

## Unconstrained Optimization I

Note Title

3/25/2022

Let's recall some basic facts about unconstrained optimization:

**Theorem 2.6 (first order optimality condition for local optima points).** Let  $f : U \rightarrow \mathbb{R}$  be a function defined on a set  $U \subseteq \mathbb{R}^n$ . Suppose that  $\mathbf{x}^* \in \text{int}(U)$  is a local optimum point and that all the partial derivatives of  $f$  exist at  $\mathbf{x}^*$ . Then  $\nabla f(\mathbf{x}^*) = 0$ .

$\nabla f(\mathbf{x}^*) = 0$  called a stationary / critical pt of  $f$ .

**Proof.** Let  $i \in \{1, 2, \dots, n\}$  and consider the one-dimensional function  $g(t) = f(\mathbf{x}^* + t\mathbf{e}_i)$ .

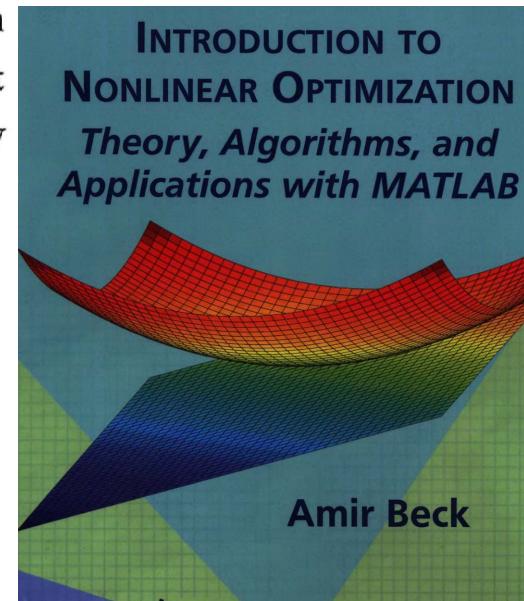
Note that  $g$  is differentiable at  $t = 0$  and that  $g'(0) = \frac{\partial f}{\partial x_i}(\mathbf{x}^*)$ . Since  $\mathbf{x}^*$  is a local optimum point of  $f$ , it follows that  $t = 0$  is a local optimum of  $g$ , which immediately implies that  $g'(0) = 0$ . The latter equality is exactly the same as  $\frac{\partial f}{\partial x_i}(\mathbf{x}^*) = 0$ . Since this is true for any  $i \in \{1, 2, \dots, n\}$ , the result  $\nabla f(\mathbf{x}^*) = 0$  follows.  $\square$

**Proposition 7.8 (sufficiency of stationarity under convexity).** Let  $f$  be a continuously differentiable function which is convex over a convex set  $C \subseteq \mathbb{R}^n$ . Suppose that  $\nabla f(\mathbf{x}^*) = 0$  for some  $\mathbf{x}^* \in C$ . Then  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $C$ .

**Proof.** Let  $\mathbf{z} \in C$ . Plugging  $\mathbf{x} = \mathbf{x}^*$  and  $\mathbf{y} = \mathbf{z}$  in the gradient inequality (7.6), we obtain that

$$f(\mathbf{z}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^T (\mathbf{z} - \mathbf{x}^*),$$

which by the fact that  $\nabla f(\mathbf{x}^*) = 0$  implies that  $f(\mathbf{z}) \geq f(\mathbf{x}^*)$ , thus establishing that  $\mathbf{x}^*$  is the global minimizer of  $f$  over  $C$ .  $\square$



**Theorem 2.27 (sufficient second order optimality condition).** Let  $f : U \rightarrow \mathbb{R}$  be a function defined on an open set  $U \subseteq \mathbb{R}^n$ . Suppose that  $f$  is twice continuously differentiable over  $U$  and that  $\mathbf{x}^*$  is a stationary point. Then the following hold:

- (a) If  $\nabla^2 f(\mathbf{x}^*) \succ 0$ , then  $\mathbf{x}^*$  is a strict local minimum point of  $f$  over  $U$ .
- (b) If  $\nabla^2 f(\mathbf{x}^*) \prec 0$ , then  $\mathbf{x}^*$  is a strict local maximum point of  $f$  over  $U$ .

**Theorem 7.12 (second order characterization of convexity).** Let  $f$  be a twice continuously differentiable function over an open convex set  $C \subseteq \mathbb{R}^n$ . Then  $f$  is convex if and only if  $\nabla^2 f(\mathbf{x}) \succeq 0$  for any  $\mathbf{x} \in C$ .

**Theorem 7.13 (sufficient second order condition for strict convexity).** Let  $f$  be a twice continuously differentiable function over a convex set  $C \subseteq \mathbb{R}^n$ , and suppose that  $\nabla^2 f(\mathbf{x}) \succ 0$  for any  $\mathbf{x} \in C$ . Then  $f$  is strictly convex over  $C$ .

From now on, assume that the objective function  $f: U \subseteq \mathbb{R}^n \xrightarrow{\text{open}} \mathbb{R}$  is at least  $C^1$ .

