

# SDS 385: Homework 3

G. Paulon

September 14, 2016



**Problem 1. Line search**

- (A) Let us recall that, in both the gradient descent and in Newton method, the update of the parameters has the following form:

$$\beta^{(k+1)} = \beta^{(k)} + \alpha^{(k)} s^{(k)},$$

where  $\alpha^{(k)}$  is a step size and  $s^{(k)}$  is a given direction. In general, the direction has the general expression

$$s^{(k)} = -B_k^{-1} \nabla f_k,$$

that is, the direction is chosen to be opposite to the gradient vector, with a correction introduced by the matrix  $B_k^{-1}$ . One can prove that, if  $B_k$  is a positive definite matrix, then the direction  $s^{(k)}$  is a descent direction. Moreover, in the case of the gradient descent  $B_k = \mathcal{I}$ , whereas in the Newton method  $B_k = H_k$ , where  $H_k$  is the Hessian matrix of the objective function evaluated with the parameters at iteration  $k$ .

Given a choice of the descent direction, the goal is now to optimize the choice of the step size  $\alpha^{(k)}$ . Ideally, give a direction  $s^{(k)}$  one would like to minimize the function

$$\phi(\alpha) = f(\beta^{(k)} + \alpha s^{(k)}).$$

However, solving a minimization problem at each iteration of the algorithm is way too computationally intense. A simpler approach consists in trying several values for the step size, discarding those who do not satisfy appropriate conditions that ensure convergence. In fact, a simple decrease in the objective function, i.e.  $f(\beta^{(k)} + \alpha^{(k)} s^{(k)}) < f(\beta^{(k)})$ , is not enough for the whole algorithm to converge. We require a sufficient decrease condition, called the *first Wolfe condition*:

$$f(\beta^{(k)} + \alpha s^{(k)}) \leq f(\beta^{(k)}) + c_1 \alpha \nabla f_k^T s^{(k)}, \quad c_1 \in (0, 1). \quad (1)$$

This expression means that  $f$  should decrease at least linearly, with a slope proportional to both  $\alpha$  and the directional derivative  $\nabla f_k^T s^{(k)}$ . In practice, this condition ensures that the step size is small enough to make the algorithm converge. However, as we discussed in the previous chapters, small step sizes often cause slow performances. Since the first Wolfe condition is always satisfied for small values of  $\alpha$ , it is mandatory to introduce a *second Wolfe condition* which has the role to penalize small step sizes, i.e.

$$\nabla f(\beta^{(k)} + \alpha s^{(k)}) \geq c_2 \nabla f_k^T s^{(k)}, \quad c_2 \in (c_1, 1). \quad (2)$$

A well known procedure that implements line search is known as *backtracking algorithm*. This algorithm allows us to get rid of (2) by implementing only the condition (1) and by choosing the candidates  $\alpha$  in a clever way. In practice, we start from a maximum value of the step size  $\alpha = \bar{\alpha}$  and we proceed backwards by shrinking it by a factor  $\rho$  until condition (1) is met ( $\alpha$  will eventually become small enough that (1) holds). In Algorithm 1 the pseudocode for the algorithm is reported.

**Algorithm 1** Backtracking line search algorithm.

---

```

1: function BACKTRACKING( $\beta, s, y, X, m, c = 0.01, \bar{\alpha} = 3, \rho = 0.5$ )
2:   Set  $\alpha = \bar{\alpha}$ ;
3:   do
4:     Shrink the step size:  $\alpha \leftarrow \rho\alpha$ ;
5:   while  $l(\beta + \alpha s) > l(\beta) + c\alpha \nabla l(\beta)^T s$ 
6:   return  $\alpha$ ;

```

---

(B) The gradient descent method coupled with backtracking has been implemented and tested on the same dataset used in the previous chapters. In Figure 1, the values of the negative log-likelihood are displayed and along with the barplot of the optimal values of the step size. The gradient descent converges, in this case, after around 5947 iterations, which represents a great improvement from the algorithm with fixed step size, which converged in around 60000 iterations.

