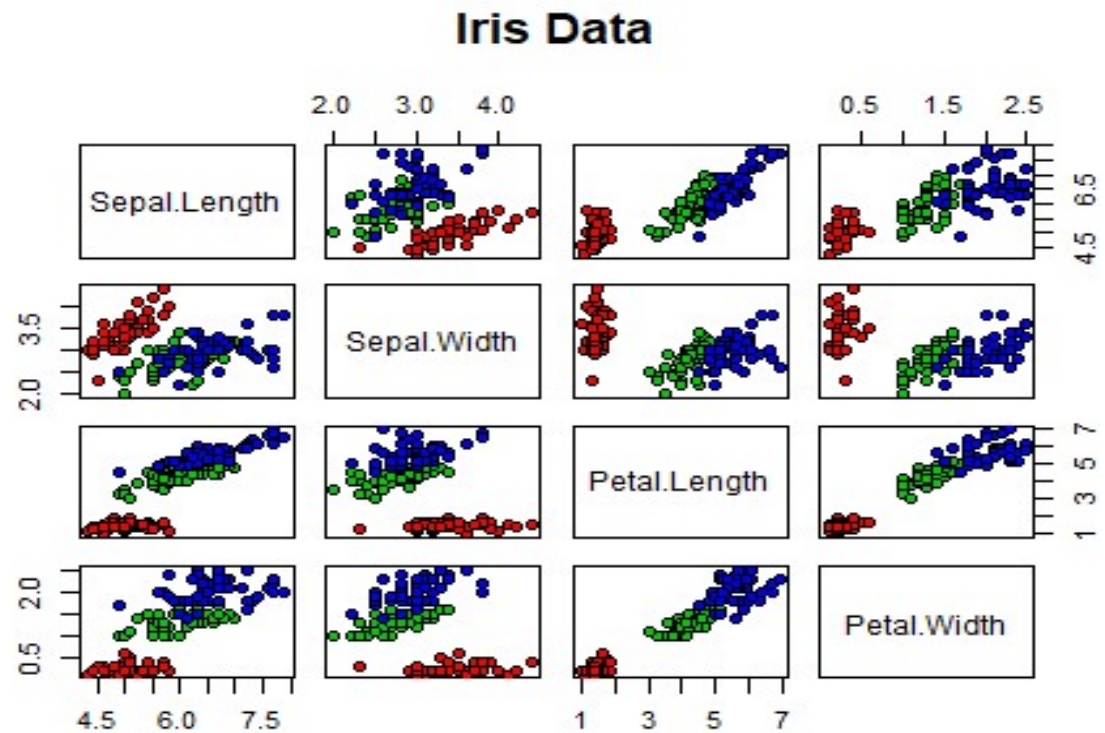


IE8990: Adv. Data Analytics for Complex Systems

- Lab 4: Classification
 - LDA & QDA
 - Logistic regression
 - SVM

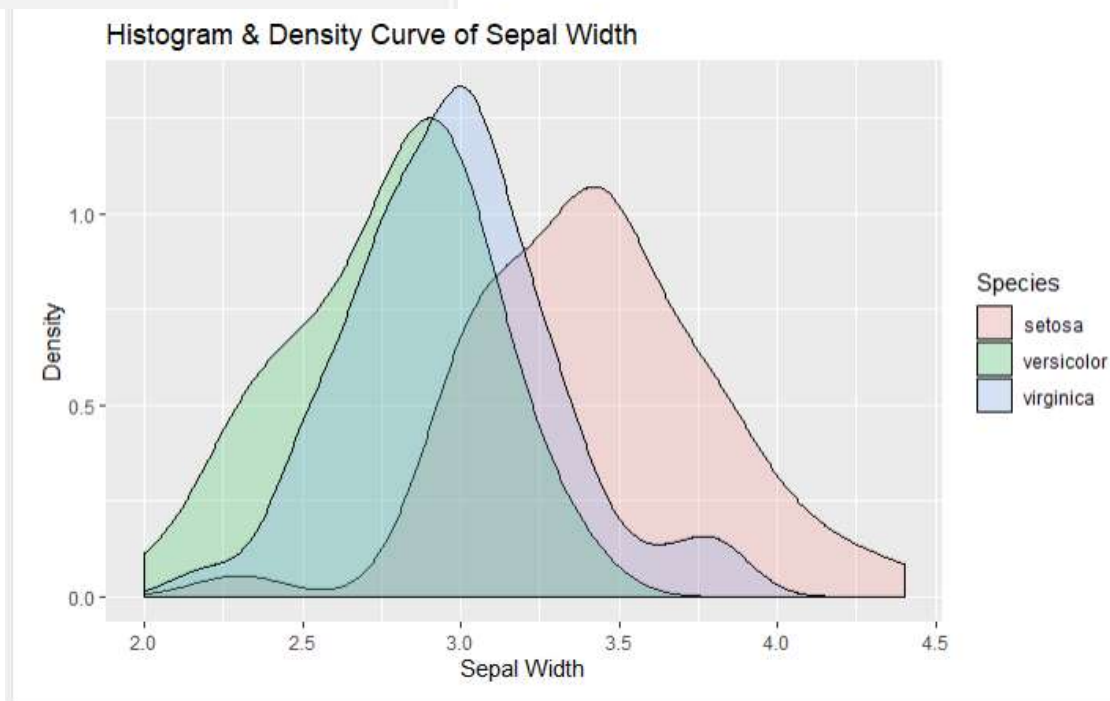
Iris Example 3: Classification

This time we try to classify the iris samples into different species using the given four inputs.



Explore functions in ggplot2 for data visualization

```
46 # try to use functions in ggplot2 to visualize data distribution
47 density2 <- ggplot(data=iris, aes(x=Sepal.width, fill=Species))|
48 density2 + geom_density(stat="density", alpha=I(0.2)) +
49   xlab("Sepal width") + ylab("Density") + ggtitle("Histogram &
   Density Curve of Sepal width")
```



LDA Model

```
50 ▾ ##{r}
51 # fit lda model
52 model.lda<- lda(Species~.,data=train.data)
53 print(model.lda)
54 ##
