# Accelerated Gradient Method

$$\min_x \ f(x) = g(x) + h(x)$$

- $g$ : convex and differentiable
- $h$ : convex

The algorithm:

$$y = (1-\theta_k) \, x_{k-1}^{\dagger} + \theta_k \, u_{k-1}$$

$$x_k = \text{prox}_{t_k h}\left(y - t_k \nabla g(y)\right)$$
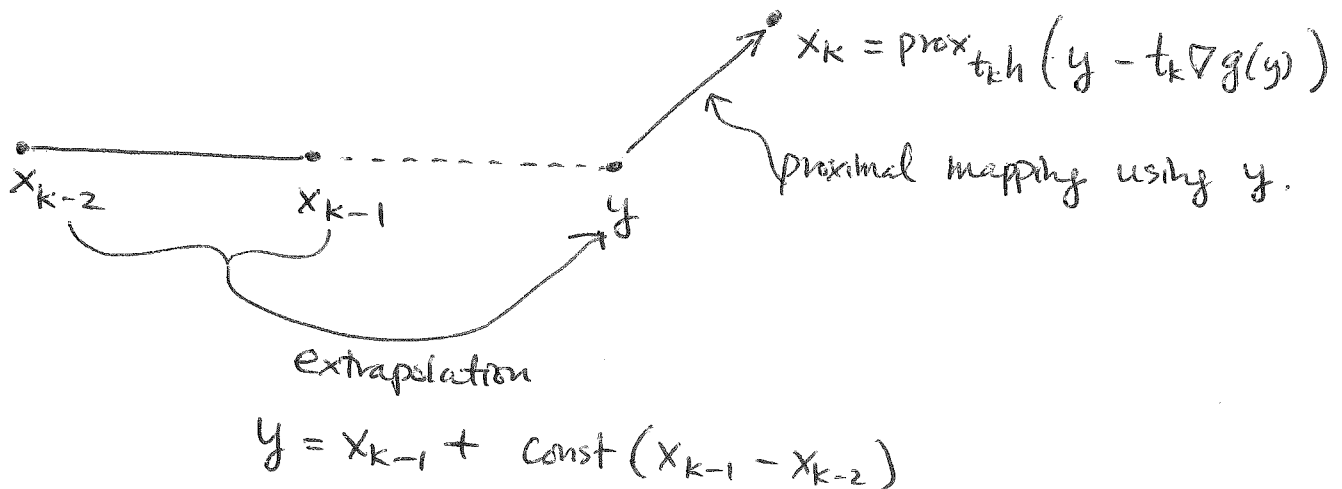
$$u_k = x_{k-1} + \frac{1}{\theta_k}(x_k - x_{k-1})$$

What is it doing?

$$y = (1-\theta_k) x_{k-1} + \theta_k u_{k-1}$$

$$= (1-\theta_k) x_{k-1} + \theta_k \left[ x_{k-2} + \frac{1}{\theta_{k-1}} (x_{k-1} - x_{k-2}) \right]$$

$$= x_{k-1} + \theta_k \left(\frac{1}{\theta_{k-1}} - 1\right)(x_{k-1} - x_{k-2}).$$

and $x_k = \text{prox}_{t_k h}\left(y - t_k \nabla g(y)\right)$.

$x_k = \text{prox}_{t_k h}\left(y - t_k \nabla g(y)\right)$

proximal mapping using $y$.

