Lab2 Group11

Group 10 28/01/2020

Question 1: Optimizing a model parameter

The aim of the first question is to perform optimization by using a data set, mortality rate.csv, consisting of information about the mortality rates of fruit flies during an observed period.

1.

First, we import the file to R, and then add a variable called LMR to the data set. The new defined variable LMR is the natural logarithm of Rate. Thereafter, we split the data into a training set and a test sets, respectively. The splitting is done using the code which is already given (see Appendix).

- 2.
- 3.
- 4.
- **5.**
- 6.

Question 2: Maximum likelihood

1.

In this task we will use the file data.RData consists of a sample coming from normal distribution with parameters μ and σ . First we load the data set into R.

2.

The sample comes from a normal distribution with parameters μ and σ , where we set $\theta = (\mu, \sigma)$. Under the assumption that the sample $\boldsymbol{x} = (x_1, ..., x_{100})$ is iid, i.e. $\boldsymbol{X_i} \stackrel{iid}{\sim} N(\mu, \sigma^2)$, for i = 1, ..., 100, then the joint density function of all n = 100 observations can be written as

$$L(\theta; \boldsymbol{x}) = f(\boldsymbol{x}|\theta) = \prod_{i=1}^{100} f(x_i|\theta).$$

Now we let the number of observations be denoted by n in the following derivations. Using the density function of a normal distribution with parameter θ we obtain the likelihood function

$$L(\theta; \boldsymbol{x}) = \prod_{i=1}^{n} \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 \right\}.$$

The log-likelihood function is given by

$$l(\theta; \mathbf{x}) = \log L(\theta; \mathbf{x}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu)^2.$$

The maximum likelihood estimators (MLEs) $\hat{\mu}_{ML}$ and $\hat{\sigma}_{ML}^2$ of μ and σ^2 are obtained by maximizing the likelihood function. This is done by differentiating the log-likelihood functions and put them to zero. In more detail, we calculate the score functions $S(\theta; \boldsymbol{x})$ w.r.t. μ and σ separately, and let them equal zero and solve for each parameter:

$$S(\mu) = \frac{\partial}{\partial \mu} l(\theta; \boldsymbol{x}) = -\frac{n(\overline{x} - \mu)}{\sigma^2} = 0,$$

where $\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$. From this we obtain $\hat{\mu}_{ML} = \overline{x}$. Further,

$$S(\sigma) = \frac{\partial}{\partial \sigma} l(\theta; \boldsymbol{x}) = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0,$$

and
$$\hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$
.

Then we use the derived formulas in order to obtain the desired parameter estimates for the loaded data. So the data set with 100 observations gives the result $\hat{\mu}_{ML} = 1.275528$ and $\hat{\sigma}_{ML} = 2.005976$.

3.

The function optim() minimizes the function by default in R. Thus we will optimize the minus log-likelihood function in order to find the maximum of the function. We will perform two types of algorithms, Conjugate Gradient and BFGS, both with gradient specified and without, to optimize the minus log-likelihood function. It is a better idea to maximize the log-likelihood than maximize the likelihood. This is due to the large values that occurs in the likelihood, which are numerically unstable, so it is preferable to take the logarithm which gives us a better scale to work with. Also, differentiating the log-likelihood function is more computationally convenient, since the product is replaced by a sum (see Question 2.2).

4.

The results of the optimization are presented in Table 1.

Algorithm	Gradient specified	$\hat{\mu}$	$\hat{\sigma}$	Function	Gradient	Time
Conjugate Gradient	No	1.275528	2.005977	297	45	0.01877809 sec
Conjugate Gradient	Yes	1.275528	2.005976	53	17	0.004215956 sec
BFGS	No	1.275528	2.005977	37	15	0.005324841 sec
BFGS	Yes	1.275528	2.005977	39	15	0.009279966 sec

Table 1: Result of the optimization using the algorithms Conjugate Gradient and BFGS, for the given data set. The algorithms, gradient, optimal values of parameters, number of function and gradient evaluations and time taken are presented.

