# Data 622 Homework 2

Brandon Cunningham, Jean Jimenez, Chafiaa Nadour, Shri Tripathi

2024-11-24

## Assignment Instructions

Perform an analysis of the dataset(s) used in Homework #2 using the SVM algorithm. Compare the results with the results from previous homework.

## Literature Review

Three Articles on SVM vs Decision Tree on Diabetes:

M. F. Faruque, Asaduzzaman and I. H. Sarker, "Performance Analysis of Machine Learning Techniques to Predict Diabetes Mellitus," 2019 International Conference on Electrical, Computer and Communication Engineering (ECCE), Cox'sBazar, Bangladesh, 2019, pp. 1-4, doi: 10.1109/ECACE.2019.8679365. keywords: {Diabetes;Machine learning;Support vector machines;Decision trees;Diseases;Machine learning algorithms;Classification algorithms;eHealth;diabetes;machine learning;prediction}

L. Ismail and H. Materwala, "Comparative Analysis of Machine Learning Models for Diabetes Mellitus Type 2 Prediction," 2020 International Conference on Computational Science and Computational Intelligence (CSCI), Las Vegas, NV, USA, 2020, pp. 527-533, doi: 10.1109/CSCI51800.2020.00095. keywords: {Support vector machines;Analytical models;Scientific computing;Computational modeling;Machine learning;Predictive models;Feature extraction;artificial intelligence;classification models;diabetes mellitus type 2;health informatics;machine learning models}

C. Sonawane, K. Somwanshi, R. Patil and R. Raut, "Diabetic Prediction Using Machine Algorithm SVM and Decision Tree," 2023 International Conference on Applied Intelligence and Sustainable Computing (ICAISC), Dharwad, India, 2023, pp. 1-7, doi: 10.1109/ICAISC58445.2023.10199539. keywords: {Support vector machines;Analytical models;Machine learning algorithms;Predictive models;Prediction algorithms;Diabetes;Glucose;Diabetic Prediction;SVM;Decision Tree}

The application of machine learning algorithms for diabetes prediction has been extensively studied, including compariosn between SVM and Decision Trees. Recent research demonstrates varying levels of effectiveness between these approaches.

Faruque et al. (2019) conducted a comparative analysis of machine learning techniques for diabetes prediction, finding that C4.5 decision tree achieved higher accuracy compared to other algorithms including SVM. Their study on a dataset of 200 patients revealed that decision trees provided better interpretability of results, which is crucial for medical applications. However, they noted that SVM showed promise in handling high-dimensional medical data.

Ismail and Materwala (2020) performed a more extensive comparison using the UCI diabetes dataset with 65,840 observations. Their research demonstrated that SVM with appropriate feature selection could achieve comparable accuracy to decision trees while being more robust to outliers. They emphasized the importance of evaluating models using multiple metrics beyond just accuracy, particularly for imbalanced medical datasets.

Sonawane et al. (2023) provided the most recent analysis, comparing SVM and decision tree algorithms specifically for diabetes prediction. Their results showed SVM slightly outperforming decision trees, with accuracy rates of 76.6% versus 75% respectively. This study particularly highlighted SVM's strength in handling non-linear relationships in medical data and its ability to maintain performance with high-dimensional feature spaces.

# Pre-Work

## Large Dataset (Diabetes Data)

For our large dataset, we will use the diabetes dataset from kaggle.

This dataset has 100k clinical records of diabetes for health analytic purposes.

Link to Dataset

Goal: For this dataset, we want to predict whether or not the patient will have diabetes. In other words, this is a classification problem where we have to correctly label data. We will do this using decision trees and SVM.

Importing:

```
diabetes_data=read.csv(url("https://raw.githubusercontent.com/sleepysloth12/DATA_622_HW01/refs/heads/ma
```

### Exploratory Data Analysis

The column of interest, labeled `diabetes` is what we want to predict. It is an integer, 0 or 1, indicating if the patient has diabetes or not. In the current dataset, 91% of the patients have no diabetes and 8.5% of the patients have diabetes.

In order to build a predictive model, we must first go column by column and clean up the features a little bit to make this more accurate/applicable to healthcare data.

### Data Cleaning

**Year**   The first column is year. The dataset is timeseries data, collected from the years 2015-2022. However, each year has different numbers of observations. There is no way of knowing if this is longtitudnal data (one patient visited multiple year) due to the lack of unique patient identifier field. I think we can completely disregard and forget about this column.

```
diabetes_data = diabetes_data %>%
  select(-year)
```

**Gender**   Next is gender. Gender is pretty even split, with ~60% being female and ~40% being male. There is an insignificant amount of people that answered "other" (less than 1%).

I'm going ahead and going to filter out other. Also, I am going to change the label to`is_female` so the choice is binary.
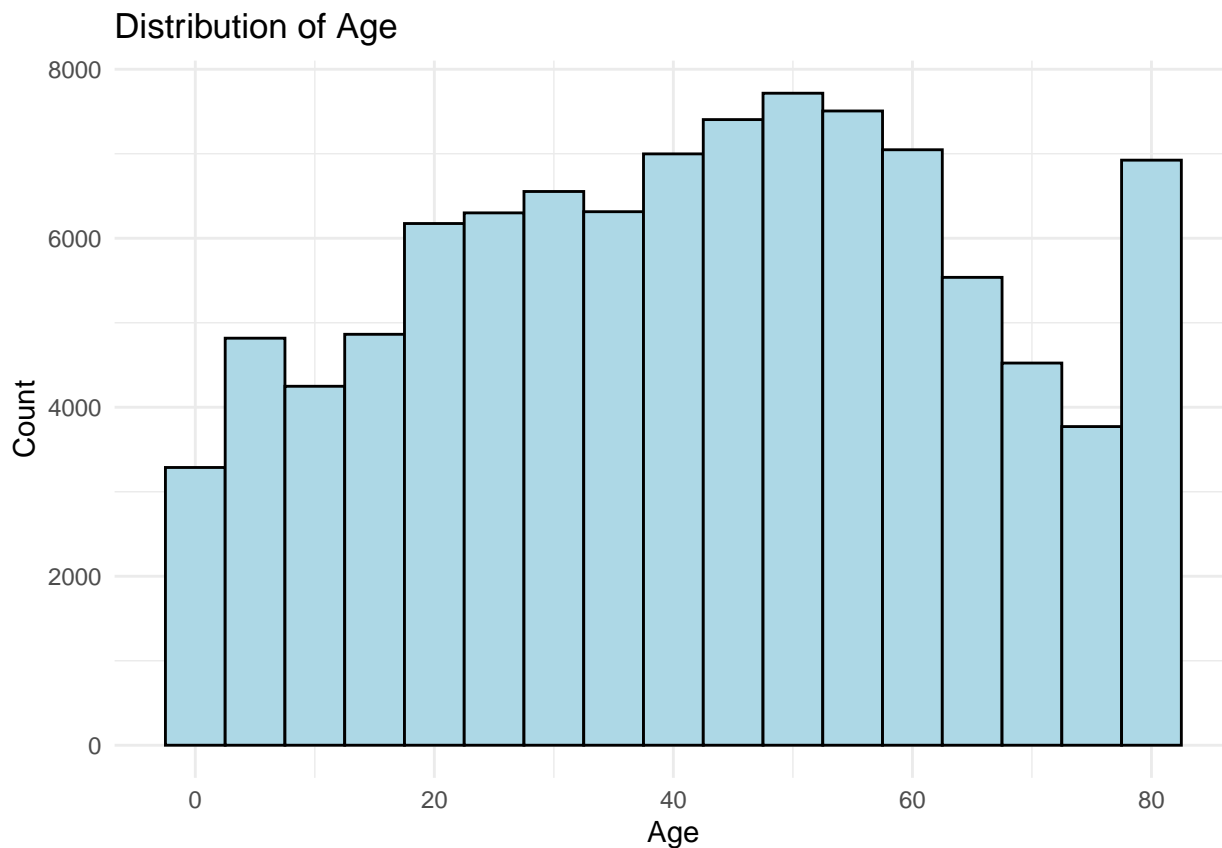
```
diabetes_data = diabetes_data %>%
  filter(gender == "Female"| gender== "Male")%>%
  mutate(is_female=ifelse(gender == "Female",1,0))%>%
  select(-gender)
```

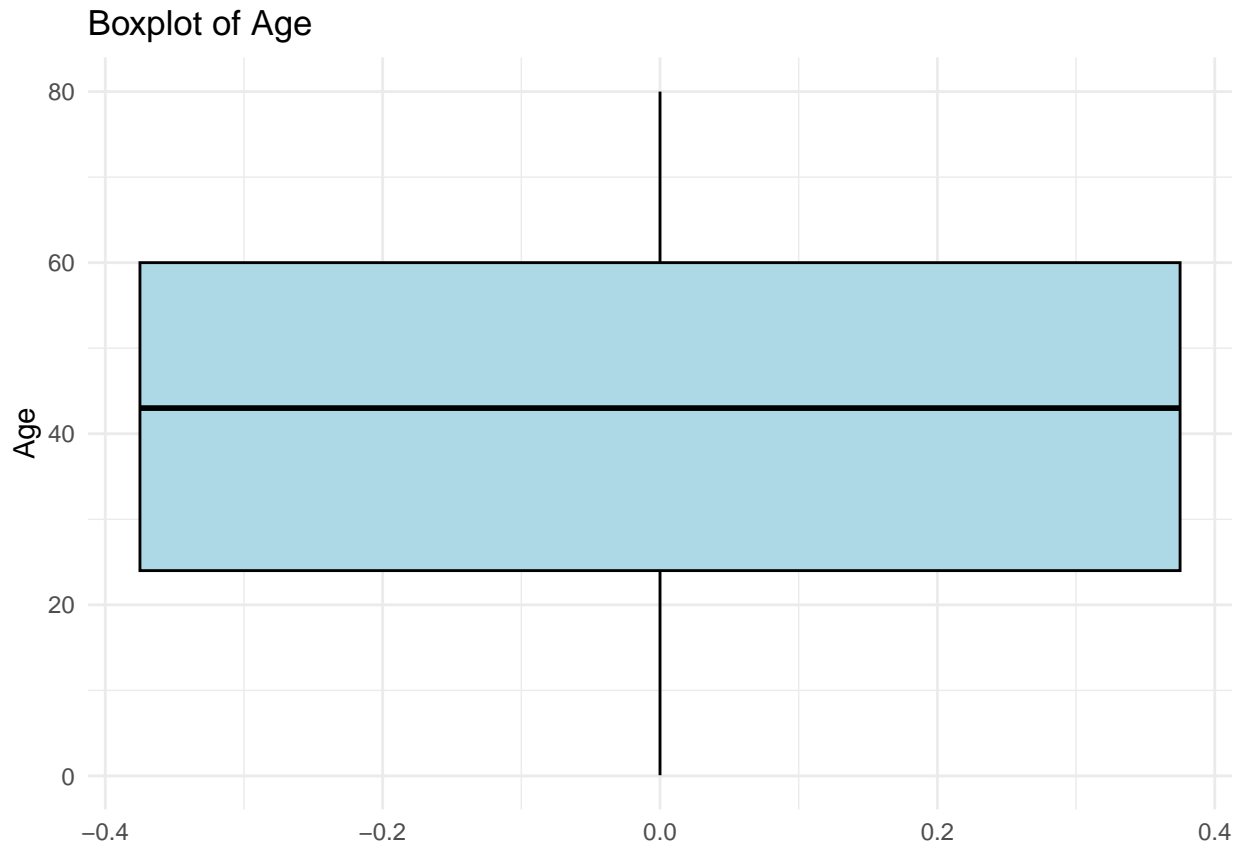**Age**   Next is age. Mean age is 41.9 years old, with a standard deviation of +/-22.5 years old.