55
56 ▾ #### 3.2 Use the testing set to evaluate the classification
   performance
57 ▾ ##{r}
58 predictions.lda = data.frame(predict(model.lda, test.data))
59 confusionM.lda<-confusionMatrix(predictions.lda$class,test.data$
   Species)
60 print(confusionM.lda$table)
61 ##
```

Prediction	Reference		
	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	1
virginica	0	0	9



QDA Model

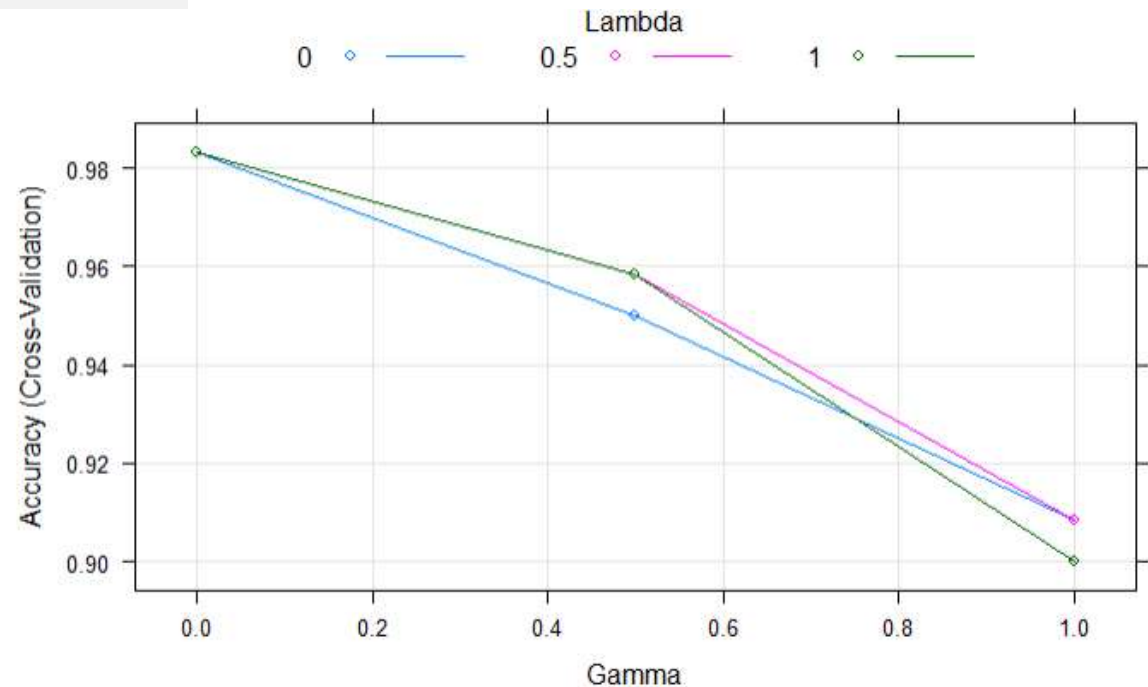
```
63 ▾ #### 3.3 Fit a QDA model and use testing set to evaluate  
    classification performance  
64 ▾ ```{r}  
65 model.qda<- qda(Species~.,data=train.data)  
66 print(model.qda)  
67 predictions.qda = data.frame(predict(model.qda, test.data))  
68 confusionM.qda<-confusionMatrix(predictions.qda$class,test.data$  
    Species)  
69 print(confusionM.qda$table)  
70 ```
```

