

Mingtao_Gao_HW6

Mingtao Gao
3/2/2020

```
# libraries that will be used for this homework
library(patchwork)
library(ggplot2)
library(e1071)
library(caret)
```

```
## Loading required package: lattice
```

```
## Loading required package: ggplot2
```

```
library(kernlab)
```

```
## Attaching package: 'kernlab'
```

```
## The following object is masked from 'package:ggplot2':
```

```
## alpha
```

Non-Linear Separation 1.

```
set.seed(1234)
# Create non-linear relation dataset
X <- rnorm(100)
Y <- 8 * X^3 + 4 * X^2 + rnorm(100) * 1

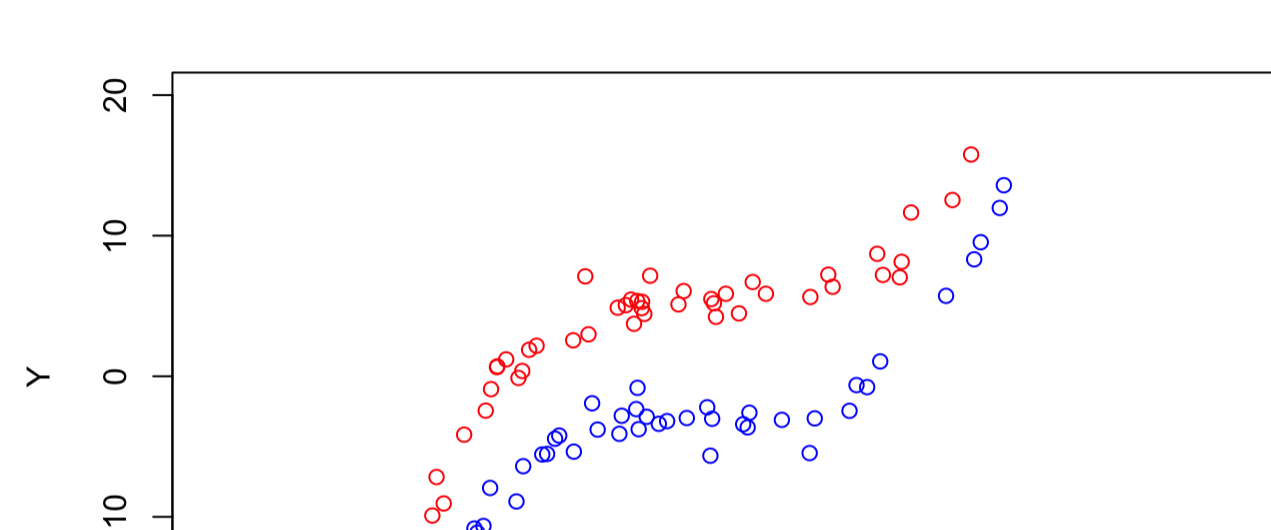
# Create separation
class = sample(1:100, 50)
Y[class] = Y[class] + 4
Y[-class] = Y[-class] - 4

# Add class data to the dataset
z = rep(-1, 100)
z[class] = 1
data = data.frame(x=X, y=Y, z=as.factor(z))

# Plot the data
plot(x[class], Y[class], col = "red", xlab = "X", ylab = "Y", ylim = c(-20, 20))
points(X[-class], Y[-class], col = "blue")

# Split dataset to training and testing
train = sample(100, 80)
data.train = data[data$train %in% train, ]
data.test = data[-train, ]

# Support Vector Classifier (Linear Kernel)
svm.linear <- svm(z ~., data=train,
  kernel="linear",
  scale=FALSE, cost=5)
plot(svm.linear, data=data.train)
```



```
# Support Vector Classifier (Linear Kernel) -- Training Error
table(predict=svm.linear, data=train, truth=data.train$z)
```

```
##      truth
## predict -1 1
##      -1 1 21 27
```

The support vector classifier with linear kernel made 35 errors on the training data, the error rate is 35/80 = 0.4375.

```
# Support Vector Classifier (Linear Kernel) -- Testing Error
table(predict=svm.linear, data=test, truth=data.test$z)
```