$x_{k-2}$  $x_{k-1}$  $y$

extrapolation

$$y = x_{k-1} + \text{const} \, (x_{k-1} - x_{k-2})$$
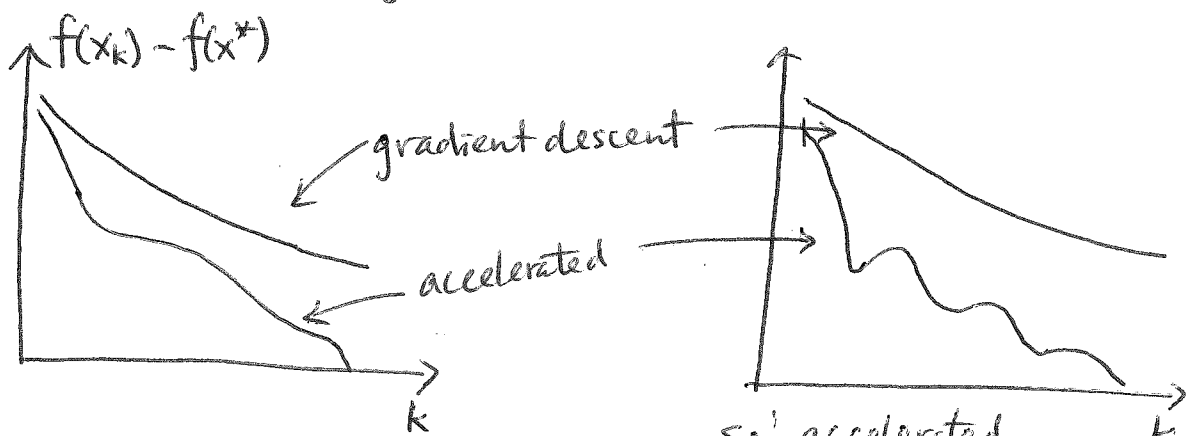
Example ( logistic regression)

$$f(x) = \underbrace{\sum_{i=1}^{n} \left( -y_i \, a_i^T x + \log(1 + \exp(a_i^T x)) \right)}_{g(x)} + \underbrace{0}_{h(x)}$$

$$\nabla g(x) = -A^T \left( y - \cancel{\text{expr}(a_i^T x)} \, P(x) \right),$$

where $P_i(x) = \dfrac{\exp(a_i^T x)}{1 + \exp(a_i^T x)}$.

$\text{prox}_h(x) = x$ because $h(x) = 0$.

Typical convergence plot



Choice of $\theta_k$:

Require:
$$\frac{\theta_k^2}{t_k} \geq (1 - \theta_k) \left( \frac{\theta_{k-1}^2}{t_{k-1}} \right).$$

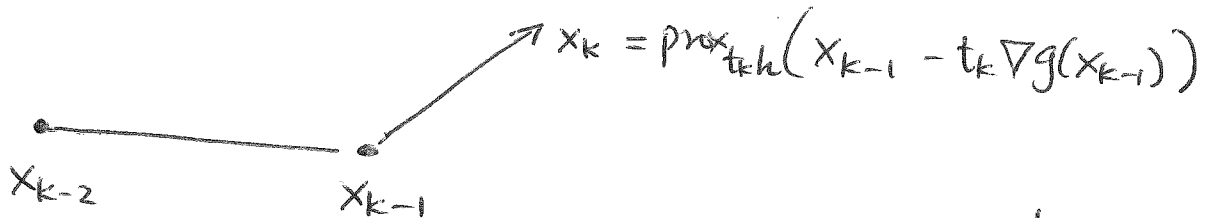$t_k$ can be fixed, e.g. $t_k = \dfrac{1}{L}$, $L = $ Lipschitz constant of $\nabla g$.

then,
$$\theta_k^2 \geq (1 - \theta_k) \, \theta_{k-1}^2$$
$$\Rightarrow \frac{1 - \theta_k}{\theta_k^2} \leq \frac{1}{\theta_{k-1}^2}$$

34

if $\theta_k = 0$, then

$$y = x_{k-1}, \quad x_k = \text{prox}_{t_k h}\left(x_{k-1} - t_k \nabla g(x_{k-1})\right)$$

is the classic gradient projection

$$\nearrow x_k = \text{prox}_{t_k h}\left(x_{k-1} - t_k \nabla g(x_{k-1})\right)$$

$x_{k-2}$ $\bullet$ —————— $\bullet$ $x_{k-1}$

$$g(x) = \frac{1}{2}\|Ax - b\|^2$$
$$\nabla g(x) = A^T(Ax - b)$$
$$h(x) = \lambda\|x\|_1$$

