

# Optimization Algorithms: Theory and Software Implementation

Prof. Thirumulanathan D

Department of Mathematics

Institute of IIT Kanpur

## Lecture: 22

Hello everyone, this is the second lecture of week five. We will continue our discussion on Newton's method.

In the last lecture, we learned the general form of the Newton's method algorithm and looked at an example with a quadratic function. We proved analytically that for a quadratic function, Newton's method finds the solution in a single step, and we also demonstrated this with a Python code.

Now, we will look at another example before moving on to discuss the properties and potential issues with Newton's method. We will again consider the function:

$$f(x_1, x_2) = x_1^2 e^{x_2} + x_2^2 e^{x_1}$$

and use the initial point  $x_0 = (\sqrt{2}, \sqrt{2})$ .

(Refer Slide Time 1:28)

The slide content is as follows:

$$f(x) = \frac{1}{2} x^T H x + b^T x + c, \text{ where } H \succ 0.$$
$$\nabla f(x) = Hx + b \Rightarrow x = -H^{-1}b.$$
$$\nabla^2 f(x) = H \quad \forall x \in \mathbb{R}^n.$$
$$x' = x^0 - H^{-1}(Hx^0 + b) = x^0 - x^0 - H^{-1}b = -H^{-1}b.$$
$$f(x) = x_1^2 e^{x_2} + x_2^2 e^{x_1}. \quad x^0 = (\sqrt{2}, \sqrt{2}).$$

We will implement this using a Python code. For this, we need to define the function  $f(x)$ , its gradient  $\text{grad}(x)$ , and its Hessian matrix  $H(x)$ .

Here is the Python code for Newton's method for this general function:

## Python

```
import numpy as np

# Define the function, its gradient, and its Hessian
def f(x):
    return x[0]**2 * np.exp(x[1]) + x[1]**2 * np.exp(x[0])

def grad(x):
    g1 = 2*x[0]*np.exp(x[1]) + x[1]**2*np.exp(x[0])
    g2 = x[0]**2*np.exp(x[1]) + 2*x[1]*np.exp(x[0])
    return np.array([g1, g2])

def H(x):
    h11 = 2*np.exp(x[1]) + x[1]**2*np.exp(x[0])
    h12 = 2*x[0]*np.exp(x[1]) + 2*x[1]*np.exp(x[0])
    h21 = 2*x[0]*np.exp(x[1]) + 2*x[1]*np.exp(x[0])
    h22 = 2*np.exp(x[0]) + x[0]**2*np.exp(x[1])
    return np.array([[h11, h12], [h21, h22]])

# Algorithm parameters
x = np.array([np.sqrt(2), np.sqrt(2)]) # Initial point
tol = 1e-6
k = 0
max_iter = 100
print(f"K = {k}: x = {x}, f(x) = {f(x)}")

# Main optimization loop
while np.linalg.norm(grad(x)) > tol and k < max_iter:
    g = grad(x)
    H_inv = np.linalg.inv(H(x))
    d = -H_inv @ g # Newton direction
    # Newton update (step size alpha = 1)
    x = x + d
```

```
k += 1
```

```
print(f"k = {k}: x = {x}, f(x) = {f(x)}")
```

```
print(f"\nSolution: {x}")
```

```
print(f"Number of iterations: {k}")
```

This code converges to the solution (0, 0) in about 7 iterations. In previous lectures, the conjugate gradient (Fletcher-Reeves) and gradient descent methods took about 30 steps for the same function. This demonstrates that Newton's method is significantly faster for many functions.

(refer Slide Time 7:24)

The first screenshot shows the following code:

```
[1] import numpy as np
def f(x):
    return x[0]**2*np.exp(x[1])+x[1]**2*np.exp(x[0])

def grad(x):
    return np.array([2*x[0]*np.exp(x[1])+x[1]**2*np.exp(x[0]), \
                    -2*x[1]*np.exp(x[0])+x[0]**2*np.exp(x[1])])

x=np.array([1,1])
d=grad(x)
k,rho,c1=0,0.8,0.75
while(np.linalg.norm(grad(x))>1e-6):
    print("%d, (%1.3e,%1.3e), %1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
    alpha=1#grad(x).dot(grad(x))/grad(x).dot(H@grad(x)) #np.dot(x,y)=x.dot(y)
    while(f(x+alpha*d)-f(x)>c1*alpha*grad(x).dot(d)):
        alpha*=rho #alpha=alpha*rho
    x_old=x
    x=x+alpha*d
    d=grad(x)+(k%2!=0)*(grad(x).dot(grad(x))/(grad(x_old).dot(grad(x_old))*d
    k+=1 #k=k+1
    print("%d, (%1.3e,%1.3e), %1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
```

