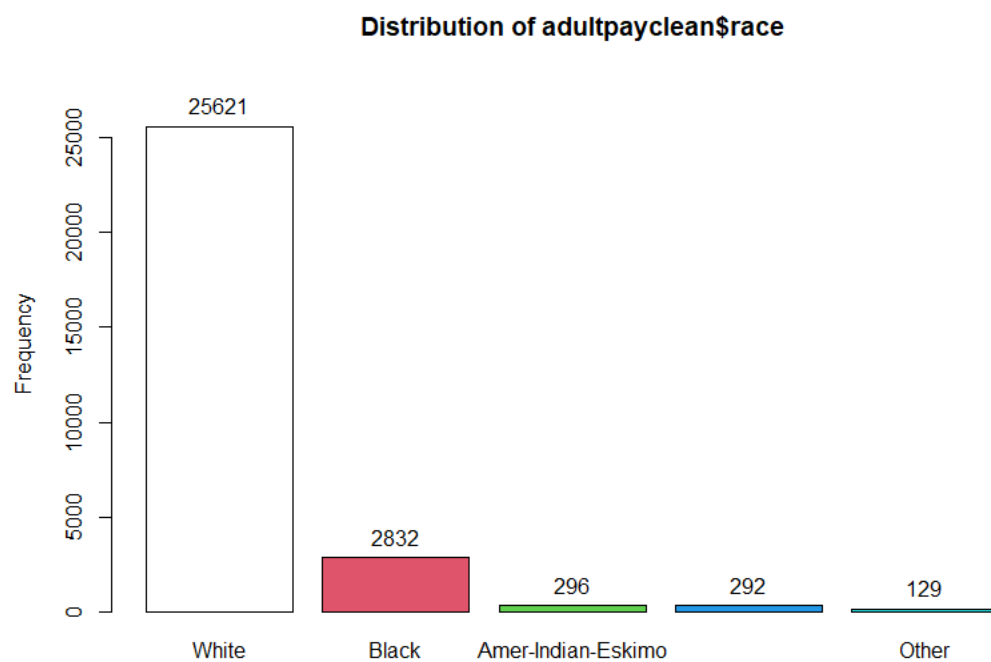


Figure 12: Race distribution

	Frequency	Percent	Cum. percent
White	25621	87.8	87.8
Black	2832	9.7	97.5
Amer-Indian-Eskimo	296	1.0	98.6
Asian-Pac-Islander	292	1.0	99.6
Other	129	0.4	100.0
Total	29170	100.0	100.0

Figure 13: Race distribution

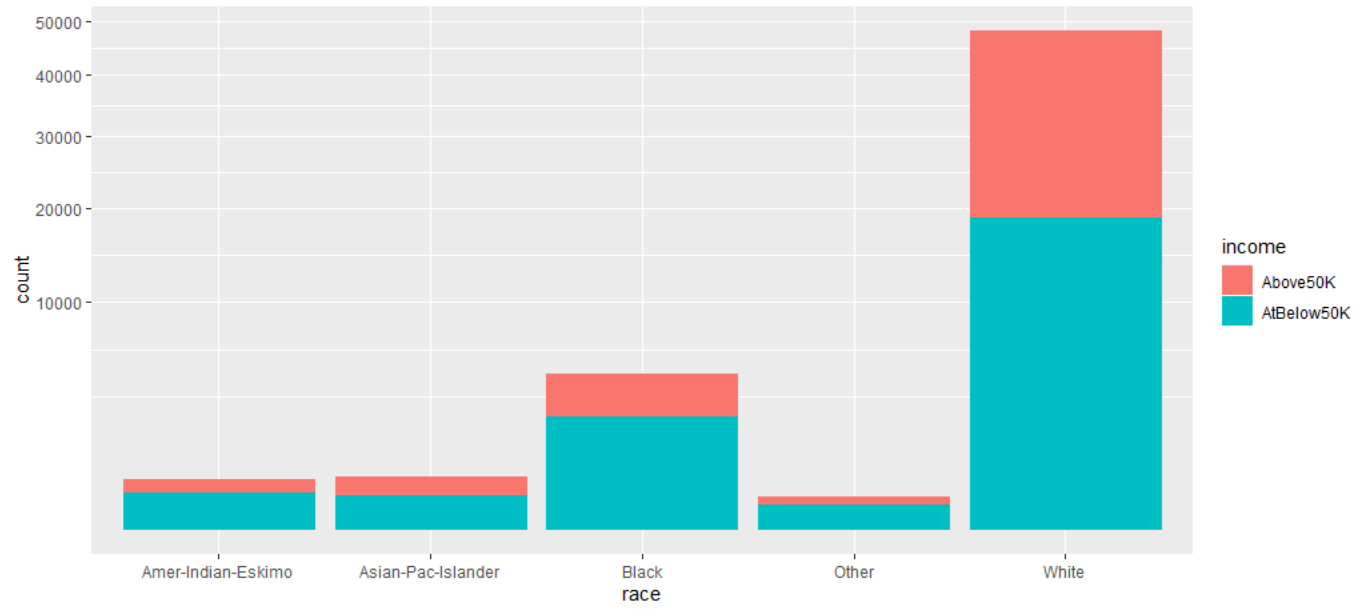


Figure 14: Race distribution

There are more whites in the labor force in both the Above-50K and at-or-below-50K categories. Blacks and other race follow next. The proportion of Above-50K earners are lesser in other races when compared to whites.

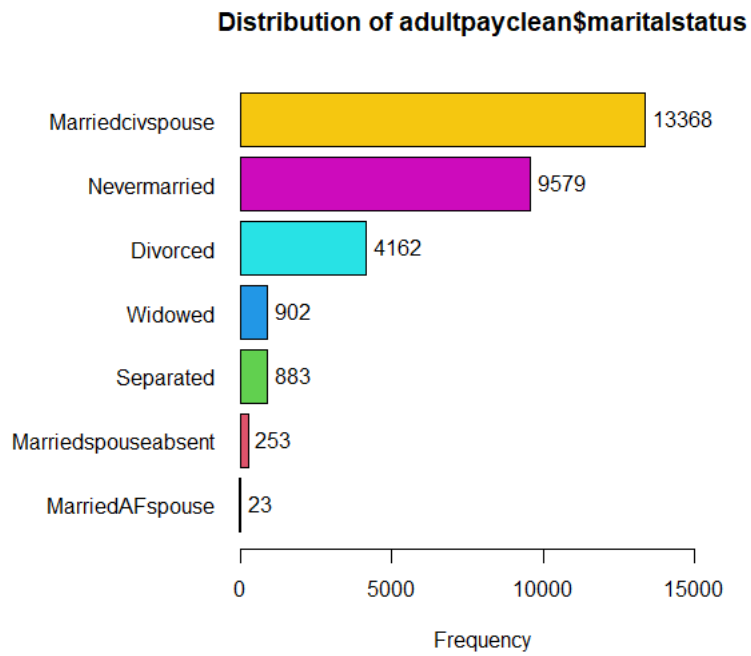


Figure 15: Marital Status frequency

	Frequency	Percent	Cum. percent
Marriedcivspouse	13368	45.8	45.8
Nevermarried	9579	32.8	78.7
Divorced	4162	14.3	92.9
Widowed	902	3.1	96.0
Separated	883	3.0	99.1
Marriedspouseabsent	253	0.9	99.9
MarriedAFspouse	23	0.1	100.0
Total	29170	100.0	100.0

Figure 16: Marital Status frequency

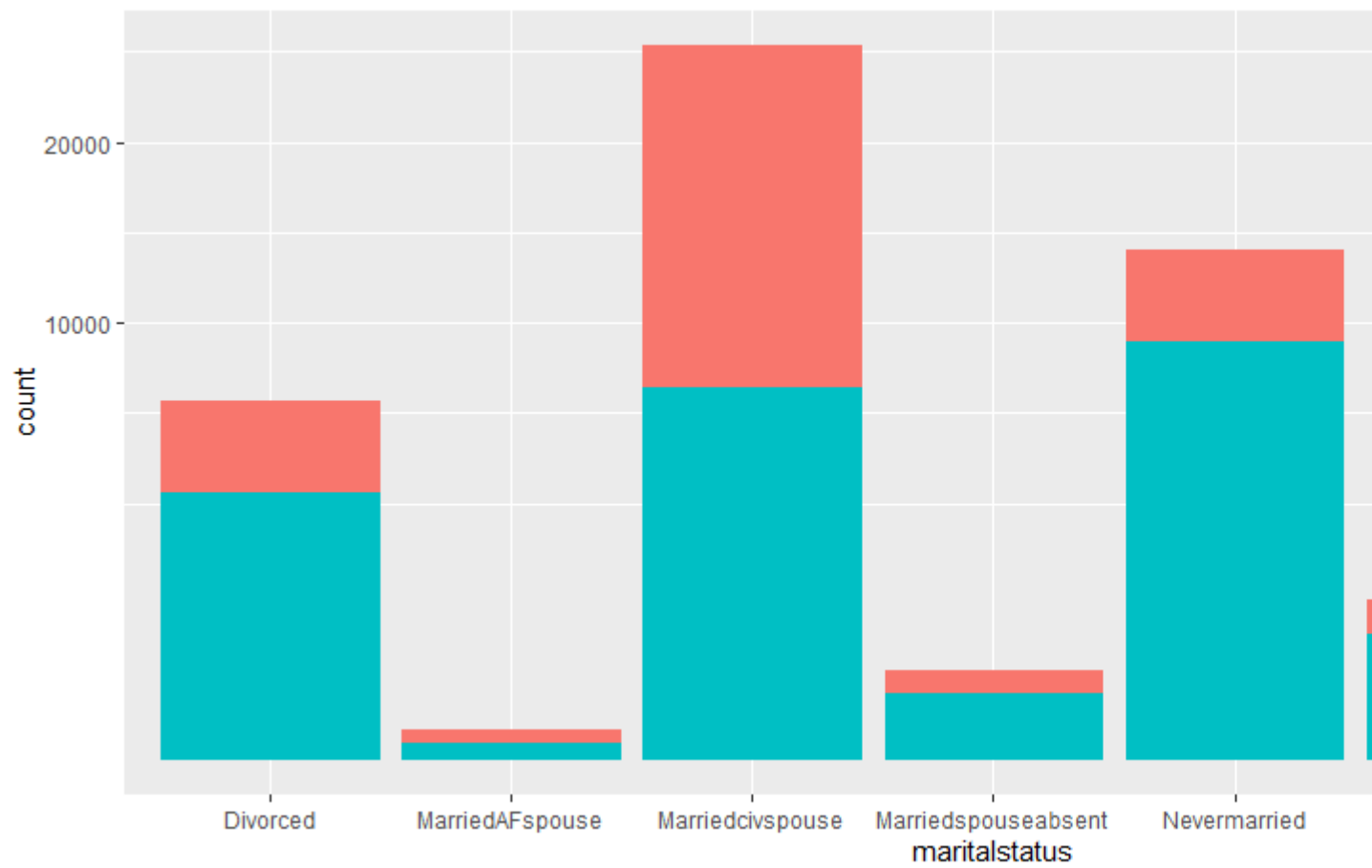


Figure 17: Marital Status distribution

The number of married with civilian spouses are more than any other groups in labor force. The number of Above-50K and at-or-below-50K earners are also more in this group when compared to others. The number of at-or-below-50K are more than the Above-50K in the singles group. This suggests most of this group is below 25 years age group. This can be correlated with the income by age group chart. The number of married with military spouses is less which may be due to lack of data.

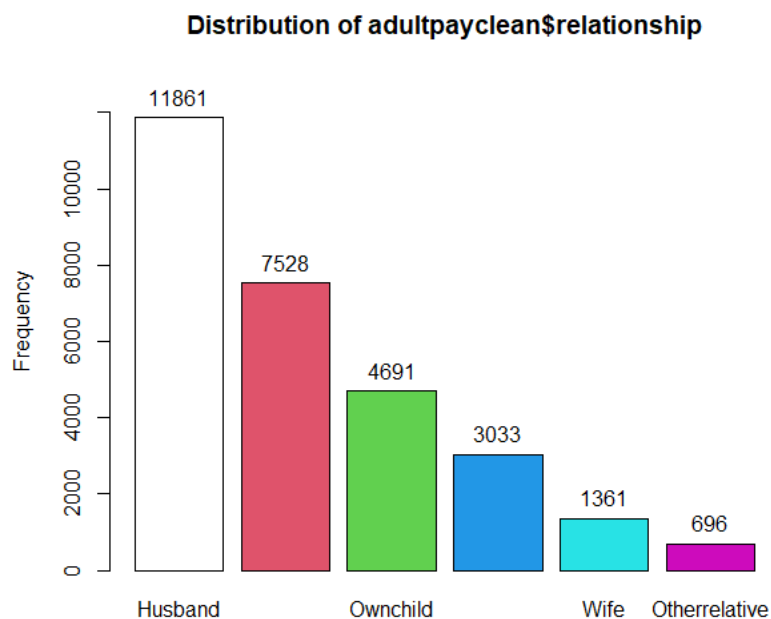


Figure 18: Relationship distribution

	Frequency	Percent	Cum. percent
Husband	11861	40.7	40.7
Notinfamily	7528	25.8	66.5
Ownchild	4691	16.1	82.6
Unmarried	3033	10.4	92.9
Wife	1361	4.7	97.6
Otherrelative	696	2.4	100.0
Total	29170	100.0	100.0

Figure 19: Relationship distribution

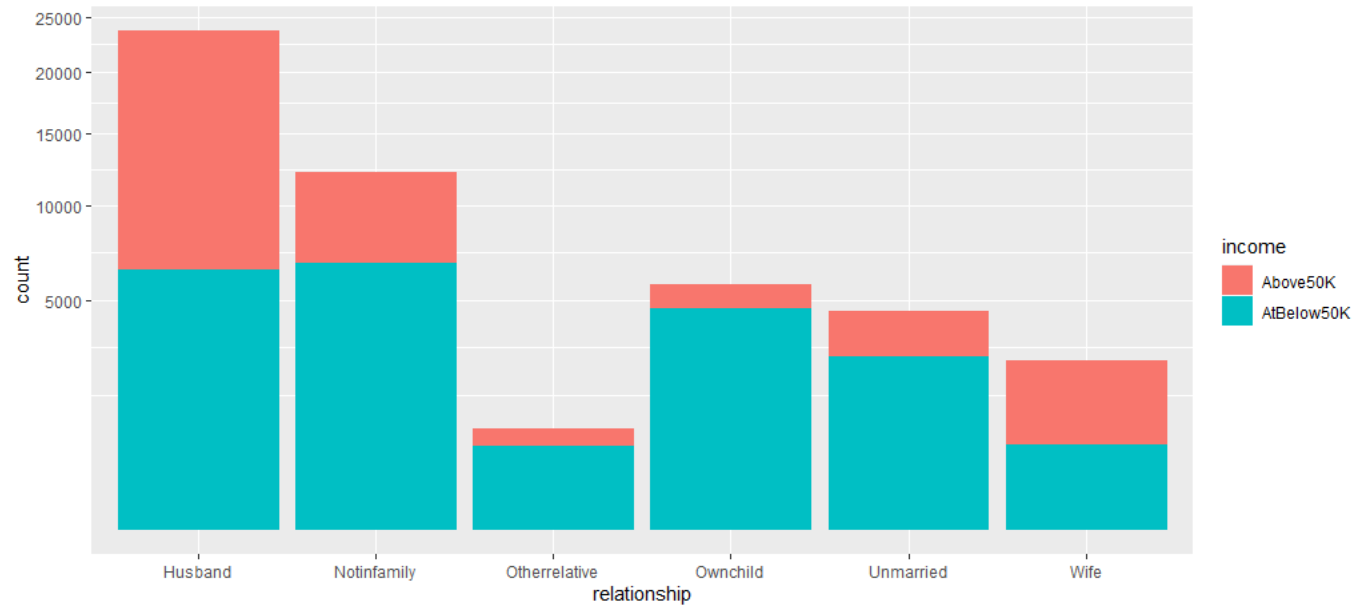


Figure 20: Relationship distribution

The number of husbands who earn more than 50K and at-or-below-50K are also more than the number of wives. Unmarried individuals also are more in the at-or-below-50K category.

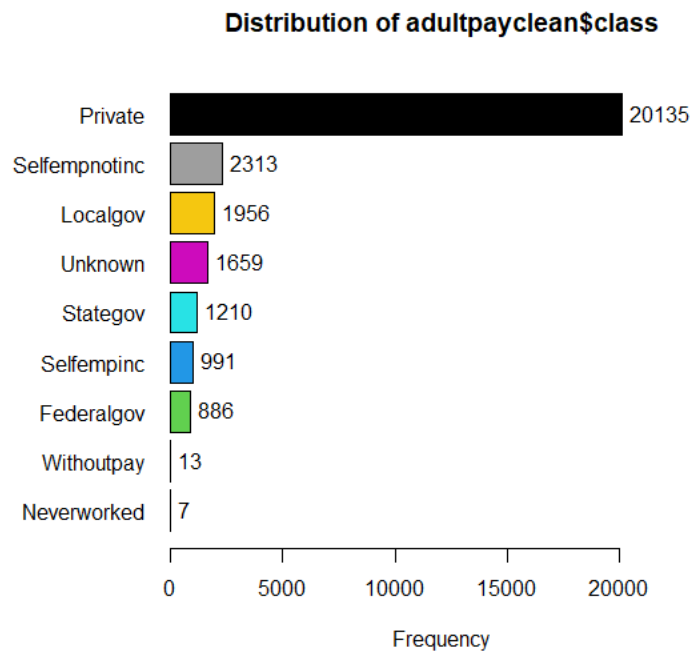


Figure 21: Class distribution

	Frequency	Percent	Cum. percent
Private	20135	69.0	69.0
Selfempnotinc	2313	7.9	77.0
Localgov	1956	6.7	83.7
Unknown	1659	5.7	89.3
Stategov	1210	4.1	93.5
Selfempinc	991	3.4	96.9
Federalgov	886	3.0	99.9
Withoutpay	13	0.0	100.0
Neverworked	7	0.0	100.0
Total	29170	100.0	100.0

Figure 22: Class distribution

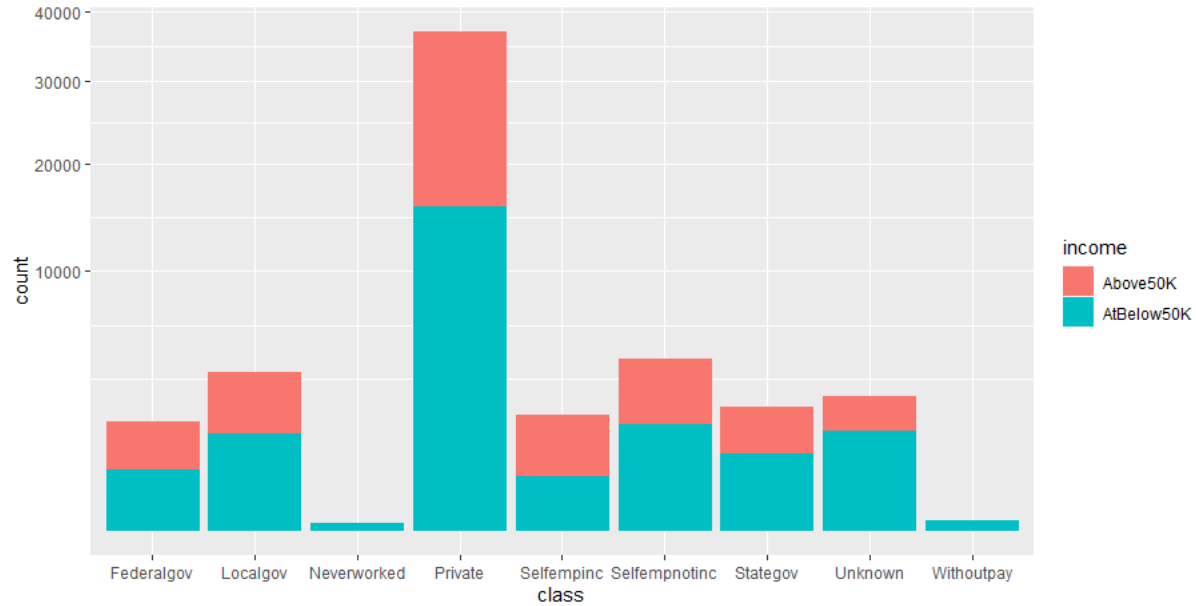


Figure 23: Class distribution

The number of private employer workers are more than other categories. They also earn more Above-50K and at-or-below-50K than other categories.

