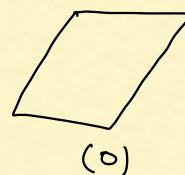
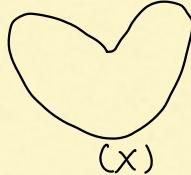
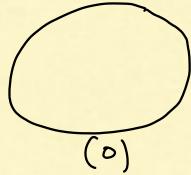


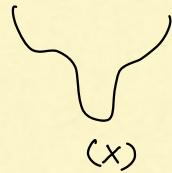
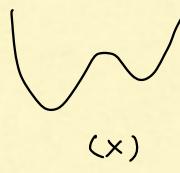
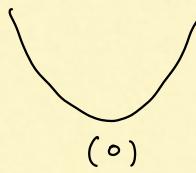
# Opt Lec 1.

## • Convexity

1. Convex set ( $C$ ) : When  $x \in C, y \in C$ , then  $\alpha x + (1-\alpha)y \in C \quad \forall \alpha \in [0, 1]$



2. Convex functions ( $f$ ) :  $\alpha f(x_1) + (1-\alpha)f(x_2) \geq f(\alpha x_1 + (1-\alpha)x_2) \quad \forall x_1, x_2, \forall \alpha \in [0, 1]$



$\Rightarrow$  linear functions ( $a^T x$ ), affine functions ( $a^T x + b$ ), norm  
exponential functions ( $e^{\alpha x}$ ), negative log functions ( $-\log \beta x$ )

★  $f$  is a convex function iff its epigraph is a convex set  
(where epigraph is  $\{(x, t) \mid f(x) \leq t\}$ )

pf) ①  $\Rightarrow$  ②  $f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2) \leq \lambda t_1 + (1-\lambda)t_2 \quad (\because f(x_1) \leq t_1, f(x_2) \leq t_2)$

②  $\Rightarrow$  ①  $f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2) \quad (\because \{(x_1, f(x_1)), (x_2, f(x_2))\} \subset \text{epi}(f))$

★ Convexity is preserved with addition and max operator between convex functions and non-negative scalar multiplication is okay, too.

3. Norms ( $p$ ) : Non-negative-valued scalar functions on a given vector space  $V$ .

Satisfying the following conditions ( $p: V \rightarrow [0, +\infty)$ )

①  $p(x) + p(y) \geq p(x+y)$  for all  $x, y \in V$  (triangular inequality)

②  $p(ax) = |a|p(x)$  (homogeneous or scalable)

③  $p(x) = 0$  if and only if  $x = 0$

(semi-norm satisfies only ① and ②)

\* Frequently used norms

$$\left. \begin{array}{l} \rightarrow 0\text{-norm} : \|x\|_0 = \# \text{ of nonzero element} \\ \rightarrow 1\text{-norm} : \|x\|_1 = \sum_{i=1}^d |x_i| \\ \rightarrow 2\text{-norm} : \|x\|_2 = \sqrt{\sum_{i=1}^d x_i^2} \\ \rightarrow \max\text{-norm} : \|x\|_\infty = \max_i |x_i| \end{array} \right\} \Rightarrow p\text{-norm} : \|x\|_p = \left( \sum_{i=1}^d |x_i|^p \right)^{\frac{1}{p}}$$

★ Every norm is convex

$$\begin{aligned} \text{pf)} \quad & \alpha p(x) + (1-\alpha)p(y) = p(\alpha x) + p((1-\alpha)y) \geq p(\alpha x + (1-\alpha)y) \\ & \therefore \alpha p(x) + (1-\alpha)p(y) \geq p(\alpha x + (1-\alpha)y) \end{aligned}$$

4. Jensen's inequality.

$\Rightarrow$  Let  $f$  is convex,  $x_1, \dots, x_m \in \text{dom}(f)$ ,  $\lambda_1, \dots, \lambda_m \in \mathbb{R}_+$  s.t.  $\sum_{i=1}^m \lambda_i = 1$

$$\text{Then, } f\left(\sum_{i=1}^m \lambda_i x_i\right) \leq \sum_{i=1}^m \lambda_i f(x_i) \quad (\Leftrightarrow f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)])$$

pf) Mathematical Induction

Step 1. Show  $m=2$  is correct

$$\Rightarrow f(\lambda_1 x_1 + \lambda_2 x_2) \leq \lambda_1 f(x_1) + \lambda_2 f(x_2) \quad (\because f \text{ is convex})$$

Step 2. Show  $m=k+1$  is correct assuming  $m=k$  is correct

$$\begin{aligned} \Rightarrow f\left(\sum_{i=1}^{k+1} \lambda_i x_i\right) &= f\left((1-\lambda_{k+1}) \sum_{i=1}^k \frac{\lambda_i x_i}{1-\lambda_{k+1}} + \lambda_{k+1} x_{k+1}\right) \\ &\leq (1-\lambda_{k+1}) f\left(\sum_{i=1}^k \frac{\lambda_i x_i}{1-\lambda_{k+1}}\right) + \lambda_{k+1} f(x_{k+1}) \\ &\leq (1-\cancel{\lambda_{k+1}}) \sum_{i=1}^k \frac{\lambda_i}{1-\cancel{\lambda_{k+1}}} f(x_i) + \cancel{\lambda_{k+1}} f(x_{k+1}) \\ &= \sum_{i=1}^{k+1} \lambda_i f(x_i) \end{aligned}$$

$$\therefore f\left(\sum_{i=1}^m \lambda_i x_i\right) \leq \sum_{i=1}^m \lambda_i f(x_i)$$

## • Young's inequality

When  $1 < p < \infty$ ,  $a > 0, b > 0$ , then  $ab \leq \frac{p-1}{p}a^{\frac{p}{p-1}} + \frac{1}{p}b^p$

(When  $p=2$ , it turns to Cauchy-Schwarz inequality)

$$\begin{aligned} \text{pf)} \quad ab &= \exp(\log a + \log b) \\ &= \exp\left(\frac{p-1}{p} \cdot \frac{p}{p-1} \log a + \frac{1}{p} \cdot p \log b\right) \\ &\leq \frac{p-1}{p} \exp\left(\frac{p}{p-1} \log a\right) + \frac{1}{p} \exp(p \log b) \quad (\because \text{convexity}) \\ &= \frac{p-1}{p} \cdot a^{\frac{p}{p-1}} + \frac{1}{p} b^p \end{aligned}$$

## • Hölder's inequality

When  $x, y \in V$ , then  $x^T y \leq \|x\|_q \|y\|_p$  where  $\frac{1}{p} + \frac{1}{q} = 1$

$$\begin{aligned} \text{pf)} \quad \sum_{i=1}^{\infty} \frac{x_i}{\|x\|_q} \cdot \frac{y_i}{\|y\|_p} &\leq \sum_{i=1}^{\infty} \frac{|x_i|}{\|x\|_q} \cdot \frac{|y_i|}{\|y\|_p} \\ &\leq \sum_{i=1}^{\infty} \left( \frac{1}{q} \frac{|x_i|^q}{\|x\|_q^q} + \frac{1}{p} \cdot \frac{|y_i|^p}{\|y\|_p^p} \right) \quad (\because \text{Young's inequality}) \\ &= \frac{1}{q} + \frac{1}{p} = 1 \end{aligned}$$

$$\Rightarrow \sum_{i=1}^{\infty} x_i \cdot y_i \leq \|x\|_q \cdot \|y\|_p$$

## • Minkowski's inequality

When  $1 < p < \infty$ ,  $x, y \in V$ , then  $\|x+y\|_p \leq \|x\|_p + \|y\|_p$

$$\begin{aligned} \text{pf)} \quad (\|x+y\|_p)^p &= \sum_i |x_i + y_i|^p \\ &= \sum_i |x_i + y_i| \cdot |x_i + y_i|^{p-1} \\ &= \sum_i (|x_i| + |y_i|) |x_i + y_i|^{p-1} \\ &= \sum_i \left( |x_i| \cdot |x_i + y_i|^{p-1} + |y_i| \cdot |x_i + y_i|^{p-1} \right) \\ &\leq \left\{ \left( \sum_i |x_i|^p \right)^{\frac{1}{p}} + \left( \sum_i |y_i|^p \right)^{\frac{1}{p}} \right\} \left( \sum_i |x_i + y_i|^{p-1} \right)^{\frac{p-1}{p}} \quad (\because \text{Hölder's inequality}) \\ &= (\|x\|_p + \|y\|_p) (\|x+y\|_p)^{p-1} \end{aligned}$$

