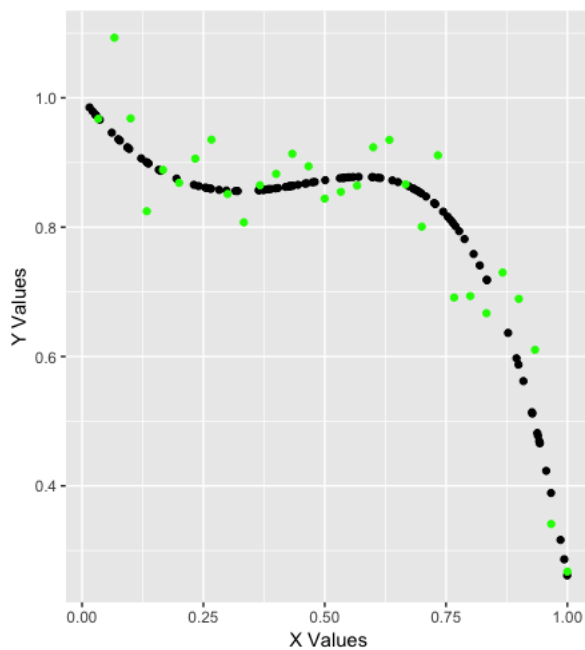# Statistical Learning Problem Set # 2

John Gaebler – `joga5948@colorado.edu`

February 10, 2020

---

1. Creation of Datasets

   (a) Write a function that generates samples - see appendix code

   (b) Display the true function f(x) using 100 uniformly spaces samples on [0,1] and the 30 noisy samples on the same plot.



*The black points correspond to the true function, and the green points are the noisy data sample*

1

2. Estimation of the Polynomial

   (a) Write a function that computes the optimal estimator $\hat{y}$ given a fixed order $p$. Give the coefficients generated for p = 1,4, and 15

   *Coefficients: p = 1*

   | Intercept | x |
   |-----------|-----------|
   | 1.0395355 | -0.4407709 |

   *Coefficients: p = 4*

   | Intercept | x | $x^2$ | $x^3$ | $x^4$ |
   |-----------|-----------|-----------|-----------|-----------|
   | 1.05401141 | -1.24916710 | 2.23042888 | -0.02250314 | -1.71335123 |

   *Coefficients: p = 15*

   | Intercept | x | $x^2$ | $x^3$ | $x^4$ | - |
   |-----------|-----------|-----------|-----------|-----------|---|
   | -5.228157e-01 | 8.330017e+01 | -1.528828e+03 | 1.263023e+04 | -4.689995e+04 | - |
   | $x^5$ | $x^6$ | $x^7$ | $x^8$ | $x^9$ | - |
   | 4.022791e+03 | 6.963936e+05 | -3.157657e+06 | 7.470297e+06 | -1.087213e+07 | - |
   | $x^{10}$ | $x^{11}$ | $x^{12}$ | $x^{13}$ | $x^{14}$ | $x^{15}$ |
   | 9.938791e+06 | -5.401311e+06 | 1.415607e+06 | NA | -5.829778e+04 | NA |

   *Note: The $x^{13}$ and $x^{15}$ coefficients in the p = 15 model are NA because they are linearly related to other components of the model. Past degree 14, NAs appear for the coefficents with very high frequency.*

   (b) Display the three optimal polynomials as well as the true function and the noisy data on the same plot



   *Blue line: degree 1, Red line: degree 4, Pink line: degree 15*

(c) Based on visual inspection, what is the optimal degree polynomial for:

    i. The best fit to the noisy data?
       *Answer: The degree 15 polynomial best fits the noisy data points.*

   ii. The best fit to the true function?
       *Answer: The degree 4 polynomial best fits the path of the true function.*