As shown above the occupation and class columns have more than 1600 Unknown data. Adding this to the predictor columns could skew the predictions and introduce errors. We will exclude this from our list of predictors. The final list of predictors/features used to predict the outcome income level are age, education, years of education (a.k.a. eduyears), maritalstatus, relationship, race, sex, and hoursperweek.

Process

The code `DatasetProcessingCode.R` downloads the dataset from a convenient location in github (originally downloaded from kaggle) and then unzips it and converts the dataset into a R data frame. It then does extensive cleaning, processing and munging of the data to make it tidy and meaningful for analysis. It then splits the original dataset into training and validation datasets. All analysis is done on this training dataset. Once the training is done the validation is done on the test dataset.

The `caret` package includes the function `createDataPartition` that helps us generate indexes for randomly splitting the data into training and test sets. The argument `times` is used to define how many random samples of indexes to return, the argument `p` is used to define what proportion of the data is represented by the index, and the argument `list` is used to decide if we want the indexes returned as a list or not. We can use the `test_index` from the `createDataPartition` function call to define the training and test sets.

```
# set.seed(1, sample.kind = 'Rounding')

# test_index <-

# createDataPartition( y = adu1tpayclean$income, times = 1, p = 0.1, list =
# FALSE )

# adu1tpayclean_train <- adu1tpayclean[-test_index, ]

# adu1tpayclean_validation <- adu1tpayclean[test_index, ]

# dim(adu1tpayclean)

# dim(adu1tpayclean_train)

# dim(adu1tpayclean_validation)

# > dim(adu1tpayclean)

# [1] 29170 13

# > dim(adu1tpayclean_train)

# [1] 26252 13

# > dim(adu1tpayclean_validation)

# [1] 2918 13
```

We then develop an algorithm using only the training set. Once we are done developing the algorithm, we then freeze it and evaluate it using the test set.

The simplest way to evaluate the algorithm for the `adultpay` dataset is by simply reporting the proportion of cases that were correctly predicted in the test set. This metric is referred to as overall accuracy. Some times accuracy is skewed due to prevalence or bias of one class or the other, therefore we will weigh in on sensitivity and specificity (check true/false positives and true/false negatives) when choosing the final model.

We will also use F1 measure to validate the accuracy of the sensitivity and specificity and see if we have a good precision. F1 score is a combination of two important error metrics: Precision and Recall. Thus, it can be considered as the Harmonic mean of Precision and Recall error metrics for an imbalanced dataset with respect to binary classification of data. As income level has only 30% data for Above-50K class this score is

useful. The highest possible value of an F-score is 1.0, indicating perfect precision and recall, and the lowest possible value is 0. Higher the precision the better.

We will also use a widely used plot receiver operating characteristic (ROC) curve. Receiver Operating Characteristic curves are a popular way to visualize the tradeoffs between sensitivity and specificity in a binary classifier. The ROC curve plots sensitivity (TPR) versus 1 - specificity or the false positive rate (FPR). This also gives us Area under the curve (a.k.a. AUC). The AUC summarizes ROC in a single value. The probabilistic interpretation is that if you randomly choose a positive case and a negative case, the probability that the positive case outranks the negative case according to the classifier is given by the AUC. Higher the AUC score the better. This indicates that the model has higher true positive rate and low false positive rate.

Model Creation

Simplest model using random sampling

We are now going to evaluate various algorithms progressively

We first start with a simple model. We sample randomly for the desired outcome. We then compare it with the actual outcomes and take a mean of the results. The result is 50% which is expected for guessing a binary outcome. Its akin to tossing a coin and getting head or tail. The chances are 50/50.

```
# > seat_of_the_pants <- sample(c('Above50K', 'AtBelow50K'),
# length(test_index), replace = TRUE)

# %>% factor(levels = levels(adultpayclean_validation$income))

# > mean(seat_of_the_pants == adultpayclean_validation$income)

# [1] 0.5010281

# build a confusion matrix for this simple model

# cm <- confusionMatrix(data =seat_of_the_pants , reference =
# adultpayclean_validation$income)

## Confusion Matrix and Statistics Reference

## Prediction Above50K AtBelow50K

## Above50K 347 1087

## AtBelow50K 371 1113 Accuracy : 0.5003

## 95% CI : (0.482, 0.5186)

## No Information Rate : 0.7539

## P-Value [Acc > NIR] : 1 Kappa : -0.0081 McNemar's Test P-Value : <2e-16
## Sensitivity : 0.4833

## Specificity : 0.5059

## Pos Pred Value : 0.2420
```



```

## Neg Pred Value : 0.7500 Prevalence : 0.2461

## Detection Rate : 0.1189

## Detection Prevalence : 0.4914

## Balanced Accuracy : 0.4946 'Positive' Class : Above50K

# p <- 0.1

# n <- length(test_index)

# y_hat <- sample(c('Above50K', 'AtBelow50K'), n, replace = TRUE, prob=c(p,
# 1-p)) %>%

# factor(levels = levels(adultpayclean_validation$income))

# mean(y_hat == adultpayclean_validation$income)

# [1] 0.7076765

# p <- 0.9

# n <- length(test_index)

# y_hat <- sample(c('Above50K', 'AtBelow50K'), n, replace = TRUE, prob=c(p,
# 1-p)) %>%

# factor(levels = levels(adultpayclean_validation$income))

# mean(y_hat == adultpayclean_validation$income)

# [1] 0.2964359

# > f1_guess

# [1] 0.3224907

# Area under the curve

# > auc_guess <- auc(ifelse(adultpayclean_validation$income == 'Above50K',1,2),
# ifelse(seat_of_the_pants == 'Above50K',1,2))

# > auc_guess

# Area under the curve: 0.4946

```

We then construct the confusion matrix, which basically tabulates each combination of prediction and actual value. We see that the above-50K and at-or-below-50K are almost evenly distributed with a slightly higher prevalence in the at-or-below-50K income level.

We can verify this by adjusting the probability of our sampling to skew towards above-50K vs at-or-below-

50K and vice-versa.

Prevalence can result in skewed results. We will keep an eye on the other metrics like sensitivity and specificity in addition to accuracy. In this case low prevalence matches with the expected accuracy.

The F1 score for this guessing is 0.3224 and Area under the curve is 0.4946. This is in line with our expectations for mere guessing of the outcome.

Our goal is to improve the accuracy > 80% while keeping sensitivity and specificity under check. Hence we further analyze the impact of other features on the income levels.

Logistic regression using linear models

We will start with a simple logistic model - linear model. We use the features age, eduyears, sex, race, maritalstatus, relationship and hoursperweek.

Both linear and logistic regressions provide an estimate for the conditional expectation:

$$E(Y \mid X = x)$$

which in the case of binary data is equivalent to the conditional probability:

$$Pr(Y = 1 \mid X = x)$$

We can use this to arrive at `y_hat_logit`. Since we have 7 features $X = x$ is more like $X_i = x_i$ where i is from 1 to 7.

```
# lm_fit <- adultpayclean_train %>%  
  
# mutate(y = as.numeric(income == 'Above50K')) %>%  
  
# lm(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +  
# relationship, data=.) p_hat_logit <- predict(lm_fit, newdata =  
# adultpayclean_validation) y_hat_logit <- ifelse(p_hat_logit > 0.5,  
# 'Above50K', 'AtBelow50K') %>% factor accuracy_lm <-  
# confusionMatrix(y_hat_logit, adultpayclean_validation$income)$overall[['Accuracy']]  
  
# accuracy_lm  
  
# [1] 0.8167  
  
## Confusion Matrix and Statistics Reference  
  
## Prediction Above50K AtBelow50K  
  
## Above50K 344 161  
  
## AtBelow50K 374 2039 Accuracy : 0.8167  
  
## 95% CI : (0.8021, 0.8305)  
  
## No Information Rate : 0.7539
```

```
## P-Value [Acc > NIR] : 2.788e-16 Kappa : 0.451 McNemar's Test P-Value : <
## 2.2e-16 Sensitivity : 0.4791

## Specificity : 0.9268

## Pos Pred Value : 0.6812

## Neg Pred Value : 0.8450

## Prevalence : 0.2461

## Detection Rate : 0.1179

## Detection Prevalence : 0.1731

## Balanced Accuracy : 0.7030 'Positive' Class : Above50K

# > f1_lm

# [1] 0.5625511

# Area under the curve: 0.703
```

The accuracy for this model is 0.8167 however, sensitivity is 0.4791 and specificity is 0.9268. This indicates this model has higher ratio of negative outcomes than positive outcomes. It does have better accuracy, confidence interval, and F1 score than plain old guessing. Prevalence remains the same. AUC has improved as well. Lets see if we can do better

We will now experiment with the general linear model (glm). We will be consistent with the features used across all algorithms.

```
# glm_fit <- adultpayclean_train %>%

# mutate(y = as.numeric(income == 'Above50K')) %>%

# glm(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
# relationship, data=., family = 'binomial')

# p_hat_logit <- predict(glm_fit, newdata = adultpayclean_validation)

# y_hat_logit <- ifelse(p_hat_logit > 0.5, 'Above50K', 'AtBelow50K') %>% factor

# accuracy_glm <- confusionMatrix(y_hat_logit,
# adultpayclean_validation$income)$overall[['Accuracy']]

# > accuracy_glm [1] 0.8098

## Confusion Matrix and Statistics Reference

## Prediction Above50K AtBelow50K

## Above50K 279 116

## AtBelow50K 439 2084 Accuracy : 0.8098
```

```
## 95% CI : (0.7951, 0.8239)

## No Information Rate : 0.7539

## P-Value [Acc > NIR] : 3.442e-13 Kappa : 0.3958 McNemar's Test P-Value : <
## 2.2e-16 Sensitivity : 0.38858

## Specificity : 0.94727

## Pos Pred Value : 0.70633

## Neg Pred Value : 0.82600

## Prevalence : 0.24606

## Detection Rate : 0.09561

## Detection Prevalence : 0.13537

## Balanced Accuracy : 0.66793 'Positive' Class : Above50K

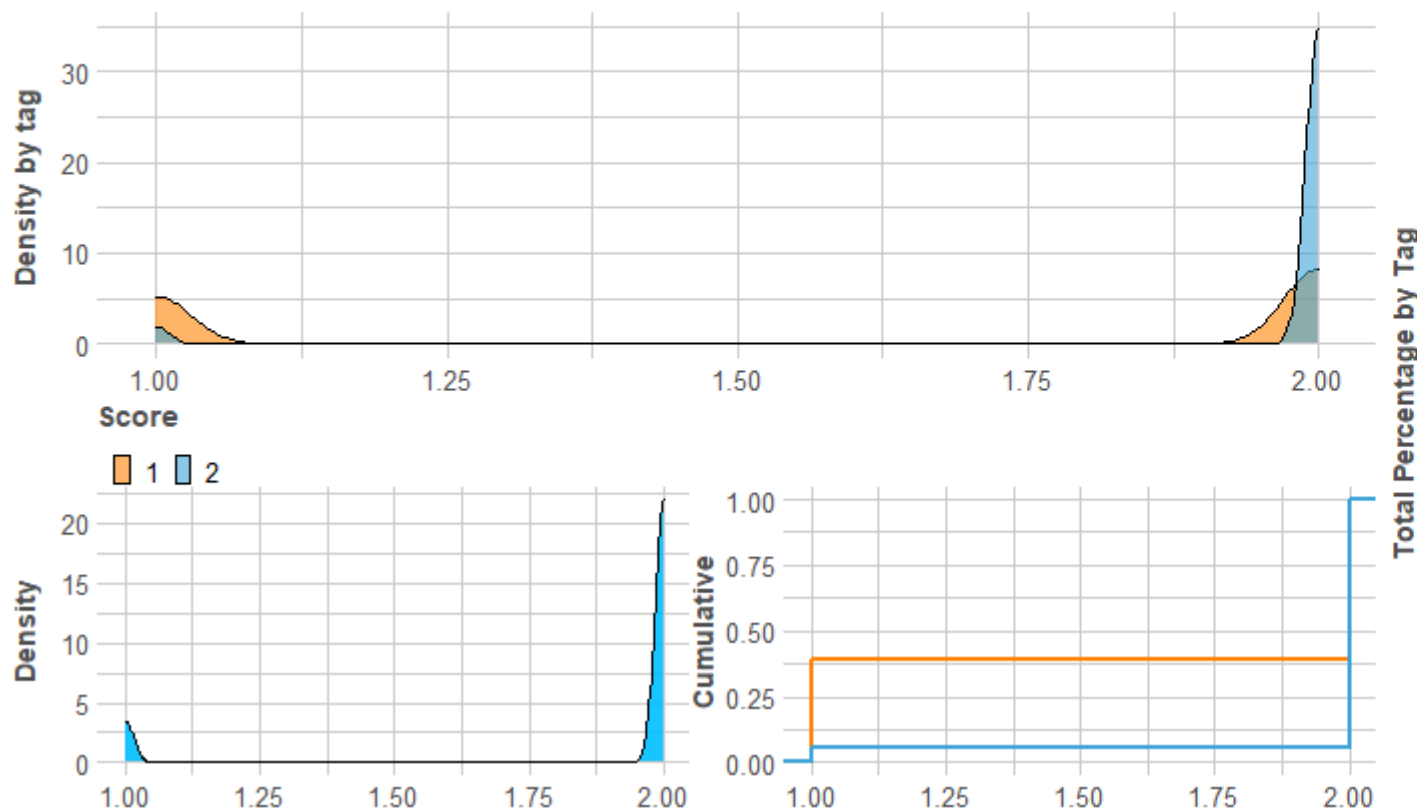
# f1_glm

# [1] 0.5013477

# Area under the curve: 0.6679
```

glm produces accuracy of 0.8098, however, specificity is still higher than sensitivity and prevalence about the same. CI and F1 scores can be better. AUC can also do better. We can do better!

Classification Model Results



Score cuts (10 quantiles)