Our goal is to develop algorithms, and the analysis of such algorithms, for solving

$$\min_{x \in U} f(x).$$

Two strategies :

- (i) **line search** : (I) choose a descent direction  $p_k$  at  $x_k$  (ie. a  $p_k \in \mathbb{R}^n$  st.  $\frac{d}{d\alpha} f(x_k + \alpha p_k) < 0$ )  
(II) then, choose a step size  $\alpha_k$  that approximately solves  
$$\min_{\alpha > 0} f(x_k + \alpha p_k)$$
  
(III) update  $x_{k+1} = x_k + \alpha_k p_k$ , and iterate.

- (ii) **trust region** : (I) information gathered about  $f$  near  $x_k$  is used to construct a model function  $m_k$ , so that  $f(x) \approx m_k(x)$  for  $x \approx x_k$ .  
[Typically  $m_k$  is a quadratic approximation.]

(II) solve  $\min_p m_p(x_k + p)$  st.  $\|p\| \leq \Delta$  approximately  $\rightarrow p_k$

Note: this "trust region subproblem" is a constraint optimization problem. (Something surprising here...)

(III) If  $f(x_k + p_k)$  is significantly smaller than  $f(x_k)$ ,  
update  $x_{k+1} = x_k + p_k$ , and iterate.

else

reduces  $\Delta$ , goto (II).

In effect, a trust region method chooses the descent direction and step size simultaneously.

$f(x) = 10(x_2 - x_1^2)^2 + (1 - x_1)^2$ . At  $x_k = (0, 1)$ , its 2nd order Taylor approximation is:

$$f(x) \approx 11 + \begin{bmatrix} -2 \\ 20 \end{bmatrix}^T \begin{bmatrix} x_1 - 0 \\ x_2 - 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} x_1 - 0 \\ x_2 - 1 \end{bmatrix}^T \begin{bmatrix} -38 & 0 \\ 0 & 20 \end{bmatrix} \begin{bmatrix} x_1 - 0 \\ x_2 - 1 \end{bmatrix}$$

*a sum of squares, and non-convex*

$m_p(x) = \nabla f(x_k)$

$\nabla^2 f(x_k)$

$m_p(x)$

Note:  $\min f(x) = 0$ ,  $\operatorname{argmin} f(x) = (1, 1)$

$\min m_p(x) = -\infty$  (we definitely need to restrict  $x$  to a compact set in order to get a meaningful minimizer.)

Search direction for line search methods

How to choose  $p \in \mathbb{R}^n$  s.t.  $\frac{d}{d\alpha} f(x + \alpha p) = \nabla f(x)^T p < 0$  ?

The most obvious choice is  $p = -\nabla f(x)$ . If  $\nabla f(x) = 0$ ,  $x$  is already a stationary point. Otherwise

$$\nabla f(x)^T (-\nabla f(x)) = -\|\nabla f(x)\|^2 < 0.$$

moreover, it is ostensibly the best choice due to the Cauchy-Schwartz inequality:

$$|\langle \nabla f(x), p \rangle| \leq \|\nabla f(x)\|_2 \|p\|_2 \text{ with equality holds } \Leftrightarrow \|\nabla f(x)\|_2 \|p\|_2$$

so  $\operatorname{argmin}_{\|p\|_2=1} \langle \nabla f(x), p \rangle = -\nabla f(x) / \|\nabla f(x)\|_2$

For this reason,  $p = -\nabla f(x)$  is also called the "steepest descent direction".

If  $\|u\|_2 = \|v\|_2 = 1$   
 $\langle u, v \rangle = -1$   
 $\Leftrightarrow u = -v$

But I want to convince you that this nomenclature can be misleading, and I prefer the name "gradient descent".

Here is one way to see it: For any non-singular matrix  $A$ ,

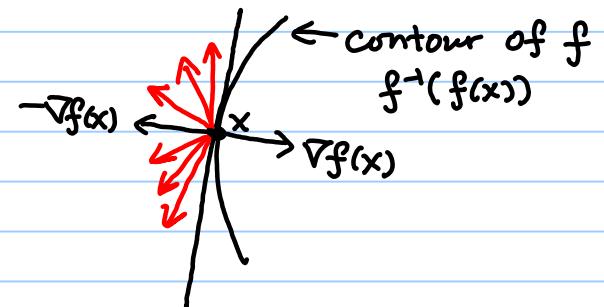
$$\frac{d}{d\alpha} f(x + \alpha p) = \nabla f(x)^T A A^T p = \langle A^T \nabla f(x), A^{-1} p \rangle$$

so  $\min_{\|A^T p\|=1} \frac{d}{d\alpha} f(x + \alpha p) = -\|A^T \nabla f(x)\|_2$ , the min is attained when and only when  $A^T p = -A^T \nabla f(x)$  and  $\|A^T p\|=1$

★  $\underset{P: \|A^T p\|=1}{\operatorname{argmin}} \frac{d}{d\alpha} f(x + \alpha p) = \frac{-A A^T \nabla f(x)}{\|A^T \nabla f(x)\|_2}$

Note :  $\{-A A^T \nabla f(x) : A \text{ nonsingular}\}$   
 $= \{d : d^T \nabla f(x) < 0\}$   
when  $\nabla f(x) \neq 0$ .

