

Problem 1

(a)

$$R_{tr}(\hat{(\beta)}) = \frac{1}{N} \sum_{i=1}^n (y_i - \beta^T x_i)^2 < \frac{1}{N-p} \sum_{i=1}^n (y_i - \beta^T x_i)^2 = \text{MSE}$$

$$\Rightarrow E(R_{tr}) < E(\text{MSE}) = \sigma^2$$

(b)

$$E(R_{te}(\hat{(\beta)})) = \frac{1}{M} E \sum_{i=1}^M (\tilde{y}_i - \hat{\beta}^T \tilde{x}_i)^2$$

$$= \frac{1}{M} E \left[\sum_{i=1}^M \left((\tilde{y}_i - \hat{\beta}^T \tilde{x}_i) + (\beta^T \tilde{x}_i - \beta^T \tilde{x}_i) \right)^2 \right]$$

$$= \frac{1}{M} E \left[\sum_{i=1}^M (\tilde{y}_i - \beta^T \tilde{x}_i)^2 + \sum_{i=1}^M (\hat{\beta}^T \tilde{x}_i - \beta^T \tilde{x}_i)^2 \right]$$

$$= \frac{1}{M} E \sum_{i=1}^M (\tilde{y}_i - \beta^T \tilde{x}_i)^2 + \frac{1}{M} E \sum_{i=1}^M (\hat{\beta}^T \tilde{x}_i - \beta^T \tilde{x}_i)^2$$

$$= \sigma^2 + (\text{something squared}) > \sigma^2$$

- (c) On average, the training set error will be *less* than the testing set error. This has many implications in both prediction and inference.
- (d) The attached code generates random noise data sets, calculates training and test set risks, then makes a histogram based on 100 simulations. The blue line is $\sigma^2 = 1$, and the red line is the average risk for training and test set.

```

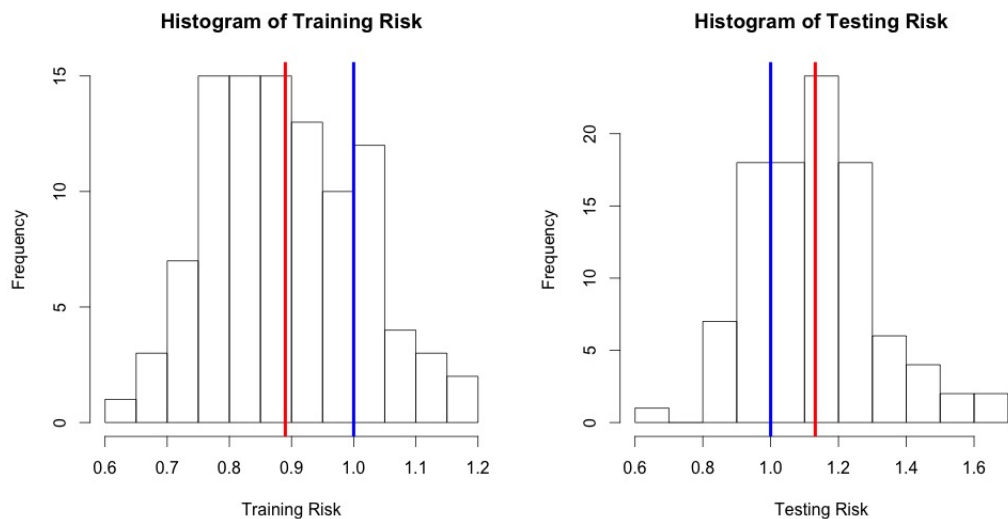
1 generateData = function(n, m, p){
2   #Takes in training and testing sample sizes and number of predictors.
3   #Generates simple noise data
4   yTrain = rnorm(n) #Responses
5   xTrain = matrix(rnorm(n*p), ncol=p) #Predictors
6
7   yTest = rnorm(m) #Responses
8   xTest = matrix(rnorm(m*p), ncol=p) #Predictors
9   list(yTrain=yTrain, xTrain=xTrain, yTest=yTest, xTest=xTest)
10 }
11
12 calcRisks = function(trainTest){
13   # Takes in training and testing data in a list.
14   #Unrolls into matrices, then calculates
15   # the training and testing errors (scalars). Returns a 2-vector.
16   yTrain = trainTest$yTrain
17   xTrain = trainTest$xTrain

