

Line Search Methods Analysis

Saurav Samantaray

Department of Mathematics

Indian Institute of Technology Madras

February 14, 2024



Step Length

- In computing the step length we face a trade-off.
- We want to choose α_k to give a substantial reduction of f , but we don't want to spend too much time making the choice.
- Off-course the ideal choice would be the global minimiser of the univariate function $\phi(\cdot)$ defined by

$$\phi(\alpha) = f(x_k + \alpha p_k), \alpha > 0. \quad (1)$$

- But in general, it is too expensive to identify this value.
- It requires too many evaluations of the objective function and/or the gradient to even find a local minimiser to moderate precision.

Step Length

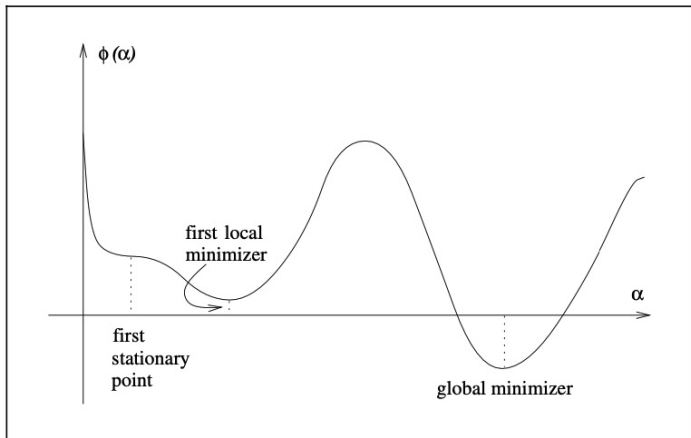


Figure: The ideal step length is the global minimiser

Step Length

- Practically, strategies perform an inexact line search to identify a step length that achieves adequate reductions in f at minimal cost.
- We will discuss these search strategies a little later.
- We will now discuss various termination conditions for line search algorithms and show that effective step lengths need not lie near minimisers of the univariate function $\phi(\alpha)$.
- Is $f(x_k + \alpha_k p_k) < f(x_k)$ good enough to get convergence??
- for example consider the function

$$f(x) = x^2 - 1$$

it has the global minima at $x = 0$, $f = -1$.

Step Length

- Consider a sequence $\{x_k\}$ s.t.

$$f(x_k) = \frac{5}{k}, \quad k = 1, 2, 3, \dots$$

$$\implies f(x_k) > f(x_{k+1})$$

- The reduction in f at each step is not enough to get it to converge to the minimiser.

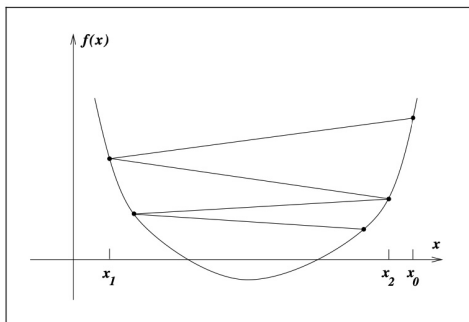


Figure: Insufficient reduction

The Wolfe Condition

Armijo Condition (Sufficient Decrease Condition):

α_k should be chosen such that

$$f(x_k + \alpha p_k) \leq f(x_k) + c_1 \alpha \nabla f_k^T p_k \quad (2)$$

for some constant $c_1 \in (0, 1)$.

- Since p_k is a descent direction and $c_1 > 0$ and $\alpha > 0$ the first thing that the **Armijo condition** asserts that there is a reduction in f from x_k to $x_{k+1} = x_k + \alpha p_k$.
- The reduction in f is at least

$$c_1 \alpha \nabla f_k^T p_k$$

therefore it also says the reduction in f must be proportional to both the step length α_k and the directional derivative $\nabla f_k^T p_k$

The Wolfe Condition

- The right hand side of (2) is a linear function in α (say) $l(\alpha)$.

$$l(\alpha) = f(x_\alpha) + c_1 \alpha \nabla f_k^T p_k$$

- The function $l(\cdot)$ has a negative slope $c_1 \nabla f_k^T p_k$ but $c_1 \in (0, 1)$.
- Therefore, it lies above the graph of ϕ for small positive values of α .
- The sufficient decrease condition states that α is acceptable only if

$$\phi(\alpha) \leq l(\alpha).$$

- In practice, c_1 is chosen to be quite small, say

$$c_1 = 10^{-4}$$

The Wolfe Condition

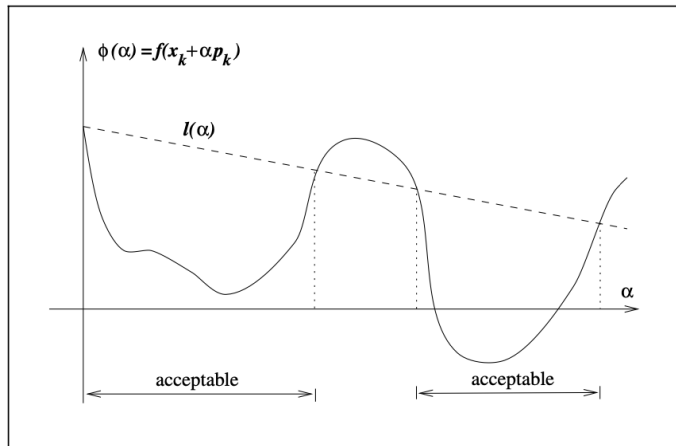


Figure: The intervals on which the Armijo condition is satisfied is shown

The Wolfe Condition

- The sufficient decrease condition is not enough by itself to ensure that the algorithm makes reasonable progress.
- As it is satisfied for all sufficiently small values of α

The Wolfe Condition

- To rule out unacceptable short steps we introduce a second requirement.