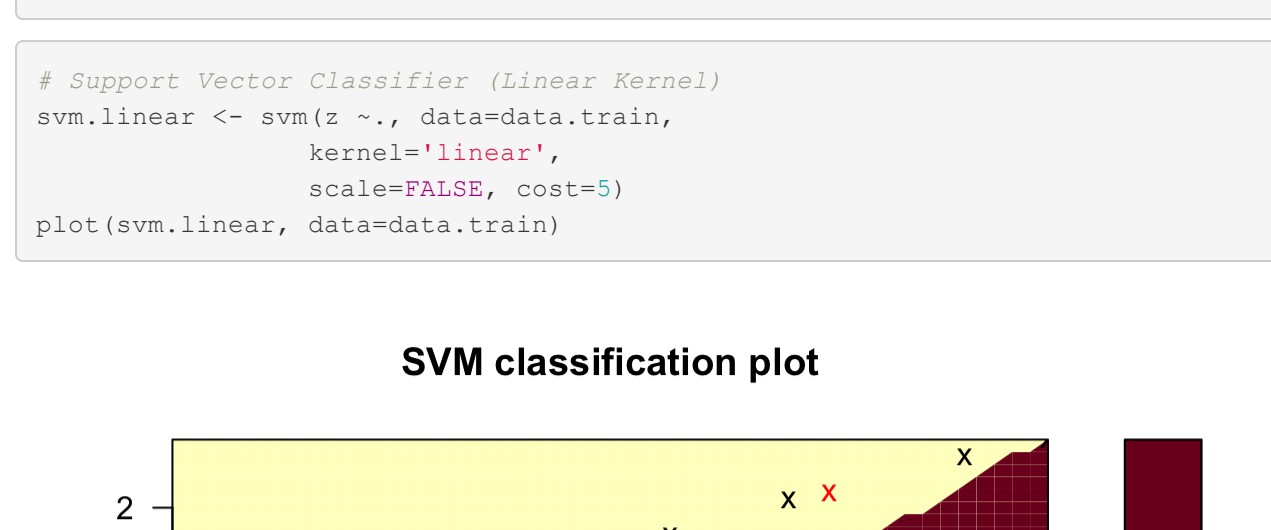
```
##      truth
## predict -1 1
##      -1 0 2
##      1 6 7
```

The support vector classifier with linear kernel made 6 errors on the training data, the error rate is 8/20 = 0.4.

The high error rate of SVC with linear kernel shows that it performed very badly with data that has a non-linear separation between two features. Next, let's move to SVC with a radial kernel and evaluate its performance.

```
# Support Vector Classifier (Radial Kernel)
svm.radial <- svm(z ~., data=train,
  kernel="radial",
  scale=FALSE, cost=5)
plot(svm.radial, data=train)
```

SVM classification plot



```
# Support Vector Classifier (Radial Kernel) -- Training Error
table(predict=svm.radial, data=train, truth=data.train$z)
```

```
##      truth
## predict -1 1
##      -1 36 1
##      1 6 40
```

The support vector classifier with radial kernel only made 4 errors on the training data, the error rate is 4/60 = 0.05.

```
# Support Vector Classifier (Radial Kernel) -- Testing Error
table(predict=svm.radial, data=test, truth=data.test$z)
```

```
##      truth
## predict -1 1
##      -1 8 3
##      1 3 6
```

