

HW 3

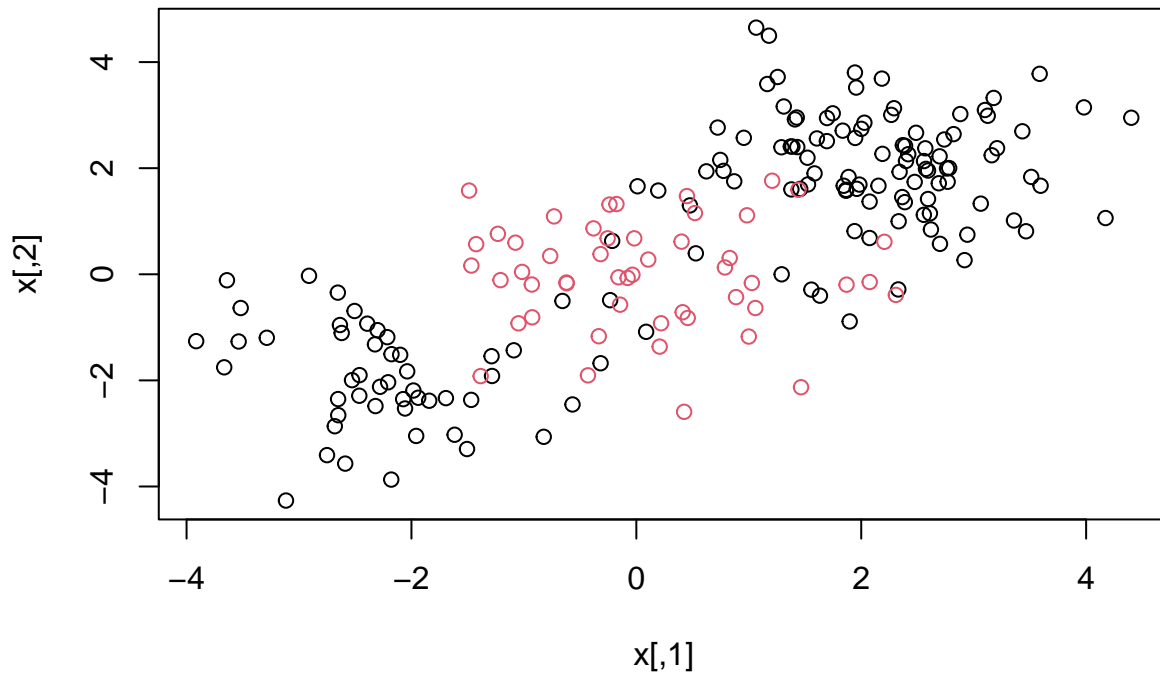
Marin Mato

9/24/2024

Let $E[X] = \mu$. Show that $Var[X] := E[(X - E[X])^2] = E[X^2] - (E[X])^2$. Note, all you have to do is show the second equality (the first is our definition from class).

In the computational section of this homework, we will discuss support vector machines and tree-based methods. I will begin by simulating some data for you to use with SVM.

```
library(e1071)
set.seed(1)
x=matrix(rnorm(200*2),ncol=2)
x[1:100,]=x[1:100,]+2
x[101:150,]=x[101:150,]-2
y=c(rep(1,150),rep(2,50))
dat=data.frame(x=x,y=as.factor(y))
plot(x, col=y)
```



