



Continuous Optimization

CO 466



Levent Tunçel

Preface

Disclaimer Much of the information on this set of notes is transcribed directly/indirectly from the lectures of CO 466 during Fall 2020 as well as other related resources. I do not make any warranties about the completeness, reliability and accuracy of this set of notes. Use at your own risk.

$$\mathbb{Z}_+ = \{0, 1, 2, 3, \dots\}; \mathbb{Z}_{++} = \{1, 2, 3, \dots\}$$

For any questions, send me an email via <https://notes.sibeliusp.com/contact/>.

You can find my notes for other courses on <https://notes.sibeliusp.com/>.

Sibelius Peng

Contents

Preface	1
1 Introduction	3
1.1 Conic Form	4
1.2 Derivatives	5
1.3 Fixed Point	6
1.4 Other	7
2 Unconstrained Continuous Optimization	13
2.1 Affine Subspace Constraints	15
2.2 Applications	16
2.3 Prototype low-rank approximation problem	16
2.4 Classical Algorithmic Approaches	17
2.5 Convergence Properties of Descent Algorithms	19
2.6 A General Conversation about Convergence	21
2.7 Quasi-Newton Methods	25
2.7.1 Convergence Results	28

Introduction: Formulations, fundamental background and definitions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}, g : \mathbb{R}^n \rightarrow \mathbb{R}^m, h : \mathbb{R}^n \rightarrow \mathbb{R}^p$ all continuous.

$$\begin{array}{ll} \inf & f(x) \\ \text{s.t.} & g(x) \leq 0 \\ & h(x) = 0 \end{array} \quad (\text{P})$$

$$S := \{x \in \mathbb{R}^n : g(x) \leq 0, h(x) = 0\}$$

is called the feasible solution set of (P), equivalently feasible region of (P).

Definition 1: global minimizer

$\bar{x} \in \mathbb{R}^n$ is a **global minimizer** of (P) if $x \in S$ and $f(x) \geq f(\bar{x})$ for all $x \in S$.

Sometimes we simply say \bar{x} is a minimizer of (P).

Definition 2: local minimizer

$\bar{x} \in \mathbb{R}^n$ is a **local minimizer** of (P) if $\bar{x} \in S$ and there exists a neighborhood U of \bar{x} such that

$$f(x) \geq f(\bar{x}) \quad \forall x \in S \cap U$$

$\bar{x} \in \mathbb{R}^n$ is a **strict local minimizer** of (P) if $\bar{x} \in S$ and there exists a neighborhood U of \bar{x} such that

$$f(x) > f(\bar{x}) \quad \forall x \in (S \cap U) \setminus \{\bar{x}\}$$

$\bar{x} \in \mathbb{R}^n$ is an **isolated local minimizer** of (P), if $\bar{x} \in S$ and there exists a neighborhood U of \bar{x} such that \bar{x} is the only local minimizer of (P) in $(S \cap U)$.

continuous optimization problem

A **continuous optimization problem** is a problem of optimizing (minimizing or maximizing) a continuous function of finitely many real variables subject to finitely many equations and inequalities on continuous functions of these variables.

What kind of problems can be formulated as Continuous Optimization problems? Almost everything.

Example 3: Fermat's Last Theorem

"There do not exist positive integers x, y, z and an integer $n \geq 3$ such that $x^n + y^n = z^n$."

Consider

$$\begin{aligned} \inf \quad & f(x) := \left(x_1^{x_4} + x_2^{x_4} - x_3^{x_4}\right)^2 + \sum_{i=1}^4 \left(\sin(\pi x_i)\right)^2 \\ \text{s.t.} \quad & g_1(x) := 1 - x_1 \leq 0 \\ & g_2(x) := 1 - x_2 \leq 0 \\ & g_3(x) := 1 - x_3 \leq 0 \\ & g_4(x) := 3 - x_4 \leq 0 \end{aligned} \tag{P}$$

The optimal objective value of (P) is zero and attained if and only if FLT is false.

We can show that (P) has a sequence of feasible solutions $\{x^{(k)}\}$ such that $f(x^{(k)}) \searrow 0$. Since $f(x) \geq 0$ for all $x \in \mathbb{R}^4$, the optimal value of (P) is zero.

FLT is true if and only if (P) does not attain its optimal value (of zero).

Even when the number of variables in a continuous optimization problem is very small (e.g., 4) the optimization problem may be notoriously hard. Even discrete structures can be formulated in our framework. $\sin(\pi x_1) = 0 \iff x_1 \in \mathbb{Z}$. In Example 3, we have functions that are "highly nonlinear".

Example 4: Combinatorial Optimization, 0,1 Integer Programming

Let m, n be positive integers, $A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m, c \in \mathbb{R}^n$ be given. Consider the 0,1 Integer Programming problem.

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & Ax \leq b \\ & x \in \{0, 1\}^n \end{aligned} \tag{IP}$$

The first condition can be written as

$$g(x) := Ax - b \leq 0.$$

The second condition can be written as

$$x_j(x_j - 1) = 0 \quad \forall j \in [n] \leftrightarrow h(x) = 0$$

Our continuous optimization problem is only mildly nonlinear.

Some conclusions from Example 3 and 4

Continuous Optimization problems can be very hard even when the number of variables and constraints are both small, the nonlinearity in f, g, h is very mild.

To successfully solve Continuous Optimization Problems, we must study the problem class at hand, discover special properties and structures and then exploit these special properties & structures.

1.1 Conic Form**Definition 5: cone**

A set $K \subseteq \mathbb{R}^n$ is a **cone** if $\forall x \in K, \forall k \in \mathbb{R}_+, \lambda x \in K$.

Definition 6: convex

A set $S \subseteq \mathbb{R}^n$ is **convex** if for every pair of points in S , the line segment joining them lies entirely in S .

That is S is convex if $\forall u, v \in S, \forall \lambda \in [0, 1], [\lambda u + (1 - \lambda)v] \in S$.

Definition 7: convex cone

A set $K \subseteq \mathbb{R}^n$ is a **convex cone** if it is convex and is a cone.

Let $g : \mathbb{R}^n \rightarrow \mathbb{R}^m, f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuous functions. Consider

$$\begin{aligned} \inf \quad & f(x) \\ \text{s.t.} \quad & g(x) \preceq_K 0 \end{aligned}$$

where $K \subseteq \mathbb{R}^m$ is a convex cone and for every $u, v \in \mathbb{R}^m$, $u \succeq_K v$ means $(u - v) \in K$.

This is at least as general as our original (P). Consider $\mathbb{R}^p \ni K := \mathbb{R}_+^m \oplus \{0\} \dots$

1.2 Derivatives

Definition 8: directional derivative

The **directional derivative** of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at $\bar{x} \in \mathbb{R}^n$ along the direction $d \in \mathbb{R}^n$ is

$$f'(\bar{x}; d) := \lim_{\alpha \searrow 0} \frac{f(\bar{x} + \alpha d) - f(\bar{x})}{\alpha}$$

(Gâteaux (directional) derivative)

Exercise:

What is the directional derivative of $f : \mathbb{R}^n \rightarrow \mathbb{R}, f(x) := \|x\|_\infty$, for every $\bar{x}, d \in \mathbb{R}^n$?

Definition 9: differentiable

$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **differentiable** at $\bar{x} \in \mathbb{R}^n$ if $\exists \mathcal{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, linear, such that

$$\lim_{\substack{h \rightarrow 0 \\ (h \in \mathbb{R}^n)}} \frac{\|f(\bar{x} + h) - [f(\bar{x}) + \mathcal{A}(h)]\|}{\|h\|} = 0$$

Such \mathcal{A} is called the derivative of f at \bar{x} and is denoted by $Df(\bar{x})$ or $f'(\bar{x})$ (matrix representation of $Df(\bar{x})$). We will also use $\nabla f(\bar{x}) := [f'(x)]^T$.

Suppose $f : \mathbb{E}_1 \rightarrow \mathbb{E}_2$, we have

$$\begin{aligned} Df(\bar{x}) &\in \mathcal{L}(\mathbb{E}_1, \mathbb{E}_2), & Df : \mathbb{E}_1 &\rightarrow \mathcal{L}(\mathbb{E}_1, \mathbb{E}_2) \\ D^2 f(\bar{x}) &\in \mathcal{L}(\mathbb{E}_1, \mathcal{L}(\mathbb{E}_1, \mathbb{E}_2)), & D^2 f : \mathbb{E}_1 &\rightarrow \mathcal{L}(\mathbb{E}_1, \mathcal{L}(\mathbb{E}_1, \mathbb{E}_2)) \end{aligned}$$

If $f : \mathbb{R}^n \rightarrow \mathbb{R}$, then $D^k f(\bar{x})[h^{(1)}, h^{(2)}, \dots, h^{(k)}] : k^{\text{th}}$ differential (derivative) along the directions $h^{(1)}, h^{(2)}, \dots, h^{(k)} \in \mathbb{R}^n$.

Theorem 10: Taylor's Theorem

Let $U \subseteq \mathbb{R}^n$ be open, $f : U \rightarrow \mathbb{R}$ be a \mathcal{C}^r function on U . Let $x, d \in \mathbb{R}^n$. If $x, (x + d)$, and the line segment joining x and $(x + d)$ lie in U , then there exists $z \in (x, x + d)$ such that

$$f(x + d) = f(x) + \sum_{k=1}^{r-1} \frac{1}{k!} D^k f(x) \underbrace{[d, d, \dots, d]}_{k\text{-times}} + \frac{1}{r!} D^r f(z) \underbrace{[d, d, \dots, d]}_{r\text{-times}}$$

Definition 11: contraction mapping

Let $U \subseteq \mathbb{R}^n$ be a closed set. $f : U \rightarrow U$ is called a **contraction mapping** if there exists $\lambda \in [0, 1)$ such that

$$\|f(x) - f(y)\| \leq \lambda \|x - y\| \quad \forall x, y \in U$$

1.3 Fixed Point

Theorem 12: Banach Fixed Point Theorem (1922)

Let $U \subseteq \mathbb{R}^n$ be a closed set and let $f : U \rightarrow U$ be a contraction mapping, then

(i) (Existence and uniqueness of solution - fixed point)

the mapping f has a unique fixed point $\bar{x} \in U$.

(ii) (Algorithm and convergence)

For all $x^{(0)} \in U$, the sequence $\{x^{(k)}\}$ generated by $x^{(k+1)} := f(x^{(k)})$, $k \in \{0, 1, 2, \dots\}$ (fixed point iteration) converges to \bar{x} . In particular,

$$\|x^{(k)} - \bar{x}\| \leq \lambda^k \|x^{(0)} - \bar{x}\| \quad \forall k \in \{0, 1, 2, \dots\}$$

Proof:

Suppose $U \subseteq \mathbb{R}^n$ is a nonempty closed set, and $f : U \rightarrow U$ is a contraction mapping with $\lambda \in [0, 1)$.