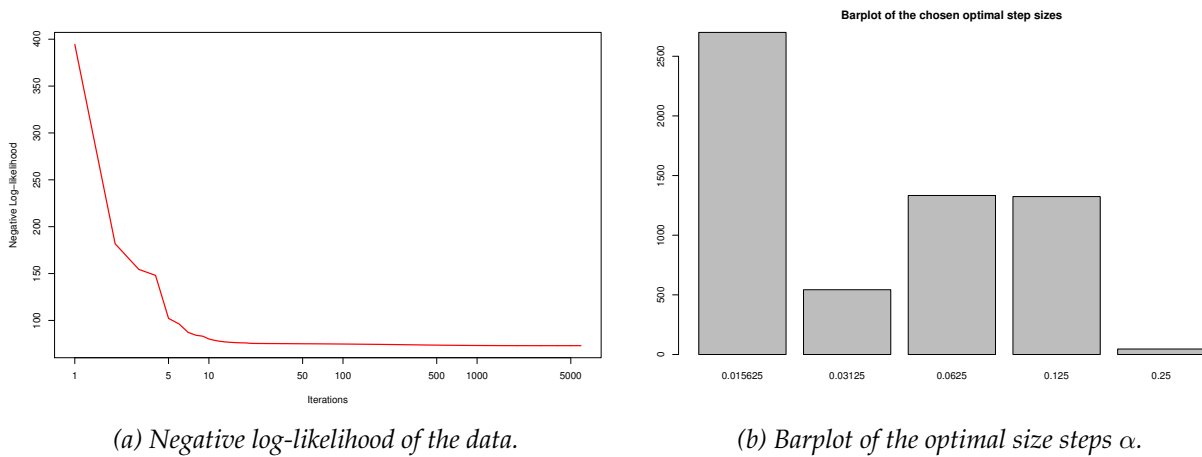


Figure 1: Results of the gradient descent coupled with backtracking.

**Problem 2. Quasi Newton**

(A) Quasi-Newton methods are an alternative to Newton's method in which the Hessian  $\nabla^2 f_k$  is approximated by  $B_k$ , a matrix which is updated at every step. In particular, since the Hessian is the "derivative of the first-order derivative", one can use the finite approximation of the derivative, i.e.

$$\nabla^2 f_k(\beta^{(k+1)} - \beta^{(k)}) \approx \nabla f_{k+1} - \nabla f_k.$$

The *secant condition* is just the equation above applied to the matrix  $B_k$  that should mimic the behaviour of the Hessian. Therefore

$$B_{k+1} s_k = y_k,$$

where  $s_k = \beta^{(k+1)} - \beta^{(k)}$  and  $y_k = \nabla f_{k+1} - \nabla f_k$ . Moreover, we typically impose also that the approximation of the Hessian matrix should be symmetric. The two conditions coupled together yield to a recursive formula to compute  $B_k$ , that is

$$B_{k+1} = B_k - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \frac{y_k y_k^T}{y_k^T s_k}.$$

A practical implementation of BFGS consists in finding an iterative formula to compute directly the inverse of the Hessian, which is

$$B_{k+1}^{-1} = (\mathcal{I} - \rho_k s_k y_k^T) B_k^{-1} (\mathcal{I} - \rho_k y_k s_k^T) + \rho_k s_k s_k^T,$$

where  $\rho_k = \frac{1}{y_k^T s_k}$ . Then, the descent direction  $s^{(k)}$  is simply computed by  $s^{(k)} = -B_k^{-1} \nabla f_k$ .

In Algorithm 2 the pseudocode for the quasi-Newton method coupled with backtracking is reported.

