

Homework #3 DATA622

Group 1

9/23/2021

Group 1 Members:

- David Moste
- Vinayak Kamath
- Kenan Sooklall
- Christian Thieme
- Lidiia Tronina

Introduction

We have been provided a dataset containing the loan approval status of 614 loan applicants. In addition to the status, we have also been given 12 additional columns of data describing different characteristics of each loan applicant. Our task is to explore the data, and then use the knowledge gained to create models using Linear Discriminant Analysis (LDA), K-Nearest Neighbors (KNN), Decision Trees, and Random Forests. Finally, we'll compare and contrast each method and its associated accuracy metrics.

The data has the following variables:

Variable	Description
Loan_ID	Unique Loan ID
Gender	Male/Female
Married	Applicant married (Y/N)
Dependents	Number of dependents
Education	Applicant Education (Graduate/Undergraduate)
Self_Employed	Self employed (Y/N)
ApplicantIncome	Applicant income
CoapplicantIncome	Coapplicant income
LoanAmount	Loan amount in thousands
Loan_Amount_Term	Term of loan in months
Credit_History	credit history meets guidelines
Property_Area	Urban/Semi Urban/Rural
Loan_Status (Target)	Loan approved (Y/N)

As can be seen from the table above, the dataset has a wide variety of numeric and categorical data. We'll dive deeper into this data now.

Exploratory Data Analysis

Let's begin exploring by taking a high level look at our dataset:

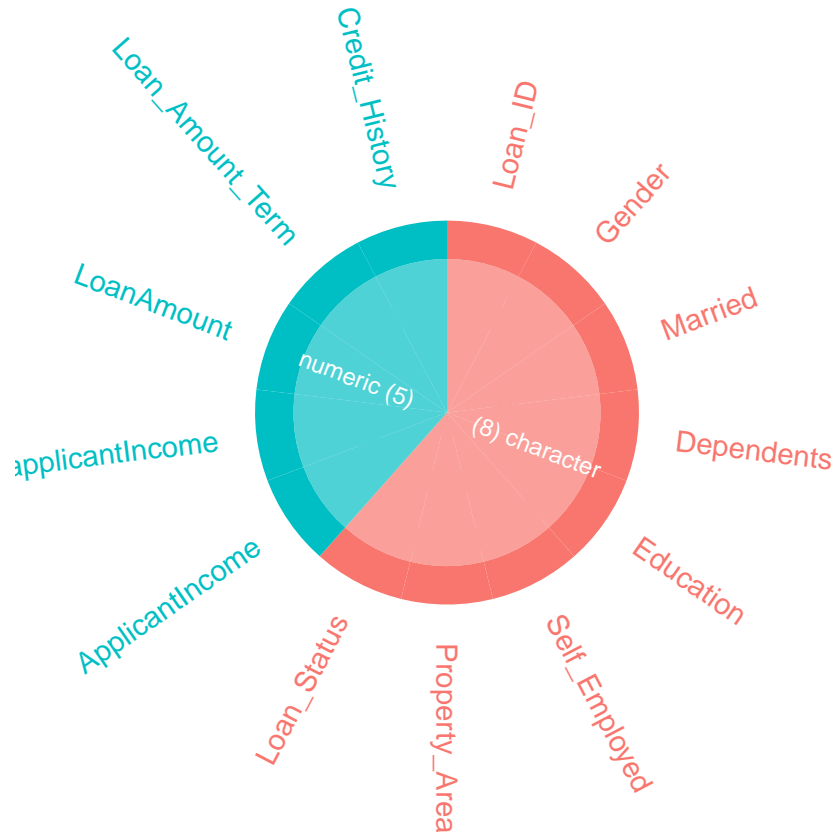
```
glimpse(loans)

## Rows: 614
## Columns: 13
## $ Loan_ID      <chr> "LP001002", "LP001003", "LP001005", "LP001006", "LP0~
## $ Gender       <chr> "Male", "Male", "Male", "Male", "Male", "Male", "Mal~
## $ Married      <chr> "No", "Yes", "Yes", "Yes", "No", "Yes", "Yes", "Yes"~
## $ Dependents   <chr> "0", "1", "0", "0", "0", "2", "0", "3+", "2", "1", "~
## $ Education    <chr> "Graduate", "Graduate", "Graduate", "Not Graduate", ~
## $ Self_Employed <chr> "No", "No", "Yes", "No", "No", "Yes", "No", "No", "N~
## $ ApplicantIncome <dbl> 5849, 4583, 3000, 2583, 6000, 5417, 2333, 3036, 4006~
## $ CoapplicantIncome <dbl> 0, 1508, 0, 2358, 0, 4196, 1516, 2504, 1526, 10968, ~
## $ LoanAmount   <dbl> NA, 128, 66, 120, 141, 267, 95, 158, 168, 349, 70, 1~
## $ Loan_Amount_Term <dbl> 360, 360, 360, 360, 360, 360, 360, 360, 360, 36~
## $ Credit_History <dbl> 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, NA, ~
## $ Property_Area <chr> "Urban", "Rural", "Urban", "Urban", "Urban", "Urban"~
## $ Loan_Status  <chr> "Y", "N", "Y", "Y", "Y", "Y", "Y", "N", "Y", "N", "Y~
```

From this view we can see that this dataset has 614 rows and 13 columns. As mentioned previously, there is a mix of both numeric and categorical data, although it appears that at least one numeric column should actually be categorical. As it stands, the split between numeric and categorical data looks like this:

	n	min	mean	median	max	sd
ApplicantIncome	614	150	5403.4593	3812.5	81000	6109.04167
CoapplicantIncome	614	0	1621.2458	1188.5	41667	2926.24837
LoanAmount	592	9	146.4122	128.0	700	85.58733
Loan_Amount_Term	600	12	342.0000	360.0	480	65.12041

```
inspectdf::inspect_types(loans) %>%
  show_plot
```



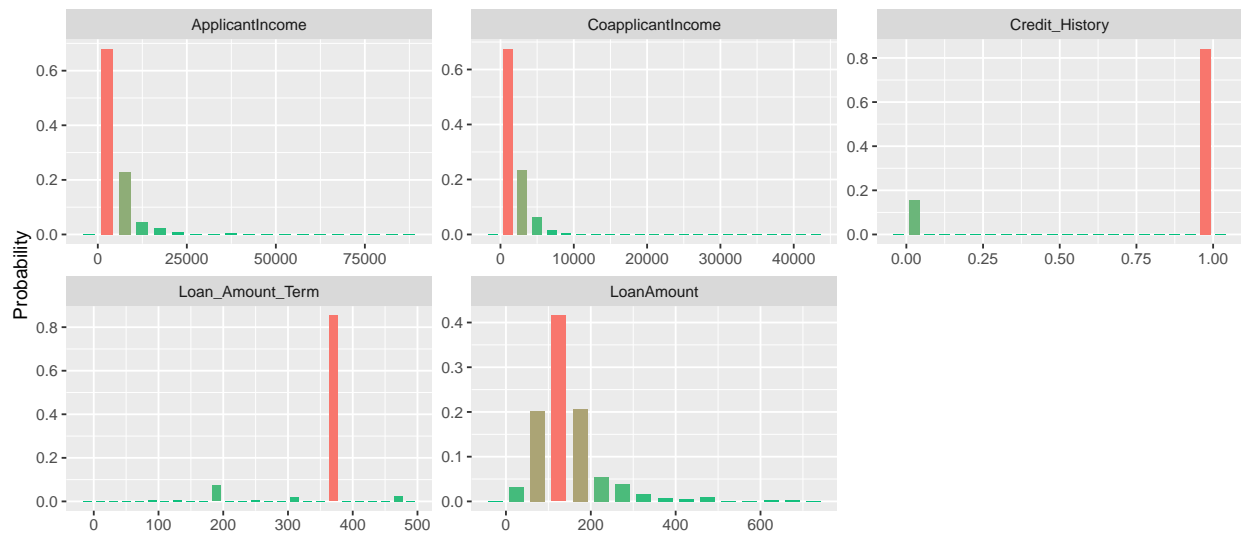
Continuous and categorical variables will be analyzed separately for the sake of clarity. We'll begin our analysis by looking at our numeric features.