Let $x^{(k+1)} := f(x^{(k)})$ for all $k \in \mathbb{Z}_+$. Then for all $k \in \mathbb{Z}_+$,

$$\|x^{(k+1)} - x^{(k)}\| = \|f(x^{(k)}) - f(x^{(k-1)})\| \leq \lambda \|x^{(k)} - x^{(k-1)}\|$$

By induction on k , ... we obtain

$$\|x^{(k)} - x^{(k-1)}\| \leq \lambda^k \|x^{(1)} - x^{(0)}\| \quad \forall k \in \mathbb{Z}_+ \quad (*)$$

$\forall m \in \mathbb{Z}_{++}, \forall k \in \mathbb{Z}_{++}$,

$$\begin{aligned} \|x^{(m+k)} - x^{(m)}\| &= \|x^{(m+k)} - x^{(m+k-1)} + x^{(m+k-1)} - x^{(m+k-2)} + \dots + x^{(m+1)} - x^{(m)}\| \\ &\leq \sum_{i=1}^k \|x^{(m+i)} - x^{(m+i-1)}\| && \triangle\text{-ineq} \\ &\leq (\lambda^{m+k-1} + \lambda^{m+k-2} + \dots + \lambda^m) \|x^{(1)} - x^{(0)}\| && \text{by } (*) \\ &= \lambda^m (1 + \lambda + \lambda^2 + \dots + \lambda^{k-1}) \|x^{(1)} - x^{(0)}\| \\ &= \frac{\lambda^m (1 - \lambda^k)}{1 - \lambda} \|x^{(1)} - x^{(0)}\| \\ &\leq \frac{\lambda^m}{1 - \lambda} \|x^{(1)} - x^{(0)}\| \rightarrow 0 \text{ as } m \rightarrow +\infty \text{ (independent of } k) \end{aligned}$$

Therefore $\{x^{(k)}\}$ is a Cauchy sequence and hence it converges. Let \bar{x} be its limit. $\bar{x} \in U$ (U is closed).

$\forall k \in \mathbb{Z}_+$, we have

$$\|f(\bar{x}) - \bar{x}\| \leq \|f(\bar{x}) - x^{(k)}\| + \|x^{(k)} - \bar{x}\| \leq \lambda \underbrace{\|\bar{x} - x^{(k-1)}\|}_{\rightarrow 0} + \underbrace{\|x^{(k)} - \bar{x}\|}_{\rightarrow 0}$$

As $k \rightarrow +\infty$, RHS $\rightarrow 0$. Thus $f(\bar{x}) = \bar{x}$. This proves the existence. Now we prove the uniqueness.

Suppose $\exists \bar{x}, \bar{y} \in U$ such that $f(\bar{x}) = \bar{x}$ and $f(\bar{y}) = \bar{y}$. Then

$$\|\bar{x} - \bar{y}\| = \|f(\bar{x}) - f(\bar{y})\| \leq \lambda \|\bar{x} - \bar{y}\| \implies (1 - \lambda) \|\bar{x} - \bar{y}\| = 0 \xRightarrow{\lambda \in [0,1)} \bar{x} = \bar{y}$$

Now that we have established existence and uniqueness of \bar{x} , for a proof of convergence rate claim, we proceed as in the beginning of the proof. However, we use \bar{x} .

$$\|x^{(1)} - \bar{x}\| = \|f(x^{(0)}) - f(\bar{x})\| \leq \lambda \|x^{(0)} - \bar{x}\| \implies \|x^{(2)} - \bar{x}\| \leq \lambda^2 \|x^{(0)} - \bar{x}\|$$

By induction on k , we have

$$\|x^{(k)} - \bar{x}\| \leq \lambda^k \|x^{(0)} - \bar{x}\| \quad \forall k \in \mathbb{Z}_k$$

as desired. □

Theorem 13: Brouwer's Fixed Point Theorem (1910)

Let $U \subset \mathbb{R}^n$ be a nonempty, compact and convex set; let $f : U \rightarrow U$ continuous such that $f(U) = U$. Then there exists $\bar{x} \in U$ such that $f(\bar{x}) = \bar{x}$.

See the application in <https://n.sibp.ro/co/456>.

Theorem 14: Kakutani's Fixed Point Theorem (1941)

Let $U \subset \mathbb{R}^n$ be a nonempty, compact convex set and $f : U \rightarrow 2^U$ be a set valued map on U . If $\text{Graph}(f) := \left\{ \begin{pmatrix} x \\ v \end{pmatrix} \in U \oplus U : v \in f(x) \right\}$ is closed and $f(x)$ is nonempty and convex for every $x \in U$, then there exists $\bar{x} \in U$ such that $\bar{x} \in f(\bar{x})$.

Theorem 15: Borsuk-Ulam Theorem (1930-1933)

Let $f : \{x \in \mathbb{R}^{n+1} : \|x\|_2 = 1\} \rightarrow \mathbb{R}^n$ be continuous. Then there exists $\bar{x} \in \mathbb{R}^{n+1}$ such that $\|\bar{x}\|_2 = 1$ and $f(\bar{x}) = f(-\bar{x})$.

Example:

Let $n := 2$. Assuming temperature and barometric air pressure are continuous functions on the Earth's surface, and Earth's surface is homeomorphic to a sphere, there always exists an antipodal pair of points on Earth with the same temperature & the same air pressure.

1.4 Other

$S^n := n \times n$ symmetric matrices with real entries.

Theorem 16: Spectral Decomposition Theorem

For every $A \in S^n$, there exists $Q \in \mathbb{R}^{n \times n}$ orthogonal ($Q^T Q = I$) such that $A = QDQ^T$, where $D \in \mathbb{R}^{n \times n}$ is a diagonal matrix.

In the above theorem, the diagonal matrix D contains all eigenvalues of A , and the columns of Q are the corresponding eigenvectors of A .

Definition 17: positive definite

$A \in \mathbb{R}^{n \times n}$ is **positive semidefinite** if $h^T A h \geq 0$ for all $h \in \mathbb{R}^n$; such A is **positive definite** if $h^T A h > 0$ for all $h \in \mathbb{R}^n \setminus \{0\}$.

If $A \in \mathbb{R}^{n \times n}$ is skew-symmetric ($A = -A^T$), then $h^T A h = (h^T A h)^T = -h^T A h = 0$ for all $h \in \mathbb{R}^n$. Therefore, such A is positive semidefinite but not positive definite.

$$S_+^n := \text{positive semidefinite matrices in } S^n,$$

$$S_{++}^n := \text{positive definite matrices in } S^n.$$

In fact, $S_{++}^n = \text{int}(S_+^n)$.

Theorem 18: Choleski Decomposition Theorem

Let $A \in S^n$, then

- (a) A is positive semidefinite if and only if there exists $L \in \mathbb{R}^{n \times n}$ lower triangular such that $A = LL^T$;
- (b) A is positive definite if and only if there exists $L \in \mathbb{R}^{n \times n}$ lower triangular and nonsingular such that $A = LL^T$.

Note that Taylor's Theorem (Theorem 10) cannot be completely generalized to functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $m \geq 2$, even for $r = 1$. However, we have

Theorem 19

Let $U \subseteq \mathbb{R}^n$ be an open set and $f : U \rightarrow \mathbb{R}^m$ be a C^1 on U . Suppose for $\bar{x}, d \in \mathbb{R}^n$, $[\bar{x}, \bar{x} + d] \subset U$. Then

$$f(\bar{x} + d) - f(\bar{x}) = \int_0^1 Df(\bar{x} + \alpha d) d(\partial \alpha)$$

A consequence of this result is obtained when $Df(\cdot)$ is Lipschitz continuous on U (in a neighborhood of $[\bar{x}, \bar{x} + d]$ suffices). Let L denote the Lipschitz constant. Then

$$\|Df(x) - Df(y)\| \leq L\|x - y\| \quad \forall x, y \in U$$

Then we have

$$\begin{aligned} \|f(\bar{x} + d) - f(\bar{x}) - Df(\bar{x})d\|_2 &= \left\| \int_0^1 [Df(\bar{x} + \alpha d) - Df(\bar{x})] d(\partial \alpha) \right\|_2 \\ &\leq \int_0^1 \|Df(\bar{x} + \alpha d) - Df(\bar{x})\|_2 \cdot \|d\|_2(\partial \alpha) \quad \text{see aside} \\ &\leq \int_0^1 L\|d\|_2 \cdot \|d\|_2(\partial \alpha) \quad \begin{array}{c} \uparrow \\ \text{operator} \\ \text{2-norm} \end{array} \quad \begin{array}{c} \uparrow \\ \text{2-norm} \\ \text{on } \mathbb{R}^n \end{array} \\ &= \frac{1}{2} L \|d\|_2^2 \end{aligned}$$

So, if $\|d\|_2 < \epsilon$, then this error in this first-order estimate of $f(\bar{x} + d)$ is bounded above by $\frac{1}{2}L\epsilon^2$.

Aside:

Let $h := \int_0^1 [Df(\bar{x} + \alpha d) - Df(\bar{x})] d(\partial\alpha)$, then

$$\begin{aligned} \|h\|_2^2 &= h^T h = h^T \int_0^1 [Df(\bar{x} + \alpha d) - Df(\bar{x})] d(\partial\alpha) \\ &= \int_0^1 h^T [Df(\bar{x} + \alpha d) - Df(\bar{x})] d(\partial\alpha) \\ &\leq \int_0^1 \|h\|_2 \| [Df(\bar{x} + \alpha d) - Df(\bar{x})] d\|_2 d\alpha \quad \text{Cauchy-Schwarz} \\ &\implies \|h\|_2 \leq \int_0^1 \| [Df(\bar{x} + \alpha d) - Df(\bar{x})] d\|_2 d\alpha \end{aligned}$$

Note that we may replace f in Theorem 19 by $Df^r(\cdot)$ (assuming $f \in \mathcal{C}^{r+1}$) and apply the same reasoning. Indeed, Theorem 19 can be very useful in the design and analysis of continuous optimization algorithms.

Theorem 20: Inverse Function Theorem

Let $U \subseteq \mathbb{R}^n$ be open, $f : U \rightarrow \mathbb{R}^n$ be \mathcal{C}^1 , $\bar{x} \in U$, $\det(\nabla f(\bar{x})) \neq 0$. Then there exists an open neighborhood V of \bar{x} in U and an open neighborhood W of $f(\bar{x})$ such that

- $f(V) = W$,
- f has a local \mathcal{C}^1 inverse $f^{-1} : W \rightarrow V$,
- $\forall y \in W$, with $x = f^{-1}(y)$, we have $Df^{-1}(y) = [Df(x)]^{-1}$.

In the above, if f is \mathcal{C}^r , then there exists such an $f^{-1} \in \mathcal{C}^r$. Theorem 20 can be proved by utilizing Theorem 12 (in showing that the inverse is well-defined, i.e., one-to-one).

Theorem 21: Implicit Function Theorem

Let $h : \mathbb{R}^n \rightarrow \mathbb{R}^p$, $h \in \mathcal{C}^1$ in a neighborhood of $\bar{x} \in \mathbb{R}^n$ where $h(\bar{x}) = 0$. Suppose $h'(\bar{x})$ has full row rank ($\text{rank}(h'(\bar{x})) = p \leq n$). Define a partition $[B|N]$ of columns of $h'(\bar{x})$:

$$h' =: \begin{bmatrix} h'_B(\bar{x}) & | & h'_N(\bar{x}) \end{bmatrix}$$

\uparrow
 $\in \mathbb{R}^{p \times p}$,
 nonsingular,
 partition

\bar{x} and x with respect to the same $[B|N]$. Then there exist neighborhoods U_B of \bar{x}_B and U_N of \bar{x}_N and a \mathcal{C}^1 function $f : U_N \rightarrow U_B$ such that

- $f(\bar{x}_N) = \bar{x}_B$,
- $h \begin{pmatrix} x_B \\ x_N \end{pmatrix} = 0 \iff x_B = f(x_N)$ for all $x_B \in U_B, x_N \in U_N$.

Moreover, $f'(x_N) = -[h'_B(\bar{x})]^{-1}h'_N(\bar{x})$.