$$\therefore \|x+y\|_p \leq \|x\|_p + \|y\|_p$$

# Opt Lec 2.

- Convex optimization

- Convex optimization problem

$$\Rightarrow \min f(x) \quad \text{s.t. } x \in C$$

↑  
convex function
↑  
convex set

- Convergence rate ( $g(t)$ )

$$\Rightarrow \frac{f(x_t) - f(x^*)}{g(t)} \leq \text{constant}$$

$$\text{ex1) } g(t) = \frac{1}{t} \rightarrow f(x_t) - f(x^*) \leq \varepsilon \text{ when } t \geq \frac{1}{\varepsilon}$$

$$\text{ex2) } g(t) = \frac{1}{t^2} \rightarrow f(x_t) - f(x^*) \leq \varepsilon \text{ when } t \geq \frac{1}{\varepsilon^2}$$

$$\text{ex3) } g(t) = e^{-t} \rightarrow f(x_t) - f(x^*) \leq \varepsilon \text{ when } t \geq -\ln \varepsilon$$

- Differentiable

①  $f$  is differentiable at  $x_0$

$\Rightarrow$  there exists a linear map  $J: \mathbb{R}^m \rightarrow \mathbb{R}^n$  such that

$$\lim_{h \rightarrow 0} \frac{\|f(x_0 + h) - f(x_0) - J(h)\|}{\|h\|} = 0$$

②  $f$  is differentiable

$\Rightarrow \nabla f(x) = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right)^T$  exists at every point  $x \in \text{dom}(f)$

③  $f$  is twice differentiable

$$\Rightarrow \nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_d \partial x_d} \end{bmatrix} \text{ exists at every point } x \in \text{dom}(f)$$

4. Tangent hyperplane of  $f$  at  $x_0$  ( $g(x)$ )

$$\Rightarrow g(x) = f(x_0) + \nabla f(x_0)^T (x - x_0)$$

5. Local minima ( $x$ ): there exist  $\varepsilon > 0$  s.t.  $f(x) \leq f(y) \forall y$  s.t.  $\|y - x\| \leq \varepsilon$

6. Critical point ( $x$ ):  $\nabla f(x) = 0$

(When  $f$  is convex, a critical point is a global minima)

★ 1st order characterization of convexity

$\Rightarrow$  Let  $\text{dom}(f)$  is open,  $f$  is differentiable

Then,  $f$  is convex iff  $\text{dom}(f)$  is convex & 
$$f(y) \geq f(x) + \nabla f(x)^T (y - x) \quad \forall x, y \in \text{dom}(f)$$

(if  $f$  is non-differentiable,  $\nabla f(x) \approx g(x) \in \partial f(x)$ )

( $f$  is lower-bounded by its tangent hyperplane for every point in domain)

$$\text{pf) } ① \Rightarrow ② \quad \lambda f(y) + (1-\lambda) f(x) \geq f(\lambda y + (1-\lambda)x) = f(\lambda(y-x) + x)$$

$$\Rightarrow \lambda f(y) \geq \lambda f(x) + f(x + \lambda(y-x)) - f(x)$$

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

$$\Rightarrow f(y) \geq f(x) + \nabla f(x)^T (y-x) \quad \text{when } \lambda \rightarrow 0$$

$$② \Rightarrow ① \quad \text{Let } z = (1-\lambda)x + \lambda y$$

$$f(x) \geq f(z) + \nabla f(z)^T (x-z) \quad \times (1-\lambda)$$

$$+ \left( f(y) \geq f(z) + \nabla f(z)^T (y-z) \right) \quad \times \lambda$$

$$(1-\lambda)f(x) + \lambda f(y) \geq f(z) + \nabla f(z)^T ((1-\lambda)x + \lambda y - z)$$

$$\Rightarrow (1-\lambda)f(x) + \lambda f(y) \geq f((1-\lambda)x + \lambda y)$$

## ★ 2nd order characterization of convexity

⇒ Let  $\text{dom}(f)$  is open,  $f$  is twice differentiable

Then,  $f$  is convex iff  $\text{dom}(f)$  is convex &  $\nabla^2 f(x)$  is positive semidefinite  $\forall x \in \text{dom}(f)$

$\left( \begin{array}{l} A \text{ is symmetric \& positive semidefinite : } x^T A x = x^T P D P^T x = y^T D y \geq 0 \quad \forall x \in \mathbb{R}^d \\ \Rightarrow \text{eigenvalues are all non-negative} \end{array} \right) \downarrow \text{diagonalizable with orthogonal matrix } (P) \rightarrow P^T = P^{-1}$

pf) ① ⇒ ②  $g(y) := f(y) - f(x) - \nabla f(x)^T (y-x) \rightarrow \text{convex + affine} = \text{convex}$

$$\Rightarrow \nabla g(y) = \nabla f(y) - \nabla f(x) \rightarrow \nabla g(x) = 0 \Rightarrow \left( \begin{array}{l} x \text{ is a global minimizer of } g \\ \nabla^2 g(x) \text{ is positive semi definite} \\ (\because \text{2nd order necessity condition for a minimizer}) \end{array} \right)$$

$$\Rightarrow \nabla^2 g(y) = \nabla^2 f(y) \rightarrow \nabla^2 g(x) = \nabla^2 f(x)$$

$\Rightarrow \nabla^2 f(x)$  is positive semi definite  $\forall x$  since  $x$  was arbitrary

$$\begin{aligned} ② \Rightarrow ① \quad f(y) &= f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(x+t(y-x)) (y-x) \quad \text{for some } t \in [0,1] \\ &\Rightarrow f(y) \geq f(x) + \nabla f(x)^T (y-x) \quad (\geq 0 \because \nabla^2 f \text{ is ps.d.}) \end{aligned}$$

## ★ Taylor theorem

$$\Rightarrow f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(c) (y-x)$$

where  $c = x + (y-x)t$ ,  $t \in [0,1]$

$$(\frac{1}{2}\alpha \|y-x\|^2 \leq \frac{1}{2} (y-x)^T \nabla^2 f(c) (y-x) \leq \frac{1}{2}\beta \|y-x\|^2)$$

$$pf) \quad f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T M (y-x)$$

$$\text{Let } g(t) = -f(y) + f(t) + \nabla f(t)^T (y-t) + \frac{1}{2} (y-t)^T M (y-t)$$

$$\Rightarrow g(x) = g(y) = 0$$

$$\Rightarrow \exists t \in [0,1] \text{ s.t. } c = x + t(y-x), \nabla g(c) = 0 \quad (\because \text{Rolle's theorem})$$

$$\Rightarrow M = \nabla^2 f(c)$$

## ★ Constrained minimization

$\Rightarrow$  Let  $f$  is convex,  $X$  is a convex set

Then,  $x^*$  is a minimizer of  $f$  over  $X$  iff  $\nabla f(x^*)^\top (y - x^*) \geq 0$

$$\downarrow$$

$$(f(x^*) \leq f(y) \quad \forall y \in X)$$

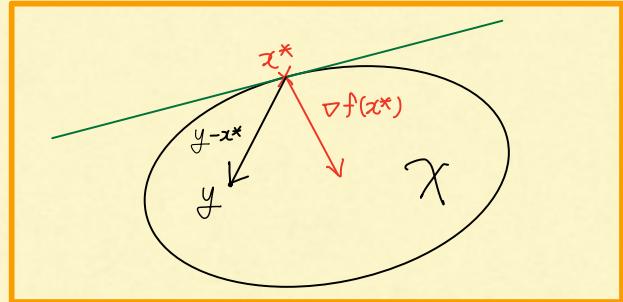
pf.)  $\textcircled{1} \Rightarrow \textcircled{2}$

Case 1.  $x^*$  is inside  $X$

$$\Rightarrow \nabla f(x^*) = 0$$

Case 2.  $x^*$  is on the boundary of  $X$

$$\Rightarrow \nabla f(x^*)^\top (y - x^*) \geq 0$$



$$\textcircled{2} \Rightarrow \textcircled{1} \quad f(y) \geq f(x^*) + \underbrace{\nabla f(x^*)^\top (y - x^*)}_{(\geq 0)} \geq f(x^*)$$