Numeric Features

Looking at the above summary table, we can see that there are outliers present in all numerical variables, with **ApplicantIncome** being the worst offender. Even at three times the standard deviation, its maximum value lies far outside of the 68-95-99.7 rule. There is also a significant difference between its mean and median, indicating there is skew present in this variable as well. Further, we'll look at the distribution of each of these features:

```
inspectdf::inspect_num(loans) %>%
  show_plot()
```

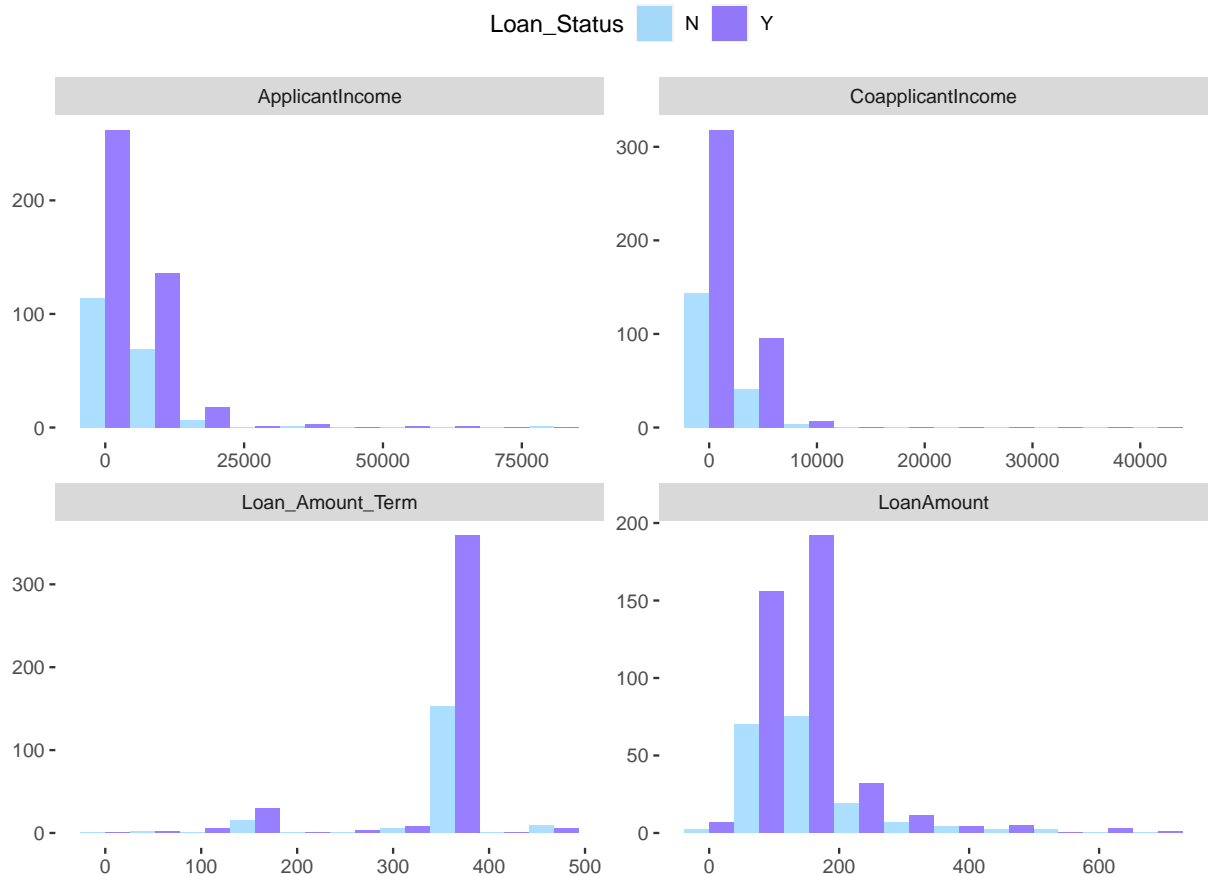
Histograms of numeric columns in df::loans



Looking at the above plots, we can see that:

- **ApplicantIncome** must be on a weekly or monthly basis as most values are fairly low
- **CoapplicantIncome** follows the same trend as applicant income
- **Credit_History** appears to be a categorical variable either 0 or 1
- **Loan_Amount_Term** appears to be 365 days (1 year) more than 80% of the time
- **LoanAmount** is fairly small, ranging from \$0-700. It appears that most loans are for somewhere between \$50 and \$200.

It will also be helpful to understand the the distributions of these variables for each application approval status, both Yes and No.



There does not seem to be a *significant* difference between those who get loan approval, and those who do not, for any of the predictors.

As part of our analysis of numeric features, let's look at the relationship between these features and `Loan_Status`:

```
loan_names <- loans %>% select_if(is.numeric) # %>% select(-Credit_History)
int_names <- names(loan_names)
myGlist <- vector('list', length(loan_names))
names(myGlist) <- int_names

for (i in int_names) {

  myGlist[[i]] <-
    ggplot(loans) +
      aes_string(x = as.factor(loans$Loan_Status), y = i) +
      geom_boxplot(color = 'steelblue',
                    outlier.color = 'firebrick',
                    outlier.alpha = 0.35) +
      labs(title = paste0(i, ' vs Loan_Status'), y = i, x = 'Loan_Status') +
      theme_minimal() +
      theme(
        plot.title = element_text(hjust = 0.45),
        panel.grid.major.y = element_line(color = "grey",
```

Table 2: Summary statistics for Categorical Variables

Dependents	Property_Area
0 :345	Rural :179
1 :102	Semiurban:233
2 :101	Urban :202
3+ : 51	NA
NA's: 15	NA