Ex : Prove it.



By choosing different  $A$ , any descent direction is "steepest" in the sense of ★.

Of course, you may find ★ contrived when we change the round sphere  $\{p : \|p\|_2=1\}$  into a ellipsoid  $\{p : \|A^{-1}p\|_2=1\}$ . But what it actually shows is that

the notion of "steepest descent direction" is not affine-invariant.

Imagine if we apply an affine change of coordinates to  $f$  and apply the "steepest descent direction" to  $g(\bar{x}) := f(A^{-1}\bar{x})$  at  $\bar{x} = Ax$ . What it means is:

In the original coordinate system:

$$x_{k+1} = x_k - s \underbrace{\nabla f(x_k)}_{\text{steepest}}$$

Chain rule:

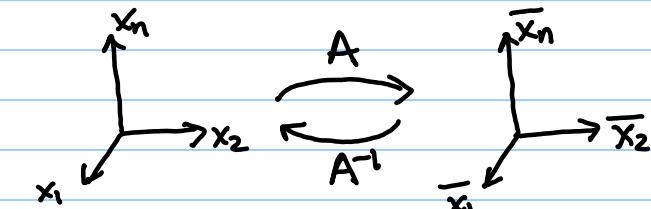
$$\begin{aligned} g &= f \circ A^{-1} \\ dg(\bar{x}) &= df(A^{-1}\bar{x}) \cdot A^{-1} \\ \nabla g(\bar{x}) &= A^{-T} \nabla f(A^{-1}\bar{x}) \end{aligned}$$

In the transformed coordinate system:

$$\bar{x}_{k+1} = \bar{x}_k - s \nabla g(\bar{x}_k)$$

$$x_{k+1} = x_k - s A^{-1} A^{-T} \nabla f(x_k)$$

Steepest in one coordinate system  
is not the steepest system in another  
coordinate system!



(assume constant size for now, other choices of step size do not really improve the situation by much as we shall see later.)

Ex: Prove that the constant step size gradient descent method is "invariant under rigid transformation of coordinates." ( $x \mapsto Ax + b$ ,  $A \in \mathbb{R}^{n \times n}$ ,  $AA^T = I$ ,  $b \in \mathbb{R}^n$ )

Here is yet another, arguably better, way to see why the "steepest descent" direction can be nowhere close to being steepest.

Consider  $f(x) = \frac{1}{2}x^T Ax + bx + c$ ,  $A \succ 0$ , which is pretty much (by Taylor's theorem) how any  $C^2$  function looks like in the vicinity of a local minimizer  $x^*$  with  $\nabla^2 f(x^*) \succ 0$ .

In virtue of the rigid motion invariance of GD and the spectral theorem

$$A = U \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} U^T$$

we may assume WLOG that  $A = \text{diag}(\lambda_1, \dots, \lambda_n)$ ,  $\lambda_i > 0$ ,  $b = 0$ ,  $c = 0$ , i.e.

$$f(x) = \frac{1}{2}x^T \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} x, \quad \lambda_i > 0.$$

In this case, the minimizer of  $f$  is at  $0$ . So at  $x$ , the "true steepest descent" direction is  $-x$ , not  $-\nabla f(x) = -Ax$ . Moreover, we can prove:

**Proposition 1.1** *The angle between  $-\nabla f(x)$  and  $-x$  can be arbitrarily close to  $90^\circ$ . In other words, the negative gradient direction can be as far from the true steepest descent direction as possible.*

**Proof:** The angle  $\theta$  between  $-\nabla f(x)$  and  $-x$  is maximized when  $\cos(\theta)$  is minimized, so we seek the solution of:

$$\min_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\langle -\nabla f(x), -x \rangle}{\|\nabla f(x)\| \|x\|} = \min_{x \in \mathbb{R}^n \setminus \{0\}} \frac{x^T A x}{\sqrt{x^T A^T A x} \sqrt{x^T x}} = \sqrt{\min_{x \in \mathbb{R}^n \setminus \{0\}} \frac{(x^T A x)^2}{(x^T A^2 x)(x^T x)}}.$$

Apply a change of variable  $y = \sqrt{A}x$ . The last minimum inside the square-root can be written as

$$\min_{y \in \mathbb{R}^n \setminus \{0\}} \frac{(y^T y)^2}{(y^T A y)(y^T A^{-1} y)}. \quad (1.8)$$

