

# EECS 227C - Homework III

Walid Krichene

March 20, 2014

1 Consider the two dimensional optimization problem

$$\text{minimize } \frac{1}{2n} \sum_{k=1}^n \left( x_1 \cos \frac{k\pi}{n} + x_2 \sin \frac{k\pi}{n} \right)^2$$

(a) Prove this is a strongly convex quadratic cost function, and compute the Lipschitz constant of the gradient and strong convexity parameters. Compute the optimal point  $x_{opt}$ .

We have

$$f(x) = \frac{1}{2} x^T Q x$$

where

$$Q = \frac{1}{n} \sum_{k=1}^n \begin{pmatrix} \cos^2 \frac{k\pi}{n} & \cos \frac{k\pi}{n} \sin \frac{k\pi}{n} \\ \cos \frac{k\pi}{n} \sin \frac{k\pi}{n} & \sin^2 \frac{k\pi}{n} \end{pmatrix}$$

the characteristic polynomial of  $Q$  is

$$\begin{aligned} \mathcal{X}(\lambda) &= \det(\lambda I - Q) \\ &= \left( \lambda - \frac{1}{n} \sum_k \cos^2 \frac{k\pi}{n} \right) \left( \lambda - \frac{1}{n} \sum_k \sin^2 \frac{k\pi}{n} \right) - \frac{1}{n^2} \left( \sum_k \cos \frac{k\pi}{n} \sin \frac{k\pi}{n} \right)^2 \\ &= \lambda^2 - \frac{1}{n} \left( \sum_k \cos^2 \frac{k\pi}{n} + \sum_k \sin^2 \frac{k\pi}{n} \right) \lambda + \frac{1}{n^2} \left( \sum_k \cos^2 \frac{k\pi}{n} \right) \left( \sum_k \sin^2 \frac{k\pi}{n} \right) - \frac{1}{n^2} \left( \sum_k \cos \frac{k\pi}{n} \sin \frac{k\pi}{n} \right)^2 \\ &= \lambda^2 - \lambda + a \end{aligned}$$

where

$$\begin{aligned} a &= \frac{1}{n^2} \left( \sum_k \cos^2 \frac{k\pi}{n} \right) \left( \sum_k \sin^2 \frac{k\pi}{n} \right) - \frac{1}{n^2} \left( \sum_k \cos \frac{k\pi}{n} \sin \frac{k\pi}{n} \right)^2 \\ &= \frac{1}{n^2} \left( \frac{1}{2} \sum_k 1 + \cos 2 \frac{k\pi}{n} \right) \left( \frac{1}{2} \sum_k 1 - \cos 2 \frac{k\pi}{n} \right) - \frac{1}{n^2} \left( \sum_k \sin 2 \frac{k\pi}{n} \right)^2 \\ &= \frac{1}{n^2} \left( \frac{n}{2} \right) \left( \frac{n}{2} \right) - 0 \\ &= \frac{1}{4} \end{aligned}$$

using the fact that  $\sum_{k=1}^n \sin 2 \frac{k\pi}{n} = \sum_{k=1}^n \cos 2 \frac{k\pi}{n} = 0$  (for example,  $\sum_{k=1}^n \cos 2 \frac{k\pi}{n} = \text{Re}(\sum_{k=1}^n e^{i 2 \frac{k\pi}{n}}) = \text{Re}(e^{i \frac{2\pi}{n}} \frac{1 - e^{i \frac{2n\pi}{n}}}{1 - e^{i \frac{2\pi}{n}}}) = 0$ ).

Thus

$$\mathcal{X}(\lambda) = \lambda^2 - \lambda + \frac{1}{4} = \left(\lambda - \frac{1}{2}\right)^2$$

therefore  $Q$  is a symmetric matrix (thus diagonalizable) with eigenvalues  $\frac{1}{2}$ , i.e.  $Q = \frac{1}{2}I$ , and  $f(x) = \frac{1}{4}\|x\|^2$ . It follows that  $f$  is strictly convex with parameter  $\ell = \frac{1}{4}$ , has Lipschitz gradient with Lipschitz constant  $\frac{1}{4}$ , and the minimizer is  $x_{opt} = 0$ .