Table 3: Summary statistics for Binary Categorical Variables

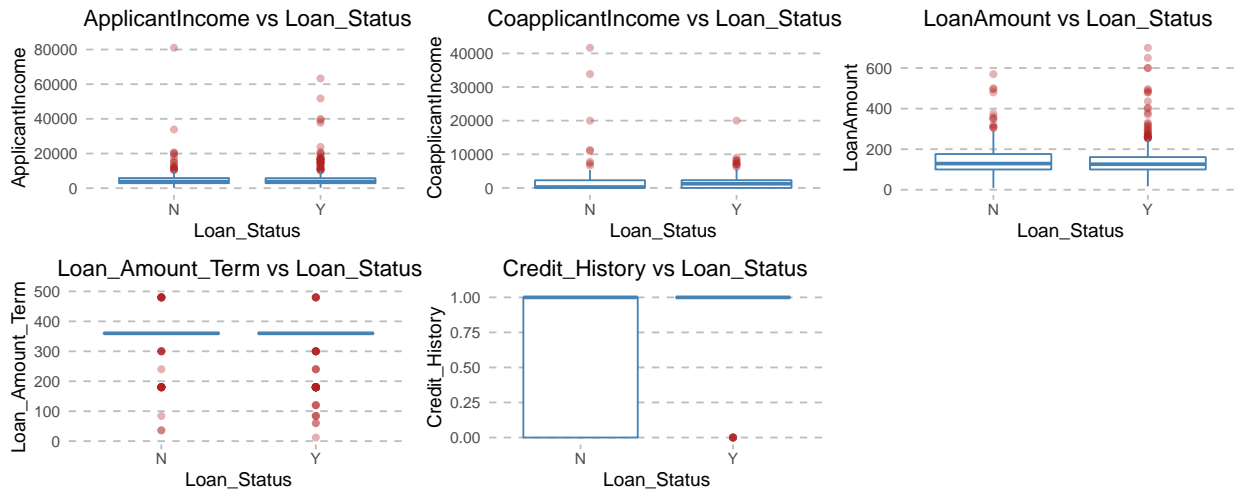
Gender	Married	Education	Self_Employed	Credit_History	Loan_Status
Female:112	No :213	Graduate :480	No :500	0 : 89	N:192
Male :489	Yes :398	Not Graduate:134	Yes : 82	1 :475	Y:422
NA's : 13	NA's: 3	NA	NA's: 32	NA's: 50	NA

```

                                linetype = "dashed"),
  panel.grid.major.x = element_blank(),
  panel.grid.minor.y = element_blank(),
  panel.grid.minor.x = element_blank(),
  axis.ticks.x = element_line(color = "grey")
)

}
myGlist <- within(myGlist, rm(target_name))
gridExtra::grid.arrange(grobs = myGlist, ncol = 3)

```



These charts are pretty hard to read, and again, we don't see any significant relationships between these variables and `Loan_Status`. We note that we need to change `Credit_History` to a categorical variable.

Categorical Features

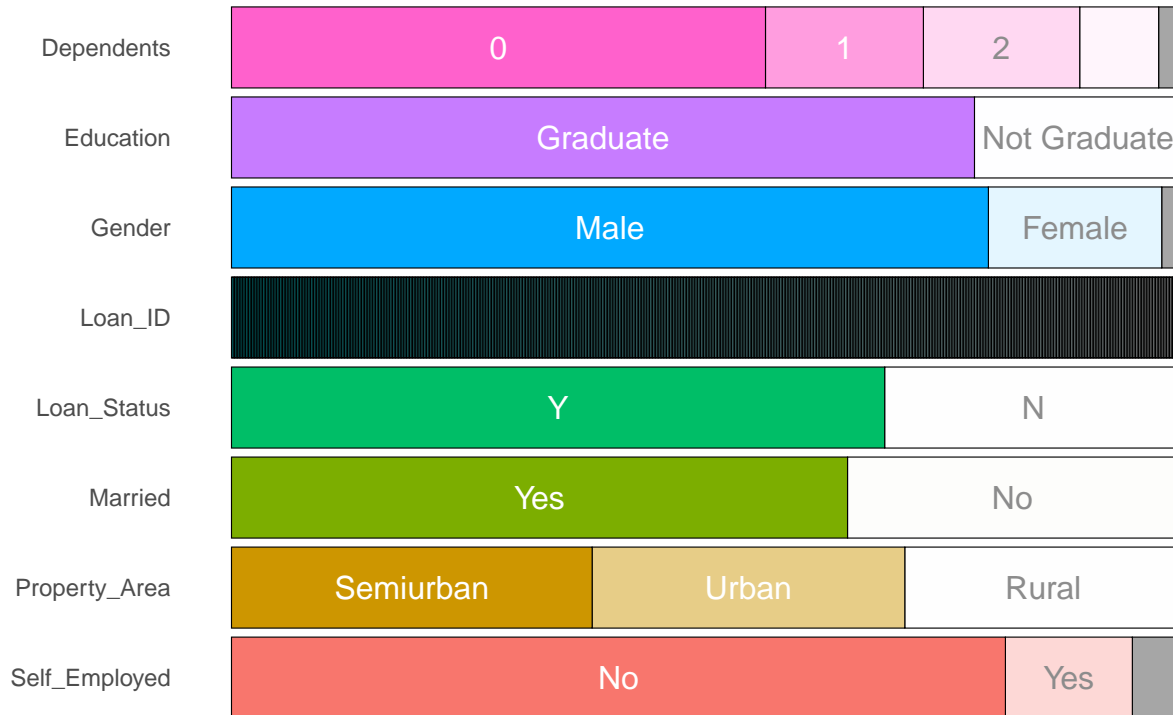
Next, let's turn our attention to our categorical variables. We have both binary categorical variables, and variables with 3 or more classes.

Dependents and Property_Area each comprise multiple categories. On the other hand, Gender, Married, Education, Self_Employed, Credit_History, Loan_Status are all binaries. These summary tables are helpful, although sometimes its easier to view this data in a visual format:

```
inspect_cat(loans) %>%
  show_plot()
```