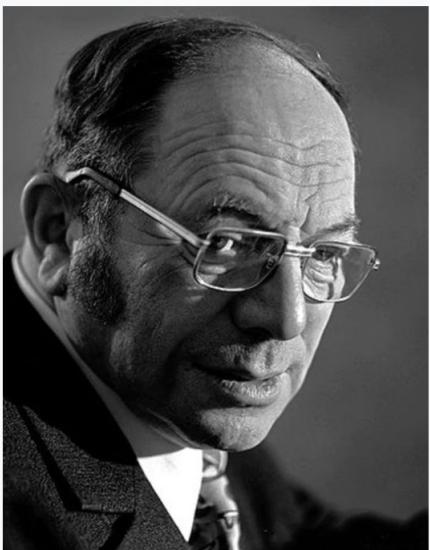
By the Kantorovich inequality, the minimum value above is  
*(see next page)*

$$\frac{4\lambda_1\lambda_n}{(\lambda_1 + \lambda_n)^2} = 1 - \frac{(\lambda_1 - \lambda_n)^2}{(\lambda_1 + \lambda_n)^2},$$

where  $\lambda_1$  and  $\lambda_n$  are the largest and smallest eigenvalues of  $A$ . This shows that  $\cos(\theta)$  gets arbitrarily close to 1, or equivalently  $\theta$  gets arbitrarily close to  $90^\circ$ , when  $A$  gets *increasingly ill-conditioned*, and  $x$  is chosen to be  $\sqrt{A}y$  where  $y$  is a direction that realizes the minimum value in (1.8). ■

Kantorovich inequality : Let  $A$  be a positive definite  $n \times n$  matrix. Then  
 $\forall 0 \neq x \in \mathbb{R}^n$ ,

Leonid Kantorovich  
Леонід Канторович



Leonid Kantorovich in 1975

$$\frac{(x^T x)^2}{(x^T A x)(x^T A^T x)} \geq \frac{4 \lambda_{\max}(A) \lambda_{\min}(A)}{(\lambda_{\max}(A) + \lambda_{\min}(A))^2}. \quad -(K)$$

WLOG, can assume  $A = \begin{matrix} \text{diag}(\lambda_1, \dots, \lambda_n) \\ \parallel \\ \lambda_{\min}(A) \end{matrix}$ ,  $\lambda_i > 0$ .  
 $\parallel \\ \lambda_{\max}(A)$

Note that if

- (i)  $x_i \neq 0 \Rightarrow \lambda_i = \lambda_1$  or  $\lambda_n$  and
- (ii)  $\sum_{i: \lambda_i = \lambda_1} x_i^2 = \sum_{i: \lambda_i = \lambda_n} x_i^2 (= A)$

$$\text{then, } \frac{(x^T x)^2}{(x^T A x)(x^T A^T x)} = \frac{(A+A)^2}{(\lambda_1 A + \lambda_n A)(\lambda_1^T A + \lambda_n^T A)} \\ = \frac{4}{(\lambda_1 + \lambda_n)(\lambda_1^T + \lambda_n^T)} = \frac{4 \lambda_1 \lambda_n}{(\lambda_1 + \lambda_n)^2}$$

The following proof of (K) shows that (i)+(ii) are not only sufficient for equality in (K), but also necessary. I learned this proof from:

P. Henrici. Two remarks on the Kantorovich inequality. *Amer. Math. Monthly*, 68:904–906, 1961.

Proof of the Kantorovich inequality :

Besides the diagonalization simplification, we may also assume  $x$  on the LHS of (K) has unit length. And (K) is equivalent to :

$$(\xi_1 \lambda_1 + \dots + \xi_n \lambda_n)(\xi_1 \lambda_1^{-1} + \dots + \xi_n \lambda_n^{-1}) \leq (\lambda_1 + \lambda_n)^2 / 4 \lambda_1 \lambda_n \text{ when } \sum \xi_i = 1, \xi_i \geq 0$$

If  $\lambda_1 = \lambda_n$ , then equality holds trivially in (K). So, let's assume  $\lambda_1 < \lambda_n$ .

A key trick is to observe that  $\forall i$ , we can always find  $p_i, q_i \geq 0$  st.

$$\begin{array}{c} \lambda_1 \quad \lambda_i \quad \lambda_n \\ \hline + \quad + \quad + \end{array} \quad \lambda_i = p_i \lambda_1 + q_i \lambda_n \quad \text{and} \quad \lambda_i^{-1} = p_i \lambda_1^{-1} + q_i \lambda_n^{-1}. \quad \begin{array}{c} \lambda_n^{-1} \quad \lambda_i^{-1} \quad \lambda_1^{-1} \\ \hline + \quad + \quad + \end{array}$$

(Just solve the  $2 \times 2$  linear system and see that the solution is non-negative.)

Then  $1 = \lambda_i \lambda_i^{-1} = (p_i \lambda_1 + q_i \lambda_n)(p_i \lambda_1^{-1} + q_i \lambda_n^{-1}) = (p_i + q_i)^2 + p_i q_i (\lambda_n - \lambda_1)^2 / (\lambda_1 \lambda_n)$ , which also implies  $p_i + q_i \leq 1$ .