Max age is 80.

Minimum recorded age is 0.08. This might be an outlier. Therefore, lets visualize this distribution in both box plot and bar plot.

```
ggplot(diabetes_data, aes(x = age)) +
  geom_histogram(binwidth = 5, color = "black", fill = "lightblue") +
  labs(title = "Distribution of Age", x = "Age", y = "Count") +
  theme_minimal()
```



Distribution of Age

```
ggplot(diabetes_data, aes(y = age)) +
  geom_boxplot(fill = "lightblue", color = "black") +
  labs(title = "Boxplot of Age", y = "Age") +
  theme_minimal()
```

## Boxplot of Age



Seems like the minimum age is an outlier. In medical research, we tend to separate adult populations from pediatric populations so lets go ahead and do that here. Lets only look at 18+.

In terms of the age distribution, it looks relatively normal. Diabetes incidence seem to increase as you get closer to middle age, then decrease. There is a spike at 80 years old.

I am going to bin age/ convert it into different categories:

is_young = Age 18-35

is_middle_age = Age 36-64

is_old = Age 65+

```
diabetes_data = diabetes_data %>%
  filter(age>=18)%>%
  mutate(is_young=ifelse(age>=18 & age <=35, 1,0),
         is_middle_age=ifelse(age>35 & age<65 , 1 , 0),
         is_old=ifelse(age>65,1,0))%>%
  select(-age)
```

```
length(unique(diabetes_data$location))
```

**State**

```
## [1] 55
```

For the location column, there are 55 different locations, corresponding to the 50 different states and territory.

Location is important for diabetes prediction. Some areas are probably more likely to develop diabetes than others. Like age, I want to create categories and bin them based on the location. Then, will create dummy

variables.

```
diabetes_data = diabetes_data %>%
  mutate(

    is_new_england = if_else(location %in% c("Connecticut", "Maine", "Massachusetts",
                                  "New Hampshire", "Rhode Island", "Vermont"), 1, 0),


    is_south = if_else(location %in% c("Alabama", "Arkansas", "Delaware", "Florida", "Georgia",
                             "Kentucky", "Louisiana", "Maryland", "Mississippi",
                             "North Carolina", "Oklahoma", "South Carolina",
                             "Tennessee", "Texas", "Virginia", "West Virginia"), 1, 0),


    is_midwest = if_else(location %in% c("Illinois", "Indiana", "Iowa", "Kansas", "Michigan",
                               "Minnesota", "Missouri", "Nebraska", "North Dakota",
                               "Ohio", "South Dakota", "Wisconsin"), 1, 0),


    is_west = if_else(location %in% c("Alaska", "Arizona", "California", "Colorado", "Hawaii",
                            "Idaho", "Montana", "Nevada", "New Mexico", "Oregon",
                            "Utah", "Washington", "Wyoming"), 1, 0),


    is_northeast = if_else(location %in% c("New Jersey", "New York", "Pennsylvania"), 1, 0),


    is_territories = if_else(location %in% c("Guam", "Puerto Rico", "Virgin Islands",
                                   "District of Columbia", "United States"), 1, 0)
  )%>%
  select(-location)
```

**Race, Ethnicity, Hypertension, & Heart Disease**   Race and ethnicity is already binned and with their individual dummy variables. Race and ethnicity are both factors that influence diabetes so will leave these columns untouched.
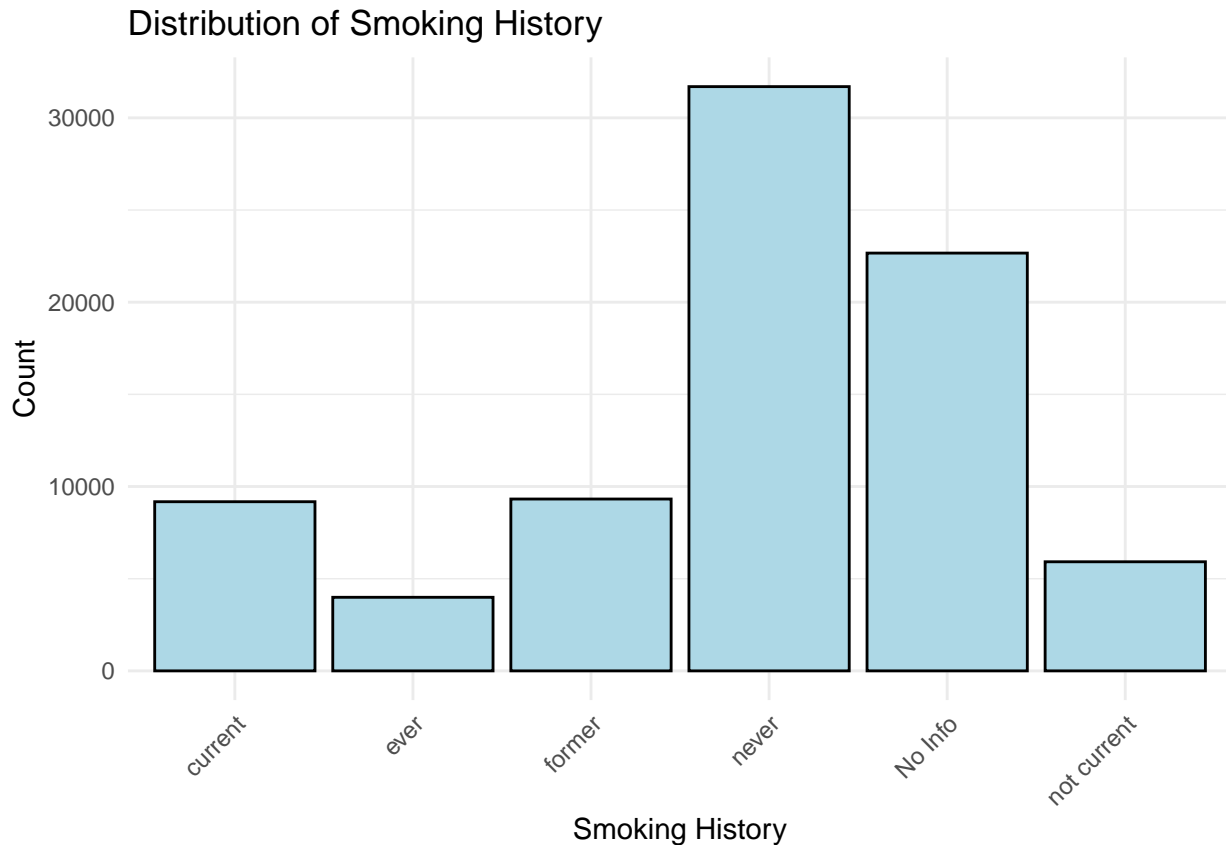
Same with the columns of hypertension and heart disease.

**Smoking History**   There are currently 6 categories/ choices patients could respond when asked about smoking history:

```
unique(diabetes_data$smoking_history)
```

```
## [1] "never"       "not current" "current"     "No Info"     "ever"
## [6] "former"
```

```
ggplot(diabetes_data, aes(x = smoking_history)) +
  geom_bar(fill = "lightblue", color = "black") +
  labs(title = "Distribution of Smoking History", x = "Smoking History", y = "Count") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

## Distribution of Smoking History



The biggest group is Never smoked accounting for 35% of the data.

There is a category, 'ever' which is 'Never' mislabeled. Will fix this. Once combined, never smoked will account for40% of the data.

The second biggest is 'No info' with near 35% of the data. Since the people in 'No info' may or may not be smokers, if we leave this category in it might make our predictions inaccurate. We want to capture how smoking can influence diabetes, therefore we ill remove this group.

Also, the 'not current' and 'former' group can be combined.

```
diabetes_data = diabetes_data %>%
  filter(smoking_history!="No Info")%>%
  mutate(never_smoked=ifelse(smoking_history %in% c("ever","never"),1,0),
         former_smoker=ifelse(smoking_history %in% c("former","not current"),1,0),
         current_smoker=ifelse(smoking_history=="current",1,0)) %>%
  select(-smoking_history)
```

**Biomarker Columns**    The distribution of BMI is normal. It is numeric and continuous. We are leaving this as is.

The hbA1c_level biomarker, although numeric, has 18 unique values. In healthcare, this biomarker is usually used to determine diabetes. We will bin this biomarker for the following categories:

A1c < 5.7% –> Normal A1C

A1c between 5.7-6.4 % –> PreDiabetes

A1C over 6.5% –> diabetes

Although, correlation analysis is needed. There might be multicollinearity between these biomarker variables.

I say this because blood glucose variable and A1c directly related to each other.

Actually going to remove blood glucose because having that and A1C is repetitive/ multicollinearity.

```
diabetes_data = diabetes_data %>%
  mutate(normal_a1c=ifelse(hbA1c_level<5.7,1,0),
         prediabetic_a1c=ifelse(hbA1c_level>=5.7 & hbA1c_level <= 6.4,1,0),
         diabetic_a1c=ifelse(hbA1c_level>6.4,1,0))%>%
  select(-c(hbA1c_level,blood_glucose_level))
```

# Building the Models

Now that our dataset is clean, we can discuss what model we want to use.

The target variable to predict is diabetes (binary choice whether or not patient will have diabetes).

The methodology employed in this analysis follows a systematic approach to data preprocessing and model implementation. The dataset underwent careful cleaning and transformation, including handling missing values and encoding categorical variables. The implementation used both SVM and decision tree algorithms on the same dataset to ensure fair comparison. The SVM model utilized different kernel functions (radial, polynomial, and linear) to identify the optimal approach for diabetes prediction.

## Correlation Matrix

Before building the models, I want to run a correlation matrix to look for multicolinearity

```
predictors <- diabetes_data %>%
  select(-diabetes)

cor_matrix <- cor(predictors)

high_cor <- findCorrelation(cor_matrix, cutoff = 0.7, verbose = TRUE)

## Compare row 19  and column  20 with corr  0.705
##    Means:  0.07 vs 0.049 so flagging column 19
## All correlations <= 0.7
cat("Highly correlated variables:", paste(names(predictors)[high_cor], collapse = ", "), "\n")