# Opt Lec 3.

- Gradient descent

1. Gradient descent algorithm.

$$\Rightarrow x_{t+1} = x_t - \gamma \nabla f(x_t)$$

( $\gamma$  could be a function of  $t \Rightarrow \gamma_t$ )

2. Convergence rate ( $t$ )

$\Rightarrow$  find  $t$  such that  $f(x_t) - f(x^*) \leq \varepsilon$

3.  $L$ -Lipschitz continuous ( $f$ )

$$\Rightarrow \|f(x) - f(y)\| \leq L \|x - y\| \quad \forall x, y$$

(When  $f$  is differentiable, then  $\|\nabla f(x)\| \leq L \quad \forall x$ )

4.  $\beta$ -smooth ( $f$ )  $\xrightarrow{\text{convex}}$

$$\Rightarrow \|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|$$

$$(\Rightarrow f(x) - f(y) \geq \nabla f(x)^T (x - y) - \frac{\beta}{2} \|x - y\|^2 \quad \forall x, y)$$

(When  $f$  is twice differentiable, then  $\|\nabla^2 f(x)\| \leq \beta \quad \forall x$ )

★  $\nabla^2 f(x) \preceq \beta I \Rightarrow$  eigenvalues of Hessian are less than equal to  $\beta$

$$\|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|$$

$$\Rightarrow (\nabla f(x) - \nabla f(y))^T (x - y) \leq \|\nabla f(x) - \nabla f(y)\| \|x - y\| \leq \beta \|x - y\|^2$$

$$\text{Let } g(t) = f(x + (y-x)t)$$

$$\Rightarrow g'(t) - g'(0) = (\nabla f(x + (y-x)t) - \nabla f(x))^T (y - x) \leq \frac{1}{t} \beta \|y - x\|^2 = \beta t \|y - x\|^2$$

$$\Rightarrow f(y) = g(1) = g(0) + \int_0^1 g'(t) dt$$

$$\leq f(x) + \int_0^1 g'(0) + \beta t \|y - x\|^2 dt = f(x) + \int_0^1 \nabla f(x)^T (y - x) + \beta t \|y - x\|^2 dt$$

$$= f(x) + \nabla f(x)^T (y - x) + \frac{\beta}{2} \|y - x\|^2$$

5.  $\alpha$ -strongly convex ( $f$ )  $\xrightarrow{\text{convex}}$

$$\Rightarrow f(x) - f(y) \leq \nabla f(x)^T(x-y) - \frac{\alpha}{2} \|x-y\|^2 \quad \forall x, y$$

$\text{pf}$   $(\Rightarrow \|\nabla f(x) - \nabla f(y)\| \geq \alpha \|x-y\|)$

(When  $f$  is twice differentiable, then  $\|\nabla^2 f(x)\| \geq \alpha \quad \forall x$ )

$\star \nabla^2 f(x) \succeq \alpha I \Rightarrow$  eigenvalues of Hessian are greater than equal to  $\alpha$

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

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

$$(\nabla f(x) - \nabla f(y))^T(x-y) \geq \alpha \|x-y\|^2 \Rightarrow \|\nabla f(x) - \nabla f(y)\| \geq \alpha \|x-y\|$$

$\star \alpha I \preceq \nabla^2 f(x) \preceq \beta I$  if  $f$  is  $\alpha$ -strongly convex and  $\beta$ -smooth

$$\text{pf) } f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{\alpha}{2} \|y-x\|^2$$

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

$$(\nabla f(x) - \nabla f(y))^T(x-y) \geq \alpha \|x-y\|^2$$

$$f(y) \leq f(x) + \nabla f(x)^T(y-x) + \frac{\beta}{2} \|y-x\|^2$$

$$+ ) f(x) \leq f(y) + \nabla f(y)^T(x-y) + \frac{\beta}{2} \|x-y\|^2$$

$$(\nabla f(x) - \nabla f(y))^T(x-y) \leq \beta \|x-y\|^2$$

$$\Rightarrow \alpha \|x-y\|^2 \leq (\nabla f(x) - \nabla f(y))^T(x-y) \leq \beta \|x-y\|^2$$

Let  $x = y + ht$

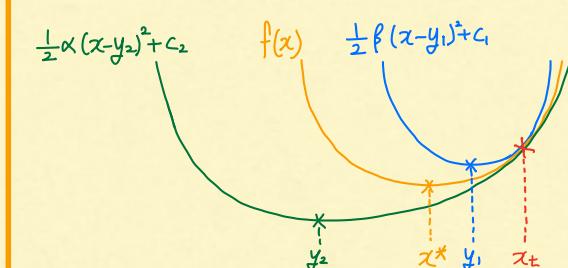
$$\Rightarrow \alpha \|ht\|^2 \leq (\nabla f(y+ht) - \nabla f(y))^T ht \leq \beta \|ht\|^2$$

$$\Rightarrow \alpha \|h\|^2 \leq \left( \frac{\nabla f(y+ht) - \nabla f(y)}{t} \right)^T h \leq \beta \|h\|^2$$

$$\Rightarrow h^T \alpha I h \leq h^T \nabla^2 f(y) \cdot h \leq h^T \beta I h \text{ as } t \rightarrow 0$$

$\rightarrow \nabla^2 f(y) - \alpha I$  is positive semi-definite

$\rightarrow \beta I - \nabla^2 f(y)$  is positive semi-definite



$\star y_2 \leq x^* \leq y_1$

$$(y_1 = x_t - \frac{1}{\beta} \nabla f(x_t))$$

$$(y_2 = x_t - \frac{1}{\alpha} \nabla f(x_t))$$

## 6. Polyak-Łojasiewicz inequality

$$\Rightarrow f(x) - f(x^*) \leq \frac{1}{2\alpha} \|\nabla f(x)\|^2 \quad \forall x \in X$$

( $\alpha$ -strongly convex function satisfies.)

$$\begin{aligned} \text{pf)} \quad f(x) - f(x^*) &\leq \nabla f(x)^T (x - x^*) - \frac{\alpha}{2} \|x - x^*\|^2 \\ &\leq \|\nabla f(x)\| \|x - x^*\| - \frac{\alpha}{2} \|x - x^*\|^2 \quad (\because \text{Hölder's inequality}) \\ &= -\frac{\alpha}{2} \left( \|x - x^*\| - \frac{1}{\alpha} \|\nabla f(x)\| \right)^2 + \frac{1}{2\alpha} \|\nabla f(x)\|^2 \\ &\leq \frac{1}{2\alpha} \|\nabla f(x)\|^2 \end{aligned}$$

## 7. Linear convergence (Sequence of $f(x_0), f(x_1), \dots$ )

$\Rightarrow$  there exists a constant  $c \in (0, 1)$  such that

$$\lim_{t \rightarrow \infty} \frac{f(x_{t+1}) - f(x^*)}{f(x_t) - f(x^*)} = c$$

(When  $\frac{f(x_{t+1}) - f(x^*)}{f(x_t) - f(x^*)} \leq c \quad \forall t$ , then  $f(x_t) - f(x^*) \leq c^t (f(x_0) - f(x^*))$ )

### ① Convergence rate for gradient descent

#### \* Vanilla analysis

$$\begin{aligned} \Rightarrow f(x_t) - f(x^*) &\leq \nabla f(x_t)^T (x_t - x^*) = \frac{1}{\gamma} (x_t - x_{t+1})^T (x_t - x^*) \\ &\leq \frac{1}{2\gamma} \left( \|x_t - x_{t+1}\|^2 + \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 \right) \\ &= \frac{\gamma}{2} \|\nabla f(x_t)\|^2 + \frac{1}{2\gamma} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2) \\ \Rightarrow \sum_{t=0}^{T-1} (f(x_t) - f(x^*)) &\leq \sum_{t=0}^{T-1} \frac{\gamma}{2} \|\nabla f(x_t)\|^2 + \frac{1}{2\gamma} (\underbrace{\|x_0 - x^*\|^2}_{(:= R)} - \underbrace{\|x_T - x^*\|^2}_{(:= R)}) \end{aligned}$$