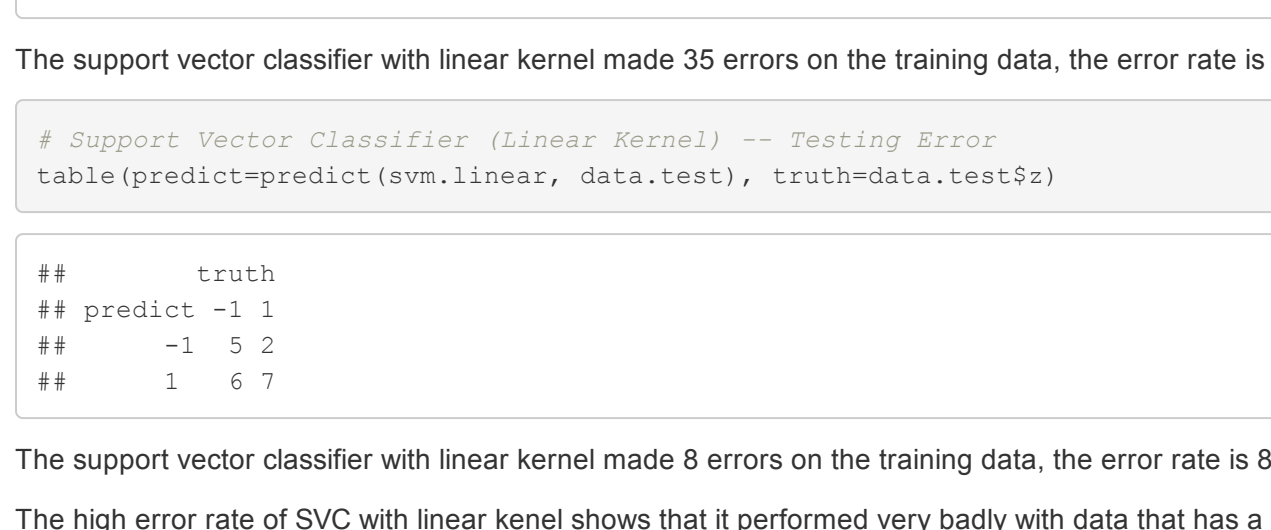
The support vector classifier with radial kernel made 6 errors on the training data, the error rate is 6/20 = 0.3.

We can see there is a good reduction in training and testing error rate with using radial kernel for support vector classifier when there is a non-linear separation between two classes. A support vector machine with a radial kernel indeed outperformed a support vector classifier on the training data in this setting.

SVM vs. Logistic Regression 2.

```
# Generate dataset
x1 <- runif(500) * 0.5
x2 <- runif(500) * 0.5
y <- ifelse(x1 < x2, 1, 0)
data = data.frame(x1=x1, x2=x2, y=y)

# Plot the observations
plot(x1[y==1], x2[y==1], col="blue", xlab="X1", ylab="X2")
points(x1[y==0], x2[y==0], col="red")
```



```
# Fit a logistic regression model with linear relation between x1 and x2
logre.linear <- glm(y ~ x1 + x2, family="binomial")
summary(logre.linear)
```

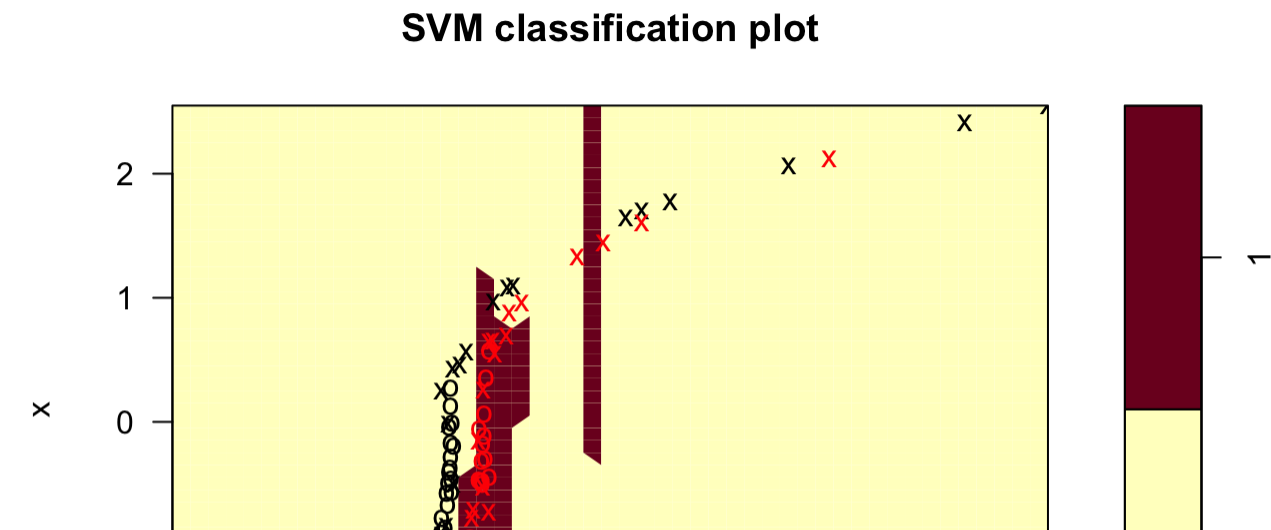
```
## Call:
## glm(formula = y ~ x1 + x2, family = "binomial")
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.296 -1.150 -1.040 -0.990 -0.922
```

```
## Coefficients:
## (Intercept)      Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -0.04575    0.03888   -1.178    0.241
## x1              -0.36555    0.30393   -1.203    0.229
## x2              -0.31404    0.31718   -0.990    0.322
```

```
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 692.95 on 499 degrees of freedom
## Residual deviance: 690.41 on 497 degrees of freedom
## AIC: 696.41
##
## Number of Fisher Scoring iterations: 3
```