## Example

$$f(x) = \frac{1}{2}\|Ax - b\|^2 + \lambda\|x\|_1$$

algorithm:

$$y = x_{k-1} + \theta_k\left(\frac{1}{\theta_{k-1}} - 1\right)(x_{k-1} - x_{k-2})$$

$$x_k = \text{prox}_{t_k h}\left(y - t_k \nabla g(x_k)\right)$$

$$= \text{argmin}_v \left\{ \lambda t_k\|v\|_1 + \frac{1}{2}\|v - (y - t_k \nabla g(x_k))\|^2 \right\}$$

$$= S_{\lambda t_k}\left(y - t_k \nabla g(x_k)\right)$$

$$= \max\left\{|y - t_k \nabla g(x_k)| - \lambda t_k, \, 0\right\} \, \text{sgn}\left(y - t_k \nabla g(x_k)\right).$$

Example of $\theta_k$:

(i) $\theta_k = \dfrac{2}{k+1}$

$$\frac{1-\theta_k}{\theta_k^2} \overset{?}{\leq} \frac{1}{\theta_{k-1}^2} \iff \frac{1}{\theta_k^2} - \frac{1}{\theta_k} \overset{?}{\leq} \frac{1}{\theta_{k-1}^2}$$

Put $\theta_k = \dfrac{2}{k+1}$, i.e. $\dfrac{1}{\theta_k} = \dfrac{k+1}{2}$, then

$$\left(\frac{k+1}{2}\right)^2 - \frac{k+1}{2} = \frac{k^2+2k+1-2k-2}{4} = \frac{k^2-1}{4}$$

$$\leq \frac{k^2}{4} = \frac{1}{\theta_{k-1}^2}.$$

if we choose $\theta_k = \dfrac{2}{k+1}$, then the algorithm becomes

$$\begin{cases} y = x_{k-1} + \dfrac{k-2}{k+1}(x_{k-1} - x_{k-2}) \\ x_k = \text{prox}_{th}\left(y - t\nabla g(y)\right) \end{cases}$$

(ii) $\theta_k = \dfrac{2}{1+\sqrt{1+\frac{4}{\theta_{k-1}^2}}}$   (FISTA).

A constructive proof:

let $z_k = \dfrac{1}{\theta_k}$. Then

$$\frac{1-\theta_k}{\theta_k^2} = \frac{1}{\theta_{k-1}^2} \iff z_k^2 - z_k = z_{k-1}^2$$

$$\implies z_k = \frac{1 + \sqrt{1+4z_{k-1}^2}}{2}$$

36

# Convergence of Accelerated Gradient Method

Suppose we want to solve

$$\min_{x} f(x) = g(x) + h(x)$$

- $g$ is convex, differentiable, $\nabla g$ is Lipschitz with constant $L > 0$.
- $h$ is convex, and proximal function can be evaluated.

## Algorithm:

$$x_0 = u_0$$

$$y = (1 - \theta_k) x_{k-1} + \theta_k u_{k-1} \qquad \theta_k = \frac{2}{k+1}$$

$$x_k = \text{prox}_{th}(y - t \nabla g(y)) \qquad , \quad t \le \frac{1}{L}.$$

$$u_k = x_{k-1} + \frac{1}{\theta_k}(x_k - x_{k-1})$$

## Theorem:

$$f(x_k) - f(x^*) \le \frac{2 \| x_0 - x^* \|^2}{t(k+1)^2} .$$

That means, the algorithm has $O(\frac{1}{k^2})$ rate of convergence, and this is a first-order method!

## Lemma 1:

The function $g$ satisfies
$$g(x) \leq g(y) + (x-y)^T \nabla g(y) + \frac{L}{2}\|x-y\|^2$$

Proof:
$$g(x) = g(y) + \int_0^1 (x-y)^T \nabla g(y + t(x-y)) \, dt$$
$$= g(y) + \int_0^1 \left[ (x-y)^T \nabla g(y + t(x-y)) - \nabla g(y) + \nabla g(y) \right] dt$$
$$= g(y) + (x-y)^T \nabla g(y) + \int_0^1 (x-y)^T \left[ \nabla g(y + t(x-y)) - \nabla g(y) \right] dt$$
$$\leq g(y) + (x-y)^T \nabla g(y) + \int_0^1 \|x-y\| \| \nabla g(y + t(x-y)) - \nabla g(y)\| \, dt$$
$$\leq g(y) + (x-y)^T \nabla g(y) + \int_0^1 \|x-y\| \, Lt \|x-y\| \, dt$$
$$= g(y) + (x-y)^T \nabla g(y) + \frac{L}{2}\|x-y\|^2.$$

## Lemma 2:

The function $h$ satisfies
$$h(x_{k+1}) \leq h(z) + \frac{1}{t}(x_{k+1} - y)^T (z - x_{k+1}) + \nabla g(y)^T (z - x_{k+1})$$
for any $z$.