The second screenshot shows the following code and output:

```
def H(x):
    return np.array([[2*np.exp(x[1])+x[1]**2*np.exp(x[0]), \
                    2*x[0]*np.exp(x[1])+2*x[1]*np.exp(x[0])], \
                    [2*x[0]*np.exp(x[1])+2*x[1]*np.exp(x[0]), \
                    2*np.exp(x[0])+x[0]**2*np.exp(x[1])])

x=np.array([np.sqrt(2),np.sqrt(2)])
k=0
print("%d, (%1.3e,%1.3e), %1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
while(np.linalg.norm(grad(x))>1e-6):
    x=x- np.linalg.inv(H(x))@grad(x)
    k+=1
    print("%d, (%1.3e,%1.3e), %1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
```

Output:

```
0, (1.414e+00, 1.414e+00), 2.809e+01
1, (9.142e-01, 9.142e-01), 9.400e+00
2, (5.039e-01, 5.039e-01), 2.953e+00
3, (2.084e-01, 2.084e-01), 8.015e-01
4, (4.842e-02, 4.842e-02), 1.472e-01
5, (3.254e-03, 3.254e-03), 9.249e-03
6, (1.580e-05, 1.580e-05), 4.469e-05
7, (3.744e-10, 3.744e-10), 1.059e-09
```

The reason for this speed is the order of convergence. The algorithms we learned before (gradient descent, conjugate gradient) have a linear rate of convergence (order  $p = 1$ ). Newton's method has a quadratic rate of convergence (order  $p = 2$ ).

The order of convergence  $p$  is defined by the following limit behavior for large  $k$ :

$$\|x_{k+1} - x^*\| \leq \beta * \|x_k - x^*\|^p$$

where  $\beta$  is a constant and  $x^*$  is the true solution.

\* **Linear Convergence (p=1):** The error decreases exponentially. If  $\|x_k - x^*\| = 10$  and  $\beta=0.5$ , the errors become 5, 2.5, 1.25, 0.625, ....

\* **Quadratic Convergence (p=2):** The error decreases much faster. If  $\|x_k - x^*\| = 1$  and  $\beta=0.5$ , the errors become 0.5, 0.125, 0.0078, .... The number of correct digits roughly doubles with each step.

(Refer Slide Time 16:28)

$f(x) = \frac{1}{2} x^T H x + b^T x + c$ , where  $H \succ 0$ .  
 $\nabla f(x) = Hx + b \Rightarrow x = -H^{-1}b$ .  
 $\nabla^2 f(x) = H \quad \forall x \in \mathbb{R}^n$ .  
 $x^1 = x^0 - H^{-1}(Hx^0 + b) = x^0 - x^0 - H^{-1}b = -H^{-1}b$ .  
 $f(x) = x_1^2 e^{x_2} + x_2^2 e^{x_1}$ .  $x^0 = (\sqrt{e}, \sqrt{e})$ .

An algorithm is said to have an order of convergence of  $p$ , if  

$$\lim_{k \rightarrow \infty} \frac{\|x^{k+1} - x^*\|}{\|x^k - x^*\|^p} = \beta \in (0, \infty).$$

\* When  $p=1$ ,  $\beta \in (0, 1)$ .  
 Then  $\|x^{k+1} - x^*\| \approx \beta \|x^k - x^*\|$   
 If  $\|x^0 - x^*\| = 10$ , and  $\beta = \frac{1}{2}$ , then  $\|x^1 - x^*\| = 5$ ,  $\|x^2 - x^*\| = 2.5$ ,  
 and  $\|x^k - x^*\| \approx \frac{10}{2^k}$ .

\* When  $p=2$ ,  $\beta = \frac{1}{2}$ . Then  $\|x^{k+1} - x^*\| \approx \beta \|x^k - x^*\|^2$   
 If  $\|x^0 - x^*\| = 1$ , then  $\|x^1 - x^*\| = 0.5$ ,  $\|x^2 - x^*\| = \frac{1}{2} \left(\frac{1}{2}\right)^2 = \frac{1}{8}$ ,  
 $\|x^3 - x^*\| = \frac{1}{2} \left(\frac{1}{8}\right)^2 = \frac{1}{128}$ ,  $\|x^4 - x^*\| = \frac{1}{2} \left(\frac{1}{128}\right)^2 = \frac{1}{32768}$ .