Curvature Conditions

α_k should satisfy

$$\nabla f(x_k + \alpha_k p_k)^T p_k \geq c_2 \nabla f_k^T p_k \quad (3)$$

for some constant $c_2 \in (c_1, 1)$.

- The left-hand side is simply the derivative $\phi'(\alpha_k)$.
- So the curvature condition ensures that the slope of ϕ at α_k is greater than c_2 times the initial slope $\phi'(0)$.
- If the **slope $\phi'(\alpha)$ is strongly negative**, we have an indication that **we can reduce f significantly** by moving further along the chosen direction.

The Wolfe Condition

- On, the other hand if $\phi'(\alpha)$ is only **slightly negative** or **even positive**, it is a sign that we cannot expect much more decrease in f in this direction.
- So it makes sense to terminate the line search. (See Figure 6)
- Typical values of c_2 are 0.9 when the search direction p_k is chosen by a Newton or quasi-Newton method, and 0.1 when p_k is obtained from a non-linear conjugate gradient method.
- The sufficient decrease and curvature conditions are known collectively as the Wolfe conditions.

The Wolfe Condition

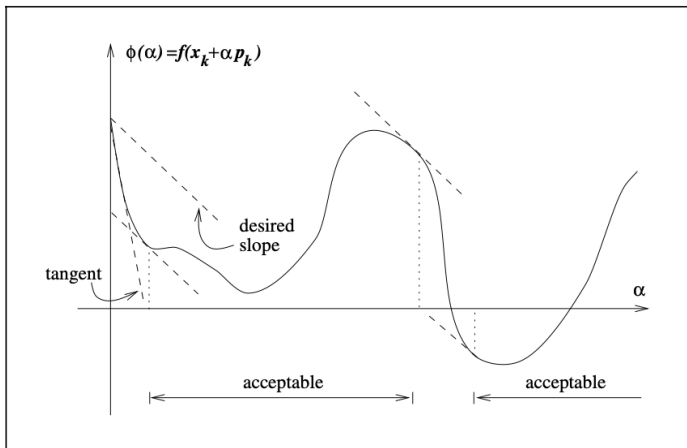


Figure: Insufficient Reduction

Wolfe Conditions

$$\begin{aligned} f(x_k + \alpha_k p_k) &\leq f(x_k) + c_1 \alpha_k \nabla f_k^T p_k \\ \nabla f(x_k + \alpha_k p_k)^T p_k &\geq c_2 \nabla f_k^T p_k. \end{aligned} \tag{4}$$

with $0 < c_1 < c_2 < 1$.

- A step length may satisfy the Wolfe conditions without being particularly close to a minimiser of ϕ . (See previous figure)
- The curvature conditions can be modified to force α_k to lie in at least a broad neighbourhood of a local minimiser or stationary point of ϕ .

The Strong Wolfe Conditions

α_k is required to satisfy

$$\begin{aligned} f(x_k + \alpha_k p_k) &\leq f(x_k) + c_1 \alpha_k \nabla f_k^T p_k \\ |\nabla f(x_k + \alpha_k p_k)^T p_k| &\leq c_2 |\nabla f_k^T p_k|. \end{aligned} \tag{5}$$

with $0 < c_1 < c_2 < 1$.

- The only difference with the Wolfe conditions is that we no longer allow the derivative $\phi'(\alpha)$ to be too positive.
- It excludes points that are far from stationary points of ϕ .
- Is it always possible to find step lengths that satisfy Wolfe conditions ?

Existence of α satisfying Wolfe conditions

Lemma

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable. Let p_k be a descent direction at x_k , and assume that f is bounded below along the ray

$$\{x_k + \alpha p_k \mid \alpha > 0\}$$

Then if $0 < c_1 < c_2 < 1$, there exist intervals of step lengths satisfying Wolfe conditions and the strong Wolfe conditions.

Sketch of Proof



$$\phi(\alpha) = f(x_k + \alpha p_k)$$

is bounded below for all $\alpha > 0$

- Let $l(\alpha) = f(x_k) + \alpha c_1 \nabla f_k^T p_k$, the line is unbounded below and must intersect the graph of ϕ at least once.

Existence of α satisfying Wolfe conditions

- Note that for very small values if α (we can find such α)

$$\begin{aligned}
 l(\alpha) &= f(x_k) + \alpha c_1 \nabla f_k^T p_k \\
 &> f(x_k) + \alpha \nabla f_k^T p_k \quad \text{as } \nabla f_k^T p_k < 0 \text{ and } c_1 < 1 \\
 &\approx f(x_k + \alpha p_k) = \phi(\alpha).
 \end{aligned}$$

Therefore, to start with, the graph of $l(\alpha)$ stays above $\phi(\alpha)$.