Frequency of categorical levels in df::loans

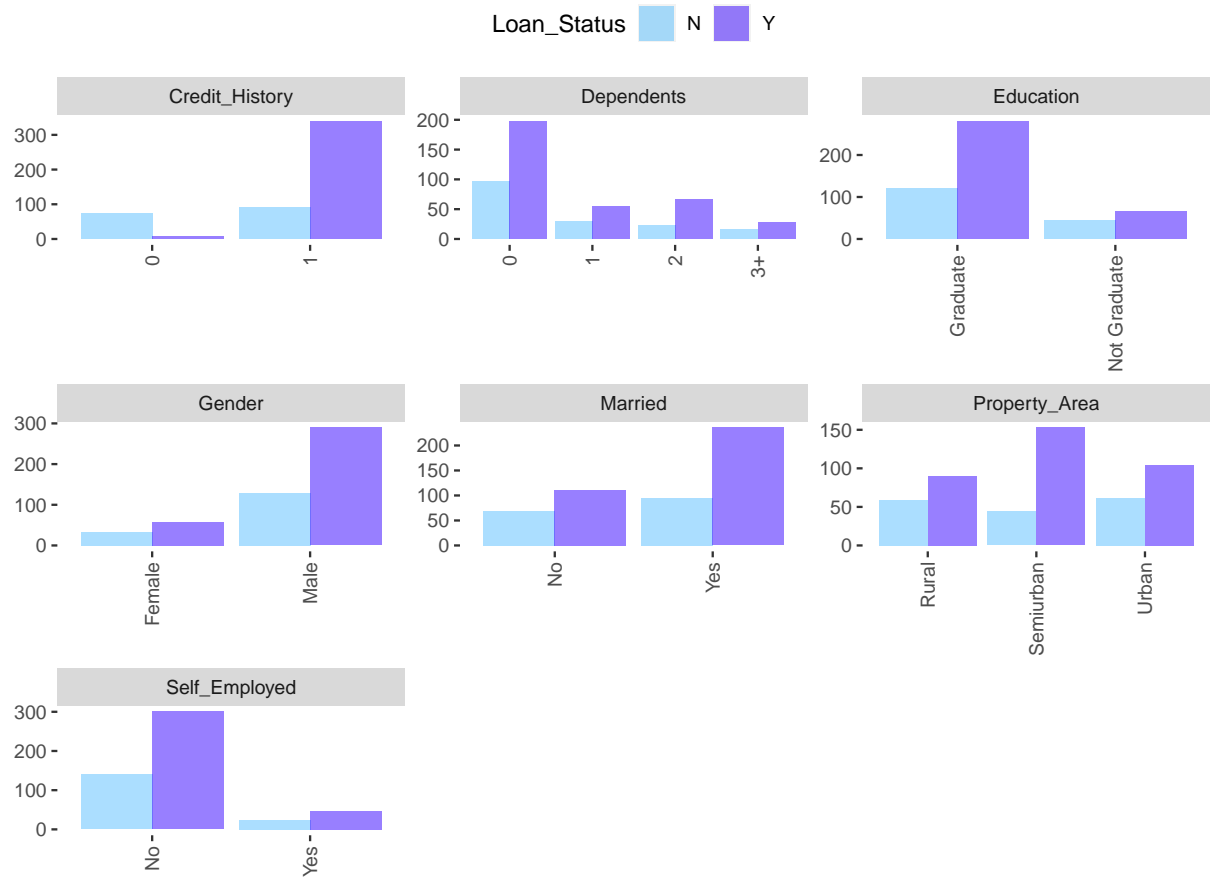
Gray segments are missing values



In looking at the chart and the summary tables, we see some very interesting things about the demographic of this dataset:

- Half of the population do not have any dependents
- Most people in this dataset are graduates
- Over 75% of the population is male
- 65% of the population is married
- Most of the individuals are not self-employed

We can also examine the dispersion of approval status between these variables.

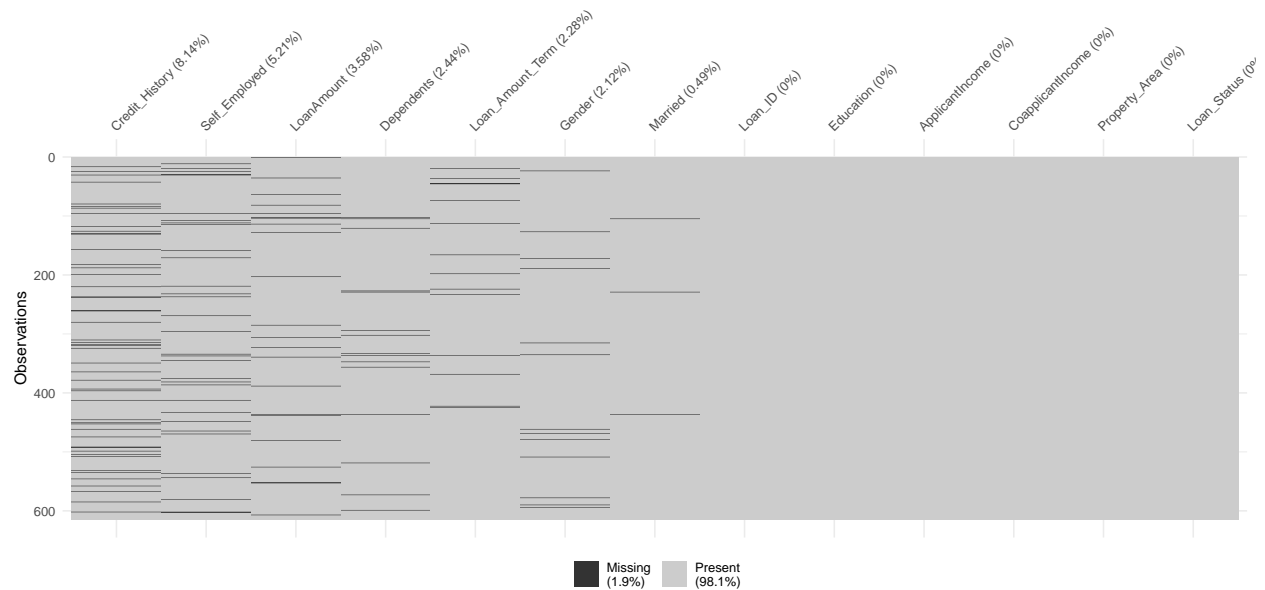


It looks like likelihoods are higher for applicants who have a credit history and who are married and live in an semi-urban area; as well as for those with a graduate degree.

Missing Data

The data appears to have some missing values that we'll have to deal with as we move along. Since we know there are missing values, let's see how pervasive the issue is:

```
visdat::vis_miss(loans, sort_miss = TRUE)
```

7 columns have missing data, ranging from 0.49% to 8.14%. What's interesting is there are 22 instances where `LoanAmount` is missing. In total, ~2% of the data has missing values.

Having completed our EDA, we'll now turn our attention to modeling.

Modeling

It is always valuable to try multiple modeling approaches to determine which model produces the most accurate results. We will proceed by building the following predictive models: LDA, KNN, a decision tree, and a random forest. Finally, we'll compare these models to determine which model is the most accurate.

LDA

First, we need to split the data into a training and test dataset. We'll do this using an 80/20 split.

```
set.seed(123)
training.samples <- loans2$Loan_Status %>%
  createDataPartition(p = 0.8, list = FALSE)
train.data <- loans2[training.samples, ]
test.data <- loans2[-training.samples, ]
```

The purpose of linear discriminant analysis is to find a combination of the variables that give the best possible separation between groups in our data set.