	Reference		
Prediction	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	1
virginica	0	0	9



Regularized Discriminant Analysis

```
71- #### 3.4 Regularized discriminant analysis
72- ```{r}
73- library(klar)
74- cv_5_grid = trainControl(method = "cv", number = 5)
75- fit_rda_grid = train(Species ~ ., data = train.data,
76-   method = "rda", trcontrol = cv_5_grid)
77- fit_rda_grid
78- plot(fit_rda_grid)
79- ```
```



Logistic Regression

```
80  
81- ### 4. Fit a logistic regression model to predict iris  
   species  
82- #### 4.1 Use the same training set in 3.1 to fit a  
   logistic regression model  
83- ```{r}  
84 library(nnet)  
85 # Fit the model  
86 model.logReg <- nnet::multinom(species ~., data =  
   train.data)  
87 # Summarize the model  
88 summary(model.logReg)  
89 # Make predictions  
90 predicted.logReg <- model.logReg %>% predict(test.data)  
91 head(predicted.logReg)  
92 # Model accuracy  
93 mean(predicted.logReg == test.data$Species)  
94  
95 confusionM.logReg<-confusionMatrix(predicted.logReg,  
   test.data$Species)  
96
```

call:

```
nnet::multinom(formula = Species ~ ., data = train.data)
```

Coefficients:

	(Intercept)	Sepal.Length	Sepal.width	Petal.Length	Petal.width
versicolor	44.88584	-50.14588	-125.0783	191.7161	76.71824
virginica	-132.54466	-62.24794	-142.1385	224.9369	155.84279

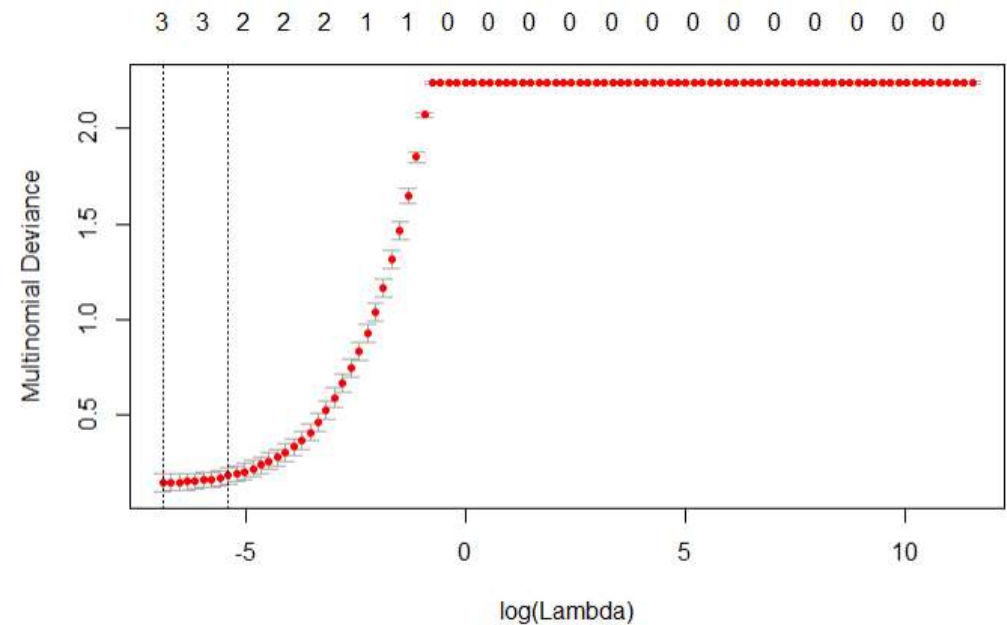
Do you remember how to interpret your coefficients in logistics regression?