- (b) Consider the stochastic gradient method with step size 1 initialized at the point  $(0, 1)$ . Compute

$$\mathbb{E} \|x^n - x_{opt}\|^2$$

when the increments are chosen with replacement.

At step  $n$ , if we draw  $k$ , then the stochastic gradient is

$$\begin{aligned} g^n &= \nabla_x \frac{1}{2} \left( \cos \frac{k\pi}{n} x_1 + \sin \frac{k\pi}{n} x_2 \right)^2 \\ &= \left( \cos \frac{k\pi}{n} x_1 + \sin \frac{k\pi}{n} x_2 \right) \begin{pmatrix} \cos \frac{k\pi}{n} \\ \sin \frac{k\pi}{n} \end{pmatrix} \\ &= z_k z_k^T x^n \end{aligned}$$

where  $z_k = \begin{pmatrix} \cos \frac{k\pi}{n} \\ \sin \frac{k\pi}{n} \end{pmatrix}$ . It is easy to check that

$$\begin{aligned} \mathbb{E}[g^n | x^n] &= \mathbb{E}[z_k z_k^T x^n | x^n] = \left( \frac{1}{n} \sum_k z_k z_k^T \right) x^n \\ &= Q x^n \end{aligned}$$

which is the gradient of  $\frac{1}{2}x^T Q x$ . Now we have

$$x^{n+1} = x^n - g^n = (I - z_k z_k^T) x^n$$

therefore we have

$$\begin{aligned} \mathbb{E} \|x^{n+1}\|^2 &= \mathbb{E}[\mathbb{E}[\|x^{n+1}\|^2 | x_n]] \\ &= \mathbb{E}[\mathbb{E}[(x^n)^T (I - z_k z_k^T)^2 x^n | x_n]] \\ &= \mathbb{E}[(x^n)^T \left( \frac{1}{n} \sum_k (I - z_k z_k^T)^2 \right) x^n] \end{aligned}$$

where  $(I - z_k z_k^T)^2 = I - 2z_k z_k^T + z_k z_k^T z_k z_k^T = I - z_k z_k^T$  since  $\|z_k\|^2 = 1$ . So  $\frac{1}{n} \sum_k (I - z_k z_k^T)^2 = \frac{1}{n} \sum_k (I - z_k z_k^T) = I - \frac{1}{n} \sum_k z_k z_k^T = I - Q = \frac{1}{2}I$ . Therefore

$$\mathbb{E} \|x^{n+1}\|^2 = \frac{1}{2} \mathbb{E} \|x^n\|^2$$

and by induction

$$\mathbb{E} \|x^n\|^2 = \frac{1}{2^n} \|x^0\|^2 = \frac{1}{2^n}$$

- (c) Consider the incremental gradient method with step size 1 initialized at the point  $(0, 1)$ . Compute

$$\|x_n - x_{opt}\|^2$$

when the increments are chosen in order, deterministically. That is, we first take a step with respect to the summand  $k = 1$ , then  $k = 2$ , all the way to  $k = n$ .

In this case, writing the consecutive updates yields

$$x^{n+1} = (I - z_n z_n^T)(I - z_{n-1} z_{n-1}^T) \dots (I - z_1 z_1^T) x^n$$

where each update  $I - z_k z_k^T$  corresponds to a projection on the orthogonal of  $z_k$ . This can be verified by writing