Proof: (Requires sub-gradients)

## Proof of Theorem:

First of all, we observe these two equations using lemmas:

$$g(x_{k+1}) \leq g(y) + \nabla g(y)^T(x_{k+1}-y) + \frac{1}{2t}\|x_{k+1}-y\|^2$$

$$h(x_{k+1}) \leq h(z) + \frac{1}{t}(x_{k+1}-y)^T(z-x_{k+1}) + \nabla g(y)^T(z-x_{k+1})$$

Summing the two equations yields

$$f(x_{k+1}) = g(x_{k+1}) + h(x_{k+1})$$

$$\leq g(y) + \nabla g(y)^T\left[(x_{k+1}-y) + (z-x_{k+1})\right] + \frac{1}{2t}\|x_{k+1}-y\|^2$$
$$+ \frac{1}{t}(x_{k+1}-y)^T(z-x_{k+1}) + h(z)$$

$$= \underbrace{g(y) + \nabla g(y)^T\left[z-y\right]}_{g \text{ convex}} + \frac{1}{2t}\|x_{k+1}-y\|^2 + \frac{1}{t}(x_{k+1}-y)^T(z-x_{k+1})$$
$$+ h(z)$$

$$\leq g(z) + h(z) + \frac{1}{2t}\|x_{k+1}-y\|^2 + \frac{1}{t}(x_{k+1}-y)^T(z-x_{k+1})$$

$$= f(z) + \frac{1}{2t}\|x_{k+1}-y\|^2 + \frac{1}{t}(x_{k+1}-y)^T(z-x_{k+1}).$$

Since $z$ is arbitrary we put $z = x_k$ and $z = x^*$:

$$\begin{cases} f(x_{k+1}) \leq f(x^*) + \frac{1}{t}(x_{k+1}-y)^T(x^*-x_{k+1}) + \frac{1}{2t}\|x_{k+1}-y\|^2 \\ f(x_{k+1}) \leq f(x_k) + \frac{1}{t}(x_{k+1}-y)^T(x_k - x_{k+1}) + \frac{1}{2t}\|x_{k+1}-y\|^2 \end{cases}$$

Subtracting the two equations yields

$$0 \leq f(x^*) - f(x_k) + \frac{1}{t}(x_{k+1}-y)^T(x^*-x_k)$$

$$\iff f(x_k) - f(x^*) \leq \frac{1}{t}(x_{k+1}-y)^T(x^*-x_k)$$

So if we consider
$$f(x_{k+1}) - f(x^*) - (1-\theta_k)\left[f(x_k) - f(x^*)\right]$$

$$\leq \frac{1}{t}(x_{k+1}-y)^T(x^* - x_{k+1}) + \frac{1}{2t}\|x_{k+1}-y\|^2 - (1-\theta_k)\frac{1}{t}(x_{k+1}-y)^T(x^*-x_k)$$

$$= \frac{1}{t}(x_{k+1}-y)^T\left[x^* - x_{k+1} - (1-\theta_k)(x^*-x_k)\right] + \frac{1}{2t}\|x_{k+1}-y\|^2$$

$$= \frac{1}{t}(x_{k+1}-y)^T\left[\theta_k x^* + (1-\theta_k)x_k - x_{k+1}\right] + \frac{1}{2t}\|x_{k+1}-y\|^2$$

Recall that
$$\begin{cases} u_{k+1} = x_k + \frac{1}{\theta_k}(x_{k+1} - x_k) \Rightarrow x_{k+1} = \theta_k(u_{k+1}-x_k) + x_k \\ \text{and} \\ \quad y = (1-\theta_k)x_k + \theta_k u_k \end{cases}$$