- Now since $\phi(\alpha)$ is bounded below \exists a minimum value and since $l(\alpha)$ is unbounded below it will (for large values of α) attain values lesser than the minimum value of $\phi(\alpha)$.
Therefore, both the graphs will intersect atleast once.
- Let $\alpha' > 0$ be the smallest intersecting value of α that is

$$f(x_k + \alpha' p_k) = f(x_k) + \alpha' c_1 \nabla f_k^T p_k.$$

Existence of α satisfying Wolfe conditions

- α' is the point where the line $l(\alpha)$ meets $\phi(\alpha)$ for the first time . Therefore for all $\alpha < \alpha'$ the sufficient decrease condition holds good.
- Now by applying the mean value theorem on $\phi(\alpha)$ in the interval $[0, \alpha']$ we get

$$\begin{aligned} \frac{\phi(\alpha') - \phi(0)}{\alpha' - 0} &= \phi'(\alpha'') \quad \alpha'' \in (0, \alpha') \\ \implies f(x_k + \alpha' p_k) - f(x_k) &= \alpha' \nabla f(x_k + \alpha'' p_k)^T p_k \\ \implies f(x_k + \alpha' p_k) &= f(x_k) + \alpha' \nabla f(x_k + \alpha'' p_k)^T p_k \\ \nabla f(x_k + \alpha'' p_k)^T p_k &= c_1 \nabla f_k^T p_k > c_2 \nabla f_k^T p_k \end{aligned} \tag{6}$$

since $c_2 > c_1$ and $\nabla f_k^T p_k < 0$.

- α'' satisfies the Wolfe conditions and the inequalities hold strictly for both the condition.

Existence of α satisfying Wolfe conditions

- Hence, by our smoothness assumption on f , there is an interval around α'' for which the Wolfe conditions hold.
- Moreover, since the left-hand side term in the curvature condition is negative, the strong Wolfe condition also holds in the same interval.

The Goldstein Conditions

The Goldstein are stated as a pair of inequalities, in the following way:

$$f(x_k) + (1-c)\alpha_k \nabla f_k^T p_k \leq f(x_k + \alpha_k p_k) \leq f(x_k) + c\alpha_k \nabla f_k^T p_k, \quad (7)$$

with $0 < c < \frac{1}{2}$.

- The second inequality is the sufficient decrease condition.
- Whereas the first inequality is introduced to control the step length from below.
- A disadvantage of the Goldstein conditions vis-a-vis the Wolfe conditions is that the first inequality in (7) may exclude all minimizers of ϕ .
- However, the Goldstein and Wolfe conditions have much in common, and their convergence theories are quite similar.

The Goldstein Conditions

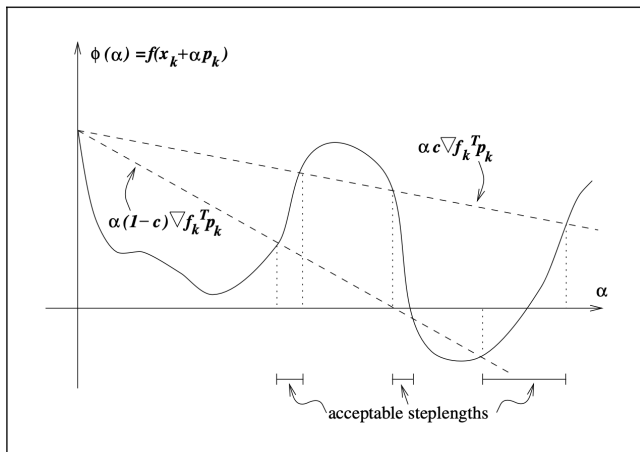


Figure: The Goldstein conditions.

Sufficient Decrease and Backtracking

- The sufficient decrease condition alone is not sufficient to ensure that the algorithm makes reasonable progress along the given search direction.
- However, the extra curvature condition can be dispensed off by using a so-called backtracking approach to choose candidate step length.

Backtracking Line Search

- 1 Choose $\bar{\alpha} > 0$, $\rho \in (0, 1)$, $c \in (0, 1)$;
- 2 Set $\alpha = \bar{\alpha}$
- 3 While $f(x_k + \alpha p_k) > f(x_k) + c\alpha \nabla f_K^T p_k$
- 4 $\alpha = \rho\alpha$;
- 5 end.

Terminate with $\alpha_k = \alpha$.

Sufficient Decrease and Backtracking

- The initial step length $\bar{\alpha}$ is chosen to be 1 in Newton and quasi-Newton methods, but can have different values in other algorithms, such as steepest descent or conjugate gradient.
- An acceptable step length will be found in a finite number of steps as α_k will eventually become small enough to satisfy the sufficient decrease condition.
- In practice the contraction factor " ρ " is allowed to vary at each iteration of the line search.
- One may need to ensure that $\rho \in [\rho_{lo}, \rho_{hi}]$ for some fixed constants $0 < \rho_{lo} < \rho_{hi} < 1$.