$$\begin{aligned} I - z_k z_k^T &= I - \begin{pmatrix} \cos^2 \frac{k\pi}{n} & \sin \frac{k\pi}{n} \cos \frac{k\pi}{n} \\ \sin \frac{k\pi}{n} \cos \frac{k\pi}{n} & \sin^2 \frac{k\pi}{n} \end{pmatrix} \\ &= \begin{pmatrix} \sin^2 \frac{k\pi}{n} & -\sin \frac{k\pi}{n} \cos \frac{k\pi}{n} \\ -\sin \frac{k\pi}{n} \cos \frac{k\pi}{n} & \cos^2 \frac{k\pi}{n} \end{pmatrix} \\ &= \begin{pmatrix} \cos^2 \frac{k\pi}{n} + \frac{\pi}{2} & \cos \frac{k\pi}{n} + \frac{\pi}{2} \sin \frac{k\pi}{n} \\ \cos \frac{k\pi}{n} + \frac{\pi}{2} \sin \frac{k\pi}{n} & \sin^2 \frac{k\pi}{n} + \frac{\pi}{2} \end{pmatrix} \\ &= \bar{z}_k \bar{z}_k^T \end{aligned}$$

where  $\bar{z}_k = \begin{pmatrix} \cos \frac{k\pi}{n} + \frac{\pi}{2} \\ \sin \frac{k\pi}{n} + \frac{\pi}{2} \end{pmatrix}$ .

Now

$$x^{m+1} = \bar{z}_n \bar{z}_n^T \dots \bar{z}_1 \bar{z}_1^T x_m$$

where for all  $k$ ,  $\bar{z}_{k+1}^T \bar{z}_k$  is the dot product of two unit vectors that form a  $\frac{\pi}{n}$  angle, i.e.  $\bar{z}_{k+1}^T \bar{z}_k = \cos \frac{\pi}{n}$ , and

$$x^{m+1} = \left(\cos \frac{\pi}{n}\right)^{n-1} \bar{z}_n \bar{z}_1^T x_m$$

and  $\bar{z}_n = \begin{pmatrix} 0 \\ 1 \end{pmatrix} = x_0$ . So  $x^m = \|x^m\| x_0$ , and

$$x^{m+1} = \|x^m\| \left(\cos \frac{\pi}{n}\right)^{n-1} (\bar{z}_1^T x_0) x_0 = \|x^m\| \left(\cos \frac{\pi}{n}\right)^n x_0$$

so  $\|x^{m+1}\| = \|x^m\| \left(\cos \frac{\pi}{n}\right)^n$  thus by induction

$$\|x^m\| = \left(\cos \frac{\pi}{n}\right)^{mn}$$

**2** In logistic regression, one is given a data set of pairs  $(x_i, y_i)$  where  $x_i \in \mathbb{R}^d$  and  $y_i \in \{0, 1\}$ . The goal is to find a vector  $w$  and scalar  $b$  such that

$$\sum_{i=1}^n -y_i(b + w^T x_i) + \log(1 + \exp(b + w^T x_i)) + \lambda \|w\|_2^2 \quad (1)$$

is minimized.

Program BFGS, LBFGS, and the Barzilai-Borwein methods. Test these methods using the function(1) on the adult data set with  $\lambda = 10^{-3}$ .

Use Weak Wolfe Line Search for the line search in BFGS and LBFGS. Use the stopping criterion

$$\|\nabla f(x)\|^2 \leq \text{toler} * (1 + |f(x)|)$$

Describe the optimal parameter settings for each algorithms to achieve the smallest number of function and gradient evaluations. Devise a methodology for getting the lowest possible computation time with some hybrid of these three methods. Submit your code as printout with comments.

## BFGS

The BFGS iteration can be summarized as

- $p_k = -H_k \nabla f(x_k)$
- $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is chosen to satisfy the Wolfe condition: for some  $0 < c_1 < c_2 < 1$ :

$$f(x_{k+1}) \leq f(x_k) + c_1 \alpha_k \nabla f(x_k)^T p_k \quad (\text{Armijo})$$

$$\nabla f(x_{k+1})^T p_k \geq c_2 \nabla f(x_k)^T p_k \quad (\text{Wolfe})$$

- $s_k = x_{k+1} - x_k$
- $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
- $H_{k+1} = (I - \rho_k s_k y_k^T)^T H_k (I - \rho_k s_k y_k^T) + \rho_k s_k s_k^T$  where  $\rho_k = \frac{1}{y_k^T s_k}$

The Weak Wolfe line search can be implemented as:  $t = 1$ ,  $\alpha = 0$ ,  $\beta = \infty$

- if Armijo condition fails,  $\beta = t$ ,  $t = \frac{1}{2}(\alpha + \beta)$
- else if Wolfe condition fails,  $\alpha = t$ ,  $t = 2\alpha$  if  $\beta = \infty$  and  $\frac{1}{2}(\alpha + \beta)$  otherwise.