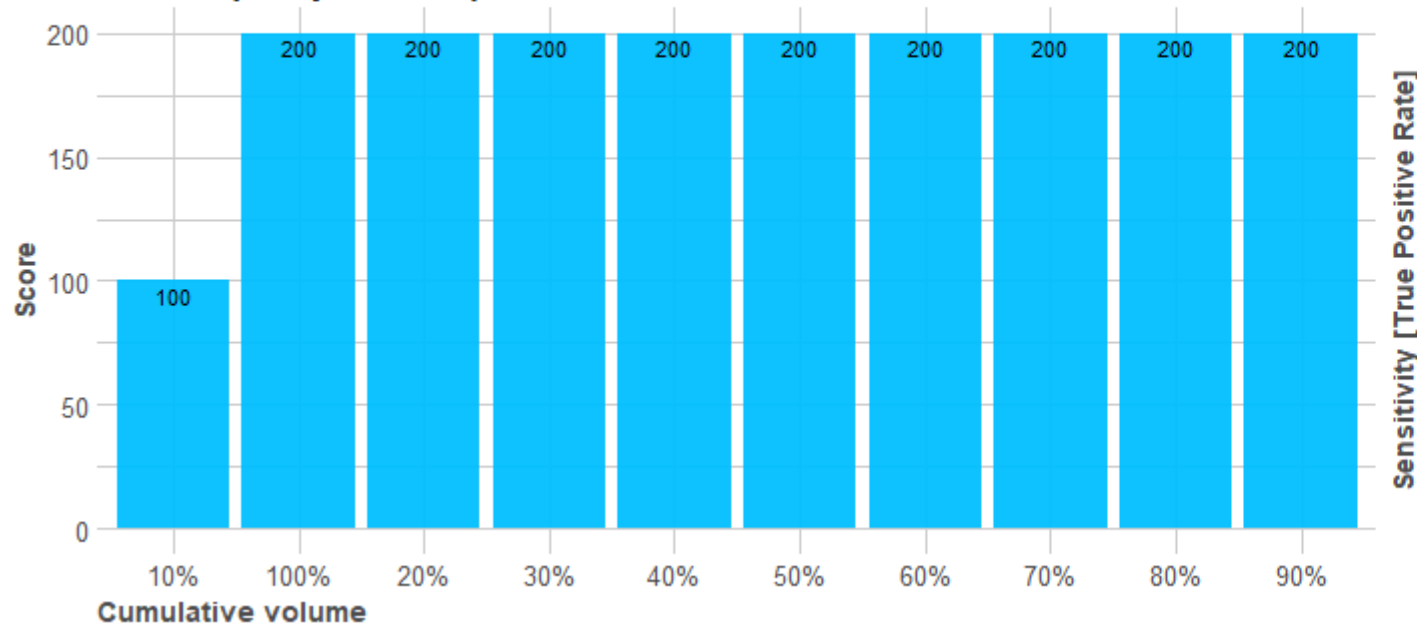


Figure 24: GLM classification summary

naive bayes

We will now experiment with generative models like naive bayes classification. It assumes the features that go into the model is independent of each other. That is changing the value of one feature, does not directly influence or change the value of any of the other features used in the algorithm. When we analyzed the dataset it didn't appear to have independent features. An example of dependent feature would be age and education years, age and marital status etc. However, let's see how naive bayes performs with this assumption.

```
## train_nb <- adultpayclean_train %>%

## mutate(y = as.factor(income == 'Above50K')) %>%

## naiveBayes(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
## relationship, data = .)

## y_hat_nb <- predict(train_nb, newdata = adultpayclean_validation)

## cm_tab <- table(adultpayclean_validation$income == 'Above50K', y_hat_nb)

## cm_nb <- confusionMatrix(cm_tab)

## Accuracy : 0.7718

## Confusion Matrix and Statistics y_hat_nb

## FALSE TRUE

## FALSE 1698 502

## TRUE 164 554 Accuracy : 0.7718

## 95% CI : (0.7561, 0.7869)

## No Information Rate : 0.6381

## P-Value [Acc > NIR] : < 2.2e-16 Kappa : 0.469 McNemar's Test P-Value : <
## 2.2e-16 Sensitivity : 0.9119

## Specificity : 0.5246

## Pos Pred Value : 0.7718

## Neg Pred Value : 0.7716

## Prevalence : 0.6381

## Detection Rate : 0.5819

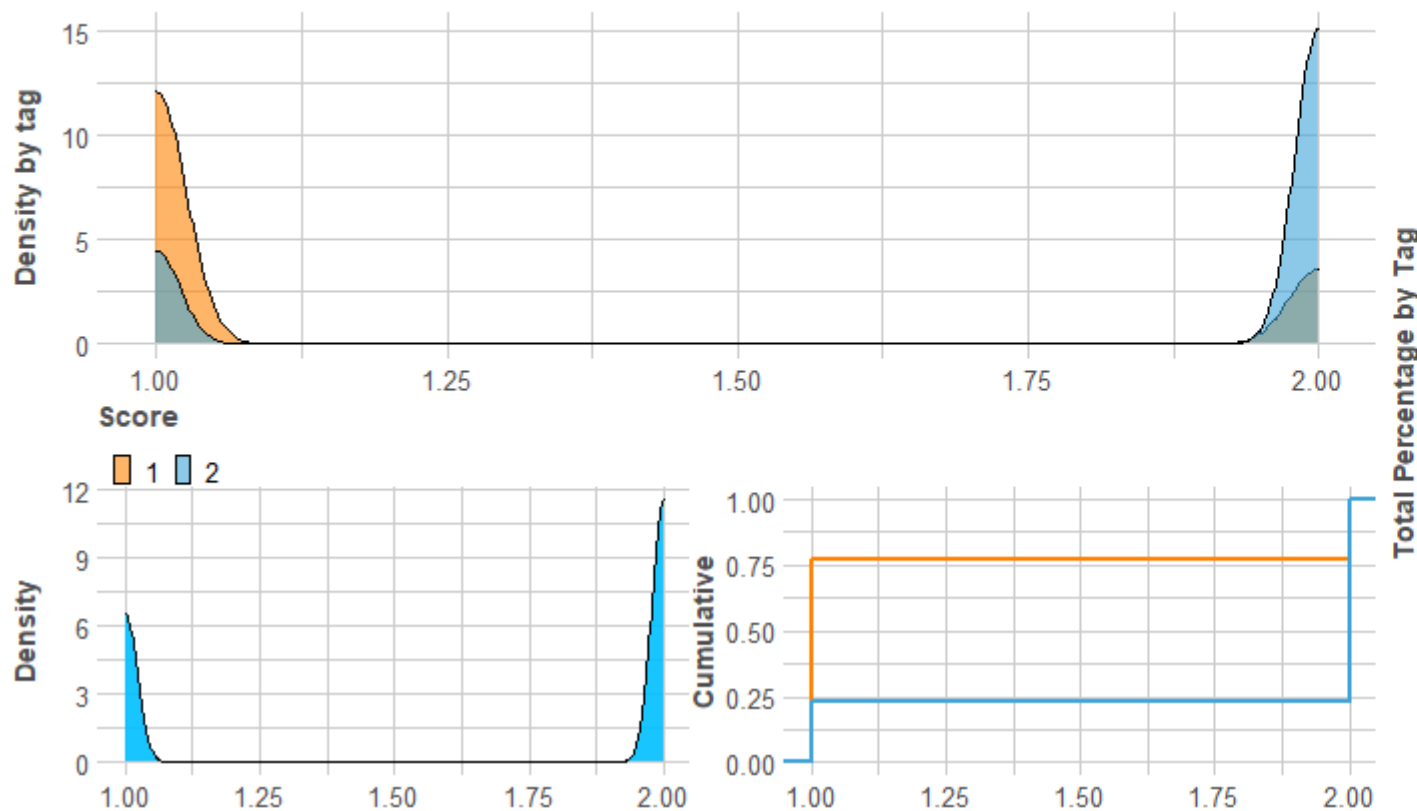
## Detection Prevalence : 0.7539

## Balanced Accuracy : 0.7183 'Positive' Class : FALSE
```

```
# > f1_nb  
# [1] 0.8360414  
# Area under the curve: 0.7717
```

Naive bayes gives us an accuracy of 0.7718 with a lower prevalence than other algorithms but higher than plain guessing. However the confidence interval is lower and accuracy of the prediction is lower too. There are more false negatives as well. F1 score and AUC has increased. This indicates there are more true positives as well. However, given that this is a imbalanced dataset, we can try other advanced algorithms and compare this later.

Classification Model Results



Score cuts (10 quantiles)

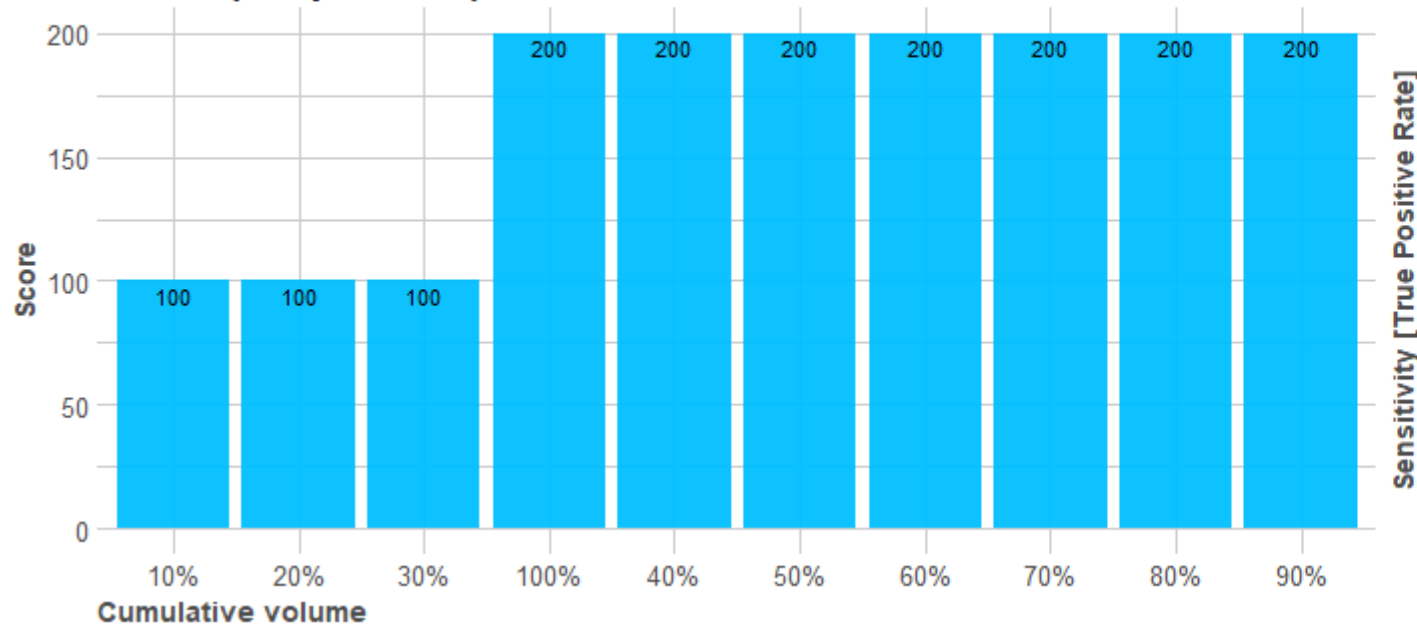


Figure 25: naive bayes classification summary

Lets explore k-nearest model now with same feature set. We will use cross-validation to tune the k parameter. By default, the cross validation is performed by taking 25 bootstrap samples comprised of 25% of the observations. For the kNN method, the default is to try k=5,7,9. We change this using the tuneGrid parameter. We will try the k values in the following sequence k = seq(3, 71, 2). Running this code will take several seconds. This is because when we run the algorithm, we will have to compute a distance between each observation in the test set and each observation in the training set. There are a lot of computations. Therefore, we use the trainControl function to make the code above go a bit faster by using, 10-fold cross validation. This means we have 10 samples using 10% of the observations each. We set the seed because cross validation is a random procedure and we want to make sure the result here is reproducible

k-nearest neighbors

```
# temp <- adultpayclean_train %>%  
  
# mutate(y = as.factor(income == 'Above50K'))  
  
# set.seed(2008)  
  
# control <- trainControl(method = 'cv', number = 10, p = .9)  
  
# train_knn <- train(y ~ age + eduyears + sex + race + hoursperweek +  
# maritalstatus + relationship, method = 'knn', data = temp, tuneGrid =  
# data.frame(k = seq(3, 71, 2)),trControl = control)  
  
# train_knn$bestTune  
  
# y_hat_knn <- predict(train_knn,adultpayclean_validation, type = 'raw')  
  
# accuracy_knn <- confusionMatrix(y_hat_knn,  
# as.factor(adultpayclean_validation$income ==  
# 'Above50K'))$overall[['Accuracy']]  
  
# ggplot(train_knn, highlight = TRUE)  
  
# > accuracy_knn [1] 0.8081
```

The k parameter that lead to maximum accuracy can be obtained by bestTune. The plot for that is shown in the figure “KNN tuning”.

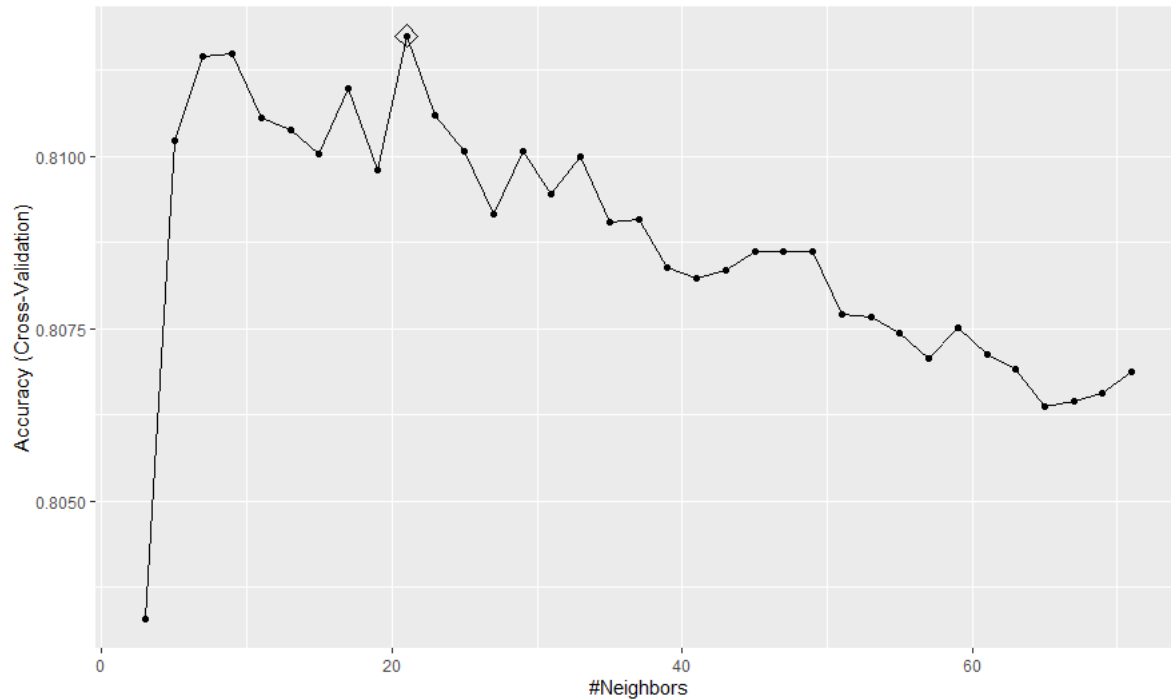


Figure 26: KNN tuning

Here is the confusion matrix for the knn tuned raw model

```
## cm_knn

## Confusion Matrix and Statistics Reference

## Prediction FALSE TRUE

## FALSE 1988 348

## TRUE 212 370 Accuracy : 0.8081

## 95% CI : (0.7933, 0.8222)

## No Information Rate : 0.7539

## P-Value [Acc > NIR] : 1.777e-12 Kappa : 0.4475 McNemar's Test P-Value :
## 1.165e-08 Sensitivity : 0.9036

## Specificity : 0.5153

## Pos Pred Value : 0.8510

## Neg Pred Value : 0.6357

## Prevalence : 0.7539

## Detection Rate : 0.6813
```

```
## Detection Prevalence : 0.8005

## Balanced Accuracy : 0.7095 'Positive' Class : FALSE

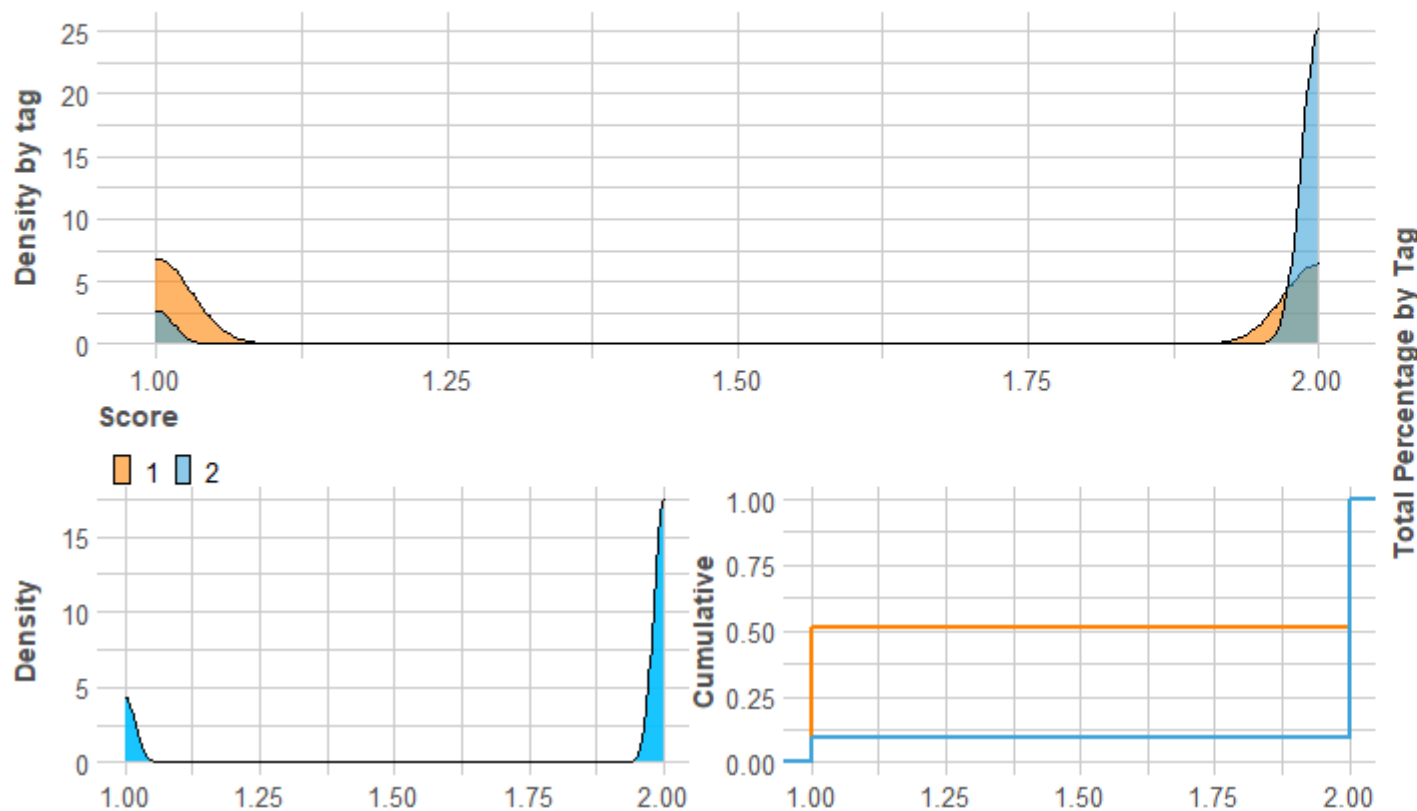
# f1_knn

# [1] 0.8765432

# Area under the curve: 0.7095
```

knn raw model produces same accuracy of 0.8081, however, specificity is lower than sensitivity and prevalence is now 0.753. We now see the positive prediction is better. F1 score has improved. However, the AUC score is less than naive bayes. This indicates that the number of true positives are lesser than before. Balanced accuracy is better. Can we do better with accuracy and precision?