From the results in Table 1, we can see that the algorithm converged in all cases, where all the obtained optimal values are the same and correspond to the estimated parameters in Question 2.2. One explanation of why all algorithms find the optimum is because the likelihood function is convex for the normal distribution. Therefore, it is guaranteed that these algorithms will find the optimal values of the parameters.

Without specifying the gradient, the Conjugate Gradient method required 297 function and 45 gradient evaluations for the algorithm to converge, while the BFGS only required 37 function and 15 gradient evaluations. Also, the time until convergence was measured by using Sys.time(), and we can notice that the BFGS is somewhat faster than Conjugagte Gradient. Even though the difference in time between the algorithms is small, it may have a bigger impact when having a larger data set. All this support the recommendation of choosing the algorithm BFGS in this situation. When we specified the gradient, the number of evaluations was remarkably reduced and time decreased for the Conjugate Gradient algorithm. In this case, it could be reasonable to specify the gradient and not use a finite-difference approximation. On the other hand, there are no big difference when specify the gradient for the BFGS algorithm.

In summary, the gradient adds more information to the optimization and therefore we recommend to specify the gradient if possible. If the gradient is not specified when optimizing, then the BFGS is a better choice in this case.

Appendix

```
knitr::opts chunk$set(echo = TRUE)
# R version
RNGversion('3.5.1')
# 2.1
# Load the data
# Sample from normal distribution with parameters \mu and \sigma
load("data.Rdata")
x <- data
# 2.2
# Derived formulas for the MLE
mu ml <- function(x){</pre>
  mean(x)
sigma_ml <- function(x){</pre>
  term \langle (x - mean(x))^2 \rangle
  sigma <- sum(term)/length(x)</pre>
  return(sqrt(sigma))
}
# Value of MLE
mu \leftarrow mu_ml(x)
sigma <- sigma_ml(x)</pre>
# 2.3
# Minus log-likelihood function
minus_logL <- function(x, par){</pre>
  mu <- par[1]</pre>
  sigma <- par[2]
  (length(x)/2)*log(2*pi*sigma^2) + 1/(2*sigma^2) * sum((x - mu)^2)
}
# Gradient (also minus)
gradient <- function(x, par){</pre>
```

```
mu <- par[1]</pre>
  sigma <- par[2]
  -c((1/sigma)^2 * sum(x - mu),
     -(\operatorname{length}(x)/\operatorname{sigma}) + (1/\operatorname{sigma})^3 * \operatorname{sum}((x - \operatorname{mu})^2))
}
# Optimize with initial parameters mu = 0, sigma = 1.
# Set seed
set.seed(123456)
# Conjugate Gradient method
# Start time
start.time <- Sys.time()</pre>
# With a finite-difference approximation
optim(par = c(0, 1), fn = minus_logL, gr = NULL, method = "CG", x = x)
# End time
end.time <- Sys.time()</pre>
time.taken <- end.time - start.time</pre>
#-----
# Conjugate Gradient method
# Start time
start.time <- Sys.time()</pre>
# With a specified gradient
optim(par = c(0, 1), fn = minus_logL, gr = gradient, method = "CG", x = x)
# End time
end.time <- Sys.time()</pre>
time.taken <- end.time - start.time</pre>
# BFGS - With a finite-difference approximation
# Start time
start.time <- Sys.time()</pre>
# With a finite-difference approximation
optim(par = c(0, 1), fn = minus_logL, gr = NULL, method = "BFGS", x = x)
# End time
end.time <- Sys.time()</pre>
time.taken <- end.time - start.time</pre>
# BFGS - With a specified gradient
# Start time
start.time <- Sys.time()</pre>
# Optimize
optim(par = c(0, 1), fn = minus_logL, gr = gradient, method = "BFGS", x = x)
```

```
# End time
end.time <- Sys.time()
time.taken <- end.time - start.time</pre>
```