However, this immense speed comes with significant drawbacks:

1. **No Global Convergence:** Newton's method is only locally convergent. This means it will converge quadratically only if the initial guess  $x_0$  is already sufficiently close to the solution  $x^*$ . If started too far away, it may not converge at all, or even diverge, as the quadratic model becomes a poor approximation of the function.

2. **Descent Direction Not Guaranteed:** The Newton direction  $d_k = -H_k^{-1}g_k$  is a descent direction only if the Hessian  $H_k$  is positive definite. For non-convex functions, the Hessian may not be positive definite at some points, causing the algorithm to potentially increase the function value.

3. **Computational Expense:** Inverting the Hessian matrix  $H_k$  (or solving the linear system  $H_k d_k = -g_k$ ) in each iteration is computationally expensive, especially for high-dimensional problems (e.g., in machine learning with thousands of variables). The Hessian might also be singular (non-invertible).

To illustrate the first issue (no global convergence), consider the function:

$$f(x) = \sqrt{(x_1^2 + 1)} + \sqrt{(x_2^2 + 1)}$$

Its minimum is at (0, 0). Its gradient and Hessian are:

$$\nabla f(x) = [x_1/\sqrt{(x_1^2+1)}, x_2/\sqrt{(x_2^2+1)}]$$

$$H(x) = \text{diag}(1/(x_1^2+1)^{3/2}, 1/(x_2^2+1)^{3/2}) \text{ (a diagonal matrix)}$$

(Refer Slide Time 29:10)

The order of convergence of Newton's method is TWO.

If  $\|x^0 - x^*\| = 10$ , then  $\|x^1 - x^*\| = \frac{1}{2}(10)^2 = 50$ ,  $\|x^2 - x^*\| = \frac{1}{2}(50)^2 = 1250$ ,  
 $\|x^3 - x^*\| = \frac{1}{2}(1250)^2 =$

Issues:

- \* No global convergence. Newton's method is only locally convergent.  
 (i.e.)  $\exists \delta > 0$  for which when  $\|x^0 - x^*\| < \delta$ , the algorithm converges.
- \*  $(-H^k)^{-1}g^k$  need not be a descent direction if  $H^k \neq 0$ . In such a case, we would have  $f(x^{k+1}) > f(x^k)$ .
- \* Matrix inversion is a computationally expensive process. We may also end up with singular  $H^k$ .

$$f(x) = \sqrt{x_1^2 + 1} + \sqrt{x_2^2 + 1}$$

$$\nabla f(x) = \left[ \frac{x_1}{\sqrt{x_1^2 + 1}}, \frac{x_2}{\sqrt{x_2^2 + 1}} \right]^T$$

$$\frac{\partial^2 f}{\partial x_1^2} = \frac{\sqrt{x_1^2 + 1} - x_1 \left( \frac{x_1}{\sqrt{x_1^2 + 1}} \right)}{x_1^2 + 1} = \frac{1}{(x_1^2 + 1)^{3/2}}$$

$$\text{So } \nabla^2 f(x) = \begin{bmatrix} \frac{1}{(x_1^2 + 1)^{3/2}} & 0 \\ 0 & \frac{1}{(x_2^2 + 1)^{3/2}} \end{bmatrix}$$

If you run Newton's method on this function:

## Python

- \* Starting from (0.5, 0.5) (close to the solution), it converges quickly to (0, 0).
- \* Starting from (2, 2) (farther away), the algorithm can diverge, sending the iterates to very large values, demonstrating the lack of global convergence.

