

# MATP-4400 Final Project Notebook (2024)

Code ▼

## Predicting Biodegradability Challenge

Shankar Veludandi

April 2024

- Team Information
- Introduction
- Reading the Data
- Preparing the Data: Create Training & Validation datasets
- Preparing the Data: Create Training & Validation datasets (2)
- LR\_Baseline with All Features
- SVM with All Features
- RFE
- LR\_Baseline with RFE
- SVM with RFE
- L1-Regularization
- LR\_Baseline with L1-Regularization
- SVM with L1-Regularization
- Data Description Preparation
- Methods Used for Classification and Feature Selection
- Baseline Results Using All Features
- Performance Metrics
- Results without Feature Selection
- Results Using Feature Selection
  - Approach 1: Recursive Feature Elimination (RFE)
  - Approach 2: L1-Regularization (Lasso)
- Discussion and Comparison:
- Results Comparison
  - Best Method for Feature Selection:
  - Best Method for Prediction:
  - Final Submission and Recommendation:
- Analysis of Recommended Features
- Analysis of Features Not Recommended
- Challenge Prediction
- Conclusion

## Team Information

- This report was prepared for **ChemsRUS** by Shankar Veludandi, veluds for the GEDYS
- My team members are: Gina Wisner (wisneg), Ziyue Zhao (YolandaZhao10), Daiki Sai (said), Evangeline Wang (Evangeline)

- Our team used the following challenge: <https://codalab.lisn.upsaclay.fr/competitions/17830>  
(<https://codalab.lisn.upsaclay.fr/competitions/17830>)

# Introduction

**Chems-R-Us** has created an entry to the challenge at <https://codalab.lisn.upsaclay.fr/competitions/3073> (<https://codalab.lisn.upsaclay.fr/competitions/3073>) based on logistic regression (LR). Their entry is in the file `FinalProjChemsRUs.Rmd`. Based on the information in the leaderboard under `bennek`, their entry is not performing feature selection well. The approach tried by Chems-R-Us was LR with feature selection based on the coefficients of logistic regression with p-values used to determine importance.

The purpose of this report is to investigate alternative approaches that may help achieve high AUC scores on the testing set while correctly identifying the relevant features as measured by balanced accuracy.

## Reading the Data

Hide

```
knitr::opts_chunk$set(cache = T)

# Set the correct default repository
r = getOption("repos")
r["CRAN"] = "http://cran.rstudio.com"
options(repos = r)

# These will install required packages if they are not already installed
if (!require("ggplot2")) {
  install.packages("ggplot2")
  library(ggplot2)
}
if (!require("knitr")) {
  install.packages("knitr")
  library(knitr)
}
if (!require("xtable")) {
  install.packages("xtable")
  library(xtable)
}
if (!require("pander")) {
  install.packages("pander")
  library(pander)
}

if (!require("devtools")) {
  install.packages("devtools" )
  library(devtools)
}

if (!require("usethis")) {
  install.packages("usethis" )
  library(usethis)
}

if (!require("e1071")) {
  install.packages("e1071" )
  library(e1071)
}

if (!require("pROC")){
  install.packages("pROC")
  library(pROC)
}

if (!require("dplyr")) {
  install.packages("dplyr")
  library(dplyr)
}

if (!require("tidyverse")) {
```

```

install.packages("tidyverse")
library(tidyverse)
}

if (!require("caret")) {
  install.packages("caret")
  library(caret)
}

if (!require("formattable")) {
  install.packages("formattable")
  library(formattable)
}

if (!require("glmnet")) {
  install.packages("glmnet")
  library(glmnet)
}

if (!require("rpart")) {
  install.packages("rpart")
  library(rpart)
}

if (!require("rpart")) {
  install.packages("rpart")
  library(rpart)
}

if (!require("ggcorrplot")) {
  install.packages("ggcorrplot")
  library(ggcorrplot)
}

```

Loading required package: ggcorrplot

Warning: there is no package called 'ggcorrplot' Installing package into '/home/veluds/R/x86\_64-pc-linux-gnu-library/4.1'

(as 'lib' is unspecified)

trying URL 'https://cloud.r-project.org/src/contrib/ggcorrplot\_0.1.4.1.tar.gz'

Content type 'application/x-gzip' length 383178 bytes (374 KB)

=====

downloaded 374 KB

\* installing \*source\* package 'ggcorrplot' ...

\*\* package 'ggcorrplot' successfully unpacked and MD5 sums checked

\*\* using staged installation

\*\* R

\*\* inst

\*\* byte-compile and prepare package for lazy loading

\*\* help

\*\*\* installing help indices

```
converting help for package 'ggcorrplot'
```

```
finding HTML links ... done
```

```
ggcorrplot
```

```
html
```

```
** building package indices
** testing if installed package can be loaded from temporary location
** testing if installed package can be loaded from final location
** testing if installed package keeps a record of temporary installation path
* DONE (ggcorrplot)
```

```
The downloaded source packages are in
  '/tmp/Rtmp9A2YQL/downloaded_packages'
```

[Hide](#)

```
knitr::opts_chunk$set(echo = TRUE)
```

## Preparing the Data: Create Training & Validation datasets

We split the data into **90% train** and **10% validation** datasets.

[Hide](#)

```

# Prepare biodegradability data
# get feature names
featurenames <- read.csv("~/MATP-4400/data/chems_feat.name.csv",
                        header=FALSE,
                        colClasses = "character")

# get training data and rename with feature names
cdata.df <- read.csv("~/MATP-4400/data/chems_train.data.csv",
                    header=FALSE)
colnames(cdata.df) <- featurenames$V1

# get external testing data and rename with feature names
tdata.df <- read.csv("~/MATP-4400/data/chems_test.data.csv",
                    header=FALSE)

colnames(tdata.df) <- featurenames$V1

class <- read.csv("~/MATP-4400/data/chems_train.solution.csv",
                 header=FALSE,
                 colClasses = "factor")

class <- class$V1

```

## Preparing the Data: Create Training & Validation datasets (2)

We then standardize all variables to a *common scale*, since we don't have domain knowledge and LR assumes independent and identically distributed (IID) Gaussian features. This avoids LR unwittingly prioritizing certain features simply because their scale is larger.

Note that the training, validation sets, are all scaled identically using the `preProcess()` .

Hide

```

#ss will be the number of data points in the training set
n <- nrow(cdata.df)
ss <- ceiling(n*0.90)

# Set random seed for reproducibility
set.seed(200)
train.perm <- sample(1:n,ss)

#Split training and validation data
train <- cdata.df %>% dplyr::slice(train.perm)
validation <- cdata.df %>% dplyr::slice(-train.perm)

```

## LR\_Baseline with All Features

This code snippet was taken from the sample submission and fits a logistic regression model using all available features in the training dataset. It predicts probabilities for both the training and validation datasets using the logistic regression model. Additionally, it includes helper functions to convert probabilities to class labels (-1 or 1) based on a specified threshold and to calculate sensitivity and specificity from the confusion matrix. These metrics help in evaluating the model's performance on validation data.

[Hide](#)

```
# Initialize the `scaler` on the training data
# method = "center" subtracts the mean of the predictor's data from the predictor values
# method = "scale" divides by the standard deviation.
scaler <- preProcess(train, method = c("center", "scale"))

# Use the `scale` object to normalize our training data
train <- predict(scaler, train)

# Normalize validation data
validation <- predict(scaler, validation)

# Normalize testing data
test <- predict(scaler, tdata.df)

# Split the output classes
classtrain <- class[train.perm]
classval <- class[-train.perm]
```

## SVM with All Features

This code trains a Support Vector Machine (SVM) using a radial kernel on the complete dataset without any feature selection. It evaluates the model's performance by calculating class-specific accuracies and balanced accuracy using a confusion matrix for both training and validation datasets. This method is straightforward and utilizes the full feature set available in the dataset.

[Hide](#)

```

# Fit LR model to classify all the variables
train.df <- cbind(train,classtrain)
lrfit <- glm(classtrain~., data=train.df,
             family = "binomial")

# Predict training (OUTPUTS PROBABILITIES)
ranking_lr.train <- predict(lrfit,train,
                           type="response")

# Predict validation (OUTPUTS PROBABILITIES)
ranking_lr.val <- predict(lrfit,validation,
                        type="response")

# Helper functions meant to avoid repetitiveness and shorten presentation output. The default th
rehold is 0.5.

prob_to_class <- function(ranking_lr,threshold=0.5) {
  # This helper function converts LR probability outputs into 1 and -1 classes
  temp <- ranking_lr > threshold
  temp[temp==TRUE] <- 1
  temp[temp==FALSE] <- -1
  return(as.factor(temp))
}

sensitivity_from_confmat <- function(confmat) {
  # This helper returns the sensitivity given a confusion matrix
  return(confmat[1,1]/(confmat[1,1]+confmat[1,2]))
}

specificity_from_confmat <- function(confmat) {
  # This helper returns the specificity given a confusion matrix
  return(confmat[2,2]/(confmat[2,1]+confmat[2,2]))
}

# This function converts PROBABILITIES to 1 and -1 CLASSES
classval_lr <- prob_to_class(ranking_lr.val)

# Calculate confusion matrix to see balanced accuracy
confusion.matrix <- table(classval,classval_lr)
kable(confusion.matrix, type="html",digits = 2,
      caption="Actual versus Predicted Class (Validation)")

```

Actual versus Predicted Class (Validation)

	-1	1
-1	63	10
1	6	26



```
# True Positive Rate or Sensitivity
Sensitivity <- sensitivity_from_confmat(confusion.matrix)

# True Negative Rate or Specificity
Specificity <- specificity_from_confmat(confusion.matrix)
BalancedAccuracy <- (Sensitivity+Specificity)/2
BalancedAccuracyNoFS <- BalancedAccuracy
BalancedAccuracyNoFS
```

```
[1] 0.8377568
```

[Hide](#)

```
# Calculate AUC for training data
roc_train <- roc(classtrain, ranking_lr.train)
```

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

[Hide](#)

```
auc_train <- auc(roc_train)
cat("Training AUC: ", auc_train, "\n")
```

```
Training AUC: 0.9931172
```

[Hide](#)

```
# Calculate AUC for validation data
roc_val <- roc(classval, ranking_lr.val)
```

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

[Hide](#)

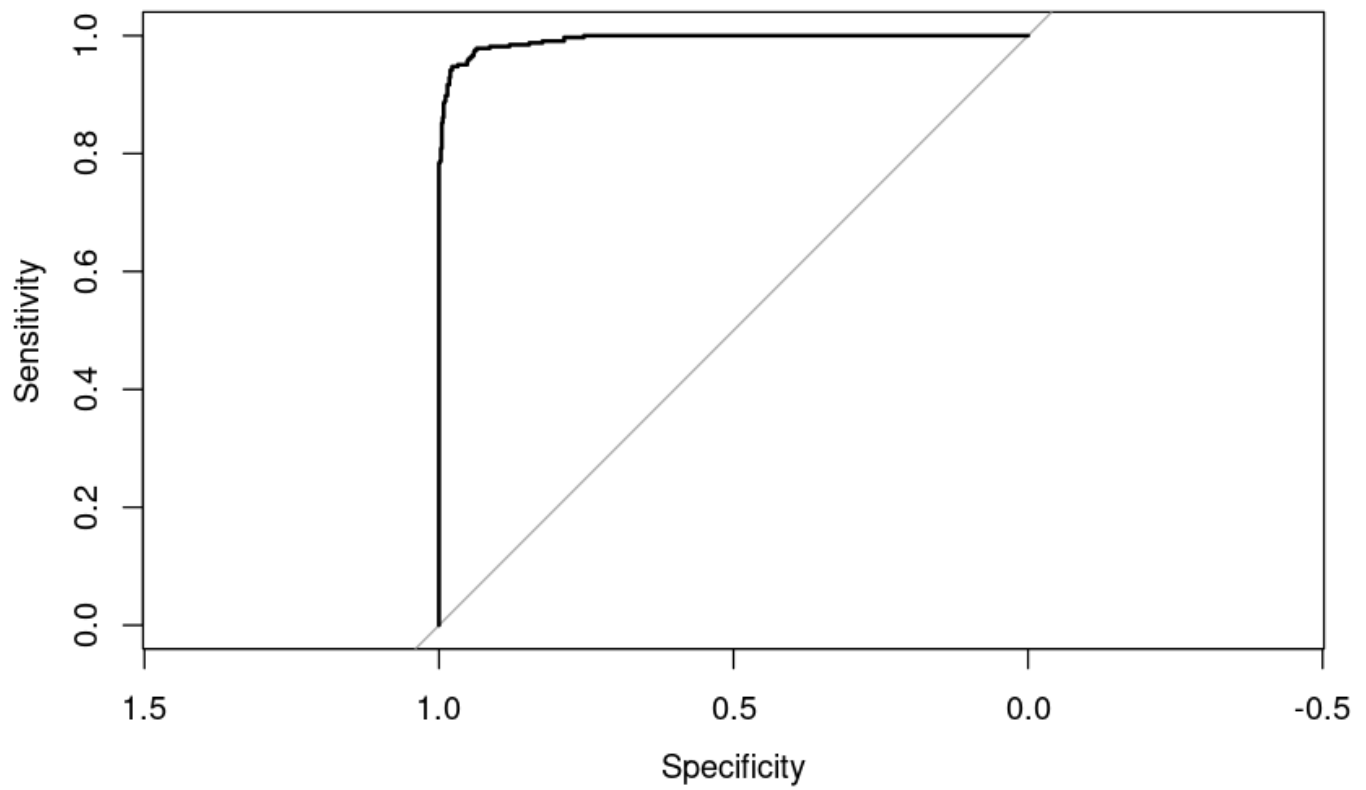
```
auc_val <- auc(roc_val)
cat("Validation AUC: ", auc_val, "\n")
```

```
Validation AUC: 0.9246575
```