```

```

18 yTest = trainTest$yTest
19 xTest = trainTest$xTest
20
21 #Compute training OLS estimate, get residuals/coeffs from it.
22 trainLM = lm(yTrain ~ xTrain)
23 trainResid = trainLM$resid
24 trainCoeffs = trainLM$coefficients
25
26 #Use training OLS estimates to make testing fits. Compute errors.
27 testFits = cbind(rep(1,nrow(xTest)), xTest) %*% trainCoeffs
28 testResid = yTest - testFits
29
30 riskTrain = 1/length(yTrain) * sum(trainResid^2)
31 riskTest = 1/length(yTest) * sum(testResid^2)
32 c(riskTrain, riskTest)
33 }
34
35 n = 100 #Train sample size
36 m = 100 #Test sample size
37 p = 10 #Number of predictors
38 nsim = 100 #Number of simulations
39
40 simErrors = t(sapply(1:nsim, function(i) calcRisks(generateData(n, m, p))))
41 par(mfrow=1:2)
42 hist(simErrors[,1], main='Histogram of Training Risk', xlab='Training Risk')
43 abline(v = mean(simErrors[,1]), col='red', lwd=4)
44 abline(v = 1, col='blue', lwd=4)
45 hist(simErrors[,2], main='Histogram of Testing Risk', xlab='Testing Risk')
46 abline(v = mean(simErrors[,2]), col='red', lwd=4)
47 abline(v = 1, col='blue', lwd=4)

```



Problem 2

- (a) We want a confidence interval for $\sum_{j=0}^3 \beta_j x_0^j$. This is accomplished by noting that (at least asymptotically) the OLS estimates are normal, so a linear combination is normal as well. We do not know σ^2 however, and estimating it, we incur a small penalty. Because we must estimate σ^2 , the linear combination, when suitably standardized, will have a t_{n-p} distribution instead of a normal. The point estimate will be $\hat{f}(x_0) = x_0 \hat{\beta}$, $\beta \in \mathbb{R}^4$,

and where, forgive the abuse of notation, $x_0 = (1, x_0, x_0^2, x_0^3)$. The variance is

$$V(\hat{f}(x_0)) = x_0' V(\hat{\beta}) x_0 = \sigma^2 \cdot x_0' \cdot (X'X)^{-1} \cdot x_0.$$

Not knowing σ^2 , we use the plug-in estimator, MSE, and argue $\left[V(\hat{f}(x_0))\right]^{-1/2} \hat{f}(x_0)$ is t-distributed with $n - p$ degrees of freedom. A confidence interval for $f(x_0)$ is

$$\hat{f}(x_0) \pm \sqrt{\text{MSE} \cdot x_0' \cdot (X'X)^{-1} \cdot x_0} \cdot t_{n-p, 1-\alpha/2}.$$

- (b) The confidence set method follows by applying the Scheffé Method, which creates simultaneous intervals for β , then using this interval to come up with all possible linear combinations for $\sum_{i=0}^3 \beta_j x_0^j$, $\forall x_0$. This interval is

$$\hat{f}(x_0) \pm \sqrt{\text{MSE} \cdot x_0' \cdot (X'X)^{-1} \cdot x_0 \cdot \chi_{4, 1-\alpha}^2}.$$