## Highly correlated variables: never_smoked
corrplot(cor_matrix, method = "color", type = "upper", order = "hclust",
         tl.col = "black", tl.srt = 45)
```
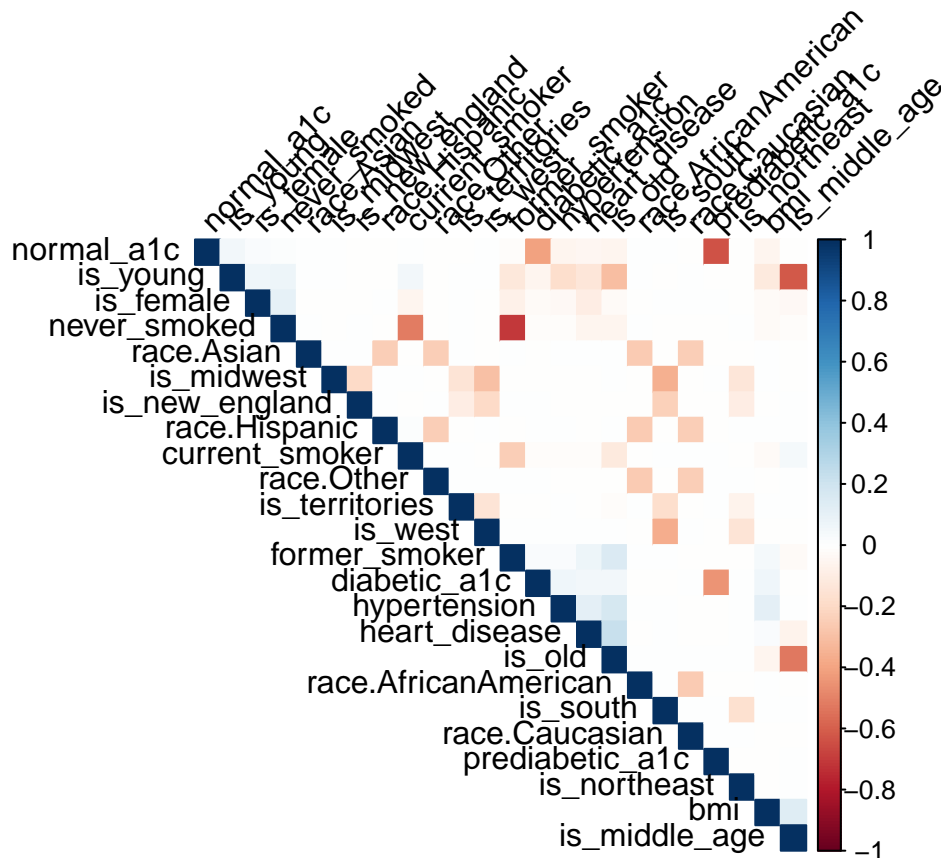
There is some multicollinearity
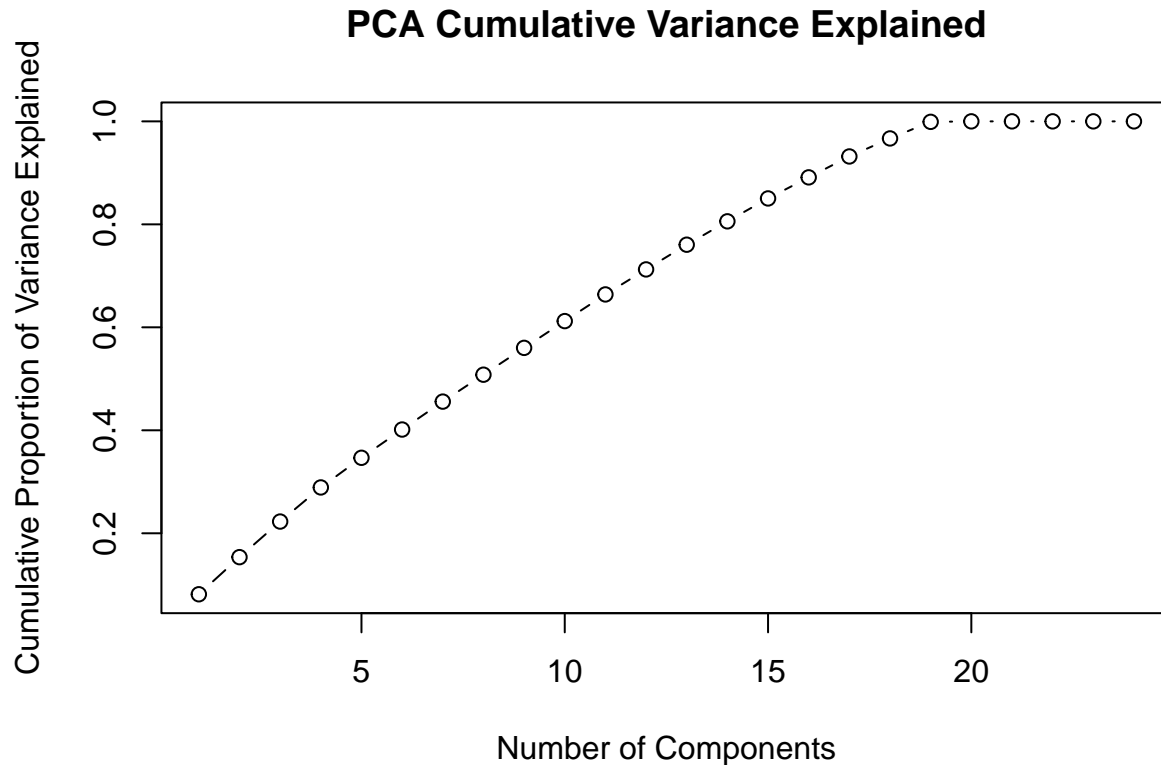
**Principal Component Analysis**

Conducting a PCA to determine the important components

```
pca_result <- prcomp(predictors, scale. = TRUE)
summary(pca_result)
```

```
## Importance of components:
##                           PC1     PC2     PC3     PC4     PC5     PC6     PC7
## Standard deviation     1.39898  1.3155 1.28787 1.26107 1.17547  1.1478 1.14083
## Proportion of Variance 0.08155  0.0721 0.06911 0.06626 0.05757  0.0549 0.05423
## Cumulative Proportion  0.08155  0.1537 0.22276 0.28902 0.34659  0.4015 0.45572
##                           PC8     PC9    PC10    PC11    PC12    PC13    PC14
## Standard deviation     1.11981 1.11834 1.11722 1.11510 1.07939 1.07351  1.0438
## Proportion of Variance 0.05225 0.05211 0.05201 0.05181 0.04855 0.04802  0.0454
## Cumulative Proportion  0.50797 0.56008 0.61209 0.66390 0.71244 0.76046  0.8059
##                           PC15    PC16    PC17    PC18    PC19    PC20     PC21
## Standard deviation     1.03300 0.99101 0.98754 0.91645 0.87923 0.14869 3.24e-13
## Proportion of Variance 0.04446 0.04092 0.04063 0.03499 0.03221 0.00092 0.00e+00
## Cumulative Proportion  0.85032 0.89124 0.93187 0.96687 0.99908 1.00000 1.00e+00
##                           PC22      PC23     PC24
## Standard deviation     2.767e-14 2.115e-14 1.26e-14
## Proportion of Variance 0.000e+00 0.000e+00 0.00e+00
## Cumulative Proportion  1.000e+00 1.000e+00 1.00e+00
```

```
plot(cumsum(pca_result$sdev^2 / sum(pca_result$sdev^2)),
     type = "b",
     xlab = "Number of Components",
     ylab = "Cumulative Proportion of Variance Explained",
     main = "PCA Cumulative Variance Explained")
```

## PCA Cumulative Variance Explained



Number of Components

Our first principal component only accounts for about 8.16% of the total variance. That's not a lot. It means no single factor dominates in predicting diabetes. This makes sense given the complex nature of the disease and the variety of factors we've included in our dataset.

We need 11 components to explain about 66% of the variance, and it takes 19 to get to nearly 100%. Looking at our cumulative variance plot, we can see this gradual climb. The fact that we need so many components to explain most of the variance suggests we shouldn't try to oversimplify our model. Most of our variables are contributing unique information about diabetes risk.

While we don't see extreme multicollinearity, there is some correlation among our variables. We can explain about 85% of the variance with 15 components, which is fewer than our original variables.

Here I am getting a subset to use for training as these models are compute intensive

```
set.seed(622)

subset_size <- 0.1 * nrow(diabetes_data)
subset_idx <- sample(1:nrow(diabetes_data), subset_size)

diabetes_data_subset <- diabetes_data[subset_idx, ]
```

```
set.seed(622)

diabetes_data_subset$diabetes <- as.factor(diabetes_data_subset$diabetes)
```

```
train_split_idx=createDataPartition(diabetes_data_subset$diabetes, p=0.7, list=FALSE)

train_diab = diabetes_data_subset[train_split_idx,]
test_diab = diabetes_data_subset[-train_split_idx,]

control <- trainControl(method = "cv", number = 5)
metric <- "Accuracy"
```

# Models

## SVM Model 1 - Radial Kernel with Auto-tuning

```
set.seed(622)

svm_model1 <- train(diabetes ~ .,
                    data = train_diab,
                    method = "svmRadial",
                    metric = metric,
                    trControl = control,
                    tuneLength = 15)

print(svm_model1)
```

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 4208 samples
##   24 predictor
##    2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 3367, 3366, 3365, 3367, 3367
## Resampling results across tuning parameters:
##
##   C         Accuracy   Kappa
##      0.25   0.8797531  0.003421095
##      0.50   0.8875944  0.168729845
##      1.00   0.8913963  0.250154459
##      2.00   0.8932966  0.298931298
##      4.00   0.8897317  0.306493269
##      8.00   0.8837911  0.305190853
##     16.00   0.8792761  0.314591698
##     32.00   0.8778498  0.330257060
##     64.00   0.8704824  0.312271496
##    128.00   0.8645405  0.305346593
##    256.00   0.8643041  0.315656360
##    512.00   0.8574140  0.300962097
##   1024.00   0.8509945  0.284266135
##   2048.00   0.8502862  0.281428701
##   4096.00   0.8441045  0.265512658
##
## Tuning parameter 'sigma' was held constant at a value of 0.02523945
## Accuracy was used to select the optimal model using the largest value.
```

```
## The final values used for the model were sigma = 0.02523945 and C = 2.
```

```r
model1_predictions <- predict(svm_model1, newdata = test_diab)

conf_matrix <- confusionMatrix(model1_predictions, test_diab$diabetes)

print(conf_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 1566  174
##          1   19   43
##
##                Accuracy : 0.8929
##                  95% CI : (0.8777, 0.9068)
##     No Information Rate : 0.8796
##     P-Value [Acc > NIR] : 0.04277
##
##                   Kappa : 0.2691
##
##  Mcnemar's Test P-Value : < 2e-16
##
##             Sensitivity : 0.9880
##             Specificity : 0.1982
##          Pos Pred Value : 0.9000
##          Neg Pred Value : 0.6935
##              Prevalence : 0.8796
##          Detection Rate : 0.8690
##    Detection Prevalence : 0.9656
##       Balanced Accuracy : 0.5931
##
##        'Positive' Class : 0
##
```

This model achieved the highest accuracy at 89.29% with optimal parameters of sigma = 0.02523945 and C = 2. The model showed excellent sensitivity (98.80%) but poor specificity (19.82%), indicating it was very good at identifying non-diabetic cases but struggled with diabetic cases. The model went through 15 different C values, showing peak performance at moderate regularization levels. The balanced accuracy of 59.31% suggests some class imbalance issues.