The higher degree polynomial can do a better job of fitting the variation in the data points, but it fails to capture the actual movement of the function. The lower degree polynomial may actually better approximate the true function even if it doesn't appear to follow the noisy data as well. We call the effect of matching the noisy data too closely and as a result loosing the true function "over-fitting the model."
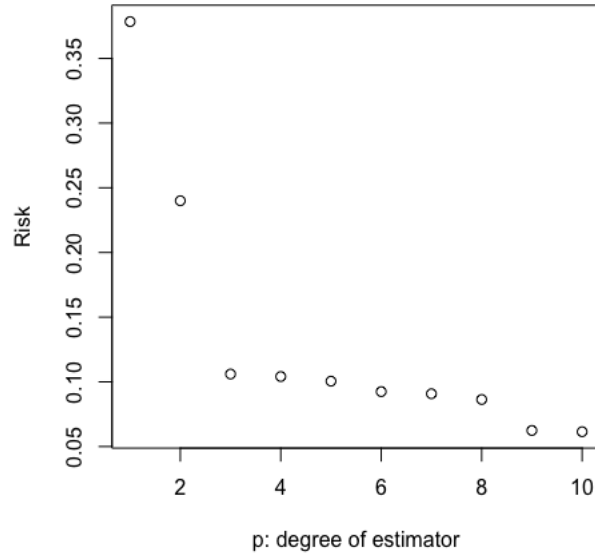
3. Choice of the Optimal Degree

(a) For p = 1, ... , 10 compute the values of $\hat{R}$. Display $\hat{R}$ as a function of p, and find the value of p that minimizes $\hat{R}$.

Values of $\hat{R}$ for each value of p:

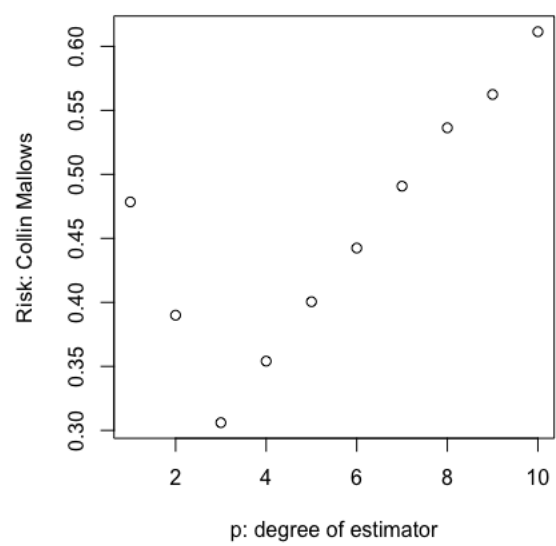| p = 1 | p = 2 | p = 3 | p = 4 | p = 5 |
|---|---|---|---|---|
| 0.37845020 | 0.23997598 | 0.10603381 | 0.10410272 | 0.10051470 |
| p = 6 | p = 7 | p = 8 | p = 9 | p = 10 |
| 0.09244547 | 0.09079719 | 0.08636445 | 0.06241269 | 0.06144891 |

The value of p that minimizes $\hat{R}$ is p = 10.



(b) To penalize overly complex models, use the Collin Mallows statistic $C_p$ to choose the optimal model. Display $C_p$ as a function of , and find the value of p that minimizes $C_p$.

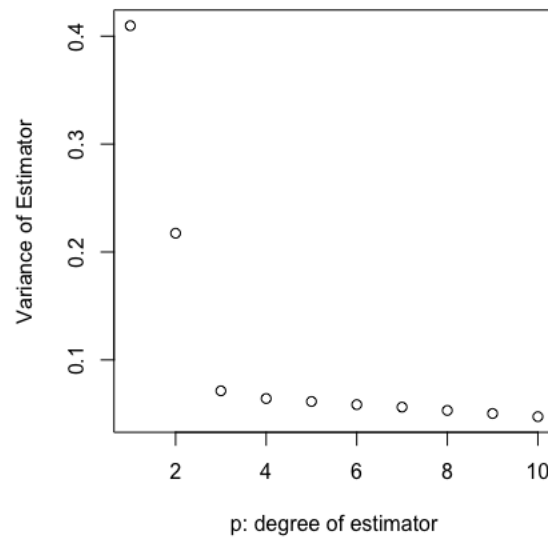| p = 1 | p = 2 | p = 3 | p = 4 | p = 5 |
|---|---|---|---|---|
| 0.4784502 | 0.3899760 | 0.3060338 | 0.3541027 | 0.4005147 |
| p = 6 | p = 7 | p = 8 | p = 9 | p = 10 |
| 0.4424455 | 0.4907972 | 0.5363644 | 0.5624127 | 0.6114489 |