Set  $p := \sum \xi_i p_i$ ,  $q = \sum \xi_i q_i$ . So  $p+q = \sum \xi_i (p_i + q_i) \leq \sum \xi_i = 1$

$$\begin{aligned} p \lambda_1 + q \lambda_n &= \sum \xi_i \lambda_i \\ p \lambda_1^{-1} + q \lambda_n^{-1} &= \sum \xi_i \lambda_i^{-1} \end{aligned}$$

$$\begin{aligned}
 & \sum \xi_i \lambda_i \sum \xi_i \lambda_i^{-1} = (p\lambda_1 + q\lambda_n)(p\lambda_1^{-1} + q\lambda_n^{-1}) = (p+q)^2 + pq \frac{(\lambda_n - \lambda_1)^2}{(\lambda_1 \lambda_n)}, \\
 & \leq (p+q)^2 \left[ 1 + (\lambda_n - \lambda_1)^2 / 4\lambda_1 \lambda_n \right] = (p+q)^2 \left[ (\lambda_n + \lambda_1)^2 / 4\lambda_1 \lambda_n \right] \stackrel{\text{Q.E.D.}}{\leq} (\lambda_n + \lambda_1)^2 / 4\lambda_1 \lambda_n.
 \end{aligned}$$

Ex: See from the proof above that equality in (K) holds iff  
Conditions (i) and (ii) (stated before the proof) hold.

Takeaway message : the "steepest descent" direction is far from steepest  
when  $\nabla^2 f(x)$  is ill-conditioned.  
(We shall establish rate-of-convergence results related  
to this.)

Newton direction is the direction that solves :  $\min_{p \in \mathbb{R}^n} m_k(p)$ , where  
 $m_k(p) = f(x_k) + \nabla f(x_k)^T p + \frac{1}{2} p^T [\nabla^2 f(x_k)] p$  is the 2nd order Taylor approximation  
to  $f(x_k + p)$ .

$m_k(p)$  has a unique minimizer when and only when  $\nabla^2 f(x_k) > 0$  (so the Newton  
direction is well-defined if  $f$  is  $C^2$  and strictly convex, or when  $x_k$  is close  
enough to a local minimizer  $x^*$  with  $\nabla^2 f(x^*) > 0$ .)

shorthand notation:

$$f_R = f(x_R)$$

$$\nabla f_R = \nabla f(x_R)$$

$$\nabla^2 f_R = \nabla^2 f(x_R)$$

etc.

$$\text{In this case, } p_R^N := \underset{p}{\operatorname{argmin}} m_R(p) = -\nabla^2 f_R^{-1} \nabla f_R$$

Facts about  $p_R^N$ :

$$(1) \text{ It is a descent direction: } \nabla f_R^T p_R^N = -\nabla f_R^T [\nabla^2 f_R]^{-1} \nabla f_R$$

Again we need the assumption  $\nabla^2 f_R \succ 0$ , this time to guarantee that  $p_R^N$  is a descent direction.

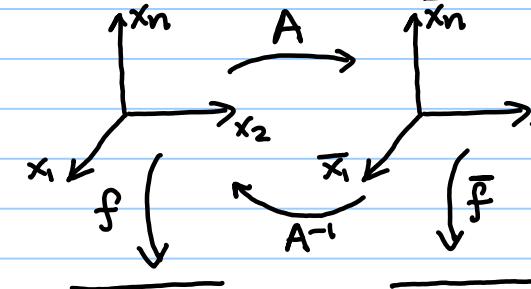
(2) It is invariant under affine change of coordinates!

In the original coordinate system:

$$x_{R+1} = x_R - s \nabla^2 f_R^{-1} \nabla f_R$$

In the transformed coordinate system:

$$\bar{x}_{R+1} = \bar{x}_R - s \nabla^2 \bar{f}_R^{-1} \nabla \bar{f}_R$$



Chain rule:

$$\bar{f} = f \circ A^{-1}$$

$$d\bar{f}(x) = df(A^{-1}x) \cdot A^{-1}$$

$$\nabla \bar{f}(x) = A^{-T} \nabla f(A^{-1}x)$$

$$d\nabla \bar{f}(x) = A^{-T} d\nabla f(A^{-1}x) \cdot A^{-1}$$

$$\nabla^2 \bar{f}(x)$$

$$\begin{aligned} x_{R+1} &= x_R - S A^{-1} A \nabla^2 f(x_R) A^T A^{-T} \nabla f_R \\ &= x_R - S \nabla^2 f_R^{-1} \nabla f_R \end{aligned}$$

$$\begin{aligned} \bar{f}(x) &= f(A^{-1}x) \\ \text{or } \bar{f} &= f \circ A^{-1} \end{aligned}$$

## GD

## Newton

- well-defined descent direction at all  $x$  (unless if  $x$  is already a critical point)
  - invariant under rigid transformation but not affine transformation
  - Step size determined from approx. solving  $\min_{\alpha} f(x + \alpha p)$
  - Computation requires only  $\nabla f(x)$   
 $n \times 1$
  - rate of convergence : linear error  $\rho = O(p^k)$   
depends on choice of  $\alpha_k$  and "conditioning" of  $f$
- not well-defined / not a descent dir when  $\nabla^2 f(x) \not\succeq 0$
  - affine invariant
  - "natural" step size 1 (more to say...)
  - Computation requires  $\nabla f(x)$ ,  $\nabla^2 f(x)$ , and solving  $-\nabla^2 f(x)p = \nabla f(x)$ .
  - rate of convergence : quadratic (a lot more to say...)