The Matlab implementation is given below

```
1 function [fxp, xp]=BFGS(f, g, x0)
2
3 global toler;
4 global nbFcalls maxFcalls;
5
6 n = size(x0, 1);
7 I = eye(n);
8 nbFcalls = 0;
9
10 % initialize H and x
11 H = I;
12 x = x0;
13 % iterate
14 while (1)
15     % compute the descent direction
16     gx = g(x);
17     p = -H*gx;
18
19     % line search for alpha to satisfy the weak Wolfe conditions
20     alpha = WW(f, g, x, p);
21
22     % compute next iterate
23     xp = x + alpha*p;
24
25     % update the H matrix
26     s = xp - x;
27     gxp = g(xp);
28     y = gxp - gx;
29     rho = 1/(y'*s);
30     Hu = (I - rho*s*y');
31     H = Hu'*H*Hu + rho*(s*s');
32
33     % compute the function value and update x
34     fxp = f(xp)
35     x = xp;
36
37     % check stopping criterion
38     e = norm(gxp)^2 - toler*(1+abs(fxp))
39     if(e < 0 || (nbFcalls > maxFcalls))
40         break
41     end
42 end
```

## Weak Wolfe line search

```
1 % f is the function handle, g the gradient handle, x the current point, d
2 % the search direction
3 function t = WW(f, g, x, p) %% WW for weak Wolfe
4
5 global c1 c2;
6 global fs;
7 t0 = 1;
8
9
10 beta = Inf;
11 alpha = 0;
12 t = t0;
13 found = 0;
14
15 % f(x)
16 ff = f(x);
17 % g(x)
18 gg = g(x);
19 while(~found)
20     % next point
21     xp = x+t*p;
22     ffp = f(xp);
23     fs = [fs min(ff, ffp)]; % one gradient or function call per iteration of the line search
24     %display({ff, ffp})
25     if(ffp > ff + t*c1*gg'*p) % if the Armijo condition fails
26         beta = t;
27         t = (alpha + beta)/2;
28     elseif (g(xp) '*p < c2*gg'*p) % if the Wolfe condition fails
29         alpha = t;
30         if(beta == Inf)
31             t = 2*alpha;
32         else
33             t = (alpha+beta)/2;
34         end
35     else
36         found = 1; % found a t which satisfies both conditions
37     end
38 end
39 t % return t
```

## LBFGS

The LBFGS iteration can be summarized as: instead of maintaining (and storing) the full matrix  $H_k$ , choose a horizon  $\tau$ , then

- $q = \nabla f(x_k)$
- for  $i \in \{k-1, \dots, k-\tau\}$ ,  $\alpha_i = \frac{s_i^T q}{s_i^T y_i}$ ,  $q = q - \alpha_i y_i$
- $r = \lambda_k q$  where  $\lambda_k = \frac{s_{k-1}^T s_{k-1}}{y_{k-1}^T s_{k-1}}$
- for  $i \in \{k-\tau, \dots, k-1\}$ ,  $\beta_i = \frac{y_i^T r}{y_i^T s_i}$ ,  $r = r + s_i(\alpha_i - \beta_i)$
- use  $r$  as descent direction.