Classification Model Results



Score cuts (10 quantiles)

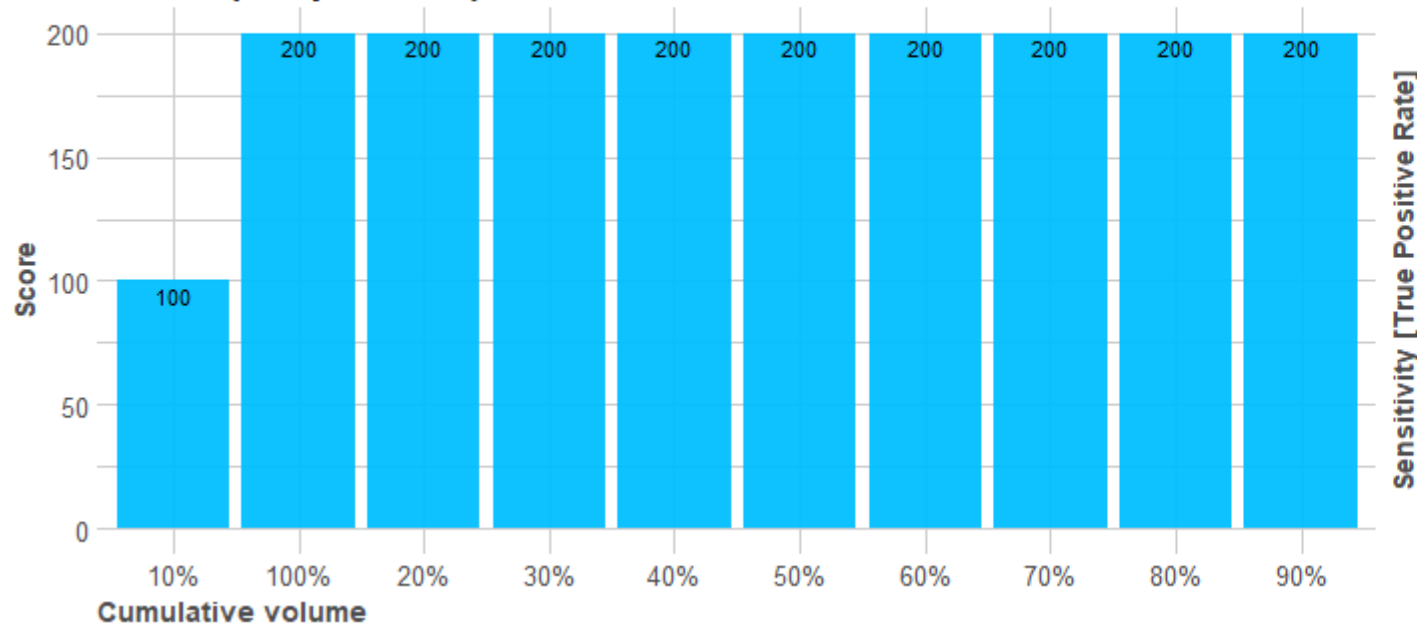


Figure 27: knn classification summary

We will now try to use classification with knn3. We will again use different values of k but using map_df function to repeat the above for each one. Running this classification model is going to be slow as it has to iterate through all the k values and find the one that is the highest.

```
# ks <- seq(3, 251, 2)

# knntune <- map_df(ks, function(k){

# temp <- adultpayclean_train %>%

# mutate(y = as.factor(income == 'Above50K'))

# temp_test <- adultpayclean_validation %>%

# mutate(y = as.factor(income == 'Above50K'))

# knn_fit <- knn3(y ~ age + eduyears + sex + race + hoursperweek +
# maritalstatus + relationship, data = temp, k = k)

# y_hat <- predict(knn_fit, temp, type = 'class')

# cm_train <- confusionMatrix(y_hat, temp$y)

# train_error <- cm_train$overall['Accuracy']

# y_hat <- predict(knn_fit, temp_test, type = 'class')

# cm_test <- confusionMatrix(y_hat, temp_test$y)

# test_error <- cm_test$overall['Accuracy']

# tibble(train = train_error, test = test_error) })

# accuracy_knntune <- max(knntune$test)

# get the confusion matrix for that k

# knn_fit <-

# knn3(

# y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
# relationship,

# data = temp,

# k = ks[which.max(knntune$test)]

# )

# y_hat <- predict(knn_fit, adultpayclean_validation, type = 'class')

# cm_knntune <- confusionMatrix(y_hat,
# as.factor(adultpayclean_validation$income == 'Above50K'))
```

```

## cm_knntune

## Confusion Matrix and Statistics Reference

## Prediction FALSE TRUE

## FALSE 1987 342

## TRUE 213 376 Accuracy : 0.8098

## 95% CI : (0.7951, 0.8239)

## No Information Rate : 0.7539

## P-Value [Acc > NIR] : 3.442e-13 Kappa : 0.4544 McNemar's Test P-Value :
## 5.532e-08 Sensitivity : 0.9032

## Specificity : 0.5237

## Pos Pred Value : 0.8532

## Neg Pred Value : 0.6384

## Prevalence : 0.7539

## Detection Rate : 0.6809

## Detection Prevalence : 0.7981

## Balanced Accuracy : 0.7134 'Positive' Class : FALSE

# > f1_knntune

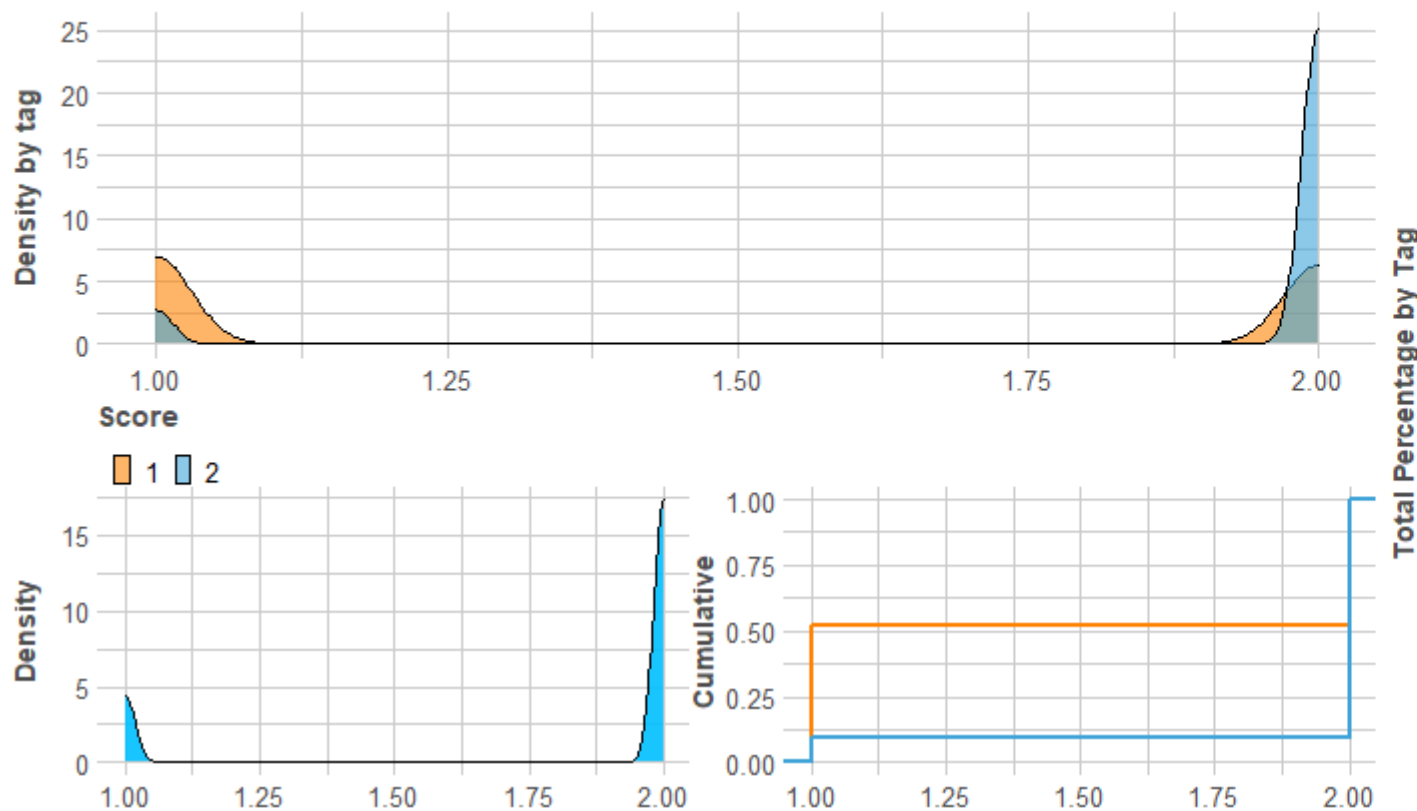
# [1] 0.8774564

# Area under the curve: 0.7134

```

The accuracy of the knn tuned classification is 0.8098 with a sensitivity of 0.9032 and specificity of 0.5237. This is better than the previous models. Prevalence is just about the same as raw knn. Balanced accuracy has improved. F1 score is also at ~ 88% and AUC has improved a bit.

Classification Model Results



Score cuts (10 quantiles)

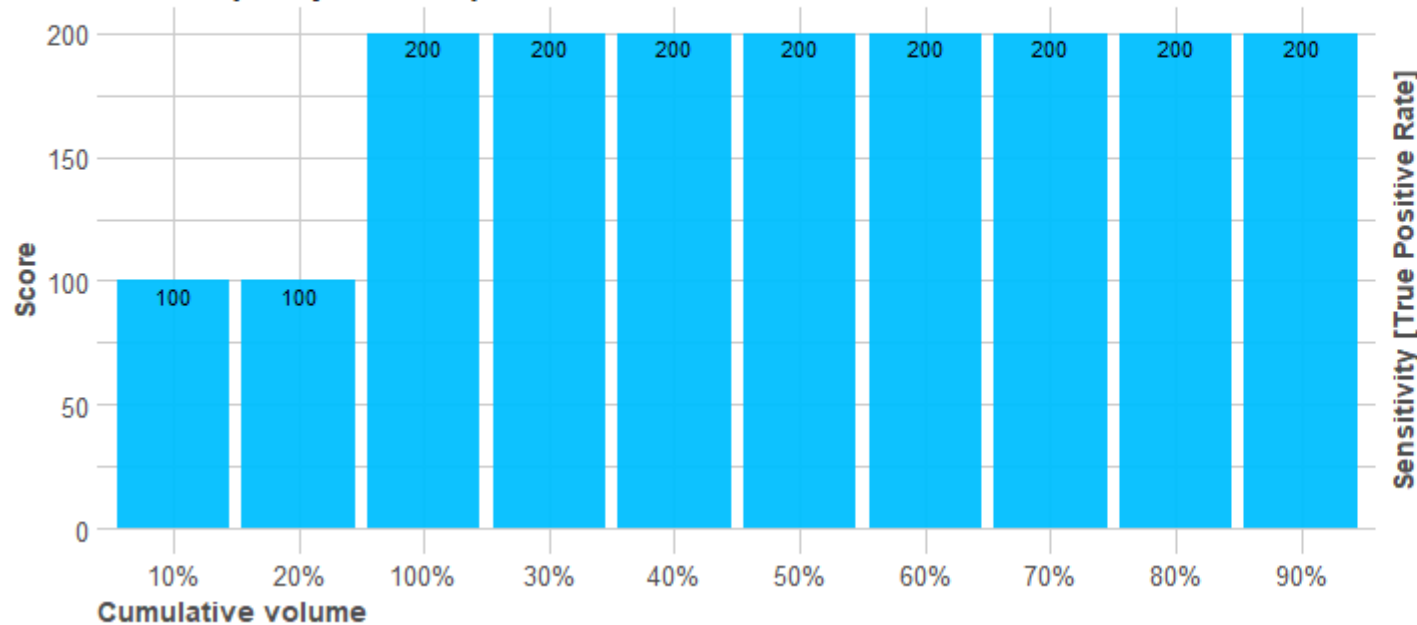


Figure 28: knn tune classification summary

Recursive partitioning with rpart

Next we use classification trees, or decision trees. We use the recursive partitioning library rpart for this.

```
# train_rpart <- train(y ~ age + eduyears + sex + race + hoursperweek +
# maritalstatus + relationship, method = 'rpart', tuneGrid = data.frame(cp =
# seq(0.0, 0.1, len = 25)), data = temp)

# y_hat <- predict(train_rpart, adu1tpayclean_validation)

# accuracy_rpart <- confusionMatrix(y_hat,
# as.factor(adu1tpayclean_validation$income == 'Above50K'))$overall['Accuracy']

# accuracy_rpart

# Accuracy

# 0.8211

## Confusion Matrix and Statistics Reference

## Prediction FALSE TRUE

## FALSE 2002 324

## TRUE 198 394 Accuracy : 0.8211

## 95% CI : (0.8067, 0.8349)

## No Information Rate : 0.7539

## P-Value [Acc > NIR] : < 2.2e-16 Kappa : 0.4876 McNemar's Test P-Value :
## 4.472e-08 Sensitivity : 0.9100

## Specificity : 0.5487

## Pos Pred Value : 0.8607

## Neg Pred Value : 0.6655

## Prevalence : 0.7539

## Detection Rate : 0.6861

## Detection Prevalence : 0.7971

## Balanced Accuracy : 0.7294 'Positive' Class : FALSE

# > f1_rpart

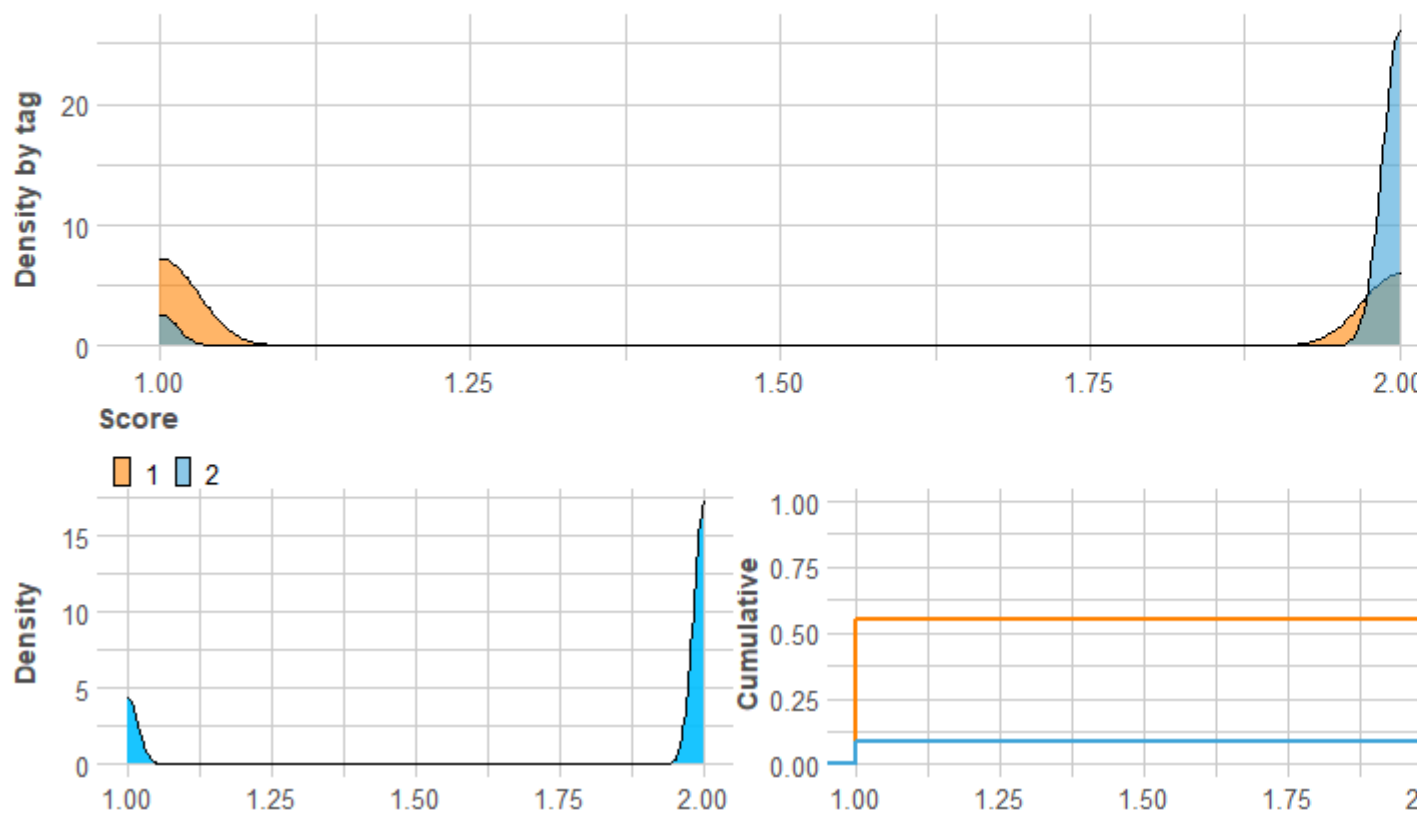
# [1] 0.8846664
```



```
# Area under the curve: 0.7294
```

The accuracy of recursive partitioning is 0.8211, sensitivity is higher at 0.91 and specificity is at 0.5487 with prevalence about the same. Balanced accuracy has improved. Confidence intervals for positive and negative has improved as well. The F1 score for rpart is 0.88467 and AUC is 0.7294 which is a positive sign.

Classification Model Results



Score cuts (10 quantiles)

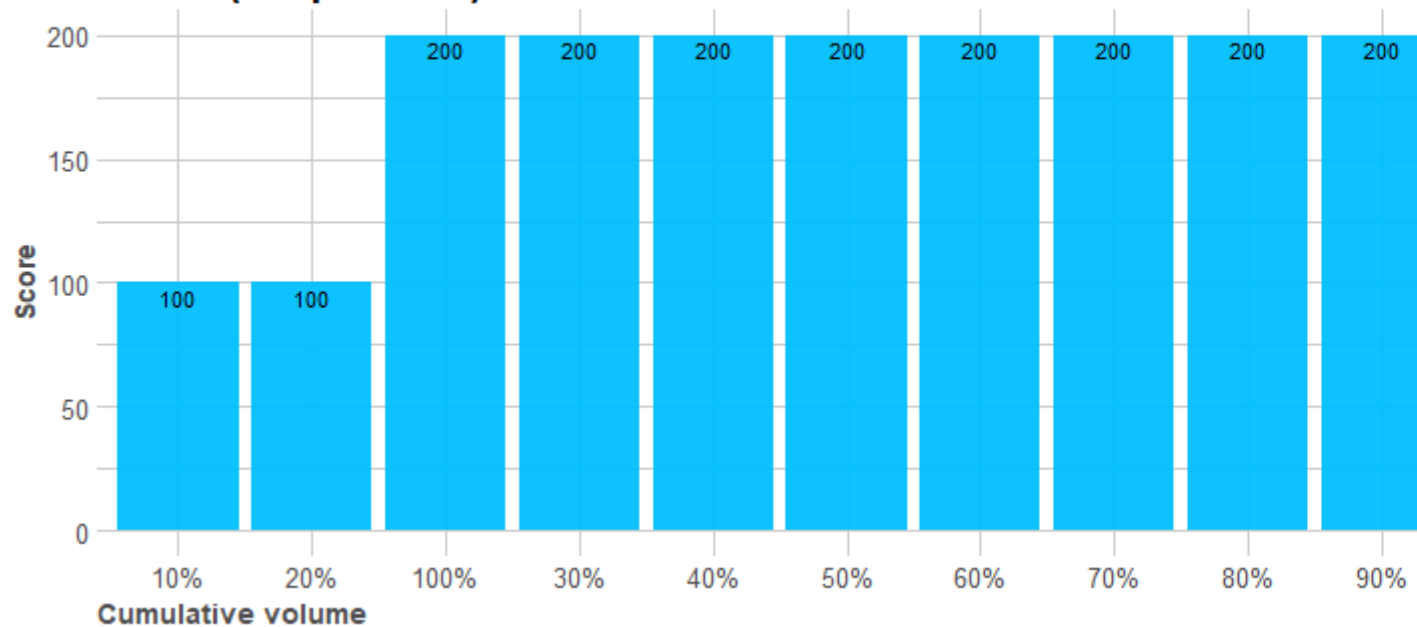


Figure 29: rpart classification summary

We will next look at random forest algorithm.

