

Theorem 2.1 If $f \in C^1$, then the followings are equivalent:

1. f is convex
2. $f(y) \geq f(x) + \nabla^T f(x)(y-x)$ for $\forall x, y \in \text{dom } f$, i.e., f lies above the tangent line.

Proof. 1. From the definition for convexity,

$$f(y) - f(x) \geq \frac{f(\lambda x + (1-\lambda)y) - f(x)}{1-\lambda}$$

Letting $\lambda \rightarrow 1$, the RHS becomes a direction derivative:

$$f(y) - f(x) \geq \nabla^T f(x)(y-x)$$

2. To show the converse, we let $z = \lambda x + (1-\lambda)y$. By applying the inequality in (2.1) twice, we have

$$f(x) \geq f(z) + \nabla^T f(z)(x-z) \quad (2.1)$$

$$f(y) \geq f(z) + \nabla^T f(z)(y-z) \quad (2.2)$$

Letting Eq.(2.1) times λ add Eq.(2.2) times $(1-\lambda)$, we derive that f is convex.

Theorem 2.2 If $f \in C^2$, then the followings are equivalent:

1. f is convex
2. $\nabla^2 f(x) \succeq 0$ for $\forall x \in \text{dom } f$.

Proof. We rewrite $f(y)$ by applying Taylor expansion:

$$f(y) = f(x) + \nabla^T f(x)(y-x) + \frac{1}{2}(y-x)^T \nabla^2 f(x + t(y-x))(y-x), \quad (2.3)$$

for some $t \in [0, 1]$.

1. If f is convex, from Theorem(2.1) and Eq.(2.3), we derive

$$(y-x)^T \nabla^2 f(x + t(y-x))(y-x) \geq 0 \implies \frac{(y-x)^T \nabla^2 f(x + t(y-x))(y-x)}{\|y-x\|^2} \geq 0$$

Set $d := \frac{y-x}{\|y-x\|}$ and let $y \rightarrow x$, we derive

$$d^T \nabla^2 f(x) d \geq 0,$$

which implies $\nabla^2 f(x) \succeq 0$ since d could have an arbitrary direction.

2. To show the converse, due to the semidefiniteness of $\nabla^2 f(x)$, we obtain a new inequality from Eq.(2.3):

$$f(y) \geq f(x) + \nabla^T f(x)(y-x)$$

From Theorem(2.1) we imply f is convex.

Definition 2.4 [Epigraph] The Epigraph of f is given by:

$$\text{Epi}(f) := \{(x, t) \in \mathbb{R}^n \times \mathbb{R} \mid x \in \text{dom } f, t \geq f(x)\} \subseteq \mathbb{R}^{n+1}$$

Theorem 2.3 f is convex iff $\text{Epi}(f)$ is convex.

Proof. 1. Suppose f is convex. For any $(x, t), (y, s) \in \text{Epi}(f)$, it suffices to show

$$(\lambda x + (1-\lambda)y, \lambda t + (1-\lambda)s) \in \text{Epi}(f) \implies \lambda t + (1-\lambda)s \geq f(\lambda x + (1-\lambda)y).$$

Put $\Gamma = \text{epi}(f)$. Suppose first that f is convex, and let $(x_1, t_1), \dots, (x_n, t_n) \in \Gamma$. For any $\lambda_1, \dots, \lambda_n \in [0, 1]$ with $\sum \lambda_i = 1$, the point $(x, t) = (\sum \lambda_i x_i, \sum \lambda_i t_i)$ has

$$t = \sum \lambda_i t_i \geq \sum \lambda_i f(x_i) \geq f(\sum \lambda_i x_i) = f(x).$$

Hence $(x, t) \in \Gamma$, and Γ is convex. The converse is entirely similar.

Choice of Step Length. To get the optimal step length α , differentiating the RHS w.r.t. α leads to

$$(\nabla f(x), d) + \alpha \|d\|^2 = 0 \implies \alpha = -\frac{(\nabla f(x), d)}{\|d\|^2} > 0,$$

which seems a reasonable choice. If d is the steepest descent direction, the step-length becomes a constant:

$$\alpha = \frac{1}{L}.$$

Theorem 2.4 If f is convex in C^1 , then $\nabla f(x) = 0$ is the necessary and sufficient condition for global minimum.