#### ① L-Lipschitz continuous

$$\begin{aligned} \bullet f(\bar{x}) - f(x^*) &\leq \frac{1}{T} \sum_{t=0}^{T-1} f(x_t) - f(x^*) \quad (\because \text{Jensen's inequality}) \\ &\leq \frac{\gamma}{2} \cdot L^2 + \frac{1}{2\gamma T} R^2 \quad (\because \text{Vanilla analysis}) \\ &\leq \frac{RL}{\sqrt{T}} \quad (\because \gamma^* = \frac{R}{L\sqrt{T}}) \quad \text{minimization} \end{aligned}$$

$\therefore O(\frac{1}{\varepsilon^2})$  steps

## ② $\beta$ -smooth

- $f(x_t) - f(x_{t+1}) \geq \nabla f(x_t)^T (x_t - x_{t+1}) - \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \quad (\because \beta\text{-smooth})$

$$= \left( \gamma - \frac{\beta \gamma^2}{2} \right) \|\nabla f(x_t)\|^2$$

(non-increasing)  $\geq \frac{1}{2\beta} \|\nabla f(x_t)\|^2 \quad (\because \gamma^* = \frac{1}{\beta})$  maximization

(Since  $f$  is non-increasing, set  $\gamma$  to maximize the amount of decrease)

- $\sum_{t=0}^{T-1} (f(x_t) - f(x^*)) \leq \sum_{t=0}^{T-1} \frac{1}{2\beta} \|\nabla f(x_t)\|^2 + \frac{\beta}{2} R^2 \quad (\because \text{Vanilla analysis})$

$$\leq \sum_{t=0}^{T-1} (f(x_t) - f(x_{t+1})) + \frac{\beta}{2} R^2$$

$$= f(x_0) - f(x_T) + \frac{\beta}{2} R^2$$

- $f(x_T) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{\beta}{2T} R^2 \quad (\because \text{non-increasing})$

$\therefore O(\frac{1}{\epsilon})$  steps

## ③ $\alpha$ -strongly convex & $\beta$ -smooth

- $\nabla f(x_t)^T (x_t - x^*) \geq f(x_t) - f(x^*) + \frac{\alpha}{2} \|x_t - x^*\|^2 \quad (\because \alpha\text{-strongly convex})$

- $\nabla f(x_t)^T (x_t - x^*) \leq \frac{\gamma}{2} \|\nabla f(x_t)\|^2 + \frac{1}{2\gamma} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2) \quad (\because \text{vanilla analysis})$

$$\therefore f(x_t) - f(x^*) - \frac{\gamma}{2} \|\nabla f(x_t)\|^2 \leq -\frac{\alpha}{2} \|x_t - x^*\|^2 + \frac{1}{2\gamma} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2)$$

- $RHS \geq LHS \geq f(x_t) - f(x_{t+1}) - \frac{1}{2\beta} \|\nabla f(x_t)\|^2 \quad (\because \text{non-increasing } f)$

$$\geq \frac{1}{2\beta} \|\nabla f(x_t)\|^2 - \frac{1}{2\beta} \|\nabla f(x_t)\|^2 = 0 \quad (\because \beta\text{-smooth})$$

$$\therefore \|x_{t+1} - x^*\|^2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|x_t - x^*\|^2 \leq \left(1 - \frac{\alpha}{\beta}\right)^{t+1} \|x_0 - x^*\|^2$$

- $f(x_T) - f(x^*) \leq \cancel{\nabla f(x^*)^T (x_T - x^*)} + \frac{\beta}{2} \|x_T - x^*\|^2 \leq \frac{\beta}{2} \left(1 - \frac{\alpha}{\beta}\right)^T R^2$

$\therefore O(\log \frac{1}{\epsilon})$  steps

- $f(x_t) - f(x_{t+1}) \geq \frac{1}{2\beta} \|\nabla f(x_t)\|^2 \quad (\because \beta\text{-smooth})$

$$\geq \frac{1}{2\beta} \cdot 2\alpha (f(x_t) - f(x^*)) \quad (\because \text{Polyak Lojasiewicz inequality})$$

$$\Rightarrow f(x_{t+1}) - f(x^*) \leq \left(1 - \frac{\alpha}{\beta}\right) (f(x_t) - f(x^*)) \quad (\because \text{subtract } f(x^*) \text{ on both side})$$

$$\Rightarrow f(x_T) - f(x^*) \leq \left(1 - \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*))$$

# Opt Lec 4.

- Projected Gradient descent

1. Constrained Optimization problem  $\rightarrow$  Convex constrained problem

$$\Rightarrow \text{minimize } f(x) \text{ s.t. } x \in X.$$

(when  $f$  is a convex function  
 $X$  is a closed convex set)

$\begin{cases} \textcircled{1} \text{ Projected Gradient Descent} \\ \textcircled{2} \text{ Transform it into unconstrained version} \end{cases}$

2. Projected Gradient Descent

$$\Rightarrow x_{t+1} = \pi_X(x_t - \gamma \nabla f(x_t)) = \arg \min_{x \in X} \|x - (x_t - \gamma \nabla f(x))\|^2$$

3. Property of projection.

$$\textcircled{1} (x - \pi_X(y))^T (y - \pi_X(y)) \leq 0 \text{ where } x \in X, y \in \mathbb{R}^d$$

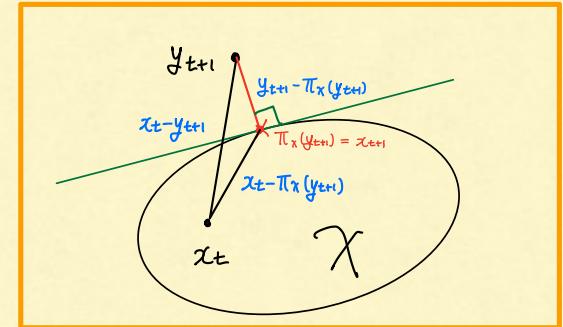
pf) Let  $g(x) = \|x - y\|^2 \rightarrow$  convex

$\Rightarrow \pi_X(y)$  is a minimizer of  $g(x)$

$$\Rightarrow \nabla g(\pi_X(y))^T (x - \pi_X(y)) \geq 0$$

$$\Rightarrow 2(\pi_X(y) - y)^T (x - \pi_X(y)) \geq 0$$

$$\therefore (x - \pi_X(y))^T (y - \pi_X(y)) \leq 0$$



$$\textcircled{2} \|x - \pi_X(y)\|^2 + \|y - \pi_X(y)\|^2 \leq \|x - y\|^2$$

$$\text{pf) } (x - \pi_X(y))^T (y - \pi_X(y)) = \frac{1}{2} (\|x - \pi_X(y)\|^2 + \|y - \pi_X(y)\|^2 - \|x - y\|^2) \leq 0$$

$$\therefore \|x - \pi_X(y)\|^2 + \|y - \pi_X(y)\|^2 \leq \|x - y\|^2$$

① Convergence rate for projected gradient descent

\* Vanilla analysis

$$\begin{aligned}
 \Rightarrow f(x_t) - f(x^*) &\leq \nabla f(x_t)^T (x_t - x^*) = \frac{1}{\gamma} (x_t - y_{t+1})^T (x_t - x^*) \\
 &\leq \frac{1}{2\gamma} \left( \|x_t - y_{t+1}\|^2 + \|x_t - x^*\|^2 - \|y_{t+1} - x^*\|^2 \right) \\
 &= \frac{\delta}{2} \cdot \|\nabla f(x_t)\|^2 + \frac{1}{2\gamma} \left( \|x_t - x^*\|^2 - \|x^* - y_{t+1}\|^2 \right) \\
 &\leq \frac{\delta}{2} \|\nabla f(x_t)\|^2 + \frac{1}{2\gamma} \left( \|x_t - x^*\|^2 - \|x^* - x_{t+1}\|^2 - \|y_{t+1} - x_{t+1}\|^2 \right) \quad (\because \text{2nd property})
 \end{aligned}$$

$$\Rightarrow \sum_{t=0}^{T-1} (f(x_t) - f(x^*)) \leq \sum_{t=0}^{T-1} \frac{\delta}{2} \|\nabla f(x_t)\|^2 + \frac{1}{2\gamma} \left( \|x_0 - x^*\|^2 - \|x_T - x^*\|^2 - \sum_{t=0}^{T-1} \|y_{t+1} - x_{t+1}\|^2 \right) \quad (\because R^2)$$

