

MTH3045: Statistical Computing

Dr. Ben Youngman
b.youngman@exeter.ac.uk
Laver 817; ext. 2314

17/3/2025

Week 10 lecture 1

Week 10 lecture 1

```
y0 <- c(3.52, 1.95, 0.62, 0.02, 5.13, 0.02, 0.01, 0.34, 0.43, 15.5,  
        4.99, 6.01, 0.28, 1.83, 0.14, 0.97, 0.22, 0.02, 1.87, 0.13, 0.01,  
        4.81, 0.37, 8.61, 3.48, 1.81, 37.21, 1.85, 0.04, 2.32, 1.06)
```

```
weib_d1 <- function(pars, y, mult = 1) {  
  # Function to evaluate first derivative of Weibull log-likelihood  
  # pars is a vector  
  # y can be scalar or vector  
  # mult is a scalar defaulting to 1; so -1 returns neg. gradient  
  # returns a vector  
  n <- length(y)  
  z1 <- y / pars[1]  
  z2 <- z1^pars[2]  
  out <- numeric(2)  
  out[1] <- (sum(z2) - n) * pars[2] / pars[1] # derivative w.r.t. lambda  
  out[2] <- n * (1 / pars[2] - log(pars[1])) +  
    sum(log(y)) - sum(z2 * log(z1)) # w.r.t k  
  mult * out  
}
```

```
weib_d2 <- function(pars, y, mult = 1) {  
  # Function to evaluate second derivative of Weibull log-likelihood  
  # pars is a vector  
  # y can be scalar or vector  
  # mult is a scalar defaulting to 1; so -1 returns neg. Hessian  
  # returns a matrix  
  n <- length(y)
```

Newton's method in R I

- We've just put together some simple code that implemented Newton's method
- There are various ways of performing variants of Newton's method in R, but not Newton's method itself
- So here we'll look at `nlminb()`, which is described as 'Unconstrained and box-constrained optimization using PORT routines'
- We can use our functions `weib_d1` and `weib_d2` from earlier for the first and second derivatives of the negative log-likelihood w.r.t. λ and k

Newton's method in R II

- We now just need a function to evaluate the negative log-likelihood itself
- We'll call this `weib_d0`
- Note, though, that it's important to ensure that invalid parameters, i.e. $\lambda \leq 0$ and/or $k \leq 0$, are avoided
- Below we achieve this by setting the log-likelihood to be extremely low (-10^8) for such parts of parameter space

```
weib_d0 <- function(pars, y, mult = 1) {  
  # Function to evaluate Weibull log-likelihood  
  # pars is a vector  
  # y can be scalar or vector  
  # mult is a scalar defaulting to 1; so -1 returns neg. log likelihood  
  # returns a scalar  
  n <- length(y)  
  if (min(pars) <= 0) {  
    out <- -1e8  
  } else {  
    out <- n * (log(pars[2]) - pars[2] * log(pars[1])) +  
      (pars[2] - 1) * sum(log(y)) - sum((y / pars[1])^pars[2])  
  }  
  mult * out  
}
```

Newton's method in R II

```
nlminb(c(1.6, .6), weib_d0, weib_d1, weib_d2, y = y0, mult = -1)
```

```
## $par  
## [1] 1.8900689 0.5375279  
##  
## $objective  
## [1] 54.95316  
##  
## $convergence  
## [1] 0  
##  
## $iterations  
## [1] 5  
##  
## $evaluations  
## function gradient  
##          6          6  
##  
## $message  
## [1] "relative convergence (4)"
```

Newton's method in R III

- We see that `nlminb`'s output is a list comprising...
 - `par`, the parameter estimates
 - `objective`, the final value of the negative log-likelihood
 - `convergence`, whether the algorithm has converged (where 0 indicates successful convergence)
 - `iterations`, the number iterations before convergence was achieved
 - `evaluations`, how many times the function and gradient were evaluated
 - `message` provides further details on the type of convergence achieved

Challenges I

- Go to Challenges I of the week 10 lecture 1 challenges at <https://byoungman.github.io/MTH3045/challenges>

Quasi-Newton methods I

- Between Newton's method and steepest descent lie **quasi-Newton** methods
- These essentially employ Newton's method, but with some approximation to the Hessian matrix
- Instead of the Newton's method search direction

$$\mathbf{p}_i = - [\nabla^2 f(\boldsymbol{\theta}_i)]^{-1} \nabla f(\boldsymbol{\theta}_i)$$

consider the search direction

$$\tilde{\mathbf{p}}_i = -\mathbf{H}_i^{-1} \nabla f(\boldsymbol{\theta}_i)$$

where \mathbf{H}_i is an approximation to the Hessian matrix $\nabla^2 f(\boldsymbol{\theta}_i)$ at the i th iteration