Note that convex function does not have a local minimum that is not global minimum

Proof. If $f \in C^1$ is convex, recall the Theorem(2.1) that

$$f(y) \geq f(x) + \nabla^T f(x)(y-x) \quad (2.4)$$

1. If $\nabla f(x) = 0$, then Eq.(2.3) implies $f(y) \geq f(x)$ for $\forall y$.
2. If x is the global minimum, recall the optimality condition, $\nabla f(x) = 0$.

Definition 4.1 [Q1 Factor] Restrict the attention to a convergent sequence $\{x^k\}$ with limit x^* . Define an error function $e_k = \|x^k - x^*\| \rightarrow 0$. The Q_1 factor of $\{x^k\}$ is given as:

$$Q_1 = \limsup_{k \rightarrow \infty} \frac{e_{k+1}}{e_k}$$

Convergence Rate Analysis. We apply the Lipschitzness to analysis the rate of convergence first. Setting $h(t) = f(x + t\alpha d)$, we find that

$$\begin{aligned} f(x + \alpha d) - f(x) &= h(1) - h(0) = \int_0^1 h'(t) dt \\ &= \int_0^1 \langle \nabla f(x + t \cdot \alpha d), \alpha d \rangle dt \\ &= \int_0^1 \langle \nabla f(x + t \cdot \alpha d), \alpha d \rangle - \langle \nabla f(x), \alpha d \rangle + \langle \nabla f(x), \alpha d \rangle dt \\ &= \langle \nabla f(x), \alpha d \rangle + \int_0^1 \langle \nabla f(x + t \cdot \alpha d) - \nabla f(x), \alpha d \rangle dt \\ &\leq \langle \nabla f(x), \alpha d \rangle + \int_0^1 \|\nabla f(x + t \cdot \alpha d) - \nabla f(x)\| \cdot \|\alpha d\| dt \\ &\leq \langle \nabla f(x), \alpha d \rangle + L \int_0^1 t \alpha^2 \|d\|^2 dt \\ &= \langle \nabla f(x), \alpha d \rangle + \frac{L \alpha^2 \|d\|^2}{2} \end{aligned}$$

Here we want to study the performance of e_k . In our case, we compare $\{e_k\}$ with the geometric progression

$$\beta^k, \quad k = 0, 1, \dots$$

- If there exists $\beta \in (0, 1)$ such that

$$Q_1 \leq \beta,$$

then we can show $e_k \leq \beta^k e_0$ for some $\beta > 0$. In this case $\{e_k\}$ is said to be Q_1 -linear convergent.

- If $Q_1 = 0$, then we say $\{e_k\}$ is Q -super-linear convergent

- If $Q_1 = 1$, then we say $\{e_k\}$ is Q -sub-linear convergent

Definition 4.2 [Q_2 Factor]

$$Q_2 = \limsup_{k \rightarrow \infty} \frac{e_{k+1}^2}{e_k^2}$$

If $Q_2 = M < +\infty$, i.e., $e_{k+1}^2 = O(e_k^2)$, then $\{x^k\} \rightarrow x^*$ Q -quadratically.

To get the optimal solution, it suffices to solve (4.1) for d :

$$d = -(\nabla^2 f(x))^{-1} \nabla f(x),$$

and hence update the solution to be $x \leftarrow x + \alpha d$.

Proposition 4.1 Newton's method guarantees the Q -quadratic convergence.