The matlab implementation is given below

```

1 function [fxp, xp]=LBFGS(f, g, x0)
2
3 global toler;
4 global nbFcalls maxFcalls;
5 global tau;
6 n = size(x0, 1);
7
8 nbFcalls = 0;
9
10 % initialize the memory
11 %tau = 5; % remember the last 5 iterates
12 ss = zeros(n, tau);
13 ys = zeros(n, tau);
14 alphas = zeros(1, tau);
15 betas = zeros(1, tau);
16
17 % initialize x and lambda
18 x = x0;
19 lambda = 1;
20 k = 0;
21 % iterate
22 while (1)
23     gx = g(x);
24     % compute the descent direction
25     q = gx;
26     T = min(tau, k);
27     for i=1:T;
28         si = ss(:, i);
29         yi = ys(:, i);
30         alphas(i) = si'*q/(si'*yi);
31         q = q - alphas(i)*yi;
32     end
33     r = lambda*q;
34     for ii=1:T
35         i = T-ii+1;
36         yi = ys(:, i);
37         si = ss(:, i);
38         betas(i) = yi'*r/(yi'*si);
39         r = r+si*(alphas(i) - betas(i));
40     end
41
42     % -r is the descent direction
43
44     % line search for alpha to satisfy the weak Wolfe conditions
45     alpha = WW(f, g, x, -r);
46
47     % compute next iterate
48     xp = x - alpha*r;
49
50     % update the memory
51     s = xp - x;
52     gxp = g(xp);
53     y = gxp - gx;
54
55     lambda = s'*s / (y'*s);
56
57     % shift all tables
58     for i=0:tau-2
59         ss(:, tau-i) = ss(:, tau-i-1);
60         ys(:, tau-i) = ys(:, tau-i-1);
61     end
62     ss(:, 1) = s;
63     ys(:, 1) = y;
64

```

```

65     % compute the function value and update x
66     fxp = f(xp)
67     x = xp;
68
69     % check stopping criterion
70     e = norm(gxp)^2 - toler*(1+abs(fxp))
71     if(e < 0 || (nbFcalls > maxFcalls))
72         break
73     end
74 end

```

## BB

The Barzilai-Borwein method is simply

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

where  $\alpha_k = \frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}}$ .  $\alpha_0$  is a parameter.

The Matlab implementation is given below

```

1  function [fxp, xp]=BB(f, g, x0)
2
3  global nbFcalls maxFcalls;
4
5  global toler;
6  global fs;
7  global alpha0;
8  nbFcalls = 0;
9  fs = [fs, f(x0)];
10
11 % initialize alpha and x
12 alpha = alpha0;
13 x = x0;
14 % iterate
15 while (1)
16     % compute the descent direction
17     p = -g(x);
18
19     % compute next iterate
20     s = alpha*p;
21     xp = x + s;
22
23     % update alpha
24     gxp = g(xp);
25     y = gxp + p;
26     alpha = s'*y/(y'*y);
27
28     % compute the function value and update x
29     fxp = f(xp)
30     fs = [fs fxp];
31     x = xp;
32
33     % check stopping criterion
34     e = norm(gxp)^2 - toler*(1+abs(fxp))
35     if(e < 0 || (nbFcalls > maxFcalls))
36         break
37     end
38 end

```



## Expression of the gradient

by augmenting all vectors  $X_i$  with  $X_{i,0} = 1$ , and writing  $w_0 = b$ , we have

$$f(x) = \sum_{i=1}^n -y_i w^T X_i + \log(1 + e^{w^T X_i}) + \lambda \|w\|^2$$

so

$$\begin{aligned} \frac{\partial f}{\partial w_j} &= \sum_{i=1}^n -y_i X_{i,j} + X_{i,j} \frac{e^{w^T X_i}}{1 + e^{w^T X_i}} + 2\lambda w_j \\ &= 2\lambda w_j + \sum_{i=1}^n X_{i,j} \left( \frac{e^{w^T X_i}}{1 + e^{w^T X_i}} - y_i \right) \end{aligned}$$