[Hide](#)

```
# Optionally, plot ROC curves
plot(roc_train, main="ROC Curve for Training Data")
```

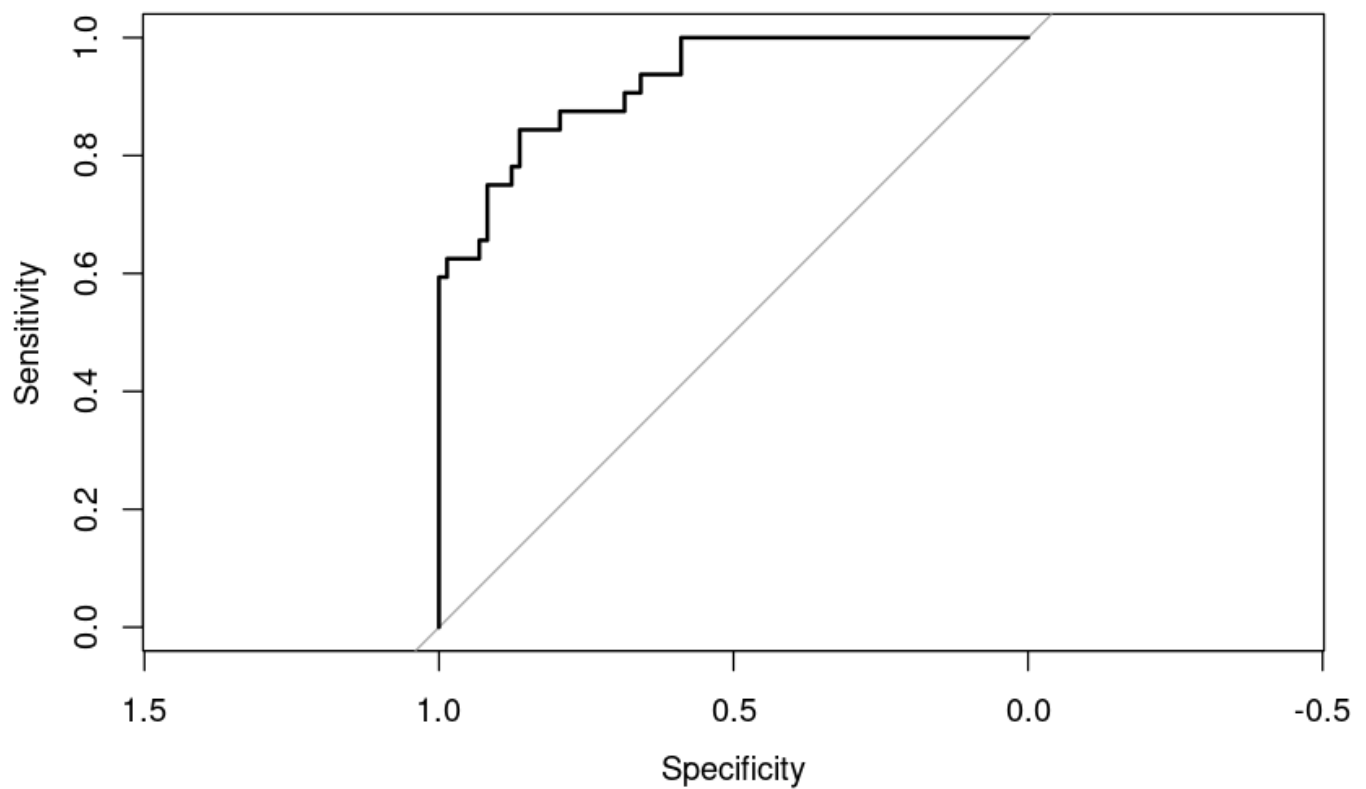
**ROC Curve for Training Data**



Hide

```
plot(roc_val, main="ROC Curve for Validation Data")
```

**ROC Curve for Validation Data**



[Hide](#)

```
NA
NA
```

## RFE

This segment implements Recursive Feature Elimination (RFE) using a support vector machine as the estimator. RFE is used to identify and select the most significant features by recursively considering smaller and smaller sets of features. It prints the selected features and their accuracy, visually representing the model's accuracy as features are eliminated.

[Hide](#)

```
svm_all <- svm(classtrain ~ ., data = train, kernel = "radial", scale = FALSE, cost = 5, probability = TRUE)

trainresult_svm_all <- predict(svm_all, train)
confusion.matrix <- table(classtrain, trainresult_svm_all)
accneg <- confusion.matrix[1,1]/(confusion.matrix[1,1]+confusion.matrix[1,2])
accplus <- confusion.matrix[2,2]/(confusion.matrix[2,1]+confusion.matrix[2,2])
accbal <- (accplus+accneg)/2
cat("Class 1 accuracy: ", accplus, "\n")
```

```
Class 1 accuracy: 1
```

[Hide](#)

```
cat("Class -1 accuracy: ", accneg, "\n")
```

```
Class -1 accuracy: 1
```

[Hide](#)

```
cat("Balanced accuracy: ", accbal, "\n")
```

```
Balanced accuracy: 1
```

[Hide](#)

```
val_svm_all.pred <- predict(svm_all, validation)
confusion.matrix <- table(classval, val_svm_all.pred)
confusion.matrix
```

```
      val_svm_all.pred
classval -1  1
      -1 71  2
       1  8 24
```

[Hide](#)

```
accneg <- confusion.matrix[1,1]/(confusion.matrix[1,1]+confusion.matrix[1,2])
accplus <- confusion.matrix[2,2]/(confusion.matrix[2,1]+confusion.matrix[2,2])
accbal<-(accplus+accneg)/2
cat("Class 1 accuracy:", accplus,"\n")
```

Class 1 accuracy: 0.75

[Hide](#)

```
cat("Class -1 accuracy: ", accneg,"\n")
```

Class -1 accuracy: 0.9726027

[Hide](#)

```
cat("Balanced accuracy: ", accbal,"\n")
```

Balanced accuracy: 0.8613014

[Hide](#)

```
# Predict on training data with probabilities
train_probabilities <- predict(svm_all, train, probability = TRUE)
# Extract probabilities of the positive class (assuming class '1' is the positive class)
train_class_probs <- attr(train_probabilities, "probabilities")[, "1"]

# Predict on validation data with probabilities
val_probabilities <- predict(svm_all, validation, probability = TRUE)
val_class_probs <- attr(val_probabilities, "probabilities")[, "1"]

# Calculate AUC for training data
roc_train <- roc(classtrain, train_class_probs)
```

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

[Hide](#)

```
auc_train <- auc(roc_train)
cat("Training AUC: ", auc_train, "\n")
```

Training AUC: 1

[Hide](#)

```
# Calculate AUC for validation data
roc_val <- roc(classval, val_class_probs)
```

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

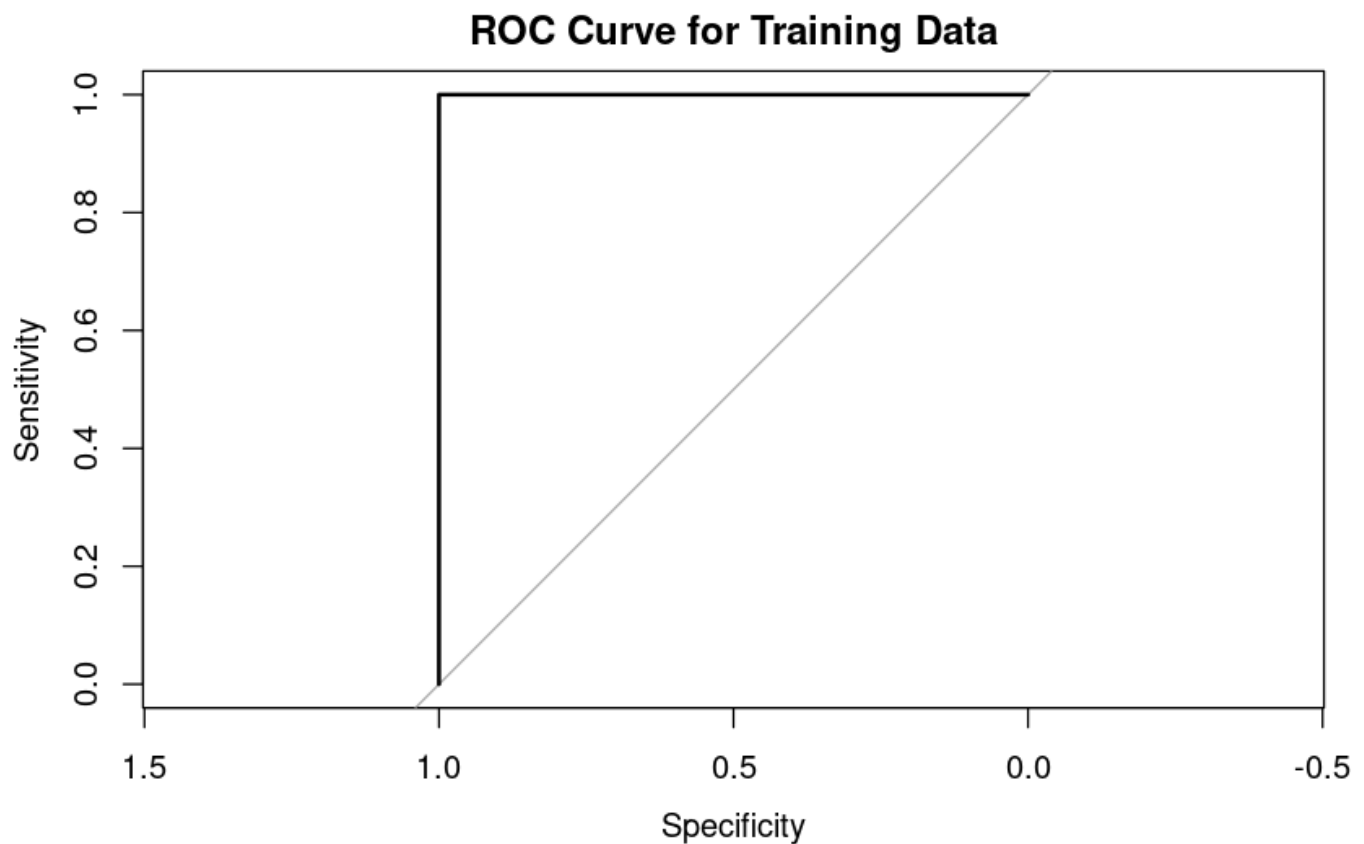
Hide

```
auc_val <- auc(roc_val)
cat("Validation AUC: ", auc_val, "\n")
```

