

Individual_assignment4

Yuhan_Xu_474154

2019/9/18

Problem 10

This question should be answered using the Weekly data set, which is part of the ISLR package. This data is similar in nature to the Smarket data from this chapter's lab, except that it contains 1,089 weekly returns for 21 years, from the beginning of 1990 to the end of 2010.

```
library("ISLR")
attach(Weekly)
```

for context: (d) Now fit the logistic regression model using a training data period from 1990 to 2008, with Lag2 as the only predictor. Compute the confusion matrix and the overall fraction of correct predictions for the held out data (that is, the data from 2009 and 2010).

```
train = (Year < 2009)
Weekly.test = Weekly[!train,]
Direction.test = Direction[!train]
```

(f) Repeat (d) using QDA.

```
library("MASS")
qda.fit = qda(Direction ~ Lag2, data = Weekly, subset = train)
qda.class = predict(qda.fit, Weekly.test)$class
table(qda.class, Direction.test)

##           Direction.test
## qda.class Down Up
##      Down    0  0
##      Up     43 61

mean(qda.class == Direction.test)

## [1] 0.5865385
```

Problem 8

We will now perform cross-validation on a simulated data set.

(a)

Q: Generate a simulated data set as follows:

```
#> set.seed(1)
```

```
#> x=rnorm (100)
```

```
#> y=x-2* x^2+ rnorm (100)
```

In this data set, what is n and what is p? Write out the model used to generate the data in equation form.

```
set.seed(1)
```

```
x = rnorm(100)
```

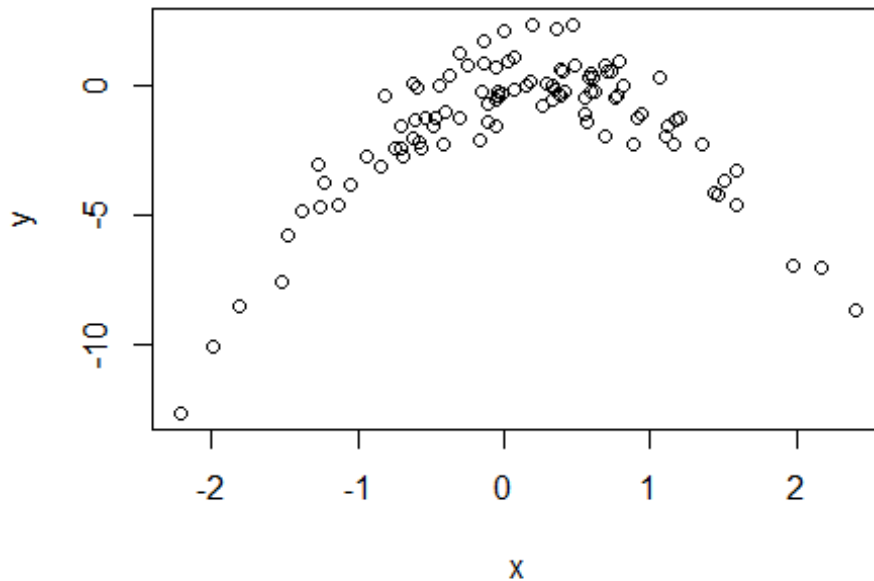
```
y = x-2*x^2 + rnorm(100)
```

A: $n = 100$, $p = 2$, the model is $Y = X - 2X^2 + \epsilon$

(b)

Q: Create a scatterplot of X against Y. Comment on what you find.

```
plot(x, y)
```



A: Y and X has a non-linear relationshp. The shape of the function seems like a parabola.

(c)

Q: Set a random seed, and then compute the LOOCV errors that result from fitting the following four models using least squares:

i. $Y = \beta_0 + \beta_1 X + \epsilon$

ii. $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$

iii. $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \varepsilon$

iv. $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \varepsilon$

Note you may find it helpful to use the `data.frame()` function to create a single data set containing both X and Y.