3.

```
# Predict the data and plot the predictions
l1.pred = predict(logre.linear, newdata=data, type="response")
l1.pred = ifelse(l1.pred > 0.5, 1, 0)
plot(data[l1.pred == 1, $x1, data[l1.pred == 1, $x2],
  xlab = "X1", ylab = "X2", col = "blue")
points(data[l1.pred == 0, $x1, data[l1.pred == 0, $x2], col = "red")
```



We can observe the predicted

decision boundary looks like a linear relation.

6.

```
# Fit a logistic regression model with non-linear relation between x1 and x2
logre.nonlinear <- glm(y ~ poly(x1, 2) + poly(x2, 2) + I(x1 * x2), family="binomial")
summary(logre.nonlinear)
```

```
## Warning: glm.fit: algorithm did not converge
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
summary(logre.nonlinear)
```

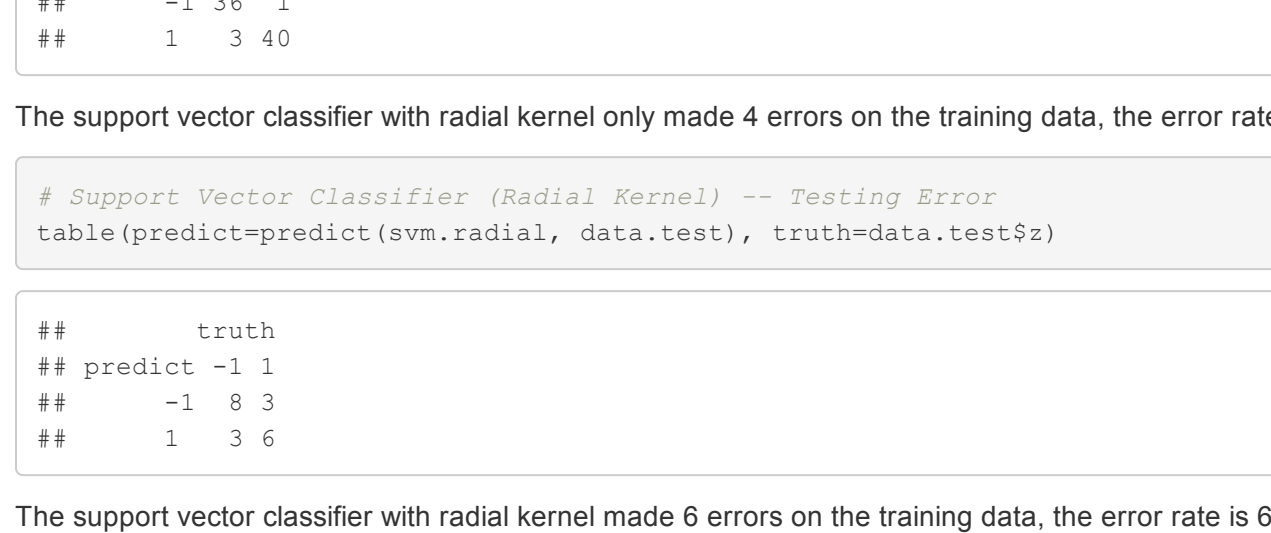
```
## Call:
## glm(formula = y ~ poly(x1, 2) + poly(x2, 2) + I(x1 * x2), family = "binomial")
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.399e-03 -2.000e-08 -2.000e-08  2.000e-08  1.294e-03
```

```
## Coefficients:
## (Intercept)      Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -0.04575    0.03888   -1.178    0.241
## x1              -0.36555    0.30393   -1.203    0.229
## x2              -0.31404    0.31718   -0.990    0.322
```

```
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 6.9295e+02 on 499 degrees of freedom
## Residual deviance: 4.3736e-06 on 494 degrees of freedom
## AIC: 12
##
## Number of Fisher Scoring iterations: 25
```