so

$$\nabla_w f = 2\lambda w + X^T z$$

where  $z_i = \frac{e^{w^T X_i}}{1 + e^{w^T X_i}} - y_i$

The Matlab implementation of the loss function and the gradient function is given below

```
1 function l = loss(w)
2
3 lambda = .001;
4 global X;
5 global y;
6 global nbFcalls;
7
8 nbFcalls = nbFcalls+1;
9 wo = w;
10 wo(1) = 0;
11 l = lambda*norm(wo)^2;
12 n = size(X, 2);
13 for i=1:n
14     wXi = w'*X(:, i);
15     l = l - y(i)*wXi + log(1+exp(wXi));
16 end
```

```
1 function g = grad(w)
2
3 lambda = .001;
4 global X y nbGcalls;
5
6 nbGcalls = nbGcalls+1;
7
8 wo = w;
9 wo(1) = 0;
10 n = size(X, 2);
11 z = zeros(n, 1);
12 for i=1:n
13     Xi = X(:, i);
14     wXi = w'*Xi;
15     z(i) = exp(wXi)/(1+exp(wXi)) - y(i);
16 end
17
18 g = 2*lambda*wo + X*z;
```

## Results

To find the optimal parameter setting for each method, we run, for different values of the parameters, the algorithm for a fixed number of steps (a step being one function or one gradient evaluation).

The code is given at the end.

Optimal parameters:

- BFGS:  $c_1 = .4$ ,  $c_2 = .5$
- LBFGS:  $c_1 = .4$ ,  $c_2 = .5$ , memory:  $\tau = 4$ .
- Barzalai-Borwein:  $\alpha_0 = 10^{-4}$

The figure below shows the losses as a function of step number for the optimal parameter choice for each method. For `toler` =  $10^{-3}$ , we have the following convergence (number of function evaluations)

- BFGS: after 506 iterations,  $f(x_{\text{best}}) = 1.0515e + 04$
- LBFGS: after 514 iterations,  $f(x_{\text{best}}) = 1.0512e + 04$
- BB: 269 iterations,  $f(x_{\text{best}}) = 1.0510e + 04$

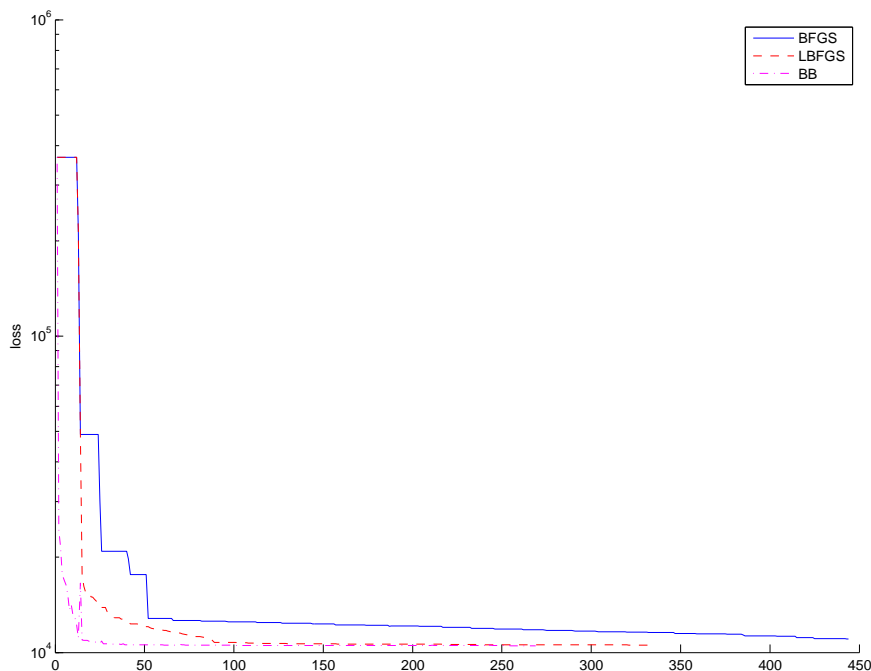


Figure 1: Loss as a function of iteration number for the optimal parameter choice. BB is not a descent method (no line search).

## Combining the algorithms

It seems that the BB method is fast to converge within a ball of the optimum, and then yields slower convergence. We can use it first, and when we get closer to the optimum, switch to either BFGS or LBFGS to have a better approximation of the function around the equilibrium.

An algorithm combining both is given below, where the switch happens at  $\|\nabla f(x)\|^2 \leq \text{toler} * (1 + |f(x)|)$ , then we decrease `toler` to  $10^{-4}$  and run LBFGS.