Sufficient Decrease and Backtracking

- The backtracking approach either choose $\alpha_k = \bar{\alpha}$ the initial choice or else α_k is short enough to satisfy the sufficient decrease condition.
- Still α_k is not very small as, $\frac{\alpha_k}{\rho}$ doesn't satisfy the sufficient decrease condition.
- It is only by a factor of $\frac{1}{\rho}$ that α_k is shorter from the previous choice of α_k which doesn't work.
- It is a very simple and quite a popular strategy to terminate line search algorithms.
- Well suited for Newton methods but less appropriate for quasi-Newton and conjugate gradient methods.

Convergence of Line Search Methods

Global Convergence

$$\|\nabla f_k\| \rightarrow 0 \quad \text{as } k \rightarrow \infty$$

i.e. convergence to a stationary point for any starting point x_0 .

To obtain global convergence:

- 1 Need to choose step lengths well;
 - 2 Choose search directions p_k appropriately as well.
- Let p_k be a chosen direction at the k th iteration of the line search method.
 - We define θ_k to be the angle between p_k and the steepest descent direction $-\nabla f_k$ given by

$$\cos \theta = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \|p_k\|} \quad (8)$$

Global Convergence

Theorem (Zoutendijk)

Consider any iteration of the form

$$x_{k+1} = x_k + \alpha_k p_k$$

where p_k is a descent direction and α_k satisfies the Wolfe conditions. Suppose that f is bounded below in \mathbb{R}^n and that f is continuously differentiable in an open set \mathcal{N} containing the level set

$$\mathcal{L} =^{\text{def}} \{x : f(x) \leq f(x_0)\}$$

where x_0 is the starting point of the iteration. Assume also that the gradient " ∇f " is Lipschitz continuous on \mathcal{N} , i.e. there exists a constant $L > 0$ s.t.

$$\|\nabla f(x) - \nabla f(\tilde{x})\| \leq L\|x - \tilde{x}\|, \quad \text{for all } x, \tilde{x} \in \mathcal{N}$$

Then

$$\sum_{k \geq 0} \cos^2 \theta_k \|\nabla f_k\|^2 < \infty$$

Global Convergence

Proof:

- Consider the second Wolfe condition,

$$\begin{aligned} \nabla f(x_k + \alpha_k p_k)^T p_k &\geq c_2 \nabla f_k^T p_k \\ \text{or, } \nabla f(x_{k+1})^T p_k &\geq c_2 \nabla f_k^T p_k \\ \text{or, } \nabla f(x_{k+1})^T p_k - \nabla f(x_k)^T p_k &\geq (c_2 - 1) \nabla f_k^T p_k \\ \text{or, } (\nabla f(x_{k+1})^T - \nabla f(x_k)^T) p_k &\geq (c_2 - 1) \nabla f_k^T p_k \end{aligned} \tag{9}$$

- For every descent direction, iteration lives in the level set.
- From the Lipschitz condition we have:

$$(\nabla f(x_{k+1}) - \nabla f(x_k))^T p_k \leq \alpha_k L \|p_k\|^2.$$

Global Convergence

- By combining the two relation i.e. the last equation in (9) and the one above we obtain

$$\alpha_k \geq \frac{(c_2 - 1) \nabla f_k^T p_k}{L \|p_k\|^2} \quad (10)$$

- Now consider the first Wolfe condition

$$f(x_k + \alpha_k p_k) \leq f(x_k) + c_1 \alpha_k \nabla f_k^T p_k$$

or, $f_{k+1} \leq f_k + c_1 \alpha_k \nabla f_k^T p_k$ (as $\nabla f_k^T p_k < 0$)

(11)

or, $f_{k+1} \leq f_k + c_1 \frac{(c_2 - 1) (\nabla f_k^T p_k)^2}{L \|p_k\|^2}$ using (10)

- Note that

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \|p_k\|} \implies \cos^2 \theta_k \|\nabla f_k\|^2 = \frac{(\nabla f_k^T p_k)^2}{\|p_k\|^2} \quad (12)$$

Global Convergence

- Therefore, $f_{k+1} \leq f_k - \frac{c_1(1-c_2)}{L} \cos^2 \theta_k \|\nabla f_k\|^2$
- Let $c = \frac{c_1(1-c_2)}{L}$.
- By summing this expression over all indices less than or equal to k , we obtain:

$$f_{k+1} \leq f_0 - c \sum_{j=0}^k \cos^2 \theta_j \|\nabla f_j\|^2$$

- Since f is bounded below, we have $f_0 - f_{k+1}$ is less than some positive constant, for all k .
- Therefore, by taking limits in the above we obtain

$$\sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f_k\|^2 < \infty.$$

which concludes the proof.

Global Convergence

- Similar results also hold for the Goldstein conditions or the strong Wolfe conditions.
- For all these strategies, the step length selection implies the inequality

$$\sum_{k \geq 0} \cos^2 \theta_k \|\nabla f_k\|^2 < \infty$$

which is called the Zoutendijk condition.