Random forests

Classification trees have certain advantages that make them very useful. They are highly interpretable, even more so than linear models. They are easy to visualize (if small enough). Finally, they can model human decision processes and don't require use of dummy predictors for categorical variables. On the other hand, the approach via recursive partitioning can easily over-train and is therefore a bit harder to train than, for example, linear regression or kNN. Furthermore, in terms of accuracy, it is rarely the best performing method since it is not very flexible and is highly unstable to changes in training data. Random forests, explained next, improve on several of these shortcomings.

The first step is bootstrap aggregation or bagging. The general idea is to generate many predictors, each using regression or classification trees, and then forming a final prediction based on the average prediction of all these trees.

```
# train_rf <- randomForest(y ~ age + eduyears + sex + race + hoursperweek +
# maritalstatus + relationship, data=temp)

# accuracy_rf <- confusionMatrix(predict(train_rf, adu1tpayclean_validation),
# as.factor(adu1tpayclean_validation$income == 'Above50K'))$overall['Accuracy']

# > accuracy_rf

# Accuracy

# 0.8235

## Confusion Matrix and Statistics Reference

## Prediction FALSE TRUE

## FALSE 2016 331

## TRUE 184 387 Accuracy : 0.8235

## 95% CI : (0.8092, 0.8372)

## No Information Rate : 0.7539

## P-Value [Acc > NIR] : < 2.2e-16 Kappa : 0.4891 McNemar's Test P-Value :
## 1.247e-10 Sensitivity : 0.9164

## Specificity : 0.5390

## Pos Pred Value : 0.8590

## Neg Pred Value : 0.6778

## Prevalence : 0.7539
```

```

## Detection Rate : 0.6909

## Detection Prevalence : 0.8043

## Balanced Accuracy : 0.7277 'Positive' Class : FALSE

# With 'Occupation' as one of the feature

# Confusion Matrix and Statistics

# Reference

# Prediction FALSE TRUE

# FALSE 2011 305

# TRUE 189 413

# Accuracy : 0.8307

# 95% CI : (0.8166, 0.8441)

# No Information Rate : 0.7539

# P-Value [Acc > NIR] : < 2.2e-16

# Kappa : 0.5175

# Mcnemar's Test P-Value : 2.29e-07

# Sensitivity : 0.9141

# Specificity : 0.5752

# Pos Pred Value : 0.8683

# Neg Pred Value : 0.6860

# Prevalence : 0.7539

# Detection Rate : 0.6892

# Detection Prevalence : 0.7937

# Balanced Accuracy : 0.7446

# 'Positive' Class : FALSE

# f1_rf

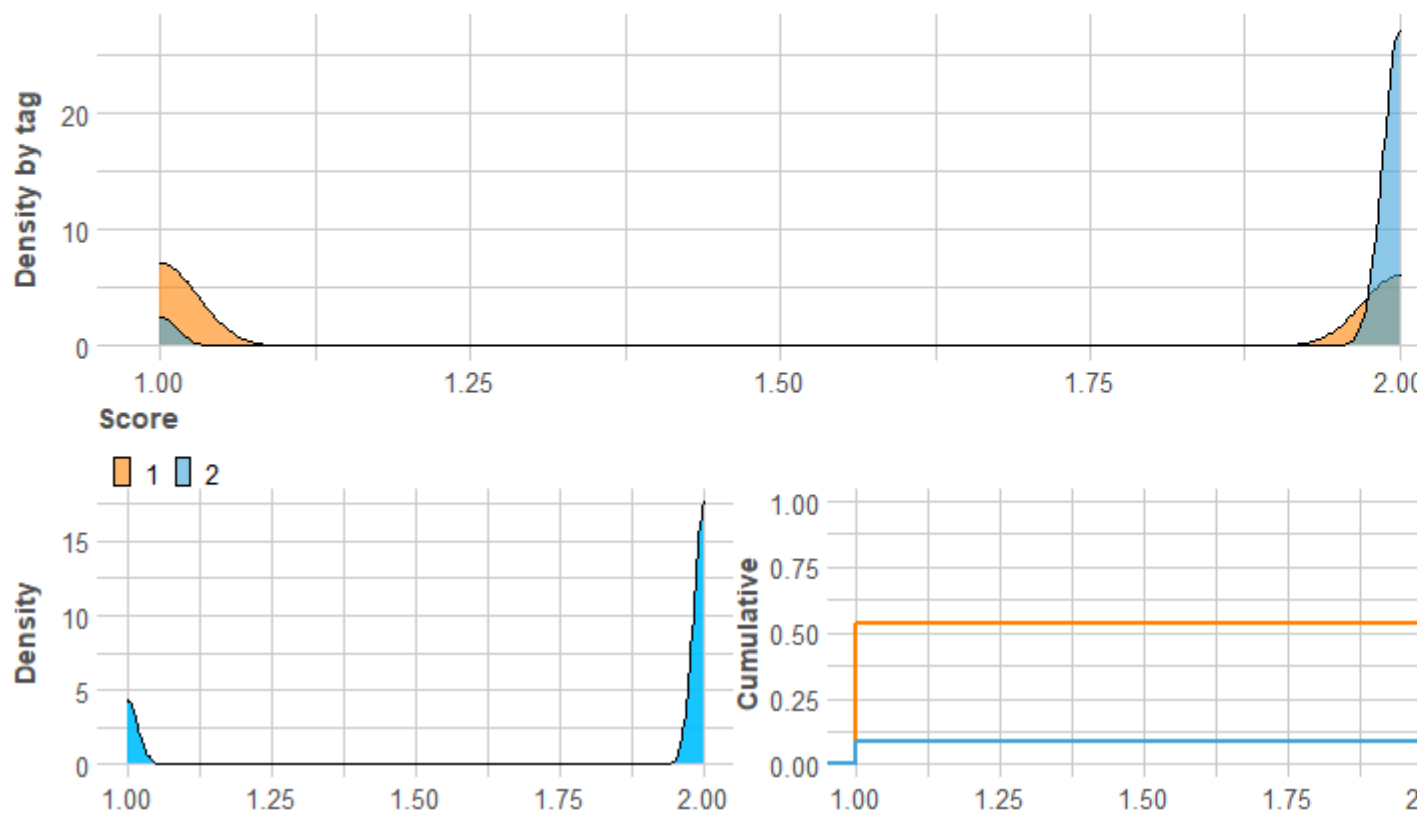
# [1] 0.8867385

# Area under the curve: 0.7277

```

The accuracy of random forest is at 0.8235, sensitivity is at 0.9164 and specificity is at 0.5390. Prevalence is about the same when compared to knn and rpart classification models. The confidence intervals is about the same as rpart. The F1 score is 0.88674 is the best so far. AUC also on the positive side

Classification Model Results



Score cuts (10 quantiles)

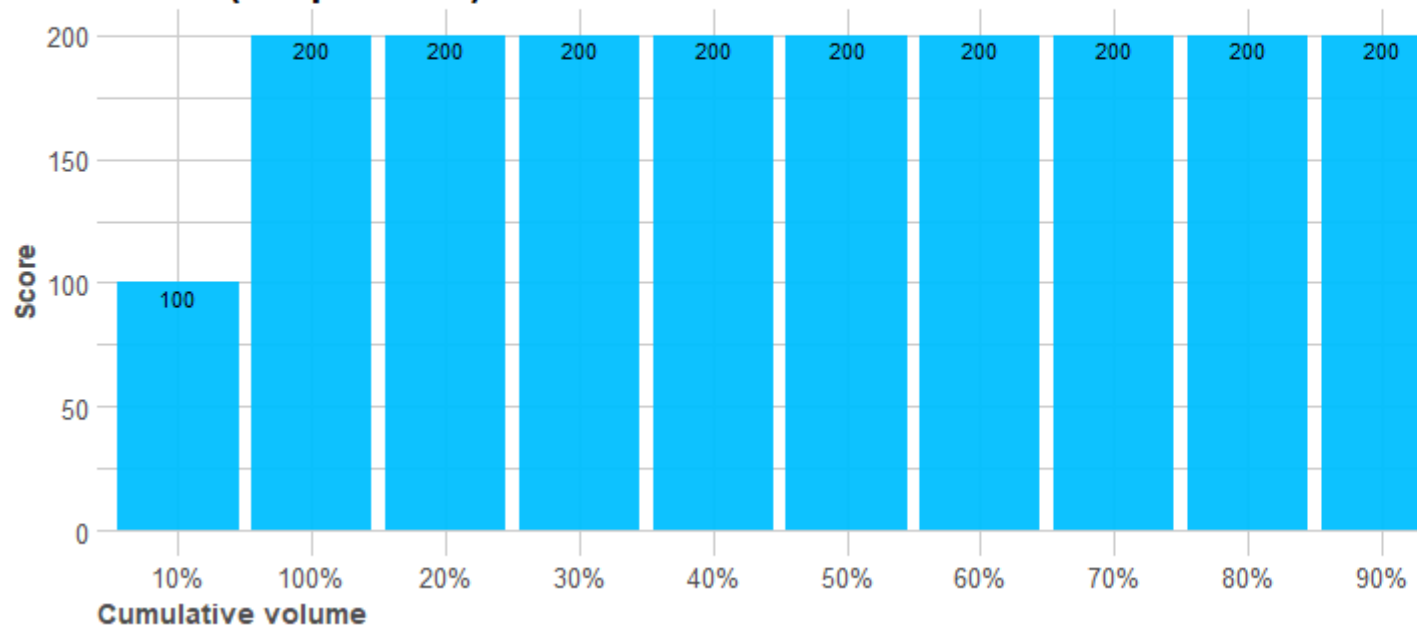


Figure 30: random forest classification summary

When we add “Occupation” as one the feature (with some missing occupation data (5%)) we see that the accuracy is 0.8307 and sensitivity and specificity are 0.9141 and 0.5752 respectively with a F1 score of 0.89061.

When we add “class” as one the feature (with some missing occupation data (5%)) we see that the accuracy is 0.8283 and sensitivity and specificity are 0.9086 and 0.5822 respectively with a F1 score of 0.88864.

Adding class has a negative impact on the accuracy and precision when compared to occupation. We can further investigate into how often a specific feature is used in the predictions using variable importance

```
# varImp(train_rf)

# Overall

# age 1431.1949

# eduyears 1124.9818

# sex 117.8895

# race 152.8045

# hoursperweek 844.8170

# maritalstatus 921.5902

# relationship 1162.1617

# occupation 1038.7262

# class 430.0374
```

We can see that age, years of education, and relationship are the most used and sex is the least used feature. We can also see that occupation did have a greater impact on decisioning even though we had 5% bad data. However, it does impact the precision of the prediction.

Lets tune this model just like the KNN3 classification and see if we can do better.

```
# nodesize <- seq(1, 51, 10)

# acc <- sapply(nodesize, function(ns){ train(y ~ age + eduyears + sex + race +
# hoursperweek + maritalstatus + relationship, method = 'rf', data = temp,
# tuneGrid = data.frame(mtry = 2), nodesize = ns)$results$Accuracy })

# qplot(nodesize, acc)

# train_rf_2 <- randomForest(y ~ age + eduyears + sex + race + hoursperweek +
# maritalstatus + relationship, data=temp,nodesize = nodesize[which.max(acc)])

# y_hat_rf2 <- predict(train_rf_2, adu1tpayclean_validation)

# accuracy_rftune <- confusionMatrix(predict(train_rf_2,
# adu1tpayclean_validation), as.factor(adu1tpayclean_validation$income ==
# 'Above50K'))$overall['Accuracy']
```

```

# > accuracy_rftune

# Accuracy

# 0.8241947

## Confusion Matrix and Statistics Reference

## Prediction FALSE TRUE

## FALSE 2031 348

## TRUE 169 370 Accuracy : 0.8228

## 95% CI : (0.8085, 0.8365)

## No Information Rate : 0.7539

## P-Value [Acc > NIR] : < 2.2e-16 Kappa : 0.4787 McNemar's Test P-Value :
## 4.94e-15 Sensitivity : 0.9232

## Specificity : 0.5153

## Pos Pred Value : 0.8537

## Neg Pred Value : 0.6865

## Prevalence : 0.7539

## Detection Rate : 0.6960

## Detection Prevalence : 0.8153

## Balanced Accuracy : 0.7193 'Positive' Class : FALSE

# Confusion Matrix and Statistics

# Reference

# Prediction FALSE TRUE

# FALSE 2029 295

# TRUE 171 423

# Accuracy : 0.8403

# 95% CI : (0.8265, 0.8534)

# No Information Rate : 0.7539

# P-Value [Acc > NIR] : < 2.2e-16

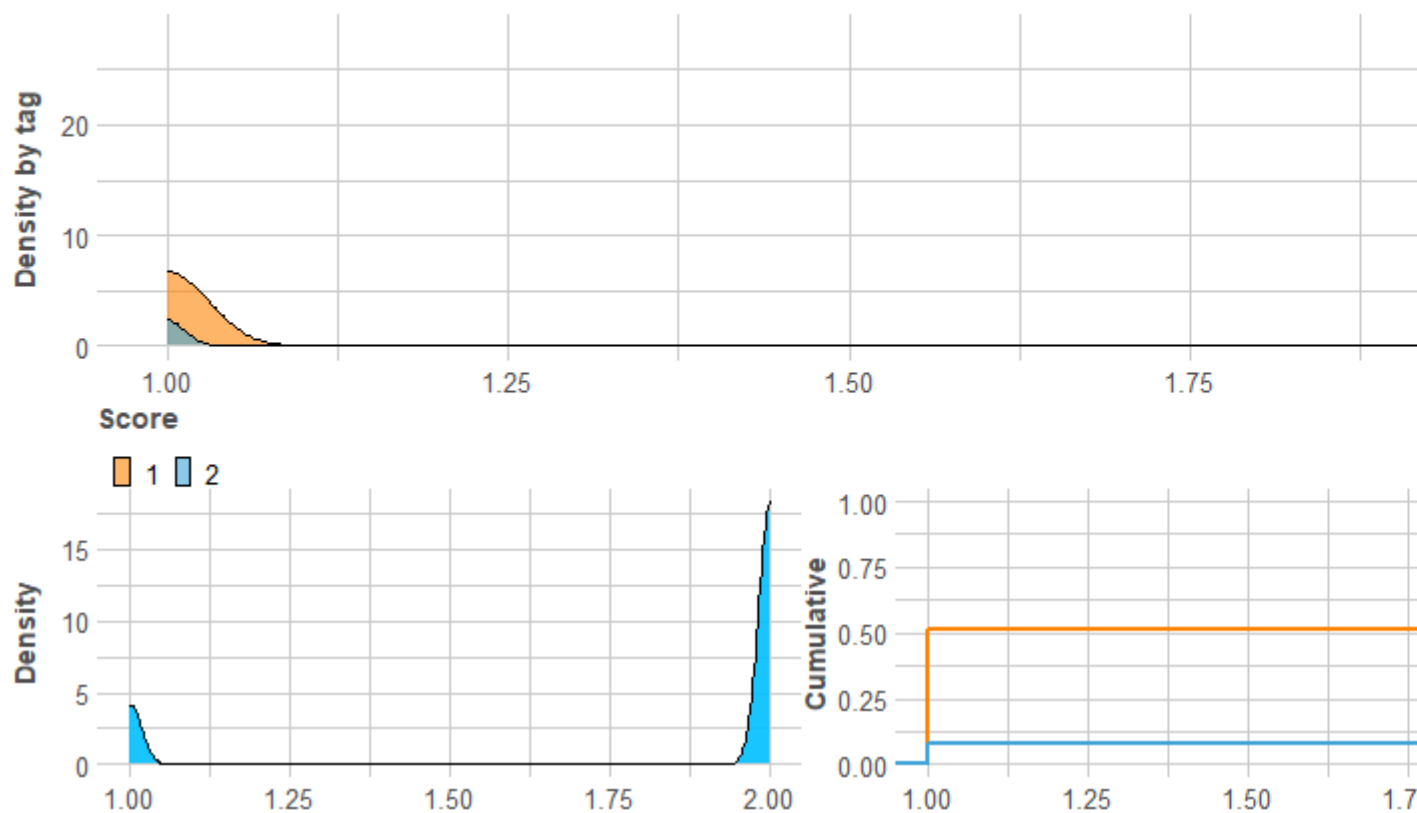
```



```
# Kappa : 0.543
# McNemar's Test P-Value : 1.213e-08
# Sensitivity : 0.9223
# Specificity : 0.5891
# Pos Pred Value : 0.8731
# Neg Pred Value : 0.7121
# Prevalence : 0.7539
# Detection Rate : 0.6953
# Detection Prevalence : 0.7964
# Balanced Accuracy : 0.7557
# 'Positive' Class : FALSE
# f1_rf2
# [1] 0.8870933
# Area under the curve: 0.7186
```

The accuracy of random forest with tuning is at 0.8228, sensitivity is at 0.9232 and specificity is at 0.5153. Prevalence is about the same when compared to knn and rpart classification models. The confidence intervals is just about the same as rpart and random forest without tuning. Interestingly enough, the accuracy is slightly lesser than the random forest test without tuning. Overall there was no further improvement in accuracy. The tuning resulted in overfitting and didn't provide additional value to precision or accuracy. AUC also decreased slightly.