So, $\theta_k x^* + (1-\theta_k)x_k - x_{k+1}$

$$= \theta_k x^* + (y - \theta_k u_k) - x_{k+1}$$

$$= \theta_k(x^* - u_k) + y - x_{k+1}$$

and hence
$$\to \frac{1}{t}(x_{k+1}-y)^T\left[\theta_k x^* + (1-\theta_k)x_k - x_{k+1}\right]$$

$$= \frac{1}{t}(x_{k+1}-y)^T\left[\theta_k(x^*-u_k) + y - x_{k+1}\right]$$

$$= \frac{\theta_k}{t}(x_{k+1}-y)^T(x^*-u_k) - \frac{1}{t}\|x_{k+1}-y\|^2.$$

Adding $\frac{1}{2t}\|x_{k+1}-y\|^2$, we have
$$\frac{\theta_k}{t}(x_{k+1}-y)^T(x^*-u_k) - \frac{1}{2t}\|x_{k+1}-y\|^2$$

$$= \frac{-1}{2t}\left\{ \|x_{k+1}-y\|^2 - 2\theta_k(x_{k+1}-y)^T(x^*-u_k) + \theta_k^2\|x^*-u_k\|^2 - \theta_k^2\|x^*-u_k\|^2 \right\}$$

40

$$= \frac{-1}{2t} \left\{ \| x_{k+1} - y - \theta_k(x^* - u_k) \|^2 - \theta_k^2 \| x^* - u_k \|^2 \right\}$$

$$\curvearrowright \quad \theta_k(u_{k+1} - x_k) + x_k$$

So $x_{k+1} - y - \theta_k(x^* - u_k)$

$$= \theta_k(u_{k+1} - x_k) + x_k - y - \theta_k(x^* - u_k)$$

$$= \theta_k u_{k+1} + \underbrace{(1 - \theta_k) x_k - y + \theta_k u_k - \theta_k x^*}_{= 0}$$

$$= \theta_k(u_{k+1} - x^*).$$

$$= \frac{-\theta_k^2}{2t} \left\{ \| u_{k+1} - x^* \|^2 - \| u_k - x^* \|^2 \right\}$$

$$= \frac{\theta_k^2}{2t} \left\{ \| u_k - x^* \|^2 - \| u_{k+1} - x^* \|^2 \right\}$$

Thus, we have shown that

$$f(x_{k+1}) - f(x^*) - (1 - \theta_k) \left[ f(x_k) - f(x^*) \right]$$

$$\leq \frac{\theta_k^2}{2t} \left\{ \| u_k - x^* \|^2 - \| u_{k+1} - x^* \|^2 \right\}$$

$$\Longleftrightarrow \frac{t}{\theta_k^2} \left[ f(x_{k+1}) - f(x^*) \right] + \frac{1}{2} \| u_{k+1} - x^* \|^2 \leq \frac{(1 - \theta_k) t}{\theta_k^2} \left[ f(x_k) - f(x^*) \right]$$

$$+ \frac{1}{2} \| u_k - x^* \|^2$$

$$\Longrightarrow \frac{t}{\theta_k^2} \left[ f(x_{k+1}) - f(x^*) \right] + \frac{1}{2} \| u_{k+1} - x^* \|^2$$

$$\leq \frac{t}{\theta_{k-1}^2} \left[ f(x_k) - f(x^*) \right] + \frac{1}{2} \| u_k - x^* \|^2 , \quad \frac{1 - \theta_k}{\theta_k^2} \leq \frac{1}{\theta_{k-1}^2}$$

$$\vdots \qquad \qquad \qquad \qquad \overbrace{\qquad}^{u_0 = x_0}$$

$$\leq t \underbrace{\left( \frac{1 - \theta_1}{\theta_1} \right)}_{= 0 \text{ because } \theta_1 = 1.} \left[ f(x_0) - f(x^*) \right] + \frac{1}{2} \| u_0 - x^* \|^2$$

$$= \frac{1}{2} \| x_0 - x^* \|^2$$

$$\Rightarrow \quad \frac{t}{\theta_k^2}\left[f(x_{k+1}) - f(x^*)\right] \leq \frac{1}{2}\|x_0 - x^*\|^2 - \frac{1}{2}\|u_k - x^*\|^2$$

$$\leq \frac{1}{2}\|x_0 - x^*\|^2$$

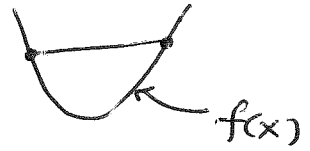$$\Rightarrow \quad f(x_{k+1}) - f(x^*) \leq \frac{\theta_k^2}{2t}\|x_0 - x^*\|^2.$$