- The assumptions of the theorem are not too restrictive.
- f needs to be bounded below for the optimisation problem to be well defined.
- The smoothness assumption - Lipschitz continuity of the gradient - is implied by many of the smoothness conditions that are used in local convergence theorems and are often satisfied in practice.

Global Convergence

- The Zoutendijk's condition implies that

$$\cos^2 \theta_k \|\nabla f_k\|^2 \rightarrow 0$$

- If the choice of the search direction p_k is made so that it ensures that the angle θ_k is bounded away from 90° , then there is a positive constant δ s.t.

$$\cos \theta_k > \delta > 0, \quad \text{for all } k.$$

- It now follows immediately that

$$\lim_{k \rightarrow \infty} \|\nabla f_k\| = 0$$

- In other words the gradient norm $\|\nabla f_k\| \rightarrow 0$, provided that the search directions are never too close to orthogonality with the gradient.

Global Convergence

- Line Search + Steepest descent (for which the search direction p_k is parallel to the negative gradient) + Wolfe or Goldstein conditions \implies Produces a gradient that converges to zero.
- For line search methods the Zoutendijk condition is the strongest global convergence result that can be obtained.
- It cannot be guaranteed that the method converges to a minimiser (let alone global minimiser).
- Only insight we get is the algorithm, is attracted to stationary points.

However, by making additional requirements on the search direction p_k

-> by introducing negative curvature information from the Hessian $\nabla^2 f(x_k)$

we can strengthen these results to include convergence to a local minimiser.

Condition Number

Consider any norm on \mathbb{R}^n , and let A be a $n \times n$ matrix

- Let $M = \|A\| = \max \frac{\|Ax\|}{\|x\|}$ (maximum stretching)
- Let $m = \|A\| = \min \frac{\|Ax\|}{\|x\|}$ (minimum stretching)
- The reciprocal of m is the norm of the inverse of A

$$m = \min \frac{\|Ax\|}{\|x\|} = \min \frac{\|y\|}{\|A^{-1}y\|} = \frac{1}{\max \frac{\|A^{-1}y\|}{\|y\|}} = \frac{1}{\|A^{-1}\|}$$

Condition number for inversion

The ratio of maximum to minimum stretching is the condition number for inversion:

$$\kappa(A) = \frac{M}{m}$$

- An equivalent definition is $\kappa(A) = \|A\| \|A^{-1}\|$.
- A finite large condition number means that the matrix is close to being singular.

Convergence for Newton-Like Methods

- Consider a Newton-like method and assume that the matrices B_k are positive definite with a uniformly bounded condition number.
- That is, there is a constant M such that

$$\|B_k\| \|B_k^{-1}\| \leq M, \quad \text{for all } k.$$

- Since B_k is symmetric and positive definite matrix, we have that the matrices $B_k^{1/2}$ and $B_k^{-1/2}$ exist and

$$\|B_k^{1/2}\| = \|B_k\|^{1/2} \quad \text{and} \quad \|B_k^{-1/2}\| = \|B_k^{-1}\|^{1/2}$$

$$\begin{aligned}
\cos \theta_k &= -\frac{\nabla f_k^T p_k}{\|\nabla f_k\| \cdot \|p_k\|} \\
&= \frac{p_k^T B_k p_k}{\|B_k p_k\| \cdot \|p_k\|} && (p_k = -B_k^{-1} \nabla f_k) \\
&\geq \frac{p^T B_k p}{\|B_k\| \|p_k\|^2} && \|B_k p_k\| \leq \|B_k\| \|p_k\| \\
&= \frac{p_k^T B_k^{1/2} B_k^{1/2} p_k}{\|B_k\| \|p_k\|^2} = \frac{\|B_k^{1/2} p_k\|^2}{\|B_k\| \|p_k\|^2} \\
&\geq \frac{\|p_k\|^2}{\|B_k^{-1/2}\|^2 \|B_k\| \|p_k\|^2} = \frac{1}{\|B_k^{-1}\| \|B_k\|} \geq \frac{1}{M}
\end{aligned} \tag{13}$$

By combining this bound with Zoutendijk condition we get

$$\lim_{k \rightarrow \infty} \|\nabla f_k\| = 0$$

Rate of Convergence

- One of the key measures of performance of an algorithm is its rate of convergence.

Q-linear Convergence

Let $\{x_k\}$ be a sequence in \mathbb{R}^n that converges to x^* . We say that the convergence is **Q-linear** if there is a constant $r \in (0, 1)$ such that

$$\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} \leq r, \quad \text{for all } k \text{ sufficiently large.}$$

That is the distance to the solution x^* decreases at each iteration by at least a constant factor bounded away from 1

Example

$\{x_k\} = 1 + (0.5)^k$ converges Q-linearly to 1, with $r = 0.5$.

Rate of Convergence

Q-superlinear

The convergence is said to be Q-superlinear if

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 0.$$

Example

For example, the sequence $1 + k^{-k}$ converges superlinearly to 1.

Q-quadratic

Q-quadratic convergence, an even more rapid convergence rate, is obtained if

$$\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^2} \leq M, \quad \text{for all } k \text{ sufficiently large.}$$

where M is a positive constant, not necessarily less than 1.