Proof. Given a nonlinear system $F(x) = 0$, suppose the sequence $\{x^k\}$ is generated by Newton with limit x^* and $F(x^*) = 0$. By Newton's iteration,

$$x^{k+1} = x^k - [F'(x^k)]^{-1} F(x^k), \quad (4.2)$$

which follows that

$$\begin{aligned} x^{k+1} - x^* &= x^k - x^* - [F'(x^k)]^{-1} (F(x^k) - F(x^*)) \\ &= [F'(x^k)]^{-1} (F(x^*) - F(x^k) - F'(x^k)(x^* - x^k)) \end{aligned} \quad (4.3)$$

Note that $F(x^*) = F(x^k) + F'(x^k)(x^* - x^k) + O(\|x^k - x^*\|^2)$, which implies

$$\|x^{k+1} - x^*\| \leq \| [F'(x^k)]^{-1} \| O(\|x^k - x^*\|^2) = O(\|x^k - x^*\|^2).$$

$$x^{k+1} = x^k - \alpha^k \nabla f(x^k) = (I - \alpha^k Q) x^k,$$

$$\begin{aligned} \|x^{k+1}\|^2 &= (x^k)^T (I - \alpha^k Q)^2 x^k \leq \lambda_1 \|(I - \alpha^k Q)^2\| \|x^k\|^2 \\ &= \|x^k\|^2 \cdot \max\{(1 - \alpha^k m)^2, (1 - \alpha^k M)^2\}, \end{aligned}$$

The idea for feasible direction method is that we generate a sequence of $\{x^k\} \subseteq X$ such that $f(x^{k+1}) \leq f(x^k)$, or equivalently, find $\tilde{x}^k \in X$ such that the $\tilde{x}^k - x^k$ is the descent direction:

$$\langle \nabla f(x^k), (\tilde{x}^k - x^k) \rangle \leq 0 \quad (7.7a)$$

and therefore

$$x^{k+1} = x^k + \alpha_k (\tilde{x}^k - x^k), \quad \alpha_k \in (0, 1] \quad (7.7b)$$

The feasibility of x^{k+1} is guaranteed since x^{k+1} is a convex combination of feasible points x^k and \tilde{x}^k . Here the problem remains to find \tilde{x}^k

Projection Gradient Method. One way of finding \tilde{x}^k is to compute $x^* - s_k \nabla f(x^k)$ and project it back into X , as \tilde{x}^k :

$$\tilde{x}^k = [x^* - s_k \nabla f(x^k)]^+, \quad s_k > 0.$$

Conditional Gradient. Another way is to linearize the objective function and solve for \tilde{x}^k :

$$x^k \approx \arg \min_{x \in X} f(x) \implies \tilde{x}^k = \arg \min_{x \in X} f(x^k) + \langle \nabla f(x^k), (x - x^k) \rangle.$$

$$\frac{\partial f^T x}{\partial x} = \frac{\partial f}{\partial x} x + \frac{\partial f}{\partial x}$$

which implies

$$\frac{\|x^{k+1}\|}{\|x^k\|} \leq \max\{|1 - \alpha^k m|, |1 - \alpha^k M|\}$$

Choosing $\alpha^k = \frac{2}{M^2 + m^2}$ s.t. $\max\{|1 - \alpha^k m|, |1 - \alpha^k M|\}$ is maximized, we have

$$\frac{\partial^T A x}{\partial x} = (A + A^T) x$$

$$\frac{\|x^{k+1}\|}{\|x^k\|} \leq \frac{M - m}{M + m}$$

$$\frac{\partial \|y - Ax\|_2^2}{\partial x} = 2 \frac{\partial y - Ax}{\partial x} (y - Ax) = -2A^T(y - Ax)$$

GD with fixed step-size.

$$\|x^{k+1} - x^*\|^2 = \|x^k - x^*\|^2 - 2\alpha \langle g^k, x^k - x^* \rangle + \alpha^2 \|g^k\|^2$$

non-increasing when $0 < \alpha < \frac{2}{L}$.

$$\| \nabla f(x^k) \| = O(\frac{1}{k})$$

$f(x^k), \| \nabla f(x^k) \| \searrow$

convergence rate be q -quadratic if $\nabla^2 L(x, \lambda)$ continuous nonsingular Lipschitz

$$\| \nabla^2 L(x, \lambda) - \nabla^2 L(y, \lambda) \|_F \leq K \| (x, \lambda) - (y, \lambda) \|$$

$$\begin{aligned} \text{MVT: } r(x) - r(x^*) &= (x - x^*)^T \int_0^1 J(x^* + t(x - x^*)) dt \\ x^+ - x^* &= (J^T J)^{-1} J^T \cdot 0 \\ &\leq (J^T J)^{-1} J^T \cdot \int_0^1 L(t) dt \\ &= \frac{1}{2} [J^T J]^{-1} J^T \|x - x^*\|^2 \end{aligned}$$