Recall the very special case (e.g., equality constraints in an LP problem): $A \in \mathbb{R}^{p \times n}$, $\text{rank}(A) = p$ given

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & Ax = b \\ & x \geq 0 \end{aligned}$$

$$\begin{aligned}
h(x) &:= Ax - b \implies h'(x) = A \\
\bar{x}_B &= A_B^{-1}b - A_B^{-1}A_N\bar{x}_N \\
x_B &= A_B^{-1}b - A_B^{-1}A_Nx_N \\
f(x_N) &:= A_B^{-1}b - A_B^{-1}A_Nx_N
\end{aligned}$$

In this setting $U_B := \mathbb{R}^p, U_N := \mathbb{R}^{n-p}$.

Lemma 22: Chain Rule

Let $U \subseteq \mathbb{R}^n, V \subseteq \mathbb{R}^m$ be both open sets. $f_1 : U \rightarrow \mathbb{R}^m, f_2 : V \rightarrow \mathbb{R}^p$ be differentiable on U and V respectively such that $f_1(U) \subseteq V$. Then $(f_2 \circ f_1)$ is differentiable on U and

$$D(f_2 \circ f_1)(\bar{x}) = Df_2(f_1(\bar{x})) \circ Df(\bar{x}) \quad \forall \bar{x} \in U$$

Example: Line search, directional derivative

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable on \mathbb{R}^n . Also, given are a current point $\bar{x} \in \mathbb{R}^n$ and a “search direction” $d \in \mathbb{R}^n$. We define $\phi : \mathbb{R} \rightarrow \mathbb{R}$ by $\phi(\alpha) := f(\bar{x} + \alpha d)$. Then $\phi'(\alpha) = \langle \nabla f(\bar{x} + \alpha d), d \rangle$. If f is C^2 , then $\phi''(\alpha) = d^T \nabla^2 f(\bar{x} + \alpha d) d$. Note $\phi'(0) = \langle \nabla f(\bar{x}), d \rangle, \phi''(0) = d^T \nabla^2 f(\bar{x}) d$.

Corollary 23

Suppose h and \bar{x} are as in Theorem 21 (Implicit Function Theorem). Also assume $Z \in \mathbb{R}^{n \times p}$ ($q \leq n - p$) such that $h'(\bar{x})Z = 0$. Then there exists a neighborhood U of $0 \in \mathbb{R}^q$ and a C^1 function $t : U \rightarrow \mathbb{R}^n$ such that

- $t(0) = 0$,
- $t'(0) = 0$,
- $h(\bar{x} + Zd_Z + t(d_Z)) = 0$ for all $d_Z \in U$.

So the function t gives us a way of moving away from \bar{x} (a solution of $h(x) = 0$) in a way that keeps feasible with respect to $h(x) = 0$.

Proof:

Let h, \bar{x} and Z be as in the assumptions. Using the partition $[B|N]$, define $z := \begin{bmatrix} z_B \\ z_N \end{bmatrix}$ (recall $h'(\bar{x}) = [h'_B(\bar{x}) | h'_N(\bar{x})]$). Let $U := \{d_Z \in \mathbb{R}^q : (\bar{x}_N + Z_N d_Z) \in U_N\}$. Define t by

$$\begin{aligned}
t_N(d_Z) &:= 0 \\
t_B(d_Z) &:= f(\bar{x}_N + Z_N d_Z) - \bar{x}_B - Z_B d_Z
\end{aligned}$$

\uparrow
 neighborhood
 of \bar{x}_N from
 Theorem 21

Thus,

$$h(\bar{x} + Zd_Z + t(d_Z)) = h \begin{bmatrix} \bar{x}_B + Z_B d_Z + f(\bar{x}_N + Z_N d_Z) - \bar{x}_B - Z_B d_Z \\ \bar{x}_N + Z_N d_Z + 0 \end{bmatrix} = h \begin{bmatrix} f_N(\bar{x}_N + Z_N d_Z) \\ \bar{x}_N + Z_N d_Z \end{bmatrix} \stackrel{\uparrow}{=} 0$$

By Theorem 21

Also,

$$\begin{aligned}
t(0) &= f(\bar{x}_N) - \bar{x}_B = 0, & t'_N(0) &= 0 \\
t'_B(0) &= f'(\bar{x}_N)Z_N - Z_B = -[h'_B(\bar{x})]^{-1}h'_N(\bar{x})Z_N - Z_B = [h'_B(\bar{x})]^{-1} \underbrace{[-h'_N(\bar{x})Z_N - h'_B(\bar{x})Z_N]}_{=-h'(\bar{x})Z=0} = 0
\end{aligned}$$

\uparrow
 Chain rule (Lemma 22)

□

What does the size of the neighborhood depend on?

Note in LPs $t(d_Z) := 0$ for all $d_Z \in \mathbb{R}^q$.

Corollary 24

Assume h and \bar{x} are as described in Theorem 21. Let $d \in \mathbb{R}^n$ such that $h'(\bar{x})d = 0$. Then there exists $\bar{\lambda} > 0$ and a C^1 arc (directed curve) \hat{t} with properties:

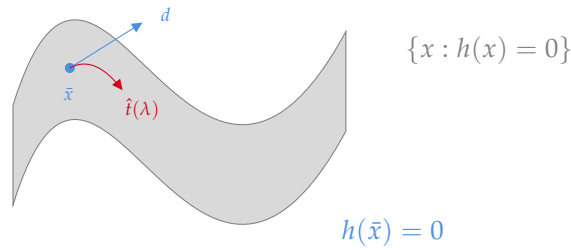
$$\begin{cases} \hat{t}(0) = \bar{x} \\ h(\hat{t}(\lambda)) = 0 \quad \forall \lambda \in [0, \bar{\lambda}) \\ \hat{t}'(0) = d \end{cases}$$

Proof:

In the statement of Corollary 23, plug in $Z := d$ and then using the resulting t ,

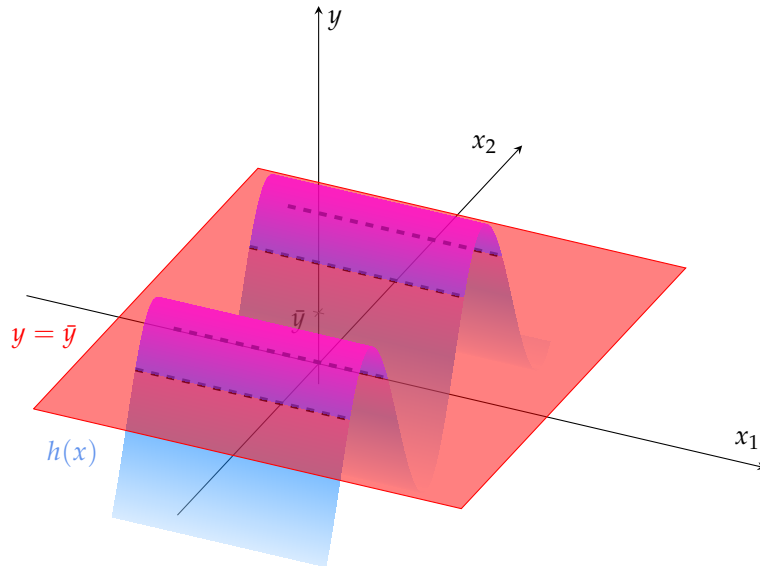
$$\hat{t}(\lambda) := \bar{x} + \lambda \underset{d_Z}{\uparrow} \underset{Z}{\uparrow} d + t(\lambda)$$

□



How applicable are the Theorems 20, 21 and their Corollaries?

Let $h : \mathbb{R}^n \rightarrow \mathbb{R}^p$, where $p \leq n$. Call $\bar{x} \in \mathbb{R}^n$ **regular** if $\text{rank}(h'(\bar{x})) = p$; call $\bar{y} \in \mathbb{R}^p$ a **regular value** if $\forall x \in h^{-1}(\bar{y})$ are regular.



Union of **these curves** is the set $h^{-1}(\bar{y})$.

If h is affine, then $h(x) = Ax - b$ for some given $A \in \mathbb{R}^{p \times n}$, $b \in \mathbb{R}^p$. Let $\bar{y} \in \mathbb{R}^p$ be given. Then $h^{-1}(\bar{y}) = \{x \in \mathbb{R}^n : Ax = \bar{y} + b\}$.

Theorem 25: Sard's Theorem, Morse-Sard Theorem

Let $h : \mathbb{R}^n \rightarrow \mathbb{R}^p$, where $p \leq n$, $h \in \mathcal{C}^r$ with $r \geq n - p + 1$. Then the p -dimensional Lebesgue measure of $\{y \in \mathbb{R}^p : y \text{ is not a regular value}\}$ is zero.

Morse (1939) proved the $p = 1$ case, Sard (1942) proved the generalization above. Smale (1965) proved an infinite dimensional version.

Unconstrained Continuous Optimization

$f : \mathbb{R}^n \rightarrow \mathbb{R}, g : \mathbb{R}^n \rightarrow \mathbb{R}^m, h : \mathbb{R}^n \rightarrow \mathbb{R}^p.$

$$\begin{aligned} \inf \quad & f(x) \\ \text{s.t.} \quad & g(x) \leq 0 \\ & h(x) = 0 \end{aligned} \tag{P}$$

$S := \{x \in \mathbb{R}^n : g(x) \leq 0, h(x) = 0\}$. Here, we assume $S = \mathbb{R}^n$.

Theorem 26: First-order necessary conditions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be \mathcal{C}^1 and $S = \mathbb{R}^n$. Then, $\bar{x} \in \mathbb{R}^n$ is a local minimizer for (P) $\implies f'(\bar{x}) = 0$.

\bar{x} is a **stationary point** of f .

Proof:

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is \mathcal{C}^1 , $S = \mathbb{R}^n$, and $\bar{x} \in \mathbb{R}^n$ is a local minimizer for (P). For the sake of seeking a contradiction, suppose $f'(\bar{x}) \neq 0$. Then, there exists $d \in \mathbb{R}^n$ such that $\langle f'(\bar{x}), d \rangle < 0$ (e.g., let $A \in \mathbb{S}_{++}^n$, and set $d := -Af'(\bar{x})$). Consider $\phi : \mathbb{R} \rightarrow \mathbb{R}$, $\phi(\alpha) := f(\bar{x} + \alpha d)$. Then, $\phi'(0) = \langle f'(\bar{x}), d \rangle < 0$. Thus, for all sufficiently small, positive α , $f(\bar{x} + \alpha d) < f(\bar{x})$. Therefore, \bar{x} is not a local minimizer for (P). \square

Optimality conditions are widely used in algorithm design. E.g., for many software $\|\nabla f(x^{(k)})\| < \epsilon$ is a part of the stopping criteria.

Definition 27

$d \in \mathbb{R}^n$ is a **decent direction** for f at $\bar{x} \in \mathbb{R}^n$, if $\langle f'(\bar{x}), d \rangle < 0$.

$d \in \mathbb{R}^n$ is an **improving direction** for f at \bar{x} , if $f(\bar{x} + \alpha d) < f(\bar{x}) \forall \alpha > 0$ and sufficiently small.

Theorem 28: Second-order necessary conditions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be \mathcal{C}^2 and $S = \mathbb{R}^n$. If $\bar{x} \in \mathbb{R}^n$ a local minimizer for (P), then $f'(\bar{x}) = 0$ and $\nabla^2 f(\bar{x}) \in \mathbb{S}_+^n$.