① L-Lipschitz continuous

$$\begin{aligned}
 \bullet f(\bar{x}) - f(x^*) &\leq \frac{1}{T} \sum_{t=0}^{T-1} (f(x_t) - f(x^*)) \quad (\because \text{Jensen's inequality}) \\
 &\leq \frac{\delta}{2} \cdot L^2 + \frac{1}{2\gamma T} R^2 \quad (\because \text{Vanilla analysis}) \\
 &\leq \frac{RL}{\sqrt{T}} \quad (\because \gamma^* = \frac{R}{L\sqrt{T}}) \quad \text{minimization}
 \end{aligned}$$

$\therefore O(\frac{1}{\varepsilon^2})$  steps

②  $\beta$ -smooth

$$\begin{aligned}
 \bullet f(x_t) - f(x_{t+1}) &\geq \nabla f(x_t)^T (x_t - x_{t+1}) - \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \\
 &= \beta (x_t - y_{t+1})^T (x_t - x_{t+1}) - \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \quad (\because \gamma := \frac{1}{\beta}) \\
 &= \frac{\beta}{2} \left( \|x_t - y_{t+1}\|^2 + \|x_t - x_{t+1}\|^2 - \|y_{t+1} - x_{t+1}\|^2 \right) - \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \\
 &= \frac{\beta}{2} \left( \|x_t - y_{t+1}\|^2 - \|y_{t+1} - x_{t+1}\|^2 \right) \\
 &= \frac{1}{2\beta} \|\nabla f(x_t)\|^2 - \frac{\beta}{2} \|y_{t+1} - x_{t+1}\|^2 \quad (\text{yet, can not guarantee non-increasing}) \\
 &\geq \frac{1}{2\beta} \|\nabla f(x_t)\|^2 - \frac{\beta}{2} \|y_{t+1} - x_t\|^2 + \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \quad (\because \text{2nd property}) \\
 &= \frac{1}{2\beta} \|\nabla f(x_t)\|^2 - \frac{1}{2\beta} \|\nabla f(x_t)\|^2 + \frac{\beta}{2} \|x_t - x_{t+1}\|^2
 \end{aligned}$$

$$(\text{non-increasing}) = \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \geq 0$$

$$\begin{aligned}
 \bullet \sum_{t=0}^{T-1} (f(x_t) - f(x^*)) &\leq \sum_{t=0}^{T-1} \frac{1}{2\beta} \|\nabla f(x_t)\|^2 + \frac{\beta}{2} \left( R^2 - \sum_{t=0}^{T-1} \|y_{t+1} - x_{t+1}\|^2 \right) \quad (\because \text{vanilla analysis}) \\
 &\leq \sum_{t=0}^{T-1} (f(x_t) - f(x_{t+1})) + \frac{\beta}{2} R^2 \\
 &= f(x_0) - f(x_T) + \frac{\beta}{2} R^2
 \end{aligned}$$

$$\bullet f(x_T) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{\beta}{2T} R^2 \quad (\because \text{non-increasing})$$

$\therefore O(\frac{1}{\epsilon})$  steps

③  $\alpha$ -strongly convex &  $\beta$ -smooth

- $\nabla f(x_t)^T(x_t - x^*) \geq f(x_t) - f(x^*) + \frac{\alpha}{2} \|x_t - x^*\|^2 \quad (\because \alpha\text{-strongly convex})$
- $\nabla f(x_t)^T(x_t - x^*) \leq \frac{\beta}{2} \|\nabla f(x_t)\|^2 + \frac{1}{2\beta} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 - \|y_{t+1} - x_{t+1}\|^2)$

$$\therefore f(x_t) - f(x^*) - \frac{\beta}{2} \|\nabla f(x_t)\|^2 \leq -\frac{\alpha}{2} \|x_t - x^*\|^2 + \frac{1}{2\beta} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 - \|y_{t+1} - x_{t+1}\|^2)$$

- RHS  $\geq$  LHS  $\geq f(x_t) - f(x_{t+1}) - \frac{1}{2\beta} \|\nabla f(x_t)\|^2 \quad (\because \text{non-increasing } f)$

$$\begin{aligned} &\geq \frac{1}{2\beta} \|\nabla f(x_t)\|^2 - \frac{\beta}{2} \|y_{t+1} - x_{t+1}\|^2 - \frac{1}{2\beta} \|\nabla f(x_t)\|^2 \\ &= -\frac{\beta}{2} \|y_{t+1} - x_{t+1}\|^2 \end{aligned}$$

$$\therefore \|x_{t+1} - x^*\|^2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|x_t - x^*\|^2 \leq \left(1 - \frac{\alpha}{\beta}\right)^{t+1} \|x_0 - x^*\|^2$$

- $f(x_T) - f(x^*) \leq \nabla f(x^*)^T(x_T - x^*) + \frac{\beta}{2} \|x_T - x^*\|^2$

$$\begin{aligned} &\leq \|\nabla f(x^*)\| \cdot \|x_T - x^*\| + \frac{\beta}{2} \|x_T - x^*\|^2 \quad (\nabla f(x^*) \text{ may not be } 0) \\ &\leq \|\nabla f(x^*)\| \cdot \left(1 - \frac{\alpha}{\beta}\right)^{\frac{T}{2}} R + \frac{\beta}{2} \left(1 - \frac{\alpha}{\beta}\right)^T R^2 \end{aligned}$$

$\therefore O(\log \frac{1}{\epsilon})$  (but slower than gradient descent)

- $f(x_T) - f(x^*) \leq \left(1 - \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*)) \text{ no longer holds}$

( $\because f(x_t) - f(x_{t+1}) \geq \frac{1}{2\beta} \|\nabla f(x_t)\|^2$  no longer holds)

# Opt Lec 5.

- Lagrange Multiplier Theory

## 1. Inequality constrained problem

$\Rightarrow$  minimize  $f(x)$   $\rightarrow$  convex function

s.t.  $h_i(x) = 0$  for  $1 \leq i \leq m$   $\rightarrow$  affine function

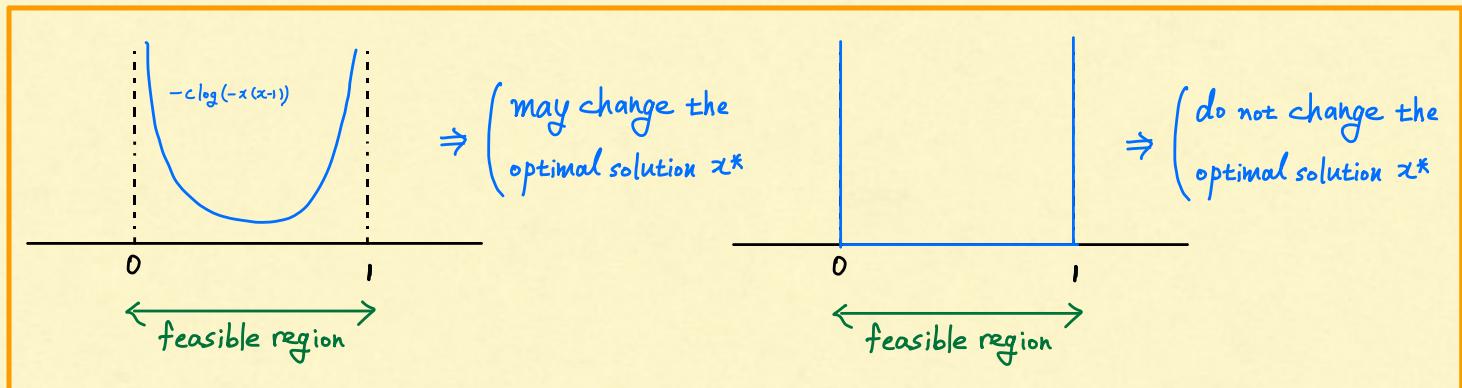
$g_j(x) \leq 0$  for  $1 \leq j \leq r$   $\rightarrow$  convex function

## 2. Barrier method

$\Rightarrow$  add a high cost to infeasibility or approaching boundary from the interior