Example

An example is the sequence $1 + (0.5)^{2^k}$

Rate of Convergence

- The speed of convergence depends on r and (more weakly) on M , whose values depend not only on the algorithm but also on the properties of the particular problem.
- Regardless of these values, however, a quadratically convergent sequence will always eventually converge faster than a linearly convergent sequence.
- Obviously, any sequence that converges Q-quadratically also converges Q-superlinearly, and any sequence that converges Q-superlinearly also converges Q-linearly.
- Higher rates of convergence (cubic, quartic, and so on) can also be defined

Q-order of convergence is p

We say that the Q-order of convergence is p (with $p > 1$) if there is a positive constant M such that

$$\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^p} \leq M, \quad \text{for all } k \text{ sufficiently large.}$$

Convergence of Line Search Methods

- Designing optimization algorithms with good convergence properties may seem to be very easy.
 - Since all we need to ensure is that the search direction p_k does not tend to become orthogonal to the gradient ∇f_k , or that steepest descent steps are taken regularly.
 - One could also simply compute $\cos \theta_k$ at every iteration and turn p_k toward the steepest descent direction if $\cos \theta_k$ is smaller than some preselected constant $\delta > 0$
 - Angle tests of this type ensure global convergence, but they are undesirable for two reasons.
- 1 First, they may impede a fast rate of convergence
 - 2 Second, angle tests destroy the invariance properties of quasi-Newton methods.

CONVERGENCE RATE OF STEEPEST DESCENT

- Consider the ideal case, in which the objective function is quadratic and the line searches are exact.
- Let us suppose

$$f(x) = \frac{1}{2}x^T Qx - b^T x,$$

where Q is symmetric and positive definite.

- The gradient is given by

$$\nabla f(x) = Qx - b$$

- The minimiser x^* is the unique solution of the linear system

$$Qx = b.$$

CONVERGENCE RATE OF STEEPEST DESCENT

- To find the step length α_k at each iteration x_k one can exactly minimise the univariate function

$$\phi(\alpha) = f(x_k - \alpha \nabla f_k)$$

- Denote ∇f_k by g_k (gradient at x_k)

- We have
$$f(x_k - \alpha g_k) = \frac{1}{2}(x_k - \alpha g_k)^T Q(x_k - \alpha g_k) - b^T(x_k - \alpha g_k)$$

- Differentiating the above w.r.t α we get

$$g_k^T Q \alpha g_k - \frac{1}{2} g_k^T Q x_k - \frac{1}{2} x_k^T Q g_k + b^T g_k$$

- Equating the above to 0 we get

$$g_k^T Q \alpha g_k - x_k^T Q g_k + b^T g_k = 0$$

$$\implies g_k^T Q \alpha g_k = x_k^T Q g_k - b^T g_k = (x_k^T Q - b^T) g_k = \nabla f_k^T g_k$$

$$\implies \alpha_k = \frac{\nabla f_k^T \nabla f_k}{\nabla f_k^T Q \nabla f_k}$$

CONVERGENCE RATE OF STEEPEST DESCENT

- By using this exact minimiser α_k , we get the steepest descent iteration for the quadratic function f as

$$x_{k+1} = x_k - \left(\frac{\nabla f_k^T \nabla f_k}{\nabla f_k^T Q \nabla f_k} \right) \nabla f_k$$

- The above expression yields a closed form expression for x_{k+1} in terms of x_k .
- To quantify the rate of convergence let us introduce the weighted norm

$$\|x\|_Q^2 = x^T Q x$$

- We know $Qx^* = b$, x^* being the unique minimiser we get

$$\frac{1}{2} \|x - x^*\|_Q^2 = f(x) - f(x^*)$$

- So this norm measures the difference between the current objective value and the optimal value.

CONVERGENCE RATE OF STEEPEST DESCENT

- By using the closed form expression for x_{k+1} and noting the fact that $\nabla f_k = Q(x_k - x^*)$, we can derive the following identity

$$\|x_{k+1} - x^*\|_Q^2 = \left\{ 1 - \frac{(\nabla f_k^T \nabla f_k)^2}{(\nabla f_k^T Q \nabla f_k)(\nabla f_k^T Q^{-1} \nabla f_k)} \right\} \|x_k - x^*\|_Q^2$$

- This expression describes the exact decrease in f at each iteration.
- But since the term inside the brackets is difficult to interpret.
- It would be more useful to bound it (may be in terms of the condition number of the problem).

CONVERGENCE RATE OF STEEPEST DESCENT

Theorem

When the steepest descent method with exact line searches is applied to the strongly convex quadratic function the error norm satisfies

$$\|x_{k+1} - x^*\|_Q^2 \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 \|x_k - x^*\|_Q^2$$

where $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ are eigenvalues of Q .

- The above inequality show that the function values f_k converge to the minimum f^* at a linear rate.
- A special case is when all the eigenvalues are equal
$$\lambda_1 = \lambda_2 = \dots = \lambda_n$$

Then the convergence is achieved in just one step.

- In general, as the condition number $\kappa(Q) = \frac{\lambda_n}{\lambda_1}$ increases, the convergence degrades.