Proof:

Suppose \bar{x} is a local minimizer for (P). Since f is \mathcal{C}^2 by Theorem 27, $f'(\bar{x}) = 0$. Suppose for the sake of contradiction that $\nabla^2 f(\bar{x}) \notin \mathbb{S}_+^n$. Since $f \in \mathcal{C}^2$, $\nabla^2 f(\bar{x}) \in \mathbb{S}^n$. Therefore, there exists $d \in \mathbb{R}^n$

such that $d^T \nabla^2 f(\bar{x})d < 0$. Define $\phi : \mathbb{R} \rightarrow \mathbb{R}$ by $\phi(\alpha) := f(\bar{x} + \alpha d)$. Then $\phi'(0) = \langle \nabla f(\bar{x}), d \rangle = 0$, $\phi''(0) = d^T \nabla^2 f(\bar{x})d < 0$. Therefore, for all $\epsilon > 0$ and sufficiently small $f(\bar{x} + \epsilon d) < f(\bar{x})$ which contradicts the fact that \bar{x} is a local minimizer for (P). \square

Definition 29: direction of negative curvature

$d \in \mathbb{R}^n$ is called a **direction of negative curvature** for f at \bar{x} if $d^T \nabla^2 f(\bar{x})d < 0$.

Theorem 30: Taylor's Theorem - implicit remainder version

Let $U \subseteq \mathbb{R}^n$ be open, $f : U \rightarrow \mathbb{R}$ be C^r on U . Let $\bar{x}, d \in \mathbb{R}^n$, assume $[\bar{x}, \bar{x} + d] \subset U$. Then,

$$f(\bar{x} + d) = f(\bar{x}) + \sum_{k=1}^r \frac{1}{k!} D^k f(\bar{x})[\underbrace{d, \dots, d}_{k\text{-times}}] + \mathcal{R}(\bar{x}, d),$$

where $\mathcal{R}(\bar{x}, \cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$\lim_{h \rightarrow 0} \frac{\mathcal{R}(\bar{x}, h)}{\|h\|^r} = 0.$$

Theorem 31: Second order sufficient conditions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f \in C^2$, $S = \mathbb{R}^n$. Let $\bar{x} \in \mathbb{R}^n$. If $f'(\bar{x}) = 0$ and $\nabla^2 f(\bar{x}) \in \mathbb{S}_{++}^n$, then \bar{x} is a strict local minimizer for (P).

Proof:

Let $\bar{x} \in \mathbb{R}^n$ such that $f'(\bar{x}) = 0$ and $\nabla^2 f(\bar{x}) \in \mathbb{S}_{++}^n$,

$$\delta := \min_{\substack{\uparrow \\ \lambda_n(\nabla^2 f(\bar{x}))}} \{d^T \nabla^2 f(\bar{x})d : \|d\|_2 = 1\} > 0$$

By Theorem 30, for all $d \in \mathbb{R}^n$, $\|d\|_2 = 1$, and $\alpha > 0$ and small enough, we have

$$f(\bar{x} + \alpha d) = f(\bar{x}) + \underbrace{\alpha \langle \nabla f(\bar{x}), d \rangle}_{=0} + \frac{\alpha^2}{2} d^T \nabla^2 f(\bar{x})d + o(\alpha^2) \geq f(\bar{x}) + \frac{\delta}{2} \alpha^2 + o(\alpha^2)$$

Choose a neighborhood U of \bar{x} such that $\frac{\delta}{2} \alpha^2 > |o(\alpha^2)|$. Then for all $x \in U \setminus \{\bar{x}\}$, $f(x) > f(\bar{x})$. Therefore, \bar{x} is a strict local minimizer for (P). \square

How applicable is this last theorem?

Proposition 32

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be C^2 and consider $\tilde{f}(x) := f(x) + c^T x$, where $c \in \mathbb{R}^n$ is given. Then for almost all $c \in \mathbb{R}^n$, $\tilde{f}(\bar{x}) = 0 \implies \nabla^2 f(\bar{x})$ is nonsingular.

Proof:

Apply Sard's Theorem (Theorem 25) to $g(x) := f'(x)$, with $r := 1$ and $p := n$. \square

What if f has some nice structure, can we say more?

Definition 33: convex function

$f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is **convex** if $\text{epi}(f) := \left\{ \begin{pmatrix} \mu \\ x \end{pmatrix} \in \mathbb{R} \oplus \mathbb{R}^n : f(x) \leq \mu \right\}$ is convex.

Here $\text{epi}(f)$ denotes the epi graph of f .

Theorem 34

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function and $S := \mathbb{R}^n$. Then every local minimizer of (P) is a global minimizer of (P). If in addition, f is differentiable on \mathbb{R}^n , then every stationary point of f is a global minimizer of (P).

2.1 Affine Subspace Constraints

One of the most popular form of continuous optimization problems is

$$\begin{aligned} \inf \quad & f(x) \\ \text{s.t.} \quad & Ax = b \end{aligned}$$

where $A \in \mathbb{R}^{p \times n}$, $b \in \mathbb{R}^p$ are given.

At a first glance (and strictly speaking), (P) does not belong to the class of unconstrained continuous optimization problems. We may assume $\text{rank}(A) = p$; otherwise

- we easily prove $Ax = b$ has no solution, which implies (P) is infeasible, or
- we easily find all redundant equations and $\bar{x} \in \mathbb{R}^n$ such that $A\bar{x} = b$.

So, $\text{rank}(A) = p$. Find a basis B of A and form the partitions

$$[A_B | A_N] := A, \quad \begin{bmatrix} x_B \\ x_N \end{bmatrix} := x.$$

Then,

$$Ax = b \iff x_B = A_B^{-1}b - A_B^{-1}A_Nx_N.$$

Therefore, for every $x \in S$,

$$f(x) = f \left(\begin{array}{c} A_B^{-1}b - A_B^{-1}A_Nx_N \\ x_N \end{array} \right).$$

We define $\tilde{f} : \mathbb{R}^{n-p} \rightarrow \mathbb{R}$ by

$$\tilde{f}(x_N) := f \left(\begin{array}{c} A_B^{-1}b - A_B^{-1}A_Nx_N \\ x_N \end{array} \right).$$

Thus (P) is equivalent to

$$\inf_{x \in \mathbb{R}^{n-p}} \tilde{f}(x) \tag{\tilde{P}}$$

and we can start any algorithm from any starting point $x^{(0)} \in \mathbb{R}^{n-p}$.

Another equivalent approach:

Let $\bar{x} \in S$ (i.e., $A\bar{x} = b$). Then, $S = \{\bar{x} + u : u \in \text{Null}(A)\}$.

Let columns of $Z \in \mathbb{R}^{n \times (n-p)}$ form a basis for $\text{Null}(A)$. Then (P) is also equivalent to

$$\inf_{v \in \mathbb{R}^{n-p}} \hat{f}(v),$$

where $\hat{f} : \mathbb{R}^{n-p} \rightarrow \mathbb{R}$ is defined as $\hat{f}(v) := f(\bar{x} + Zv)$.

In applications, with either of these two approaches, we must be very careful about exploiting sparsity as well as making sure we can efficiently and accurately evaluate all ingredients of the algorithms we choose to use on such problems.

2.2 Applications

Some other ways of dealing with constrained optimization problems using unconstrained optimization algorithms: Form the Lagrangian for (P):

$$\mathcal{L}(x, v) := f(x) + v^T(b - Ax),$$

where $v \in \mathbb{R}^p$ represents the Lagrange multipliers (dual variables corresponding to the constraints).

Use a **penalty function** (penalizing any violation of the constraints):

$$\rho(x, \eta) := f(x) + \eta \|Ax - b\|_{\beta}^{\gamma},$$

where $\beta, \gamma \in \mathbb{R}$ suitably defined, $\eta \in \mathbb{R}_{++}$ a penalty parameter.

In compressed sensing and related applications, one seeks a solution of

$$\inf \{f(x) + \eta \|x\|_0 : Ax = b\},$$

where $\|x\|_0 :=$ number of nonzero entries of x . As an approximation, many researches and practitioners work with

$$\inf \{f(x) + \eta_1 \|x\|_1 + \eta_2 \|Ax - b\|_2^{\gamma}\},$$

where $\eta_1, \eta_2, \gamma \in \mathbb{R}$, usually fixed.

We can generalize such approaches to matrix variables. Very many interesting applications in Machine Learning, AI and modern Data Science. In many of these applications, we want to find a low-rank solution.

Example:

$\min\{\text{rank}(X) : A(X) = b\}$, where $A : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$ linear, $b \in \mathbb{R}^p$, both A, b are given.

2.3 Prototype low-rank approximation problem

Given $A \in \mathbb{R}_+^{m \times n}$ (both m and n are huge). We want to find matrices $U \in \mathbb{R}_+^{m \times k}$, $V \in \mathbb{R}_+^{n \times k}$ such that $A = UV^T$ and k is as small as possible.

If we do not require U and V to be nonnegative, the problem is solved by Singular Value Decomposition (SVD) and optimal k is the rank of A .

$A = Q_1 D Q_2^T$ where $Q_1 \in \mathbb{R}^{m \times m}$, $Q_2 \in \mathbb{R}^{n \times n}$ are orthogonal and $D \in \mathbb{R}^{m \times n}$ diagonal. Let's assume $m \leq n$, then

$$D = \begin{bmatrix} \sigma_1(A) & & & 0 & 0 & \dots & 0 \\ & \sigma_2(A) & & & 0 & \dots & 0 \\ & & \ddots & & & & \vdots \\ 0 & & & \sigma_m(A) & 0 & \dots & 0 \end{bmatrix}$$

where $\sigma_1(A) \geq \sigma_2(A) \geq \dots \geq \sigma_m(A) \geq 0$ are the singular values of A .

Theorem 35

Every $A \in \mathbb{R}^{m \times n}$ has a SVD.

Requiring U, V to be nonnegative, makes the problem hard. Let p be an upper bound on k (taking p as $(mn + 1)$ suffices, but in practice, better guesses can help). Suppose our guess for the minimum

nonnegative rank of A is p . Then let $U \in \mathbb{R}^{m \times p}$ and $V \in \mathbb{R}^{n \times p}$ denote the variable matrices and consider

$$\inf f(U, V) := \eta_1 \|A - UV^T\| + \eta_2 \|U_-\| + \eta_3 \|V_-\|,$$

where $\eta_1, \eta_2, \eta_3 \in \mathbb{R}^+$ are parameters that we can fix, and U_- denotes the $\mathbb{R}^{m \times p}$ matrix with only negative entries of U .

2.4 Classical Algorithmic Approaches

I. Search direction + line-search strategies

Pick a search direction $d^{(k)}$, pick a step-size $\alpha_k > 0$. $x^{(k+1)} := x^{(k)} + \alpha_k d^{(k)}$. Repeat.

- $d^{(k)} := -\nabla f(x^{(k)})$ steepest-descent direction
- any $d^{(k)}$ with $\langle \nabla f(x^{(k)}), d^{(k)} \rangle < 0$
- $d^{(k)} := -[\nabla^2 f(x^{(k)})]^{-1} \nabla f(x^{(k)})$ Newton direction

↑
Assuming $\nabla^2 f(x^{(k)}) \in \mathbb{S}_{++}^n$

For convex optimization problems and also near local minimizers of nonconvex problems, we want $\alpha_k \approx 1$ with this direction \rightarrow superlinear or quadratic convergence

Exact Line-Search Find $\alpha > 0$ such that $\phi(\alpha) := f(x^{(k)} + \alpha d^{(k)})$ is minimized. Typically *not* practical.

