#### Line Search Methods

#### Saurav Samantaray

Department of Mathematics

Indian Institute of Technology Madras

#### February 16, 2024



<ロ > < 母 > < 臣 > < 臣 > 臣 の Q で 1/20

#### Line Search Method

- In each iteration of a line search method a search direction <u>pk</u> is computed, and
- then its decided how far to move along that direction.

An iteration is given by

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \tag{1}$$

where  $\alpha_k > 0$  (scalar) called the step length

. The success of a line search method depends on effective choices of both:

- **1** the direction  $p_k$
- **2** the step length  $\alpha_k$

#### The Steepest Descent Direction

- The steepest descent direction −∇*f<sub>k</sub>* is the most obvious choice for search direction for a line search method.
- Among all Directions once could from x<sub>k</sub>, along −∇f<sub>k</sub>, f decreases most rapidly.

#### Justification

• Consider any search direction p and step-length  $\alpha$ , we have

$$f(x_k+\alpha p) = f_k+\alpha p^T \nabla f_k + \frac{\alpha^2 p^T \nabla^2 f(x_k+tp)p}{2} \text{ for some } t \in (0,\alpha)$$

 Let α << 1 (small) and we consider the first-order approximation of f at x<sub>k</sub> + αp around x<sub>k</sub> as:

$$f(x_k + \alpha p) \approx f(x_k) + \alpha p^T \nabla f_k$$

• Change in f moving from  $x_k$  to  $x_k + \alpha p$  is  $f(x_k + \alpha p) - f(x_k)$ 

Line Search Methods The Steepest Descent Direction

# The Steepest Descent Direction

 As the distance moved in the direction is α, therefore, the rate of change of f along the direction p at x<sub>k</sub> is

$$\frac{f(x_k + \alpha p) - f(x_k)}{\alpha}$$

• which is coefficient of  $\alpha$ , i.e.

$$p^T \nabla f_k$$

- This implies smaller the above value is, more descent can be achieved.
- Hence, the unit direction *p* of most rapid decrease is the solution to the problem

$$\min_{p} p^{T} \nabla f_{k}, \qquad \text{sub}$$

subject to ||p|| = 1.

Line Search Methods The Steepest Descent Direction

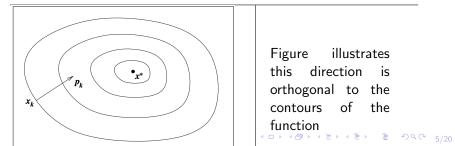
#### The Steepest Descent Direction

$$p^T \nabla f_k = ||p|| ||\nabla f_k|| \cos \theta = ||\nabla f_k|| \cos \theta$$

where  $\theta$  is the angle between p and  $\nabla f_k$ .

• The minimiser is attained when  $\cos heta = -1$  and

$$p = -\frac{\nabla f_k}{||\nabla f_k||}$$



## The Steepest Descent Direction

• At every step (iteration) in the steepest descent method the search direction is chosen along

$$p = -\nabla f_k$$

- $\alpha_k$  can be chosen in a variety of ways.
- One advantage of this method is it requires only the calculation of gradient (∇f<sub>k</sub>), but not second derivatives.
- Line search methods may use search directions other than the steepest descent direction.

# **Descent Direction**

#### **Descent Direction**

Any direction that makes an angle of strictly less than  $\frac{\pi}{2}$  radians with  $-\nabla f_k$  is guaranteed to produce a decrease in f, provided the step length is sufficiently small and is called a descent direction.

