

Lab 8: Support Vector Machines

```
library(tidyverse) ## data manipulation
library(tidymodels) ## models
library(knitr) ## tables

## reproducible
set.seed(445)
```

0.1 Data Preparation

We will make some simulated data to see how support vector classifiers and SVMs work.

Run the following code to create two datasets.

```
n1 <- 20
n2 <- 200
p <- 2

## training data sets
x_small <- matrix(rnorm(n1 * p), ncol = p)
x_large <- matrix(rnorm(n2 * p), ncol = p)
y_small <- c(rep(-1, n1/2), rep(1, n1/2))
y_large <- c(rep(1, n2/4*3), rep(2, n2/4))

## shift data farther apart
x_small[y_small == 1,] <- x_small[y_small == 1,] + 1
x_large[1:100,] <- x_large[1:100,] + 2
x_large[101:150,] <- x_large[101:150,] - 2

## put data into dataframes
df_small <- data.frame(x_small, y = as.factor(y_small))
df_large <- data.frame(x_large, y = as.factor(y_large))
```

1. Make two scatterplots to inspect the small and large training data sets. Describe what you see.

0.2 Support Vector Classifier

We will use the `svm_poly` and `svm_rbf` functions to fit the support vector classifier and

the SVM.

Here is an example model specification for a support vector classifier (linear decision boundary):

```
svm_linear_spec <- svm_poly(degree = 1) %>%  
  set_mode("classification") %>%  
  set_engine("kernlab", scaled = FALSE)
```

The `cost` argument allows us to specify the cost of violation to the margin. When the `cost` argument is small, margins will be wide. An example of fitting an SVM with a specified cost is:

```
svm_linear_fit <- svm_linear_spec %>%  
  set_args(cost = 10) %>%  
  fit(y ~ ., data = df)
```

1. Fit a support vector classifier on the small data with $C = 10$ (use `scaled = FALSE` to indicate your data should not be scaled.)
2. How many support vectors were used to fit your classifier?
3. Predict a grid of \mathbf{X} values between the range of X_1 and X_2 . Plot these predictions using `geom_tile()` to visualize the decision boundary and add a scatterplot of training data on top, colored by training label. Describe what you see.
4. Alter your plot from 2 to change the shape of the support vectors. [*Hint*: You can extract the fit engine from your object using `extract_fit_engine` and then access the index of the support vectors from your stored object as `object@alphaindex`]
5. Perform CV on the cost parameter. Which value of C would you choose?
6. Repeat 3. and 4. using your chosen C value. Describe what you see.

0.3 Support Vector Machines

1. Split the large data frame into 50% training and 50% test.
2. Fit a linear SVM, radial SVM with $\gamma = 1$, and polynomial SVM with $d = 3$ using CV to choose the appropriate cost for each model.
3. Predict a grid of \mathbf{X} values between the range of X_1 and X_2 . Plot these predictions using `geom_tile()` to visualize the decision boundary and add a scatterplot of

training data on top, colored by training label. Describe what you see.

4. Predict your test data with your three models. Which model would you choose?