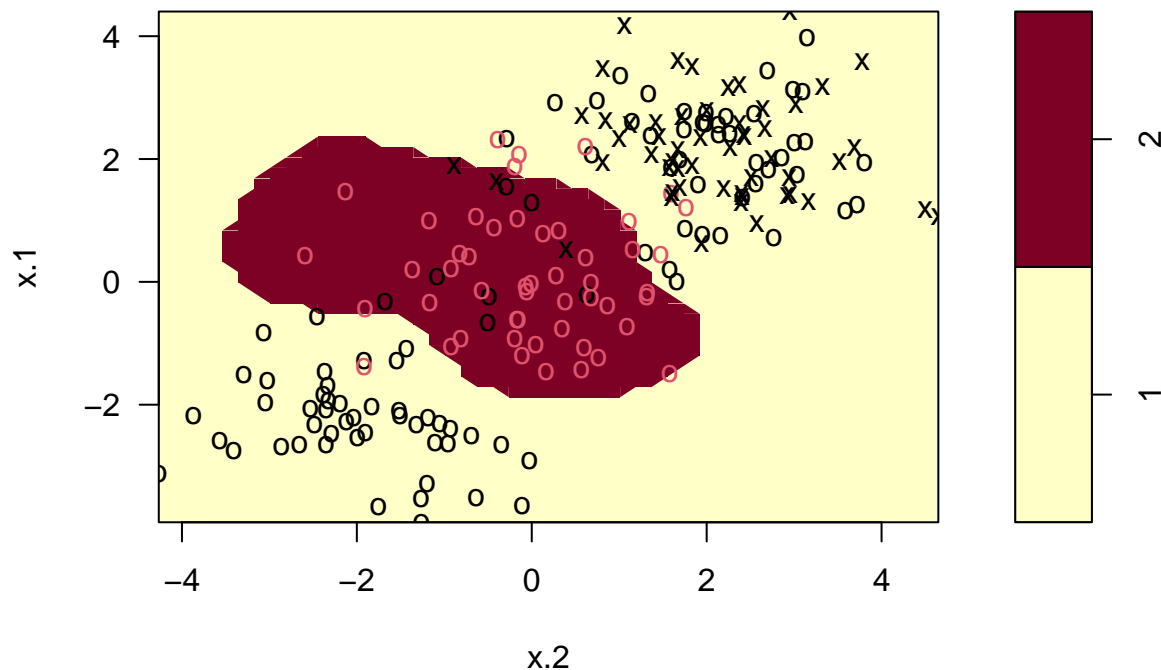
Quite clearly, the above data is not linearly separable. Create a training-testing partition with 100 random observations in the training partition. Fit an svm on this training data using the radial kernel, and tuning parameters $\gamma = 1$, cost = 1. Plot the svm on the training data.

```
set.seed(1)

train = sample(1:nrow(dat), 100)
training_dataset = dat[train,]
test_dataset = dat[-train,]

#now we plot svm
svm_fit = svm(y ~ ., data = training_dataset, kernel = "radial", cost = 1, scale = FALSE, gamma = 1)
plot(svm_fit, dat)
```