Quasi-Newton equation

$$B^{k+1} g^k = p^k \quad p^k = x^{k+1} - x^k, \quad g^k = \nabla f(x^{k+1}) - \nabla f(x^k)$$

$$U = V^T [BS] U$$

$$D = \text{diag}(-1/STBS, 1/y^TS)$$

constraint:

$$\langle \nabla f(x^*), (x - x^*) \rangle \geq 0 \quad \forall x \in X.$$

GN: γ_i 's are Lipschitz continuous differentiable

$J(x)$ non-singular near x^*

$$\begin{aligned} x^+ - x^* &= x - x^* - [J^T J]^{-1} J^T r(x) \\ &= [J^T J]^{-1} J^T [J(x - x^*) - (r(x) - r(x^*))] \\ &= [J^T J]^{-1} J^T \int_0^1 J(x^* + t(x - x^*)) dt (x - x^*) \\ \|J(x) - J(x^* + t(x - x^*))\| &\leq L(1-t)\|x - x^*\| \end{aligned}$$

$$\|J(x) - J(x^* + t(x - x^*))\| \leq L(1-t)\|x - x^*\|$$

$f \in \text{convex}$

$\forall (x_i, \lambda_i) \in E_f$ for $y_i = f(x_i)$

$$f(t x_1 + (1-t)x_2) \leq \dots$$

the quadratic (local) rate of convergence is guaranteed if the following conditions hold: From this Lagrange function we define a dual function

hold:

1. there exists x^* such that $F(x^*) = 0$
2. $[F'(x^*)]^{-1}$ exists
3. F is Lipschitz continuous near x^* .

Theorem 5.4 — Convergence Rate for invariant step-size. Given a convex function $f \in C^2$ with Lipschitz gradient, i.e.,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|,$$

the iteration $x^{k+1} = x^k - \frac{1}{L}\nabla f(x^k)$ with local minimum point x^* gives sub-linear convergence rate, i.e.,

$$f(x^k) - f(x^*) \leq \frac{L\|x^0 - x^*\|^2}{k+1}$$

Proposition 5.1 A convex function $f \in C^2$ with Lipschitz gradient constant L has a bounded Hessian matrix:

$$\nabla^2 f(x) \leq LI$$

Proof for proposition(5.1). Otherwise $\exists x_0, v$ such that

$$\|\nabla^2 f(x_0)v\| > L\|v\|.$$

Thus we apply Taylor expansion near x_0 for $x = x_0 + v$:

$$\nabla f(x) = \nabla f(x_0) + \nabla^2 f(x_0)(x - x_0) + o(1)(x - x_0)$$

It follows that

$$\|\nabla f(x) - \nabla f(x_0)\| = \|\nabla^2 f(x_0)(x - x_0) + o(1)(x - x_0)\| \leq L\|x - x_0\|$$

Thus for sufficiently small v , we have $\|\nabla f(x) - \nabla f(x_0)\| \leq L\|x - x_0\|$, which is a contradiction. ■

Step 1: Apply Lipschitz condition. We do the Taylor expansion of f^k for the point x^{k+1} :

$$\begin{aligned} f(x^{k+1}) &= f(x^k) + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{1}{2} \langle \nabla^2 f(x^k)(x^{k+1} - x^k), x^{k+1} - x^k \rangle + o(\|x^{k+1} - x^k\|) \\ &\leq f(x^k) + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L}{2} \|x^{k+1} - x^k\|^2 \\ &= f(x^k) - \frac{1}{2L} \|\nabla f(x^k)\|^2 \end{aligned}$$

implying that

$$\|\nabla f(x^k)\|^2 \leq 2L(f(x^k) - f(x^{k+1}))$$

and therefore

$$\sum_{k=0}^t \|\nabla f(x^k)\|^2 \leq \sum_{k=0}^t 2L(f(x^k) - f(x^{k+1})) \leq 2L(f(x^0) - f(x^*)) \leq L^2\|x^0 - x^*\|^2 \quad (5.5)$$