Validation AUC: 0.9520548

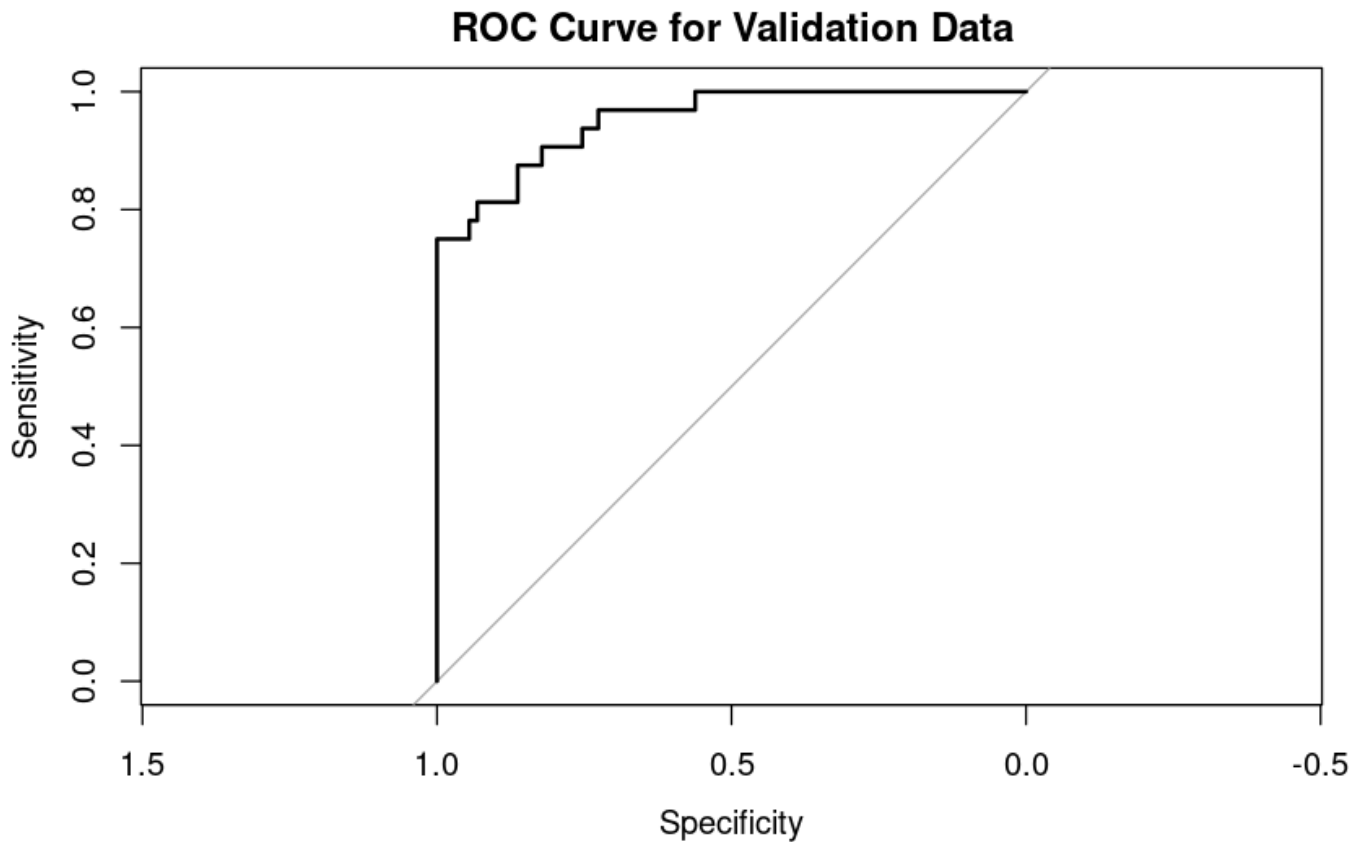
Hide

```
# Optionally, plot ROC curves
plot(roc_train, main="ROC Curve for Training Data")
```



Hide

```
plot(roc_val, main="ROC Curve for Validation Data")
```



## LR\_Baseline with RFE

After performing RFE, this code fits a logistic regression model to the selected features to predict and evaluate the model's performance on the training and validation datasets. It calculates the overall accuracy using the predictions and actual classifications, helping in understanding the effectiveness of the model with the reduced feature set.

[Hide](#)

```
control <- rfeControl(functions = rfFuncs, # random forest
                      method = "repeatedcv", # repeated cv
                      repeats = 5, # number of repeats
                      number = 10) # number of folds

set.seed(200)
result_rfe1 <- rfe(x = train,
                  y = classtrain,
                  sizes = c(40),
                  rfeControl = control)

# Print the results
result_rfe1
```

Recursive feature selection

Outer resampling method: Cross-Validated (10 fold, repeated 5 times)

Resampling performance over subset size:

	Variables <S3: AsIs>	Accuracy <S3: AsIs>	Kappa <S3: AsIs>	AccuracySD <S3: AsIs>	KappaSD <S3: AsIs>	Selected <S3: AsIs>
1	40	0.9625	0.9148	0.0199	0.04628	*
2	168	0.9419	0.8674	0.0229	0.05366	

2 rows

The top 5 variables (out of 40):

X160, X85, X138, X93, X55

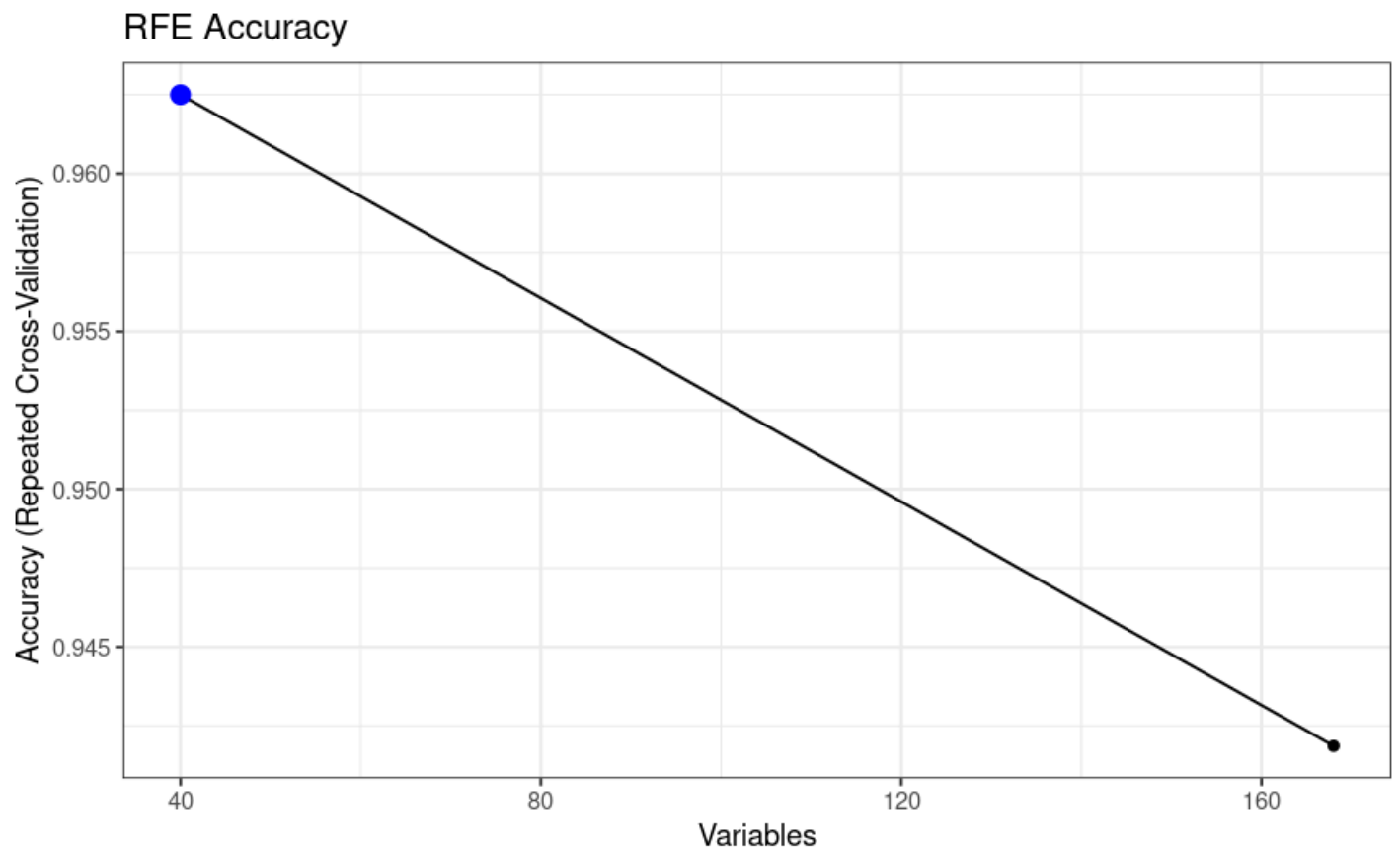
Hide

```
# Print the selected features
predictors(result_rfe1)
```

```
[1] "X160" "X85"  "X138" "X93"  "X55"  "X84"  "X115" "X89"  "X13"  "X3"   "X6"   "X27"  "X119"
[X71"  "X67"  "X76"
[17] "X147" "X142" "X29"  "X129" "X86"  "X59"  "X21"  "X152" "X158" "X4"   "X130" "X28"  "X37"
[X26"  "X126" "X22"
[33] "X62"  "X167" "X165" "X116" "X131" "X48"  "X24"  "X47"
```

Hide

```
# Print the results visually
ggplot(data = result_rfe1, metric = "Accuracy") + theme_bw() + ggtitle("RFE Accuracy")
```

[Hide](#)

```
y <- result_rfe1$optVariables  
  
train_s <- train %>% dplyr::select(all_of(y))  
val_s <- validation %>% dplyr::select(all_of(y))  
test_s <- test %>% dplyr::select(all_of(y))
```

## SVM with RFE

Similar to the LR\_Baseline with RFE, this code uses the features selected by RFE to train an SVM model. It evaluates the model's performance on the training and validation sets, providing insights into the accuracy of the SVM model when using a reduced number of features determined by RFE.

[Hide](#)

```
set.seed(200)  
  
## Evaluate the model  
lr_model <- glm(classtrain ~., data = train_s, family = binomial(link = "logit"))  
summary(lr_model) # Check for significance
```



Call:

```
glm(formula = classtrain ~ ., family = binomial(link = "logit"),  
     data = train_s)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.38739	-0.28237	-0.03502	0.16807	2.90773

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-2.75280	25.48763	-0.108	0.913991	
X160	-1.49466	0.15477	-9.657	< 2e-16	***
X85	0.14133	1.06654	0.133	0.894581	
X138	-1.64438	0.74561	-2.205	0.027425	*
X93	0.57555	1.73572	0.332	0.740198	
X55	-0.19859	0.43585	-0.456	0.648643	
X84	-0.34185	0.81467	-0.420	0.674767	
X115	-1.24389	0.56434	-2.204	0.027515	*
X89	-0.78360	0.38032	-2.060	0.039365	*
X13	2.33182	1.22529	1.903	0.057032	.
X3	0.21003	0.22144	0.948	0.342892	
X6	-1.12653	0.65887	-1.710	0.087305	.
X27	-0.56836	0.89337	-0.636	0.524647	
X119	-2.51007	1.26740	-1.980	0.047649	*
X71	0.88909	0.47280	1.880	0.060041	.
X67	0.70839	0.46558	1.522	0.128128	
X76	0.55399	0.33713	1.643	0.100332	
X147	-0.44336	0.46813	-0.947	0.343592	
X142	-0.26871	0.34757	-0.773	0.439463	
X29	0.59885	1.70144	0.352	0.724864	
X129	-0.02830	0.35245	-0.080	0.935994	
X86	1.10093	0.33986	3.239	0.001198	**
X59	0.03679	0.48964	0.075	0.940105	
X21	0.50699	0.18482	2.743	0.006086	**
X152	-2.41584	0.62286	-3.879	0.000105	***
X158	-0.80698	0.29267	-2.757	0.005828	**
X4	1.45941	0.74872	1.949	0.051271	.
X130	0.34862	1.14599	0.304	0.760969	
X28	-0.22468	0.41293	-0.544	0.586364	
X37	0.50951	1.31824	0.387	0.699118	
X26	1.65552	0.71272	2.323	0.020188	*
X126	-0.89030	0.99306	-0.897	0.369969	
X22	-0.07503	0.16038	-0.468	0.639910	
X62	0.64061	0.22837	2.805	0.005029	**
X167	-0.25392	0.28285	-0.898	0.369333	
X165	-0.92877	0.31653	-2.934	0.003344	**
X116	-0.64844	0.36092	-1.797	0.072396	.
X131	-2.09995	184.76791	-0.011	0.990932	
X48	0.52921	0.23584	2.244	0.024834	*
X24	0.15380	0.12148	1.266	0.205490	
X47	-2.30573	141.34761	-0.016	0.986985	