MISSISSIPPI STATE
UNIVERSITY™

Regularized Logistic Regression

```
97 ▾ #### 4.2 Regularized logistic regression model
98
99 ▾ ```{r}
100 library(glmnet)
101 library(fastDummies)
102 x<-train.data[,1:4] %>% as.matrix()
103
104 y_dummy<-dummy_columns(y)
105 y_dummy1<-y_dummy[,2:4] %>% as.matrix()
106
107 ▾ # Perform 10-fold cross-validation to select lambda
-----
108 lambdas_to_try <- 10^seq(-3, 3, length.out = 100)
109 # Setting alpha = 1 implements lasso regression
110 lasso_cv <- cv.glmnet(x, y_dummy1, family="multinomial",
111                      alpha=1, lambda = lambdas_to_try,
112                      standardize = TRUE, nfolds = 10)
112 plot(lasso_cv)|
113 #lasso_model<-glmnet(x, y_dummy1, family="multinomial",
114                      alpha = 1, standardize = TRUE)
114
115 #plot(lasso_model,xvar = "lambda")
116 #legend("bottomright", lwd = 1, col = 1:6, legend =
117       colnames(x), cex = .7)
117 ...
```



HW 2: Q2

- Use data set [kyphosis](#) from R package “gam”, randomly divide the data into training (75%) and testing (25%) set.
 - 2.1 Fit a regularized discriminant analysis model, find the best tuning parameter combination, and evaluate the model using the testing set and report the confusion matrix
 - 2.2 Fit a logistic regression model, and interpret the coefficients you estimated
 - 2.3 Fit a regularized logistic regression model, and find the best model from your regularization

Logistic Regression – Model Interpretation

```
80
81- ### 4. Fit a logistic regression model to predict iris
   species
82- #### 4.1 Use the same training set in 3.1 to fit a
   logistic regression model
83- ```{r}
84 library(nnet)
85 # Fit the model
86 model.logReg <- nnet::multinom(species ~., data =
   train.data)
87 # Summarize the model
88 summary(model.logReg)
89 # Make predictions
90 predicted.logReg <- model.logReg %>% predict(test.data)
91 head(predicted.logReg)
92 # Model accuracy
93 mean(predicted.logReg == test.data$Species)
94
95 confusionM.logReg<-confusionMatrix(predicted.logReg,
   test.data$Species)
96
```

call:

```
nnet::multinom(formula = Species ~ ., data = train.data)
```

Coefficients:

	(Intercept)	Sepal.Length	Sepal.width	Petal.Length	Petal.width
versicolor	44.88584	-50.14588	-125.0783	191.7161	76.71824
virginica	-132.54466	-62.24794	-142.1385	224.9369	155.84279