Quasi-Newton methods II

- We might, for example, want to avoid explicitly calculating $\nabla^2 f(\theta_i)$ because it's a matrix that's sufficiently more difficult to calculate than the gradient (e.g. mathematically, or just in terms of time)
 - so that using an approximation to the Hessian matrix (provided it is an adequate approximation) gives a more efficient approach to optimisation than using the Hessian matrix itself
- We should typically expect quasi-Newton methods to converge slower than Newton's method, but provided that convergence isn't too much slower or less reliable, then we may prefer this over analytically forming the Hessian matrix

BFGS (Broyden–Fletcher–Goldfarb–Shanno) I

- In MTH3045 we shall consider the so-called **BFGS** (shorthand for Broyden–Fletcher–Goldfarb–Shanno) quasi-Newton algorithm
- Put simply, at iteration i , the BFGS algorithm assumes that

$$\nabla^2 f(\theta_i) \simeq \mathbf{H}_i = \mathbf{H}_{i-1} + \frac{\mathbf{y}_i \mathbf{y}_i^T}{\mathbf{y}_i^T \mathbf{s}_i} - \frac{(\mathbf{H}_{i-1})^{-1} \mathbf{s}_i \mathbf{s}_i^T (\mathbf{H}_{i-1})^{-T}}{\mathbf{s}_i^T (\mathbf{H}_{i-1})^{-1} \mathbf{s}_i},$$

where $\mathbf{s}_i = \theta_i - \theta_{i-1}$ and $\mathbf{y}_i = \nabla f(\theta_i) - \nabla f(\theta_{i-1})$

- Hence the BFGS algorithm uses differences in the gradients of successive iterations to approximate the Hessian matrix
- We now note that we use \mathbf{H}_i in $\mathbf{p}_i = -[\mathbf{H}_i]^{-1} \nabla f(\theta_i)$
- We can avoid solving this system of linear equations by instead directly obtaining $[\mathbf{H}_i]^{-1}$ through

$$[\mathbf{H}_i]^{-1} = \left(\mathbf{I}_p - \frac{\mathbf{s}_i \mathbf{y}_i^T}{\mathbf{y}_i^T \mathbf{s}_i} \right) [\mathbf{H}_{i-1}]^{-1} \left(\mathbf{I}_p - \frac{\mathbf{y}_i \mathbf{s}_i^T}{\mathbf{s}_i^T \mathbf{y}_i} \right) + \frac{\mathbf{s}_i \mathbf{s}_i^T}{\mathbf{y}_i^T \mathbf{y}_i}$$

BFGS (Broyden–Fletcher–Goldfarb–Shanno) II

- The following R function updates $[\mathbf{H}_{i-1}]^{-1}$ to $[\mathbf{H}_i]^{-1}$ given θ_i , θ_{i-1} , $\nabla f(\theta_{i-1})$ and $\nabla f(\theta_i)$, which are the arguments x0, x1, g0 and g1, respectively

```
iH1 <- function(x0, x1, g0, g1, iH0) {  
  # Function to update Hessian matrix  
  # x0 and x1 are p-vectors of second to last and last estimates, respectively  
  # g0 and g1 are p-vectors of second to last and last gradients, respectively  
  # iH0 is previous estimate of p x p Hessian matrix  
  # returns a p x p matrix  
  s0 <- x1 - x0  
  y0 <- g1 - g0  
  denom <- sum(y0 * s0)  
  I <- diag(rep(1, 2))  
  pre <- I - tcrossprod(s0, y0) / denom  
  post <- I - tcrossprod(y0, s0) / denom  
  last <- tcrossprod(s0) / denom  
  pre %*% iH0 %*% post + last  
}
```

Example: Weibull maximum likelihood: BFGS I

- Repeat Example 5.5 using the BFGS method
- Comment on how it compares to Newton's method