(transform constrained problem to unconstrained version)

$$\text{ex1) } \min f(x) \text{ s.t. } x(x-1) \leq 0 \rightarrow \underbrace{\min f(x) - c \log(-x(x-1))}_{\text{penalty function}}$$



### 3. Lagrange Dual Problem

$$\Rightarrow \max_{\mu, \lambda} L(\mu, \lambda) \rightarrow \text{Lagrange dual function}$$

s.t.  $\lambda_j \geq 0$  for  $1 \leq j \leq r$

where  $L(\mu, \lambda) := \min_x \Delta(x, \mu, \lambda)$

$\Delta(x, \mu, \lambda) = f(x) + \sum_{i=1}^m \mu_i h_i(x) + \sum_{j=1}^r \lambda_j g_j(x)$

$(\sum_{i=1}^m \mu_i h_i(x) + \sum_{j=1}^r \lambda_j g_j(x) \text{ is a penalty function})$

$\Rightarrow$  updates

$$\mu_i \leftarrow \mu_i + \tau \cdot \frac{\partial \Delta(x^*(\mu, \lambda), \mu, \lambda)}{\partial \mu_i} = \mu_i + \tau h_i(x)$$

$$\lambda_j \leftarrow \lambda_j + \tau \cdot \frac{\partial \Delta(x^*(\mu, \lambda), \mu, \lambda)}{\partial \lambda_j} = \lambda_j + \tau g_j(x)$$

★  $x^*(\mu, \lambda)$  may not be in  $X$ , but  $x^*(\mu^*, \lambda^*)$  is in  $X$ .

★ Lagrange dual function is concave.

$$\begin{aligned} \text{pf)} & L(\alpha \mu^{(1)} + (1-\alpha) \mu^{(2)}, \alpha \lambda^{(1)} + (1-\alpha) \lambda^{(2)}) \\ &= \min_x (f(x) + \sum_{i=1}^m (\alpha \mu_i^{(1)} + (1-\alpha) \mu_i^{(2)}) h_i(x) + \sum_{j=1}^r (\alpha \lambda_j^{(1)} + (1-\alpha) \lambda_j^{(2)}) g_j(x)) \\ &\geq \alpha \left( \min_x f(x) + \sum_{i=1}^m \mu_i^{(1)} h_i(x) + \sum_{j=1}^r \lambda_j^{(1)} g_j(x) \right) + (1-\alpha) \left( \min_x f(x) + \sum_{i=1}^m \mu_i^{(2)} h_i(x) + \sum_{j=1}^r \lambda_j^{(2)} g_j(x) \right) \end{aligned}$$

★ Important remarks on optimal solution

$$\textcircled{1} \quad \nabla_x L(\mu^*, \lambda^*) = 0$$

$$\textcircled{2} \quad h_i(x^*(\mu^*, \lambda^*)) = 0 \quad \forall i$$

$$\textcircled{3} \quad \lambda_j^* g_j(x^*(\mu^*, \lambda^*)) = 0 \quad \forall j$$

$$\begin{cases} g_j(x) < 0 \Rightarrow \lambda_j^* = 0 \\ g_j(x) > 0 \Rightarrow \lambda_j^* = \infty \\ g_j(x) = 0 \Rightarrow \lambda_j^* \geq 0 \end{cases}$$

★ Weight decay can be expressed as Lagrange dual problem

$$\text{pf)} \quad x^* = \arg \min_x f(x) \text{ s.t. } \|x\| \leq c$$

$$\Leftrightarrow x^* = \arg \min_x f(x) + \lambda (\|x\| - c) \text{ s.t. } \lambda \geq 0$$

$$= \arg \min_x f(x) + \lambda \|x\| \text{ s.t. } \lambda \geq 0$$

( $\therefore$  For all  $c$ , there is a corresponding value of  $\lambda$ .)

- Proximal Gradient

1. Proximal gradient descent

$\Rightarrow$  decompose the original objective function:  $f(x) = g(x) + h(x)$

( $g$  is a nice function and  $h$  is a simple additional term)

$$\Rightarrow x_{t+1} = \arg \min_y \left( \frac{1}{2\gamma} \|y - (x_t - \gamma \nabla g(x_t))\|^2 + h(y) \right)$$

$$:= \text{prox}_{h,\gamma}(x_t - \gamma \nabla g(x_t))$$

$$(\text{Generalized Gradient } (G_{h,\gamma}(x_t)) = \frac{1}{\gamma} (x_t - x_{t+1}) = \frac{1}{\gamma} (x_t - \text{prox}_{h,\gamma}(x_t - \gamma \nabla g(x_t)))$$

$$(x_{t+1} = \text{prox}_{h,\gamma}(x_t - \gamma \nabla g(x_t)) \xrightarrow{\text{differentiation}} \nabla h(x_{t+1}) = G_{h,\gamma}(x_t) - \nabla g(x_t))$$

★ G.D. and P.G.D can be generalized to proximal gradient descent.

① Gradient descent  $\rightarrow h(x) = 0$

② Projected gradient descent  $\rightarrow h(x) = 0$  (for  $x \in \mathcal{X}$ ) or  $h(x) = \infty$  (for  $x \notin \mathcal{X}$ )

## ★ Convergence rate for proximal gradient descent

\*  $\beta$  - smooth

$$\bullet g(x_{t+1}) \leq g(x_t) + \nabla g(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2 \quad (\because \beta\text{-smooth})$$

$$= g(x_t) + \nabla g(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|\nabla G_{h,r}(x_t)\|^2$$

$$= g(x_t) + \nabla g(x_t)^T (x_{t+1} - x_t) + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2 \quad (\because r := \frac{1}{\beta})$$

$$\bullet f(x_{t+1}) = g(x_{t+1}) + h(x_{t+1})$$

$$\leq g(x_t) + \nabla g(x_t)^T (x_{t+1} - x_t) + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2 + h(x_{t+1})$$

$$(\because g \text{ is convex}) \leq g(z) + \nabla g(x_t)^T (x_t - z) + \nabla g(x_t)^T (x_{t+1} - x_t) + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2 + h(x_{t+1})$$

$$(\because h \text{ is convex}) \leq g(z) + \nabla g(x_t)^T (x_{t+1} - z) + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2 + h(z) + \nabla h(x_{t+1})^T (x_{t+1} - z)$$

$$= f(z) + \nabla g(x_t)^T (x_{t+1} - z) + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2 + (G_{h,r}(x_t) - \nabla g(x_t))^T (x_{t+1} - z)$$

$$= f(z) + G_{h,r}(x_t)^T (x_{t+1} - z) + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2$$

$$= f(z) + G_{h,r}(x_t)^T (x_t - z + x_{t+1} - x_t) + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2$$

$$= f(z) + G_{h,r}(x_t)^T (x_t - z) - \frac{1}{\beta} \|G_{h,r}(x_t)\|^2 + \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2$$

$$= f(z) + G_{h,r}(x_t)^T (x_t - z) - \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2$$

$\begin{cases} z = x_t \\ z = x^* \end{cases} \therefore f(x_{t+1}) \leq f(x_t) - \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2 \quad (\text{non-increasing})$

$$\bullet f(x_{t+1}) \leq f(x^*) + G_{h,r}(x_t)^T (x_t - x^*) - \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2$$

$$= f(x^*) + \beta \cdot (x_t - x_{t+1})^T (x_t - x^*) - \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2$$

$$= f(x^*) + \frac{\beta}{2} (\|x_t - x_{t+1}\|^2 + \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2) - \frac{1}{2\beta} \|G_{h,r}(x_t)\|^2$$

$$= f(x^*) + \frac{\beta}{2} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2)$$

$$\bullet f(x_T) - f(x^*) \leq \frac{1}{T} \sum_{t=0}^{T-1} (f(x_{t+1}) - f(x^*)) \leq \frac{\beta}{2T} (\|x_0 - x^*\|^2 - \underbrace{\|x_T - x^*\|^2}_{\text{red}}) \leq \frac{\beta}{2T} R^2$$