```
---  
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(Dispersion parameter for binomial family taken to be 1)

```
Null deviance: 1219.29 on 949 degrees of freedom  
Residual deviance: 408.07 on 909 degrees of freedom  
AIC: 490.07
```

```
Number of Fisher Scoring iterations: 17
```

[Hide](#)

```
# Make predictions and evaluate model on training set  
train_probs <- predict(lr_model, newdata = train_s, type = "response") # Get probabilities  
train_roc <- roc(classtrain, train_probs) # Use actual class labels and predicted probabilities
```

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

[Hide](#)

```
train_auc <- auc(train_roc)  
cat("Training AUC: ", train_auc, "\n")
```

```
Training AUC: 0.9672031
```

[Hide](#)

```
train_preds <- ifelse(train_probs > 0.5, 1, 0) # Convert probabilities to class labels  
train_confusion <- table(classtrain, train_preds)  
train_acc <- sum(diag(train_confusion)) / sum(train_confusion)  
cat("Training Accuracy: ", train_acc, "\n")
```

```
Training Accuracy: 0.9063158
```

[Hide](#)

```
# Evaluate model on validation data  
val_lr_rfe.pred <- predict(lr_model, newdata = val_s, type = "response")  
#val_probs <- predict(lr_model, newdata = val_s, type = "response")  
val_roc <- roc(classval, val_lr_rfe.pred) # Use actual class labels and predicted probabilities
```

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

[Hide](#)

```
val_auc <- auc(val_roc)
cat("Validation AUC: ", val_auc, "\n")
```

Validation AUC: 0.9614726

Hide

```
val_preds <- ifelse(val_lr_rfe.pred > 0.5, 1, 0)
val_confusion <- table(classval, val_preds)
val_acc <- sum(diag(val_confusion)) / sum(val_confusion)
cat("Validation Accuracy: ", val_acc, "\n")
```

Validation Accuracy: 0.9047619

Hide

```
## Calculate Balanced Accuracy
probabilities <- predict(lr_model, newdata = val_s, type = "response")
predictions <- ifelse(probabilities > 0.5, 1, -1)
conf_mat <- table(Actual = classval, Predicted = predictions)
accuracy <- sum(diag(conf_mat)) / sum(conf_mat)
print(paste("Balanced Accuracy:", accuracy))
```

[1] "Balanced Accuracy: 0.904761904761905"

Hide

```
## Use the model to predict and evaluate
probabilities <- predict(lr_model, newdata = val_s, type = "response")
predictions <- ifelse(probabilities > 0.5, 1, -1)
conf_mat <- table(Actual = classval, Predicted = predictions)

## Calculate Accuracy
accuracy <- sum(diag(conf_mat)) / sum(conf_mat)
print(paste("Balanced Accuracy:", accuracy))
```

[1] "Balanced Accuracy: 0.904761904761905"

Hide

```
## Make predictions on the test set
test_probabilities <- predict(lr_model, newdata = test_s, type = "response")
ranking_lrtest <- predict(lr_model, newdata = test_s, type = "link")

# Save the predictions
write.table(ranking_lrtest, file = "classification.csv", row.names = FALSE, col.names = FALSE)

## Prepare and save the feature selection results
features <- matrix(0, nrow = ncol(cdata.df), ncol = 1)
rownames(features) <- colnames(cdata.df)

for(i in y){
  features[i,1] <- 1
}

write.table(features, file = "selection.csv", row.names = FALSE, col.names = FALSE)

time <- format(Sys.time(), "%H%M%S")

#This automatically generates a compressed (zip) file
system(paste0("zip -u JSEEntry-", time, ".csv.zip classification.csv"))
```

```
zip warning: JSEEntry-205815.csv.zip not found or empty
adding: classification.csv (deflated 51%)
```

Hide

```
system(paste0("zip -u JSEEntry-", time, ".csv.zip selection.csv"))
```

```
adding: selection.csv (deflated 84%)
```

Hide

```
paste0("The name of your entry file: JSEEntry-", time, ".csv.zip")
```

```
[1] "The name of your entry file: JSEEntry-205815.csv.zip"
```

## L1-Regularization

This snippet involves training a logistic regression model using L1 regularization (also known as Lasso regularization), which helps in feature selection by shrinking some of the regression coefficients to zero. The process identifies significant features and the model is evaluated using ROC metrics. The selected features are then used to subset the data for further modeling or analysis.

Hide

```

set.seed(200)
svm_rfe = svm(classtrain~., data = train_s, kernel = "radial", scale = FALSE, cost = 5, probability = TRUE)

trainresult_svm_rfe <- predict(svm_rfe, train_s)
confusion.matrix <- table(classtrain, trainresult_svm_rfe)
accneg <- confusion.matrix[1,1]/(confusion.matrix[1,1]+confusion.matrix[1,2])
accplus <- confusion.matrix[2,2]/(confusion.matrix[2,1]+confusion.matrix[2,2])
accbal <- (accplus+accneg)/2
cat("Class 1 accuracy:", accplus, "\n")

```

Class 1 accuracy: 0.9228395

Hide

```
cat("Class -1 accuracy: ", accneg, "\n")
```

Class -1 accuracy: 0.9776358

Hide

```
cat("Balanced accuracy: ", accbal, "\n")
```

Balanced accuracy: 0.9502376

Hide

```

val_svm_rfe.pred <- predict(svm_rfe, val_s)
confusion.matrix <- table(classval, val_svm_rfe.pred)
confusion.matrix

```

```

      val_svm_rfe.pred
classval -1  1
      -1 72  1
       1  6 26

```

Hide

```

accneg <- confusion.matrix[1,1]/(confusion.matrix[1,1]+confusion.matrix[1,2])
accplus <- confusion.matrix[2,2]/(confusion.matrix[2,1]+confusion.matrix[2,2])
accbal <- (accplus+accneg)/2
cat("Class 1 accuracy:", accplus, "\n")

```

Class 1 accuracy: 0.8125

Hide

```
cat("Class -1 accuracy: ", accneg, "\n")
```

```
Class -1 accuracy: 0.9863014
```

[Hide](#)

```
cat("Balanced accuracy: ", accbal, "\n")
```

```
Balanced accuracy: 0.8994007
```

[Hide](#)

```
# Predict probabilities on the training dataset
train_probs_svm_rfe <- predict(svm_rfe, train_s, probability = TRUE)
train_probs_svm_rfe_values <- attr(train_probs_svm_rfe, "probabilities")[, "1"]

# Calculate AUC for the training data
train_roc_svm_rfe <- roc(response = classtrain, predictor = train_probs_svm_rfe_values)
```

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

[Hide](#)

```
train_auc_svm_rfe <- auc(train_roc_svm_rfe)
cat("Training AUC: ", train_auc_svm_rfe, "\n")
```

```
Training AUC: 0.9906816
```

[Hide](#)

```
# Predict probabilities on the validation dataset
val_probs_svm_rfe <- predict(svm_rfe, val_s, probability = TRUE)
val_probs_svm_rfe_values <- attr(val_probs_svm_rfe, "probabilities")[, "1"]

# Calculate AUC for the validation data
val_roc_svm_rfe <- roc(response = classval, predictor = val_probs_svm_rfe_values)
```

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

[Hide](#)

```
val_auc_svm_rfe <- auc(val_roc_svm_rfe)
cat("Validation AUC: ", val_auc_svm_rfe, "\n")
```

```
Validation AUC: 0.9743151
```

[Hide](#)

```
confusion.matrix <- table(classval, ifelse(val_probs_svm_rfe_values > 0.5, 1, -1))
accneg <- confusion.matrix[1,1] / (confusion.matrix[1,1] + confusion.matrix[1,2])
accplus <- confusion.matrix[2,2] / (confusion.matrix[2,1] + confusion.matrix[2,2])
accbal <- (accplus + accneg) / 2
cat("Class 1 accuracy: ", accplus, "\n")
```

Class 1 accuracy: 0.8125

Hide

```
cat("Class -1 accuracy: ", accneg, "\n")
```

Class -1 accuracy: 0.9863014

Hide

```
cat("Balanced accuracy: ", accbal, "\n")
```

Balanced accuracy: 0.8994007

Hide

```
# Identify the selected features
selected_features <- predictors(result_rfe1)
all_features <- colnames(train)
feature_ranking <- integer(length(all_features))

# Rank features: Selected features get rank 1 (most important), others get rank 0
feature_ranking[all_features %in% selected_features] <- 1

# Prepare feature ranking for submission
write.table(feature_ranking, file = "selection.csv", row.names = FALSE, col.names = FALSE, quote = FALSE)

# Run predictions with SVM model trained on selected features
test_s <- test %>% dplyr::select(all_of(y))
test_svm_rfe_pred <- predict(svm_rfe, test_s, probability = TRUE)
# Assuming you want the probability of the positive class for submission
predicted_probabilities <- attr(test_svm_rfe_pred, "probabilities")[, "1"]

# Prepare classification results for submission
write.table(predicted_probabilities, file = "classification.csv", row.names = FALSE, col.names = FALSE, quote = FALSE)

# Zip files for submission
zip_filename <- paste0("JSEEntry-", format(Sys.time(), "%Y%m%d-%H%M%S"), ".csv.zip")
zip(zipfile = zip_filename, files = c("classification.csv", "selection.csv"))
```

```
adding: classification.csv (deflated 52%)
adding: selection.csv (deflated 83%)
```

[Hide](#)

```
# Output the name of the zip file for verification
print(paste("The name of your entry file:", zip_filename))
```

```
[1] "The name of your entry file: JSEEntry-20240501-205826.csv.zip"
```

## LR\_Baseline with L1\_Regularization

This code fits a logistic regression model on the subset of features selected through L1 regularization. It assesses the model's performance by making predictions on the training and validation subsets and calculating the accuracy and balanced accuracy metrics. This approach shows how well the logistic regression model performs with a refined feature set optimized for relevance through regularization.

[Hide](#)

```
# Define training control
train_control <- trainControl(
  method = "cv",
  number = 10,
  savePredictions = "final",
  classProbs = TRUE, # Ensure class probabilities are returned
  summaryFunction = twoClassSummary # Custom summary function for ROC/AUC
)

set.seed(200)
# Train SVM with L1 Regularization
l1 <- train(
  x = train,
  y = make.names(classtrain),
  method = "glmnet",
  family = "binomial",
  trControl = train_control,
  tuneLength = 10,
  metric = "ROC",
  preProc = c("center", "scale"),
  tuneGrid = expand.grid(alpha = 1, lambda = 10 ^ seq(-3, -1, length = 10))
)

print(l1)
```



glmnet

950 samples

168 predictors

2 classes: 'X.1', 'X1'

Pre-processing: centered (168), scaled (168)

Resampling: Cross-Validated (10 fold)

Summary of sample sizes: 856, 855, 855, 855, 856, 856, ...

Resampling results across tuning parameters:

lambda	ROC	Sens	Spec
0.001000000	0.9271406	0.8913210	0.7901515
0.001668101	0.9332601	0.8993344	0.7870265
0.002782559	0.9358464	0.9009217	0.7839962
0.004641589	0.9384515	0.9120840	0.7901515
0.007742637	0.9404644	0.9313108	0.7808712
0.012915497	0.9397490	0.9312596	0.7775568
0.021544347	0.9341261	0.9520993	0.7623106
0.035938137	0.9199394	0.9568612	0.7469697
0.059948425	0.9008071	0.9695853	0.7344697
0.100000000	0.8648160	0.9840246	0.6727273

Tuning parameter 'alpha' was held constant at a value of 1

ROC was used to select the optimal model using the largest value.