Now consider

$$f(x_k + \epsilon p_k) = f(x_k) + \epsilon p_k^T \nabla f_k + \mathscr{O}(\epsilon^2).$$

When p<sub>k</sub> is a downhill (descent) direction, the angle θ<sub>k</sub> between p<sub>k</sub> and ∇f<sub>k</sub> has cos θ<sub>k</sub> < 0, so that</li>

$$p_k^T \nabla f_k = ||p_k|| ||\nabla f_k|| \cos \theta_k < 0$$
$$\implies f(x_k + \epsilon p_k) < f(x_k)$$

 Most line search algorithms require pk to be descent direction, because this property guarantees that the function are 7/20

• This direction is derived from the second-order Taylor series approximation to  $f(x_k + p)$ , which is

$$f(x_{k+p}) \approx f_k + p^T \nabla f_k + \frac{1}{2} p^T \nabla^2 f_k p = {}^{def} m_k(p).$$
 (2)

- Assuming for the moment that  $\nabla^2 f_k$  is positive definite, we obtain the Newton direction by finding the vector p that minimizes  $m_k(p)$ .
- By simply setting the derivative of m<sub>k</sub>(p) to zero, we obtain the following explicit formula:

$$p_k^N = -(\nabla^2 f_k)^{-1} \nabla f_k.$$
(3)

 The Newton direction is reliable when the difference between the true function f(x<sub>k</sub> + p) and its quadratic model m<sub>k</sub>(p) is not too large.

- If ∇<sup>2</sup>f is sufficiently smooth, this difference introduces a perturbation of only O(||p||<sup>3</sup>).
- Therefore, when ||p|| is small

 $f(x_k + p) \approx m_k(p)$  quite accurately

• The Newton direction can be used in a line search method when  $\nabla^2 f_k$  is positive definite, as in this case we have

$$abla f_k^T p_k^N = -p_k^{N^T} 
abla^2 f_k p_k^N 
onumber \\ < -\sigma_k ||p_k^N||^2$$

for some  $\sigma_k > 0$  (+ve definiteness of  $\nabla^2 f_k$ )

Unless the gradient ∇f<sub>k</sub> (and therefore the step p<sup>N</sup><sub>k</sub>) is zero, we have

$$abla f^{T} p_{k}^{N} < 0$$
 , the set of th

- Unlike the steepest descent direction, there is a <u>"natural"</u> step length of 1 associated with the Newton direction.
- Adjust  $\alpha$  only when it does not produce a satisfactory reduction in the value of f.
- Note that when  $\nabla^2 f_k$  is not positive definite the Newton direction may not exist, since  $(\nabla^2 f_k)^{-1}$  may not exist.
- Even when it is defined, it may not satisfy the descent property, and therefore is unsuitable.
- Methods that use Newton direction have fast rate of local convergence (more on this later).
- After a neighbourhood of the solution is reached, convergence to high accuracy often occurs in just a few iterations.

- The main drawback is the need to calculate the Hessian  $\nabla^2 f(x)$ .
- Explicit computation of this matrix of second derivatives can sometimes be a cumbersome, error-prone, and expensive process.
- Finite-difference and automatic differentiation techniques come useful in avoiding the need to calculate second derivatives by hand.

#### Quasi-Newton Search Direction

- <u>Quasi-Newton</u> search directions provide an attractive alternative to Newton's method.
- They do not require computation of the Hessian and yet still attain a superlinear rate of convergence.
- In place of the true Hessian ∇<sup>2</sup> f<sub>k</sub>, they use an approximation B<sub>k</sub>, which is updated after each step to take account of the additional knowledge gained during the step.
- The updates make use of the fact that changes in the gradient g provide information about the second derivative of f along the search direction.

#### Quasi-Newton Search Direction

• Using the Taylor's expansion we have

$$\nabla f(x+p) = \nabla f(x) + p^T \nabla^2 f(x) + \mathcal{O}(||p||^2)$$

• By setting  $x = x_k$  and  $p = x_{k+1} - x_k$ , we obtain

$$\nabla f_{k+1} = \nabla f_k + \nabla^2 f_k (x_{k+1} - x_k) + \mathscr{O}(||x_{k+1} - x_k||^2).$$

<□ ▶ < □ ▶ < ■ ▶ < ■ ▶ < ■ ▶ ■ の へ · 13/20

#### Quasi-Newton Search Direction

• When  $x_k$  and  $x_{k+1}$  lie in a region near the solution  $x^*$ , within which  $\nabla^2 f$  is positive definite, the final term in this expansion is eventually dominated by the  $\nabla^2 f_k(x_{k+1} - x_k)$  term, and we an write.

