

# Optimization Algorithms: Theory and Software Implementation

Prof. Thirumulanathan D

Department of Mathematics

Institute of IIT Kanpur

Lecture: 26

Hello everyone. This is the first lecture of the sixth week. We will be starting with quasi-Newton methods this week.

Until now, we have covered line search methods, both exact and backtracking line searches. To find the descent directions, we have discussed the gradient descent method, the popular Fletcher-Reeves conjugate gradient method, and the eigenvalue/eigenvector method. Last week, we discussed Newton's methods and variants like the modified Newton method and the damped Newton method.

We were able to rectify most of the issues that Newton's methods faced, except for one particular issue: the inversion of the Hessian matrix. The inversion of the Hessian matrix is computationally very expensive; it takes a lot of time if  $n$  is large. It is not scalable. That is the reason we move to quasi-Newton methods, which is what we will be learning this week.

You can see that variants of Newton's methods are named along with Newton itself. We saw damped Newton methods and modified Newton method. Now we have a few different quasi-Newton methods. We will see the popular quasi-Newton methods this week.

Where do we start? Please recall that for Newton's method, we considered the second-order Taylor's approximation of  $f$ :

$$f(x) \approx f(x_k) + \nabla f(x_k)^T(x - x_k) + \frac{1}{2}(x - x_k)^T \nabla^2 f(x_k)(x - x_k)$$

Notation: Hessian =  $\nabla^2 = \nabla^2$

We might have written this as  $g_k^T$  and  $H_k$ .

From this, if you recall what we did in Newton's method, we differentiated this with respect to  $x$  and evaluated it at  $x^*$ , which is the minimizer of  $f$ . Differentiating the above equation and evaluating at  $x = x^*$ , we get:

$$\nabla f(x^*) = 0 = g_k + H_k(x^* - x_k)$$

From here, we actually got the direction to be:

$$d_k = -H_k^{-1}g_k$$

This was the basis of Newton's method.

Now, how do we actually avoid the inverse in this step? If taking the inverse is costly and computationally expensive, how do we avoid this step  $H_k^{-1}$ ? The idea is as follows. We consider that  $H_k^{-1}$  and  $H_{k+1}^{-1}$  are close to each other; they are not very different because  $x_{k+1} = x_k + \alpha_k d_k$ . The difference between  $x_{k+1}$  and  $x_k$  is not too much, so we expect the difference between  $H_k^{-1}$  and  $H_{k+1}^{-1}$  not to be too much either. We try to update  $H_k^{-1}$  to  $H_{k+1}^{-1}$ . That is how the quasi-Newton method comes into practice.

How to avoid inversion of the Hessian? The idea is to update  $H_k^{-1}$  to  $H_{k+1}^{-1}$ , assuming that they are not very different. I am just giving the heuristics; this is not a well-defined statement. I am saying "very different" and "update" to help you understand what is actually happening in the quasi-Newton method so that you can follow the different methods and algorithms.

Now, let us define  $B_k = H_k^{-1}$ . Then we have the second-order approximation using this equation:

$$f(x) \approx f(x_k) + g_k^T(x - x_k) + \frac{1}{2}(x - x_k)^T B_k^{-1}(x - x_k)$$

It should have been  $H_k$ , but instead I am writing  $B_k^{-1}$  now.

With this definition, in the idea above, I need to update  $B_{k+1}$  using  $B_k$ . I need to update  $B_k$  to  $B_{k+1}$ , just like how I update  $x_{k+1}$  as  $x_k + \alpha_k d_k$ . How do we do it? I will differentiate this particular equation (the last equation I have written) and evaluate it at  $x = x_{k+1}$ .

Differentiating the above equation and evaluating at  $x = x_{k+1}$ , we get:

$$\nabla f(x_{k+1}) = g_{k+1} = g_k + B_k^{-1}(x_{k+1} - x_k)$$

This is:

$$g_{k+1} - g_k = B_k^{-1}(x_{k+1} - x_k)$$

Define  $\gamma_k = g_{k+1} - g_k$  and  $\delta_k = x_{k+1} - x_k$ . So we have:

$$\gamma_k = B_k^{-1} \delta_k$$

Or, equivalently:

$\delta_k = B_k \gamma_k$ ? Let's be careful. From  $g_{k+1} - g_k = B_k^{-1} \delta_k$ , multiplying both sides by  $B_k$  gives:

$$B_k (g_{k+1} - g_k) = \delta_k$$

So,  $\delta_k = B_k \gamma_k$ . But we want to update  $B_k$  to  $B_{k+1}$ . The condition we want for the new matrix is that it satisfies a similar secant equation. The standard quasi-Newton condition is:

$$\delta_k = B_{k+1} \gamma_k$$

This is called the **quasi-Newton condition**. When you update  $B_k$  to  $B_{k+1}$ , you must make sure that this quasi-Newton condition is satisfied. That is basically how we construct different quasi-Newton methods.