---

**Algorithm 2** BFGS coupled with backtracking line search algorithm.

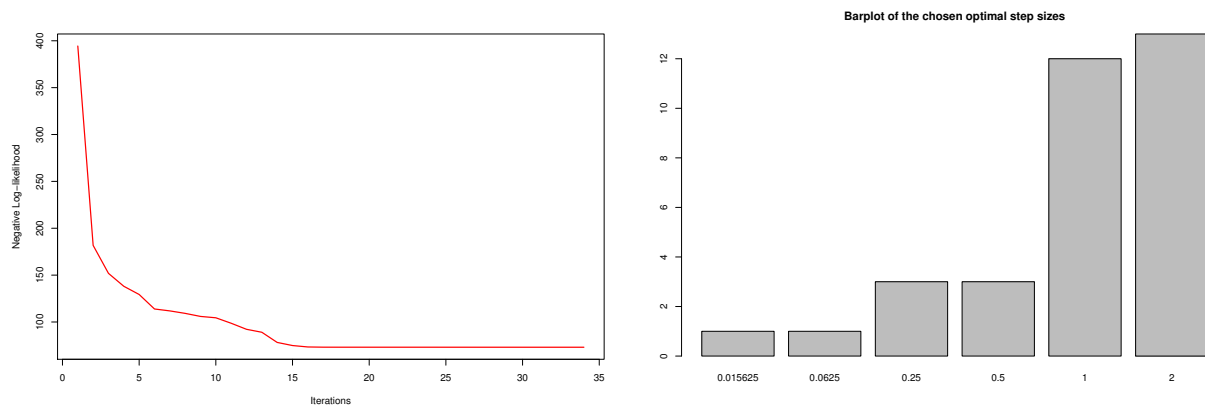
---

```

1: function BFGS( $y, X, m, \beta_0, T$ )
2:   Initialize the data structures that contain  $\beta^{(t)}, l(\beta)^{(t)}, \nabla l(\beta)^{(t)}$ ;
3:   Initialize the Hessian matrix  $B^{-1} = \mathcal{I}$ ;
4:   for  $t$  in  $2 : T$  do
5:     Store the previous gradient:  $gradient_{old} \leftarrow gradient_{new}$ ;
6:     Find the descent direction:  $s^{(t-1)} \leftarrow B^{-1} * gradient_{old}$ ;
7:     Compute the step size:  $\alpha^{(t-1)} \leftarrow \text{BACKTRACKING}(\beta^{(t-1)}, s^{(t-1)}, y, X, m)$ ;
8:     Update the parameters:  $\beta^{(t)} \leftarrow \beta^{(t-1)} + \alpha^{(t-1)} * s^{(t-1)}$ ;
9:     Update the gradient of the log-likelihood:  $gradient_{new} \leftarrow \text{GRAD.LOGLIK}(\beta^{(t)}, y, X, m)$ ;
10:    Compute the quantities:  $s_t \leftarrow \alpha^{(t-1)} s^{(t-1)}$ ;  $y_t \leftarrow gradient_{new} - gradient_{old}$ ;  $\rho_t \leftarrow \frac{1}{y_t^T s_t}$ ;
11:    Update the inverse of the Hessian matrix:  $B^{-1} = (\mathcal{I} - \rho_t s_t y_t^T) B^{-1} (\mathcal{I} - \rho_t y_t s_t^T) + \rho_t s_t s_t^T$ ;
12:  return  $\beta^{(T)}$ ;
```

---

- (B) The BFGS method coupled with backtracking has been implemented and tested on the same dataset used in the previous chapters. In Figure 2, the values of the negative log-likelihood are displayed and along with the barplot of the optimal values of the step size. The gradient descent converges, in this case, after around 34 iterations, which can seem a worse performance than the one obtained with the algorithm with fixed step size, which converged in around 11 iterations. However, in the BFGS version we avoid to compute the Hessian matrix at each iteration, which in high dimensional problems can be cumbersome. The Hessian approximation reduces the per-iteration computational time and it makes the problem more stable (sometimes the Hessian cannot be computed analytically or it could not exist).



(a) Negative log-likelihood of the data.

(b) Barplot of the optimal size steps  $\alpha$ .

Figure 2: Results of the BFGS method coupled with backtracking.

# Appendix A

## R code

```
1 # Compute the probabilities associated with the logit model
2 comp.wi <- function(X, beta){
3   wi <- 1 / (1 + exp(-X %*% beta))
4   return(wi)
5 }
6
7 # Compute the negative log-likelihood of the logit model
8 log.lik <- function(beta, y, X, mi){
9   ll <- array(NA, dim = N)
10  wi <- comp.wi(X, beta)
11  ll <- -sum(y * log(wi + 1E-10) + (mi - y) * log(1 - wi + 1E-10))
12  return(ll)
13 }
14
15 # Compute the gradient of the negative log-likelihood of the logit model
16 grad.loglik <- function(beta, y, X, mi){
17   grad.ll <- array(NA, dim = length(beta))
18   wi <- comp.wi(X, beta)
19   grad.ll <- -apply(X * as.numeric(y - mi * wi), 2, sum)
20   return(grad.ll)
21 }
22
23 # Compute the hessian of the negative log-likelihood of the logit model
24 hessian.loglik <- function(beta, y, X, mi){
25   wi <- comp.wi(X, beta)
26   W.12 <- sqrt(mi * wi * (1 - wi))
27   hes.ll <- crossprod(as.numeric(W.12) * X)
28   return(hes.ll)
29 }
```

*Listing A.1: General functions implementing the computation of log-likelihood, its gradient and its Hessian.*

```
1 # Compute the line search
2 optimal.step <- function(param.act, direct, ll.act, grad.ll.act, y, X, mi, c=0.01, max.alpha=2,
3   rho=0.5){
4   opt.alpha <- max.alpha
5   while(log.lik(param.act + opt.alpha * direct, y, X, mi) > ll.act + c * opt.alpha * crossprod(
6     grad.ll.act, direct)){
7     opt.alpha <- opt.alpha * rho
8   }
9 }
```

```

7   return(opt.alpha)
8 }