Example: Weibull maximum likelihood: BFGS II

- The following code implements five iterations of the BFGS method

```
## This build of rgl does not include OpenGL functions. Use
## rglwidget() to display results, e.g. via options(rgl.printRglwidget = TRUE)
iterations <- 5
xx <- matrix(0, iterations + 1, 2)
dimnames(xx) <- list(paste('iter', 0:iterations), c('lambda', 'k'))
xx[1, ] <- c(1.6, .6)
g <- iH <- list()
for (i in 2:(iterations + 1)) {
  g[[i]] <- weib_d1(xx[i - 1, ], y0, mult = -1)
  if (sqrt(sum(g[[i]]^2)) < 1e-6)
    break
  if (i == 2) {
    iH[[i]] <- diag(1, 2)
  } else {
    iH[[i]] <- iH1(xx[i - 2, ], xx[i - 1, ], g[[i - 1]], g[[i]], iH[[i - 1]])
  }
  search_dir <- -(iH[[i]] %*% g[[i]])
  alpha <- line_search(xx[i - 1, ], search_dir, weib_d0, y = y0, mult = -1)
  xx[i, ] <- xx[i - 1, ] + alpha * search_dir
}
```

Example: Weibull maximum likelihood: BFGS III

- Our estimates at each iteration are

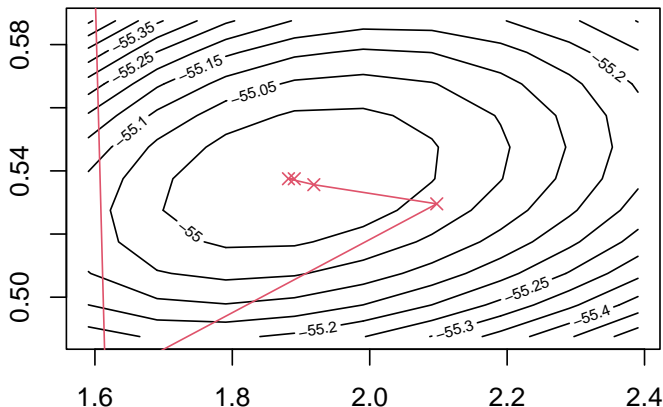
```
xx
```

```
##          lambda          k
## iter 0 1.600000 0.6000000
## iter 1 1.615241 0.4736904
## iter 2 2.097571 0.5295654
## iter 3 1.918661 0.5356343
## iter 4 1.881726 0.5375145
## iter 5 1.890167 0.5374986
```

and we see that we need two more iterations than Newton's method to reach convergence to three decimal places

Example: Weibull maximum likelihood: BFGS IV

- Finally, we'll plot the course of the iterations



and we can see the slightly less direct route we've taken

Quasi-Newton methods in R I

- There are various options for performing quasi-Newton methods in R
- For these, we just need to supply the function to be minimised and its gradient
- The first option is to use `nlminb()` again
 - if we don't supply a function to evaluate the Hessian, then `nlminb()` uses a quasi-Newton approach
- The alternative, and possibly preferred option, is to use `optim()` with option `method = 'BFGS'`
- We'll repeat Example 5.5 using BFGS instead

Quasi-Newton methods in R II

- To use `nlminb()` we can use the following.

```
nlminb(c(1.6, .6), weib_d0, weib_d1, y = y0, mult = -1)
```

```
## $par  
## [1] 1.8900689 0.5375279  
##  
## $objective  
## [1] 54.95316  
##  
## $convergence  
## [1] 0  
##  
## $iterations  
## [1] 7  
##  
## $evaluations  
## function gradient  
##          9          8  
##  
## $message  
## [1] "relative convergence (4)"
```

Quasi-Newton methods in R III

- To use `optim()` we can use the following.

```
optim(c(1.6, .6), weib_d0, weib_d1, y = y0, mult = -1, method = 'BFGS')
```

```
## $par
## [1] 1.8900632 0.5375283
##
## $value
## [1] 54.95316
##
## $counts
## function gradient
##      14      6
##
## $convergence
## [1] 0
##
## $message
## NULL
```

Quasi-Newton methods in R IV

- We note `nlminb()` and `optim()` have essentially given the same value of `par`, i.e. for the $\hat{\lambda}$ and \hat{k} , which is reassuring
- Note that `nlminb()` has used fewer function evaluations than `optim()`
 - we won't go into the details of the cause of this, but it is worth noting that the functions use different stopping criteria, and slightly different variants of the BFGS algorithm
- Note also that `nlminb()` has used three more function evaluations with the BFGS method than with Newton's method
 - this is typically the case, and reflects the improved convergence achieved by using the actual Hessian matrix with Newton's method, as opposed to the approximation that's used with the BFGS approach

Challenges I

- Go to Challenges I of the week 10 lecture 3 challenges at <https://byoungman.github.io/MTH3045/challenges>