

732A90: Lab 2

Computational Statistics, Group 6

Chih-Yuan Lin, Sarah Walid Alsaadi, Carles Sans Fuentes, Joshua Burrata

February 10, 2017

Question 1

Exercise 1

In this exercise, the datafile **mortality rate.csv** is used, which contains information about mortality rates of the fruit flies during a certain period.

On the data, an extra column called "lograte" is added to it (being the logarithm of our data column "Rate") and a partition is performed with the code provided by the activity in train and test data.

Then, a function `myMSE()` (with parameters "pars" and "lambda") is implemented to predict the performance of the loess model on our test data given different lambdas (from 0.1 to 40 by 0.1). A plot of the MSE with given different lambdas can be seen below in Figure 1:

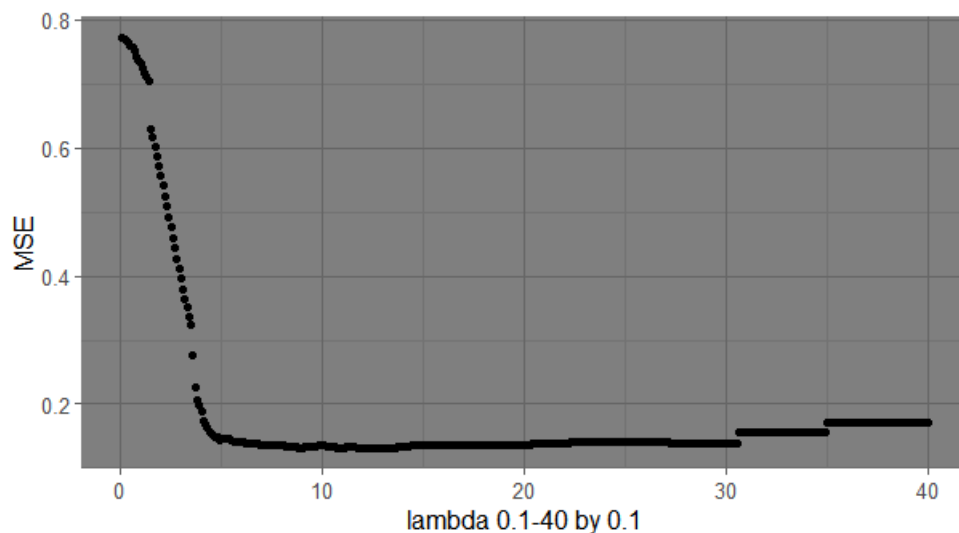


Figure 1: Plot of the MSE of the LOESS model for lambdas from 0.1 to 40 by 0.1

The goal of it is to find the optimal value of lambda for which our model is best. It can be seen that the best lambda used is 11.7 with an MSE of 0.131047. In order to find this solution, the number of iterations needed is the length of the vector lambda (in this case 400 different values for lambda were evaluated).

In order to improve time on optimization, the `optimize()` and the `optim()` formula from R are used.

When using the `optimize()` function for accuracy equals to 0.01, the number of iterations performed are 18 giving the following result:

```
1 $minimum
2 [1] 10.69361
3
4 $objective
5 [1] 0.1321441
```

When using the `optim()` function starting from lambda equals to 35 and using the BFGS method, the number of iterations performed are only 3, getting the following result:

```
1 $par
2 [1] 35
3
4 $value
5 [1] 0.1719996
6
7 $counts
8 function gradient
9      1          1
10
11 $convergence
12 [1] 0
13
14 $message
15 NULL
```

In both cases the function has found some local optimum, but not the best one because the function is a polynomial of more than one degree. For this reason, none of both functions is good to find a local optima if you run it once, because it can get stuck in local optima. For this reason, if this kind of optimization is to be used, then it would be good to try different starting points when optimizing the function for the parameter.

Question 2

In this exercise, the datafile **data.RData**, which contains a sample from the normal distribution with some parameters μ , σ , is used.

We derive the maximum likelihood estimates for these parameters as follows. We start from the probability density function of the normal distribution:

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\sigma^2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (1)$$

So the likelihood function for N observations is:

$$L(\mu, \sigma|x) = \prod_{n=1}^N \frac{1}{\sqrt{2\sigma^2\pi}} e^{-\frac{(x_n-\mu)^2}{2\sigma^2}} \quad (2)$$

Taking the natural logarithm gives us the log-likelihood function:

$$l(\mu, \sigma|x) = -\frac{N}{2} \log(2\sigma^2\pi) - \frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \quad (3)$$