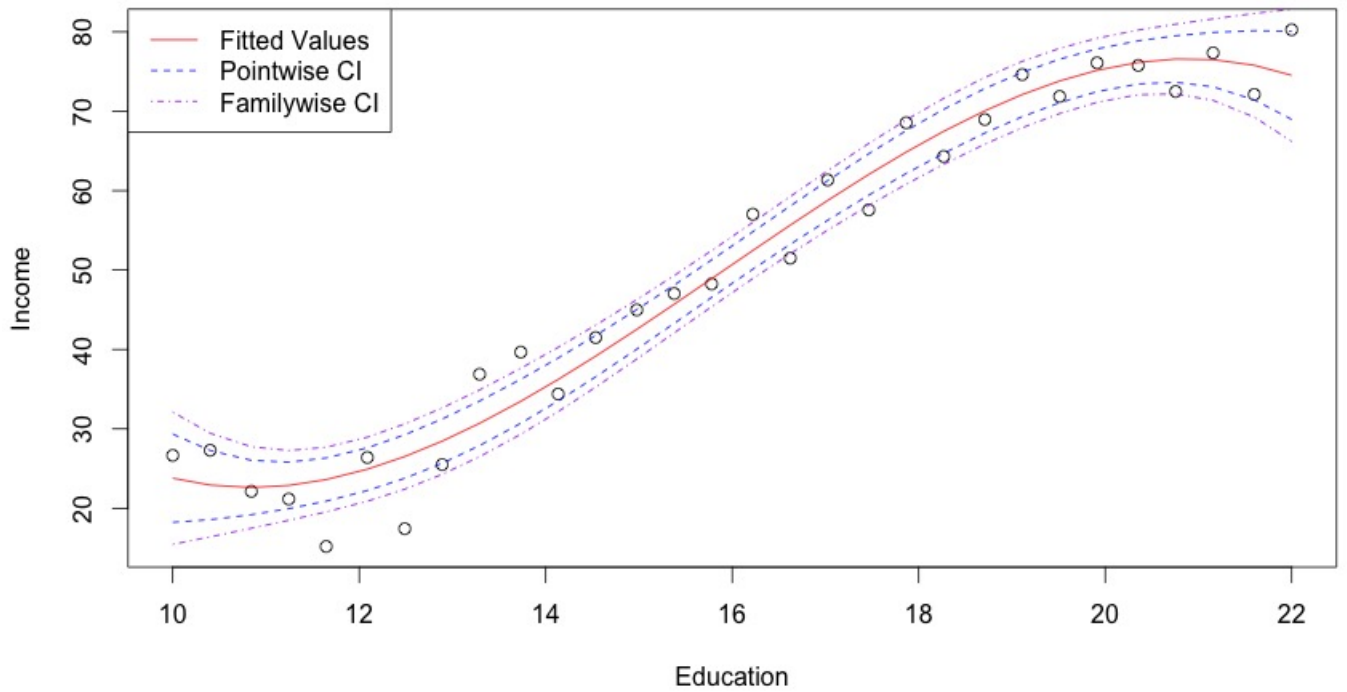
There was some note about using $pF_{p, n-p, 1-\alpha}$. The χ^2 method is asymptotic. These methods are equivalent for large n and small p/n , though the F method is based on non-asymptotic normality where we need to estimate σ^2 , and the χ^2 relies on Central Limit Theorem.

The method in (a) needs only give coverage to the point $\hat{f}(x_0)$, whereas the method in (b) gives proper coverage for any linear combination of the form $\sum_{j=0}^3 \beta_j x_0^j$ we choose. This means we have to be a bit more conservative than if we could make statements about a particular x_0 .

```

1 income = read.csv("http://www-bcf.usc.edu/~gareth/ISL/Income1.csv")[,-1]
2 incomeModel = lm(Income ~ Education + I(Education^2) + I(Education^3), data=income)
3 pointwiseCI = predict(incomeModel, interval = "confidence")
4
5 multiplier = sqrt(qchisq(.95, 4))
6 familywiseCI = predict(incomeModel, se.fit=TRUE)
7 scheffeInterval = data.frame(fit = familywiseCI$fit,
8                               lwr = familywiseCI$fit - multiplier * familywiseCI$se.fit,
9                               upr = familywiseCI$fit + multiplier * familywiseCI$se.fit)
10
11 plot(Income ~ Education, data=income)
12 legend("topleft", c("Fitted Values", "Pointwise CI", "Familywise CI"),
13        lty = c(1, 2, 4), col = c("red", "blue", "purple")) #Add a legend
14 #Plot Pointwise CI
15 lines(income$Education, pointwiseCI[,1], col = "red")
16 lines(income$Education, pointwiseCI[,2], lty = 2, col = "blue")
17 lines(income$Education, pointwiseCI[,3], lty = 2, col = "blue")
18 #Plot Scheffe
19 lines(income$Education, scheffeInterval[,2], lty = 4, col = "purple")
20 lines(income$Education, scheffeInterval[,3], lty = 4, col = "purple")