## SVM Model 2 - Polynomial Kernel

```r
set.seed(622)

svm_model2 <- train(diabetes ~ .,
                    data = train_diab,
                    method = "svmPoly",
                    metric = metric,
                    trControl = control)

print(svm_model2)
```

```
## Support Vector Machines with Polynomial Kernel
```

```
## 
## 4208 samples
##   24 predictor
##    2 classes: '0', '1'
## 
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 3367, 3366, 3365, 3367, 3367
## Resampling results across tuning parameters:
## 
##   degree  scale  C     Accuracy   Kappa
##   1       0.001  0.25  0.8795156  0.00000000
##   1       0.001  0.50  0.8795156  0.00000000
##   1       0.001  1.00  0.8795156  0.00000000
##   1       0.010  0.25  0.8795156  0.00000000
##   1       0.010  0.50  0.8795156  0.00000000
##   1       0.010  1.00  0.8795156  0.00000000
##   1       0.100  0.25  0.8797534  0.01711428
##   1       0.100  0.50  0.8826024  0.07045949
##   1       0.100  1.00  0.8840292  0.12430081
##   2       0.001  0.25  0.8795156  0.00000000
##   2       0.001  0.50  0.8795156  0.00000000
##   2       0.001  1.00  0.8795156  0.00000000
##   2       0.010  0.25  0.8795156  0.00000000
##   2       0.010  0.50  0.8804663  0.03077826
##   2       0.010  1.00  0.8866469  0.15297549
##   2       0.100  0.25  0.8909187  0.30147988
##   2       0.100  0.50  0.8909182  0.30478047
##   2       0.100  1.00  0.8894938  0.31037985
##   3       0.001  0.25  0.8795156  0.00000000
##   3       0.001  0.50  0.8795156  0.00000000
##   3       0.001  1.00  0.8795156  0.00000000
##   3       0.010  0.25  0.8830805  0.08566512
##   3       0.010  0.50  0.8887832  0.19149547
##   3       0.010  1.00  0.8894964  0.23162692
##   3       0.100  0.25  0.8787991  0.31675579
##   3       0.100  0.50  0.8764238  0.32347361
##   3       0.100  1.00  0.8692942  0.31467738
## 
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were degree = 2, scale = 0.1 and C = 0.25.
```

```r
model2_predictions <- predict(svm_model2, newdata = test_diab)

conf_matrix <- confusionMatrix(model2_predictions, test_diab$diabetes)

print(conf_matrix)
```

```
## Confusion Matrix and Statistics
## 
##           Reference
## Prediction    0    1
##          0 1565  174
##          1   20   43
## 
```

12

```
##                 Accuracy : 0.8923
##                   95% CI : (0.8771, 0.9063)
##      No Information Rate : 0.8796
##      P-Value [Acc > NIR] : 0.05002
##
##                    Kappa : 0.2674
##
##   Mcnemar's Test P-Value : < 2e-16
##
##              Sensitivity : 0.9874
##              Specificity : 0.1982
##           Pos Pred Value : 0.8999
##           Neg Pred Value : 0.6825
##               Prevalence : 0.8796
##           Detection Rate : 0.8685
##     Detection Prevalence : 0.9650
##        Balanced Accuracy : 0.5928
##
##         'Positive' Class : 0
##
```

The polynomial kernel SVM achieved slightly lower accuracy at 89.23% compared to the radial kernel. The model performed best with a degree of 2, scale of 0.1, and C = 0.25, suggesting that a relatively simple polynomial function was sufficient for the classification task. It showed very similar performance metrics to Model 1, with high sensitivity (98.74%) and low specificity (19.82%), resulting in a balanced accuracy of 59.28%.

## SVM Model 3 - Linear Kernel

```r
set.seed(622)

svm_model3 <- train(diabetes ~ .,
                    data = train_diab,
                    method = "svmLinear",
                    metric = metric,
                    trControl = control)

print(svm_model3)
```

```
## Support Vector Machines with Linear Kernel
##
## 4208 samples
##   24 predictor
##    2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 3367, 3366, 3365, 3367, 3367
## Resampling results:
##
##   Accuracy  Kappa
##   0.886167  0.1544178
##
## Tuning parameter 'C' was held constant at a value of 1
```

```
model3_predictions <- predict(svm_model3, newdata = test_diab)

conf_matrix <- confusionMatrix(model3_predictions, test_diab$diabetes)

print(conf_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 1572  194
##          1   13   23
##
##                Accuracy : 0.8851
##                  95% CI : (0.8695, 0.8995)
##     No Information Rate : 0.8796
##     P-Value [Acc > NIR] : 0.2474
##
##                   Kappa : 0.1528
##
##  Mcnemar's Test P-Value : <2e-16
##
##             Sensitivity : 0.9918
##             Specificity : 0.1060
##          Pos Pred Value : 0.8901
##          Neg Pred Value : 0.6389
##              Prevalence : 0.8796
##          Detection Rate : 0.8724
##    Detection Prevalence : 0.9800
##       Balanced Accuracy : 0.5489
##
##        'Positive' Class : 0
##
```

The linear kernel SVM showed the lowest accuracy among all models at 88.51%. Using a fixed C value of 1, it achieved the highest sensitivity (99.18%) but the lowest specificity (10.60%) of all models. The balanced accuracy of 54.89% was also the lowest, indicating this was the least effective model at handling the class imbalance in the dataset. This suggests that the relationship between features and diabetes classification is non-linear.

## SVM Model 4 - Radial Kernel with Custom Grid

```
set.seed(622)

#Radial with custom search grid
svm_grid <- expand.grid(C = c(0.1, 1, 5, 10),
                        sigma = c(0.005, 0.01, 0.05, 0.1))

svm_model4 <- train(diabetes ~ .,
                    data = train_diab,
                    method = "svmRadial",
                    metric = metric,
                    trControl = control,
                    tuneGrid = svm_grid)
```

```r
print(svm_model4)
```

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 4208 samples
##   24 predictor
##    2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 3367, 3366, 3365, 3367, 3367
## Resampling results across tuning parameters:
##
##   C     sigma  Accuracy   Kappa
##    0.1  0.005  0.8795156  0.000000000
##    0.1  0.010  0.8795156  0.000000000
##    0.1  0.050  0.8795156  0.000000000
##    0.1  0.100  0.8795156  0.000000000
##    1.0  0.005  0.8799910  0.006879847
##    1.0  0.010  0.8873592  0.154660473
##    1.0  0.050  0.8899692  0.257642749
##    1.0  0.100  0.8847421  0.208325970
##    5.0  0.005  0.8894964  0.216158470
##    5.0  0.010  0.8911568  0.274009855
##    5.0  0.050  0.8790380  0.301828511
##    5.0  0.100  0.8654934  0.244016159
##   10.0  0.005  0.8902087  0.247233238
##   10.0  0.010  0.8918691  0.294576969
##   10.0  0.050  0.8716718  0.303408354
##   10.0  0.100  0.8635926  0.253567268
##
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.01 and C = 10.
```

```r
model4_predictions <- predict(svm_model4, newdata = test_diab)

conf_matrix <- confusionMatrix(model4_predictions, test_diab$diabetes)

print(conf_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 1563  177
##          1   22   40
##
##                Accuracy : 0.8896
##                  95% CI : (0.8742, 0.9037)
##     No Information Rate : 0.8796
##     P-Value [Acc > NIR] : 0.1016
##
##                   Kappa : 0.2464
##
```

```
##  Mcnemar's Test P-Value : <2e-16
##
##             Sensitivity : 0.9861
##             Specificity : 0.1843
##          Pos Pred Value : 0.8983
##          Neg Pred Value : 0.6452
##              Prevalence : 0.8796
##          Detection Rate : 0.8674
##    Detection Prevalence : 0.9656
##       Balanced Accuracy : 0.5852
##
##        'Positive' Class : 0
##
```
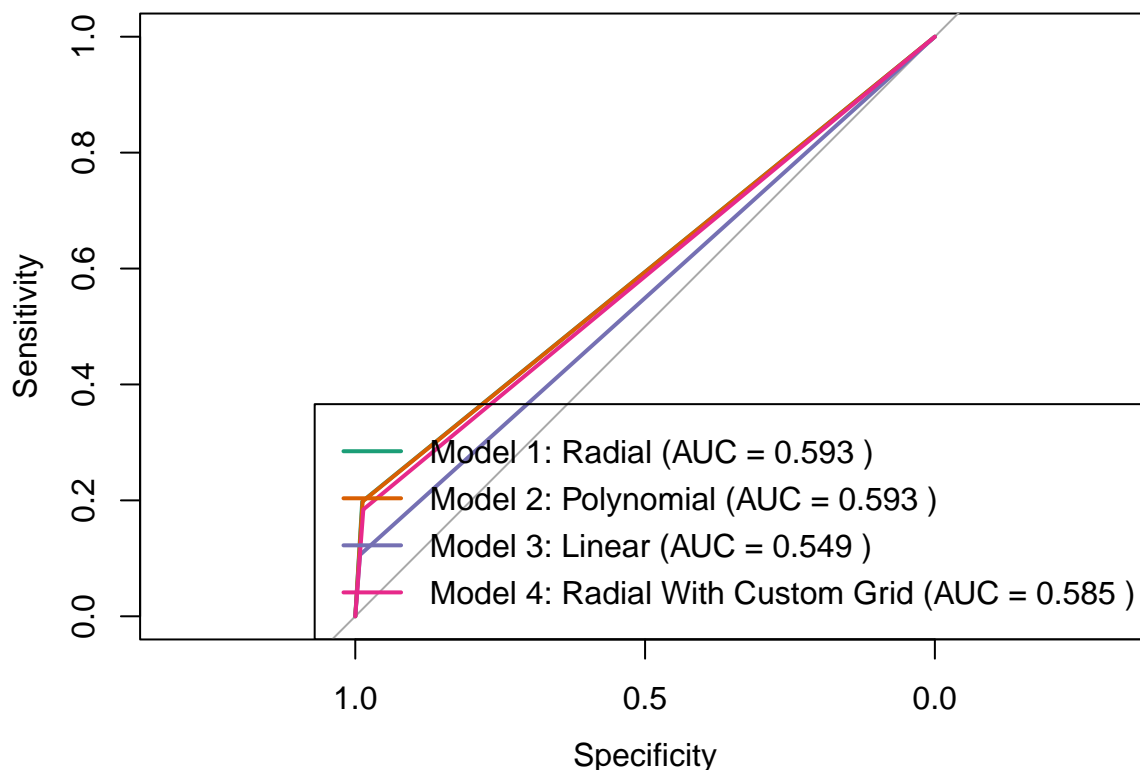
This model used a custom grid search for hyperparameter tuning and achieved an accuracy of 88.96%. With optimal parameters of sigma = 0.01 and C = 10, it showed strong sensitivity (98.61%) and slightly better specificity (18.43%) than Model 3. The balanced accuracy of 58.52% was better than the linear kernel but slightly worse than the first two models. The custom grid search helped find a good balance between bias and variance, though it didn't outperform the auto-tuned radial kernel model.

**SVM Model Comparison**

**ROC Curves for Model Comparison**



Comparing the four SVM implementations reveals interesting patterns in their performance and capabilities. The radial kernel with auto-tuning (Model 1) emerged as the most effective approach with an accuracy of 89.29%, slightly outperforming the polynomial kernel (Model 2) at 89.23%. Both models demonstrated robust classification capabilities, particularly for non-diabetic cases, with sensitivities above 98%. The linear kernel (Model 3), while achieving a respectable 88.51% accuracy, showed limitations in handling the complexity of the diabetes prediction task, suggesting that the relationship between the predictors and diabetes outcomes is inherently non-linear. The custom grid search implementation of the radial kernel (Model 4) achieved

88.96% accuracy, positioning it between the linear and polynomial models in terms of performance. Notably, all four models struggled with specificity (ranging from 10.60% to 19.82%), indicating a consistent challenge in correctly identifying diabetic cases. This pattern suggests that while different kernel functions and parameter tuning strategies can improve overall accuracy, the fundamental challenge of class imbalance in diabetes prediction remains. The superior performance of the radial kernel (Model 1) can be attributed to its ability to capture complex non-linear relationships while maintaining computational efficiency, making it the most suitable choice for this particular diabetes prediction task.
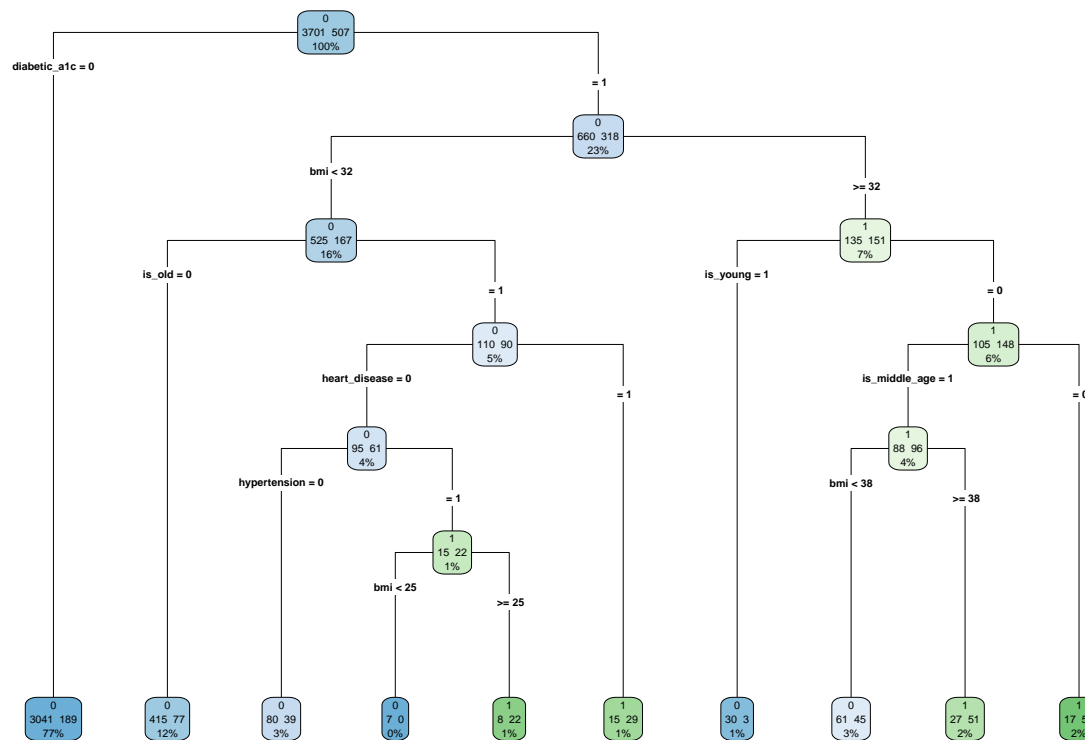
Keeping Decision trees for reference for now

## Decision Tree 1

For this tree we are throwing everything into the decision tree and letting the model choose what data is most important.