The value of p that minimizes $C_p$ is p = 3.
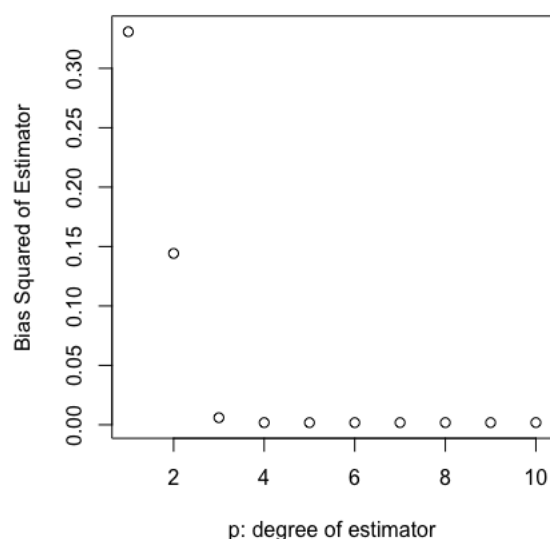
4

4. Estimate of the Bias and Variance of the Estimator

(a) Let S $= 100$. Generate S random realizations of $y^{(s)}$ using the code that you wrote in the first problem. For each $y^{(s)}$, compute the optimal estimate $\hat{y}^{(s)}$. See appendix code Finally, compute the Variance and the Bias of the estimator, and combine theses statistics to to compute the risk:

$$\tilde{R} = Var(\hat{y}) + bias^2(\hat{y})$$
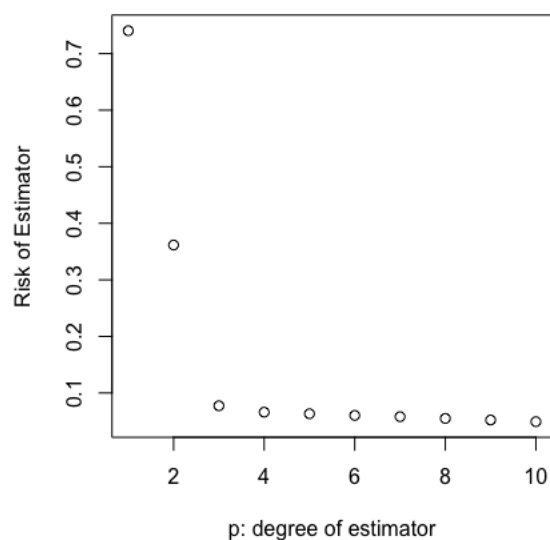
for $p = 1,...,10$. Then display $Var(\hat{y})$ , $bias^2(\hat{y})$, and $\tilde{R}$ as functions of p. Describe and explain the behavior of the curves as p increases. What is the value of p that minimizes $\tilde{R}$?



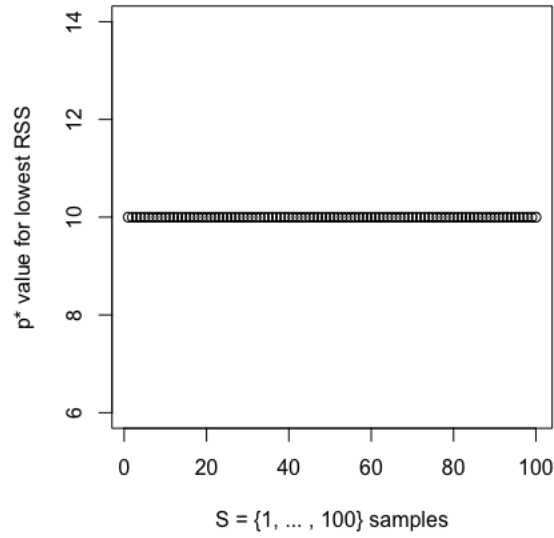*The value of p that minimizes variance is p = 10*

*The value of p that minimizes bias squared is p = 10*



*The value of p that minimizes Risk is p = 10*
For all of the statistics, the values decreases as p increases. This is because a function with higher degree is going to be able to better approximate the sample data. For Bias squared in particular, the value rapidly approaches zero as early as p = 3. However, the higher degree estimator leads to an overly complex model that is likely to over fit the data and produce predictions that are less accurate than a more simple model.

(b) For each of the S data sets, compute the value $p^*$ that minimizes the residual sum of squares. Additionally, compute the value $p^+$ that minimize the the statistic $C_p$.

*For all the S = 100 samples, the degree 10 polynomial estimator produces the least value of Residual Sum of Squares*



*For all the S = 100 samples, the degree 3 polynomial estimator produces the least value of the Collin Mallows statistic.*

The differences between these measures of risk demonstrates the importance of penalizing excessively complex models. There is an apparent large difference in the optimal degree of estimator between these two metrics of risk (3 vs 10). These estimates are very consistent across the S = 100 realizations of testing data, considering there is not a single deviation from p = 10 or p = 3 respectively.