Inexact Line-Search Armijo-Goldstein (1966-67) conditions (or Wolfe (1969) conditions):

Choose $\alpha > 0$ so that

$$f(x^{(k)} + \alpha d^{(k)}) \leq f(x^{(k)}) + c_1 \cdot \alpha \langle \nabla f(x^{(k)}), d^{(k)} \rangle$$

(sufficiently good rate for the decrease in the objective function) and “**curvature condition**”

$$\langle \nabla f(x^{(k)} + \alpha d^{(k)}), d^{(k)} \rangle \geq c_2 \langle \nabla f(x^{(k)}), d^{(k)} \rangle$$

(step size should not be too small) where constants c_1, c_2 satisfy $0 < c_1 < c_2 < 1$.

Strong Wolfe Conditions

$$f(x^{(k)} + \alpha d^{(k)}) \leq f(x^{(k)}) + c_1 \cdot \alpha \langle \nabla f(x^{(k)}), d^{(k)} \rangle$$

and

$$|\langle \nabla f(x^{(k)} + \alpha d^{(k)}), d^{(k)} \rangle| \leq c_2 |\langle \nabla f(x^{(k)}), d^{(k)} \rangle|$$

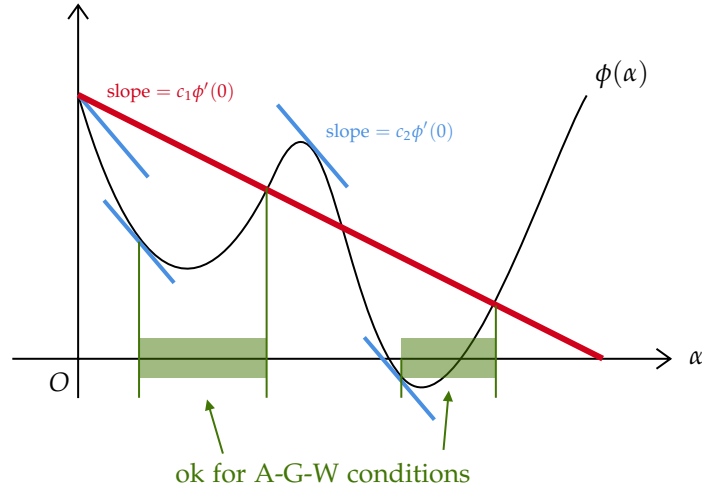
↑
The second conditions disallows
this being too large and positive

Lemma 36

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be \mathcal{C}^1 , and $d \in \mathbb{R}^n$ be a descent direction at $\bar{x} \in \mathbb{R}^n$ for f . Suppose f is bounded from below on the ray $\{\bar{x} + \alpha d : \alpha \in \mathbb{R}_+\}$. Then $\forall 0 < c_1 < c_2 < 1$, there exists step lengths $\alpha > 0$ satisfying Armijo-Goldstein-Wolfe as well as Strong Wolfe conditions.

With $\phi(\alpha) := f(\bar{x} + \alpha d)$, $0 < c_1 < c_2 < 1$, choose $\alpha > 0$ such that

$$\text{Armijo-Goldstein-Wolfe} \begin{cases} \phi(\alpha) \leq \phi(0) + c_1 \alpha \phi'(0) \\ \phi'(\alpha) \geq c_2 \phi'(0) \end{cases} \quad \text{Strong Wolfe} \begin{cases} \phi(\alpha) \leq \phi(0) + c_1 \alpha \phi'(0) \\ |\phi'(\alpha)| \leq c_2 |\phi'(0)| \end{cases}$$



Proof of Lemma 36:

Suppose the stated assumptions hold. We adopt the above mentioned notion with ϕ . Then $\phi(\alpha)$ is bounded from below on $\{\alpha \in \mathbb{R} : \alpha \geq 0\}$. Since $c_1 \in (0, 1)$ and $\phi'(0) = \langle \nabla f(\bar{x}), d \rangle < 0$ (d is a descent direction for f), the ray $\{\phi(0) + (c_1 \cdot \phi'(0))\alpha : \alpha \geq 0\}$ is unbounded below and therefore, intersects the graph of ϕ at least once for $\alpha > 0$. Let $\bar{\alpha} > 0$ denote the smallest value of α for which the ray intersects the graph of ϕ . Then,

$$\phi(\bar{\alpha}) = \phi(0) + \bar{\alpha}c_1\phi'(0) \quad (*)$$

Thus, the first condition of A-G-W holds on $(0, \bar{\alpha}]$.

By the Mean Value Theorem, there exists $\hat{\alpha} \in (0, \bar{\alpha})$ such that $\phi(\bar{\alpha}) - \phi(0) = \bar{\alpha}\phi'(\hat{\alpha})$. Therefore,

$$\phi(\bar{\alpha}) - \phi(0) \stackrel{(*)}{=} \bar{\alpha}\phi'(\hat{\alpha}) \stackrel{(*)}{=} \bar{\alpha}c_1\phi'(0) > c_2\bar{\alpha}\phi'(0).$$

\uparrow
 $c_2 > c_1, \phi'(0) < 0$

Thus, A-G-W conditions strictly hold at $\hat{\alpha}$. Since $\phi'(\hat{\alpha}) < 0$, Strong Wolfe conditions also hold at $\hat{\alpha}$ as well as in a sufficiently small neighborhood of $\hat{\alpha}$. \square

In the textbook, Backtracking line search.

II. Trust-Region Strategies

Use the information gathered about f so far and construct an approximation ("model") m_k of the function f . Then solve

$$\begin{aligned} \min \quad & m_k(d) \\ \text{s.t.} \quad & d \in \text{Trust Region (around } x^{(k)}) \end{aligned}$$

$x^{(k)} \in \mathbb{R}^n$ is our current iterate. Let B_k denote $\nabla^2 f(x^{(k)})$ or an approximation of it. Choose $\delta_k > 0$, and solve

$$\begin{aligned} \min \quad & m_k(d) := f(x^{(k)}) + \langle \nabla f(x^{(k)}), d \rangle + \frac{1}{2}d^T B_k d \\ \text{s.t.} \quad & \|d\|_2 \leq \delta_k \end{aligned}$$

Let \bar{d} denote an optimal solution of this trust-region subproblem. If $x^{(k)} + \bar{d}$ satisfies certain criteria, then set $x^{(k+1)} := x^{(k)} + \bar{d}$; otherwise either modify δ_k , or the step size, ...

Depending on how well we did with the latest δ_k choose a suitable value for δ_{k+1} and repeat. (size of the Trust-Region is being adjusted.)

2.5 Convergence Properties of Descent Algorithms

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$. For every $\beta \in \mathbb{R}$,

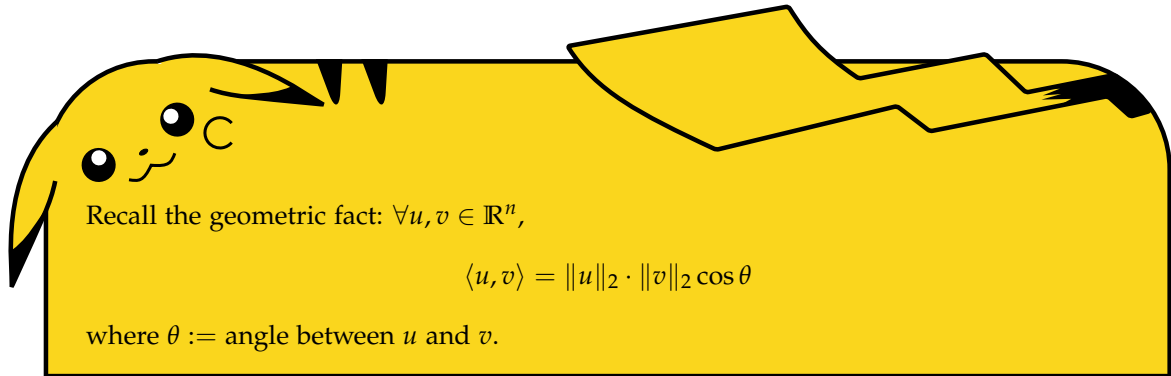
$$\{x \in \mathbb{R}^n : f(x) \leq \beta\}$$

is called a **sublevel set** of f (some literature use level set).

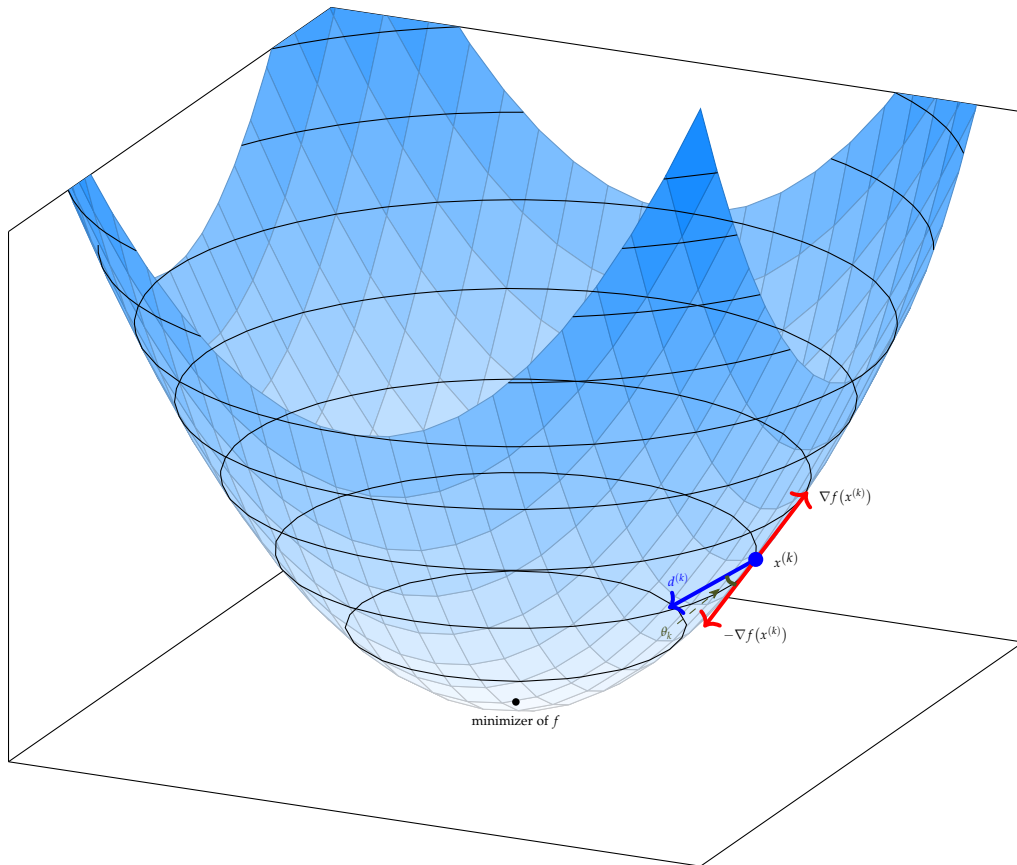
$$\{x \in \mathbb{R}^n : f(x) = \beta\}$$

is called a **level set** of f (some also call it a contour of f)

Consider a descent algorithm: Start with $x^{(0)} \in \mathbb{R}^n$, at each iteration k , choose $d^{(k)} \in \mathbb{R}^n$ such that $\langle \nabla f(x^{(k)}), d^{(k)} \rangle < 0$ and choose $\alpha_k > 0$, $x^{(k+1)} := x^{(k)} + \alpha_k d^{(k)}$.