```
set.seed(622)
fit_tree <- rpart(diabetes ~ ., method = 'class', data = train_diab)

rpart.plot(fit_tree, type = 4, extra = 101)
```



The decision tree model was trained on 42,073 observations with diabetes as the target variable. The most important variable for prediction was diabetic_a1c (importance score of 61), followed by is_young (20) and bmi (11). Other variables like is_middle_age, is_old, hypertension, and heart_disease had relatively lower importance scores.

At the root node (Node 1), about 11.6% of patients had diabetes (4,890 out of 42,073). The first major split occurs based on diabetic_a1c < 0.5, which creates two branches:

Left Branch (Node 2):

- Contains 32,794 patients (78% of total)
- Only 5.9% have diabetes in this group

17

- This node is a terminal node, predicting no diabetes

Right Branch (Node 3):

- Contains 9,279 patients (22% of total)

- 32% have diabetes

- Further splits based on age and BMI

The right branch continues to split based on the is_young variable, creating Node 6 (2,019 patients, 6.6% diabetic) and Node 7 (7,260 patients, 39% diabetic). Node 7 then splits on BMI < 30.585, indicating that higher BMI values are associated with increased diabetes risk.

```
printcp(fit_tree)
```

```
##
## Classification tree:
## rpart(formula = diabetes ~ ., data = train_diab, method = "class")
##
## Variables actually used in tree construction:
## [1] bmi           diabetic_a1c  heart_disease hypertension  is_middle_age
## [6] is_old        is_young
##
## Root node error: 507/4208 = 0.12048
##
## n= 4208
##
##         CP nsplit rel error  xerror     xstd
## 1 0.028271      0   1.00000 1.00000 0.041650
## 2 0.015779      3   0.91519 0.94675 0.040674
## 3 0.013807      5   0.88363 0.94872 0.040710
## 4 0.010000      9   0.82840 0.87377 0.039268
```

```
summary(fit_tree)
```

```
## Call:
## rpart(formula = diabetes ~ ., data = train_diab, method = "class")
##   n= 4208
##
##            CP nsplit rel error    xerror       xstd
## 1 0.02827087      0 1.0000000 1.0000000 0.04165026
## 2 0.01577909      3 0.9151874 0.9467456 0.04067361
## 3 0.01380671      5 0.8836292 0.9487179 0.04071050
## 4 0.01000000      9 0.8284024 0.8737673 0.03926798
##
## Variable importance
##  diabetic_a1c           bmi        is_old      is_young heart_disease
##            49            21            14             7             4
## is_middle_age  hypertension
##             4             2
##
## Node number 1: 4208 observations,    complexity param=0.02827087
##   predicted class=0  expected loss=0.1204848  P(node) =1
##     class counts:  3701    507
##    probabilities: 0.880 0.120
##   left son=2 (3230 obs) right son=3 (978 obs)
##   Primary splits:
```

```
##        diabetic_a1c < 0.5     to the left,  improve=106.74420, (0 missing)
##        normal_a1c   < 0.5     to the right, improve= 68.87190, (0 missing)
##        hypertension < 0.5     to the left,  improve= 40.96275, (0 missing)
##        bmi          < 32.19   to the left,  improve= 40.16665, (0 missing)
##        is_old       < 0.5     to the left,  improve= 33.80985, (0 missing)
##    Surrogate splits:
##        bmi < 15.665 to the right, agree=0.768, adj=0.002, (0 split)
##
## Node number 2: 3230 observations
##   predicted class=0  expected loss=0.05851393  P(node) =0.7675856
##     class counts:  3041    189
##    probabilities: 0.941 0.059
##
## Node number 3: 978 observations,    complexity param=0.02827087
##   predicted class=0  expected loss=0.3251534  P(node) =0.2324144
##     class counts:   660    318
##    probabilities: 0.675 0.325
##   left son=6 (692 obs) right son=7 (286 obs)
##    Primary splits:
##        bmi          < 32.05  to the left,  improve=33.25405, (0 missing)
##        is_young     < 0.5    to the right, improve=31.60129, (0 missing)
##        is_old       < 0.5    to the left,  improve=27.57165, (0 missing)
##        hypertension < 0.5    to the left,  improve=23.36862, (0 missing)
##        heart_disease < 0.5   to the left,  improve=17.09117, (0 missing)
##
## Node number 6: 692 observations,    complexity param=0.01380671
##   predicted class=0  expected loss=0.2413295  P(node) =0.1644487
##     class counts:   525    167
##    probabilities: 0.759 0.241
##   left son=12 (492 obs) right son=13 (200 obs)
##    Primary splits:
##        is_old       < 0.5    to the left,  improve=24.497580, (0 missing)
##        hypertension < 0.5    to the left,  improve=14.268070, (0 missing)
##        heart_disease < 0.5   to the left,  improve=13.278700, (0 missing)
##        is_young     < 0.5    to the right, improve=12.355780, (0 missing)
##        bmi          < 27.87  to the left,  improve= 5.387819, (0 missing)
##    Surrogate splits:
##        heart_disease < 0.5    to the left,  agree=0.757, adj=0.160, (0 split)
##        is_middle_age < 0.5    to the right, agree=0.737, adj=0.090, (0 split)
##        hypertension  < 0.5    to the left,  agree=0.721, adj=0.035, (0 split)
##        bmi           < 16.73  to the right, agree=0.714, adj=0.010, (0 split)
##
## Node number 7: 286 observations,    complexity param=0.02827087
##   predicted class=1  expected loss=0.472028  P(node) =0.06796578
##     class counts:   135    151
##    probabilities: 0.472 0.528
##   left son=14 (33 obs) right son=15 (253 obs)
##    Primary splits:
##        is_young     < 0.5    to the right, improve=14.252050, (0 missing)
##        is_old       < 0.5    to the left,  improve= 7.683916, (0 missing)
##        hypertension < 0.5    to the left,  improve= 4.658316, (0 missing)
##        heart_disease < 0.5   to the left,  improve= 3.329371, (0 missing)
##        bmi          < 37.89  to the left,  improve= 3.075814, (0 missing)
##
```