Appendix: Code

```r
1   #=============================================================================
2   ##------------------------Homework #2------------------------------------------
3   library(ggplot2)
4   #Generating data
5   x_Data <- seq(1/30,1, by = 1/30)
6   true_fun_in <- runif(100, 0 ,1)
7
8   #function to generate the samples
9   generate_Data <- function(x, sd =0) {
10      data <- 1 - x + (2 * x ^2) - (.8 *  x ^ 3) + (.06 * x ^4) - (x ^ 5)
11      data <- data + rnorm(30, 0, sd)
12      return(data)
13  }
14
15
16  y_Sample <- generate_Data(x_Data, .05)
17  true_fun_out <- generate_Data(true_fun_in)
18  both_in <- append(true_fun_in, x_Data)
19  both_out<- append(true_fun_out, y_Sample)
20  y_Data
21
22  bigplot <- ggplot()+
23      geom_point(aes(x = true_fun_in, y = true_fun_out))+
24      geom_point(aes(x = x_Data, y = y_Sample), color = "green")+
25      labs(x = "X Values", y = "Y Values")
26
27  ##------------------------------------------------------------------
28  |
29  estimator_degree <- function(n, data = x_Data) {
30      i = 1
31      data_matrix <-  array(data = rep(x_Data, n), dim = c(length(x_Data),n))
32      estimator <- c(rep(x_Data, n))
33      while (i <= n) {
34          data_matrix[,i] <- data_matrix[,i]^i
35          i = i +1
36      }
37      return(data_matrix)
38  }
39
40  #-------------------------------
41  p1_data <- estimator_degree(1)
42  p1 <- lm(y_Sample ~ p1_data)
43
44  degree1 <- function(x) 1.0395 + -0.4408 *x
45
```

```r
44  degree1 <- function(x) 1.0395 + -0.4408 *x
45
46  p2_data <- estimator_degree(2)
47  p2 <- lm(y_Sample ~ p2_data)
48
49  p3_data <- estimator_degree(3)
50  p3 <- lm(y_Sample ~ p3_data)
51
52  p4_data <- estimator_degree(4)
53  p4 <- lm(y_Sample ~ p4_data)
54
55  p5_data <- estimator_degree(5)
56  p5 <- lm(y_Sample ~ p5_data)
57
58  p6_data <- estimator_degree(6)
59  p6 <- lm(y_Sample ~ p6_data)
60
61  p7_data <- estimator_degree(7)
62  p7 <- lm(y_Sample ~ p7_data)
63
64  p8_data <- estimator_degree(8)
65  p8 <- lm(y_Sample ~ p8_data)
66
67  p9_data <- estimator_degree(9)
68  p9 <- lm(y_Sample ~ p9_data)
69
70  p10_data <- estimator_degree(10)
71  p10 <- lm(y_Sample ~ p10_data)
72
73  data_vec <- c(p1_data, p2_data, p3_data, p4_data, p5_data, p6_data, p7_data, p8_data, p9_data, p10_data)
74
75
76
77  degree4 <- function(x) 1.05401 -1.24917*x +2.23043 *x^2 -0.02250 * x^3 -1.71335 *x^4
78
79  p15_data <- estimator_degree(15)
80  p15 <- lm(y_Sample ~ p15_data)
81
82  degree4 <- function(x) {
83      output <- 0
84      for(i in 1:length(coef(p4))) {
85        output <- output + as.numeric(coef(p4)[i]) * x ^(i-1)
86      }
87      return(output)
88  }
```

```r
82 ▾ degree4 <- function(x) {
83      output <- 0
84 ▾    for(i in 1:length(coef(p4))) {
85         output <- output + as.numeric(coef(p4)[i]) * x ^(i-1)
86      }
87      return(output)
88   }
89 ▾ degree15 <- function(x) {
90      output <- 0
91 ▾    for(i in 1:length(coef(p15))) {
92         output <- output + ifelse(is.na(as.numeric(coef(p15)[i]) * x ^(i-1)), 0,  as.numeric(coef(p15)[i]) * x ^(i-1))
93      }
94      return(output)
95   }
96
97 ▾ degree_n <- function(x, degree) {
98      output <- 0
99      reg <- eval(parse(text = paste("p", degree, sep = "")))
100 ▾    for(i in 1:length(coef(reg))) {
101         output <- output + ifelse(is.na(as.numeric(coef(reg)[i]) * x ^(i-1)), 0,  as.numeric(coef(reg)[i]) * x ^(i-1))
102      }
103      return(output)
104   }
105
106   bigplot + stat_function(fun = degree1, aes(x = x), data = data.frame(x = 0), color = "blue") +
107      stat_function(fun = degree4, aes(x = x), data = data.frame(x = 0), color = "red") +
108      stat_function(fun = degree15, aes(x = x), data = data.frame(x = 0), color = "pink") +
109      ylim(.2,1)
110
111
112
113 ▾ Risk <- function(deg) {
114      r_hat <- sum((y_Sample - degree_n(x_Data, deg))^2)
115      return(r_hat)
116   }
117
118   risk_vec <- sapply(c(1:10), Risk)
119   plot(x = 1:10, y = risk_vec,