To find the maximum likelihood estimates for each parameter, we first find the partial derivatives of this log-likelihood function with respect to each parameter:

$$\begin{aligned} \frac{\partial l}{\partial \mu} &= -\frac{1}{2\sigma^2} \sum_{n=1}^N \frac{\partial}{\partial \mu} (x_n - \mu)^2 \\ &= -\frac{1}{2\sigma^2} \sum_{n=1}^N (-2)(x_n - \mu) \\ &= \frac{1}{\sigma^2} \left(-N\mu + \sum_{n=1}^N x_n \right) \end{aligned} \quad (4)$$

$$\begin{aligned}
\frac{\partial l}{\partial \sigma} &= -\frac{N}{2} \frac{\partial}{\partial \sigma} (\log(2\sigma^2\pi)) - \sum_{n=1}^N (x_n - \mu)^2 \frac{\partial}{\partial \sigma} \left(\frac{1}{2\sigma^2}\right) \\
&= -\frac{N}{2} \frac{4\sigma\pi}{2\sigma^2\pi} - \left(-\frac{4\sigma}{(2\sigma^2)^2}\right) \sum_{n=1}^N (x_n - \mu)^2 \\
&= -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{n=1}^N (x_n - \mu)^2
\end{aligned} \tag{5}$$

Setting the partial derivatives to zero and re-arranging gives us the following ML estimators for μ and σ :

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \tag{6}$$

$$\sigma_{ML} = \sqrt{\frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2} \tag{7}$$

A function that calculates these maximum likelihood estimators for a normal distribution is implemented. The following results for the data are provided below:

```

1 > loglikelihood(newdata)
2 $mu
3 [1] 1.275528
4
5 $var
6 [1] 4.023942
7
8 $sd
9 [1] 2.005976
10
11 $loglikelihood
12 [1] 211.5069

```

Then, we are asked to optimize the minus loglikelihood function with initial parameters $\mu = 0, \sigma = 1$ using the Conjugate Gradient (CG)) and "BFGS" method with gradient specified and without.

For specifying the gradient, we created a function which takes the same arguments as our minus loglikelihood function and returns the gradient of this function. We know that the gradient of the log-likelihood function here is given by:

$$\nabla l = \left(\frac{\partial l}{\partial \mu}, \frac{\partial l}{\partial \sigma} \right)^T \tag{8}$$

Therefore we can use the partial derivative expressions obtained in equations 4 and 5 and just take the negative as we are optimizing the minus log-likelihood function here.

The usage of the log-likelihood in front of the likelihood is just a matter of convenience: (1) the logarithm is a monotonically increasing function, which means that it achieves its maximum value at the same points as the function itself. (2) having everything in logarithm eases the derivative work when try to optimize and find maximum or minimums. (3) it helps on improving the quality of the calculations computationally speaking reducing the problems that arise when using large numbers.

Here we provide with the results of each data:

```

1 > myframe
2
3 mu          CG with gradient CG without gradient BFGS with gradient BFGS without gradient
4 sigma          2.005976          2.005976          2.005977          2.005977
5 loglikelihood  211.506949          211.506949          211.506949          211.506949
6 counts.function  53.000000          210.000000          38.000000          37.000000
7 counts.gradient  17.000000          35.000000          15.000000          15.000000
8 convergence     0.000000          0.000000          0.000000          0.000000

```

The algorithm converges in all cases giving the same output (as provided above) which is the same as in the case when calculating the maximum likelihood with the expected parameters. It is necessary to mention that the counter for the optim function without gradient does not count well iterations since it excludes those calls needed to compute the Hessian, if requested, and any calls to the function to compute a finite-difference approximation to the gradient. Thus, the real number of times that the function was called has been computed being 97 times for BFGS and 350 times for th CG method.

Nevertheless, it can be seen that when the gradient function is provided, the optimization is performed much faster and the number of iterations is lowered a lot. The fastest model then is the BFGS model, being called 38 times for the loglikelihood function and 15 times for the gradient function.

Contributions

All results and comments presented have been developed and discussed together by the members of the group.