7.

```
# Predict the data and plot the predictions
nl1.pred = predict(logre.nonlinear, newdata=data, type="response")
nl1.pred = ifelse(nl1.pred > 0.5, 1, 0)
plot(data[nl1.pred == 1, $x1, data[nl1.pred == 1, $x2],
  xlab = "X1", ylab = "X2", col = "blue")
points(data[nl1.pred == 0, $x1, data[nl1.pred == 0, $x2], col = "red")
```

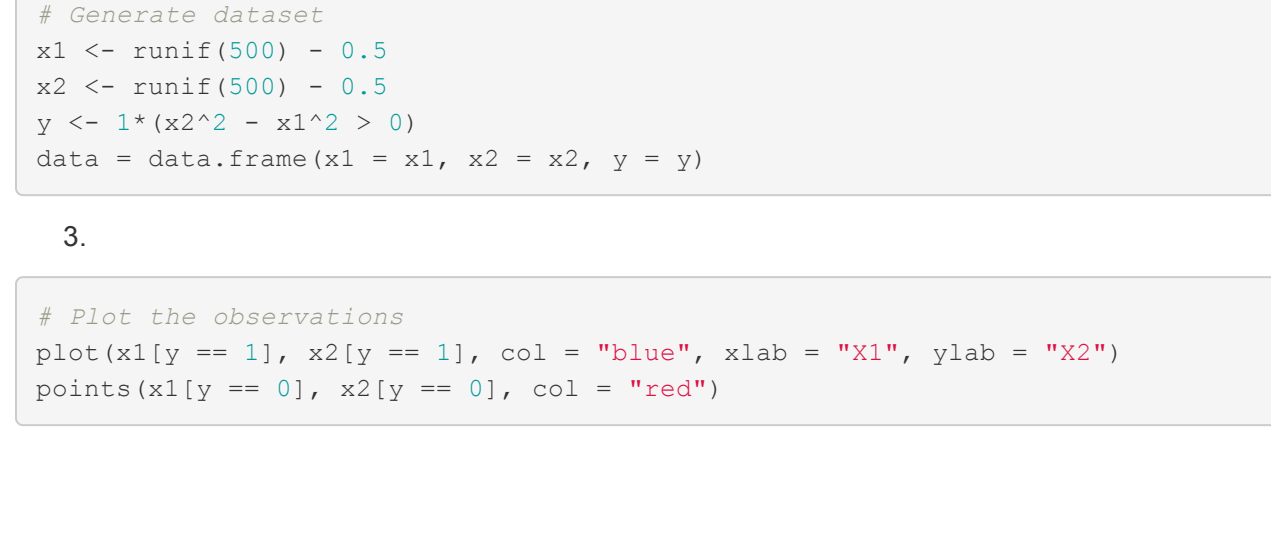


We can see the non-linear logistic

regression model performs much better than the linear one and its non-linear decision boundary is very close to the true decision boundary.

8.

```
# Fit Support Vector Classifier (Linear Kernel)
svm.linear <- svm(x.factor(data$y) ~ x1 + x2, data, kernel="linear", cost=0.01)
svm1.pred <- predict(svm.linear, data)
plot(data[svm1.pred == 0, $x1, data[svm1.pred == 0, $x2],
  xlab = "X1", ylab = "X2", col="red")
points(data[svm1.pred == 1, $x1, data[svm1.pred == 1, $x2], col="blue")
```