Classification Model Results



Score cuts (10 quantiles)

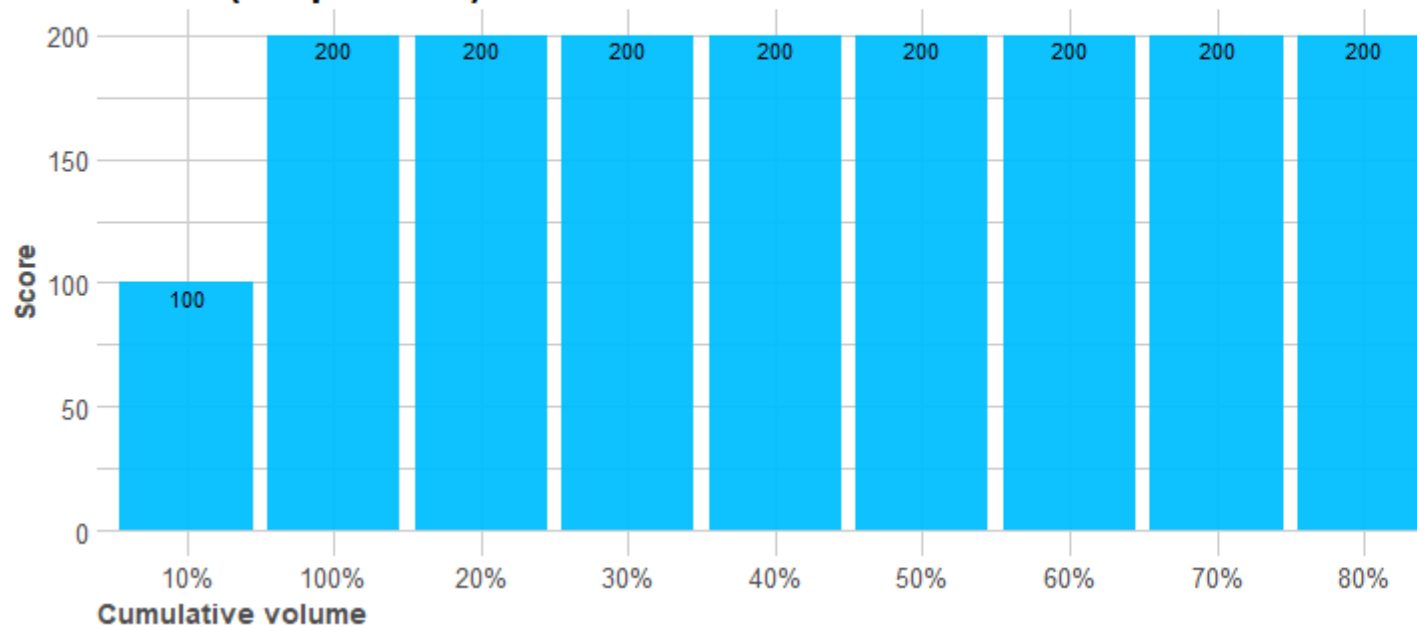


Figure 31: random forest tuned classification summary

When we added the experimental occupation and class data into the dataset for rf tuned test we see that the accuracy improves to 0.8403 and sensitivity is 0.9223 and specificity is 0.5891 with prevalence being the same as before and balanced accuracy improved to 0.7557. The F1 score for this is 0.896

When we compare our final model with naive bayes, the AUC for naive bayes seems to be far better than random forest. However, if you look closer at the classification models, there is a high degree of imbalance in the dataset and naive bayes doesn't do very well when the dataset is imbalanced. F1 score of random forest is higher than naive bayes. For this reason I chose random forest over naive bayes.

Results

This project created a machine learning model that predicts income level of adults based on 7 shortlisted attributes and 2 experimental attributes. This model was run on a test dataset that included an income level indicator, allowing us to compare the predicted and actual value.

The final model that was chosen after comparison was random forest model. The model returned an accuracy of 82 to 83% while predicting annual income more than \$50,000 annually. If occupation and class (with 5% unknown data) are included as features then accuracy changes to 83 to 84%. A summary table of results is given below:

	Method	Accuracy	Sensitivity	Specificity	Prevalence	F1
1.	Plain old guess	0.50034	0.48329	0.50591	0.24606	0.32249
2.	linear model	0.81666	0.47911	0.92682	0.24606	0.56255
3.	General linear model	0.8098	0.38858	0.94727	0.24606	0.50135
4.	naive bayes	0.77176	0.91192	0.52462	0.63811	0.83604
5.	knn	0.80809	0.90364	0.51532	0.75394	0.87654
6.	knn tune	0.81151	0.90318	0.52368	0.75394	0.87746
7.	rpart	0.82111	0.91	0.54875	0.75394	0.88467
8.	rf	0.82351	0.91636	0.539	0.75394	0.88674
9.	rf tune	0.82248	0.92318	0.51532	0.75394	0.88709

Figure 32: Final Results

The accuracy of the predictions is backed by sensitivity of 0.9163 and specificity of 0.539. The F1 score 0.8867 vouches for the sensitivity and specificity. AUC is around 0.72 and is considered as acceptable considering some degree of imbalance in the dataset.

This model can be used to determine income levels for adults in any year with similar attribute sets and achieve comparable accuracies. The following attributes were determined to influence annual adult incomes in the US:

- age
- eduyears
- sex
- race
- hoursperweek
- maritalstatus

- relationship
- occupation and class (experimental with 5% unknown data each)

Conclusion

This machine learning model was able to predict annual incomes of persons in US based (1994) on 7 parameters with an accuracy of 82 to 83%. This model can be applied to data from other census years as well. The model will perform better if the training set is updated with new data that is confirmed for correctness, that is, the label value is the real life value, and not the predicted value. Additional data for certain features like marital status, occupation, and class will make the model better. The model can continuously learn from changing data in the training set to adapt to new parameters, thus improving its accuracy and other metrics

Further improvements could be made by use random oversampling and/or undersampling techniques to fill the void created by imbalanced datasets.

Another approach for improving accuracy would be to create an ensemble based on multiple algorithms like rpart, random forest and glm.

Appendix A - Complete code

```
## @knitr CensusPayR
```

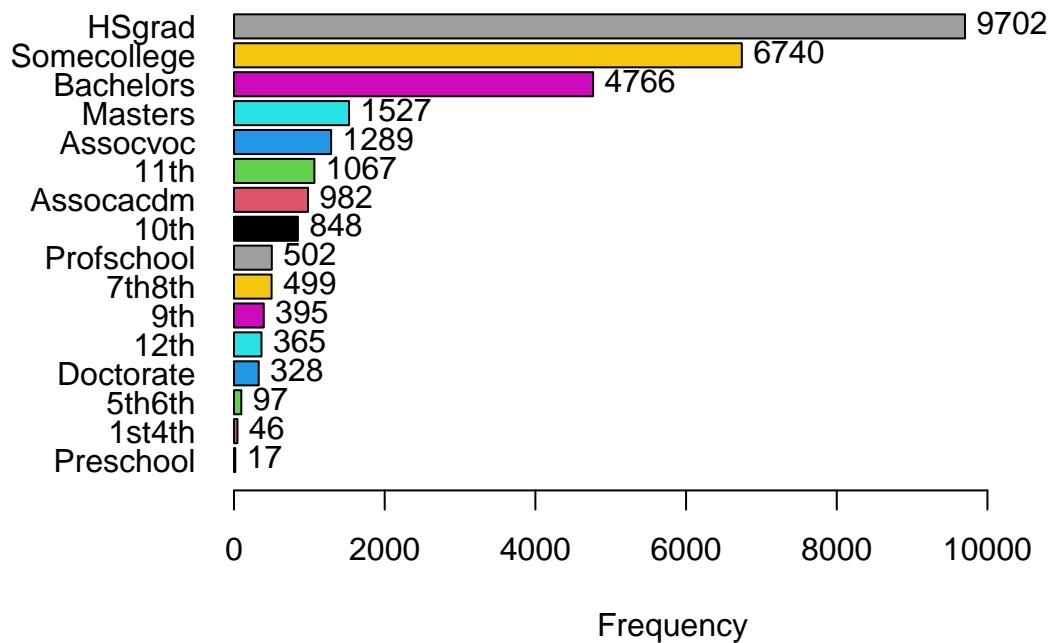
```
# Note: This script will take a while to run. In particular the knn and random  
# forest algorithms with tuning grids will take more time. please be patient if  
# you happen to execute it. The execution report is available in the github  
# location as well
```

```
# Execute the given source code for the project  
source("DatasetProcessingCode.R")
```

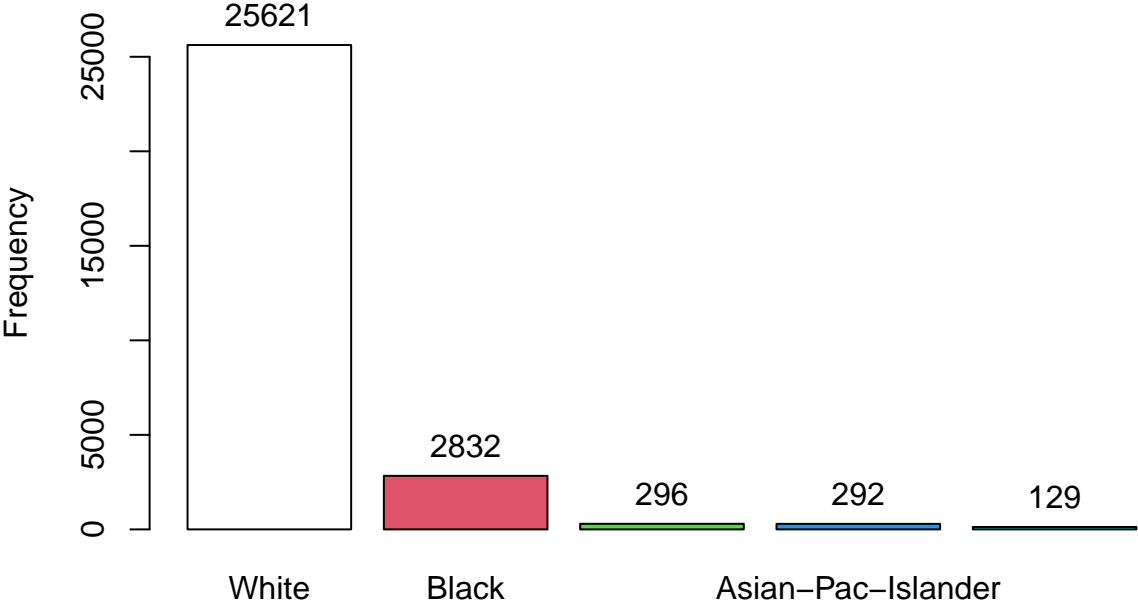
```
## Rows: 32,561  
## Columns: 15  
## $ age          <int> 90, 82, 66, 54, 41, 34, 38, 74, 68, 41, 45, 38, 52, 32, ~  
## $ workclass    <chr> "?", "Private", "?", "Private", "Private", "Private", "~  
## $ fnlwgt       <int> 77053, 132870, 186061, 140359, 264663, 216864, 150601, ~  
## $ education    <chr> "HS-grad", "HS-grad", "Some-college", "7th-8th", "Some--  
## $ education.num <int> 9, 9, 10, 4, 10, 9, 6, 16, 9, 10, 16, 15, 13, 14, 16, 1~  
## $ marital.status <chr> "Widowed", "Widowed", "Widowed", "Divorced", "Separated~  
## $ occupation   <chr> "?", "Exec-managerial", "?", "Machine-op-inspct", "Prof~  
## $ relationship <chr> "Not-in-family", "Not-in-family", "Unmarried", "Unmarri~  
## $ race         <chr> "White", "White", "Black", "White", "White", "White", "~  
## $ sex          <chr> "Female", "Female", "Female", "Female", "Female", "Fema~  
## $ capital.gain  <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~  
## $ capital.loss  <int> 4356, 4356, 4356, 3900, 3900, 3770, 3770, 3683, 3683, 3~  
## $ hours.per.week <int> 40, 18, 40, 40, 40, 45, 40, 20, 40, 60, 35, 45, 20, 55, ~  
## $ native.country <chr> "United-States", "United-States", "United-States", "Uni~  
## $ income       <chr> "<=50K", "<=50K", "<=50K", "<=50K", "<=50K", "<=50K", "~  
## Rows: 29,170  
## Columns: 13  
## $ age          <int> 90, 82, 66, 54, 41, 34, 38, 74, 68, 45, 38, 52, 32, 51, ~  
## $ fnlwgt       <int> 77053, 132870, 186061, 140359, 264663, 216864, 150601, 8~
```

```
## $ education      <fct> HSgrad, HSgrad, Somecollege, 7th8th, Somecollege, HSgrad~
## $ eduyears       <int> 9, 9, 10, 4, 10, 9, 6, 16, 9, 16, 15, 13, 14, 16, 15, 7,~
## $ maritalstatus  <fct> Widowed, Widowed, Widowed, Divorced, Separated, Divorced~
## $ occupation     <fct> Unknown, Execmanagerial, Unknown, Machineopinspct, Profs~
## $ relationship   <fct> Notinfamily, Notinfamily, Unmarried, Unmarried, Ownchild~
## $ race           <fct> White, White, Black, White, White, White, White, White, ~
## $ sex            <fct> Female, Female, Female, Female, Female, Female, Male, Fe~
## $ hoursperweek   <int> 40, 18, 40, 40, 40, 45, 40, 20, 40, 35, 45, 20, 55, 40, ~
## $ native         <chr> "UnitedStates", "UnitedStates", "UnitedStates", "UnitedS~
## $ income         <fct> AtBelow50K, AtBelow50K, AtBelow50K, AtBelow50K, AtBelow5~
## $ class          <fct> Unknown, Private, Unknown, Private, Private, Private, Pr~
```

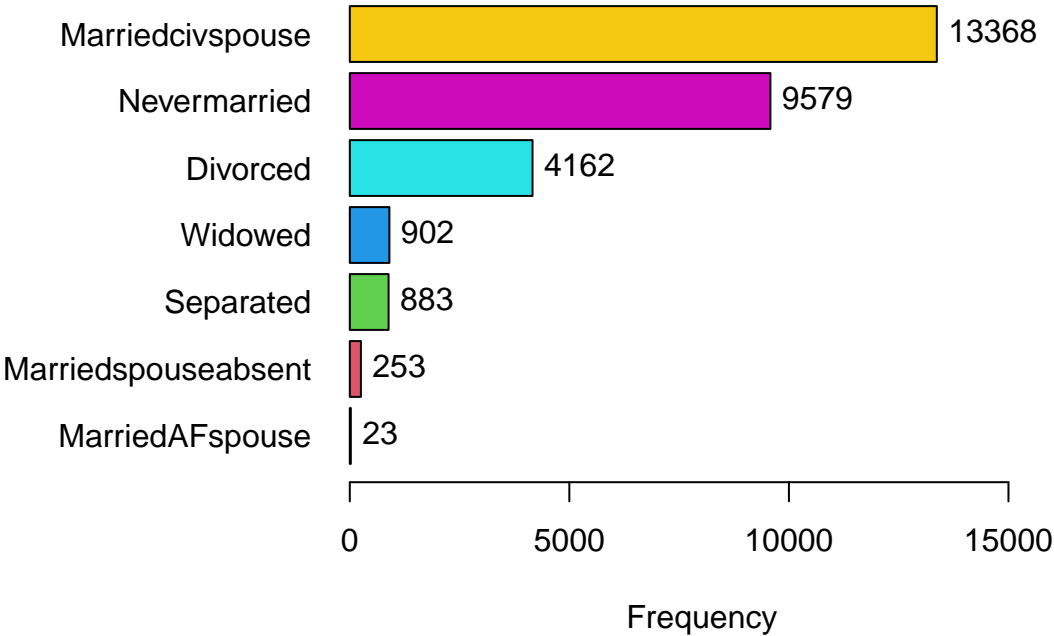
Distribution of adultpayclean\$education