```



Problem 3

- (a) A possible rule would be to classify observations to $Y = 1$ if $P(Y = 1) = \text{sigmoid}(\hat{\beta}^T X) \geq 0.5$. This is equivalent to classifying to $Y = 1$ if $\hat{\beta}^T X \geq 0$.

```

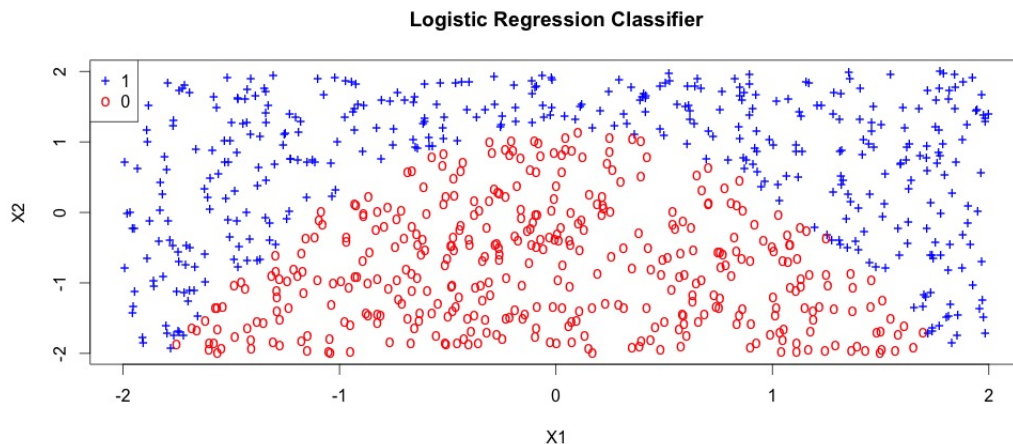
1 load("hw1prob3.Rdata")
2 logisticY1 = glm(y1~x, family='binomial')
3 plot(x[,1], x[,2], col=c("red","blue")[y1+1], pch = c('o', '+')[y1+1],
4       ylab="X2", xlab="X1", main='Logistic Regression Classifier')
5 legend('topleft', legend=unique(y1), col=c('blue', 'red'), pch=c('+', 'o'))
6 coef(logisticY1) #Classification Rule: b'X >= 0, classify to Y=1.
7
8 truePred = cbind(y1, logisticY1$fitted.values >= 0.5) #Obs Y1 and Pred Y1
9 table(Actual=truePred[,1], Predicted=truePred[,2])
10 (101+87)/length(y1)

```

Using our logistic regression, classify to $Y = 1$ if $0.127 + 0.063 \cdot x_1 + 1.390 \cdot x_2 \geq 0$. The confusion matrix:

	Predict 0	Predict 1
Actual 0	308	87
Actual 1	101	304

The misclassification rate is $(101 + 87)/800 = .235$.



(b) Classification rule drawn on the plot, corresponds to the line

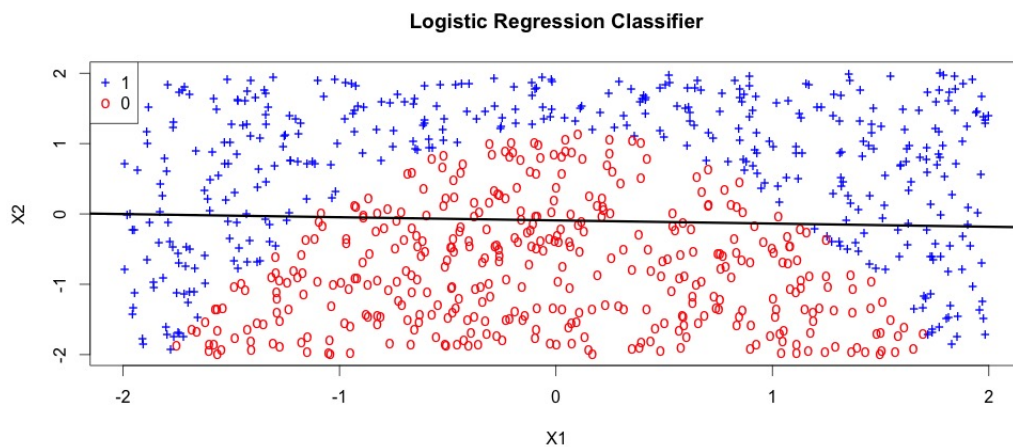
$$x_2 = -0.063x_1/1.390 - 0.127/1.390.$$

This is obtained by taking the classification rule and converting it to slope-intercept form.

```

1 plot(x[,1], x[,2], col=c("red","blue")[y1+1], pch = c('o', '+')[y1+1],
2      ylab="X2", xlab="X1", main='Logistic Regression Classifier')
3 legend('topleft', legend=unique(y1), col=c('blue', 'red'), pch=c('+', 'o'))
4
5 logisticVSlope = -coef(logisticY1)[2] / coef(logisticY1)[3]
6 logisticCVInt = -coef(logisticY1)[1] / coef(logisticY1)[3]
7 abline(a = logisticCVInt, b = logisticVSlope, col='black', lty=1, lwd=3)

```



The classifier is clearly nonlinear (it looks like it could be quadratic or periodic with a restricted support), and yet we fit a linear rule to the data set. Not smart!