Step 2: Applying Convexity of f . By the convexity of f and the bound on its

gradient, we can estimate the total error $\sum_{k=0}^t (f(x^k) - f(x^*))$:

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - \frac{1}{L}\nabla f(x^k) - x^*\|^2 \\ &= \|x^k - x^*\|^2 - \frac{2}{L} \langle \nabla f(x^k), x^k - x^* \rangle + \frac{1}{L^2} \|\nabla f(x^k)\|^2 \\ &\leq \|x^k - x^*\|^2 - \frac{2}{L} (f(x^k) - f(x^*)) + \frac{1}{L^2} \|\nabla f(x^k)\|^2 \end{aligned}$$

which implies

$$\begin{aligned} \sum_{k=0}^t f(x^k) - f(x^*) &\leq \sum_{k=0}^t \left[\|x^k - x^*\|^2 - \frac{2}{L} (f(x^k) - f(x^*)) + \frac{1}{L^2} \|\nabla f(x^k)\|^2 \right] \\ &\leq \sum_{k=0}^t \left[\|x^k - x^*\|^2 - \frac{2}{L} (f(x^k) - f(x^*)) \right] + \frac{1}{L^2} \sum_{k=0}^t \|\nabla f(x^k)\|^2 \\ &\leq \sum_{k=0}^t \|x^k - x^*\|^2 \leq L\|x^0 - x^*\|^2 \end{aligned}$$

Step 3: applying monotonicity of $f(x^k) - f(x^*)$. By the monotonicity of $f(x^k) - f(x^*)$,

$$\begin{aligned} f(x^k) - f(x^*) &\leq \frac{1}{r+1} \sum_{i=0}^r f(x^i) - f(x^*) \\ &\leq \frac{L\|x^0 - x^*\|^2}{r+1} \end{aligned}$$

The standard form of linear programming is

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

Define the Lagrange function

$$\begin{aligned} L(x, y, z) &= c^T x - y^T (Ax - b) - z^T x \\ &= (c - A^T y - z)^T x + b^T y \end{aligned}$$

1. The necessary condition for primal and dual optimal points $x^*, (\lambda, \gamma^*)$ is second order sufficient optimality condition is

$$\begin{cases} \nabla_x L(x^*, \lambda^*) = 0, & \nabla_\lambda L(x^*, \lambda^*) = 0 \\ y^T \nabla_{xx}^2 L(x^*, \lambda^*) y > 0, & \forall y \neq 0 \text{ with } \langle \nabla h(x^*), y \rangle = 0 \end{cases}$$

$$Q(y, z) = \inf_x L(x, y, z) = \begin{cases} b^T y, & \text{if } c - A^T y - z = 0 \\ -\infty, & \text{otherwise} \end{cases}$$

For any feasible x and $z \geq 0$, we always have

$$Q(y, x) = c^T x - z^T x \leq c^T x$$

• The first order necessary optimality condition is

$$\nabla_x L(x^*, \lambda^*) = 0, \iff \nabla f(x^*) + \sum_{i=1}^m \lambda_i^* \nabla h_i(x^*) = 0,$$

where λ^* is uniquely determined.

Theorem 7.1 The LP optimality condition is given by:

1. Primal-Feasibility:

$$Ax = b, x \geq 0$$

2. Dual-Feasibility:

$$A^T y + z = c, z \geq 0$$

3. Complementarity/Strong-Duality:

$$x \circ z = 0$$

$$[z]^+ = x^* = \arg \min_{x \in X} \|x - z\|$$

Given the optimization problem

$$\begin{aligned} \min \quad & f(x) \\ \text{such that} \quad & x \in X \text{ is convex} \end{aligned}$$

Recall that The necessary optimality condition for (7.2) is

$$\langle \nabla f(x^*), (x - x^*) \rangle \geq 0 \quad \forall x \in X$$

Proposition 7.4 — First order projection property. Given convex set X , the necessary and sufficient condition for the local minimum x^* for problem (7.3) is:

$$\langle x - [z]^+, z - [z]^+ \rangle \leq 0, \quad \forall x \in X \quad (7.4)$$