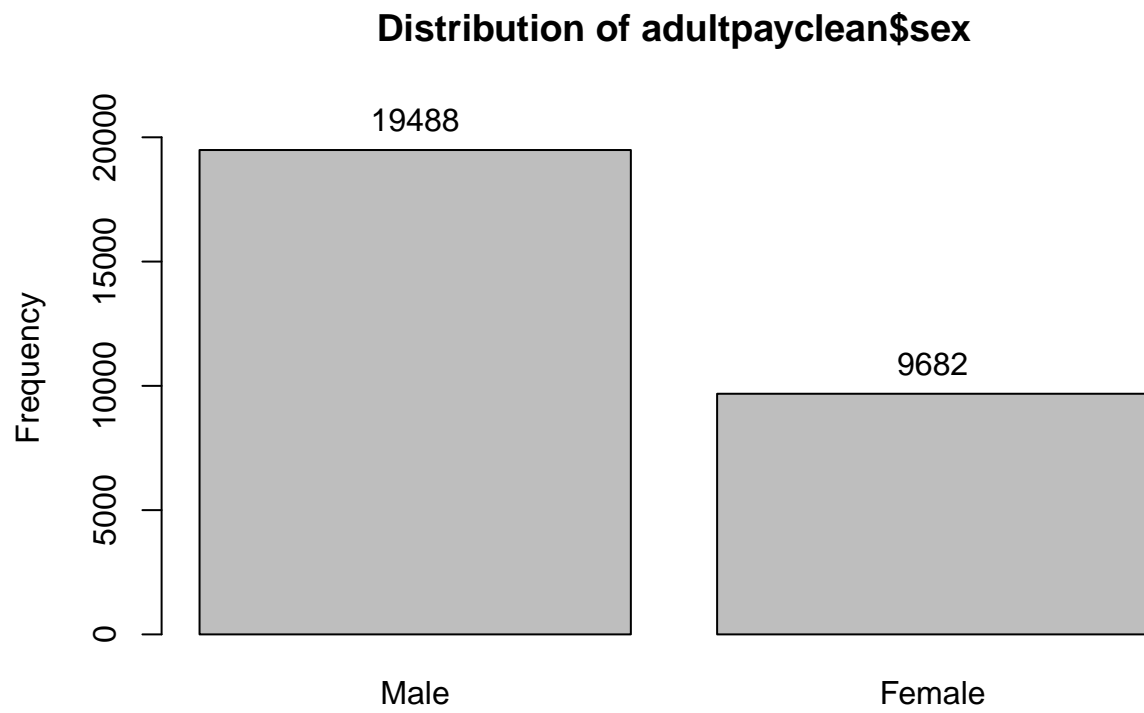


Distribution of adultpayclean\$race

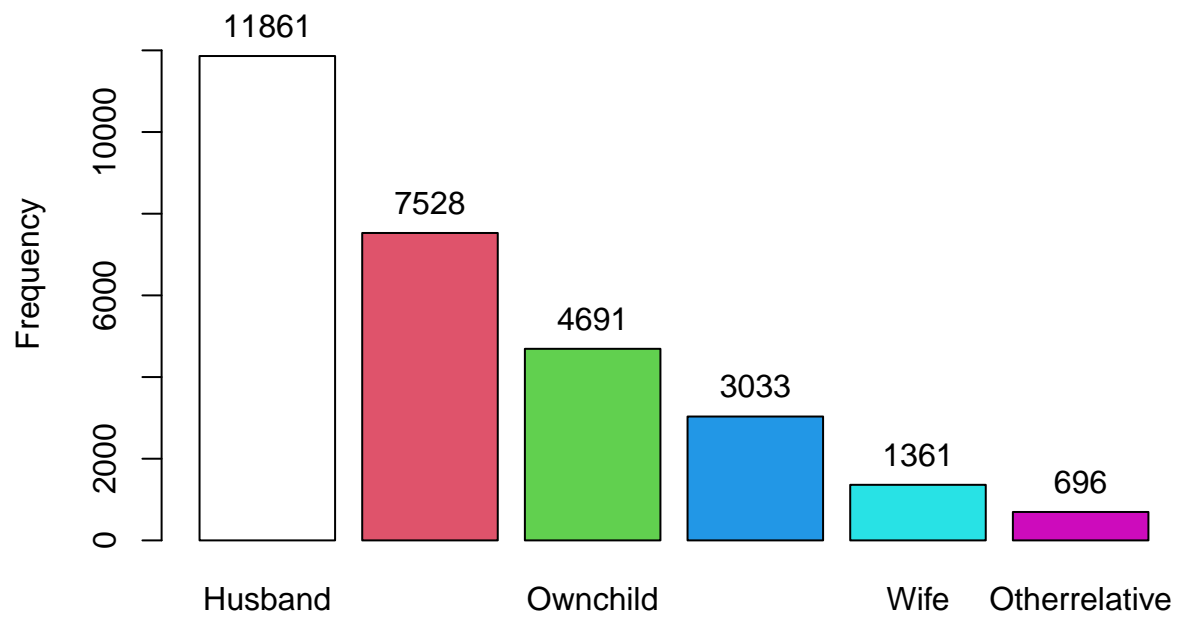


Distribution of adultpayclean\$maritalstatus

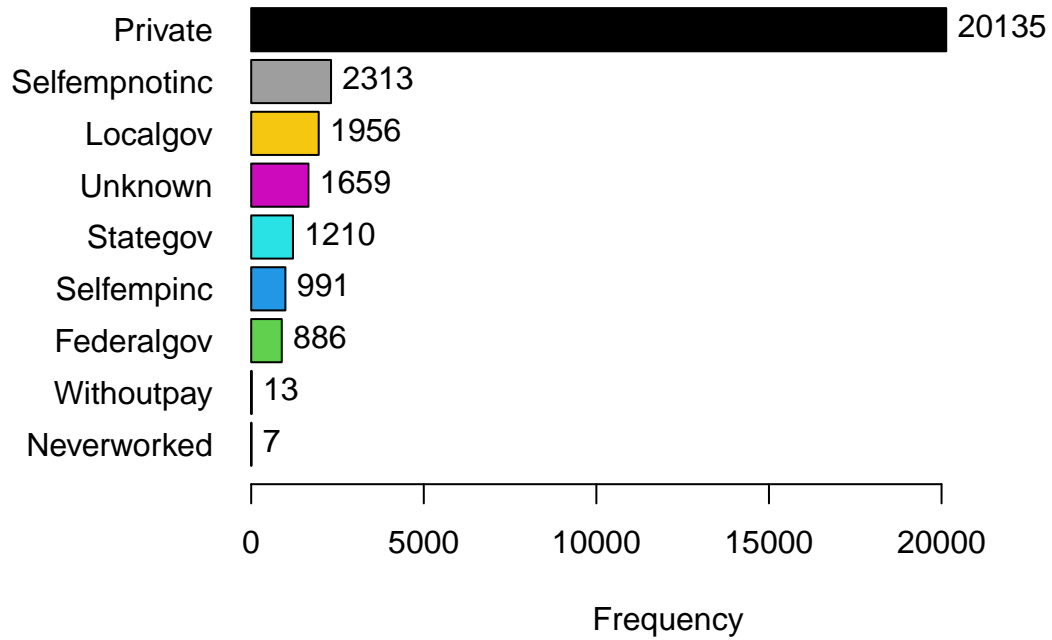




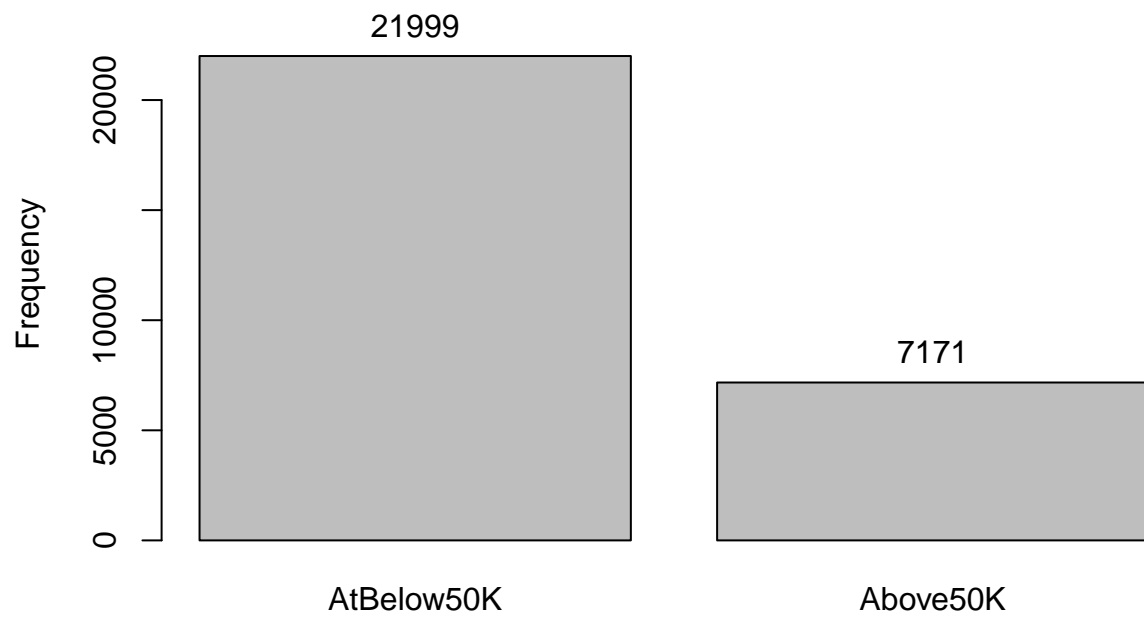
Distribution of adultpayclean\$relationship



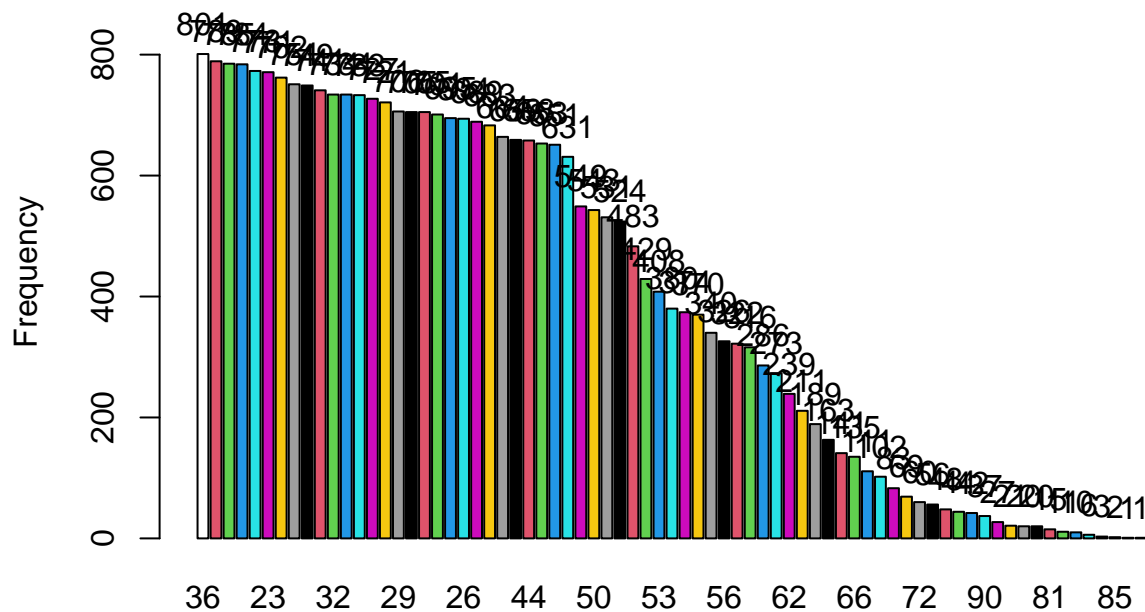
Distribution of adultpayclean\$class



Distribution of adultpayclean\$income



Distribution of adultpayclean\$age



```

if (!require(randomForest)) install.packages("randomForest", repos = "http://cran.us.r-project.org")

if (!require(purrr)) install.packages("purrr", repos = "http://cran.us.r-project.org")

if (!require(e1071)) install.packages("e1071")

if (!require(pROC)) install.packages("pRoc")

library(caret)
library(gridExtra)
library(kableExtra)
library(randomForest)
library(purrr)
library(e1071)
library(caTools)
library(pROC)

set.seed(1996, sample.kind = "Rounding")

# the simplest possible machine algorithm: guessing the outcome
seat_of_the_pants <- sample(c("Above50K", "AtBelow50K"), length(test_index), replace = TRUE) %>%
  factor(levels = levels(adultpayclean_validation$income))
accuracy_guess <- mean(seat_of_the_pants == adultpayclean_validation$income)

# build a confusion matrix for this simple model

```

```
table(predicted = seat_of_the_pants, actual = adultpayclean_validation$income)
```

```
##           actual
## predicted  Above50K AtBelow50K
## Above50K      347      1087
## AtBelow50K    371      1113
```

```
# tabulate accuracy by income levels
adultpayclean_validation %>%
  mutate(y_hat = seat_of_the_pants) %>%
  group_by(income) %>%
  summarize(accuracy = mean(y_hat == income))
```

```
## # A tibble: 2 x 2
##   income      accuracy
##   <fct>      <dbl>
## 1 Above50K    0.483
## 2 AtBelow50K  0.506
```

```
# confusion matrix using R function
cm <- confusionMatrix(data = seat_of_the_pants, reference = adultpayclean_validation$income)
cm
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  Above50K AtBelow50K
## Above50K      347      1087
## AtBelow50K    371      1113
##
##           Accuracy : 0.5003
##           95% CI : (0.482, 0.5186)
##       No Information Rate : 0.7539
##       P-Value [Acc > NIR] : 1
##
##           Kappa : -0.0081
##
##  Mcnemar's Test P-Value : <2e-16
##
##           Sensitivity : 0.4833
##           Specificity : 0.5059
##       Pos Pred Value : 0.2420
##       Neg Pred Value : 0.7500
##       Prevalence : 0.2461
##       Detection Rate : 0.1189
##       Detection Prevalence : 0.4914
##       Balanced Accuracy : 0.4946
##
##       'Positive' Class : Above50K
##
```

```

# record the sensitivity, specificity, and prevalence
sensitivity_guess <- cm$byClass[["Sensitivity"]]
specificity_guess <- cm$byClass[["Specificity"]]
prevalence_guess <- cm$byClass[["Prevalence"]]
f1_guess <- cm$byClass[["F1"]]

auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(seat_of_the_pants ==
  "Above50K", 1, 2))

## Area under the curve: 0.4946

# logistic linear model create the model
lm_fit <- adultpayclean_train %>%
  mutate(y = as.numeric(income == "Above50K")) %>%
  lm(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus + relationship,
    data = .)

# predict using test set
p_hat_logit <- predict(lm_fit, newdata = adultpayclean_validation)

# translate predicted data into factor
y_hat_logit <- ifelse(p_hat_logit > 0.5, "Above50K", "AtBelow50K") %>%
  factor

# compare the predicted vs observed values and use confusionMatrix to get the
# accuracy and other metrics
cm_lm <- confusionMatrix(y_hat_logit, adultpayclean_validation$income)
accuracy_lm <- confusionMatrix(y_hat_logit, adultpayclean_validation$income)$overall[["Accuracy"]]

cm_lm

## Confusion Matrix and Statistics
##
##               Reference
## Prediction  Above50K AtBelow50K
##   Above50K      344      161
##   AtBelow50K     374     2039
##
##               Accuracy : 0.8167
##               95% CI : (0.8021, 0.8305)
##   No Information Rate : 0.7539
##   P-Value [Acc > NIR] : 2.788e-16
##
##               Kappa : 0.451
##
##   Mcnemar's Test P-Value : < 2.2e-16
##
##               Sensitivity : 0.4791
##               Specificity : 0.9268
##   Pos Pred Value : 0.6812
##   Neg Pred Value : 0.8450
##               Prevalence : 0.2461
##   Detection Rate : 0.1179

```

```
## Detection Prevalence : 0.1731
## Balanced Accuracy : 0.7030
##
## 'Positive' Class : Above50K
##
```

```
# record the sensitivity, specificity, and prevalence
sensitivity_lm <- cm_lm$byClass[["Sensitivity"]]
specificity_lm <- cm_lm$byClass[["Specificity"]]
prevalence_lm <- cm_lm$byClass[["Prevalence"]]
f1_lm <- cm_lm$byClass[["F1"]]

auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_logit) ==
  "Above50K", 1, 2))
```

```
## Area under the curve: 0.703
```

```
# general linear model create the glm model
glm_fit <- adultpayclean_train %>%
  mutate(y = as.numeric(income == "Above50K")) %>%
  glm(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus + relationship,
    data = ., family = "binomial")

# predict using validation set
p_hat_logit <- predict(glm_fit, newdata = adultpayclean_validation)

# translate the predicted data into factor
y_hat_logit <- ifelse(p_hat_logit > 0.5, "Above50K", "AtBelow50K") %>%
  factor

# compare the predicted vs observed values and use confusionMatrix to get the
# accuracy and other metrics for the glm model
cm_glm <- confusionMatrix(y_hat_logit, adultpayclean_validation$income)
accuracy_glm <- confusionMatrix(y_hat_logit, adultpayclean_validation$income)$overall[["Accuracy"]]

cm_glm
```

```
## Confusion Matrix and Statistics
```

```
##
##           Reference
## Prediction  Above50K AtBelow50K
## Above50K      279      116
## AtBelow50K    439      2084
##
##           Accuracy : 0.8098
##           95% CI : (0.7951, 0.8239)
##       No Information Rate : 0.7539
##       P-Value [Acc > NIR] : 3.442e-13
##
##           Kappa : 0.3958
##
## Mcnemar's Test P-Value : < 2.2e-16
##
```

```
##           Sensitivity : 0.38858
##           Specificity : 0.94727
##           Pos Pred Value : 0.70633
##           Neg Pred Value : 0.82600
##           Prevalence : 0.24606
##           Detection Rate : 0.09561
##           Detection Prevalence : 0.13537
##           Balanced Accuracy : 0.66793
##
##           'Positive' Class : Above50K
##
```

```
# record the sensitivity, specificity, and prevalence
sensitivity_glm <- cm_glm$byClass[["Sensitivity"]]
specificity_glm <- cm_glm$byClass[["Specificity"]]
prevalence_glm <- cm_glm$byClass[["Prevalence"]]
f1_glm <- cm_glm$byClass[["F1"]]

auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_logit) ==
  "Above50K", 1, 2))
```

```
## Area under the curve: 0.6679
```

```
# Naive bayes

train_nb <- adultpayclean_train %>%
  mutate(y = as.factor(income == "Above50K")) %>%
  naiveBayes(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus + relationship,
    data = .)

y_hat_nb <- predict(train_nb, newdata = adultpayclean_validation)
cm_tab <- table(adultpayclean_validation$income == "Above50K", y_hat_nb)
cm_nb <- confusionMatrix(cm_tab)
cm_nb
```

```
## Confusion Matrix and Statistics
##
##           y_hat_nb
##           FALSE TRUE
## FALSE  1698  502
## TRUE   164  554
##
##           Accuracy : 0.7718
##           95% CI : (0.7561, 0.7869)
##           No Information Rate : 0.6381
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.469
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9119
##           Specificity : 0.5246
```



```
##          Pos Pred Value : 0.7718
##          Neg Pred Value : 0.7716
##          Prevalence : 0.6381
##          Detection Rate : 0.5819
##          Detection Prevalence : 0.7539
##          Balanced Accuracy : 0.7183
##
##          'Positive' Class : FALSE
##
```

```
accuracy_nb <- cm_nb$overall[["Accuracy"]]
sensitivity_nb <- cm_nb$byClass[["Sensitivity"]]
specificity_nb <- cm_nb$byClass[["Specificity"]]
prevalence_nb <- cm_nb$byClass[["Prevalence"]]
f1_nb <- cm_nb$byClass[["F1"]]

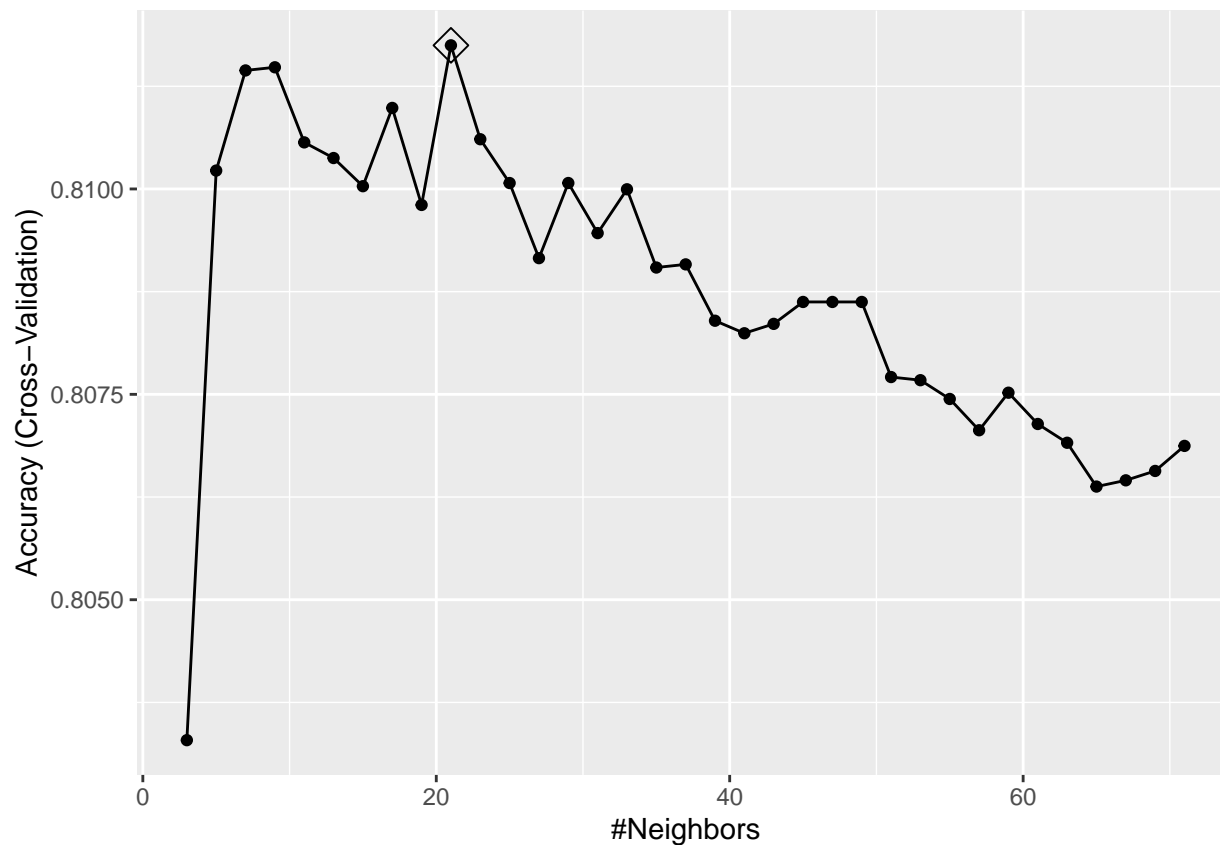
auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_nb) ==
  "TRUE", 1, 2))
```

```
## Area under the curve: 0.7717
```

```
# translate income factor into binary outcome
temp <- adultpayclean_train %>%
  mutate(y = as.factor(income == "Above50K"))

# k-nearest neighbors with a train control and tuning
set.seed(2008)
# train control to use 10% of the observations each to speed up computations
control <- trainControl(method = "cv", number = 10, p = 0.9)
# train the model using knn. choose the best k value using tuning algorithm
train_knn <- train(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
  relationship, method = "knn", data = temp, tuneGrid = data.frame(k = seq(3, 71,
    2)), trControl = control)

# plot the resulting model
ggplot(train_knn, highlight = TRUE)
```



```
# verify which k value was used
train_knn$bestTune
```

```
##      k
## 10 21
```

```
train_knn$finalModel
```

```
## 21-nearest neighbor model
## Training set outcome distribution:
##
## FALSE  TRUE
## 19799  6453
```

```
# use this trained model to predict raw knn predictions
y_hat_knn <- predict(train_knn, adu1tpayclean_validation, type = "raw")

# compare the predicted and observed values using confusionMatrix to get the
# accuracy and other metrics
cm_knn <- confusionMatrix(y_hat_knn, as.factor(adu1tpayclean_validation$income ==
  "Above50K"))
accuracy_knn <- confusionMatrix(y_hat_knn, as.factor(adu1tpayclean_validation$income ==
  "Above50K"))$overall[["Accuracy"]]

cm_knn
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction FALSE TRUE
##      FALSE  1988  348
##      TRUE   212  370
##
##           Accuracy : 0.8081
##           95% CI : (0.7933, 0.8222)
##      No Information Rate : 0.7539
##      P-Value [Acc > NIR] : 1.777e-12
##
##           Kappa : 0.4475
##
##      McNemar's Test P-Value : 1.165e-08
##
##           Sensitivity : 0.9036
##           Specificity : 0.5153
##      Pos Pred Value : 0.8510
##      Neg Pred Value : 0.6357
##           Prevalence : 0.7539
##      Detection Rate : 0.6813
##      Detection Prevalence : 0.8005
##      Balanced Accuracy : 0.7095
##
##      'Positive' Class : FALSE
##
```

```
# record the sensitivity, specificity, and prevalence
```

```
sensitivity_knn <- cm_knn$byClass[["Sensitivity"]]
specificity_knn <- cm_knn$byClass[["Specificity"]]
prevalence_knn <- cm_knn$byClass[["Prevalence"]]
f1_knn <- cm_knn$byClass[["F1"]]
```

```
auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_knn) ==
  "TRUE", 1, 2))
```

```
## Area under the curve: 0.7095
```

```
# k-nearest classification using tuning function
```

```
set.seed(2008)
```

```
# train the model using knn3 classification
```

```
ks <- seq(3, 251, 2)
```

```
knntune <- map_df(ks, function(k) {
  temp <- adultpayclean_train %>%
    mutate(y = as.factor(income == "Above50K"))
  temp_test <- adultpayclean_validation %>%
    mutate(y = as.factor(income == "Above50K"))
```

```
# create the knn3 model
```

```
knn_fit <- knn3(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
  relationship, data = temp, k = k)
```

```
# predict the model for the current k
```

```

y_hat <- predict(knn_fit, temp, type = "class")
# get the confusionmatrix for the current k
cm_train <- confusionMatrix(y_hat, temp$y)
train_error <- cm_train$overall["Accuracy"]
# do the same for test model
y_hat <- predict(knn_fit, temp_test, type = "class")
cm_test <- confusionMatrix(y_hat, temp_test$y)
test_error <- cm_test$overall["Accuracy"]

tibble(train = train_error, test = test_error)
})
# get the accuracy for the k with maximum accuracy
accuracy_knntune <- max(knntune$test)
# get the confusion matrix for that k
knn_fit <- knn3(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
  relationship, data = temp, k = ks[which.max(knntune$test)])

y_hat_knntune <- predict(knn_fit, adultpayclean_validation, type = "class")
cm_knntune <- confusionMatrix(y_hat_knntune, as.factor(adultpayclean_validation$income ==
  "Above50K"))

cm_knntune