So if $\theta_k = \frac{2}{k+1}$, then

$$f(x_{k+1}) - f(x^*) \leq \frac{2}{t(k+1)^2}\|x_0 - x^*\|^2.$$

# Subgradients

1. Convex Review

A function is _convex_ if
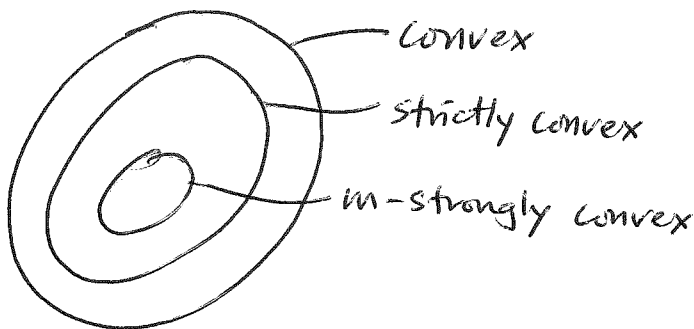
$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y), \quad 0 \leq t \leq 1.$$

_Strictly convex_:    the "$\leq$" becomes "$<$".

$$f(tx + (1-t)y) < tf(x) + (1-t)f(y)$$

m-_Strongly convex_:

$$f(x) \,\,- \frac{m}{2}\|x\|^2 \text{ is convex.}$$

i.e. $f(x)$ is "at least quadratic".

convex

strictly convex

m-strongly convex

Claim: if $f$ is m-strongly convex, then $f$ is convex.

Proof:

Let $g(x) = f(x) - \frac{m}{2}\|x\|^2$.

Then,
$$g(tx + (1-t)y) = f(tx+(1-t)y) - \frac{m}{2}\|tx+(1-t)y\|^2$$
$$\leq tg(x) + (1-t)g(y) = t\{f(x) - \frac{m}{2}\|x\|^2\} + (1-t)\{f(y) - \frac{m}{2}\|y\|^2\}$$
$$= [tf(x) + (1-t)f(y)] - \frac{m}{2}\{t\|x\|^2 + (1-t)\|y\|^2\}$$

So $f(tx+(1-t)y) \leq tf(x) + (1-t)f(y) + \boxed{\frac{m}{2}\{\|tx+(1-t)y\|^2 - (t\|x\|^2 + (1-t)\|y\|^2)\}} \leq 0$

Note that
$$\|tx + (1-t)y\|^2 = t^2\|x\|^2 + 2t(1-t)x^Ty + (1-t)^2\|y\|^2$$
$$- (t\|x\|^2 + (1-t)\|y\|^2) \qquad\qquad - t\|x\|^2 - (1-t)\|y\|^2$$
$$= -t(1-t)\|x\|^2 + 2t(1-t)x^Ty - t(1-t)\|y\|^2$$
$$= -t(1-t)\|x-y\|^2 \leq 0.$$

43

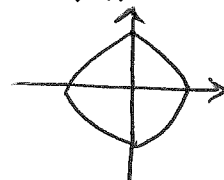# Examples of Convex Functions

(i) $f(x) = a^T x + b$    affine function

(ii) $f(x) = \frac{1}{2} x^T A x + b^T x + c$,    $A \geq 0$.

(iii) $f(x) = \frac{1}{2} \| A x - b \|^2$, because $\frac{1}{2} x^T A^T A x$ has $A^T A \geq 0$.

(iv) $f(x) = \| x \|_p = \left( \sum x_i^p \right)^{1/p}$

(v) $f(x) = \mathbb{1}_\Omega(x) = \begin{cases} 0 &, x \in \Omega \\ +\infty &, x \notin \Omega, \end{cases}$   where $\Omega$ is convex
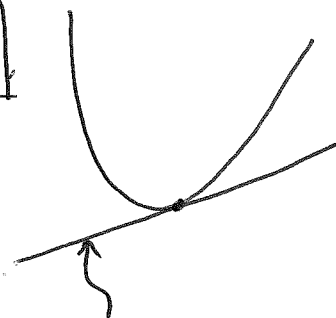
(vi) $f(x) = \max(x_1, x_2, \ldots, x_n)$.