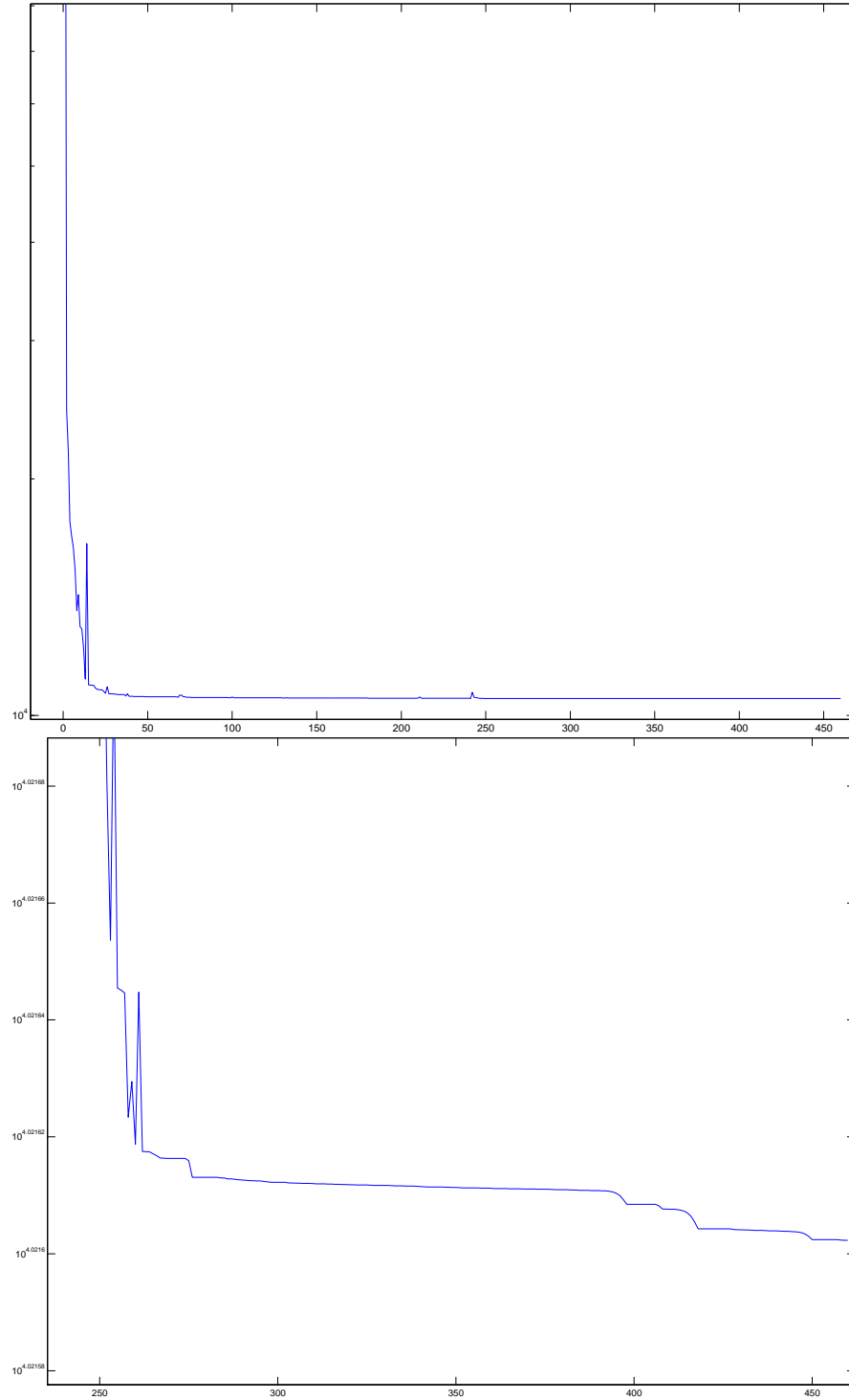


Figure 2: Loss of the combined method, and zoomed in plot around the point where the switch happens. After the switch, the decrease becomes slower, but we do have a marginal improvement by using LBFGS after BB.

```

1 n = size(X0, 2);
2 p = size(X0, 1);
3
4 % training set
5 global X y;
6 X = [ones(1, n); X0]; % augment the training samples with a 1 on their first entry
7 y = y;
8
9 % parameters of the algorithms
10 global c1 c2 toler tau alpha0;
11 toler = .001;
12 % Weak Wolfe line search parameters: c1, c2
13 % LBFGS parameter: tau
14 % BB parameter: alpha0
15 % number of function calls and gradient calls
16 global nbFcalls nbGcalls maxFcalls;
17 nbFcalls = 0;
18 nbGcalls = 0;
19
20 % sequence of functions values f(x.k)
21 global fs;
22
23 % initial point
24 w0 = ones(p+1, 1);
25
26 %%-----
27 %% compare performance of BFGS with different parameters
28 maxFcalls = 200;
29 fBest = Inf;
30 fcs = [];
31
32 c1s = [.1, .2, .3, .4, .5];
33 c2s = [.4, .5, .6, .7, .8];
34 for c1 = c1s
35     for c2 = c2s
36         if(c1 < c2)
37             [fx, x] = BFGS(@loss, @grad, w0);
38             fcs = [fcs fx];
39             if (fx < fBest)
40                 fBest = fx;
41                 cBest = [c1, c2];
42             end
43         end
44     end
45 end
46
47 %%-----
48 %% compare performance of LBFGS with different parameters
49 maxFcalls = 400;
50 fBest = Inf;
51 fcs = [];
52
53 c1s = [.2, .4];
54 c2s = [.6, .8];
55 taus = [4, 8];
56 for c1 = c1s
57     for c2 = c2s
58         for tau=taus