(Refer Slide Time 34:49)

Untitled13.ipynb - Colab

colabresearch.google.com/drive/1DTfYUgRkdeZuGeZ15MnPCiqMPiaP2ofscroIfTo-Xy5RQ2G5CdM

File Edit View Insert Runtime Tools Help

Example for no global convergence

```
import numpy as np
def f(x):
    return np.sqrt(x[0]**2+1)+np.sqrt(x[1]**2+1)

def grad(x):
    return np.array([x[0]/np.sqrt(x[0]**2+1),\
                    x[1]/np.sqrt(x[1]**2+1)])

def H(x):
    return np.array([(x[0]**2+1)**(-1.5),0],\
                    [0,(x[1]**2+1)**(-1.5)])

x=np.array([0.5,0.5])
k=0
print("%d, (%1.3e,%1.3e),%1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
while(np.linalg.norm(grad(x))>1e-6):
    x=x-np.linalg.inv(H(x))@grad(x)
    k+=1
    print("%d, (%1.3e,%1.3e),%1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
```

0, (5.000e-01, 5.000e-01), 6.325e-01  
1, (-1.250e-01, -1.250e-01), 1.754e-01  
2, (1.953e-03, 1.953e-03), 2.762e-03  
3, (-7.451e-09, -7.451e-09), 1.854e-08

32:29 / 35:45

0s completed at 16:21

Bilingual YouTube

Untitled13.ipynb - Colab

colabresearch.google.com/drive/1DTfYUgRkdeZuGeZ15MnPCiqMPiaP2ofscroIfTo-Xy5RQ2G5CdM

File Edit View Insert Runtime Tools Help

```
import numpy as np
def f(x):
    return np.sqrt(x[0]**2+1)+np.sqrt(x[1]**2+1)

def grad(x):
    return np.array([x[0]/np.sqrt(x[0]**2+1),\
                    x[1]/np.sqrt(x[1]**2+1)])

def H(x):
    return np.array([(x[0]**2+1)**(-1.5),0],\
                    [0,(x[1]**2+1)**(-1.5)])

x=np.array([1,1])
k=0
print("%d, (%1.3e,%1.3e),%1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
while(k!=100 and np.linalg.norm(grad(x))>1e-6):
    x=x-np.linalg.inv(H(x))@grad(x)
    k+=1
    print("%d, (%1.3e,%1.3e),%1.3e"%(k,x[0],x[1],np.linalg.norm(grad(x))))
```

0, (1.000e+00, 1.000e+00), 1.000e+00  
1, (-1.000e+00, -1.000e+00), 1.000e+00  
2, (1.000e+00, 1.000e+00), 1.000e+00  
3, (-1.000e+00, -1.000e+00), 1.000e+00  
4, (1.000e+00, 1.000e+00), 1.000e+00  
5, (-1.000e+00, -1.000e+00), 1.000e+00  
6, (1.000e+00, 1.000e+00), 1.000e+00  
7, (-1.000e+00, -1.000e+00), 1.000e+00

32:54 / 35:45

0s completed at 16:21

Bilingual YouTube

The screenshot shows a Google Colab notebook interface. The top bar includes the Google Colab logo, the notebook title 'Untitled13.ipynb', and various icons for file management, settings, and sharing. The main area is divided into a code editor and an output area. The code editor contains a Python script for a gradient descent algorithm. The output area shows the results of the algorithm, including a series of numerical values and two runtime warnings about overflow encountered in scalar power.

```
def H(x):  
    return np.array([[x[0]**2+1]**(-1.5), 0],  
                    [0, (x[1]**2+1)**(-1.5)])  
  
x=np.array([2,2])  
k=0  
print("%d, (%1.3e, %1.3e), %1.3e" % (k, x[0], x[1], np.linalg.norm(grad(x))))  
while (k!=100 and np.linalg.norm(grad(x))>1e-6):  
    x=-np.linalg.inv(H(x))@grad(x)  
    k+=1  
    print("%d, (%1.3e, %1.3e), %1.3e" % (k, x[0], x[1], np.linalg.norm(grad(x))))  
  
0, (2.000e+00, 2.000e+00), 1.265e+00  
1, (-8.000e+00, -8.000e+00), 1.403e+00  
2, (5.120e+02, 5.120e+02), 1.414e+00  
3, (-1.342e+08, -1.342e+08), 1.414e+00  
4, (2.418e+24, 2.418e+24), 1.414e+00  
5, (-1.413e+73, -1.413e+73), 1.414e+00  
6, (2.824e+219, 2.824e+219), 0.000e+00  
<ipython-input-10-2423383aaab0>:6: RuntimeWarning: overflow encountered in scalar power  
    return np.array([x[0]/np.sqrt(x[0]**2+1)],\n<ipython-input-10-2423383aaab0>:7: RuntimeWarning: overflow encountered in scalar power  
    x[1]/np.sqrt(x[1]**2+1)])
```

In the next lecture, we will look at examples for the other issues (non-descent directions and singular Hessians). Thank you.