SVM classification plot

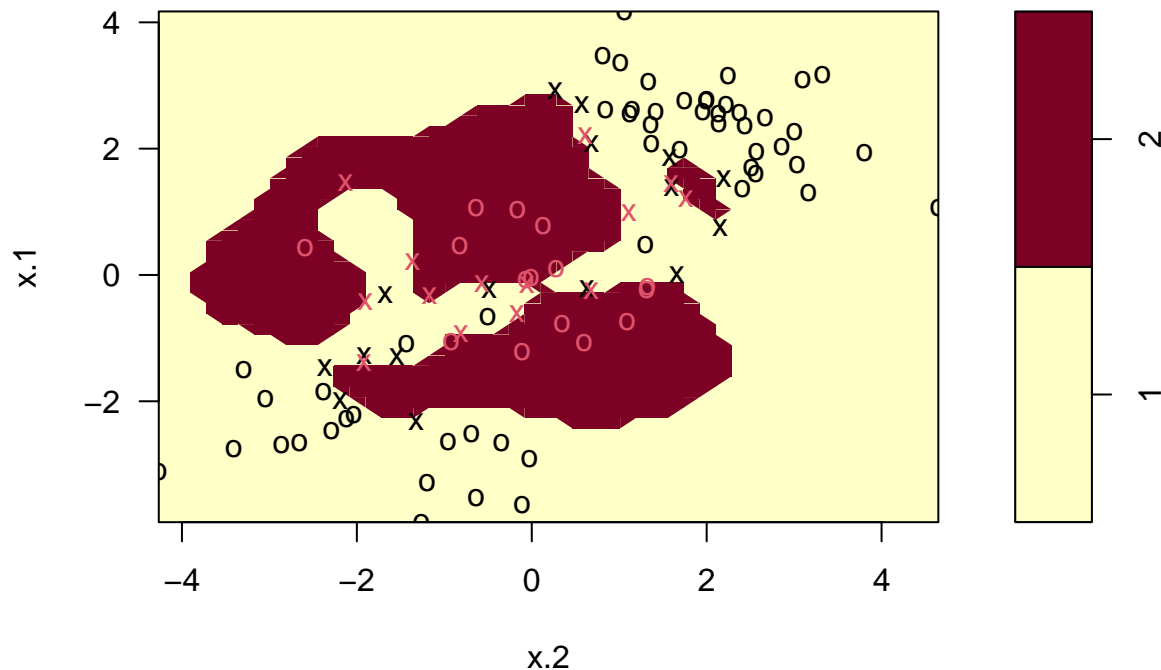


Notice that the above decision boundary is decidedly non-linear. It seems to perform reasonably well, but there are indeed some misclassifications. Let's see if increasing the cost ¹ helps our classification error rate. Refit the svm with the radial kernel, $\gamma = 1$, and a cost of 10000. Plot this svm on the training data.

```
svmfit = svm(y ~ ., data = training_dataset, kernel = "radial", cost = 10000, scale = FALSE, gamma = 1)
plot(svmfit, training_dataset)
```

¹Remember this is a parameter that decides how smooth your decision boundary should be

SVM classification plot



It would appear that we are better capturing the training data, but comment on the dangers (if any exist), of such a model.

By looking at the model closely, we can easily determine that it is overfitted, meaning that the testing data will lead to misleading results.

Create a confusion matrix by using this svm to predict on the current testing partition. Comment on the confusion matrix. Is there any disparity in our classification results?

```
#remove eval = FALSE in above  
table(true=dat[-train,"y"], pred=predict(svmfit, newdata=dat[-train,]))
```

```
##      pred  
## true  1  2  
##      1 62 17  
##      2  3 18
```

Is this disparity because of imbalance in the training/testing partition? Find the proportion of class 2 in your training partition and see if it is broadly representative of the underlying 25% of class 2 in the data as a whole.

```
prop <- nrow(training_dataset[y=="2",])/nrow(training_dataset)
```

```
prop
```

```
## [1] 0.5
```

Student Response

Let's try and balance the above to solutions via cross-validation. Using the `tune` function, pass in the training data, and a list of the following cost and γ values: $\{0.1, 1, 10, 100, 1000\}$ and $\{0.5, 1, 2, 3, 4\}$. Save the output of this function in a variable called `tune.out`.

```
set.seed(1)

tune.out = tune(svm, y ~ .,
               data = training_dataset,
               ranges = list(cost = c(0.1, 1, 10, 100, 1000),
                             gamma = c(.5, 1, 2, 3, 4)))
```

I will take `tune.out` and use the best model according to error rate to test on our data. I will report a confusion matrix corresponding to the 100 predictions.

```
table(true=dat[-train,"y"], pred=predict(tune.out$best.model, newdata=dat[-train,]))
```

Comment on the confusion matrix. How have we improved upon the model in question 2 and what qualifications are still necessary for this improved model.

Student Response

Let's turn now to decision trees.

```
library(kmed)
data(heart)
library(tree)
```

The response variable is currently a categorical variable with four levels. Convert heart disease into binary categorical variable. Then, ensure that it is properly stored as a factor.

Train a classification tree on a 240 observation training subset (using the seed I have set for you). Plot the tree.

```
set.seed(101)
```

Use the trained model to classify the remaining testing points. Create a confusion matrix to evaluate performance. Report the classification error rate.

Above we have a fully grown (bushy) tree. Now, cross validate it using the `cv.tree` command. Specify cross validation to be done according to the misclassification rate. Choose an ideal number of splits, and

plot this tree. Finally, use this pruned tree to test on the testing set. Report a confusion matrix and the misclassification rate.

```
set.seed(101)
```

Discuss the trade-off in accuracy and interpretability in pruning the above tree.

Student Input

Discuss the ways a decision tree could manifest algorithmic bias.

Student Answer