Define $\theta_k := \arccos\left(-\frac{\langle \nabla f(x^{(k)}), d^{(k)} \rangle}{\|\nabla f(x^{(k)})\|_2 \cdot \|d^{(k)}\|_2}\right)$



Theorem 37: Zoutendijk (1970), Wolfe (1969)

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a bounded from below, $x^{(0)} \in \mathbb{R}^n$, and f be \mathcal{C}^1 on

$$N := \text{nbhd}\{x \in \mathbb{R}^n : f(x) \leq f(x^{(0)})\}.$$

Assume ∇f is Lipschitz continuous on N with Lipschitz constant $L \in \mathbb{R}_{++}$. Then, every descent algorithm following Armijo-Goldstein-Wolfe conditions for stepsize selection satisfies:

$$\sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(x^{(k)})\|_2^2 < \infty.$$

Proof:

Suppose the assumptions in the statement hold. For every iteration k , due to the second A-G-W condition, we have $\langle \nabla f(x^{(k+1)}), d^{(k)} \rangle \geq c_2 \langle \nabla f(x^{(k)}), d^{(k)} \rangle$, which implies

$$\langle \nabla f(x^{(k+1)}) - \nabla f(x^{(k)}), d^{(k)} \rangle \geq (c_2 - 1) \langle \nabla f(x^{(k)}), d^{(k)} \rangle \quad (\star)$$

Due to the fact that we are working with a descent algorithm ($\langle \nabla f(x^{(k)}), d^{(k)} \rangle < 0$ for all k) and the first condition of A-G-W, $\{x^{(k)}\} \subset N$. Since ∇f is Lipschitz continuous on N with Lipschitz constant L ,

$$\langle \nabla f(x^{(k+1)}), d^{(k)} \rangle \leq \underbrace{\|\nabla f(x^{(k+1)}) - \nabla f(x^{(k)})\|_2}_{\text{Cauchy-Schwarz}} \|d^{(k)}\|_2 \leq \underbrace{\alpha_k L}_{\nabla f \in \text{Lip}(L)} \|d^{(k)}\|_2 \quad (\spadesuit)$$

By (\star) and (\spadesuit) , we have

$$\alpha_k \geq \frac{(c_2 - 1) \langle \nabla f(x^{(k)}), d^{(k)} \rangle}{L \|d^{(k)}\|_2^2}$$

Substituting this lower bound on α_k into the first A-G-W condition, we obtain

$$\begin{aligned} f(x^{(k)} + \alpha_k d^{(k)}) &\leq f(x^{(k)}) - \frac{c_1(c_2 - 1) \langle \nabla f(x^{(k)}), d^{(k)} \rangle}{L \|d^{(k)}\|_2^2} \\ \iff f(x^{(k+1)}) &\leq f(x^{(k)}) - \left(\frac{c_1(1 - c_2)}{L} \right) \cos^2 \theta_k \|\nabla f(x^{(k)})\|_2^2 \end{aligned}$$

Applying the above to pairs of consecutive iterates, we obtain:

$$f(x^{(k+1)}) \leq f(x^{(0)}) - \frac{c_1(1 - c_2)}{L} \sum_{\ell=0}^k \cos^2 \theta_\ell \|\nabla f(x^{(\ell)})\|_2^2$$

Since f is bounded from below, $[f(x^{(0)}) - f(x^{(k)})]$ is bounded from above, and

$$\frac{c_1(1 - c_2)}{L} \sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(x^{(k)})\|_2^2 < +\infty$$

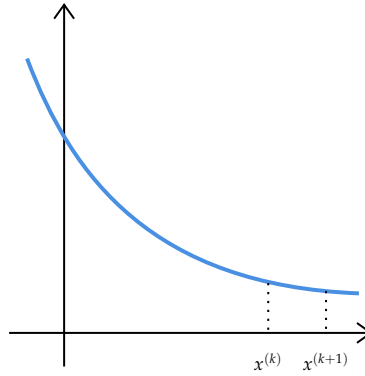
□

A consequence of Theorem 37:

$$\sum_{k=0}^{\infty} \cos^2 \theta_k \|\nabla f(x^{(k)})\|_2^2 < +\infty \implies \cos^2 \theta_k \|\nabla f(x^{(k)})\|_2^2 \rightarrow 0 \text{ as } k \rightarrow +\infty$$

Therefore, if $\cos^2 \theta_k \geq \delta > 0$ for all $k \in \mathbb{Z}_+$, then $\lim_{k \rightarrow +\infty} \|\nabla f(x^{(k)})\|_2 \rightarrow 0$. In some places, including the textbook, this criterion is used to conclude that Steepest-Descent Algorithm is “globally convergent”.

What if



i.e., what if $\|\nabla f(x^{(k)})\| < 10^{-8}$ but (even assuming convexity of f etc.) the unique minimizer \bar{x} is far away from $x^{(k)}$?

2.6 A General Conversation about Convergence

Example:

$f : \mathbb{R} \rightarrow \mathbb{R}, f(x) := \frac{1}{4}x^4 - 5x$. Then f is convex, global minimizer is unique and attained at $\bar{x} := \sqrt[3]{5}$ which is irrational (even though the data $\subseteq \mathbb{Z}^{\cdots}$)

Thus we cannot expect finite algorithms in the worst-case. We will generate a sequence $x^{(1)}, x^{(2)}, \dots$. We will hope for conclusions like:

- For every $x^{(0)}, x^{(k)} \rightarrow \bar{x}$ a global minimizer.
- For every $x^{(0)}, x^{(k)} \rightarrow \bar{x}$ a local minimizer.
- For every $x^{(0)}$, all limit points of $\{x^{(k)}\}$ are global (local) minimizers, or $f(x^{(k)}) \rightarrow -\infty$.
- For every $x^{(0)}$, all limit points of $\{x^{(k)}\}$ satisfy second-order necessary conditions.
- For every $x^{(0)}$, all limit points of $\{x^{(k)}\}$ satisfy first-order necessary conditions.
- For every $x^{(0)}, \lim_{k \rightarrow +\infty} \|\nabla f(x^{(k)})\| = 0$.

Locally, replace “every $x^{(0)} \in \mathbb{R}$ ” by “every $x^{(0)} \in B(\bar{x}, \eta) := \{x \in \mathbb{R}^n : \|x - \bar{x}\|_2 < \eta\}$,” and hope that the 2nd-order sufficient conditions hold.

How fast? $\epsilon_k := \|x^{(k)} - \bar{x}\|$.

Example:

$\epsilon_k := (0.1)^k \rightarrow 10^{-1}, 10^{-2}, \dots$	linear convergence
$\epsilon_k := (0.9)^k \rightarrow 0.9, 0.081, 0.729, \dots$	linear convergence
$\epsilon_k := (0.1)^{2^k} \rightarrow (10)^{-2}, (10)^{-4}, (10)^{-8}, \dots$	quadratic convergence
$\epsilon_k := (0.9)^{2^k} \rightarrow 0.81, 0.65, 0.43, 0.185, 0.034, 10^{-3}, \dots$	quadratic convergence
$\epsilon_k := (0.1)^{3^k} \rightarrow \dots$	cubic convergence

Definition 38

If $\epsilon_k \searrow 0$ and $\epsilon_{k+1} \leq \beta(\epsilon_k)^p$ for some constants $p \geq 1$ and β ($\beta < 1$ if $p = 1$) and for all sufficiently large k , then we say $\epsilon_k \rightarrow 0$ with Q-order (at least) p . If $\epsilon_k \searrow 0$ and $\frac{\epsilon_{k+1}}{\epsilon_k} \rightarrow 0$ (as $k \rightarrow +\infty$) then the convergence is **Q-superlinear**.

Q-linear := Q-order 1, Q-quadratic := Q-order 2

Example:

$\epsilon_k := (\frac{1}{2})^k, k \in \mathbb{Z}_{++}$. Then $\epsilon_k \searrow 0$, Q-superlinearly, but it does not have Q-order $p > 1$.

Given a sequence $\{\epsilon_k\} \subset \mathbb{R}_+$, let $\eta_i := \sup\{\epsilon_k : k \geq i\}$. Then $\limsup_{k \rightarrow +\infty} \{\epsilon_k\} := \lim_{i \rightarrow \infty} \eta_i$.

Definition 39

If $\epsilon_k \searrow 0$ and $\limsup_{k \rightarrow +\infty} \{\epsilon_k^{1/q^k}\} < 1$ for all $0 < q < p, p > 1$, then $\epsilon_k \rightarrow 0$ with R-order (at least) p .

This is the same as $\limsup_{k \rightarrow +\infty} \{\frac{1}{q^k} \ln(\epsilon_k)\} < 0$

Proposition 40

- (i) If $x^{(k)} \rightarrow \bar{x}$ with Q-order p (R-order p), so does $\{x^{(k+\ell)}\}$ for every fixed $\ell \in \mathbb{Z}_+$.
- (ii) If $\epsilon_k \searrow 0$ with Q-order p and $0 < \eta_k \leq \epsilon_k$ for $\forall k \in \mathbb{Z}_{++}$, then $\eta_k \searrow 0$ with R-order p .

Fast Local Convergence of Newton's Method

This goes back at least to Kantorovich (Nobel Prize in Economics for his work on “the theory of optimal allocation of resources”, 1975). In addition to his foundational work on the convergence theory of Newton's Method, Kantorovich also made significant contributions to functional analysis and operator theory.

Lemma 41

Let $A, B \in \mathbb{R}^{n \times n}$, A nonsingular, $\|A^{-1}\|_2 \leq \gamma$ and $\|A - B\|_2 \leq \frac{1}{3\gamma}$. Then B is nonsingular and $\|B^{-1}\|_2 \leq \frac{3\gamma}{2}$.

$$\|A^{-1}\|_2 \leq \gamma \iff \text{dist}(A, \text{singular matrices}) \geq \frac{1}{\gamma}$$

Proof:

Suppose $A, B \in \mathbb{R}^n$ satisfy the assumptions. Note that

$$B = A - (A - B) = A[I - A^{-1}(A - B)],$$

and

$$\|A^{-1}(A - B)\|_2 \leq \|A^{-1}\|_2 \|A - B\|_2 \leq \gamma \cdot \frac{1}{3\gamma} = \frac{1}{3}.$$

If $C \in \mathbb{R}^{n \times n}$ nonsingular such that $\|C\|_2 \leq \frac{1}{3}$, then $(I - C)$ is invertible and

$$(I - C)^{-1} = I + C + C^2 + \dots = \sum_{k=0}^{\infty} C^k.$$

This implies

$$\|(I - C)^{-1}\|_2 \leq \sum_{k=0}^{\infty} \left(\frac{1}{3}\right)^k = \frac{1}{2/3} = \frac{3}{2}.$$

Thus, with $C := A^{-1}(A - B)$, (then $B = A(I - C)$) B is invertible. $B^{-1} = (I - C)^{-1}A^{-1}$ and

$$\|B^{-1}\|_2 \leq \|(I - C)^{-1}\|_2 \cdot \|A^{-1}\|_2 \leq \frac{3}{2}\gamma.$$

□

Lemma 42

Let $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $g \in \mathcal{C}^1$ and $\nabla g \in \text{Lip}(L)$ on some open and convex set $D \subseteq \mathbb{R}^n$. Then

$$\|g(y) - g(x) - \nabla g(x)(y - x)\|_2 \leq \frac{L}{2} \|y - x\|_2^2,$$

$\forall x, y \in D$.