Quasi-Newton search directions provide an attractive alternative to Newton's method

- do not require computation of Hessian
- "Superlinear" rate of convergence

Let me motivate this method.

The Newton method is based on

$$\min_x f(x) \xrightarrow{\text{local quadratic approx.}} \min_p \frac{1}{2} p^T \nabla^2 f_{x_k} p + \nabla f_{x_k}^T p + f_{x_k} \quad (\text{Solve, and then iterate})$$

It is also a method for solving nonlinear system of equations

$$F(x) = 0 \xrightarrow{\text{local linear approx.}} F(x_k) + Df(x_k)(x - x_k) = 0 \quad (\text{Solve, and then iterate})$$

The latter applies to a general  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$ . In this case,  $Df$  needs not be a symmetric matrix.

When  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $F = \nabla f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ , the latter reduces to the former.

Ex : Check it. And recall that  $D\nabla f(x) = (\text{Hessian of } f \text{ at } x)$  is always symmetric.

Newton method for solving nonlinear equation in 1-D

$$x_{k+1} \text{ is the solution of } \underbrace{f_k + f'_k(x - x_k)}_{\approx f(x) \text{ near } x_k} = 0. \quad |x_{k+1} - x^*| \leq C |x_k - x^*|^2$$

Secant method :

$$x_{k+1} \text{ is the solution of } f_k + \frac{f_k - f_{k-1}}{x_k - x_{k-1}}(x - x_k) = 0. \quad |x_{k+1} - x^*| \leq C |x_k - x^*|^\alpha$$

$\alpha = (\sqrt{5} + 1)/2 = 1.618.$

Quasi-Newton methods for minimization is based on generalizing the secant method to higher-dimensions for solving  $\nabla f(x) = 0$ .

In 1-D,  $(f'_k - f'_{k-1})/(x_k - x_{k-1})$  is an approximation to  $f''_k$

In higher-dimensions, a quasi-Newton method chooses a Hessian approximation  $B_k$  to  $\nabla^2 f_k$  that satisfies the secant equation :  $B_k(x_k - x_{k-1}) = \nabla f_k - \nabla f_{k-1}$ .

In dimension  $n > 1$ , such a  $B_k$  is under-determined. (In contrast, if we have available  $x_{k-1}, \dots, x_{k-n}$ ,

$\nabla f_{k-1}, \dots, \nabla f_{k-n}$ , then (generically)  $\exists! B_k$  st.  $B_k [x_k - x_{k-1}, \dots, x_k - x_{k-n}] = [\nabla f_k - \nabla f_{k-1}, \dots, \nabla f_k - \nabla f_{k-n}]$ .

But this is too expensive to compute and quasi-Newton methods do not do that.)

We shall see that : with certain (economic) choices of  $B_k$  that satisfy the secant equation, the resulted iterative method enjoys **superlinear convergence** :

$$\|e_{k+1}\|/\|e_k\| \rightarrow 0. \quad (\text{SL})$$

Note that (SL) is weaker than saying

$$\|e_{k+1}\| \leq C \|e_k\|^{\alpha} \text{ for some } \alpha > 1. \quad (\text{Strong-SL})$$

We know that when  $n=1$ , (Strong-SL) holds with  $\alpha = 1.618$ .

When  $n>1$ , even the existence of an exponent  $\alpha > 1$  is not guaranteed by the theory (and if it exists, one would expect it to deteriorate to 1 as  $n \uparrow$ .)

In contrast, Newton's method enjoys quadratic convergence  $\|e_{k+1}\| \leq C \|e_k\|^2$  in any dimension, at the cost of much more expensive computation per iteration.

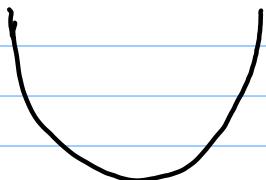
Surprisingly, quasi-Newton methods are much more effective than GD.

Moreover, a version of Quasi-Newton method under the name of **BFGS** is widely used.

We shall present Quasi-Newton methods in details, Several aspects (choice of  $B_k$ , choice of step size, etc) are quite tricky, and it is a tour de force to put all the pieces together.

---

### Choice of step size



Let's begin with  $f(x) = \frac{1}{2}x^2$  in 1-D.

$$\begin{aligned} f'(x) &= x . \quad x_{k+1} = x_k - \alpha_k x_k \quad \leftarrow \text{GD with step size } \alpha_k \\ &= (1 - \alpha_k)x_k \end{aligned}$$

$\alpha_k \in (-1, 1)$  guarantees a decrease in  $|x_k|$ , and hence also in  $f(x_k)$ , in each step.