# Properties of a convex function

(i) Suppose that $f$ is differentiable. Then $f$ is convex if and only if $\operatorname{dom} f$ is convex and

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

for all $x, y \in \operatorname{dom} f$.

(ii) Suppose that $f$ is twice differentiable. Then $f$ is convex if and only if $\operatorname{dom} f$ is convex and

$$\boxed{\nabla^2 f(x) \geq 0,}$$

for all $x \in \operatorname{dom} f$.

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

(iii) Jensen's Inequality:

If $f$ is convex, and $X$ is a Random Variable, then

$$\boxed{f(\mathbb{E}[X]) \leq \mathbb{E}(f(x)).}$$

44

**Proof**: Consider the discrete case:
$$\mathbb{E}[X] = \sum_{i=1}^{n} \pi_i x_i \quad, \quad \pi_i = \mathbb{P}(X = x_i)$$
So $f(\mathbb{E}[X]) = f\left(\sum_{i=1}^{n} \pi_i x_i\right) \leq \sum_{i=1}^{n} \pi_i f(x_i) = \mathbb{E}[f(X)]$.

$\underbrace{\qquad}$ Convexity of $f$.

## Operations Preserving Convexity

(i) Non-negative linear combination

if $f_1, f_2, \ldots, f_m$ convex,

then $\sum_{i=1}^{m} w_i f_i$ is also convex, $w_i \geq 0$.

Example: $(a_i^T x - b_i)^2$ is convex, so

$\sum_{i=1}^{n} (a_i^T x - b_i)^2$ is also convex

(ii) Affine Composition

if $f$ is convex, then $f(Ax + b)$ is also convex

Example: $-\sum_{i=1}^{n} \log(x_i)$ is convex.

So $-\sum_{i=1}^{n} \log(a_i^T x + b_i)$ is also convex

(iii) log-sum-exp

The function

$$g(x) = \log\left(\sum_{i=1}^{n} e^{a_i^T x + b_i}\right) \text{ is convex.}$$

**Proof**: Just need to show $f(x) = \log\left(\sum_{i=1}^{n} e^{x_i}\right)$ is convex.

$$(\nabla f)_i = \frac{e^{x_i}}{\sum_{i=1}^{n} e^{x_i}} \quad, \quad (\nabla^2 f)_{ij} = \frac{e^{x_i}}{\sum_{i=1}^{n} e^{x_i}} \mathbb{1}\{i = j\} - \frac{e^{x_i} e^{x_j}}{\left(\sum_{i=1}^{n} e^{x_i}\right)^2}$$

45

We can show that

$$\left|(\nabla^2 f)_{ii}\right| = \left|\frac{e^{x_i}}{\sum_{i=1}^{n} e^{x_i}} - \frac{(e^{x_i})^2}{\left(\sum_{i=1}^{n} e^{x_i}\right)^2}\right|$$

$$\sum_{j \neq i}\left|(\nabla^2 f)_{ij}\right| = \sum_{j \neq i} \frac{e^{x_i} e^{x_j}}{\left(\sum e^{x_i}\right)^2} = \left|\sum_{j=1}^{n} \frac{e^{x_i} e^{x_j}}{\left(\sum e^{x_i}\right)^2} - \frac{(e^{x_i})^2}{\left(\sum e^{x_i}\right)^2}\right|$$

$$= \left|\frac{e^{x_i}}{\sum e^{x_i}} - \frac{(e^{x_i})^2}{\left(\sum e^{x_i}\right)^2}\right|.$$

So $\left|(\nabla^2 f)_{ii}\right| = \sum_{j \neq i}\left|(\nabla^2 f)_{ij}\right|$, and this implies $\nabla^2 f$ is diagonally dominant. Then by Gershgorin Disk Theorem, $\nabla^2 f$ is positive semi-definite.