Appendix

Assignment 1

```
1 ### Assignment 2
2
3 ## Exercise 1
4
5 data<- read.csv2("C:/Users/M/Desktop/Statistics and Data Mining Master/Semester 2/Computational
  Statistics/Lecture 2/mortality_rate.csv", sep = ";", stringsAsFactors = FALSE)
6
7 # Adding logarithmic of rate
8
9 lograte<- log(data[,2])
10 #binding data
11 data<- cbind( data, lograte)
12
13
14 # Splitting data
15 n=dim(data)[1]
16 set.seed(123456)
17 id=sample(1:n, floor (n*0.5))
18 train=data[id,]
19 test=data[-id,]
20
21 mylist<- list(X = train[,1], Y = train[,3], Xtest= test[,1], Ytest= test[,3])
22
23 myseq<- seq(0.1,40,0.1)
24 ##1.3 doing the function
25 counter<- 0
26 myMSE <- function(lambda= myseq, pars = mylist){
27   mytrain<- data.frame(X =pars[["X"]],Y = pars[["Y"]])
28
29
30   pred_MSE<- integer(length(lambda)) ##variable to store values
31   counter <- counter+1
32   print(counter)
33   for(count in 1:length(lambda)){
34     model <- loess(formula = Y ~ X, data = mytrain, enp.target = lambda[count])
35     fYpred<- predict(model, pars[["Xtest"]])
36     pred_MSE[count]<- sum((pars[["Ytest"]]-fYpred)**2)/length(fYpred)
37
38   }
39   return(pred_MSE)
40 }
41
42 MSEresults<-myMSE(pars = mylist, lambda= myseq)
43
44 ##best lambda
45 best<- c(best_lambda = myseq[which.min(MSEresults)], MSE = min(MSEresults))
46
47
48 ## 1.4 plot
49 myplot<-function(vector1, vector2){
50   library(ggplot2)
51   ggplot2::ggplot(data = data.frame(X1 =vector1, X2= vector2))+
52     geom_point(aes(x = X1, y= X2))+
53     xlab("lambda 0.1-40 by 0.1")+ ylab("MSE")+
54     theme_dark()}
55 myplot(myseq, MSEresults)
56
57
58 ## 1.5
59
60
61 myoptimize<- optimize(myMSE , interval = c(0.1,40), tol = 0.01, pars = mylist) ##iterations =
  18
62
63 ## 1.6
64
65 myoptim <- optim(35 , myMSE , method = "BFGS", pars = mylist) ##iterations = 3
66
67
68
69 #####2
70
71 load("C:/Users/M/Desktop/Statistics and Data Mining Master/Semester 2/Computational Statistics/
  Lecture 2/data.RData")
72
73 newdata<- data
74
75 ##2.2
76
77 Nloglikelihood<- function(x){
78   n<- length(x)
```

```

79 mu<- sum(x)/n
80 var<- (1/n)*sum((x-mu)**2)
81 sd<- sqrt(var)
82 loglikelihood<- n*log(1/(sqrt(2*var*pi))) - sum((x-mu)**2)/(2*var)
83
84 res<- list(mu = mu, var = var, sd = sd, loglikelihood= loglikelihood)
85
86 return(res)
87
88 }
89
90 Nloglikelihood(newdata)
91
92 ##2.3
93 count<- 0
94 optimloglikelihood<- function(par, x){
95
96   n<- length(newdata)
97   loglikelihood<- -(-n/2*log(2*pi)-n/2*log(par[2]**2) - sum((x-par[1])**2)/(2*par[2]**2))
98   count<- count+1
99   print(count)
100   return(loglikelihood)
101 }
102 }
103
104
105 ###97 iterations
106 BFGS_no_grad<-optim(par = c(0, 1), optimloglikelihood, x= newdata , method = "BFGS")
107 ##350 iterations
108 CG_no_grad<-optim(par = c(0, 1), optimloglikelihood, x= newdata , method = "CG")
109
110
111 ##gradient
112 count<-0
113 mycount<- 0
114 gradient <- function(par, x){
115   n<- length(newdata)
116   gradient_mu<- -((1/(par[2]**2))*(sum(x)-n*par[1]))
117   gradient_sigma <- -(-n/par[2] + (1/par[2]**3)*sum((x-par[1])**2))
118
119   result <- c(gradient_mu = gradient_mu, gradient_sigma= gradient_sigma )
120   mycount<- mycount+1
121   print(mycount)
122   return(result)
123 }
124
125 ##38 iterations, 15 times the gradient function
126 BFGS_grad <-optim(par = c(0, 1), optimloglikelihood, x= newdata, method = "BFGS", gr = gradient
127 )
128 CG_grad <-optim(par = c(0, 1), optimloglikelihood ,x= newdata, method = "CG", gr = gradient)
129
130 ##Creating a data frame for answers
131 myframe <- data.frame(unlist(CG_grad), unlist(CG_no_grad), unlist(BFGS_grad), unlist(BFGS_no_
132 grad))
133 rownames(myframe)<- c("mu", "sigma", "loglikelihood", "counts.function", "counts.gradient", "
134 convergence")
135 colnames(myframe)<- c("CG with gradient", "CG without gradient", "BFGS with gradient", "BFGS
136 without gradient")
137
138 myframe <- rbind(myframe)
139 myframe

```