But it is insufficient to guarantee that  $x_n \xrightarrow{\parallel} 0 = \underset{\parallel}{\arg\min} f(x)$ .

$$(1 - \alpha_1)(1 - \alpha_2) \cdots (1 - \alpha_n)x_0$$

When  $\alpha_i \in [0, 1]$ ,

$$\prod_{i=1}^{\infty} (1 - \alpha_i) > 0 \iff \sum_{i=1}^{\infty} \alpha_i \text{ converges}$$

This means  $\lim_{n \rightarrow \infty} x_n \neq 0$  if  $\alpha_i$  is too small (e.g.  $\alpha_i = 1/(i+1)^2$ )

Three popular choices of step size

- constant step size  $\alpha_k = \bar{\alpha} \quad \forall k$
- exact line search  $\alpha_k \in \operatorname{argmin}_{\alpha \geq 0} f(x_k + \alpha p_k)$
- backtracking (see below)

A step size  $\alpha_k$  satisfies the **Sufficient decrease / Armijo condition** if :

$$f(x_k + \alpha p_k) \leq f(x_k) + C \alpha \nabla f(x_k)^T p_k \quad \text{for some } C \in (0,1)$$

If  $p_k$  is a descent direction, we can always find a  $\alpha > 0$  st.  $f(x_k + \alpha p_k) < f(x_k)$ .  
The sufficient decrease condition is stronger than that. Still, such a step size is guaranteed to exist :

Lemma Let  $f: U \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$  be  $C^1$ ,  $p$  is a descent direction at  $x$ . For any  $C \in (0,1)$ ,  
 $\exists \varepsilon > 0$  st.

$$f(x) - f(x + \alpha p) \geq -C \alpha \nabla f(x)^T p \quad , \quad \forall \alpha \in [0, \varepsilon].$$

Proof:  $f(x + \alpha p) = f(x) + \alpha \nabla f(x)^T p + o(\alpha \|p\|)$  (differentiability at  $x$ )

$$f(x) - f(x + \alpha p) = -c\alpha \nabla f(x)^T p - \underbrace{(1-c)\alpha \nabla f(x)^T p}_{>0 \text{ for } \alpha > 0} - o(\alpha \|p\|)$$

$\Downarrow$

$$\text{so } f(x) - f(x + \alpha p) \geq -c\alpha \nabla f(x)^T p \quad \text{for small } \alpha > 0.$$

Q.E.D.

### Algorithm 3.1 (Backtracking Line Search).

Choose  $\bar{\alpha} > 0$ ,  $\rho \in (0, 1)$ ,  $c \in (0, 1)$ ; Set  $\alpha \leftarrow \bar{\alpha}$ ;

**repeat** until  $f(x_k + \alpha p_k) \leq f(x_k) + c\alpha \nabla f_k^T p_k$

$\alpha \leftarrow \rho\alpha$ ;

**end (repeat)**

Terminate with  $\alpha_k = \alpha$ .

Start with a big  $\bar{\alpha}$ , shrink it by a factor of  $\rho$  until the C-Armijo condition is satisfied.

A subtle point : The "sufficient decrease condition" alone cannot guarantee sufficient decrease in the objective value, because, as the lemma above shows, any small enough step size satisfies it. Put differently, the condition alone does not prohibit the step size from getting arbitrarily close to 0.

But the backtracking algorithm does !

This means : In order to establish a convergence result for a line search method based on backtracking , the proof must exploit the specific feature of the backtracking algorithm instead of simply invoking the Armijo Condition.

Turns Out it can be done for GD , ie. when  $P_k = -\nabla f(x_k)$ . ( Stay tuned ! )

But for more general descent methods , it is probably easier/necessary to strengthen the Armijo condition so that the stronger condition prohibits the step size to get arbitrarily close to 0 .

One such stronger condition requires :

P. Wolfe

$$f(x_k + \alpha_k P_k) \leq f(x_k) + C_1 \alpha_k \nabla f_k^T P_k$$

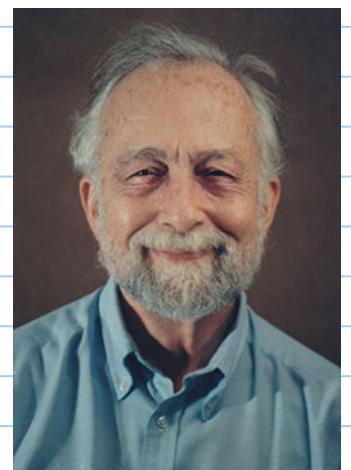
$$\nabla f(x_k + \alpha_k P_k)^T P_k \leq C_2 \nabla f_k^T P_k$$

with

$$0 < C_1 < C_2 < 1$$

$\leftarrow$  Armijo cond.  
 $\leftarrow$  "Curvature cond."

Known collectively as  
the Wolfe Conditions .