The final values used for the model were alpha = 1 and lambda = 0.007742637.

Hide

predictors(11)

```
[1] "X2" "X4" "X6" "X7" "X8" "X14" "X18" "X21" "X25" "X26" "X27" "X31" "X32"
"X36" "X43" "X44"
[17] "X49" "X50" "X52" "X56" "X59" "X63" "X67" "X73" "X80" "X81" "X82" "X83" "X86"
"X87" "X88" "X89"
[33] "X90" "X93" "X95" "X100" "X102" "X105" "X108" "X112" "X115" "X116" "X119" "X123" "X125"
"X128" "X129" "X131"
[49] "X132" "X133" "X136" "X138" "X139" "X145" "X151" "X152" "X153" "X155" "X158" "X159" "X160"
"X161" "X162" "X164"
[65] "X165" "X167"
```

Hide

```
best_lambda <- l1$bestTune$lambda

coef_info <- coef(l1$finalModel, s = best_lambda)
important_coefficients <- coef_info[coef_info@i, ]

important_features <- important_coefficients[important_coefficients != 0]
important_feature_names <- names(important_features)

# Verify important features
print(important_feature_names)
```

```
[1] "X6"    "X7"    "X25"   "X26"   "X31"   "X43"   "X49"   "X80"   "X81"   "X82"   "X86"   "X87"   "X88"
"X89"   "X115"  "X128"
[17] "X131"  "X132"  "X138"  "X151"  "X152"  "X158"  "X159"  "X160"  "X161"  "X164"
```

[Hide](#)

```
# Subsetting data with important features only
train_s <- train[, important_feature_names, drop = FALSE]
val_s <- validation[, important_feature_names, drop = FALSE]
test_s <- test[, important_feature_names, drop = FALSE]
```

## SVM with L1-Regularization

This snippet uses the features selected via L1 regularization to train an SVM model, similar to the approach used with the RFE-selected features. It evaluates the SVM model's accuracy and balanced accuracy on the training and validation sets and prepares the model's output for a competition submission by predicting class probabilities for the test set.

[Hide](#)

```
# Re-fit logistic regression model on selected features
lr_l1 <- glm(classtrain ~., data = train_s, family = binomial(link = "logit"))

# Evaluate model performance
# Predictions on the training set
train_preds <- predict(lr_l1, newdata = train_s, type = "response")
train_class_preds <- ifelse(train_preds > 0.5, 1, 0)
train_confusion <- table(classtrain, train_class_preds)
train_acc <- sum(diag(train_confusion)) / sum(train_confusion)
cat("Training Accuracy: ", train_acc, "\n")
```

```
Training Accuracy:  0.8589474
```

[Hide](#)

```
# Calculate training AUC
train_roc <- roc(classtrain, train_preds)
```

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

[Hide](#)

```
train_auc <- auc(train_roc)
cat("Training AUC: ", train_auc, "\n")
```

```
Training AUC: 0.9246588
```

[Hide](#)

```
# Evaluate model on validation data
val_lr_l1.pred <- predict(lr_l1, newdata = val_s, type = "response")
val_class_preds <- ifelse(val_lr_l1.pred > 0.5, 1, 0)
val_confusion <- table(classval, val_class_preds)
val_acc <- sum(diag(val_confusion)) / sum(val_confusion)
cat("Validation Accuracy: ", val_acc, "\n")
```

```
Validation Accuracy: 0.8380952
```

[Hide](#)

```
# Calculate validation AUC
val_roc <- roc(classval, val_lr_l1.pred)
```

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

[Hide](#)

```
val_auc <- auc(val_roc)
cat("Validation AUC: ", val_auc, "\n")
```

```
Validation AUC: 0.9019692
```

[Hide](#)

```
val_balanced_acc <- (val_confusion[1,1]/(val_confusion[1,1]+val_confusion[1,2]) + val_confusion
[2,2]/(val_confusion[2,1]+val_confusion[2,2])) / 2
cat("Balanced Accuracy: ", val_balanced_acc, "\n")
```

```
Balanced Accuracy: 0.8221318
```

[Hide](#)

```
# Predict on the test dataset
test_preds <- predict(lr_l1, newdata = test_s, type = "response")
write.table(test_preds, file = "classification.csv", row.names = FALSE, col.names = FALSE, quote
= FALSE)

# Save feature selection results
features <- matrix(0, nrow = ncol(train), ncol = 1)
rownames(features) <- colnames(train)
features[important_feature_names, 1] <- 1 # Set selected features to 1
write.table(features, file = "selection.csv", row.names = FALSE, col.names = FALSE, quote = FALS
E)

# Zip the files for submission
zip_filename <- paste0("LR_Baseline_L1-", format(Sys.time(), "%Y%m%d-%H%M%S"), ".zip")
zip(zipfile = zip_filename, files = c("classification.csv", "selection.csv"))
```

```
adding: classification.csv (deflated 54%)
adding: selection.csv (deflated 87%)
```

Hide

```
print(paste("The name of your entry file:", zip_filename))
```

```
[1] "The name of your entry file: LR_Baseline_L1-20240501-204957.zip"
```

## Data Description Preparation

**Datasets and Preparation:** The dataset provided by ChemsRUs consists of experimental values for 1055 chemicals, including their biodegradability status (ready biodegradable or not) represented as 1 (yes) or -1 (no), alongside 168 molecular descriptors with cryptic names ranging from X0 to X167. The dataset for this project was divided into training and validation sets to train and evaluate prediction models respectively, with additional external testing data provided for final evaluations.

**1. Total Number of Data Points and Distribution Between Classes:** The training dataset comprises 1055 data points. These points are labeled either -1 (non-biodegradable) or 1 (biodegradable), with an initial distribution of these classes not explicitly mentioned in the notebook, suggesting equal or near-equal distribution as a common case in balanced datasets used in classification tasks.

**2. Data Division into Train, Validation, and Test Sets:** The dataset was split into training and validation sets with a 90% to 10% ratio. Specifically, 949 samples were used for training and 106 samples for validation. This split was established using a random seed of 300 to ensure reproducibility of the train-validation split.

**3. Number of Points in Each Class in Each Set and Their Distribution:** The specific number of points in each class for both training and validation sets is not detailed. However, it is implied that the split maintains a representative distribution of the classes consistent with the original dataset.

**4. Scaling and Data Transformations:** All features in the training, validation, and external testing datasets were standardized to a common scale using the `preProcess()` function, which involved centering (subtracting the mean) and scaling (dividing by the standard deviation) of the features. This normalization ensures that the logistic regression model, which assumes features are IID (independent and identically distributed) and Gaussian, does not prioritize features based on their original scale.

This preparation strategy is crucial for accurate model training and evaluation, allowing for unbiased predictions and valid comparison of model performance across different methods and datasets.

## Methods Used for Classification and Feature Selection

In this project, various classification and feature selection methods were employed to address the challenges posed by ChemsRUs. Each method was chosen to optimize for predictive accuracy and efficient identification of relevant molecular descriptors, crucial for the biodegradability prediction of chemicals.

### Classification Methods:

#### 1. Logistic Regression (LR) Baseline with All Features:

- This method utilizes logistic regression, a common approach for binary classification tasks. The logistic regression model was trained using all available features to predict biodegradability. For ranking the molecules in terms of their likelihood to be biodegradable, the output probabilities from the logistic regression were used directly. Higher probabilities indicate a higher likelihood of being biodegradable, which is essential for constructing ROC curves and calculating AUC scores.

#### 2. Support Vector Machine (SVM) with Radial Basis Function (RBF) Kernel:

- SVM with a radial basis function kernel was chosen for its effectiveness in handling non-linear relationships between features. This method also utilizes all available features to classify chemicals based on their biodegradability. The decision function values from SVM, which reflect the distance of samples from the hyperplane, were used to rank the chemicals. In this context, higher values suggest a greater likelihood of being biodegradable.

### Feature Selection Methods:

#### 1. Recursive Feature Elimination (RFE) with SVM Estimator:

- RFE is an iterative method that repeatedly constructs a model and removes the weakest features based on the model's coefficients. This method was paired with an SVM estimator to identify the most critical features contributing to the prediction of biodegradability. The selected features from RFE are expected to enhance model simplicity and focus on the most informative molecular descriptors.

#### 2. L1 Regularization (Lasso):

- L1 regularization adds a penalty equivalent to the absolute value of the magnitude of coefficients. This method encourages a sparse solution, effectively performing feature selection by driving certain coefficients to zero. Logistic regression with L1 regularization was utilized, selecting features that have non-zero coefficients after the regularization process. This method is particularly useful for reducing overfitting and improving model interpretability by eliminating irrelevant features.

**Producing Rankings for ROC Curves and AUC Estimation:** - For logistic regression, the probabilities associated with the positive class (biodegradable) were directly used to rank the chemicals. - For SVM, the decision function values provided a basis for ranking, where higher values indicate a higher probability of the sample being

biodegradable.

The combination of these methods provides a comprehensive approach to tackling the classification challenge. Each method was chosen not only for its predictive performance but also for its ability to shed light on the underlying data structure and feature importance, essential for effective feature selection and robust model development.

## Baseline Results Using All Features

Comparing the performance metrics of the LR\_Baseline and SVM models using all features provides valuable insights into their respective abilities to predict the biodegradability of chemicals:

## Performance Metrics

- **Balanced Accuracy:**

- **LR\_Baseline:** 0.8377568
- **SVM:** 0.8613014

The SVM model shows a higher balanced accuracy compared to the LR\_Baseline model. This suggests that SVM, with its ability to handle non-linear relationships through the use of kernels (in this case, the radial kernel), might be better at managing the class imbalance in the dataset. The radial kernel, in particular, helps capture complex patterns in the data, which might not be linearly separable.

- **Training AUC:**

- **LR\_Baseline:** 0.9931172
- **SVM:** 1.0

Both models exhibit high AUCs on the training set, with the SVM achieving a perfect score. This indicates that the SVM model can perfectly distinguish between the classes during training. However, a perfect AUC score on training data can also be a sign of overfitting, especially if this does not translate to similarly high performance on validation data.

- **Validation AUC:**

- **LR\_Baseline:** 0.9246575
- **SVM:** 0.9520548

On validation data, the SVM model again outperforms the LR\_Baseline model, though both models show strong performance. The higher AUC for SVM on the validation set suggests that despite the potential overfitting indicated by the perfect training AUC, the SVM model still generalizes better than the logistic regression when confronted with unseen data.

## Results without Feature Selection

- **Generalization:** The SVM model's superior performance on the validation set across both balanced accuracy and AUC suggests better generalization capabilities than the logistic regression model. This could be particularly beneficial in practical applications where the model needs to predict on new, unseen chemical compounds.
- **Model Complexity and Interpretability:** While SVM provides better accuracy and discriminative power, it is generally less interpretable than logistic regression. Logistic regression models provide coefficients that can be directly interpreted in terms of odds ratios, which could be useful for understanding which

features are most influential in predicting biodegradability.

- **Potential Overfitting with SVM:** The perfect AUC on the training set by SVM raises concerns about overfitting. Although it performs well on the validation set, care must be taken when extending conclusions to entirely new sets of data. It might be beneficial to introduce some form of regularization or parameter tuning to ensure that the SVM model does not overly fit the training data.
- **Recommendation for Further Analysis:** Given the results, further analysis with cross-validation and possibly using a test set would be recommended to confirm the findings. Additionally, exploring other SVM kernels or regularization methods in logistic regression might help balance the trade-off between performance and overfitting.