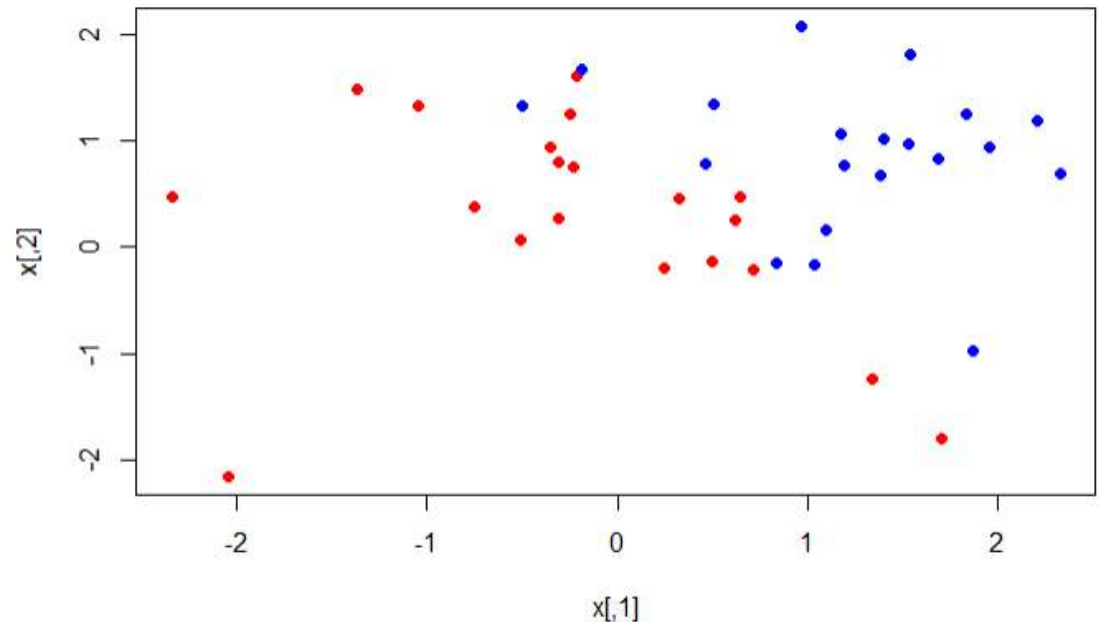


MISSISSIPPI STATE
UNIVERSITY™

Linear SVM Example

- Random data generation – will be involved in the exam

```
##{r}  
set.seed(10111)  
x = matrix(rnorm(80), 40, 2)  
y = rep(c(-1, 1), c(20, 20))  
x[y == 1,] = x[y == 1,] + 1  
plot(x, col = y + 3, pch = 19)  
  
set.seed(123)  
xt = matrix(rnorm(80), 40, 2)  
yt = rep(c(-1, 1), c(20, 20))  
xt[yt == 1,] = xt[yt == 1,] + 1
```



Linear SVM Example

```
31 - ```{r}
32 library(e1071)
33 library(caret)
34 library(dplyr)
35 library(tidyverse)
36 dat = data.frame(x, y = as.factor(y))
37 dat.test=data.frame(xt, y = as.factor(yt))
38 |
39 svmfit = svm(y ~ ., data = dat, kernel = "linear", cost = 10, scale = FALSE)
40 print(svmfit)
41 dat.test.svm<-predict(svmfit,data=dat.test)
42 ConfusionM.svm<-confusionMatrix(dat.test.svm,dat.test$y)
43 print(ConfusionM.svm)
44 ```
```

	Reference	
Prediction	-1	1
-1	19	2
1	1	18

- kernel = “linear”
- “cost” is the C in svm formulation

Sensitivity : 0.9500
Specificity : 0.9000
Pos Pred Value : 0.9048
Neg Pred Value : 0.9474
Prevalence : 0.5000
Detection Rate : 0.4750
Detection Prevalence : 0.5250
Balanced Accuracy : 0.9250