120        xlab = "p: degree of estimator",
121        ylab = "Risk")
122
123
124   c_p <- risk_vec + 2 * (c(2:11)) * .025
125
```

```r
124  c_p <- risk_vec + 2 * (c(2:11)) * .025
125
126  plot(x = 1:10, y = c_p,
127       xlab = "p: degree of estimator",
128       ylab = "Risk: Collin Mallows")
129  ##-------- Looks like the optimal degree is 3
130
131
132  ##-- Creatin the S = 100 random realization of y(s)
133  y_s <- rerun(100, generate_Data(x_Data, .05))
134
135  model_creator <- function(y, data) {
136      return(lm(y ~ data))
137  }
138  ## generating the p = {1,...10} models for each of the S = 100 realizations, 1000 total models
139  y_s_model_p1 <- lapply(y_s, model_creator, data = p1_data)
140  y_s_model_p2 <- lapply(y_s, model_creator, data = p2_data)
141  y_s_model_p3 <- lapply(y_s, model_creator, data = p3_data)
142  y_s_model_p4 <- lapply(y_s, model_creator, data = p4_data)
143  y_s_model_p5 <- lapply(y_s, model_creator, data = p5_data)
144  y_s_model_p6 <- lapply(y_s, model_creator, data = p6_data)
145  y_s_model_p7 <- lapply(y_s, model_creator, data = p7_data)
146  y_s_model_p8 <- lapply(y_s, model_creator, data = p8_data)
147  y_s_model_p9 <- lapply(y_s, model_creator, data = p9_data)
148  y_s_model_p10 <- lapply(y_s, model_creator, data = p10_data)
149
150  model_vec <- list(y_s_model_p1, y_s_model_p2, y_s_model_p3, y_s_model_p4, y_s_model_p5, y_s_model_p6,
151                    y_s_model_p7, y_s_model_p8, y_s_model_p9, y_s_model_p10)
152
153  mass_predictor <- function(model) {
154      fitted <- 0
155      reg <- model
156      for(i in 1:length(coef(reg))) {
157        fitted <- fitted + ifelse(is.na(as.numeric(coef(reg)[i]) * x_Data ^(i-1)), 0,
158                                  as.numeric(coef(reg)[i]) * x_Data ^(i-1))
159      }
160     return(fitted)
161  }
162  ## creating fitted data with the models, 10 fitted sets per S = 100 datasets, 1000 fitted datasets
163  fitted_vec_p1 <- lapply(y_s_model_p1, mass_predictor)
164  fitted_vec_p2 <- lapply(y_s_model_p2, mass_predictor)
165  fitted_vec_p3 <- lapply(y_s_model_p3, mass_predictor)
166  fitted_vec_p4 <- lapply(y_s_model_p4, mass_predictor)
167  fitted_vec_p5 <- lapply(y_s_model_p5, mass_predictor)
```

```r
162   ## creating fitted data with the models, 10 fitted sets per S = 100 datasets, 1000 fitted datasets
163   fitted_vec_p1 <- lapply(y_s_model_p1, mass_predictor)
164   fitted_vec_p2 <- lapply(y_s_model_p2, mass_predictor)
165   fitted_vec_p3 <- lapply(y_s_model_p3, mass_predictor)
166   fitted_vec_p4 <- lapply(y_s_model_p4, mass_predictor)
167   fitted_vec_p5 <- lapply(y_s_model_p5, mass_predictor)
168   fitted_vec_p6 <- lapply(y_s_model_p6, mass_predictor)
169   fitted_vec_p7 <- lapply(y_s_model_p7, mass_predictor)
170   fitted_vec_p8 <- lapply(y_s_model_p8, mass_predictor)
171   fitted_vec_p9 <- lapply(y_s_model_p9, mass_predictor)
172   fitted_vec_p10 <- lapply(y_s_model_p10, mass_predictor)
173
174   fitted_vec <- list(fitted_vec_p1, fitted_vec_p2, fitted_vec_p3, fitted_vec_p4, fitted_vec_p5,
175                      fitted_vec_p6, fitted_vec_p7, fitted_vec_p8, fitted_vec_p9, fitted_vec_p10)
176
177
178   Var_calc <- function(fit) {
179     b <- mapply('-', y_s, fit)
180     RSS <- 0
181     for(i in 1:100) {
182       RSS <- RSS + sqrt(sum(b[,i]^2))^2
183     }
184     return(RSS / 99)
185   }
186
187   Variance_vec <- lapply(fitted_vec, Var_calc)
188   # creating the sample means for the optimal estimated for each p = {1, ... , 10}
189   y_bar <- lapply(fitted_vec, function(x) Reduce("+", x)/ 99)
190
191   true_Values <- generate_Data(x_Data)
192   true_Values_vec <- list(true_Values, true_Values, true_Values, true_Values, true_Values,
193                           true_Values, true_Values, true_Values, true_Values, true_Values)
194
195   bias <- mapply('-', true_Values_vec, y_bar)
196   for(i in 1:10) {
197     bias_squared[i] <- sqrt(sum(bias[,i]^2))^2
198   }
199
200
201   Risk_Model <- mapply('+', Variance_vec, bias_squared)
202
```