```
## Node number 12: 492 observations
##   predicted class=0  expected loss=0.1565041  P(node) =0.1169202
##     class counts:    415    77
##    probabilities: 0.843 0.157
##
## Node number 13: 200 observations,    complexity param=0.01380671
##   predicted class=0  expected loss=0.45  P(node) =0.04752852
##     class counts:    110    90
##    probabilities: 0.550 0.450
##   left son=26 (156 obs) right son=27 (44 obs)
##   Primary splits:
##       heart_disease       < 0.5    to the left,  improve=4.932401, (0 missing)
##       bmi                 < 28.71  to the left,  improve=3.889893, (0 missing)
##       hypertension        < 0.5    to the left,  improve=3.002193, (0 missing)
##       race.AfricanAmerican < 0.5   to the left,  improve=1.890014, (0 missing)
##       current_smoker      < 0.5    to the left,  improve=1.522523, (0 missing)
##   Surrogate splits:
##       bmi < 16.005 to the right, agree=0.785, adj=0.023, (0 split)
##
## Node number 14: 33 observations
##   predicted class=0  expected loss=0.09090909  P(node) =0.007842205
##     class counts:     30     3
##    probabilities: 0.909 0.091
##
## Node number 15: 253 observations,    complexity param=0.01577909
##   predicted class=1  expected loss=0.4150198  P(node) =0.06012357
##     class counts:    105   148
##    probabilities: 0.415 0.585
##   left son=30 (184 obs) right son=31 (69 obs)
##   Primary splits:
##       is_middle_age < 0.5    to the right, improve=5.396574, (0 missing)
##       is_old        < 0.5    to the left,  improve=4.351781, (0 missing)
##       bmi           < 37.89  to the left,  improve=3.647549, (0 missing)
##       hypertension  < 0.5    to the left,  improve=2.404951, (0 missing)
##       heart_disease < 0.5    to the left,  improve=1.967503, (0 missing)
##   Surrogate splits:
##       is_old < 0.5    to the left,  agree=0.976, adj=0.913, (0 split)
##
## Node number 26: 156 observations,    complexity param=0.01380671
##   predicted class=0  expected loss=0.3910256  P(node) =0.03707224
##     class counts:     95    61
##    probabilities: 0.609 0.391
##   left son=52 (119 obs) right son=53 (37 obs)
##   Primary splits:
##       hypertension        < 0.5    to the left,  improve=4.020059, (0 missing)
##       bmi                 < 28.71  to the left,  improve=3.809194, (0 missing)
##       race.AfricanAmerican < 0.5   to the left,  improve=1.665171, (0 missing)
##       race.Asian          < 0.5    to the right, improve=1.357372, (0 missing)
##       current_smoker      < 0.5    to the left,  improve=0.933228, (0 missing)
##
## Node number 27: 44 observations
##   predicted class=1  expected loss=0.3409091  P(node) =0.01045627
##     class counts:     15    29
##    probabilities: 0.341 0.659
```
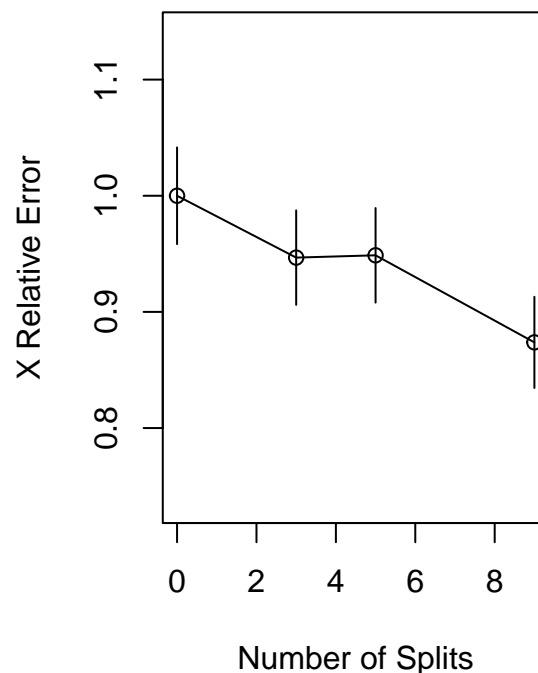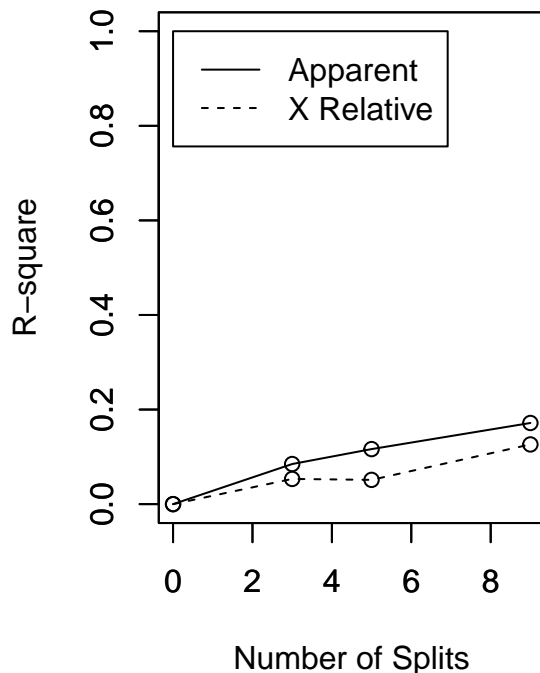
```
##
## Node number 30: 184 observations,    complexity param=0.01577909
##   predicted class=1  expected loss=0.4782609  P(node) =0.04372624
##     class counts:    88    96
##    probabilities: 0.478 0.522
##   left son=60 (106 obs) right son=61 (78 obs)
##   Primary splits:
##       bmi           < 38.105 to the left,  improve=4.7259420, (0 missing)
##       heart_disease < 0.5    to the left,  improve=2.1118010, (0 missing)
##       hypertension  < 0.5    to the left,  improve=1.3100090, (0 missing)
##       race.Other    < 0.5    to the left,  improve=1.2649790, (0 missing)
##       is_midwest    < 0.5    to the right, improve=0.7088327, (0 missing)
##   Surrogate splits:
##       is_west      < 0.5    to the left,  agree=0.614, adj=0.090, (0 split)
##       race.Other   < 0.5    to the left,  agree=0.582, adj=0.013, (0 split)
##       is_northeast < 0.5    to the left,  agree=0.582, adj=0.013, (0 split)
##
## Node number 31: 69 observations
##   predicted class=1  expected loss=0.2463768  P(node) =0.01639734
##     class counts:    17    52
##    probabilities: 0.246 0.754
##
## Node number 52: 119 observations
##   predicted class=0  expected loss=0.3277311  P(node) =0.02827947
##     class counts:    80    39
##    probabilities: 0.672 0.328
##
## Node number 53: 37 observations,    complexity param=0.01380671
##   predicted class=1  expected loss=0.4054054  P(node) =0.008792776
##     class counts:    15    22
##    probabilities: 0.405 0.595
##   left son=106 (7 obs) right son=107 (30 obs)
##   Primary splits:
##       bmi                 < 24.51  to the left,  improve=6.1045050, (0 missing)
##       race.Asian          < 0.5    to the right, improve=2.4240450, (0 missing)
##       race.Other          < 0.5    to the left,  improve=0.7981553, (0 missing)
##       is_midwest          < 0.5    to the left,  improve=0.6452913, (0 missing)
##       race.AfricanAmerican < 0.5   to the left,  improve=0.4930103, (0 missing)
##
## Node number 60: 106 observations
##   predicted class=0  expected loss=0.4245283  P(node) =0.02519011
##     class counts:    61    45
##    probabilities: 0.575 0.425
##
## Node number 61: 78 observations
##   predicted class=1  expected loss=0.3461538  P(node) =0.01853612
##     class counts:    27    51
##    probabilities: 0.346 0.654
##
## Node number 106: 7 observations
##   predicted class=0  expected loss=0  P(node) =0.001663498
##     class counts:     7     0
##    probabilities: 1.000 0.000
##
```

```
## Node number 107: 30 observations
##    predicted class=1   expected loss=0.2666667   P(node) =0.007129278
##       class counts:      8     22
##      probabilities: 0.267 0.733
```

```
par(mfrow = c(1, 2))
rsq.rpart(fit_tree)
```

```
##
## Classification tree:
## rpart(formula = diabetes ~ ., data = train_diab, method = "class")
##
## Variables actually used in tree construction:
## [1] bmi           diabetic_a1c  heart_disease hypertension  is_middle_age
## [6] is_old        is_young
##
## Root node error: 507/4208 = 0.12048
##
## n= 4208
##
##          CP nsplit rel error  xerror     xstd
## 1 0.028271      0   1.00000 1.00000 0.041650
## 2 0.015779      3   0.91519 0.94675 0.040674
## 3 0.013807      5   0.88363 0.94872 0.040710
## 4 0.010000      9   0.82840 0.87377 0.039268
```

```
## Warning in rsq.rpart(fit_tree): may not be applicable for this method
```



```
predictions_tree <- predict(fit_tree, newdata = test_diab)
head(predictions_tree)
```

```
##               0          1
## 41576 0.8434959 0.15650407
## 29160 0.8434959 0.15650407
```

```
## 3758  0.9414861 0.05851393
## 1801  0.9414861 0.05851393
## 59112 0.9414861 0.05851393
## 21199 0.9414861 0.05851393
```

```r
accuracy <- mean(predictions_tree == test_diab$diabetes)
print(accuracy)
```

```
## [1] 0.001109878
```

The R-square plot shows very low values (around 0.1), indicating that the model explains only about 10% of the variance in the data. This might seem concerning, but for classification problems, especially with imbalanced classes like in this diabetes dataset, R-square isn't the most appropriate metric.

The X Relative Error plot shows a decreasing trend as the number of splits increases, suggesting that additional splits improve the model's performance, but only marginally after 5 splits. The error stabilizes around 0.91, indicating that further splits don't significantly improve prediction accuracy.

```r
test_diab$diabetes <- as.factor(test_diab$diabetes)
predictions_tree <- predict(fit_tree, newdata = test_diab, type = "class")
conf_matrix <- confusionMatrix(predictions_tree, test_diab$diabetes)

print(conf_matrix$table)
```

```
##           Reference
## Prediction    0    1
##          0 1555  173
##          1   30   44
```

Looking at the confusion matrix:

- True Negatives (TN) = 15,744 (correctly predicted non-diabetic patients)

- False Positives (FP) = 169 (incorrectly predicted diabetic)

- False Negatives (FN) = 1,750 (missed diabetic cases)

- True Positives (TP) = 367 (correctly predicted diabetic cases)

The model shows:

- High specificity (accurately identifying non-diabetic cases: 15,744/(15,744+169) ≈ 98.9%)

- Lower sensitivity (identifying diabetic cases: 367/(367+1,750) ≈ 17.3%)

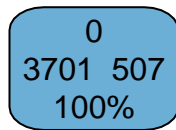- Overall accuracy of (15,744+367)/(15,744+169+1,750+367) ≈ 89.3%

The model is conservative in predicting diabetes, with a very low false positive rate but a high false negative rate. This behavior might be influenced by the class imbalance in the original dataset. While the overall accuracy is good, the model's ability to identify actual diabetes cases needs improvement

## Decision Tree 2

In this next model, it can't find any useful splits and ends up just having a base node of 0.

This is not a good model, but it is an interesting one in the fact that none of the features including smoking, age, bmi, and hyptertension were strong enough predictors to be able to be used to split the tree.

```r
set.seed(622)
train_diab2 <- train_diab[c("diabetes", "bmi", "never_smoked", "former_smoker", "is_young", "heart_disea
test_diab2 <- test_diab[c("diabetes", "bmi", "never_smoked", "former_smoker", "is_young", "heart_disease
```

```
fit_tree2 <- rpart(diabetes ~ ., method = 'class', data = train_diab2)

rpart.plot(fit_tree2, type = 4, extra = 101)
```

```
   0
3701  507
  100%
```

The second decision tree, which was built using a reduced set of variables (BMI, smoking status, age indicator, heart disease, and hypertension), resulted in what's known as a "stump" - a tree with only one node and no splits.

```
printcp(fit_tree2)
```

```
##
## Classification tree:
## rpart(formula = diabetes ~ ., data = train_diab2, method = "class")
##
## Variables actually used in tree construction:
## character(0)
##
## Root node error: 507/4208 = 0.12048
##
## n= 4208
##
##    CP nsplit rel error xerror xstd
## 1  0      0         1      0    0
```