$\therefore O(\frac{1}{\epsilon})$  steps

- subgradient

1. subgradient ( $g$ )

$$\Rightarrow f(y) \geq f(x) + g^T(y-x) \quad \forall y \in X.$$

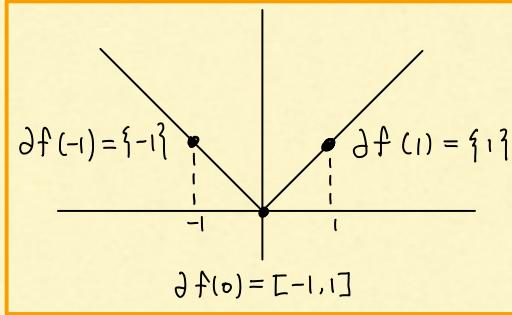
( $\partial f(x)$  := the set of subgradients of  $f$  at  $x$ )

$$\Rightarrow \forall x \in \text{dom}(f), \forall g \in \partial f(x), \|g\| \leq L \Leftrightarrow \forall x \in \text{dom}(f), \underbrace{|f(x) - f(y)| \leq L \|x - y\|}_{(\text{Lipchitz-continuous})}$$

2. subgradient descent

$$\Rightarrow x_{t+1} = x_t - \gamma_t g_t \quad \text{where } g_t \in \partial f(x_t)$$

(projected subgradient descent is also possible)



### ★ Convergence rate for subgradient descent

\* Lipchitz-continuous

$$\bullet f(x_t) - f(x^*) \leq g_t^T(x_t - x^*) = \frac{1}{2} (x_t - x_{t+1})^T (x_t - x^*)$$

$$= \frac{1}{2r} (\|x_t - x_{t+1}\|^2 + \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2)$$

$$= \frac{r}{2} \|g_t\|^2 + \frac{1}{2r} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2)$$

$$\bullet f(\bar{x}) - f(x^*) \leq \frac{1}{T} \sum_{t=0}^{T-1} (f(x_t) - f(x^*)) \quad (\because \text{Jensen's inequality})$$

$$= \frac{r}{2T} \sum_{t=0}^{T-1} \|g_t\|^2 + \frac{1}{2rT} (\|x_0 - x^*\|^2 - \underbrace{\|x_T - x^*\|^2}_{\text{↑}})$$

$$\leq \frac{r}{2} L^2 + \frac{1}{2rT} R^2 \quad (\because \|g_t\| \leq L)$$

$$\leq \frac{RL}{2\sqrt{T}} + \frac{RL}{2\sqrt{T}} = \frac{RL}{\sqrt{T}} \quad (\because r^* = \frac{R}{L\sqrt{T}})$$

$\therefore O(\frac{1}{\varepsilon^2})$  steps

- Mirror descent

1. Dual

⇒ dual space ( $V^*$ ): a vector space consisting of all linear functionals on original vector space with a vector space structure of pointwise addition & scalar multiplication

$$(\forall x, y \in V, f \in V^* \text{ s.t. } f(c_1x + c_2y) = c_1f(x) + c_2f(y))$$

(if  $V = \mathbb{R}^d$ , then  $V^* = \mathbb{R}^d$ )

→ (definition of linear functional)

$$\Rightarrow \text{dual norm} (\|\cdot\|_*): \|g\|_* = \|g\|_q = \sup_{x \in V} g^T x \text{ s.t. } \|x\|_p \leq 1 \quad (x \in V, g \in V^*)$$

★  $\frac{1}{p} + \frac{1}{q} = 1$  where  $p$  is a primal norm and  $q$  is a dual norm

$$\text{pf) } g^T x \leq \|g\|_q \|x\|_p \text{ where } \frac{1}{p} + \frac{1}{q} = 1 \quad (\because \text{Hölder's inequality})$$

$$\Rightarrow g^T x \leq \|g\|_q \quad (\because \|x\|_p \leq 1)$$

★ L-Lipschitz continuous:  $\forall x \in X, g \in \partial f(x), \|g\|_* \leq L$  ( $\|\cdot\|_1$  converge faster than  $\|\cdot\|_2$ )

★  $\beta$ -smooth:  $\forall x, y \in X, \|\nabla f(x) - \nabla f(y)\|_* \leq \beta \|x-y\|$

2. Bregman divergence ( $D_\phi(x, y)$ ):  $\phi(x) - \phi(y) - \nabla \phi(y)^T(x-y)$

$$(\nabla_x D_\phi(x, y) = \nabla \phi(x) - \nabla \phi(y))$$

$$\star (\nabla \phi(x) - \nabla \phi(y))^T(x-z) = D_\phi(x, y) + D_\phi(z, x) - D_\phi(z, y)$$

$$\text{pf) } D_\phi(x, y) + D_\phi(z, x) - D_\phi(z, y)$$

$$= \cancel{\phi(x)} - \cancel{\phi(y)} - \nabla \phi(y)^T(x-y) + \cancel{\phi(z)} - \cancel{\phi(x)} - \nabla \phi(x)^T(z-x) - \cancel{\phi(z)} + \cancel{\phi(y)} + \nabla \phi(y)^T(z-y)$$

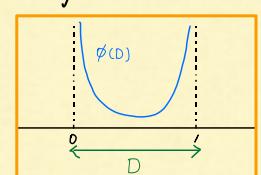
$$= -\nabla \phi(y)^T(x-y) - \nabla \phi(x)^T(z-x) + \nabla \phi(y)^T(z-y) = (\nabla \phi(x) - \nabla \phi(y))^T(z-x)$$

3. Mirror map ( $\phi$ ): a function ( $\phi: D \rightarrow \mathbb{R}^n$ ) where  $X \subset \overline{D}$  &  $D \subset \mathbb{R}^n$  is a convex, open set as follows

$\xrightarrow{\text{closed set of } D}$

①  $\phi$  is strictly convex and differentiable

② The gradient of  $\phi$  takes all possible values  $\rightarrow \nabla \phi(D) = \mathbb{R}^n$



③ The gradient of  $\phi$  diverges on the boundary of  $D \rightarrow \lim_{x \rightarrow \partial D} \|\nabla \phi(x)\| = +\infty$

#### 4. Mirror descent

$$\Rightarrow x_{t+1} = \arg \min_{x \in \mathcal{X}} r_t \nabla f(x_t)^T x + D_\phi(x, x_t)$$

①  $x_t$  is mapped to  $\nabla \phi(x_t)$   $\rightarrow$  Mirror map (primal  $\rightarrow$  dual)

② compute  $\nabla \phi(x_t) - r \nabla f(x_t)$   $\rightarrow$  Gradient descent in dual space

③ Find  $y_{t+1}$  s.t.  $\nabla \phi(y_{t+1}) = \nabla \phi(x_t) - r \nabla f(x_t)$   $\rightarrow$  Inverse mirror map (dual  $\rightarrow$  primal)

④ Projection  $x_{t+1} = \Pi_x^*(y_{t+1}) = \arg \min_{x \in \mathcal{X}} D_\phi(x, y_{t+1})$   $\rightarrow$  Projection in primal space

★ Relationship with proximal gradient descent

$$pf) x_{t+1} = \arg \min_{x \in \mathcal{X}} D_\phi(x, y_{t+1}) = \arg \min_{x \in \mathcal{X}} (\phi(x) - \underbrace{\phi(y_{t+1})}_{\textcircled{1}} - \nabla \phi(y_{t+1})^T (x - \underbrace{y_{t+1}}_{\textcircled{2}}))$$

$$= \arg \min_{x \in \mathcal{X}} (\phi(x) - \nabla \phi(y_{t+1})^T x) = \arg \min_{x \in \mathcal{X}} (\phi(x) - (\nabla \phi(x_t) - r \nabla f(x_t))^T x)$$

$$= \arg \min_{x \in \mathcal{X}} (\phi(x) - \underbrace{\phi(x_t)}_{\textcircled{1}} - \nabla \phi(x_t)^T (x - x_t) + r \nabla f(x_t)^T x)$$

$$= \arg \min_{x \in \mathcal{X}} (\underbrace{r \nabla f(x_t)^T x}_{\textcircled{1}} + \underbrace{D_\phi(x, x_t)}_{\textcircled{2}})$$

$\left( \begin{array}{l} \textcircled{1}: \text{local minimization inducing to move opposite direction of gradient} \\ \textcircled{2}: \text{additional term inducing not to move far from the previous point} \end{array} \right)$