CONVERGENCE RATE OF STEEPEST DESCENT

The rate-of-convergence behaviour of the steepest descent method is essentially the same on general nonlinear objective functions.

Theorem

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable, and that the iterates generated by the steepest-descent method with exact line searches converge to a point x^* at which the Hessian matrix $\nabla^2 f(x^*)$ is positive definite. Let r be any scalar satisfying

$$r \in \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}, 1 \right),$$

where $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ are the eigenvalues of $\nabla^2 f(x^*)$. Then for all k sufficiently large, we have

$$f(x_{k+1}) - f(x^*) \leq r^2 [f(x_k) - f(x^*)].$$

Therefore, the steepest descent method can have an unacceptably slow rate of convergence, even when the Hessian is reasonably well conditioned.

CONVERGENCE RATE OF NEWTON'S METHOD

- Consider the Newton iteration, for which the search is given by

$$p_k^N = -\nabla^2 f_k^{-1} \nabla f_k \quad (14)$$

- Since the Hessian matrix $\nabla^2 f_k$ may not always be P.D. p_k^N may not always be a descent direction

$$x_{k+1} = x_k + \alpha_k p_k \quad (15)$$

- For all x in the vicinity of a solution point x^* s.t. $\nabla^2 f(x^*)$ is P.D. then the Hessian $\nabla^2 f(x)$ will also be P.D. .
- Newton's method will be well defined in this region.

CONVERGENCE RATE OF NEWTON'S METHOD

Theorem

Suppose that f is twice differentiable and that the Hessian $\nabla^2 f$ is Lipschitz continuous in a neighbourhood of a solution x^* at which the second order sufficient conditions are satisfied. Consider the iteration (15) where p_k is given by (14). Then

- 1 if the starting point x_0 is sufficiently close to x^* , the sequence of iterates converges to x^* ;
- 2 the rate of convergence of $\{x_k\}$ is quadratic; and
- 3 the sequence of gradient norms $\{\|\nabla f_k\|\}$ converges quadratically to zero.

CONVERGENCE RATE OF NEWTON'S METHOD

Sketch of proof:

- From the definition of the Newton step and the optimality condition $\nabla f(x^*) = 0$ we have

$$\begin{aligned}x_k + p_k^N - x^* &= x_k - x^* - \nabla^2 f_k^{-1} \nabla f_k \\ &= \nabla^2 f_k^{-1} [\nabla^2 f_k(x_k - x^*) - (\nabla f_k - \nabla f_*)].\end{aligned}$$

- From Taylor's theorem we have

$$\nabla f_k - \nabla f_* = \int_0^1 \nabla^2 f(x_k + t(x^* - x_k))(x_k - x^*) dt,$$

- Therefore, we have

$$\begin{aligned}& \|\nabla^2 f(x_k)(x_k - x^*) - (\nabla f_k - \nabla f_*)\| \\ &= \left\| \int_0^1 [\nabla^2 f(x_k) - \nabla^2 f(x_k + t(x^* - x_k))] (x_k - x^*) dt \right\|\end{aligned}$$

CONVERGENCE RATE OF NEWTON'S METHOD

$$\begin{aligned} &\leq \int_0^1 \|\ [\nabla^2 f(x_k) - \nabla^2 f(x_k + t(x^* - x_k))] (x_k - x^*) \| dt \\ &\leq \int_0^1 \|\nabla^2 f(x_k) - \nabla^2 f(x_k + t(x^* - x_k))\| \|(x_k - x^*)\| dt \end{aligned}$$

Now $\nabla^2 f$ is Lipschitz around x^* therefore we have:

$$\begin{aligned} &\|\nabla^2 f(x_k) - \nabla^2 f(x_k + t(x^* - x_k))\| \\ &\leq L \|t(x^* - x_k)\| \end{aligned}$$

Therefore, integral is

$$\leq \|x_k - x^*\|^2 \int_0^1 Lt dt = \frac{1}{2} L \|x_k - x^*\|^2 \quad (16)$$

CONVERGENCE RATE OF NEWTON'S METHOD

- Now $\nabla^2 f(x^*)$ is nonsingular and continuous.
- Therefore $(\nabla^2 f(x^*))^{-1}$ is defined and is continuous in at least a small neighbourhood of x^*
- Therefore, there exists a radius $r > 0$ s.t.

$$\|\nabla^2 f_k^{-1}\| \leq 2\|\nabla^2 f(x^*)^{-1}\| \quad \text{for all } x_k \text{ with } \|x_k - x^*\| \leq r.$$

- Substituting all the above in (16) we get

$$\begin{aligned} \|x_k + p_k^N - x^*\| &\leq L\|\nabla^2 f(x^*)^{-1}\| \|x_k - x^*\|^2 \\ &= \tilde{L}\|x_k - x^*\|^2 \end{aligned} \tag{17}$$

where $\tilde{L} = L\|\nabla^2 f(x^*)^{-1}\|$.

- By choosing x_0 such that $\|x_0 - x^*\| \leq \min\left(r, \frac{1}{2\tilde{L}}\right)$ we can use the above inequality to inductively deduce that the sequence converges to x^* , and the rate of convergence is quadratic.