59             [fx, x] = LBFGS(@loss, @grad, w0);
60             fcs = [fcs fx];
61             if (fx < fBest)
62                 fBest = fx;
63                 cBest = [c1, c2, tau];
64             end

```

```

65     end
66   end
67 end
68
69 %%-----
70 %% compare performance of BB with different parameters
71 maxFcalls = 400;
72 fBest = Inf;
73 w0 = ones(p+1, 1);
74
75 alpha0s = [.0001, .000001];
76 for alpha0 = alpha0s
77     fx = BB(@loss, @grad, w0);
78     if (fx < fBest)
79         fBest = fx;
80         cBest = [alpha0];
81     end
82 end
83
84 %%-----
85 %% Compare performance of the different methods with optimal parameters
86 maxFcalls = Inf;
87 toler = .001;
88 w0 = ones(p+1, 1);
89
90 figure(1)
91 hold on
92
93 % BFGS
94 fs = [];
95 c1 = .4;
96 c2 = .5;
97 [fx, x] = BFGS(@loss, @grad, w0)
98 nbFcalls
99 plot(fs, 'b')
100
101 % LBFGS
102 fs = [];
103 c1 = .4;
104 c2 = .5;
105 tau = 4;
106 [fx, x] = LBFGS(@loss, @grad, w0)
107 nbFcalls
108 plot(fs, 'r--')
109
110 % BB
111 fs = [];
112 alpha0 = .0001;
113 [fx, x] = BB(@loss, @grad, w0)
114 nbFcalls
115 plot(fs, 'm-.')
116
117 ylabel('loss')
118 legend({'BFGS', 'LBFGS', 'BB'})
119
120 %%-----
121 %% Combined
122 maxFcalls = Inf;
123 toler = .001;
124 w0 = ones(p+1, 1);
125
126 figure(1)
127 hold on
128
129 % BB

```

```
130 fs = [];  
131 [fx, x] = BB(@loss, @grad, w0)  
132  
133 % then LBFGS  
134 toler = .0001;  
135 [fx, x] = LBFGS(@loss, @grad, x)  
136  
137 plot(fs, 'b')  
138 ylabel('loss')  
139  
140 %%  
141 set(gca, 'Yscale', 'log');  
142 savefig('losses_combined.zoom', 'pdf')
```