```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction FALSE TRUE
##      FALSE  1987  342
##      TRUE   213  376
##
##           Accuracy : 0.8098
##           95% CI : (0.7951, 0.8239)
##      No Information Rate : 0.7539
##      P-Value [Acc > NIR] : 3.442e-13
##
##           Kappa : 0.4544
##
##  McNemar's Test P-Value : 5.532e-08
##
##           Sensitivity : 0.9032
##           Specificity : 0.5237
##      Pos Pred Value : 0.8532
##      Neg Pred Value : 0.6384
##           Prevalence : 0.7539
##      Detection Rate : 0.6809
##      Detection Prevalence : 0.7981
##      Balanced Accuracy : 0.7134
##
##      'Positive' Class : FALSE
##

```

```

# record the sensitivity, specificity, and prevalence
sensitivity_knntune <- cm_knntune$byClass[["Sensitivity"]]
specificity_knntune <- cm_knntune$byClass[["Specificity"]]
prevalence_knntune <- cm_knntune$byClass[["Prevalence"]]
f1_knntune <- cm_knntune$byClass[["F1"]]

auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_knntune) ==
  "TRUE", 1, 2))

```

```
## Area under the curve: 0.7134
```

```

# recursive partitioning using rpart
set.seed(2008)
# train the model with the recursive partitioning
train_rpart <- train(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
  relationship, method = "rpart", tuneGrid = data.frame(cp = seq(0, 0.1, len = 25)),
  data = temp)
# predict the outcomes with this model
y_hat_rpart <- predict(train_rpart, adultpayclean_validation)
# confusion matrix for the rpart model
cm_rpart <- confusionMatrix(y_hat_rpart, as.factor(adultpayclean_validation$income ==
  "Above50K"))
# get the accuracy
accuracy_rpart <- confusionMatrix(y_hat_rpart, as.factor(adultpayclean_validation$income ==
  "Above50K"))$overall["Accuracy"]

cm_rpart

```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction FALSE TRUE
##      FALSE  2002  324
##      TRUE   198  394
##
##           Accuracy : 0.8211
##           95% CI : (0.8067, 0.8349)
##      No Information Rate : 0.7539
##      P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.4876
##
##  Mcnemar's Test P-Value : 4.472e-08
##
##           Sensitivity : 0.9100
##           Specificity : 0.5487
##      Pos Pred Value : 0.8607
##      Neg Pred Value : 0.6655
##           Prevalence : 0.7539
##      Detection Rate : 0.6861
##      Detection Prevalence : 0.7971
##      Balanced Accuracy : 0.7294
##

```

```
##          'Positive' Class : FALSE
##
```

```
# record the sensitivity, specificity, and prevalence
sensitivity_rpart <- cm_rpart$byClass[["Sensitivity"]]
specificity_rpart <- cm_rpart$byClass[["Specificity"]]
prevalence_rpart <- cm_rpart$byClass[["Prevalence"]]
f1_rpart <- cm_rpart$byClass[["F1"]]

auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_rpart) ==
  "TRUE", 1, 2))
```

```
## Area under the curve: 0.7294
```

```
# random forest
set.seed(2008)
# train the vanilla random forest model
train_rf <- randomForest(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
  relationship, data = temp)

y_hat_rf <- predict(train_rf, adultpayclean_validation)

# create the confusionMatrix
cm_rf <- confusionMatrix(y_hat_rf, as.factor(adultpayclean_validation$income == "Above50K"))
# get the accuracy
accuracy_rf <- confusionMatrix(y_hat_rf, as.factor(adultpayclean_validation$income ==
  "Above50K"))$overall["Accuracy"]

cm_rf
```

```
## Confusion Matrix and Statistics
##
##          Reference
## Prediction FALSE TRUE
##      FALSE  2016  331
##      TRUE   184  387
##
##          Accuracy : 0.8235
##          95% CI : (0.8092, 0.8372)
##      No Information Rate : 0.7539
##      P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.4891
##
##      McNemar's Test P-Value : 1.247e-10
##
##          Sensitivity : 0.9164
##          Specificity : 0.5390
##      Pos Pred Value : 0.8590
##      Neg Pred Value : 0.6778
##          Prevalence : 0.7539
##      Detection Rate : 0.6909
##      Detection Prevalence : 0.8043
```

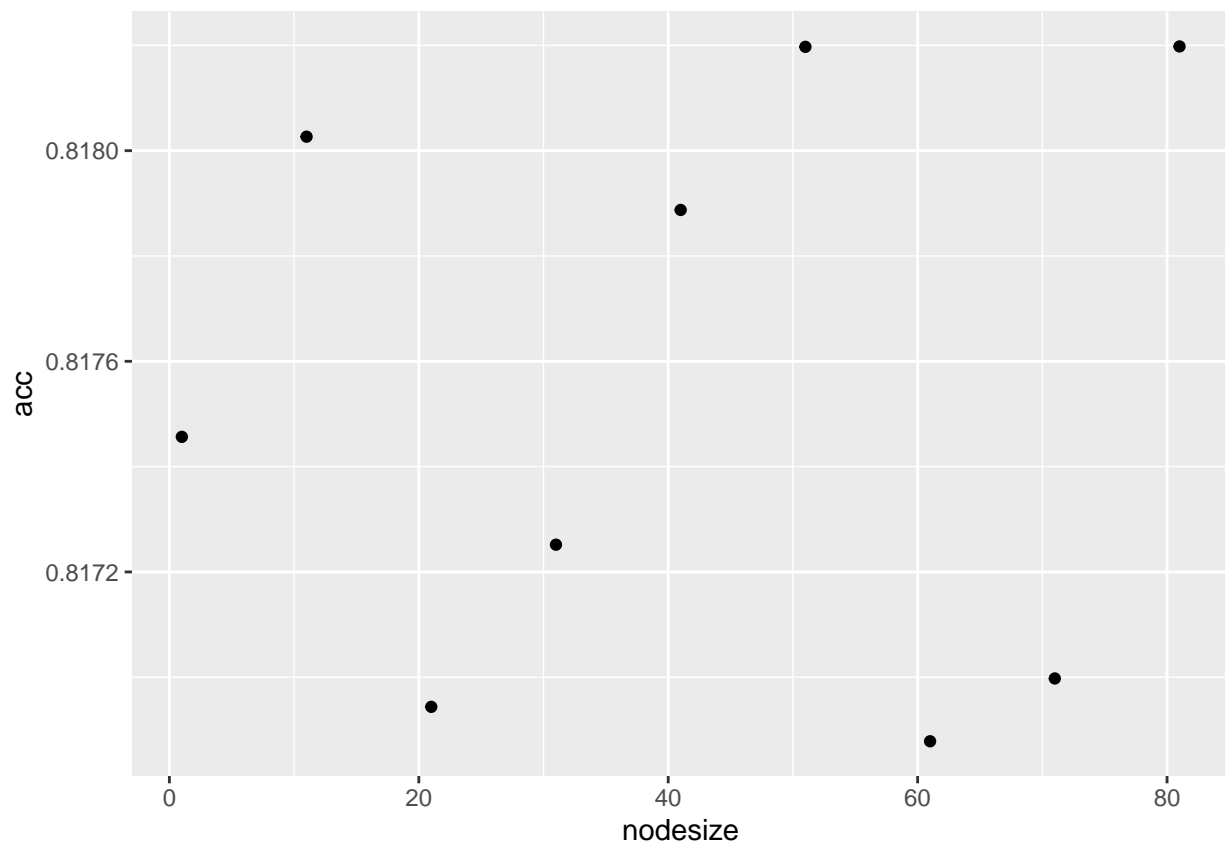
```
##      Balanced Accuracy : 0.7277
##
##      'Positive' Class : FALSE
##
```

```
# record the sensitivity, specificity, and prevalence
sensitivity_rf <- cm_rf$byClass[["Sensitivity"]]
specificity_rf <- cm_rf$byClass[["Specificity"]]
prevalence_rf <- cm_rf$byClass[["Prevalence"]]
f1_rf <- cm_rf$byClass[["F1"]]

auc(ifelse(adultpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_rf) ==
  "TRUE", 1, 2))
```

```
## Area under the curve: 0.7277
```

```
# random forest with tuning
nodesize <- seq(1, 90, 10)
acc <- sapply(nodesize, function(ns) {
  # train the model with tuning
  train(y ~ age + edueyears + sex + race + hoursperweek + maritalstatus + relationship,
    method = "rf", data = temp, tuneGrid = data.frame(mtry = 2), nodesize = ns)$results$Accuracy
})
qplot(nodesize, acc)
```



```

# get the trained model for the max node size
train_rf_2 <- randomForest(y ~ age + eduyears + sex + race + hoursperweek + maritalstatus +
  relationship, data = temp, nodesize = nodesize[which.max(acc)])
# predict the outcomes
y_hat_rf2 <- predict(train_rf_2, adu1tpayclean_validation)
# get the confusion matrix for random forest model
cm_rf2 <- confusionMatrix(predict(train_rf_2, adu1tpayclean_validation), as.factor(adu1tpayclean_validation$income == "Above50K"))
# get the accuracy
accuracy_rftune <- confusionMatrix(predict(train_rf_2, adu1tpayclean_validation),
  as.factor(adu1tpayclean_validation$income == "Above50K"))$overall["Accuracy"]

cm_rf2

```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction FALSE TRUE
##      FALSE  2031  348
##      TRUE   169  370
##
##              Accuracy : 0.8228
##              95% CI : (0.8085, 0.8365)
##      No Information Rate : 0.7539
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.4787
##
##  Mcnemar's Test P-Value : 4.94e-15
##
##      Sensitivity : 0.9232
##      Specificity : 0.5153
##      Pos Pred Value : 0.8537
##      Neg Pred Value : 0.6865
##      Prevalence : 0.7539
##      Detection Rate : 0.6960
##      Detection Prevalence : 0.8153
##      Balanced Accuracy : 0.7193
##
##      'Positive' Class : FALSE
##

```

```

# record the sensitivity, specificity, and prevalence
sensitivity_rf2 <- cm_rf2$byClass[["Sensitivity"]]
specificity_rf2 <- cm_rf2$byClass[["Specificity"]]
prevalence_rf2 <- cm_rf2$byClass[["Prevalence"]]
f1_rf2 <- cm_rf2$byClass[["F1"]]

auc(ifelse(adu1tpayclean_validation$income == "Above50K", 1, 2), ifelse(unname(y_hat_rf2) ==
  "TRUE", 1, 2))

```

```

## Area under the curve: 0.7186

```


	Method	Accuracy	Sensitivity	Specificity	Prevalence	F1
1.	Plain old guess	0.50034	0.48329	0.50591	0.24606	0.32249
2.	linear model	0.81666	0.47911	0.92682	0.24606	0.56255
3.	General linear model	0.8098	0.38858	0.94727	0.24606	0.50135
4.	naive bayes	0.77176	0.91192	0.52462	0.63811	0.83604
5.	knn	0.80809	0.90364	0.51532	0.75394	0.87654
6.	knn tune	0.81151	0.90318	0.52368	0.75394	0.87746
7.	rpart	0.82111	0.91	0.54875	0.75394	0.88467
8.	rf	0.82351	0.91636	0.539	0.75394	0.88674
9.	rf tune	0.82248	0.92318	0.51532	0.75394	0.88709

```
# tabulate all the accuracy results with sensitivity and specificity
accuracy_results <- matrix(c("Plain old guess", round(accuracy_guess, 5), round(sensitivity_guess,
5), round(specificity_guess, 5), round(prevalence_guess, 5), round(f1_guess,
5), "linear model", round(accuracy_lm, 5), round(sensitivity_lm, 5), round(specificity_lm,
5), round(prevalence_lm, 5), round(f1_lm, 5), "General linear model", round(accuracy_glm,
5), round(sensitivity_glm, 5), round(specificity_glm, 5), round(prevalence_glm,
5), round(f1_glm, 5), "naive bayes", round(accuracy_nb, 5), round(sensitivity_nb,
5), round(specificity_nb, 5), round(prevalence_nb, 5), round(f1_nb, 5), "knn",
round(accuracy_knn, 5), round(sensitivity_knn, 5), round(specificity_knn, 5),
round(prevalence_knn, 5), round(f1_knn, 5), "knn tune", round(accuracy_kntune,
5), round(sensitivity_kntune, 5), round(specificity_kntune, 5), round(prevalence_kntune,
5), round(f1_kntune, 5), "rpart", round(accuracy_rpart, 5), round(sensitivity_rpart,
5), round(specificity_rpart, 5), round(prevalence_rpart, 5), round(f1_rpart,
5), "rf", round(accuracy_rf, 5), round(sensitivity_rf, 5), round(specificity_rf,
5), round(prevalence_rf, 5), round(f1_rf, 5), "rf tune", round(accuracy_rftune,
5), round(sensitivity_rf2, 5), round(specificity_rf2, 5), round(prevalence_rf2,
5), round(f1_rf2, 5)), nrow = 9, ncol = 6, byrow = TRUE, dimnames = list(c("1.",
"2.", "3.", "4.", "5.", "6.", "7.", "8.", "9."), c("Method", "Accuracy", "Sensitivity",
"Specificity", "Prevalence", "F1")))
# style the table with knitr
accuracy_results %>%
  knitr::kable() %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed"))
```

Appendix B - Links

<https://www.edx.org/professional-certificate/harvardx-data-science->

[https://www.crcpress.com/Introduction-to-Data-Science-Data-Analysis-and-Prediction-Algorithms-with/Irizarry/p/book/9780367357986-](https://www.crcpress.com/Introduction-to-Data-Science-Data-Analysis-and-Prediction-Algorithms-with-Irizarry/p/book/9780367357986-)

<https://leanpub.com/datasciencebook->

Citations

Irizarry, Rafael A., “Introduction to Data Science: Data Analysis and Prediction Algorithms in R” <https://rafalab.github.io/dsbook/>

ML-Friendly kaggle dataset for adult census income <https://www.kaggle.com/uciml/adult-census-income>