- (c) Similar to the plotting method described in (b), convert general form to slope-intercept form.

```
1 logisticY1.X1Square = glm(y1~x + I(x[,1]^2), family='binomial')
2 truePred = cbind(y1, logisticY1.X1Square$fitted.values >= 0.5)
3 table(Actual=truePred[,1], Predicted=truePred[,2])
```

We arrive at the conclusion that we classify to $Y = 1$ if $-2322.5 - 11.6x_1 + 1955.6x_2 + 1937.1x_1^2 \geq 0$.

	Predict 0	Predict 1
Actual 0	395	0
Actual 1	0	405

There is a perfect classification rate (0% misclassification rate). This is definitely a better classifier, not just based on the classification rate, but it is clear that this classifier much more faithfully captures the relationship of y_1 on (x_1, x_2) .

(d)

```
plot(x[,1], x[,2], col=c("red","blue")[y1+1], pch = c('o', '+')[y1+1],
     ylab="X2", xlab="X1", main='Logistic Regression Classifier')
legend('topleft', legend=unique(y1), col=c('blue', 'red'), pch=c('+', 'o'))
a = -coef(logisticY1.X1Square)[1] / coef(logisticY1.X1Square)[3]
b = -coef(logisticY1.X1Square)[2] / coef(logisticY1.X1Square)[3]
c = -coef(logisticY1.X1Square)[4] / coef(logisticY1.X1Square)[3]
curve(a + b*x + c*x^2, from=min(x[,1]), to=max(x[,1]), add = TRUE, lwd=3)
```



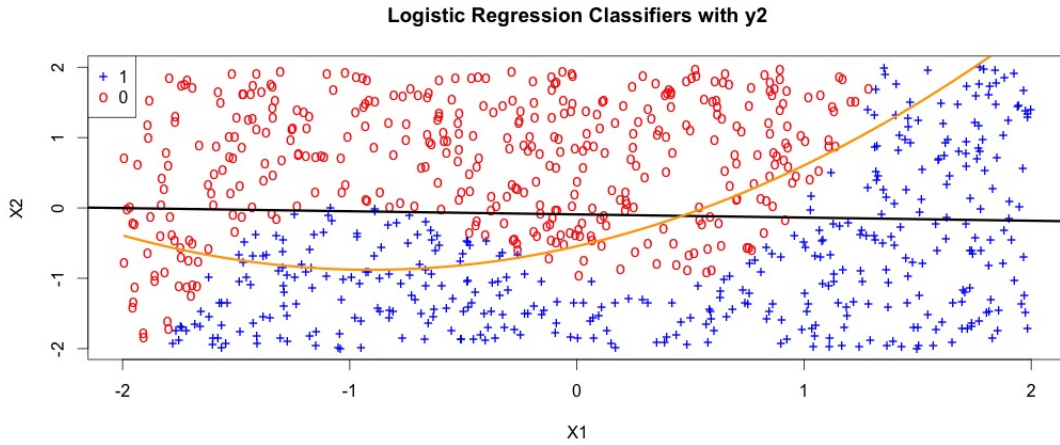
The classifier has a quadratic shape in x_1 .

- (e) What happens when we look at y_2 instead of y_1 ?

```

1 logisticY2 = glm(y2~x, family='binomial')
2 logisticY2.X1Square = glm(y2~x + I(x[,1]^2), family='binomial')
3 plot(x[,1], x[,2], col=c("red","blue")[y2+1], pch = c('o', '+')[y2+1],
4      ylab="X2", xlab="X1", main='Logistic Regression Classifiers with y2')
5 legend('topleft', legend=unique(y1), col=c('blue', 'red'), pch=c('+', 'o'))
6 a = -coef(logisticY2)[1] / coef(logisticY2)[3]
7 b = -coef(logisticY2)[2] / coef(logisticY2)[3]
8 abline(a = logisticCVInt, b = logistCVSlope, col='black', lty=1, lwd=3)
9 a = -coef(logisticY2.X1Square)[1] / coef(logisticY2.X1Square)[3]
10 b = -coef(logisticY2.X1Square)[2] / coef(logisticY2.X1Square)[3]
11 c = -coef(logisticY2.X1Square)[4] / coef(logisticY2.X1Square)[3]
12 curve(a + b*x + c*x^2, from=min(x[,1]), to=max(x[,1]), add = TRUE, lwd=3, col='orange')
13 truePredLinear = cbind(y2, logisticY2$fitted.values >= 0.5)
14 mean(truePredLinear[,1] != truePredLinear[,2])
15 truePredQuadratic = cbind(y2, logisticY2.X1Square$fitted.values >= 0.5)
16 mean(truePredQuadratic[,1] != truePredQuadratic[,2])

```



The linear rule has a misclassification rate of 18.5% while the quadratic rule has a misclassification of 13.75%.