(iv) Composition:

$$f = h \circ g.$$

$$\Rightarrow f''(x) = h''\left(g(x)\right) g'(x)^2 + h'(g(x)) g''(x).$$

So

| f | h | | g |
|---|---|---|---|
| Convex | Convex | non-decreasing | Convex |
| Convex | Convex | non-increasing | Concave |
| Concave | Concave | non-decreasing | Concave |
| Concave | Concave | non-increasing | Convex |

# 2. Subgradients

The <u>subgradient</u> of a convex function $f$ is any $g \in \mathbb{R}^n$ s.t.
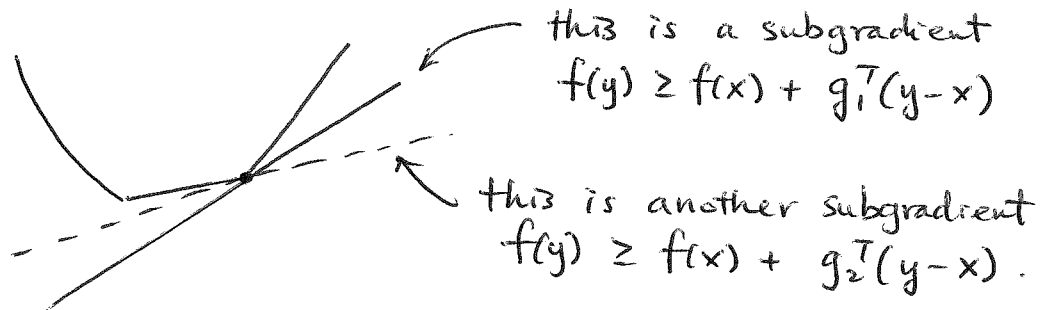
$$f(y) \geq f(x) + g^T(y-x).$$

Why this definition?

Recall that a convex function $f$ always have
$$f(y) \geq f(x) + \nabla f(x)^T(y-x).$$

So a subgradient is a generalization for non-diff. functions.

Pictorial Illustration:



this is a subgradient
$f(y) \geq f(x) + g_1^T(y-x)$

this is another subgradient
$f(y) \geq f(x) + g_2^T(y-x)$.

<u>Sub-differential</u>:

$$\partial f(x^*) = \left\{ g \mid f(y) \geq f(x^*) + g^T(y-x^*) \right\}$$

the set containing all the subgradients at $x^*$.

- $\partial f(x^*)$ is closed and convex
- $\partial f(x^*)$ is non-empty
- if $f$ is differentiable at $x^*$, then $\partial f(x^*) = \{\nabla f(x^*)\}$.

47

- For any $f$,

$$x^* = \underset{x}{\text{argmin}}\ f(x) \iff 0 \in \partial f(x^*)$$

or, in other words,

$$f(y) \geq f(x^*) + 0^T(y - x^*) = f(x^*), \quad \forall y \in \text{dom} f.$$

(this of course implies $x^*$ is the minimizer.)

## Lemma 2 of Accelerated Gradient Method

$$h(x_{k+1}) \leq h(z) + \tfrac{1}{t}(x_{k+1} - y)^T(z - x_{k+1}) + \nabla g(y)^T(z - x_{k+1}).$$

Proof: Note that

$$x_{k+1} = \underset{x}{\text{argmin}}\left\{ h(v) + \tfrac{1}{2t}\|v - (y - t\nabla g(y))\|^2 \right\}$$

$\Longrightarrow$ $0 \in \partial(\cdot)$ implies that

$$\Rightarrow \quad 0 \in \partial h(x_{k+1}) + \tfrac{1}{t}(x_{k+1} - (y - t\nabla g(y)))$$

$$\Rightarrow \quad \tfrac{1}{t}(y - t\nabla g(y) - x_{k+1}) \in \partial h(x_{k+1})$$

$\underset{\substack{\text{def} \\ \text{of subgrad}}}{\Rightarrow} \quad h(z) \geq h(x_{k+1}) + \tfrac{1}{t}(y - t\nabla g(y) - x_{k+1})^T(z - x_{k+1})$

$$\Rightarrow \quad h(x_{k+1}) \leq h(z) + \tfrac{1}{t}(x_{k+1} - (y - t\nabla g(y)))^T(z - x_{k+1})$$

$$= h(z) + \tfrac{1}{t}(x_{k+1} - y)^T(z - x_{k+1})$$

$$+ \nabla g(y)^T(z - x_{k+1}).$$