(Refer Slide Time 10:51)

Week 6 - Quasi-Newton methods

Taylor's approximation (second-order):

$$f(x) \approx f(x^k) + \nabla f(x^k)^T (x - x^k) + \frac{1}{2} (x - x^k)^T \nabla^2 f(x^k) (x - x^k)$$

Differentiating the above eqn., and evaluating at  $x = x^*$ , we get

$$0 = g^k + H^k (x^* - x^k) \Rightarrow d^k = -(H^k)^{-1} g^k.$$

How to avoid inversion of Hessian?

Idea: Update  $(H^k)^{-1}$  to  $(H^{k+1})^{-1}$ , assuming that they are not very different.

Define  $B^k = (H^k)^{-1}$ . Then we have

$$f(x) \approx f(x^k) + g^k (x - x^k) + \frac{1}{2} (x - x^k)^T (B^k)^{-1} (x - x^k)$$

Differentiating the above eqn. and evaluating at  $x = x^{k+1}$ , we get

$$\nabla f(x^{k+1}) = g^k + (B^{k+1})^{-1} (x^{k+1} - x^k) \Rightarrow (g^{k+1} - g^k) = (B^{k+1})^{-1} (x^{k+1} - x^k)$$

Define  $\gamma^k = g^{k+1} - g^k$ ,  $\delta^k = x^{k+1} - x^k$ . So,  $\boxed{\delta^k = B^{k+1} \gamma^k} \Rightarrow$  Quasi-Newton condition

This week, we will see three different quasi-Newton methods. One is called the rank one correction. The other is the DFP method (Davidon-Fletcher-Powell), named after the discoverers. The third one is the BFGS method (Broyden-Fletcher-Goldfarb-Shanno), named after the four optimizers who came up with these algorithms. We will see these three quasi-Newton methods, their conditions, differences in the algorithms, their performances, and so on.

We will start with rank one correction. We have to update  $B_{k+1}$  from  $B_k$ . As the name suggests, we will update this matrix  $B_{k+1}$  using a rank 1 matrix. Any rank 1 matrix can be written in the form  $\beta u u^T$ , where  $u$  is an  $n$ -length vector. So we update it as:

$$B_{k+1} = B_k + \beta u u^T$$

We need to find out what  $\beta$  and  $u$  are so that the quasi-Newton condition is satisfied. We want  $B_{k+1} \gamma_k = \delta_k$ .

Substituting the update:

$$(B_k + \beta u u^T) \gamma_k = \delta_k$$

$$B_k \gamma_k + \beta u (u^T \gamma_k) = \delta_k$$

Note that  $u^T \gamma_k$  is a scalar. So:

$$\beta (u^T \gamma_k) u = \delta_k - B_k \gamma_k$$

Let us choose  $u$  to be  $\delta_k - B_k \gamma_k$ . In that case, we need  $\beta (u^T \gamma_k)$  to be 1. Therefore, we have:

$$\beta = 1 / (u^T \gamma_k)$$

$$\text{And } u = \delta_k - B_k \gamma_k.$$

Therefore, we actually have:

$$B_{k+1} = B_k + \frac{(\delta_k - B_k \gamma_k)(\delta_k - B_k \gamma_k)^T}{(\delta_k - B_k \gamma_k)^T \gamma_k}$$

This is what we call the rank one correction update. Whenever you update your B matrix using this particular equation (the boxed one), that is the rank one correction method.

I will write down the algorithm fully. We first start by initializing  $x_0$ , a tolerance, and  $k = 0$ . Here we also have to initialize  $B_0$ . The starting point  $B_0$  needs to be a symmetric positive definite matrix. It cannot be any arbitrary matrix. This is an extra step compared to the algorithms we discussed until now.