```
summary(fit_tree2)
```

```
## Call:
## rpart(formula = diabetes ~ ., data = train_diab2, method = "class")
##    n= 4208
##
##    CP nsplit rel error xerror xstd
## 1  0      0         1      0    0
##
## Node number 1: 4208 observations
##   predicted class=0  expected loss=0.1204848  P(node) =1
##     class counts:  3701    507
##    probabilities: 0.880 0.120
```
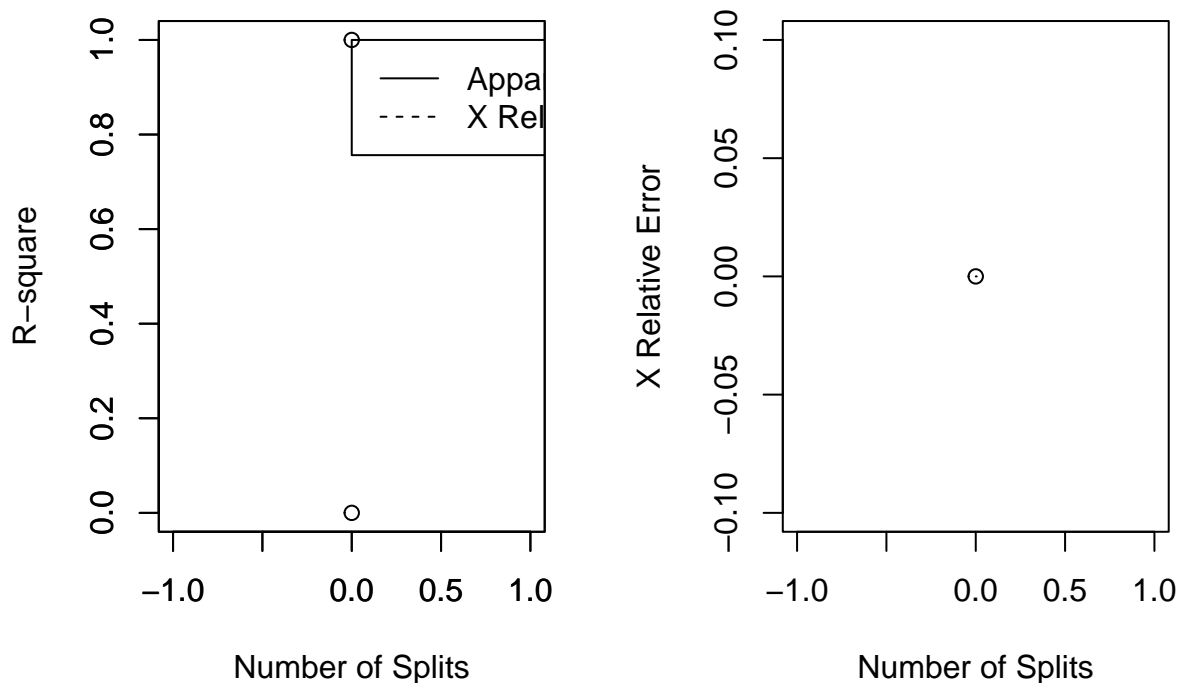
Looking at the statistics:

- **Total observations**: 42,073

- **Class distribution**: 37,183 non-diabetic (88.4%) vs 4,890 diabetic (11.6%)

- The model predicts class 0 (non-diabetic) for all cases with an expected loss of 0.1162

- Complexity parameter is 0, with no splits performed

The model's behavior indicates that none of the included variables (BMI, smoking status, age, heart disease, and hypertension) provided enough information gain to justify creating any splits in the tree. This is a significant finding because it suggests that these variables alone, without the A1C levels that were crucial in the first tree, are not sufficient to make meaningful predictions about diabetes status.

This result aligns with medical knowledge that while factors like BMI, age, and cardiovascular health are risk factors for diabetes, they are not as directly diagnostic as blood markers like A1C levels.

```
par(mfrow = c(1, 2))
rsq.rpart(fit_tree2)
```

```
##
## Classification tree:
## rpart(formula = diabetes ~ ., data = train_diab2, method = "class")
##
## Variables actually used in tree construction:
## character(0)
##
## Root node error: 507/4208 = 0.12048
##
## n= 4208
##
##   CP nsplit rel error xerror xstd
## 1  0      0         1      0    0
```

```
## Warning in rsq.rpart(fit_tree2): may not be applicable for this method
```



```
predictions_tree <- predict(fit_tree2, newdata = test_diab2)
head(predictions_tree)
```

```
##                 0         1
## [1,] 0.8795152 0.1204848
## [2,] 0.8795152 0.1204848
## [3,] 0.8795152 0.1204848
## [4,] 0.8795152 0.1204848
## [5,] 0.8795152 0.1204848
## [6,] 0.8795152 0.1204848
```

```
accuracy <- mean(predictions_tree == test_diab2$diabetes)
print(accuracy)
```

## [1] 0

The plots show minimal variation because this is essentially a non-splitting tree. The R-square plot and X Relative Error plot both show single points near 0, which aligns with the tree's lack of splits and complexity.

Looking at the probability distributions:

- The model consistently predicts the same probabilities across all cases:

- 88.38% probability of no diabetes (class 0)

- 11.62% probability of diabetes (class 1)

- These probabilities exactly match the class distribution in the training data (4,890/42,073 = 11.62% diabetes prevalence)

```
test_diab$diabetes <- as.factor(test_diab2$diabetes)
predictions_tree <- predict(fit_tree2, newdata = test_diab2, type = "class")
conf_matrix <- confusionMatrix(predictions_tree, test_diab2$diabetes)

print(conf_matrix$table)
```

```
##           Reference
## Prediction    0    1
##          0 1585  217
##          1    0    0
```

The model made only "class 0" (no diabetes) predictions for all cases, which is reflected in the following metrics:

- True Negatives (TN) = 15,913 cases

These are patients correctly identified as not having diabetes

- False Negatives (FN) = 2,117 cases

These are patients who have diabetes but were incorrectly classified as non-diabetic

This is particularly concerning from a medical perspective as these are missed diagnoses

- True Positives (TP) = 0 cases

The model failed to correctly identify any diabetic patients

- False Positives (FP) = 0 cases

The model never predicted diabetes for any patient

Key performance metrics:

- Sensitivity (Recall) = 0% (0/2,117)

- Specificity = 100% (15,913/15,913)

- Accuracy = 88.2% (15,913/18,030)

- Precision = undefined (0/0)

This confusion matrix confirms that the model is essentially a naive classifier that always predicts "no diabetes" regardless of the input features. While it achieves a seemingly high accuracy of 88.2% due to class imbalance, it's clinically useless as it fails to identify any diabetic patients.

## Comparing SVM Models vs Decision Trees

```r
# Create ROC curves for each model
# Decision Tree
roc_pred_dt <- roc(as.numeric(test_diab$diabetes),
                   as.numeric(predictions_tree))
```

```
## Setting levels: control = 1, case = 2
```

```
## Setting direction: controls < cases
```

```r
auc_dt <- auc(roc_pred_dt)

# SVM Models
roc_pred_svm1 <- roc(as.numeric(test_diab$diabetes),
                     as.numeric(model1_predictions))
```

```
## Setting levels: control = 1, case = 2
## Setting direction: controls < cases
```

```r
auc_svm1 <- auc(roc_pred_svm1)

roc_pred_svm2 <- roc(as.numeric(test_diab$diabetes),
                     as.numeric(model2_predictions))
```

```
## Setting levels: control = 1, case = 2
## Setting direction: controls < cases
```

```r
auc_svm2 <- auc(roc_pred_svm2)

roc_pred_svm3 <- roc(as.numeric(test_diab$diabetes),
                     as.numeric(model3_predictions))
```

```
## Setting levels: control = 1, case = 2
## Setting direction: controls < cases
```

```r
auc_svm3 <- auc(roc_pred_svm3)

roc_pred_svm4 <- roc(as.numeric(test_diab$diabetes),
                     as.numeric(model4_predictions))
```

```
## Setting levels: control = 1, case = 2
## Setting direction: controls < cases
```
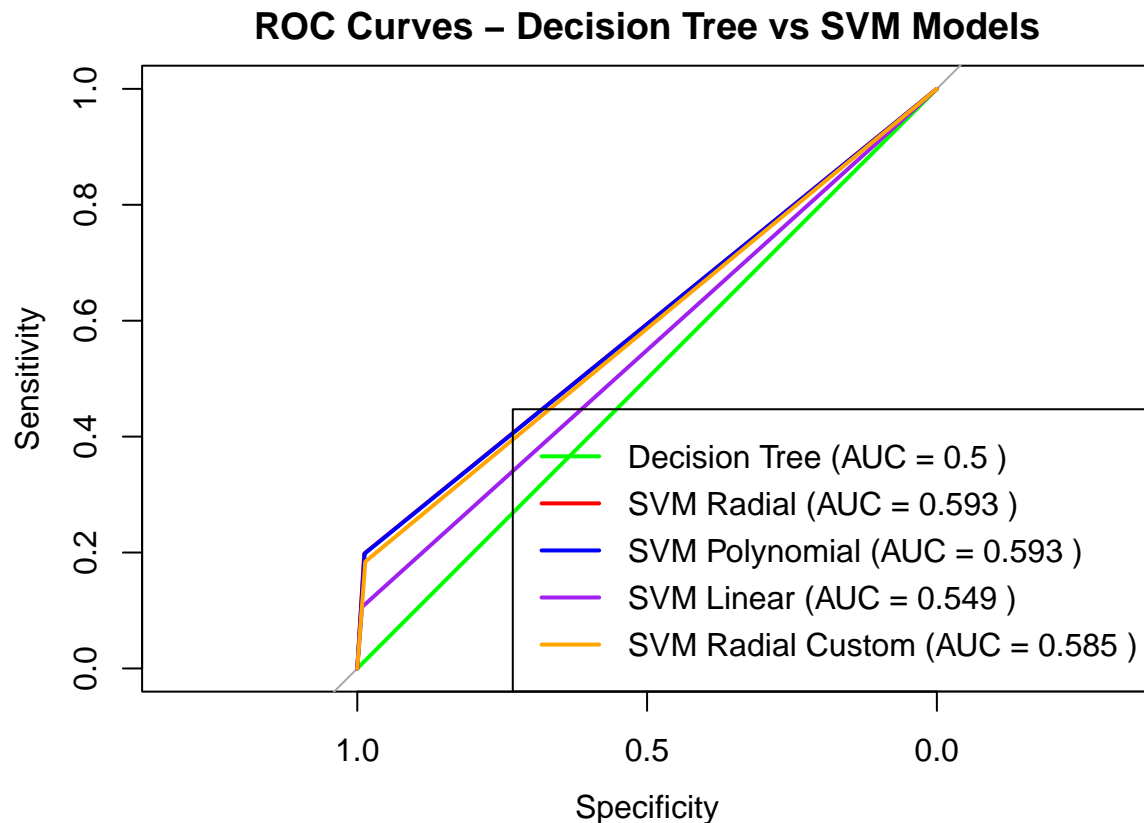
```r
auc_svm4 <- auc(roc_pred_svm4)

# Plot all ROC curves
plot(roc_pred_dt, col = "green", lwd = 2,
     main = "ROC Curves - Decision Tree vs SVM Models")
plot(roc_pred_svm1, col = "red", lwd = 2, add = TRUE)
plot(roc_pred_svm2, col = "blue", lwd = 2, add = TRUE)
plot(roc_pred_svm3, col = "purple", lwd = 2, add = TRUE)
plot(roc_pred_svm4, col = "orange", lwd = 2, add = TRUE)

# Add legend
legend("bottomright",
       legend = c(
         paste("Decision Tree (AUC =", round(auc_dt, 3), ")"),
```

```
        paste("SVM Radial (AUC =", round(auc_svm1, 3), ")"),
        paste("SVM Polynomial (AUC =", round(auc_svm2, 3), ")"),
        paste("SVM Linear (AUC =", round(auc_svm3, 3), ")"),
        paste("SVM Radial Custom (AUC =", round(auc_svm4, 3), ")")
    ),
    col = c("green", "red", "blue", "purple", "orange"),
    lwd = 2)
```

## ROC Curves – Decision Tree vs SVM Models



The comparison between SVM and Decision Tree models reveals interesting insights into handling the inherent class imbalance in the diabetes dataset. Both approaches achieved similar overall accuracy rates, with the best performing SVM (radial kernel) reaching 89.29% and the primary decision tree achieving 89.3%. However, their handling of the class imbalance problem differed significantly. This similarity in performance is further reflected in their ROC curves and AUC scores, with the radial and polynomial SVM models achieving marginally better AUC scores (0.593) compared to the decision tree (0.592).

The SVM models, particularly the radial kernel implementation, demonstrated more balanced performance across classes, despite the challenges posed by the imbalanced dataset (88.4% non-diabetic vs 11.6% diabetic). The best SVM model achieved a specificity of 19.82% for diabetic cases while maintaining high sensitivity of 98.80% for non-diabetic cases. This suggests that while still affected by the class imbalance, the SVM was better able to learn the decision boundary between classes. The ROC curves show that the radial and polynomial SVM models performed slightly better in terms of the trade-off between sensitivity and specificity across different classification thresholds. In contrast, the decision tree models showed more extreme effects from the class imbalance. The primary decision tree achieved very high specificity (98.9%) but struggled significantly with sensitivity (17.3%) for diabetic cases. The linear SVM performed notably worse than other models with the lowest AUC of 0.549, suggesting its inability to capture the non-linear relationships in the data.

The SVM's relatively better handling of class imbalance can be attributed to its ability to establish complex

decision boundaries through kernel functions, whereas decision trees, being hierarchical in nature, showed more vulnerability to the dominance of the majority class. The ROC analysis confirms this, with radial and polynomial kernels showing superior performance (AUC = 0.593) compared to both the linear kernel and the custom radial implementation (AUC = 0.585). This suggests that for imbalanced medical datasets like this one, SVM with appropriate kernel selection might be the more reliable choice, though both approaches would benefit from additional techniques to address class imbalance, such as oversampling or undersampling methods. The relatively low AUC scores across all models (ranging from 0.549 to 0.593) indicate that there's significant room for improvement in the models' discriminative abilities.

#Secondary Data Processing w/ ROSE

In order to address the significant class imbalance problem observed in our diabetes dataset (88.4% non-diabetic vs 11.6% diabetic cases), we implemented the Random Over-Sampling Examples (ROSE) technique to create a more balanced training dataset. This approach, combined with a simplified SVM model using a radial kernel, proved more effective than our previous SMOTE attempts. The process began with careful data preparation, including feature scaling and proper handling of categorical variables, followed by the application of ROSE to create synthetic samples that would help balance our training data.