Strong Wolfe conditions :

$$f(x_k + \alpha_k P_k) \leq f(x_k) + C_1 \alpha_k \nabla f_k^T P_k$$

$$|\nabla f(x_k + \alpha_k P_k)^T P_k| \leq C_2 |\nabla f_k^T P_k| , \quad 0 < C_1 < C_2 < 1$$

Note :  $\nabla f(x_k + \alpha p_k)^T p_k = \phi'(\alpha)$  where  $\phi(\alpha) = f(x_k + \alpha p_k)$ .

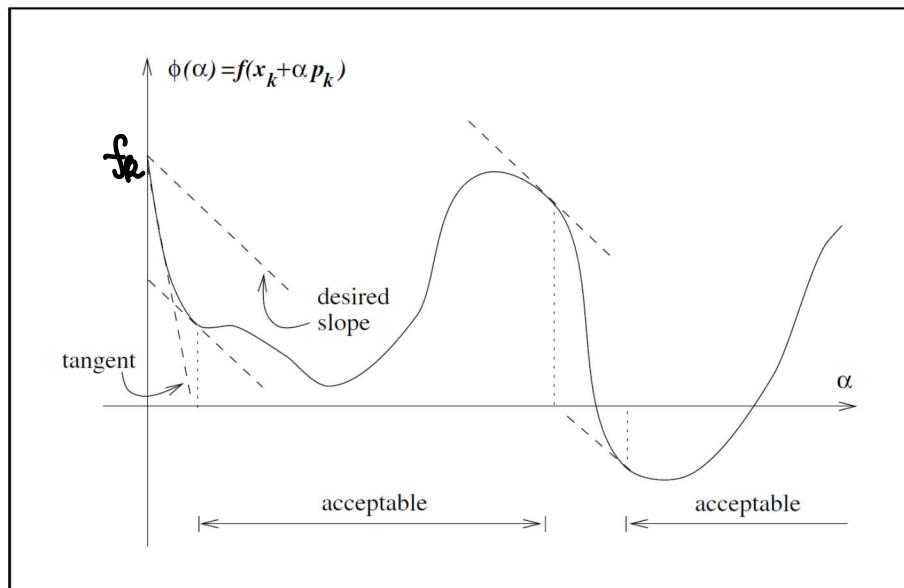


Figure 3.4 The curvature condition.

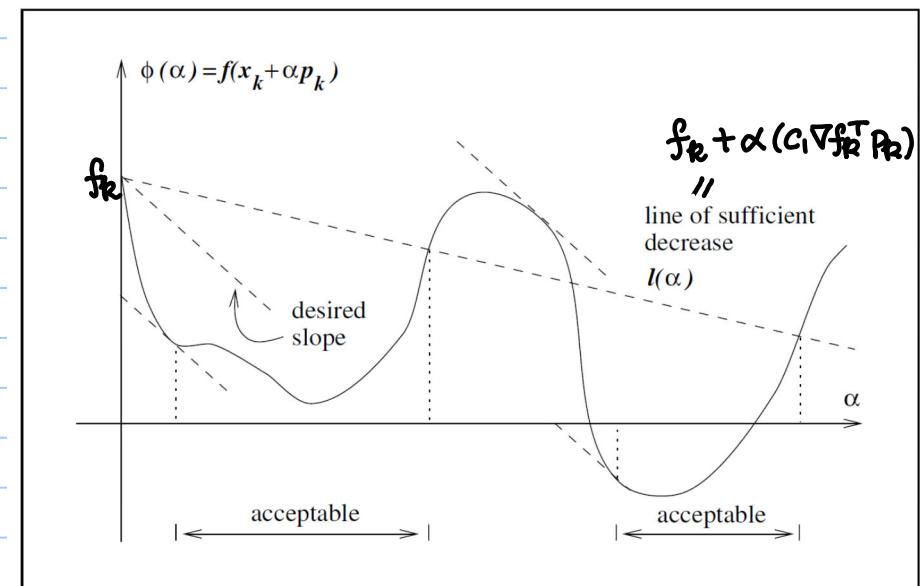


Figure 3.5 Step lengths satisfying the Wolfe conditions.

### **Lemma 3.1.**

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  $\{x_k + \alpha p_k | \alpha > 0\}$ . Then if  $0 < c_1 < c_2 < 1$ , there exist intervals of step lengths satisfying the Wolfe conditions (3.6) and the strong Wolfe conditions (3.7).

Proof : See N&W , pg 35.

The proof is not hard , but constructing an algorithm for the Wolfe conditions is tricky.

We shall establish the following results :

- Convergence of GD , without assuming "strong convexity" near a critical point , for all 3 choices of step sizes (constant , exact , backtracking)
- rate of convergence results of GD under "strong convexity" assumption .
- a convergence result of general descent method with step size satisfying the Wolfe Conditions
- various rate of convergence results for Newton and Quasi-Newton methods , after we

explained these methods. (We shall see, e.g., that backtracking does not work well with quasi-Newton methods and we need step size selection algorithms that satisfy Wolfe conditions.)