The LDA output indicates that our prior probabilities are:

- No = 0.309
- Yes = 0.690

In other words, 69% of the training observations are customers who got loan approval. It also provides the group means; these are the average of each predictor within each class, and are used by LDA as estimates. These suggest that customers that have approved applications, on average, have a lower loan amount, are

more likely to be male, to be married, have a credit history, and be from semi-urban areas. The coefficients of linear discriminant output provide the linear combination of variables that are used to form the LDA decision rule.

```
# Run LDA (from MASS library)
lda.loans <- MASS::lda(Loan_Status ~ Gender +Married +Dependents+Education +ApplicantIncome +CoapplicantIncome
                        +LoanAmount + Loan_Amount_Term + Credit_History +Property_Area +Self_Employed , data = train.data)
lda.loans

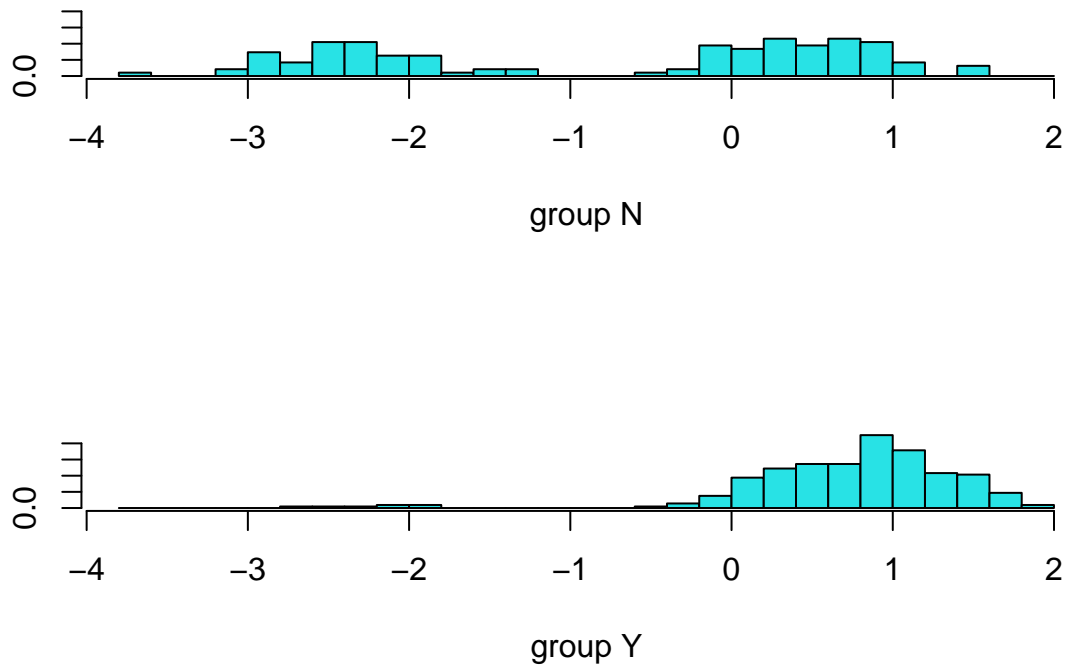
## Call:
## lda(Loan_Status ~ Gender + Married + Dependents + Education +
##      ApplicantIncome + CoapplicantIncome + LoanAmount + Loan_Amount_Term +
##      Credit_History + Property_Area + Self_Employed, data = train.data)
##
## Prior probabilities of groups:
##      N      Y
## 0.3090909 0.6909091
##
## Group means:
##      GenderMale MarriedYes Dependents1 Dependents2 Dependents3+
## N   0.7731092   0.5294118   0.1764706   0.1428571   0.09243697
## Y   0.8345865   0.6879699   0.1390977   0.2030075   0.08646617
##      EducationNot Graduate ApplicantIncome CoapplicantIncome LoanAmount
## N               0.2689076           5786.269           1768.336    150.0168
## Y               0.1842105           5018.222           1473.688    139.8083
##      Loan_Amount_Term Credit_History1 Property_AreaSemiurban Property_AreaUrban
## N           342.6555           0.5798319           0.2941176           0.3781513
## Y           340.3308           0.9736842           0.4586466           0.2894737
##      Self_EmployedYes
## N           0.1596639
## Y           0.1278195
##
## Coefficients of linear discriminants:
##                                LD1
## GenderMale                2.433286e-01
## MarriedYes                 6.018814e-01
## Dependents1               -4.065561e-01
## Dependents2                1.673478e-01
## Dependents3+              1.791003e-02
## EducationNot Graduate    -3.881465e-01
## ApplicantIncome          -5.456522e-06
## CoapplicantIncome        -5.389572e-05
## LoanAmount               -1.340982e-03
## Loan_Amount_Term         -8.934374e-04
## Credit_History1          2.965301e+00
## Property_AreaSemiurban   5.587254e-01
## Property_AreaUrban      -8.155911e-02
## Self_EmployedYes        -1.966809e-01
```

The linear discriminant function from the result above is:

$$0.24*GenderMale+0.602*MarriedYes-0.406*Dependents1+0.167*Dependents2+0.017*Dependents3-0.388*EducationNotGraduate+0.00000005456522*ApplicantIncome-0.00005389572*CoapplicantIncome-0.001340982*LoanAmount-0.0008934374*Loan_Amount_Term+2.965301*Credit_History1+0.5587254*Property_AreaSemiurban-0.08155911*Property_AreaUrban-0.1966809*Self_EmployedYes$$

We can use `plot()` to produce plots of the linear discriminants obtained by computing this formula for each training observation.

```
plot(lda.loans)
```



As you can see, when it's greater than 0, the probability increases that the customer will get approved.

Prediction accuracy of LDA compares prediction results from the model output with the actual data using the confusion matrix.

```
lda.predictions <- lda.loans %>% predict(test.data)
#Confusion Matrix
table(lda.predictions$class, test.data$Loan_Status)
```

```
##
##      N  Y
##  N 13  0
##  Y 16 66
```

```
#Prediction Accuracy
mean(lda.predictions$class == test.data$Loan_Status)
```

```
## [1] 0.8315789
```

It appears that LDA has a correct classification rate of 83%. We'll use this as a benchmark as we move forward in our algorithm selection process.

KNN

KNN requires that there is no missing data in the dataset. Based on our analysis, we know that we have NULL values in our dataset. To fill in the gaps, we can either impute or omit our missing values. Since we have already omitted these observations for LDA, we'll continue with that method here.

Following the required omission of the missing data, we need to do a little bit of pre-processing. KNN is highly susceptible to data that is on different scales (large values will be much further from each other than small values). With this in mind, we have chosen to center and scale all of our predictors. The final step of pre-processing is to remove any predictors that have near-zero variance so that there are no overlapping predictors.

```
knn_features <- subset(loans2, select=-c(Loan_Status))
knn_trans <- preProcess(knn_features,
                        method = c("center", "scale"))
knn_transformed_feat <- predict(knn_trans, knn_features)
nzv <- nearZeroVar(knn_transformed_feat, saveMetrics = TRUE)
nzv[nzv[, "nzv"] == TRUE,]
```

```
## [1] freqRatio    percentUnique zeroVar      nzv
## <0 rows> (or 0-length row.names)
```

It turns out that none of our predictors have near-zero variance, so we're good to proceed!