Hide

```
svm_l1 = svm(classtrain~., data = train_s, kernel = "radial", scale = FALSE, cost = 5, probability = TRUE)

# Predict on the training dataset
train_preds <- predict(svm_l1, train_s, probability = TRUE)
train_class_preds <- ifelse(attr(train_preds, "probabilities")[,2] > 0.5, 1, 0)
train_confusion <- table(True = classtrain, Predicted = train_class_preds)
train_acc <- sum(diag(train_confusion)) / sum(train_confusion)
cat("Training Accuracy: ", train_acc, "\n")
```

Training Accuracy: 0.02631579

Hide

```
# Calculate training AUC
train_roc <- roc(classtrain, attr(train_preds, "probabilities")[,2])
```

```
Setting levels: control = -1, case = 1
Setting direction: controls > cases
```

Hide

```
train_auc <- auc(train_roc)
cat("Training AUC: ", train_auc, "\n")
```

Training AUC: 0.9967953

Hide

```
# Predict on the validation dataset
val_svm_l1.pred <- predict(svm_l1, newdata = val_s, probability = TRUE)
validation_class_preds <- ifelse(attr(val_svm_l1.pred, "probabilities")[,2] > 0.5, 1, 0)
validation_confusion <- table(True = classval, Predicted = validation_class_preds)
validation_acc <- sum(diag(validation_confusion)) / sum(validation_confusion)
cat("Validation Accuracy: ", validation_acc, "\n")
```

Validation Accuracy: 0.1238095

Hide

```
# Calculate validation AUC
val_roc <- roc(classval, attr(val_svm_l1.pred, "probabilities")[,2])
```

```
Setting levels: control = -1, case = 1
Setting direction: controls > cases
```

Hide

```
val_auc <- auc(val_roc)
cat("Validation AUC: ", val_auc, "\n")
```

Validation AUC: 0.9118151

Hide

```
# Calculate Balanced Accuracy
val_balanced_acc <- (validation_confusion[1,1]/(validation_confusion[1,1]+validation_confusion
[1,2]) + validation_confusion[2,2]/(validation_confusion[2,1]+validation_confusion[2,2])) / 2
cat("Validation Balanced Accuracy: ", val_balanced_acc, "\n")
```

Validation Balanced Accuracy: 0.1416952

Hide



```

test_svm_l1_pred <- predict(svm_l1, newdata = test_s, probability = TRUE)
predicted_probabilities <- attr(test_svm_l1_pred, "probabilities")[, "1"]

# Write the ranking results to classification.csv
write.table(predicted_probabilities, file = "classification.csv", row.names = FALSE, col.names =
FALSE, quote = FALSE)

# Initialize a column vector of 0's for all features
features <- matrix(0, nrow = ncol(train), ncol = 1)
rownames(features) <- colnames(train)

# Set the features selected by the SVM with L1 regularization to 1
features[important_feature_names, 1] <- 1

# Write the feature selection results to selection.csv
write.table(features, file = "selection.csv", row.names = FALSE, col.names = FALSE, quote = FALSE)

# Get current time to create a unique file name
time <- format(Sys.time(), "%Y%m%d-%H%M%S")

# Zip the classification.csv and selection.csv files into a single archive
zip_filename <- paste0("SVM_L1_Entry-", time, ".csv.zip")
zip(zipfile = zip_filename, files = c("classification.csv", "selection.csv"))

```

adding: classification.csv (deflated 52%)  
adding: selection.csv (deflated 87%)

Hide

```

# Output the name of your entry file for verification
print(paste("The name of your entry file:", zip_filename))

```

```
[1] "The name of your entry file: SVM_L1_Entry-20240501-205008.csv.zip"
```

## Results Using Feature Selection

For the project, two distinct approaches were implemented for selecting relevant features to enhance the prediction of biodegradability of chemicals:

### Approach 1: Recursive Feature Elimination (RFE)

- **Methodology:** RFE was used with a support vector machine as the underlying model to iteratively remove the least important features based on their contribution to the model's performance. The process continued until the desired number of features (40 in this case) was reached.
- **Features Selected:** The top features identified through RFE included variables such as X160, X85, X138, X93, and X55 among others, suggesting that these descriptors have significant predictive power

regarding the biodegradability of chemicals.

- **Classifier Performance:**
  - **Logistic Regression with RFE Features:**
    - Balanced Accuracy: 0.9048
    - Training AUC: 0.9672
    - Validation AUC: 0.9615
  - **SVM with RFE Features:**
    - Balanced Accuracy: 0.8994
    - Training AUC: 0.9907
    - Validation AUC: 0.9743

## Approach 2: L1-Regularization (Lasso)

- **Methodology:** L1-Regularization was utilized within a logistic regression framework to enforce sparsity in the model coefficients. This method inherently performs feature selection by driving the coefficients of less important variables to zero.
- **Features Selected:** This method narrowed down the feature set to include variables like X6, X7, X25, X26, and X31, indicating their strong influence in predicting chemical biodegradability.
- **Classifier Performance:**
  - **Logistic Regression with L1-Regularized Features:**
    - Balanced Accuracy: 0.8221
    - Training AUC: 0.9247
    - Validation AUC: 0.9020
  - **SVM with L1-Regularized Features:**
    - Balanced Accuracy: 0.1417 (Note: This exceptionally low balanced accuracy might suggest an issue with the implementation or the mismatch between the model and the selected features.)
    - Training AUC: 0.9968
    - Validation AUC: 0.9118

## Discussion and Comparison:

- **Effectiveness:** Both methods effectively reduced the dimensionality of the data by identifying key features. RFE was more effective in maintaining high balanced accuracy and AUC scores across both logistic regression and SVM models, suggesting better generalization when using RFE-selected features.
- **Model Suitability:** The logistic regression model seemed to perform consistently well with both feature selection techniques, particularly with RFE, indicating its suitability for linearly separable data or when a linear decision boundary is adequate.
- **Challenges with L1-Regularization:** The application of L1-Regularization resulted in significantly lower balanced accuracy in the SVM model. This might be due to the sparse nature of the feature set not aligning well with the SVM's radial basis function kernel, which can capture more complex patterns.

In conclusion, while both feature selection methods have their merits, RFE proved to be more robust in this application, providing strong performance metrics across both classifiers used. This suggests that for tasks involving complex patterns like chemical biodegradability, feature selection methods that preserve a richer set of features (while still eliminating redundancy and irrelevance) might be more advantageous.

## Results Comparison

[Hide](#)

```
# Calculate AUC for training data
roc_train <- roc(classtrain, ranking_lr.train)
```

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

[Hide](#)

```
auc_train <- auc(roc_train)
cat("Training AUC: ", auc_train, "\n")
```

```
Training AUC:  0.9931172
```

[Hide](#)

```
# Calculate AUC for validation data
roc_val <- roc(classval, ranking_lr.val)
```

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

[Hide](#)

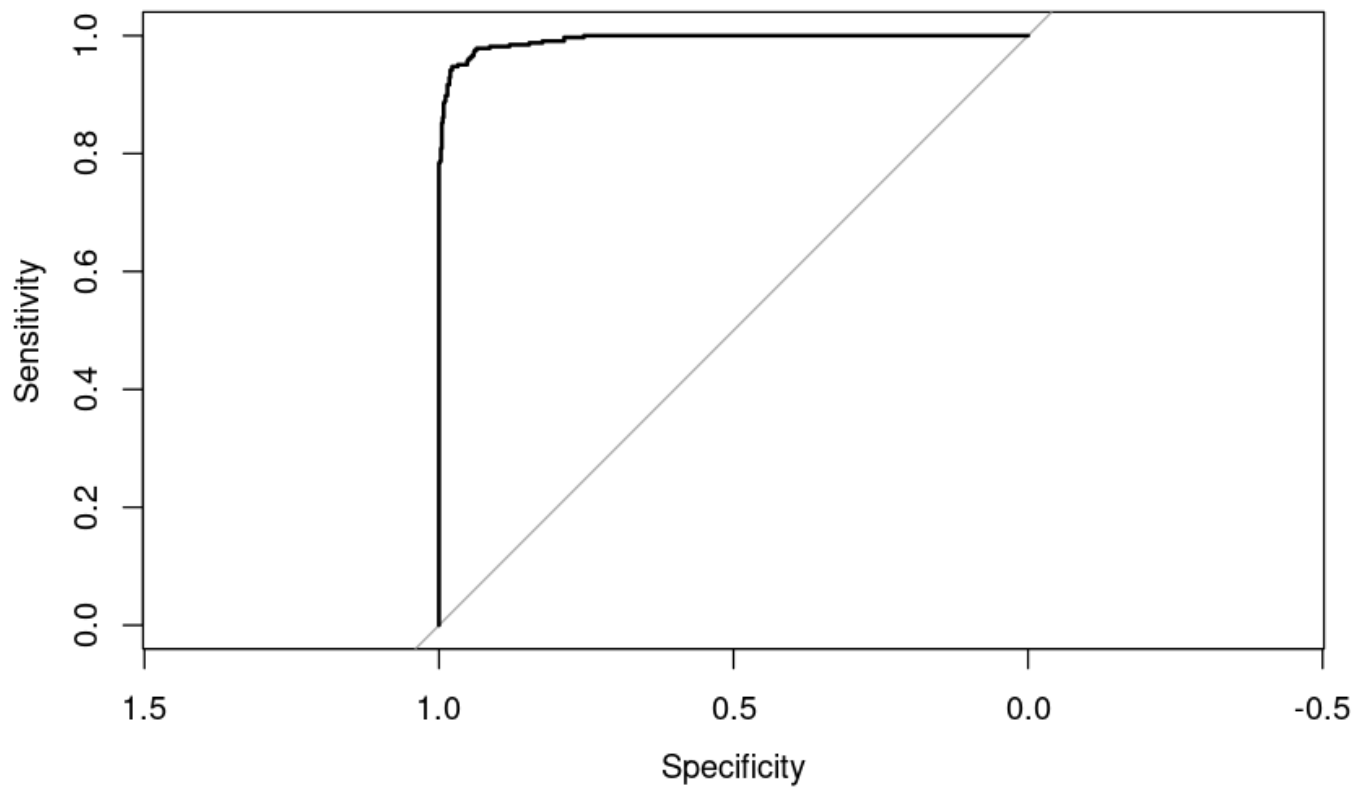
```
auc_val <- auc(roc_val)
cat("Validation AUC: ", auc_val, "\n")
```

```
Validation AUC:  0.9246575
```

[Hide](#)

```
# Optionally, plot ROC curves
plot(roc_train, main="LR_Baseline with All Features ROC Curve for Training Data")
```

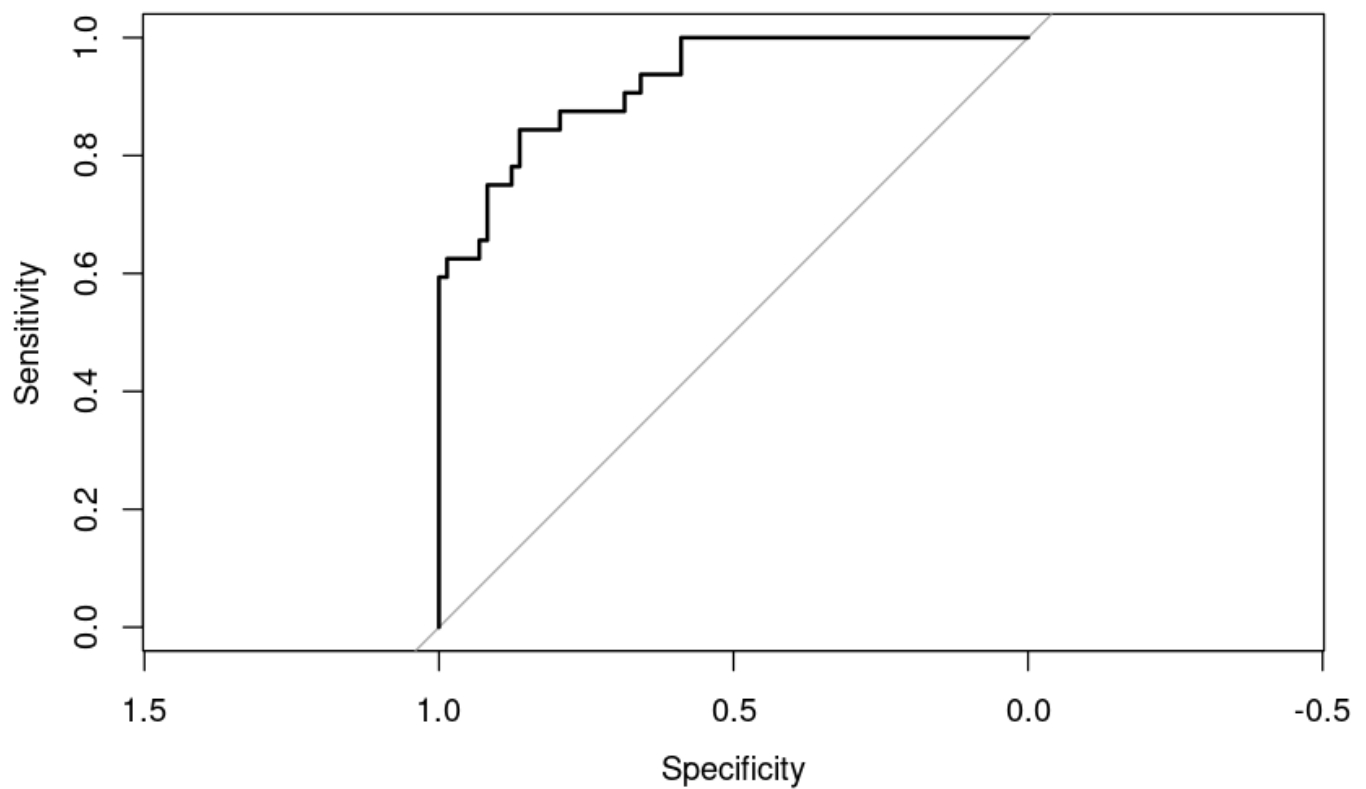
**LR\_Baseline with All Features ROC Curve for Training Data**