```

*Listing A.2: Function implementing the backtracking algorithm for the optimal step size.*

```

1  # Function for the gradient descent of the logit model coupled with backtracking
2  GD.line.search <- function(y, X, mi, beta0, maxiter, tol){
3    betas <- array(NA, dim=c(maxiter, length(beta0)))
4    betas[1, ] <- beta0 # Initial guess
5
6    ll <- array(NA, dim = maxiter)
7    ll[1] <- log.lik(betas[1, ], y, X, mi)
8
9    alpha <- array(NA, dim = maxiter)
10
11   for (iter in 2:maxiter){
12     gradient <- grad.loglik(betas[iter-1, ], y, X, mi)
13     direct <- -gradient
14     alpha[iter-1] <- optimal.step(betas[iter-1, ], direct, ll[iter-1], gradient, y, X, mi)
15     betas[iter, ] <- betas[iter-1, ] + alpha[iter-1] * direct
16     ll[iter] <- log.lik(betas[iter, ], y, X, mi)
17
18     # Convergence check
19     if (abs(ll[iter-1] - ll[iter]) / abs(ll[iter-1] + 1E-3) < tol){
20       cat('Algorithm has converged after', iter, 'iterations')
21       ll <- ll[1:iter]
22       betas <- betas[1:iter, ]
23       break;
24     }
25     else if (iter == maxiter & abs(ll[iter-1] - ll[iter]) / abs(ll[iter-1] + 1E-3) >= tol){
26       print('WARNING: algorithm has not converged')
27       break;
28     }
29   }
30   return(list("ll" = ll, "beta" = betas[iter, ], "alpha" = alpha[1:(iter-1)]))
31 }

```

*Listing A.3: Function implementing the gradient descent using the backtrack algorithm.*

```

1  # Function for the quasi Newton method of the logit model coupled with backtracking
2  quasi.newton <- function(y, X, mi, beta0, maxiter, tol){
3    p <- length(beta0)
4
5    betas <- array(NA, dim=c(maxiter, p))
6    betas[1, ] <- beta0 # Initial guess
7
8    ll <- array(NA, dim = maxiter)
9    ll[1] <- log.lik(betas[1, ], y, X, mi)
10   gradient_new <- grad.loglik(betas[1, ], y, X, mi)
11
12   alpha <- array(NA, dim = maxiter)
13
14   # We initialize the inverse of the Hessian to the identity matrix
15   id.matrix <- diag(rep(1, p))
16   Hk <- id.matrix
17
18   for (iter in 2:maxiter){

```



```

19  # The previous gradient is now the old one
20  gradient <- gradient_new
21
22  # We find the descent direction and we compute the step size
23  direct <- - Hk %*% gradient
24  alpha[iter-1] <- optimal.step(betas[iter-1, ], direct, ll[iter-1], gradient, y, X, mi)
25
26  # We update the values of the parameters, of the log-likelihood and of the gradient
27  sk <- alpha[iter-1] * direct
28  betas[iter,] <- betas[iter-1,] + sk
29  ll[iter] <- log.lik(betas[iter,], y, X, mi)
30  gradient_new <- grad.loglik(betas[iter, ], y, X, mi)
31
32  # We update the inverse of the Hessian matrix
33  yk <- gradient_new - gradient
34  rhok <- as.numeric(1 / crossprod(yk, sk))
35  acc <- tcrossprod(sk, yk)
36  Hk <- (id.matrix - rhok * acc) %*% Hk %*% (id.matrix - rhok * t(acc)) + rhok * sk %*% t(sk)
37
38  # Convergence check
39  if (abs(ll[iter-1] - ll[iter]) / abs(ll[iter-1] + 1E-3) < tol){
40    cat('Algorithm has converged after', iter, 'iterations')
41    ll <- ll[1:iter]
42    betas <- betas[1:iter, ]
43    break;
44  }
45  else if (iter == maxiter & abs(ll[iter-1] - ll[iter]) / abs(ll[iter-1] + 1E-3) >= tol){
46    print('WARNING: algorithm has not converged')
47    break;
48  }
49  }
50  return(list("ll" = ll, "beta" = betas[iter, ], "alpha" = alpha[1:(iter-1)]))
51  }

```

*Listing A.4: Function implementing the quasi-Newton method (BFGS) coupled with a backtracking algorithm.*