While  $\|g_k\| > \text{tolerance}$ :

1. Compute the descent direction:  $d_k = -B_k g_k$  (since  $H_k^{-1}$  is approximated by  $B_k$ ).
2. Choose  $\alpha_k$  by a line search algorithm (exact or backtracking).
3. Update the iterate:  $x_{k+1} = x_k + \alpha_k d_k$ .
4. Compute the differences:  $\delta_k = x_{k+1} - x_k = \alpha_k d_k$ ,  $\gamma_k = g_{k+1} - g_k$ .
5. Update the matrix:  $B_{k+1} = B_k + \frac{(\delta_k - B_k \gamma_k)(\delta_k - B_k \gamma_k)^T}{(\delta_k - B_k \gamma_k)^T \gamma_k}$ .
6. Set  $k = k + 1$ .

Finally, output  $x^* = x_k$ .

(Refer Slide Time 20:04)

Rank-one correction:

$$B^{k+1} = B^k + \beta u u^T$$

$$B^{k+1} \eta^k = \delta^k \Rightarrow B^k \eta^k + \beta u (u^T \eta^k) = \delta^k$$

$$\Rightarrow (\beta u^T \eta^k) u = \delta^k - B^k \eta^k$$

Let  $u = \delta^k - B^k \eta^k$ . Then,  $\beta = \frac{1}{(\delta^k - B^k \eta^k)^T \eta^k}$

$$\therefore B^{k+1} = B^k + \frac{(\delta^k - B^k \eta^k)(\delta^k - B^k \eta^k)^T}{(\delta^k - B^k \eta^k)^T \eta^k}$$

Algorithm:

- (i) Initialize  $x^0$ ,  $\text{tol}$ ,  $k=0$ ,  $B^0 \succ 0$ .
- (ii) While  $(\|g^k\| > \text{tol})$ :
  - \*  $d^k = -B^k g^k$
  - \* Choose  $\alpha^k$  by any line search algorithm
  - \*  $x^{k+1} = x^k + \alpha^k d^k$ ,  $\delta^k = \alpha^k d^k$ ,  $\eta^k = g^{k+1} - g^k$ .
  - \*  $B^{k+1} = B^k + \frac{(\delta^k - B^k \eta^k)(\delta^k - B^k \eta^k)^T}{(\delta^k - B^k \eta^k)^T \eta^k}$
- (iii) output  $x^* = x^k$ .

This is basically the rank one correction algorithm. The first thing we do after writing down an algorithm is to evaluate it for quadratic functions. That is how we derived the theories for many

methods earlier. You can recall that conjugate gradient methods (the eigenvector method or Fletcher-Reeves) find the answer in  $n$  steps for quadratic functions, and Newton's method finds it in one step if the Hessian of the quadratic function is positive definite.

Let us try the same with rank one correction. I claim that for quadratic functions, this converges in  $n$  steps. Let us write down the code for this particular algorithm and check whether the steps are followed.

This is rank one correction for quadratic functions. Import numpy. I would like to use a  $5 \times 5$  positive definite matrix, like we did for conjugate gradient and Newton's method. Let  $H$  be a random matrix consisting of integers between  $-9$  and  $9$ , a  $5 \times 5$  matrix. I want to make this positive definite, so I will consider  $H H^T + I$ .

Consider the quadratic function  $f(x) = \frac{1}{2} x^T H x + b^T x$ .  $b$  can be a random vector.

Define  $f(x)$  as `0.5 * np.dot(x, np.dot(H, x)) + np.dot(b, x)`.

The gradient  $\text{grad}(x)$  is `np.dot(H, x) + b`.

Now, do the initializations. Initialize  $x$  randomly. Initialize  $k = 0$ . Initialize  $B$  as the identity matrix  $I$  (which is positive definite).

While `np.linalg.norm(grad(x)) > tolerance` (say  $1e-6$ ):

1. `d = - np.dot(B, grad(x))`
2. For a quadratic function, use exact line search:  $\alpha = - (\text{np.dot}(\text{grad}(x), d)) / (\text{np.dot}(d, \text{np.dot}(H, d)))$
3. Update `x = x +  $\alpha$  * d`
4. Compute  $\delta = \alpha * d$  (which is `x_new - x_old`)
5. Compute  $\gamma = \text{grad}(x_{\text{new}}) - \text{grad}(x_{\text{old}})$
6. Update  $B$ : First compute the vector  $u = \delta - \text{np.dot}(B, \gamma)$ . Then compute the denominator `denom = np.dot(u,  $\gamma$ )`. Then `B = B + np.outer(u, u) / denom`
7. `k = k + 1`

You can also print statements to track  $k$ ,  $x$ , and the norm of the gradient.