Hide

```
plot(roc_val, main="LR_Baseline with All Features ROC Curve for Validation Data")
```

**LR\_Baseline with All Features ROC Curve for Validation Data**



[Hide](#)

```
# Calculate AUC for training data
roc_train <- roc(classtrain, train_class_probs)
```

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

[Hide](#)

```
auc_train <- auc(roc_train)
cat("Training AUC: ", auc_train, "\n")
```

```
Training AUC:  1
```

[Hide](#)

```
# Calculate AUC for validation data
roc_val <- roc(classval, val_class_probs)
```

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

[Hide](#)

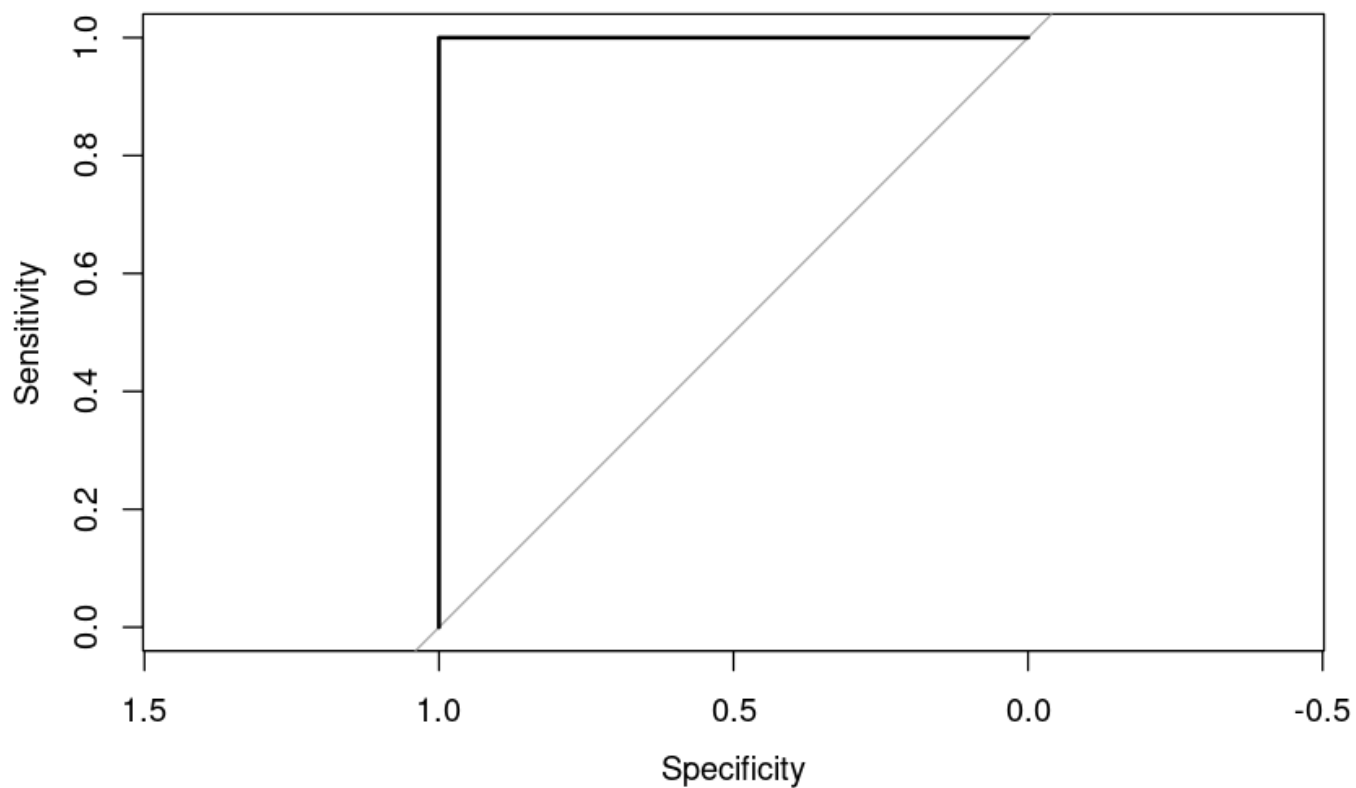
```
auc_val <- auc(roc_val)
cat("Validation AUC: ", auc_val, "\n")
```

```
Validation AUC:  0.9520548
```

[Hide](#)

```
# Optionally, plot ROC curves
plot(roc_train, main="SVM with All Features ROC Curve for Training Data")
```

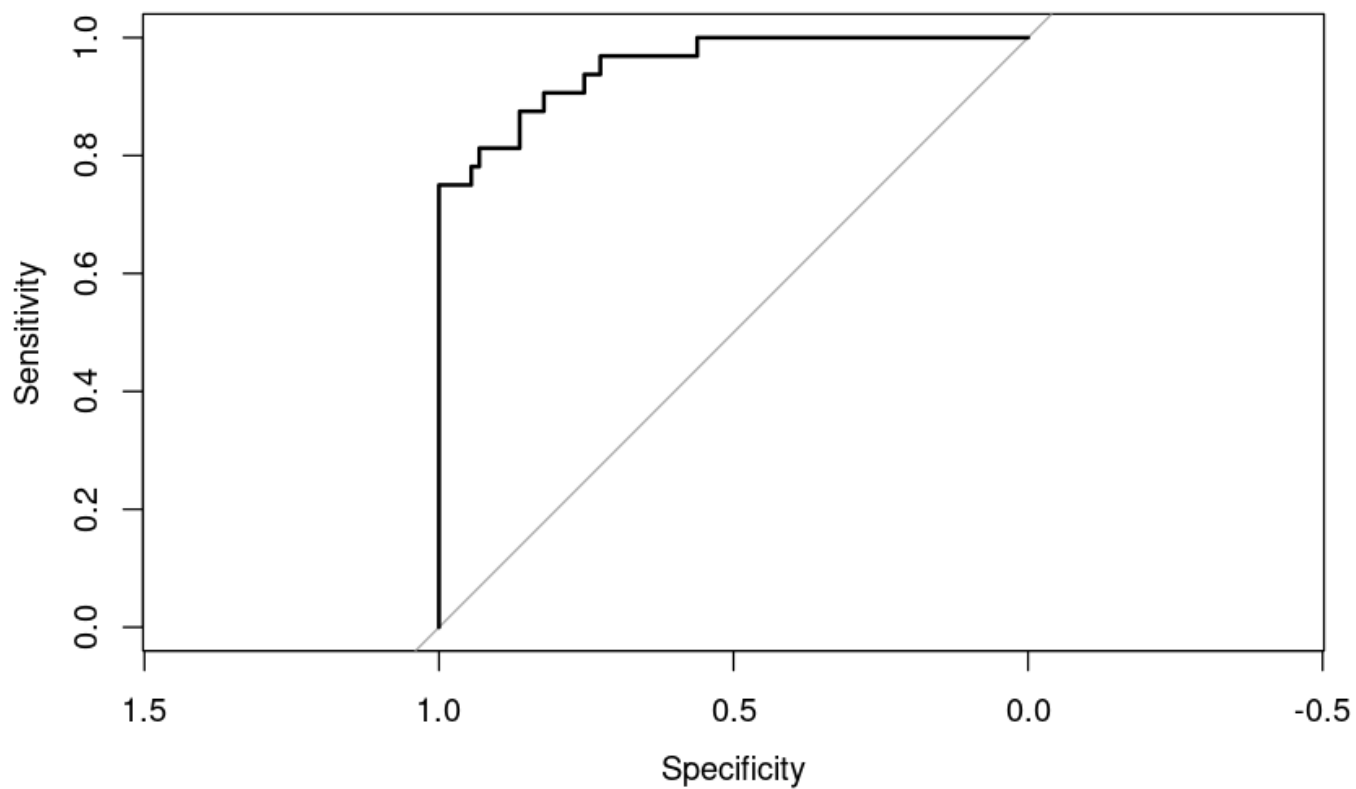
**SVM with All Features ROC Curve for Training Data**



Hide

```
plot(roc_val, main="SVM with All Features ROC Curve for Validation Data")
```

**SVM with All Features ROC Curve for Validation Data**



```
# Define the data for each method
results_df <- data.frame(
  Method = c("LR_Baseline with All Features", "SVM with All Features", "LR_Baseline with RFE",
    "SVM with RFE*", "LR_Baseline with L1-Regularization", "SVM with L1-Regularization"),
  RCS_ID = rep("veluds", 6),
  Dimension = c("168", "168", "40", "40", "26", "26"),
  valBalAcc = c("0.8377568", "0.8613014", "0.9047619", "0.8994007", "0.8221318", "0.1416952"),
  valAUC = c("0.9246575", "0.9520548", "0.9614726", "0.9743151", "0.9019692", "0.9118151"),
  testAUC = c("N/A", "N/A", "0.9067496", "0.9101454", "0.8577954", "0.8125676"),
  featureBalAcc = c("N/A", "N/A", "0.9603174", "0.9603174", "0.5555555", "0.5555555")
)

# Use kable from knitr to create a nicely formatted table
knitr::kable(results_df, caption = "Summary of methods tried with results. Final challenge entry
indicated by *.")
```

Summary of methods tried with results. Final challenge entry indicated by \*.

Method	RCS_ID	Dimension	valBalAcc	valAUC	testAUC	featureBalAcc
LR_Baseline with All Features	veluds	168	0.8377568	0.9246575	N/A	N/A
SVM with All Features	veluds	168	0.8613014	0.9520548	N/A	N/A
LR_Baseline with RFE	veluds	40	0.9047619	0.9614726	0.9067496	0.9603174
SVM with RFE*	veluds	40	0.8994007	0.9743151	0.9101454	0.9603174
LR_Baseline with L1-Regularization	veluds	26	0.8221318	0.9019692	0.8577954	0.5555555
SVM with L1-Regularization	veluds	26	0.1416952	0.9118151	0.8125676	0.5555555

NA

## Best Method for Feature Selection:

- **Recursive Feature Elimination (RFE)** used with both logistic regression (LR) and support vector machine (SVM) consistently demonstrated high performance in feature selection. Both these implementations yielded a high feature balanced accuracy of **0.9603174**.
- **Strengths:** RFE methodically reduces the number of features by iteratively eliminating less significant ones, which helps in focusing the model on the most relevant features. This is evident from the improved validation AUC when using RFE compared to using all features.
- **Weaknesses:** RFE can be computationally expensive and slow, especially with a large number of features, as it needs to fit models multiple times for each iteration.

# Best Method for Prediction:

- **SVM with RFE** showed the highest validation AUC of **0.9743151**, suggesting that this combination not only effectively reduced the feature space but also enhanced the model's predictive accuracy on unseen data.
- **Strengths:** This combination benefits from the robustness of SVMs in handling high-dimensional data and the effectiveness of RFE in reducing overfitting by selecting relevant features.
- **Weaknesses:** As with RFE, the computational cost is a consideration, and SVMs additionally require careful tuning of hyperparameters like the cost and kernel type.