A:

```
set.seed(1)
data = data.frame(x,y)

library(boot)
#i.
glm.fit1 = glm(y~x)
cv.err1 = cv.glm(data, glm.fit1)
cv.err1$delta

## [1] 7.288162 7.284744

#ii.
glm.fit2 = glm(y~poly(x,2))
cv.err2 = cv.glm(data, glm.fit2)
cv.err2$delta

## [1] 0.9374236 0.9371789

#iii.
glm.fit3 = glm(y~poly(x,3))
cv.err3 = cv.glm(data, glm.fit3)
cv.err3$delta

## [1] 0.9566218 0.9562538

#iv.
glm.fit4 = glm(y~poly(x,4))
cv.err4 = cv.glm(data, glm.fit4)
cv.err4$delta

## [1] 0.9539049 0.9534453
```

(d)

Q: Repeat (c) using another random seed, and report your results. Are your results the same as what you got in (c)? Why?

```
set.seed(2)

#i.
glm.fit5 = glm(y~x)
```

```

cv.err5 = cv.glm(data, glm.fit5)
cv.err5$delta

## [1] 7.288162 7.284744

#ii.
glm.fit6 = glm(y~poly(x,2))
cv.err6 = cv.glm(data, glm.fit6)
cv.err6$delta

## [1] 0.9374236 0.9371789

#iii.
glm.fit7 = glm(y~poly(x,3))
cv.err7 = cv.glm(data, glm.fit7)
cv.err7$delta

## [1] 0.9566218 0.9562538

#iv.
glm.fit8 = glm(y~poly(x,4))
cv.err8 = cv.glm(data, glm.fit8)
cv.err8$delta

## [1] 0.9539049 0.9534453

```

A: The results are the same as (c). Because LOOCV produces the average MSE from every single observation, splitting the data set in different ways will not affect the results.

(e)

Q: Which of the models in (c) had the smallest LOOCV error? Is this what you expected? Explain your answer.

A: The second model **ii**. $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$ had the smallest LOOCV error. This what I expected since the true function is also quadratic.

(f)

Q: Comment on the statistical significance of the coefficient estimates that results from fitting each of the models in (c) using least squares. Do these results agree with the conclusions drawn based on the cross-validation results?

```

summary(glm.fit1)

##
## Call:
## glm(formula = y ~ x)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max

```

```
## -9.5161 -0.6800 0.6812 1.5491 3.8183
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.6254      0.2619  -6.205 1.31e-08 ***
## x           0.6925      0.2909   2.380 0.0192 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 6.760719)
##
##      Null deviance: 700.85  on 99  degrees of freedom
## Residual deviance: 662.55  on 98  degrees of freedom
## AIC: 478.88
##
## Number of Fisher Scoring iterations: 2

summary(glm.fit2)

##
## Call:
## glm(formula = y ~ poly(x, 2))
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9650  -0.6254  -0.1288   0.5803   2.2700
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.5500      0.0958  -16.18 < 2e-16 ***
## poly(x, 2)1  6.1888      0.9580   6.46 4.18e-09 ***
## poly(x, 2)2 -23.9483      0.9580 -25.00 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.9178258)
##
##      Null deviance: 700.852  on 99  degrees of freedom
## Residual deviance:  89.029  on 97  degrees of freedom
## AIC: 280.17
##
## Number of Fisher Scoring iterations: 2

summary(glm.fit3)

##
## Call:
## glm(formula = y ~ poly(x, 3))
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -1.9765 -0.6302 -0.1227 0.5545 2.2843
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.55002    0.09626 -16.102 < 2e-16 ***
## poly(x, 3)1  6.18883    0.96263  6.429 4.97e-09 ***
## poly(x, 3)2 -23.94830    0.96263 -24.878 < 2e-16 ***
## poly(x, 3)3  0.26411    0.96263  0.274 0.784
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.9266599)
##
## Null deviance: 700.852 on 99 degrees of freedom
## Residual deviance: 88.959 on 96 degrees of freedom
## AIC: 282.09
##
## Number of Fisher Scoring iterations: 2

summary(glm.fit4)

##
## Call:
## glm(formula = y ~ poly(x, 4))
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.0550  -0.6212  -0.1567   0.5952   2.2267
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.55002    0.09591 -16.162 < 2e-16 ***
## poly(x, 4)1  6.18883    0.95905  6.453 4.59e-09 ***
## poly(x, 4)2 -23.94830    0.95905 -24.971 < 2e-16 ***
## poly(x, 4)3  0.26411    0.95905  0.275 0.784
## poly(x, 4)4  1.25710    0.95905  1.311 0.193
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.9197797)
##
## Null deviance: 700.852 on 99 degrees of freedom
## Residual deviance: 87.379 on 95 degrees of freedom
## AIC: 282.3
##
## Number of Fisher Scoring iterations: 2
```

A: The p-value of x and x^2 is extremely small which indicates that x , x^2 are statistical significant. This agrees with the conclusions draw based on the cross-validation result.