At this point, we are ready to build our model. We'll start by splitting our data into training and testing sets. We also need to remove Loan_ID since it will throw off our KNN model.

```
knn_processed <- cbind(loans2[,13], knn_transformed_feat)
knn_processed <- subset(knn_processed, select=-c(Loan_ID))
names(knn_processed)[1] <- ("Loan_Status")
knn_processed <- knn_processed[complete.cases(knn_processed),]
set.seed(54321)
train_ind <- sample(seq_len(nrow(knn_processed)),
                  size = floor(0.75*nrow(knn_processed)))
knn_train <- knn_processed[train_ind,]
knn_test <- knn_processed[-train_ind,]
```

Our KNN model uses the `kkn` library. With this library we are able to test different distances (Manhattan, Euclidean, etc.) as well as different weights (kernels).

```
kknn_func <- function(train_x, train_y, test_x, test_y){
  acc_df <- data.frame(matrix(nrow = 0, ncol = 4))

  weights <- c("rectangular", "triangular",
              "biweight", "triweight")

  for(d in 1:10){
    for(w in weights){
      for(i in 2:50){
        kknnModel <- kknn(train_y ~ .,
                          train_x,
                          test_x,
                          k = i,
```

```

        distance = d,
        kernel = w)

    cM <- table(test_y, fitted(kknnModel))
    accuracy <- (cM[1]+cM[4])/(cM[1]+cM[2]+cM[3]+cM[4])
    acc_df <- rbind(acc_df, c(i, accuracy, w, d))
  }
}
}
colnames(acc_df) <- c("k", "Accuracy", "Weight", "Distance")
acc_df[,1] <- as.integer(acc_df[,1])
acc_df[,2] <- as.numeric(acc_df[,2])
acc_df[,4] <- as.integer(acc_df[,4])
return(acc_df)
}
kknn_acc <- kknn_func(knn_train[,-1],
                     knn_train[,1],
                     knn_test[,-1],
                     knn_test[,1])
head(kknn_acc[order(-kknn_acc$Accuracy),], n = 10)

```

```

##      k Accuracy      Weight Distance
## 792   9 0.8250000 rectangular      5
## 793  10 0.8250000 rectangular      5
## 988   9 0.8250000 rectangular      6
## 989  10 0.8250000 rectangular      6
## 1184  9 0.8250000 rectangular      7
## 1185 10 0.8250000 rectangular      7
## 1380  9 0.8250000 rectangular      8
## 1381 10 0.8250000 rectangular      8
## 400   9 0.8166667 rectangular      3
## 401  10 0.8166667 rectangular      3

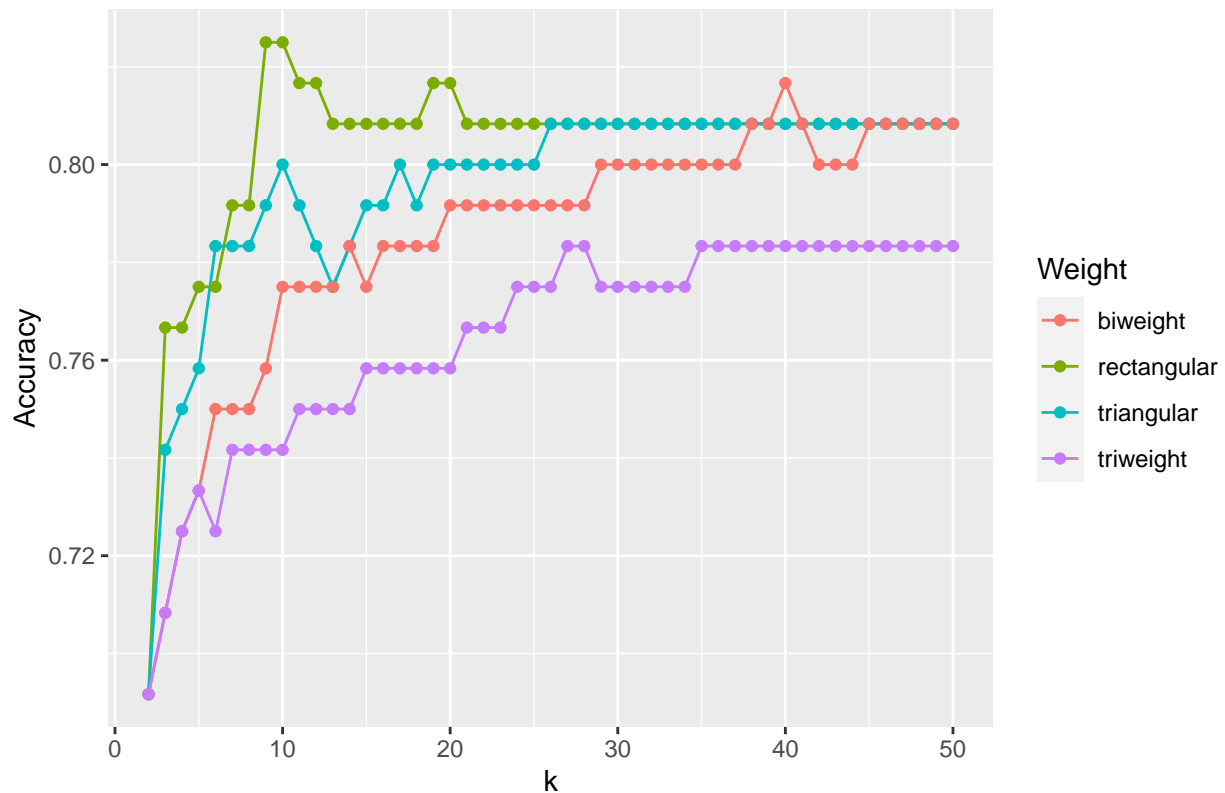
```

```

acc_plot_data <- kknn_acc[which(kknn_acc$Distance == 5),]
ggplot(data = acc_plot_data, aes(x = k, y = Accuracy, color = Weight)) +
  geom_line() +
  geom_point() +
  labs(title = "KKN: k distribution",
       x = "k",
       y = "Accuracy")

```

KKNN: k distribution



From the visual above, we can see that our model found that a k value of around 9 with a distance of 5 and a weighting function of rectangular produced the best model with an accuracy of 82.5%.

To validate these results, we'll randomize our training and testing sets one more time and check to see if the results are similar.

```
set.seed(12345)
train_ind <- sample(seq_len(nrow(knn_processed)),
                    size = floor(0.75*nrow(knn_processed)))
knn_train <- knn_processed[train_ind,]
knn_test <- knn_processed[-train_ind,]
kknnc_func <- function(train_x, train_y, test_x, test_y){
  acc_df <- data.frame(matrix(nrow = 0, ncol = 4))

  weights <- c("rectangular", "triangular",
               "biweight", "triweight")

  for(d in 1:10){
    for(w in weights){
      for(i in 2:50){
        kknncModel <- kknnc(train_y ~ .,
                           train_x,
                           test_x,
                           k = i,
                           distance = d,
                           kernel = w)
```

```

      cM <- table(test_y, fitted(kknnModel))
      accuracy <- (cM[1]+cM[4])/(cM[1]+cM[2]+cM[3]+cM[4])
      acc_df <- rbind(acc_df, c(i, accuracy, w, d))
    }
  }
}
colnames(acc_df) <- c("k", "Accuracy", "Weight", "Distance")
acc_df[,1] <- as.integer(acc_df[,1])
acc_df[,2] <- as.numeric(acc_df[,2])
acc_df[,4] <- as.integer(acc_df[,4])
return(acc_df)
}
kknn_acc <- kknn_func(knn_train[,-1],
                     knn_train[,1],
                     knn_test[,-1],
                     knn_test[,1])
head(kknn_acc[order(-kknn_acc$Accuracy),], n = 10)