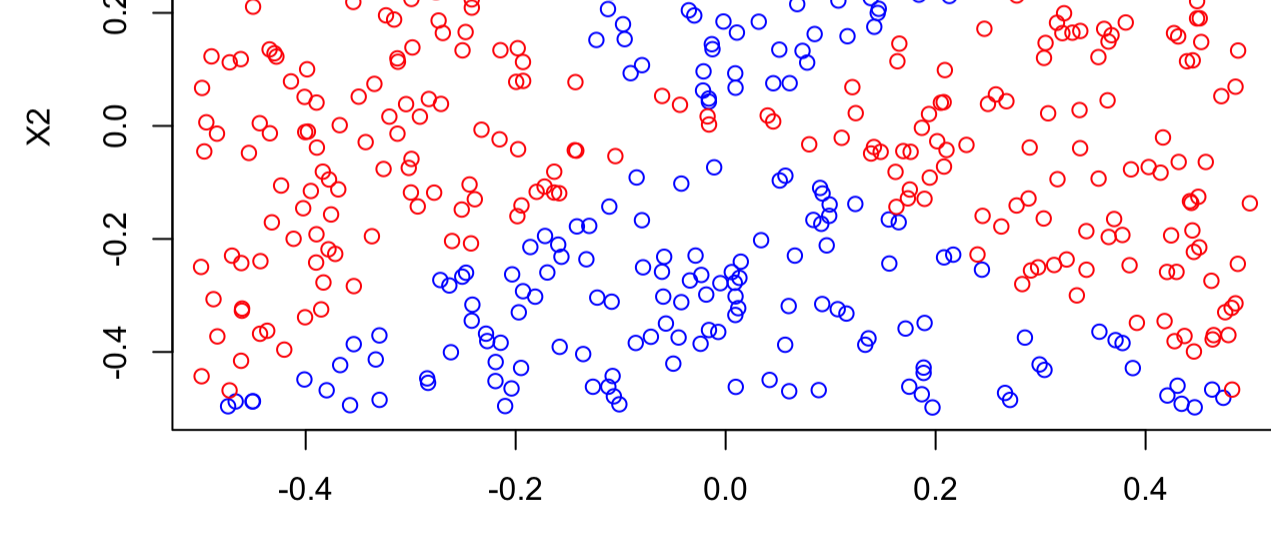


We can see SVC with linear kernel

performed very poorly. It classified all values into one single class.

9.

```
# Fit Support Vector Classifier (Non-linear/radial kernel)
svm.nonlinear <- svm(x.factor(data$y) ~ x1 + x2, data, kernel="radial", gamma=1)
svm1.pred <- predict(svm.nonlinear, data)
plot(data[svm1.pred == 0, $x1, data[svm1.pred == 0, $x2],
  xlab = "X1", ylab = "X2", col="red")
points(data[svm1.pred == 1, $x1, data[svm1.pred == 1, $x2], col="blue")
```



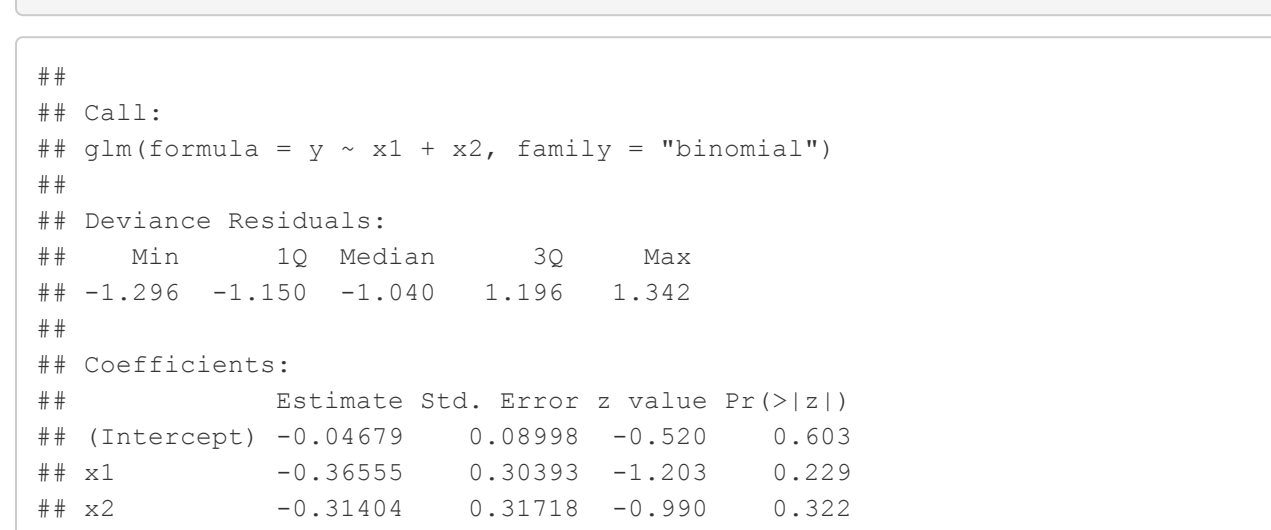
We can see the non-linear kernel

also outperform the linear kernel used for support vector classifier.