Proof. Define $f(x) := \frac{1}{2} \|x - z\|_2^2$, thus $\nabla f(x) = x - z$, with the necessary condition of local minimum x^* as:

$$\langle \nabla f(x^*), (x - x^*) \rangle \geq 0, \quad \forall x \in X \implies \langle x - [z]^+, z - [z]^+ \rangle \leq 0.$$

The sufficiency of this condition is due to the convexity of problem (7.3). ■

Armijo rule:

Let $\sigma \in (0, \frac{1}{2})$. Start with s and continue with $\beta s, \beta^2 s, \dots$, until $\beta^m s$ falls within the set of α with

$$f(x^r) - f(x^r + \alpha d^r) \geq -\sigma \alpha \nabla f(x^r)^T d^r.$$

Claim: if $d^r = -D^r \nabla f(x^r) \neq 0$ with $D^r \succ 0$, then m is finite.

$$\| [z_1]^+ - [z_2]^+ \| \leq \| z_1 - z_2 \|$$

with X to be a convex set.

Proof. Recall the first order property on z_1, z_2 , we obtain

$$\begin{cases} \langle z_1 - [z_1]^+, x - [z_1]^+ \rangle \leq 0, \forall x \in X \\ \langle z_2 - [z_2]^+, x - [z_2]^+ \rangle \leq 0, \forall x \in X \end{cases} \implies \begin{cases} \langle z_1 - [z_1]^+, [z_2]^+ - [z_1]^+ \rangle \leq 0, \forall x \in X \\ \langle z_2 - [z_2]^+, [z_1]^+ - [z_2]^+ \rangle \leq 0, \forall x \in X \end{cases}$$

Adding the inequalities above, we derive

$$\langle z_1 - z_2 + [z_2]^+ - [z_1]^+, [z_2]^+ - [z_1]^+ \rangle \leq 0 \implies \langle [z_2]^+ - [z_1]^+, [z_2]^+ - [z_1]^+ \rangle \leq \langle z_2 - z_1, [z_2]^+ - [z_1]^+ \rangle$$

Applying Cauchy-Schwarz Inequality, we obtain:

$$\| [z_2]^+ - [z_1]^+ \|^2 \leq \langle z_2 - z_1, [z_2]^+ - [z_1]^+ \rangle \leq \| z_2 - z_1 \| \| [z_2]^+ - [z_1]^+ \|$$

Or equivalently, $\| [z_2]^+ - [z_1]^+ \| \leq \| z_2 - z_1 \|$. ■

Proposition 8.4 — KKT conditions. Consider the standard optimization problem

$$\begin{aligned} \min \quad & f_0(x) \\ \text{such that} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p \end{aligned}$$

with the dual problem

$$\begin{aligned} \max \quad & g(\lambda, \gamma) \\ & \lambda_i \geq 0, \quad i = 1, \dots, m \end{aligned}$$

$$(8.7) \quad \nabla_x L(x^*, \lambda^*, \gamma^*) = 0$$

$$(8.8) \quad \begin{aligned} & \lambda_i^* = 0 \text{ for non-active} \\ & \lambda_i^* \geq 0 \end{aligned}$$

* Dual:

$$\begin{aligned} & \max g(\lambda, \gamma) \\ & \lambda_i \geq 0 \end{aligned}$$

Equality constraint:

$$\min f(x) \quad h_i(x) = 0 \quad i=1, \dots, m$$

$$\nabla_x L(x^*, \lambda^*) = 0 \quad \nabla_{\lambda} L(x^*, \lambda^*) = 0$$

$$y^T \nabla_{xx}^2 L(x^*, \lambda^*) y > 0 \quad \forall y \perp \nabla h(x^*)$$

Suff:

$$y^T \nabla_{xx}^2 L(x^*, \lambda^*) y > 0 \quad \forall y \neq 0$$

In equality

$$(8.7) \quad \nabla_x L(x^*, \lambda^*, \gamma^*) = 0$$