```

```

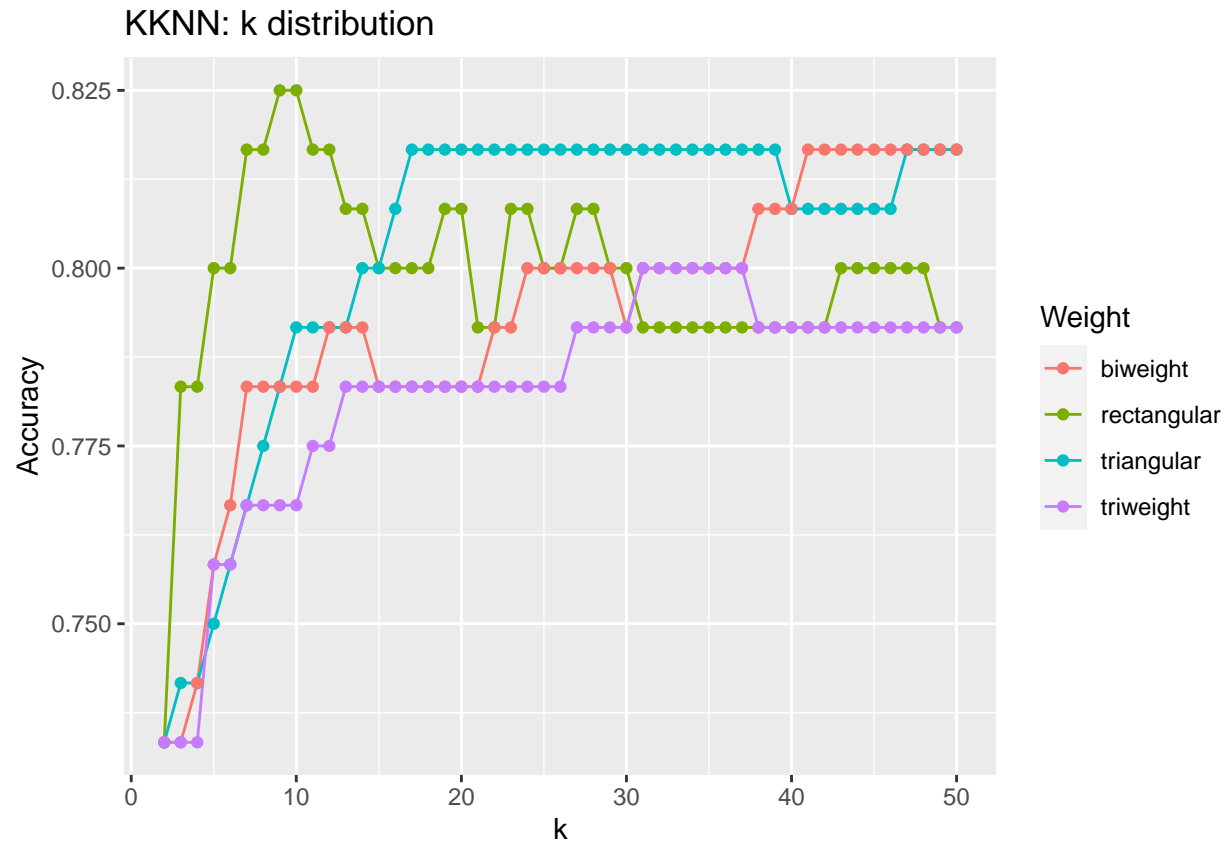
##      k Accuracy      Weight Distance
## 1380  9 0.8333333 rectangular      8
## 1381 10 0.8333333 rectangular      8
## 1576  9 0.8333333 rectangular      9
## 1577 10 0.8333333 rectangular      9
## 1772  9 0.8333333 rectangular     10
## 1773 10 0.8333333 rectangular     10
##  596  9 0.8250000 rectangular      4
##  597 10 0.8250000 rectangular      4
##  792  9 0.8250000 rectangular      5
##  793 10 0.8250000 rectangular      5

```

```

acc_plot_data <- kknn_acc[which(kknn_acc$Distance == 5),]
ggplot(data = acc_plot_data, aes(x = k, y = Accuracy, color = Weight)) +
  geom_line() +
  geom_point() +
  labs(title = "KKNN: k distribution",
       x = "k",
       y = "Accuracy")

```



After tuning our hyper parameters again, we arrived at a similar result: a k value of around 9 and a rectangular weighting provided us with the best version of a KNN model.

Decision Tree

A decision tree is a supervised machine learning algorithm that can be used for both classification and regression problems. A decision tree is simply a series of sequential decisions made to reach a specific result.

```
features <- c('Gender', 'Married', 'Dependents', 'Education', 'ApplicantIncome',
              'CoapplicantIncome', 'LoanAmount', 'Loan_Amount_Term', 'Credit_History',
              'Property_Area', 'Self_Employed')
```

```
dt <- rpart(Loan_Status~., data=train.data %>% dplyr::select(c(all_of(features), Loan_Status)), method = "class")
dt
```

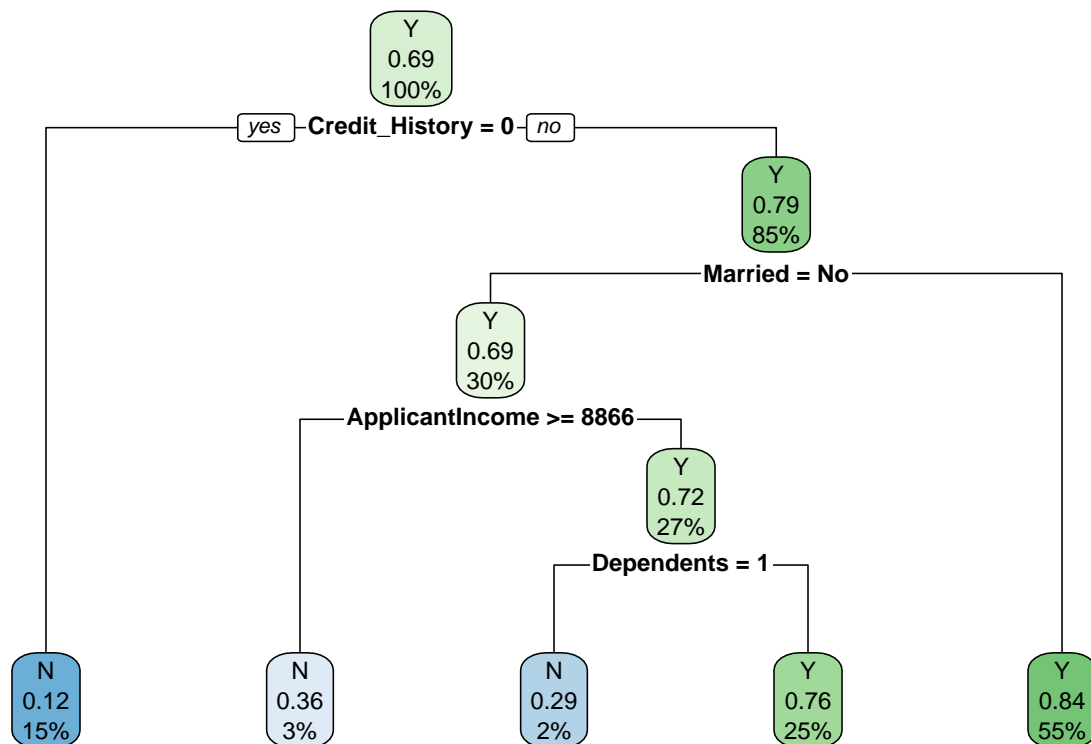
```
## n= 385
##
## node), split, n, loss, yval, (yprob)
##      * denotes terminal node
##
## 1) root 385 119 Y (0.3090909 0.6909091)
##    2) Credit_History=0 57 7 N (0.8771930 0.1228070) *
##    3) Credit_History=1 328 69 Y (0.2103659 0.7896341)
##      6) Married=No 116 36 Y (0.3103448 0.6896552)
```



```
##      12) ApplicantIncome>=8866 11    4 N (0.6363636 0.3636364) *
##      13) ApplicantIncome< 8866 105  29 Y (0.2761905 0.7238095)
##      26) Dependents=1 7    2 N (0.7142857 0.2857143) *
##      27) Dependents=0,2,3+ 98   24 Y (0.2448980 0.7551020) *
##      7) Married=Yes 212   33 Y (0.1556604 0.8443396) *
```