'Positive' class : -1

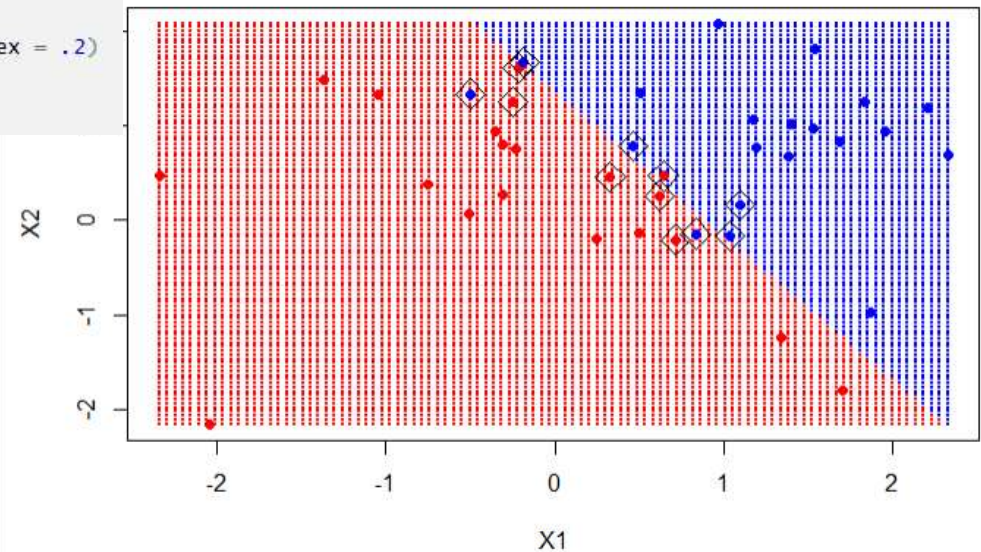


MISSISSIPPI STATE
UNIVERSITY™

Linear SVM Example

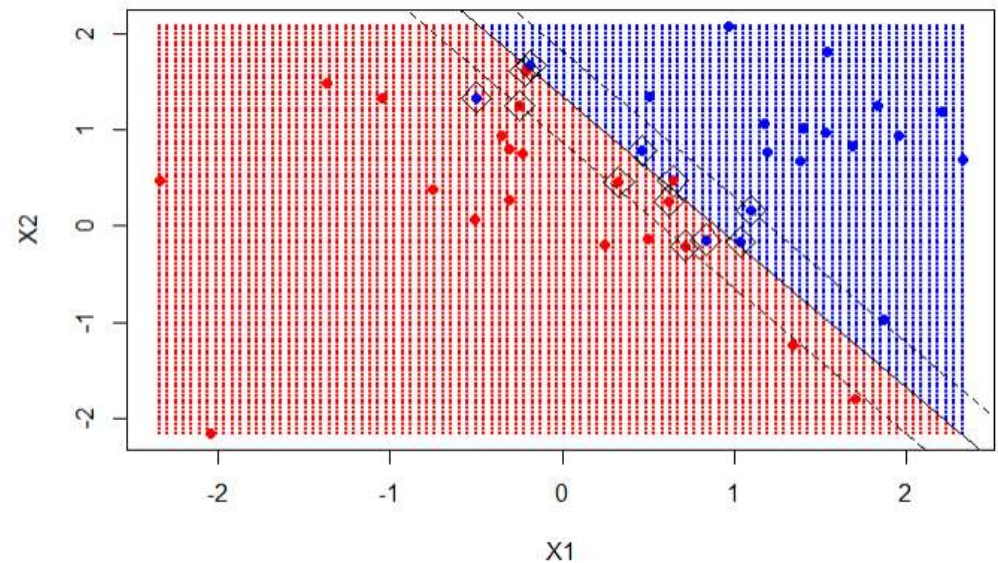
- Create a function named “make.grid”

```
46+ ...{r}  
47+ make.grid = function(x, n = 100) {  
48+   grange = apply(x, 2, range)  
49+   x1 = seq(from = grange[1,1], to = grange[2,1], length = n)  
50+   x2 = seq(from = grange[1,2], to = grange[2,2], length = n)  
51+   expand.grid(x1 = x1, x2 = x2)  
52+ }  
53+ xgrid = make.grid(x)  
54+ # xgrid[1:10,]  
55+ ygrid = predict(svmfit, xgrid)  
56+ plot(xgrid, col = c("red", "blue")[as.numeric(ygrid)], pch = 20, cex = .2)  
57+ points(x, col = y + 3, pch = 19)  
58+ points(x[svmfit$index,], pch = 5, cex = 2)  
59+ ...
```