$$\nabla^2 f_k(x_{k+1}-x_k) \approx \nabla f_{k+1} - \nabla f_k.$$
(4)

- We choose the new Hessian approximation B<sub>k+1</sub> so that it mimics the property (4) of the true Hessian.
- That is we require it to satisfy the following condition, known as the secant equation:

$$B_{k+1}s_k = y_k, \tag{5}$$

where  $s_k = x_{k+1} - x_k$ ,  $y_k = \nabla f_{k+1} - \nabla f_k$ .

# Quasi-Newton Search Direction

Typically, additional conditions are imposed on  $B_{k+1}$ , such as

- Symmetry (Motivated by symmetry of the exact Hessian).
- a requirement that the difference between successive approximations  $B_k$  and  $B_{k+1}$  have low rank.

Two of the most popular formulae for updating the Hessian approximation  $B_k$  are:

Symmetric-rank-one (SR1) formula:

$$B_{k+1} = B_k + \frac{(y_k - B_k s_k)(y_k - B_k s_k)^T}{(y_k - B_k s_k)^T s_k}.$$
 (6)

▲□▶▲圖▶▲≣▶▲≣▶ ≣ のへで 15/20

# Quasi-Newton Search Direction

#### **BFGS Formula:**

$$B_{k+1} = B_k - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \frac{y_k y_k^T}{y_k^T s_k}.$$
 (7)

• The difference between the matrices  $B_k$  and  $B_{k+1}$  is

#### • rank-one for SR1.

- rank-two in case of BFGS.
- Both updates satisfy the secant equation and both maintain symmetry.

## Quasi-Newton Search Direction

• It can be shown that the BFGS formula generates positive definite approximations whenever the initial approximation  $B_o$  is positive definite and

$$s_k^T y_k > 0$$

The quasi-Netwon search direction is obtained by using B<sub>k</sub> instead of the exact Hessian ∇<sup>2</sup>f<sub>k</sub> in the Newton direction

$$p_k = -B_k^{-1} \nabla f_k.$$

• Some practical implementations of quasi-Newton methods avoid the need to factorise  $B_k$  at each iteration by <u>updating</u> inverse of  $B_k$ , instead of  $B_k$  itself.

## Quasi-Newton Search Direction

• The equivalent formula for (SR1), applied to the inverse approximation

$$H_k : {}^{defn} = B_k^{-1} \text{ is}$$
$$H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho_k s_k s_k^T$$

where  $\rho_k = \frac{1}{y_k^T s_k}$ 

• Then *p<sub>k</sub>* is given by

$$p_k = -H_k \nabla f_k$$

◆□▶ ◆□▶ ◆ ■▶ ◆ ■ ・ ● ● ○ ○ 18/20

• This matrix-vector multiplication is simpler than the factorisation / back-substitution produce.

Line Search Methods Non-linear Conjugate Gradient Methods

# Non-linear Conjugate Gradient Methods

- The direction here are generated by non-linear conjugate gradient methods.
- They have the form

$$p_{K} = -\nabla f(x_{k}) + B_{k}p_{k-1}$$

where  $B_k$  is a scalar that ensures that  $p_k$  and  $p_{k-1}$  are conjugate (to be defined later).

• Conjugate gradient methods were originally designed to solve systems of linear equations:

$$Ax = b$$

 Line Search Methods Non-linear Conjugate Gradient Methods

# Non-linear Conjugate Gradient Methods

• The problem of solving this linear system is equivalent to the problem of minimising the convex quadratic function defined by

$$\phi(x) = \frac{1}{2}x^T A x - b^T x.$$

- In general, non-linear conjugate gradient directions are much more effective than the steepest descent direction and are almost as simple to compute.
- These methods do not attain fast convergence rates of Newton or quasi-Newton methods.
- But, they have the advantage of not requiring storage of matrices.