Running our data through a decision tree approves 70% of the loans and denies 30%.

```
rpart.plot(dt)
```



From the plot above we can see that having no credit history is a major factor in determining if the loan will be approved or not, followed by marriage, income and dependents. Upon closer examination, having a good credit history and being married gives the highest chance of being approved.

```
dt.predictions <- predict(dt, test.data, type='class')
confusionMatrix(table(dt.predictions, test.data$Loan_Status), positive='Y')
```

```
## Confusion Matrix and Statistics
##
##
## dt.predictions  N  Y
##                N 14  7
##                Y 15 59
##
##                Accuracy : 0.7684
```

```
##          95% CI : (0.6706, 0.8488)
##    No Information Rate : 0.6947
##    P-Value [Acc > NIR] : 0.07119
##
##          Kappa : 0.4083
##
##    McNemar's Test P-Value : 0.13559
##
##          Sensitivity : 0.8939
##          Specificity : 0.4828
##          Pos Pred Value : 0.7973
##          Neg Pred Value : 0.6667
##          Prevalence : 0.6947
##          Detection Rate : 0.6211
##    Detection Prevalence : 0.7789
##          Balanced Accuracy : 0.6883
##
##          'Positive' Class : Y
##
```

The confusion matrix shows we correctly predicted Yes 59 times and No 14 times for an overall accuracy of 76.84%. Our model has higher precision than recall. This means that our model is better at predicting if someone qualifies for a loan than predicting if someone should be denied a loan.

Random Forest

The decision tree algorithm is quite easy to understand and interpret. However, often a single tree is not sufficient for producing effective results. This is where the Random Forest algorithm comes into the picture.

We'll build a random forest model with 500 trees.

```
rf <- randomForest(Loan_Status~., data=train.data %>% dplyr::select(c(all_of(features), Loan_Status)), n
rf
```

```
##
## Call:
## randomForest(formula = Loan_Status ~ ., data = train.data %>%          dplyr::select(c(all_of(features)
##          Type of random forest: classification
##          Number of trees: 500
## No. of variables tried at each split: 3
##
##          OOB estimate of  error rate: 20.78%
## Confusion matrix:
##      N   Y class.error
## N 55  64  0.53781513
## Y 16 250  0.06015038
```

```
rf.predictions <- predict(rf, test.data, type='class')
confusionMatrix(table(rf.predictions, test.data$Loan_Status), positive='Y')
```

```

## Confusion Matrix and Statistics
##
##
## rf.predictions  N  Y
##                N 15  7
##                Y 14 59
##
##                Accuracy : 0.7789
##                95% CI : (0.6822, 0.8577)
##                No Information Rate : 0.6947
##                P-Value [Acc > NIR] : 0.04436
##
##                Kappa : 0.441
##
## Mcnemar's Test P-Value : 0.19043
##
##                Sensitivity : 0.8939
##                Specificity : 0.5172
##                Pos Pred Value : 0.8082
##                Neg Pred Value : 0.6818
##                Prevalence : 0.6947
##                Detection Rate : 0.6211
##                Detection Prevalence : 0.7684
##                Balanced Accuracy : 0.7056
##
##                'Positive' Class : Y
##

```

The confusion matrix shows we correctly predicted Yes 58 times and No 15 times. Our model has higher recall than precision. This means that our model is more accurate at predicting those who do not qualify for a loan than those who do.

Model Comparison

Here we summarize the accuracy of each of our models:

Model	Accuracy
LDA	0.8316
KNN	0.8250
Decision Tree	0.7684
Random Forest	0.7789

Based on the prediction accuracy, LDA appears to be the most accurate model. With the goal of building the most accurate model, We would select this model above the others.

However, if the goal of the model was accuracy as well as interpretability, we would potentially select one of the tree based methods as there are many tools available to see visually how the model is making predictions. Based on our summary table above, it appears that both the decision tree and random forest model have fairly similar accuracy. Let's break down these metrics to see if one is a clear winner over the other.

	DT_Model	RF_Model
Sensitivity	0.4827586	0.5172414
Specificity	0.8939394	0.8939394
Pos Pred Value	0.6666667	0.6818182
Neg Pred Value	0.7972973	0.8082192
Precision	0.6666667	0.6818182
Recall	0.4827586	0.5172414
F1	0.5600000	0.5882353
Prevalence	0.3052632	0.3052632
Detection Rate	0.1473684	0.1578947
Detection Prevalence	0.2210526	0.2315789
Balanced Accuracy	0.6883490	0.7055904

```
DT_Model <- confusionMatrix(dt.predictions, test.data$Loan_Status)$byClass
DT_Model <- data.frame(DT_Model)

RF_Model <- confusionMatrix(rf.predictions, test.data$Loan_Status)$byClass
RF_Model <- data.frame(RF_Model)

compare <- data.frame(DT_Model, RF_Model)

compare %>% kableExtra::kbl() %>% kable_styling(bootstrap_options = c("striped", "hover", "condensed"),
```

We can see from the above that the random forest has better performance than the decision tree model for many of the metrics like Sensitivity, Neg Pred Value, and Balanced Accuracy. Based on these metrics, we would conclude that the random forest model would generate more accurate predictions.

Conclusion and Next Steps

Based on our extensive EDA and modeling, LDA was chosen as our preferred predictive model based on accuracy. However, if interpretability was a large factor, we would consider selection the random forest model.

Next steps for this analysis may include such things as trying different imputation methods for the nulls within the dataset and trying more advanced algorithms like XGBOOST to see if we could gain additional accuracy.