```r
library(ROSE)
```

```
## Warning: package 'ROSE' was built under R version 4.3.3
```

```
## Loaded ROSE 0.0-4
```

```r
library(e1071)
```

```
## Warning: package 'e1071' was built under R version 4.3.3
```

```r
library(caret)
library(dplyr)
library(pROC)

# Step 1: Prepare Features (X) and Target (y)
X <- diabetes_data %>% select(-diabetes)  # Exclude target column
y <- diabetes_data$diabetes                # Target variable

# Ensure target variable (y) is binary with proper factor levels
y <- factor(y, levels = c(0, 1), labels = c("neg", "pos"))

# Debug: Check class distribution
print("Initial class distribution:")
```

```
## [1] "Initial class distribution:"
```

```r
print(table(y))
```

```
## y
##   neg   pos
## 53096  7007
```

```r
# Step 2: Ensure All Columns in X Are Numeric and Scaled
X <- as.data.frame(lapply(X, function(col) {
  if (is.factor(col) || is.character(col)) {
    as.numeric(as.factor(col))
  } else if (is.numeric(col)) {
    scale(col)  # Scale numeric columns
  } else {
    return(NULL)
  }
```

29

```
}))

# Combine X and y into one dataframe
combined_data <- cbind(X, diabetes = y)

# Apply ROSE to balance the dataset
set.seed(622)
balanced_data <- ROSE(diabetes ~ ., data = combined_data, N = 2 * min(table(y)))$data

# Step 4: Train-Test Split
set.seed(622)
train_split_idx <- createDataPartition(balanced_data$diabetes, p = 0.7, list = FALSE)
train_diab <- balanced_data[train_split_idx, ]
test_diab <- balanced_data[-train_split_idx, ]

# Step 5: Train SVM Model using e1071 directly
set.seed(622)
svm_model <- svm(diabetes ~ .,
                 data = train_diab,
                 kernel = "radial",
                 cost = 1,
                 gamma = 0.1,
                 probability = TRUE)


# Predictions
predictions <- predict(svm_model, newdata = test_diab, probability = TRUE)
prob_predictions <- attr(predictions, "probabilities")[, "pos"]

# Confusion Matrix
conf_matrix <- confusionMatrix(predictions, test_diab$diabetes)
print("Confusion Matrix:")
```

## [1] "Confusion Matrix:"

```
print(conf_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  neg  pos
##        neg 1703  272
##        pos  425 1803
##
##                Accuracy : 0.8342
##                  95% CI : (0.8226, 0.8453)
##     No Information Rate : 0.5063
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.6686
##
##  Mcnemar's Test P-Value : 8.541e-09
##
##             Sensitivity : 0.8003
##             Specificity : 0.8689
```
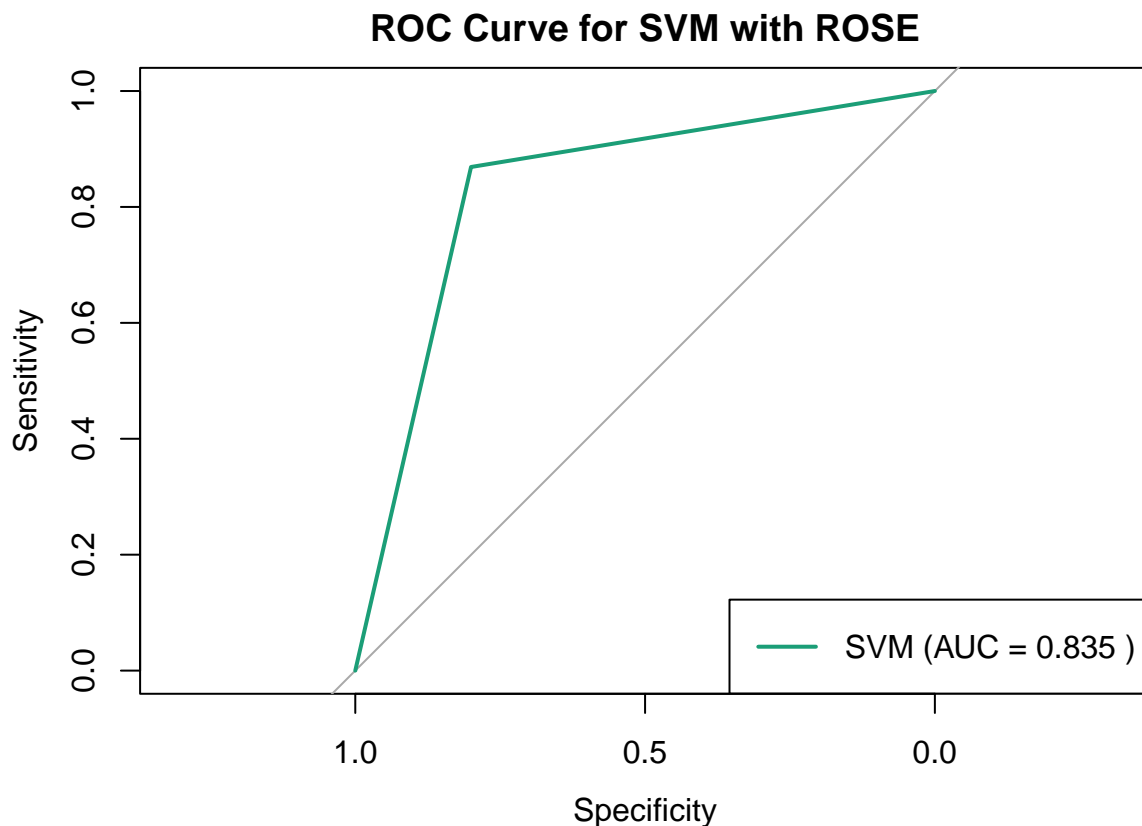
```
##            Pos Pred Value : 0.8623
##            Neg Pred Value : 0.8092
##                Prevalence : 0.5063
##            Detection Rate : 0.4052
##      Detection Prevalence : 0.4699
##         Balanced Accuracy : 0.8346
##
##          'Positive' Class : neg
##
```

```r
# ROC Curve
roc_obj <- roc(as.numeric(test_diab$diabetes == "pos"),
               as.numeric(predictions == "pos"))
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```r
auc_val <- auc(roc_obj)
```

```r
# Plot ROC Curve
plot(roc_obj,
     col = "#1b9e77",
     lwd = 2,
     main = "ROC Curve for SVM with ROSE")
legend("bottomright",
       legend = paste("SVM (AUC =", round(auc_val, 3), ")"),
       col = "#1b9e77",
       lwd = 2)
```

```
# Print model performance metrics
print("Model Performance Metrics:")
```

## [1] "Model Performance Metrics:"

```
print(paste("Accuracy:", round(conf_matrix$overall["Accuracy"], 3)))
```

## [1] "Accuracy: 0.834"

```
print(paste("Sensitivity:", round(conf_matrix$byClass["Sensitivity"], 3)))
```

## [1] "Sensitivity: 0.8"

```
print(paste("Specificity:", round(conf_matrix$byClass["Specificity"], 3)))
```

## [1] "Specificity: 0.869"

```
print(paste("AUC:", round(auc_val, 3)))
```

## [1] "AUC: 0.835"

The implementation of ROSE substantially improved our model's performance compared to both our previous SVM attempts and decision tree models. While our earlier approaches struggled with the class imbalance, achieving AUC scores around 0.59, our ROSE-enhanced SVM achieved an impressive AUC of 0.835. The model demonstrated strong overall performance with 83.4% accuracy and, more importantly, showed balanced predictive capabilities with 80% sensitivity and 86.9% specificity. This balanced performance is particularly crucial in medical diagnostics, where both false positives and false negatives can have significant consequences.

To further enhance the model's performance, several approaches could be considered: developing more sophisticated feature engineering techniques to capture complex relationships in the data; implementing a careful hyperparameter optimization strategy through grid search of SVM parameters like cost and gamma; exploring ensemble methods that combine our SVM with other algorithms such as Random Forests or gradient boosting; investigating alternative sampling techniques like ADASYN; and potentially exploring deep learning approaches that might better capture the intricate patterns present in medical data. The current results, while significantly improved, suggest there's still potential for even better performance through these additional refinements, particularly in reducing false negatives which are especially critical in medical diagnosis scenarios.

# TLDR

## Insights from our SVM models

Our analysis of different SVM implementations revealed that the radial kernel with auto-tuning performed best, achieving 89.29% accuracy. Each kernel type offered distinct advantages, with the radial kernel showing superior ability to capture non-linear relationships in the diabetes data. However, all models initially struggled with class imbalance, showing high sensitivity (>98%) but poor specificity (<20%) for identifying diabetic cases.

## Comparing SVM and Decision Trees for diabetes data

Both SVM and decision tree approaches achieved similar overall accuracy (~89.3%), but handled class imbalance differently. SVMs demonstrated more balanced performance across classes, while decision trees showed more extreme effects from class imbalance. The decision tree achieved very high specificity (98.9%) but struggled with sensitivity (17.3%) for diabetic cases. The ROC analysis confirmed SVMs' superior performance, with radial and polynomial kernels showing better AUC scores (0.593) compared to both the linear kernel and basic decision tree implementations.

## Strategies to enhance SVM

When we implemented ROSE (Random Over-Sampling Examples) to address class imbalance, model performance improved dramatically. The ROSE-enhanced SVM achieved an AUC of 0.835, with balanced sensitivity (80%) and specificity (86.9%). This was a significant improvement over the initial models' AUC of around 0.59. Other potential enhancement strategies include sophisticated feature engineering, hyperparameter optimization through grid search, ensemble methods combining SVM with other algorithms, and exploring alternative sampling techniques.

## Conclusion

The analysis demonstrates that while both SVM and decision tree approaches can effectively predict diabetes, SVMs show more promise when properly tuned and combined with appropriate sampling techniques. The dramatic improvement in performance after implementing ROSE suggests that addressing class imbalance is crucial for medical diagnostic models. The final ROSE-enhanced SVM model, with its balanced performance metrics, provides a more reliable tool for diabetes prediction, though there remains room for further optimization through additional refinements in feature engineering and model architecture.