CONVERGENCE RATE OF NEWTON'S METHOD

- Now by using $x_{k+1} - x_k = p_k^N$ and $\nabla f_k + \nabla^2 f_k p_k^N = 0$, we obtain:

$$\begin{aligned}
 \|\nabla f(x_{k+1})\| &= \|\nabla f(x_{k+1}) - \nabla f_k - \nabla^2 f(x_k) p_k^N\| \\
 &= \left\| \int_0^1 \nabla^2 f(x_k + t p_k^N) (x_{k+1} - x_k) dt - \nabla^2 f(x_k) p_k^N \right\| \\
 &= \left\| \int_0^1 \nabla^2 f(x_k + t p_k^N) p_k^N dt - \nabla^2 f(x_k) p_k^N \right\| \\
 &\leq \int_0^1 \|(\nabla^2 f(x_k + t p_k^N) - \nabla^2 f(x_k)) p_k^N\| dt \\
 &\leq \int_0^1 \|(\nabla^2 f(x_k + t p_k^N) - \nabla^2 f(x_k))\| \|p_k^N\| dt \\
 &\leq \int_0^1 L t \|p_k^N\|^2 dt \\
 &= \frac{L}{2} \|p_k^N\|^2 = \frac{L}{2} \|\nabla^2 f(x_k)^{-1} \nabla f_k\|^2 \\
 &\leq \frac{L}{2} \|\nabla^2 f_k^{-1}\|^2 \|\nabla f_k\|^2 \leq 2L \|\nabla^2 f(x_k^*)^{-1}\|^2 \|\nabla f_k\|^2
 \end{aligned}$$

CONVERGENCE RATE OF QUASI-NEWTON'S METHOD

Theorem

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable consider the iteration

$$x_{k+1} = x_k + \alpha_k p_k \quad (18)$$

where p_k is a descent direction and α_k satisfies the Wolfe conditions with $c_1 < \frac{1}{2}$. If the sequence $\{x_k\}$ converges to a point x^* such that $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive definite, and if the search direction satisfies:

$$\lim_{k \rightarrow \infty} \frac{\|\nabla f_k + \nabla^2 f_k p_k\|}{\|p_k\|} = 0 \quad (19)$$

then

- 1 the step length $\alpha_k = 1$ is admissible for all k greater than a certain index k_0 ; and
- 2 if $\alpha_k = 1$ for all $k > k_0$, $\{x_k\}$ converges to x^* superlinearly.

CONVERGENCE RATE OF QUASI-NEWTON'S METHOD

The search direction is given by

$$p_k = -B_k^{-1} \nabla f_k \quad (20)$$

where B_k is some approximation to the Hessian $\nabla^2 f_k$ and is symmetric and positive definite.

- If p_k is the quasi-Newton search direction then (19) is equivalent to

$$\lim_{k \rightarrow \infty} \frac{\|(B_k - \nabla^2 f^*)p_k\|}{\|p_k\|} = 0 \quad (21)$$

CONVERGENCE RATE OF QUASI-NEWTON'S METHOD

Theorem

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable. Consider the iteration

$$x_{k+1} = x_k + p_k$$

(that is, the step length α_k is uniformly 1) and that p_k is given by (20). Let us assume also that $\{x_k\}$ converges to a point x^* such that $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive definite. Then $\{x_k\}$ converges superlinearly if and only if (21) holds.

CONVERGENCE RATE OF QUASI-NEWTON'S METHOD

Sketch of Proof

- Consider

$$\begin{aligned} p_k - p_k^N &= \nabla^2 f_k^{-1}(\nabla^2 f_k p_k + \nabla f_k) & p_k^N &= -\nabla^2 f_k^{-1} \nabla f_k \\ &= \nabla^2 f_k^{-1}(\nabla^2 f_k p_k - B_k p_k) \\ &= \nabla^2 f_k^{-1}(\nabla^2 f_k - B_k)p_k \end{aligned}$$

- We assume that $\|\nabla^2 f_k^{-1}\|$ is bounded above for x_k sufficiently close to x^* .
- As result we have

$$p_k - p_k^N = \mathcal{O}(\|(\nabla^2 f_k - B_k)p_k\|)$$

- Asserting (21) is same as saying

$$\|(B_k - \nabla^2 f(x_k))p_k\| = o(\|p_k\|)$$

CONVERGENCE RATE OF QUASI-NEWTON'S METHOD

- Therefore,

$$p_k - p_k^N = o(\|p_k\|) \quad (22)$$

- Now if we multiply both sides of (22) by $\nabla^2 f_k$ and use the definition of the quasi-Newton direction we get (21).
- Therefore, (21) and (22) are equivalent for the quasi-Newton search direction.
- Now combining (17) and (22) we get

$$\|x_k + p_k - x^*\| \leq \|x_k + p_k^N - x^*\| + \|p_k - p_k^N\| = \mathcal{O}(\|x_k - x^*\|^2) + o(\|p_k\|) \quad (23)$$

- One can show that $\|p_k\| = \mathcal{O}(\|x_k - x^*\|)$, so we obtain

$$\|x_k + p_k - x^*\| \leq o(\|x_k - x^*\|) \quad (24)$$

giving superlinear convergence.