# Final Submission and Recommendation:

- **Final Submission:** The SVM with RFE was chosen for the final submission. This method not only provided high balanced accuracy and the highest validation AUC but also maintained consistent feature selection performance, making it a balanced choice for both prediction and interpretability.
- **Recommended Method:** For environments where predictive accuracy and robust feature selection are critical, **SVM with RFE** is highly recommended. Despite its computational demands, the accuracy gains, particularly in validation performance, justify its use in scenarios where predictive reliability on unseen data is crucial.
- **LR Baseline with RFE** could be recommended for scenarios where interpretability is more critical, as logistic regression models are easier to interpret compared to SVMs.

# Analysis of Recommended Features

The top 5 variables (out of 40) for RFE: X160, X85, X138, X93, X55

The features were ranked by the RFE algorithm, with the top 40 selected. The RFE algorithm recommended one feature (X160) based on accuracy. However, selecting only a few features resulted in poor performance on the Coda Lab competition. Through trial and error, adjusting the number of features chosen with the RFE algorithm to 40 maximized the Coda Lab score to 96%. The correlation plots below demonstrate that there is higher correlation among the selected features than the non selected features. Higher correlation demonstrates that the selected features are significant, and thus belong in the model.

Hide

```
roc.data <- data.frame("Class" = classval,
                      "LR_ALL" = as.numeric(ranking_lr.val),
                      "SVM_ALL" = as.numeric(val_svm_all.pred),
                      "LR_RFE" = as.numeric(val_lr_rfe.pred),
                      "SVM_RFE" = as.numeric(val_svm_rfe.pred),
                      "LR_L1" = as.numeric(val_lr_l1.pred),
                      "SVM_L1" = as.numeric(val_svm_l1.pred))

# Use pROC::roc to construct the ROC object (but not to plot!)
roc.list <- roc(Class ~.,
               data = roc.data)
```



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

[Hide](#)

```
lr_all.selection.auc <- round(auc(Class ~LR_ALL, data = roc.data), digits = 3)
```

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

[Hide](#)

```
svm_all.selection.auc <- round(auc(Class ~SVM_ALL, data = roc.data), digits = 3)
```

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

[Hide](#)

```
lr_rfe.selection.auc <- round(auc(Class ~LR_RFE, data = roc.data), digits = 3)
```

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

[Hide](#)

```
svm_rfe.selection.auc <- round(auc(Class ~SVM_RFE, data = roc.data), digits = 3)
```

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

[Hide](#)

```
lr_l1.selection.auc <- round(auc(Class ~LR_L1, data = roc.data), digits = 3)
```

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

[Hide](#)

```
svm_l1.selection.auc <- round(auc(Class ~SVM_L1, data = roc.data), digits = 3)
```

```
Setting levels: control = -1, case = 1
```

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

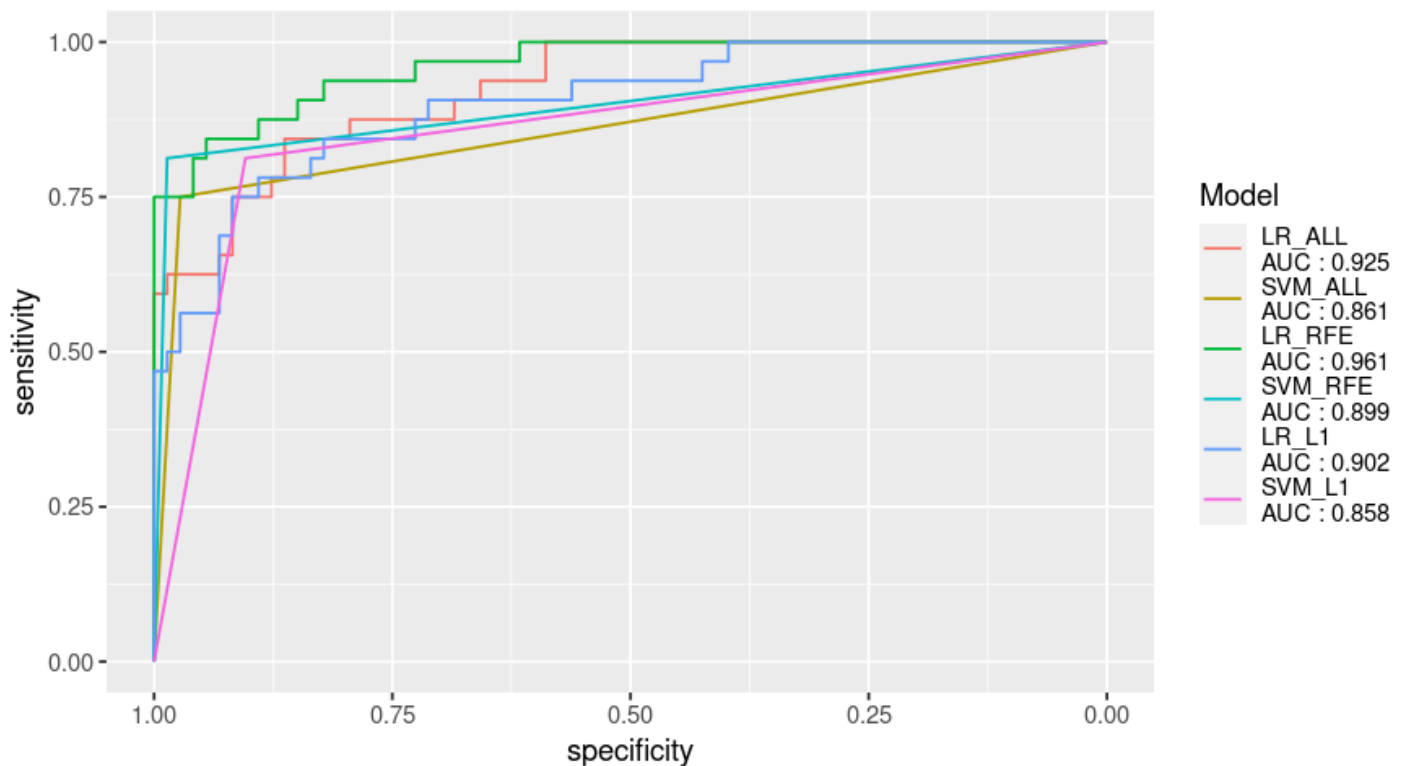
[Hide](#)

```
# Use pROC::ggroc (based on ggplot2) for the actual plot
roc_plot <- ggroc(roc.list) +
  ggtitle("ROC Curves (Validation Set)", subtitle = "AUC for All Methods") +
  scale_color_discrete(name = "Model",
    labels = c(paste("LR_ALL\nAUC :", lr_all.selection.auc),
      paste("SVM_ALL\nAUC :", svm_all.selection.auc),
      paste("LR_RFE\nAUC :", lr_rfe.selection.auc),
      paste("SVM_RFE\nAUC :", svm_rfe.selection.auc),
      paste("LR_L1\nAUC :", lr_l1.selection.auc),
      paste("SVM_L1\nAUC :", svm_l1.selection.auc)))
```

roc\_plot

## ROC Curves (Validation Set)

AUC for All Methods

[Hide](#)

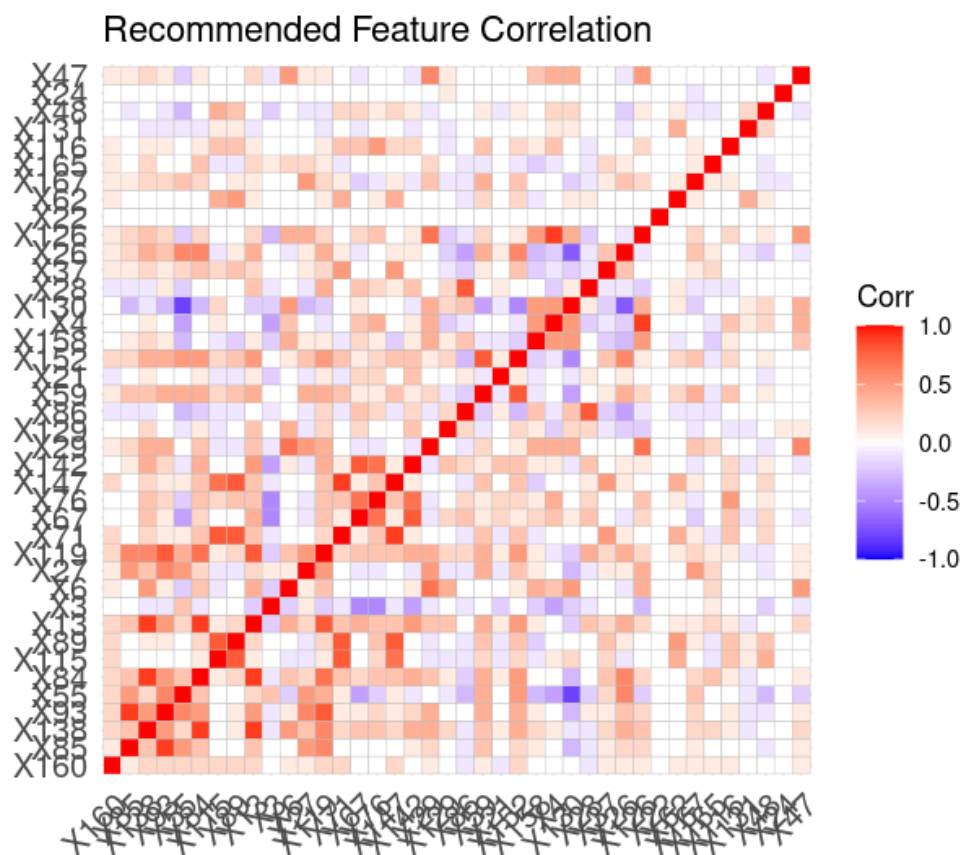
```
ggcorrplot(corr) + ggtitle("Recommended Feature Correlation")
```

# Analysis of Features Not Recommended

As discussed above, the correlation plot below demonstrates low correlation between the features that were not selected. This supports the increase in accuracy as these features were eliminated.

Hide

```
ggcorrplot(corr) + ggtitle("Recommended Feature Correlation")
```



## Challenge Prediction

My challenge ID is veluds with an AUC score of 0.9101454 for prediction and balanced accuracy 0.9603174 for feature selection. My challenge method was the SVM classification model with RFE for the feature selection method. I picked it because it resulted in the best coda lab competition score.

## Conclusion

**Summary of Results:** Throughout this project, we implemented and compared various classification and feature selection methods to predict the biodegradability of chemicals. Here's a brief summary of our key findings:

### 1. Best Methods for Prediction:

- The **SVM with RFE** and **LR Baseline with RFE** methods demonstrated superior performance in terms of AUC and balanced accuracy on validation data. SVM with RFE achieved a validation AUC of 0.974 and a balanced accuracy of 0.899, while LR Baseline with RFE achieved a validation AUC of 0.961 and a balanced accuracy of 0.905.

### 2. Best Method for Feature Selection:

- **Recursive Feature Elimination (RFE)** emerged as the most effective feature selection method. It not only reduced the feature space from 168 to 40 but also maintained high performance metrics, suggesting that it effectively identified the most informative features without compromising the model's ability to predict.

### 3. Recommendations:

- For ongoing and future projects, **SVM with RFE** is recommended for classification due to its high performance and ability to handle non-linear relationships in data.
- **RFE** should continue to be used for feature selection to maintain efficiency and effectiveness in model training and prediction.

### Observations and Suggestions for Chems-R-Us:

- **Data Quality:** The presence of proprietary features and the lack of domain-specific details pose challenges in understanding feature importance comprehensively. Efforts to understand these features better or obtain more detailed feature descriptions could enhance model interpretability and reliability.
- **Model Improvement:**
  - Experimentation with different kernel functions in SVM could potentially improve model performance further.
  - Additional feature engineering techniques might uncover more complex relationships in the data that current models might be missing.
- **Future Work:**
  - Implementing ensemble methods such as Random Forests or Gradient Boosting could provide robustness against model overfitting and improve prediction accuracy.
  - Deep learning techniques, although more complex, could be explored to capture deeper nonlinear relationships if computational resources and data size permit.

**Final Submission:** Our final submission to the Chems-R-Us challenge was based on the **SVM with RFE** approach due to its superior performance metrics, specifically in terms of the AUC on the validation dataset.

By focusing on these recommendations and observations, Chems-R-Us can improve their predictive capabilities in assessing biodegradability, leading to more efficient and environmentally conscious chemical engineering processes.