Linear SVM Example

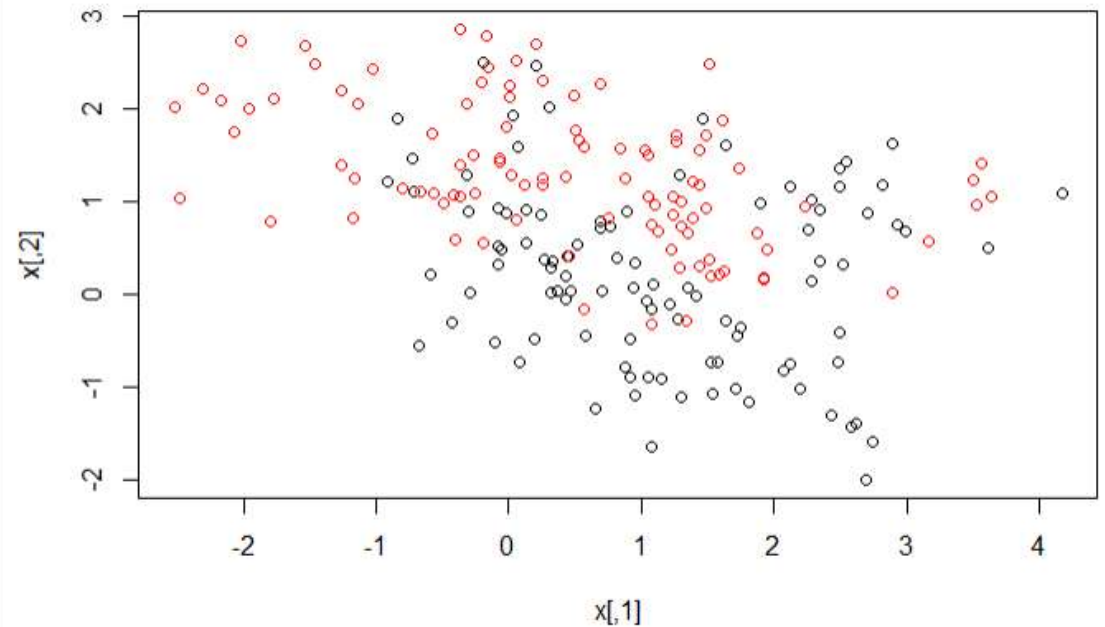
```
60 {r}
61 beta = drop(t(svmfit$coefs)%%x[svmfit$index,])
62 beta0 = svmfit$rho
63 plot(xgrid, col = c("red", "blue")[as.numeric(ygrid)], pch = 20, cex = .2)
64 points(x, col = y + 3, pch = 19)
65 points(x[svmfit$index,], pch = 5, cex = 2)
66 abline(beta0 / beta[2], -beta[1] / beta[2])
67 abline((beta0 - 1) / beta[2], -beta[1] / beta[2], lty = 2)
68 abline((beta0 + 1) / beta[2], -beta[1] / beta[2], lty = 2)
69 }
```



Nonlinear SVM

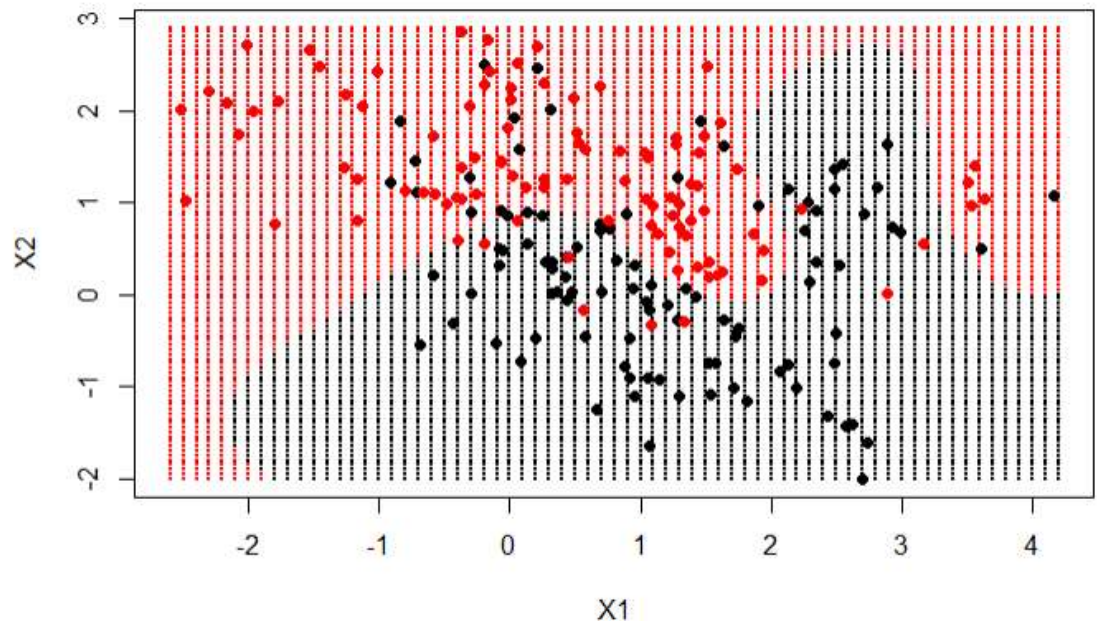
- Example from the textbook

```
80 {r}  
81 load(file = "ESL.mixture.rda")  
82 names(ESL.mixture)|  
83 rm(x, y)  
84 attach(ESL.mixture)  
85 plot(x, col = y + 1)  
86
```