13

```r
203  ##  Plots
204
205  plot(x = 1:10, y = Variance_vec,
206       xlab = "p: degree of estimator",
207       ylab = "Variance of Estimator")
208
209  plot(x = 1:10, y = bias_squared,
210       xlab = "p: degree of estimator",
211       ylab = "Bias Squared of Estimator")
212
213  plot(x = 1:10, y = Risk_Model,
214       xlab = "p: degree of estimator",
215       ylab = "Risk of Estimator")
216
217  match(min(Risk_Model), Risk_Model)
218
219  ## complute the RSS for the models p 1-10
220
221  RSS <- function( model) {
222     output <- mapply('-', y_s, model)
223     output1 <- c(1:100)
224     for(i in 1:100) {
225       output[,i] <-output[,i]^2
226       output1[i] <- sum(output[,i])
227     }
228     return(output1)
229  }
230
231  p_RSS <- lapply(fitted_vec, RSS)
232
233  p_star <- c(1:100)
234
235  for(i in 1:100) {
236     vec <- c(1:10)
237     for(j in 1:10) {
238       vec[j] <- p_RSS[[j]][i]
239     }
240     p_star[i] <- match(min(vec), vec)
241  }
242
```

```r
219  ## complute the RSS for the models p 1-10
220
221  RSS <- function( model) {
222     output <- mapply('-', y_s, model)
223     output1 <- c(1:100)
224     for(i in 1:100) {
225        output[,i] <-output[,i]^2
226        output1[i] <- sum(output[,i])
227     }
228     return(output1)
229  }
230
231  p_RSS <- lapply(fitted_vec, RSS)
232
233  p_star <- c(1:100)
234
235  for(i in 1:100) {
236     vec <- c(1:10)
237     for(j in 1:10) {
238        vec[j] <- p_RSS[[j]][i]
239     }
240     p_star[i] <- match(min(vec), vec)
241  }
242
243
244  plot(x = 1:100, y = p_star,
245       ylab = "p* value for lowest RSS",
246       xlab = "S = {1, ... , 100} samples")
247
248
249
250  c_of_p <- c(1:100)
251
252  for(i in 1:100) {
253     vec <- c(1:10)
254     for(j in 1:10) {
255        vec[j] <- p_RSS[[j]][i] + 2* Variance_vec[[j]] * (j +1)
256     }
257     c_of_p[i] <- match(min(vec), vec)
258  }
259
260  plot(x = 1:100, y = c_of_p,
261       ylab = "p-dagger value for lowest Collin Mallows Statistic",
262       xlab = "S = {1, ... , 100} samples")
263
```