10. Results and Comments: From the above analysis on linear/non-linear logistic regression and linear/non-linear support vector classifier, we can find that, in this setting with a non-linear decision boundary, SVC with non-linear kernel and logistic regression with interaction terms both perform well on fitting the data. However, SVC with linear kernel and linear logistic regression performed very poorly when the decision boundary is non-linear. The tradeoff for SVM is that we need to test many interactive relations in order to find the best performed model. The tradeoff for SVC is that we also need to tune the gamma value in order to fit the best model.

Tuning Cost 11.

```
# Generate barely linear separated dataset
x1 <- runif(1000, -4, 4)
x2 <- runif(1000, -4, 4)
y <- rep(NA, 1000)
for (i in seq(1, 1000)) {
  if (x1[i] - x2[i] > 0.5) {
    y[i] <- 1
  } else if (x1[i] - x2[i] < -0.5) {
    y[i] <- 0
  } else {
    y[i] <- sample(c(0,1), replace=TRUE, size=1)
  }
}
data <- data.frame(x1 = x1, x2 = x2, y = as.factor(y))
plot(data$x1, data$x2, col=as.integer(data$y) + 1)
```



12.

```
# Split data into training and testing set
train = sample(100, 80)
data.train = data[train, ]
data.test = data[-train, ]
```

```
# Tuning cost
costs = c(0.01, 0.1, 1, 5, 10, 100, 1000, 10000)
tune.out <- tune(svm, y ~., data=train, kernel="linear",
  ranges=list(cost = costs))
summary(tune.out)
```

```
## Parameter tuning of 'svm':
## - sampling method: 10-fold cross validation
## - best parameters:
## cost gamma
## 10
## - best performance: 0.1
```

```
## - Detailed performance results:
## cost error dispersion
## 1 le=02 0.1025 0.14731391
## 2 le=01 0.1025 0.10465919
## 3 le=00 0.1250 0.11785113
## 4 le=00 0.1125 0.10944438
## 5 le=01 0.1000 0.09601133
## 6 le=02 0.1125 0.09223310
## 7 le=02 0.1125 0.09223310
## 8 le=04 0.1125 0.09223310
```

The training errors are displayed above. We can find when the cost value starts to increase, the CV errors starts to decrease and remain stable at a certain point. The best cost value based on CV error rates is 10.

```
# Compute the training errors with different cost values
train.errs <- rep(NA, length(costs))
for (cost in costs) {
  svm.fit <- svm(y ~., data=train, kernel="linear", cost=cost)
  res <- table(predict(svm.fit, newdata=data.train), truth=data.train$y)
  train.errs[match(cost, costs)] <- (res[2,1] + res[1,2]) / sum(res)
}
train.errs
```

```
## [1] 0.2250 0.1125 0.1000 0.1000 0.1000 0.1125 0.0875 0.0875
```

The testing errors show a similar trend as the training errors. As cost values increase, the testing errors also decrease first and then remain at a stable level. The best cost value based on CV error rates is 1.

13.

```
# Compute the testing errors with different cost values
test.errs <- rep(NA, length(costs))
for (cost in costs) {
  svm.fit <- svm(y ~., data=train, kernel="linear", cost=cost)
  res <- table(predict(svm.fit, newdata=data.test), truth=data.test$y)
  test.errs[match(cost, costs)] <- (res[2,1] + res[1,2]) / sum(res)
}
test.errs
```

```
## [1] 0.3195622 0.09239130 0.0695622 0.05217391 0.0543478 0.05434783 0.05000000
## [8] 0.05000000
```

The testing errors also show a similar trend as training and CV errors. As cost values increase, the testing errors also decrease first and then remain at a stable level. The best cost value based on CV error rates is 1.

