

Solution to Homework 6

Date: April 27, 2023

Scribe: Jianliang He

Warning: This note is only used as a reference solution for the homework, and the solution to each question is not unique. The solution may contain factual and/or typographic errors and comments and criticism are kindly welcomed.

Problem 1 Get the Sloan Digital Sky Survey data from <https://www.stat.cmu.edu/~larry/all-of-nonpar/=data/galaxy.dat> Investigate the distribution of redshifts using a histogram and a kernel density estimator. Use least squares cross-validation to choose the amount of smoothing. Also consider the Normal reference rule for picking a bandwidth for the kernel. Plot the graphs of the estimated densities. Carefully read example 4.3 from the textbook.

Solution:

```
1 # Load data
2 data <- read.table
3 ("https://www.stat.cmu.edu/~larry/all-of-nonpar/=data/galaxy.dat")[,3]
4
5 estden.hist <- function(x, data, h) {
6   n <- length(data)
7   range <- max(data) - min(data)
8   block <- floor((x - min(data)) / h)
9   p.hat <- sum(data >= min(data) + range * (block - 1) * h &
10               data < min(data) + range * block * h) / n
11   return(p.hat / h)
12 }
```

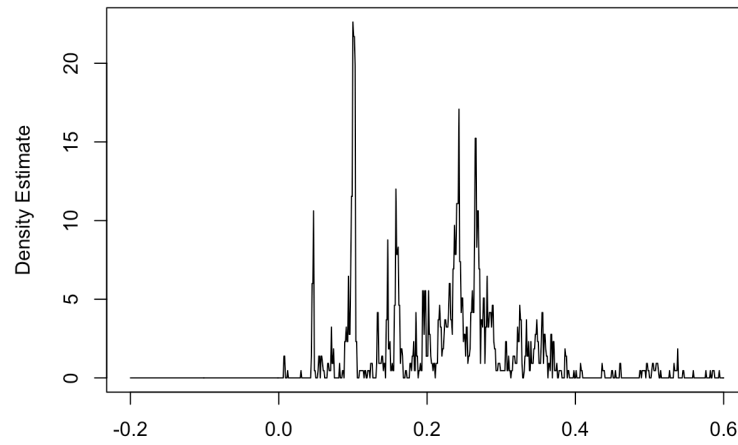
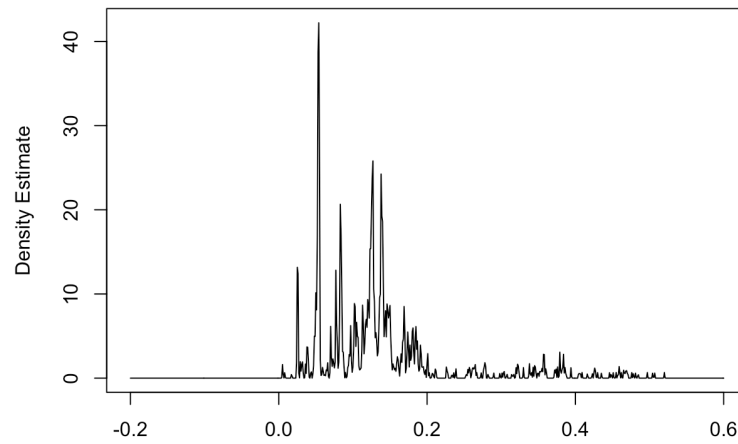
```
13
14 band.hcv <- function(data){
15   h.candidate <- seq(from = 0.00002, to = 0.01, by = 0.0001)
16   n <- length(data)
17   nh <- length(h.candidate)
18   range <- max(data)-min(data)
19   riskseq <- c()
20   for (i in 1:nh){
21     h <- h.candidate[i]
22     p.hat = c()
23     for (j in 1:n){
24       p.hat[j] = sum(data>=min(data)+range*(j-1)*h
25                     & data<min(data)+range*j*h)/n
26     }
27     # Empirical risk (6.16)
28     riskseq[i] <- 2/h/(n-1)-(n+1)/h/(n-1)*sum(p.hat^2)
29   }
30   return(h.candidate[which.min(riskseq)])
31 }
32
33 estden.kernel <- function(x, data, h){
34   n <- length(data)
35   kernel.RBF <- function(x) exp(-(x^2)/2) / sqrt(2 * pi)
36   return(sum(kernel.RBF((x - data)/h)) / (n*h))
37 }
38
39 band.kcv <- function(data){
40   h.candidate <- seq(from = 0.00002, to = 0.01, by = 0.0001)
41   n <- length(data)
42   nh <- length(h.candidate)
```

```
43 riskseq <- c()
44 for (i in 1:nh){
45   h <- h.candidate[i]
46   # Integrated square estimation
47   estdensq.kernel <- function(x) estden.kernel(x, data = data, h = h)^2
48   J1 <- integrate(Vectorize(estdensq.kernel),
49                   lower = -1, upper = 2, subdivisions = 2000)$value
50   # leave-one-out
51   J2 <- 0
52   for (j in 1:n){
53     J2 <- J2 + estden.kernel(data[j], data = data[-j], h = h)
54   }
55   # Empirical risk (6.33)
56   riskseq[i] <- J1 - (2 * J2)/n
57 }
58 return(h.candidate[which.min(riskseq)])
59 }
60
61 band.normal <-function(data){
62   n <- length(data)
63   s <- sd(data)
64   qdata <- as.vector(quantile(data,probs=c(0.25,0.75)))
65   Q <- qdata[2]-qdata[1]
66   sigma <- min(s, Q/1.34)
67   h <- 1.06*sigma*n^(-1/5)
68   return(h)
69 }
70
71 # Density estimation
72 x.grid <- seq(from = -0.2, to = 0.6, by = 0.001)
```

```
73 # Histogram : Cross-Validation
74 h.hcv <- band.hcv(data)
75 y.grid <- sapply(x.grid, estden.hist, data = data, h = h.hcv)
76 plot(x.grid, y.grid, type = "l", xlab = "", ylab = "Density_Estimate",
77      main = "Redshift:_Histogram_density_estimation_with_cross-validation")
78
79 # Kernel: Cross-Validation
80 h.kcv <- band.kcv(data)
81 y.grid <- sapply(x.grid, estden.kernel, data = data, h = h.kcv)
82 plot(x.grid, y.grid, type = "l", xlab = "", ylab = "Density_Estimate",
83      main = "Redshift:_Kernel_density_estimation_with_cross-validation")
84
85 # Kernel: Normal reference rule
86 h.normal <- band.normal(data)
87 y.grid <- sapply(x.grid, estden.kernel, data = data, h = h.normal)
88 plot(x.grid, y.grid, type = "l", xlab = "", ylab = "Density_Estimate",
89      main = "Redshift:_Kernel_density_estimation_with_normal_reference_rule")
```

Detailed plots for each density estimation can be found in next page.

□

Redshift: Histogram density estimation with cross-validation**Redshift: Kernel density estimation with cross-validation****Redshift: Kernel density estimation with normal reference rule**