Proof:

We already proved this as a part of our discussion following Theorem 19. □

Newton's Method: $x^{(0)} \in \mathbb{R}^n, f : \mathbb{R}^n \rightarrow \mathbb{R}, f \in \mathcal{C}^2$.

$$\forall k \in \mathbb{Z}_+ : \begin{cases} d^{(k)} := -[\nabla^2 f(x^{(k)})]^{-1} \nabla f(x^{(k)}) \\ x^{(k+1)} := x^{(k)} + d^{(k)} \end{cases}$$

Theorem 43

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f \in \mathcal{C}^2$, $x^{(0)} \in \mathbb{R}^n, L \geq 1$. Assume $\nabla f(\bar{x}) = 0$, $\nabla^2 f(\bar{x})$ is nonsingular, $\nabla^2 f \in \text{Lip}(L)$ in an open neighborhood of \bar{x} . Then there exists an open neighborhood N_1 of \bar{x} such that for all $x^{(0)} \in N_1$, Newton's Method converges to \bar{x} linearly and the method is **locally Q-quadratically convergence** (there exists an open neighborhood $N_2 \subseteq N_1$ of \bar{x} such that

$$\forall x^{(0)} \in N_2, \|x^{(k+1)} - \bar{x}\|_2 \leq \text{constant} \cdot \|x^{(k)} - \bar{x}\|_2^2,$$

$\forall k \in \mathbb{Z}_+$). Moreover, $\|\nabla f(x^{(k)})\|$ also converges to zero in N_1 , locally Q-quadratically

$$(\forall x^{(0)} \in N_2, \|\nabla f(x^{(k+1)})\|_2 \leq \text{constant} \cdot \|\nabla f(x^{(k)})\|_2^2, \quad \forall k \in \mathbb{Z}_+).$$

Proof:

Suppose the assumptions hold.

Let $\gamma := \|\nabla^2 f(\bar{x})\|_2^{-1}$. Choose $\eta > 0$ such that with $\mathcal{B} := B(\bar{x}, \eta) := \{x \in \mathbb{R}^n : \|x - \bar{x}\|_2 < \eta\}$, $\nabla^2 f \in \text{Lip}(L)$ on \mathcal{B} and $\eta \leq \frac{1}{3\gamma L}$. Then for all $x \in \mathcal{B}$,

$$\|\nabla^2 f(x) - \nabla^2 f(\bar{x})\|_2 \leq L\|x - \bar{x}\|_2 < L\eta \leq \frac{1}{3\gamma} \quad (\diamond)$$

Therefore, by lemma 41 (with $A := \nabla^2 f(\bar{x})$, $B := \nabla^2 f(x), x \in \mathcal{B}$), $\nabla^2 f(x)$ is nonsingular $\forall x \in \mathcal{B}$; thus, Newton's Method is well-defined for $\{x^{(k)}\} \subset \mathcal{B}$.

We will prove the theorem by induction on the iteration number k .

Let $x^{(0)} \in \mathcal{B}$ (in general, $x^{(k)} \in \mathcal{B}$), then

$$\begin{aligned}
\|x^{(k+1)} - \bar{x}\|_2 &= \|x^{(k)} - [\nabla^2 f(x^{(k)})]^{-1} \nabla f(x^{(k)}) - \bar{x}\|_2 \\
&= \|[\nabla^2 f(x^{(k)})]^{-1} - [\underset{\substack{\uparrow \\ \nabla f(\bar{x})}}{0} - \nabla f(x^{(k)}) - \nabla^2 f(x^{(k)})(\bar{x} - x^{(k)})]\|_2 \\
&\leq \|[\nabla^2 f(x^{(k)})]^{-1}\|_2 \|\nabla f(\bar{x}) - \nabla f(x^{(k)}) - \nabla^2 f(x^{(k)})(\bar{x} - x^{(k)})\|_2 \\
&\leq \frac{3\gamma}{2} \cdot \frac{L}{2} \|x^{(k)} - \bar{x}\|_2^2 \quad \text{by } (\diamond) \text{ and Lemma 41, Lemma 42} \\
&= \frac{3\gamma L}{4} \|x^{(k)} - \bar{x}\|_2^2 \quad \text{by induction on } k, \text{ we get quadratic convergence}
\end{aligned}$$

Also, if $x^{(k)} \in \mathcal{B}$, then

$$\|x^{(k+1)} - \bar{x}\|_2 \leq \frac{3\gamma L}{4} \cdot \frac{1}{3\gamma L} \|x^{(k)} - \bar{x}\|_2 = \frac{1}{4} \|x^{(k)} - \bar{x}\|_2$$

which is linear convergence.

Next, let $d := x^{(k+1)} - x^{(k)}$. Then

$$\begin{aligned}
\|\nabla f(x^{(k+1)})\|_2 &= \|\nabla f(x^{(k+1)}) - \nabla f(x^{(k)}) - \nabla^2 f(x^{(k)})d\|_2 \\
&\leq \frac{L}{2} \|d\|_2^2 \quad \text{Lemma 42} \\
&= \frac{L}{2} \|[\nabla^2 f(x^{(k)})]^{-1} \nabla f(x^{(k)})\|_2^2 \quad (\clubsuit) \\
&\leq \frac{L}{2} \|[\nabla^2 f(x^{(k)})]^{-1}\|_2^2 \|\nabla f(x^{(k)})\|_2^2 \\
&\leq \frac{9\gamma^2 L}{8} \|\nabla f(x^{(k)})\|_2^2 \quad \text{by } x^{(k)} \in \mathcal{B}, (\clubsuit) \text{ and Lemma 41}
\end{aligned}$$

Now, we have all the ingredients for an induction proof. We proved

$$\begin{aligned}
&\text{"}\forall x^{(0)} \in \mathcal{B}, \|x^{(1)} - \bar{x}\|_2 \leq \frac{1}{4} \|x^{(0)} - \bar{x}\|_2; \text{ so, } x^{(1)} \in \mathcal{B}, \text{ and } \|x^{(1)} - \bar{x}\|_2 \leq \frac{3\gamma L}{4} \|x^{(0)} - \bar{x}\|_2^2, \text{ and} \\
&\|\nabla f(x^{(1)})\|_2 \leq \frac{9\gamma^2 L}{8} \|\nabla f(x^{(0)})\|_2^2\text{"}
\end{aligned}$$

By induction on k , we establish the desired inequalities on $x^{(k)}$. For the gradient, from (\clubsuit) ,

$$\begin{aligned}
\|\nabla f(x^{(k+1)})\|_2 &\leq \frac{L}{2} \|[\nabla^2 f(x^{(k)})]^{-1} \nabla f(x^{(k)})\|_2^2 \\
&= \frac{L}{2} \|x^{(k+1)} - x^{(k)}\|_2 \cdot \|[\nabla^2 f(x^{(k)})]^{-1} \nabla f(x^{(k)})\|_2 \\
&\leq \frac{L}{2} \cdot \frac{2}{3\gamma L} \cdot \frac{3\gamma}{2} \|\nabla f(x^{(k)})\|_2 \quad x^{(k)}, x^{(k+1)} \in \mathcal{B}, \text{ Lemma 41} \\
&= \frac{1}{2} \|\nabla f(x^{(k)})\|_2
\end{aligned}$$

Therefore, for every $x^{(0)} \in \mathcal{B}$,

$$\begin{aligned}
\|x^{(k)} - \bar{x}\| &\rightarrow 0, \quad \text{Q-linearly, and locally Q-quadratically} \\
\|\nabla f(x^{(k)})\| &\rightarrow 0, \quad \text{Q-linearly, and locally Q-quadratically}
\end{aligned}$$

□

This proof also applies to the problem of solving systems of nonlinear equations.

$g : \mathbb{R}^n \rightarrow \mathbb{R}^n, g \in \mathcal{C}^1$ on open and convex set $D \subseteq \mathbb{R}^n$. $\exists \bar{x} \in D$ such that $g(\bar{x}) = 0, \nabla g(\bar{x})$ is nonsingular, $\nabla g \in \text{Lip}(L)$ on D . Let $x^{(0)} \in D$,

$$x^{(k+1)} := x^{(k)} - [\nabla g(x^{(k)})]^{-1} g(x^{(k)}) \quad \forall k \in \mathbb{Z}_+.$$

Potential Problems with Newton's Method

- (I) Fails if $\nabla^2 f(x^{(k)})$ is singular (or very ill-conditioned)
- (II) $x^{(k+1)}$ is not the minimizer of the local quadratic model \tilde{f} for f if $\nabla^2 f(x^{(k)})$ is not positive definite.
- (III) Not globally convergent in general.
- (IV) May not even provide descent in general.

Possible Remedies:

- To address (I) & (II) modify $\nabla^2 f(x^{(k)})$, if necessary, to a “nearby” symmetric positive definite matrix B_k . (do this in an efficient & numerically stable way.)
- Together with the above remedy, use A-G-W or Strong Wolfe based line searches to address (III) & (IV).

Still there are some disadvantages:

- (i) We must evaluate Hessians at every iteration,
- (ii) We must solve $n \times n$ linear systems of equations in every iteration.

For some problems evaluating the Hessian is very little extra work compared to $f, \nabla f$. Also, in some cases Automatic Differentiation via a small number of $\nabla f(\cdot)$ evaluations suffice. (Chapter 8 of the textbook)

2.7 Quasi-Newton Methods

Consider $B_k \in \mathbb{S}_{++}^n$, then $-B_k^{-1} \nabla f(x^{(k)})$ is a descent direction for f at $x^{(k)}$. Consider a quadratic model of f (near $x^{(k)}$):

$$\tilde{f}(d) := f(x^{(k)}) + \langle \nabla f(x^{(k)}), d \rangle + \frac{1}{2} d^T B_k d.$$

Since $B_k \in \mathbb{S}_{++}^n$, \tilde{f} has a unique global minimizer at $\vec{d} = -B_k^{-1} \nabla f(x^{(k)})$. Now we can do a line search and find $x^{(k+1)}$, then we have $f(x^{(k+1)})$ and $\nabla f(x^{(k+1)})$. How do we find B_{k+1} ?

Wish List for B_{k+1}

$B_{k+1} \in \mathbb{S}_{++}^n$. B_{k+1} should incorporate newly discovered information about $\nabla^2 f$.

$$\begin{aligned} s^{(k)} &:= x^{(k+1)} - x^{(k)} && \text{(primal step at iteration } k) \\ y^{(k)} &:= \nabla f(x^{(k+1)}) - \nabla f(x^{(k)}) && \text{(dual step at iteration } k) \end{aligned}$$

Magical solution: BFGS

$$B_{k+1} := \underbrace{B - \frac{1}{s^T B s} B s s^T B + \frac{1}{y^T s} y y^T}_{\text{we dropped the iteration number } k}$$

Note:

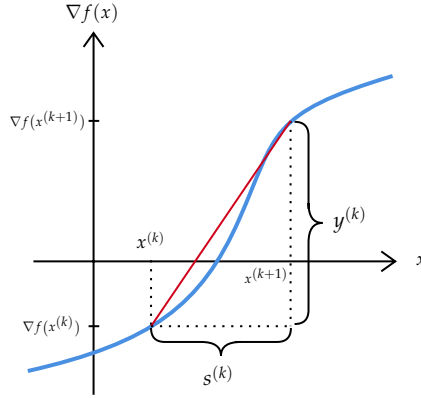
By Theorem 19, we have

$$y^{(k)} = \left[\int_0^1 \nabla^2 f(x^{(k)} - \alpha s^{(k)}) \partial \alpha \right] s^{(k)}$$

i.e., $y^{(k)}$ tells us the behavior of the “average” Hessian (along the line segment $[x^{(k)}, x^{(k+1)}]$) on the subspace $\text{span}\{s^{(k)}\}$. So, we want $B_{k+1} \in \mathbb{R}^{n \times n}$ such that $y^{(k)} = B_{k+1} s^{(k)}$ which is **secant equation**. By enforcing this equation on B_{k+1} , we can incorporate new “secant” information about $\nabla^2 f$.