14. When we have dataset that the two classes are just barely linearly separable, a small cost creates a large margin and too small would lead to higher training/testing CV error rates and poor performance. In general, a small cost creates a large margin and allows more misclassifications, with higher training errors, while a large cost creates a narrow margin and permits fewer misclassifications, with lower training errors. Thus, after reaching a certain level, the increase in cost will not lead to any changes in error rates for testing/testing data. For CV, higher, and testing error rates across all cost values, when cost is equal to 1, the error rate is minimized.

Predicting attitudes towards racist college professors 15.

```
# Load the dataset
gss.test <- read.csv("data/gss_test.csv")
gss.train <- read.csv("data/gss_train.csv")
```

```
# SVC with linear kernel
tune.out <- tune(svm, colrac ~., data=gss.train, kernel="linear",
  ranges=list(cost=c(0.01, 0.1, 1, 5, 10, 100)))
summary(tune.out)
```

```
## Parameter tuning of 'svm':
## - sampling method: 10-fold cross validation
## - best parameters:
## cost gamma
## 1 0.01
## - best performance: 0.1829203
```

```
## - Detailed performance results:
## cost error dispersion
## 1 le=02 0.1829203 0.03724689
## 2 le=01 0.2003460 0.03748506
## 3 le=00 0.2017173 0.03974464
## 4 le=00 0.2018641 0.03959338
## 5 le=01 0.2019736 0.03987396
## 6 le=02 0.2056114 0.03919998
```

Based on the CV error rates reported above, when cost is equal to 0.01, the SVC with linear kernel performs the best with the smallest CV error rate. As the cost value increases, the CV error rates started to increase.

16.

```
# SVC with polynomial kernel
tune.out <- tune(svm, colrac ~., data=gss.train, kernel="polynomial",
  ranges=list(cost=c(0.01, 0.1, 1, 5, 10, 100),
  degree=c(2, 3, 4)))
summary(tune.out)
```

```
## Parameter tuning of 'svm':
## - sampling method: 10-fold cross validation
## - best parameters:
## cost degree
## 1 3
## - best performance: 0.1510476
```

```
## - Detailed performance results:
## cost degree error dispersion
## 1 le=02 2 0.4069942 0.030528692
## 2 le=01 2 0.2642998 0.02893764
## 3 le=00 2 0.2157774 0.024515812
## 4 le=00 2 0.2844111 0.023608391
## 5 le=01 2 0.2123266 0.021914944
## 6 le=02 2 0.2173320 0.02741129
## 7 le=02 3 0.3704906 0.030983154
## 8 le=02 3 0.1794033 0.027655565
## 9 le=00 3 0.1510476 0.016197828
## 10 le=00 3 0.1691927 0.016197828
## 11 le=01 3 0.1691927 0.016197828
## 12 le=02 3 0.1691927 0.016197828
## 13 le=02 3 0.1691927 0.016197828
## 14 le=02 3 0.1691927 0.016197828
## 15 le=02 3 0.1691927 0.016197828
## 16 le=02 3 0.1691927 0.016197828
## 17 le=02 3 0.1691927 0.016197828
## 18 le=02 3 0.1691927 0.016197828
```

Based on the CV error rates reported above, when cost is equal to 1 and degree is equal to 3, the SVC with polynomial kernel performs the best with the smallest CV error rate.

```
# Prediction - SVC (linear)
svm.linear <- svm(colrac ~., data=gss.train, kernel="linear", cost=0.01)
table(predict=svm.linear, data=gss.test, colrac)
truths.test$colrac
```

```
##      truth
## predict 0 1
##      0 117 27
##      1 111 238
```

```
# Prediction - SVC (polynomial)
svm.poly <- svm(colrac ~., data=gss.train, kernel="polynomial",
  cost=1, degree=3)
table(predict=svm.poly, gss.test) > 0.5, 1, 0,
truths.test$colrac
```

```
##      truth
## predict 0 1
##      0 152
##      1 79 220
```

```
# Prediction - SVC (radial)
svm.radial <- svm(colrac ~., data=gss.train, kernel="radial",
  cost=gamma)
table(predict=svm.radial, gss.test) > 0.5, 1, 0,
truths.test$colrac
```

```
##      truth
## predict 0 1
##      0 154 34
##      1 74 231
```

We found that among all three models, the radial kernel SVM with tuned cost and gamma value outperformed the others, which shows a non-linear decision boundary between the two classes.