- (f) Neither rule is particularly good, because neither captures the underlying structure of the data, which appears to be higher degree than just quadratic. One approach we can do is to fit a cubic (third-order) term, which is what the relationship seems to be. If we try third-order, and classify $Y_2 = 1$ if

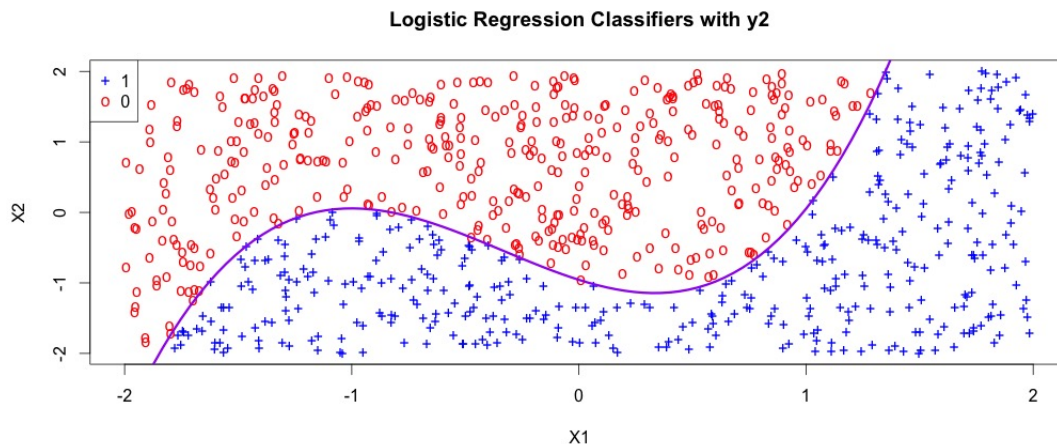
$$-1455 - 1541x_1 - 1521x_2 + 1537x_1^2 + 1537x_1^3 \geq 0,$$

we can achieve perfect classification.

```

1 logisticY2.X1Third = glm(y2~x + I(x[,1]^2) + I(x[,1]^3), family='binomial')
2 plot(x[,1], x[,2], col=c("red","blue")[y2+1], pch = c('o', '+')[y2+1],
3      ylab="X2", xlab="X1", main='Logistic Regression Classifiers with y2')
4 legend('topleft', legend=unique(y1), col=c('blue', 'red'), pch=c('+', 'o'))
5 a = -coef(logisticY2.X1Third)[1] / coef(logisticY2.X1Third)[3]
6 b = -coef(logisticY2.X1Third)[2] / coef(logisticY2.X1Third)[3]
7 c = -coef(logisticY2.X1Third)[4] / coef(logisticY2.X1Third)[3]
8 d = -coef(logisticY2.X1Third)[5] / coef(logisticY2.X1Third)[3]
9 curve(a + b*x + c*x^2 + d*x^3, from=min(x[,1]), to=max(x[,1]), add = TRUE, lwd=3, col='purple')
10 truePredCubic = cbind(y2, logisticY2.X1Third$fitted.values >= 0.5)
11 mean(truePredCubic[,1] != truePredCubic[,2])

```

- (g) Why not keep adding more and more polynomial terms? We have seen before that we run the risk of over-fitting by adding more and more terms. That is, we begin to fine-tune a model that performs well in the sample we have, but may not generalize well to data we haven't collected yet. There are a few approaches to dealing with this over-fitting problem. We've already encountered a few in this course, as well as in previous courses. One popular way is to use cross-validation, fitting several candidate models on one set of data (the training set) and validating the model's performance on a set of data that was not used in training (the testing set), and choosing the model that performs best on the test set. Other ways would involve striking some balance between classification accuracy and the size of the model. These include AIC, BIC, Mallows's C_p , and many more.