The true solution is  $x^* = -\text{np.linalg.inv}(H) @ b$ . You can verify that the algorithm finds this solution in  $n=5$  steps. Note that I chose a random  $H$  (positive definite), a random  $b$ , and started with a random  $x$ . My  $B_0$  was the identity matrix. Given that the function is quadratic, this rank one correction method found the answer in just  $n$  steps. No inversion of  $H$  has happened at all. This is just to verify that the answer is correct.

(Refer Slide Time 29:15)

The image displays two screenshots of a Jupyter Notebook interface. The top screenshot shows the initial code for a rank-one correction algorithm. The code defines a function  $f(x)$  and its gradient  $\text{grad}(x)$ , and then iteratively updates  $x$  and  $B$  until the norm of the gradient is small.

```
import numpy as np

H=np.random.randint(-9,10,(5,5))
H=H.T+np.eye(5)
b=np.random.randn(5)

def f(x):
    return 0.5*x.dot(H*x)+b.dot(x)

def grad(x):
    return H*x+b

x=np.random.randn(5)
k=0
B=np.eye(5)
while(np.linalg.norm(grad(x))>1e-6):
    d=-B@grad(x)
    alpha=-grad(x).dot(d)/(d.dot(H@d))
    x,d1=x+alpha*d,alpha*d
    gm=grad(x+alpha*d)-grad(x)
    B=B+(d1-B@gm)[:,None]@(d1-B@gm)[None,:]/(d1-B@gm).dot(d1)
```

The bottom screenshot shows the same code with numerical output. The output displays the values of  $x$  and  $B$  after the algorithm has converged.

```
x=np.random.randn(5)
k=0
B=np.eye(5)
print("%d,(%1.3e,%1.3e,%1.3e,%1.3e,%1.3e),(%1.3e,%1.3e,%1.3e,%1.3e,%1.3e)"%(k,x[0],x[1],x[2],x[3],x[4],np.linalg.norm(grad(x))))
while(np.linalg.norm(grad(x))>1e-6):
    d=-B@grad(x)
    alpha=-grad(x).dot(d)/(d.dot(H@d))
    x,d1=x+alpha*d,alpha*d
    gm=grad(x+alpha*d)-grad(x)
    B=B+(d1-B@gm)[:,None]@(d1-B@gm)[None,:]/(d1-B@gm).dot(d1)
    k+=1
print("%d,(%1.3e,%1.3e,%1.3e,%1.3e,%1.3e),(%1.3e,%1.3e,%1.3e,%1.3e,%1.3e)"%(k,x[0],x[1],x[2],x[3],x[4],np.linalg.norm(grad(x))))
y=-np.linalg.inv(H)@b
print("%1.3e,%1.3e,%1.3e,%1.3e,%1.3e"%(y[0],y[1],y[2],y[3],y[4]))
```

0,(-3.916e-01,-1.693e+00,-4.953e-01,3.579e-01,2.424e-01),7.582e+02  
1,(1.519e-01,-6.789e-01,-2.113e-01,-4.015e-01,-2.864e-01),1.651e+02  
2,(-8.053e-02,-1.469e-01,-4.398e-01,-2.950e-01,-1.706e-01),1.982e+01  
3,(-1.126e-01,-1.355e-01,-4.181e-01,-2.605e-01,-2.250e-01),1.756e+00  
4,(2.194e-01,1.294e-01,2.945e-01,3.268e-01,1.560e-01),2.912e+00  
5,(1.754e-01,2.652e-01,8.599e-01,4.751e-01,4.581e-01),1.348e-13  
1.754e-01,2.652e-01,8.599e-01,4.751e-01,4.581e-01

So, no inversion has happened, but I have managed to find the answer for any quadratic function in  $n$  steps. That is how good the rank one correction method is. In the next lecture, we will understand more about rank one correction and then move to the other methods.

Thank you.