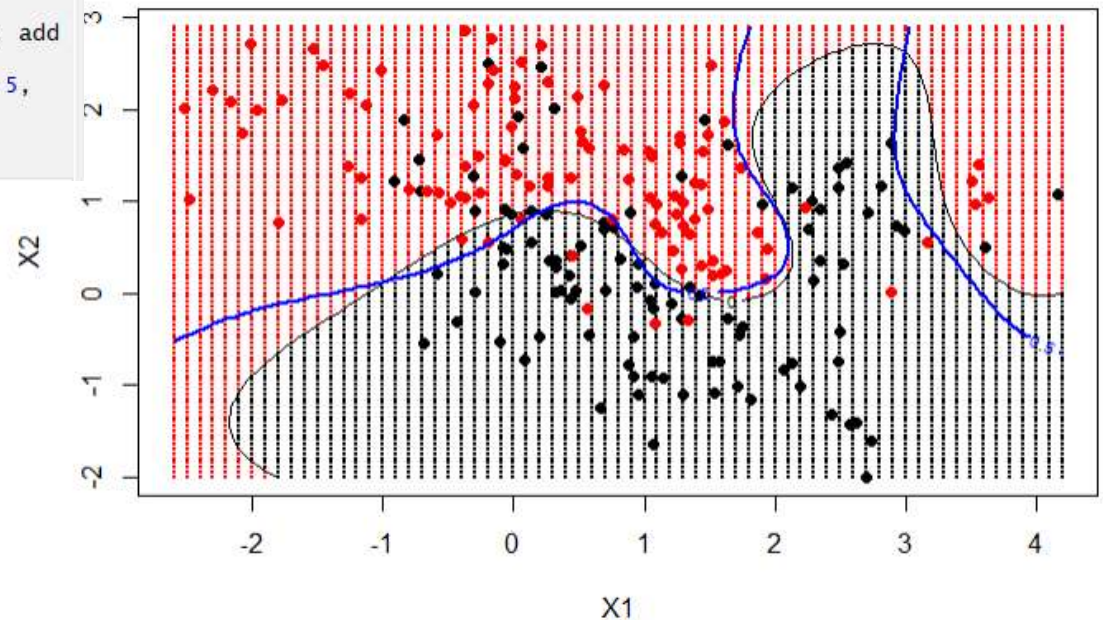
Nonlinear SVM

```
88 + ...{r}  
89 dat = data.frame(y = factor(y), x)  
90 fit = svm(factor(y) ~ ., data = dat, scale = FALSE, kernel = "radial", cost = 5)  
91 xgrid = expand.grid(x1 = px1, x2 = px2)  
92 ygrid = predict(fit, xgrid)  
93  
94 plot(xgrid, col = as.numeric(ygrid), pch = 20, cex = .2)  
95 points(x, col = y + 1, pch = 19)  
96 ...
```



Nonlinear SVM

```
97 {r}  
98 func = predict(fit, xgrid, decision.values = TRUE)  
99 func = attributes(func)$decision  
100  
101 xgrid = expand.grid(x1 = px1, x2 = px2)  
102 ygrid = predict(fit, xgrid)  
103 plot(xgrid, col = as.numeric(ygrid), pch = 20, cex =  
104 .2)  
105 points(x, col = y + 1, pch = 19)  
106 contour(px1, px2, matrix(func, 69, 99), level = 0, add  
107 = TRUE)  
108 contour(px1, px2, matrix(prob, 69, 99), level = 0.5,  
109 add = TRUE, col = "blue", lwd = 2)
```



Naïve Bayes Classifier Demonstration

```
71 ▾ ### Naïve Bayes Classifier Demonstration
72 ▾ ```{r}
73 NBclassifier=naiveBayes(y~., data=dat)
74 print(NBclassifier)
75 dat.test.NB=predict(NBclassifier, newdata = dat.test, type = "class")
76 confusionM.svm<-confusionMatrix(dat.test.NB,dat.test$y)
77 print(confusionM.svm)
78 ...
```

Confusion Matrix and Statistics

	Reference	
Prediction	-1	1
-1	17	7
1	3	13

Sensitivity : 0.8500
Specificity : 0.6500
Pos Pred Value : 0.7083
Neg Pred Value : 0.8125
Prevalence : 0.5000
Detection Rate : 0.4250
Detection Prevalence : 0.6000
Balanced Accuracy : 0.7500

'Positive' Class : -1