Example:

Consider $n := 1$.



If $B_{k+1} \succ 0$ satisfies the secant equation, then $\langle y^{(k)}, s^{(k)} \rangle = \langle B_{k+1} s^{(k)}, s^{(k)} \rangle > 0$ since $s^{(k)} \neq 0, B_{k+1} \succ 0$. Notice that $\langle y^{(k)}, s^{(k)} \rangle$ is positively proportional to:

$$\langle y^{(k)}, d^{(k)} \rangle = \langle \nabla f(x^{(k+1)}), d^{(k)} \rangle - \langle \nabla f(x^{(k)}), d^{(k)} \rangle = \phi'(\alpha_k) - \phi'(0) > 0$$

The condition $\langle y^{(k)}, s^{(k)} \rangle > 0$ is called the **curvature condition**.

↑
If we use A-G-W
or Strong Wolfe
based line-search

How do we ensure B_{k+1} is close to B_k ?

Solve the optimization problem

$$\begin{aligned} \min \quad & \|B - H\|_F \\ \text{s.t.} \quad & Bs = y, \quad B \in \mathbb{R}^{n \times n} \end{aligned} \quad (\text{P}_1)$$

for a fixed $H \in \mathbb{R}^{n \times n}$, e.g., $H := B_k$, and fixed $y, s \in \mathbb{R}^n$.

Here,

$$\|A\|_F := \left(\sum_{i=1}^n \sum_{j=1}^n A_{ij}^2 \right)^{1/2} = [\text{Tr}(A^T A)]^{1/2} = [\text{vec}(A)^T \text{vec}(A)]^{1/2}$$

↑
Frobenius norm

(P₁) always has a unique solution \bar{B} . $Z = \bar{B} - H$. Note $\bar{B}s = y \iff \bar{B}s - Hs = y - Hs =: r$.

With this change of variable and definitions, (P₁) is equivalent to

$$\begin{aligned} \min \quad & \|Z\|_F \\ \text{s.t.} \quad & Zs = r \end{aligned} \quad (\text{P}_2)$$

Suppose $s \neq 0$ (i.e., we moved!). Let $Q \in \mathbb{R}^{n \times n}$ be orthogonal such that $Qs = \beta e_1, \beta \neq 0$. $\tilde{Z} := ZQ^T$. Then (P₂) is equivalent to

$$\begin{aligned} \min \quad & \|\tilde{Z}\|_F \\ \text{s.t.} \quad & \tilde{Z}e_1 = \frac{1}{\beta}r \end{aligned} \quad (\text{P}_3)$$

which implies

$$\tilde{Z} = \begin{bmatrix} \frac{1}{\beta}r & 0 & 0 & \cdots & 0 \end{bmatrix}$$

Using our definitions, we compute

$$Z = \tilde{Z}Q = \frac{1}{\beta}re_1^T Q = \frac{1}{\beta^2}rs^T Q^T Q = \frac{1}{\beta^2}rs^T$$

$$Zs = r \implies \frac{1}{\beta^2}r(s^T s) = r \implies s^T = \beta^2$$

↑
unless $r = 0$ in
which $Z = 0$

Therefore, the unique optimal solution of (P₁) is

$$Z = \frac{rs^T}{s^T s} = \frac{1}{s^T s}(y - Hs)s^T$$

Theorem 44: Broyden (1965)

Let $s, y \in \mathbb{R}^n, s \neq 0, H \in \mathbb{R}^{n \times n}$ be given. Then the unique optimal solution of (P_1) is

$$B := H + \frac{1}{s^T s} (y - Hs)s^T.$$

Setting $v := H^T y$ and $B := H + \frac{1}{v^T s} (y - Hs)v^T$ leads to “Broyden’s Second Method”

Bad Broyden

Good Broyden

Let us modify problem (P_1) by requiring $B \in \mathbb{S}^n$ (and $H \in \mathbb{S}^n$ in the data). Consider

$$\begin{aligned} \min \quad & \|B - H\|_F \\ \text{s.t.} \quad & Bs = y \\ & B = B^T \\ & B \in \mathbb{R}^{n \times n} \end{aligned}$$

Theorem 45: Powell (1970)

The unique optimal solution ($s \neq 0$) of the above problem is given by

$$B := H + \frac{1}{s^T s} \left[(y - Hs)s^T + s(y - Hs)^T - (y - Hs)ss^T \right].$$

In the formula above, B may *not* be *positive definite* even if the curvature condition is satisfied ($y^T s > 0$) and H is symmetric positive definite.

We want B to be symmetric, positive definite, provided $H \succ 0$ and $y^T > 0$. We consider solving

$$\begin{aligned} \min \quad & \|W^{1/2}(B - H)W^{1/2}\|_F \\ \text{s.t.} \quad & Bs = y \\ & B \in \mathbb{S}^n \end{aligned} \tag{P_W}$$

where

$$W := \left[\int_0^1 \nabla^2 f(x^{(k)} + t\alpha_k d^{(k)}) dt \right]^{-1},$$

but any $W \in \mathbb{S}_{++}^n$ satisfying $Wy^{(k)} = s^{(k)}$ works.

Theorem 45

For every $H \in \mathbb{S}_{++}^n, y, s \in \mathbb{R}^n$ such that $y^T s > 0$ and $W \in \mathbb{S}_{++}^n$ such that $Wy^{(k)} = s^{(k)}$, the unique solution of (P_W) is

$$B := \left(I - \frac{ys^T}{y^T s} \right) H \left(I - \frac{sy^T}{y^T s} \right) + \frac{yy^T}{y^T s}.$$

Moreover, $B \in \mathbb{S}_{++}^n$.

This is called the Davidon–Fletcher–Powell update.

Note that

$$B^{-1} = H^{-1} - \frac{H^{-1}yy^TH^{-1}}{y^TH^{-1}y} + \frac{ss^T}{y^Ts} \leftarrow \text{DFP for the Hessian inverse}$$

Sherman–Morrison–Woodbury
Formula applied to the above

Next, consider

$$\begin{aligned} \min \quad & \|W^{-1/2}(B - H^{-1})W^{-1/2}\|_F \\ \text{s.t.} \quad & By = s \\ & B \in \mathbb{S}^n \end{aligned} \tag{P_W^{BFGS}}$$

Theorem 46

For every $H \in \mathbb{S}_{++}^n$, $y, s \in \mathbb{R}^n$ such that $y^T s > 0$ and $W \in \mathbb{S}_{++}^n$ such that $Wy^{(k)} = s^{(k)}$, the unique solution of (P_W^{BFGS}) is

$$B := \left(I - \frac{sy^T}{y^T s} \right) H \left(I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s}.$$

Moreover, $B \in \mathbb{S}_{++}^n$.

↑
Broyden–Fletcher–Goldfarb–Shanno formula for approximation Hessian inverse

To approximate the Hessian we invert the above formula and obtain (in terms of H as an approximation to the Hessian):

$$H - \frac{Hss^T H}{s^T H s} + \frac{yy^T}{y^T s}$$

In practice, this is most successful with updating its Choleski decomposition.

Summary of DFP, BFGS

$$\mathcal{P} := \{B \in \mathbb{S}^n : Bs = y, B \succ 0\}$$

$$\mathcal{D} := \{B \in \mathbb{S}^n : By = s, B \succ 0\}$$

With

$$W := \left[\int_0^1 \nabla^2 f(x^{(k)} + t\alpha_k d^{(k)}) dt \right]^{-1},$$

DFP: Solve

$$\begin{aligned} \min \quad & \|W^{1/2}(B - H)W^{1/2}\|_F \\ \text{s.t.} \quad & Bs = y \\ & B \in \mathbb{S}^n \end{aligned}$$

BFGS: Solve

$$\begin{aligned} \min \quad & \|W^{-1/2}(B - H^{-1})W^{-1/2}\|_F \\ \text{s.t.} \quad & By = s \\ & B \in \mathbb{S}^n \end{aligned}$$

then the inverse of the solution is the BFGS estimate of the Hessian $\nabla^2 f$.

\mathcal{P} and \mathcal{D} are convex sets. $U \in \mathcal{P} \iff U^{-1} \in \mathcal{D}$.

Therefore, $\forall U \in \mathcal{P}, \forall V \in \mathcal{D}, \forall \lambda \in [0, 1], [\lambda U + (1 - \lambda)V^{-1}] \in \mathcal{P}$ and $[\lambda U^{-1} + (1 - \lambda)V] \in \mathcal{D}$.

Broyden's convex class: $\{\lambda B^{\text{DFP}} + (1 - \lambda)B^{\text{BFGS}} : \lambda \in [0, 1]\}$

2.7.1 Convergence Results**(I) Global**

- (a) Powell (1972): If f is strictly convex (f is convex and $f(\lambda u + (1 - \lambda)v) < \lambda f(u) + (1 - \lambda)f(v)$, for all $\lambda \in (0, 1)$ and $u \neq v$), $\{x \in \mathbb{R}^n : f(x) \leq f(x^{(0)})\}$ is compact, $f \in \mathcal{C}^2$, and exact line search is used then quasi-Newton method based on DFP converges.
- (b) Dixon (1972): If exact line search is used then DFP, BFGS (and many others) all give identical sequence of iterates $\{x^{(k)}\}$ for the same $(x^{(0)}, B_0)$.
- (c) Powell (1976): Same assumptions on f as in (a), but line-search satisfying A-G-W conditions imply global convergence of BFGS.

- (d) Byrd, Nocedal and Yuan (1987): Result of (c) holds for all Broyden's convex class, except for DFP (i.e., $\lambda \in [0, 1)$).

It seems that DFP is worse than BFGS in practice (with inexact line search).

Local Convergence

Assume $f \in \mathcal{C}^2, x^{(k)} \rightarrow \bar{x}, \nabla f(x^{(k)}) \rightarrow 0, \nabla^2 f(\bar{x}) \in \mathbf{S}_{++}^n$.

- (a) Powell (1971): with exact line search, DFP, BFGS both attain Q-superlinear convergence.
- (b) Broyden, Dennis, Moré (1973): If we use $\alpha_k := 1$ for all $k \in \mathbb{Z}_+$ and for suitably small $\epsilon > 0, \delta > 0$, we have $\|x^{(0)} - \bar{x}\| \leq \epsilon$ and $\|B_0 - \nabla^2 f(\bar{x})\| \leq \delta$, then $x^{(k)} \rightarrow \bar{x}$ Q-superlinearly.

Index

C

cone.....	4
continuous optimization problem.....	3
contraction mapping.....	6
convex.....	5
convex cone.....	5
convex function.....	15
curvature condition.....	26

D

decent direction.....	13
differentiable.....	5
direction of negative curvature.....	14
directional derivative.....	5

G

global minimizer.....	3
-----------------------	---

I

improving direction.....	13
isolated local minimizer.....	3

L

level set.....	19
local minimizer.....	3
locally Q-quadratically convergence.....	23

P

positive definite.....	8
positive semidefinite.....	8

Q

Q-superlinear.....	22
--------------------	----

R

regular.....	11
regular value.....	11

S

secant equation.....	25
stationary point.....	13
strict local minimizer.....	3
sublevel set.....	19