★ Well used mirror map :  $\phi(x_t) = \sum_{i=1}^d x_t^{(i)} \log x_t^{(i)}$  where  $x_t^{(i)}$  is the  $i$ th element of  $x_t$

$$\textcircled{1} \quad \nabla \phi(x_t) = \begin{bmatrix} \log x_t^{(1)} + 1 \\ \vdots \\ \log x_t^{(d)} + 1 \end{bmatrix}$$

$$\textcircled{2} \quad \nabla \phi(x_t) - r \nabla f(x_t) = \begin{bmatrix} \log y_{t+1}^{(1)} + 1 \\ \vdots \\ \log y_{t+1}^{(d)} + 1 \end{bmatrix} = \nabla \phi(y_{t+1})$$

$$\textcircled{3} \quad y_{t+1}^{(i)} = x_t^{(i)} + \exp \left( -r \cdot \frac{\partial f(x_t)}{\partial x_t^{(i)}} \right)$$

$$\textcircled{4} \quad x_{t+1} = \arg \min_{x \in \mathcal{X}} D_\phi(x, y_{t+1}) = \arg \min_{x \in \mathcal{X}} (\phi(x) - \nabla \phi(y_{t+1})^T x)$$

$$= \arg \min_{x \in \mathcal{X}} \left( \sum_{i=1}^d x^{(i)} \log \frac{x^{(i)}}{y_{t+1}^{(i)}} - \|x\|_1 \right) = \frac{y_{t+1}}{\|y_{t+1}\|_1} \quad \text{where } \mathcal{X} = \{x \in \mathbb{R}^d \mid \|x\|_1 = 1, x^{(i)} \geq 0\}$$

(differentiation)

## ★ Convergence rate for mirror descent

\*  $f$  be  $L$ -Lipchitz continuous &  $\phi$  be  $\alpha$ -strongly convex

- $$\nabla_{x_{t+1}} D(x_{t+1}, y_{t+1})^T (x_{t+1} - x) \leq 0 \quad (\because x_{t+1} = \arg\min_{x \in \mathcal{X}} (D_\phi(x, y_{t+1})) \text{ and } D \text{ is convex})$$

$$\Rightarrow \nabla \phi(x_{t+1})^T (x_{t+1} - x) \leq \nabla \phi(y_{t+1})^T (x_{t+1} - x) \leq \nabla \phi(y_{t+1})^T (x_{t+1} - y_{t+1} + y_{t+1} - x)$$

$$\Rightarrow -\nabla \phi(x_{t+1})^T (x - x_{t+1}) - \nabla \phi(y_{t+1})^T (x_{t+1} - y_{t+1}) \leq -\nabla \phi(y_{t+1})^T (x - y_{t+1})$$

$$\therefore D_\phi(x, x_{t+1}) + D_\phi(x_{t+1}, y_{t+1}) \leq D_\phi(x, y_{t+1})$$
- $$D_\phi(x_t, y_{t+1}) - D_\phi(x_{t+1}, y_{t+1}) = \phi(x_t) - \phi(x_{t+1}) - \nabla \phi(y_{t+1})^T (x_t - x_{t+1})$$

$$(\because \alpha\text{-strongly convex}) \leq \nabla \phi(x_t)^T (x_t - x_{t+1}) - \frac{\alpha}{2} \|x_t - x_{t+1}\|^2 - \nabla \phi(y_{t+1})^T (x_t - x_{t+1})$$

$$= (\nabla \phi(x_t) - \nabla \phi(y_{t+1}))^T (x_t - x_{t+1}) - \frac{\alpha}{2} \|x_t - x_{t+1}\|^2$$

$$= \gamma g_t^T (x_t - x_{t+1}) - \frac{\alpha}{2} \|x_t - x_{t+1}\|^2$$

$$\leq \gamma L \|x_t - x_{t+1}\| - \frac{\alpha}{2} \|x_t - x_{t+1}\|^2$$

$$\leq -\frac{\alpha}{2} \left( \|x_t - x_{t+1}\| - \frac{\gamma L}{\alpha} \right)^2 + \frac{\alpha}{2} \cdot \frac{\gamma^2 L^2}{\alpha^2} \leq \frac{\gamma^2 L^2}{2\alpha}$$

- $$f(x_t) - f(x) \leq g_t^T (x_t - x) \quad \text{where } g_t \in \partial f(x_t)$$

$$\leq \frac{1}{\gamma} (\nabla \phi(x_t) - \nabla \phi(y_{t+1}))^T (x_t - x)$$

$$= \frac{1}{\gamma} (D_\phi(x_t, y_{t+1}) + D_\phi(x, x_t) - D_\phi(x, y_{t+1}))$$

$$\leq \frac{1}{\gamma} (D_\phi(x_t, y_{t+1}) + D_\phi(x, x_t) - D_\phi(x, x_{t+1}) - D_\phi(x_{t+1}, y_{t+1}))$$

$$\Rightarrow f(x_t) - f(x^*) \leq \frac{1}{\gamma} (D_\phi(x_t, y_{t+1}) - D_\phi(x_{t+1}, y_{t+1}) + D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1}))$$

$$\leq \frac{\gamma L^2}{2\alpha} + \frac{1}{\gamma} (D_\phi(x^*, x_t) - D_\phi(x^*, x_{t+1}))$$

- $$f(\bar{x}) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{\gamma L^2}{2\alpha} + \frac{1}{\gamma T} (D_\phi(x^*, x_1) - \underbrace{D_\phi(x^*, x_{T+1})}_{R^2})$$

$$\leq \frac{\gamma L^2}{2\alpha} + \frac{R^2}{\gamma T} \quad (R^2 := \sup_{x \in \mathcal{X}} (\phi(x) - \phi(x_1)))$$

$$\leq RL \sqrt{\frac{2}{\alpha T}} \quad (\because \gamma^* = \frac{R}{L} \sqrt{\frac{2\alpha}{T}})$$

$\therefore O(\frac{1}{\epsilon^2})$  steps

# Opt Lec 6.

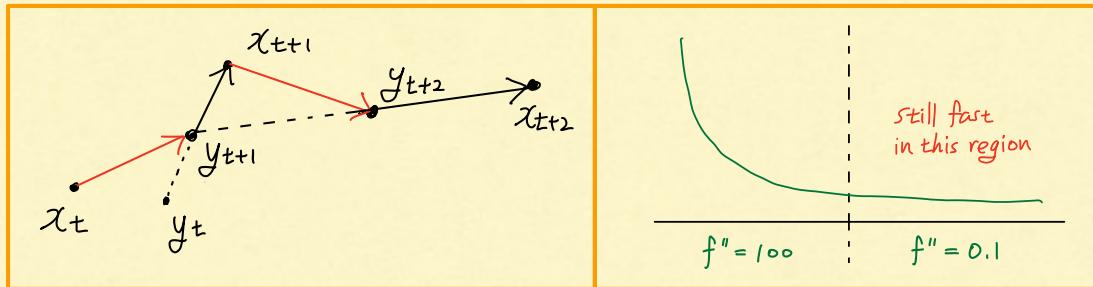
- Nesterov's accelerated gradient descent

## 1. Nesterov's accelerated gradient descent

$$\Rightarrow \begin{cases} y_{t+1} = x_t - \frac{1}{\beta} \nabla f(x_t) \\ x_{t+1} = y_{t+1} + \left( \frac{\sqrt{k}-1}{\sqrt{k+1}} \right) (y_{t+1} - y_t) \text{ where } k = \frac{\beta}{\alpha} \quad (1 \leq k \leq \infty) \end{cases}$$

K=1 K=∞

• Not informative  
• Getting harder



## 2. key properties

Let  $f$  is  $\alpha$ -strongly convex and  $\beta$ -smooth

$$\begin{aligned} \phi_s(x) &= f(x_s) + \frac{\alpha}{2} \|x - x_s\|^2 \\ \phi_{s+1}(x) &= \left(1 - \frac{1}{\sqrt{k}}\right) \phi_s(x) + \frac{1}{\sqrt{k}} \left( f(x_s) + \nabla f(x_s)^T (x - x_s) + \frac{\alpha}{2} \|x - x_s\|^2 \right) \\ \Rightarrow \phi_s(x) &= \text{affine function} + \frac{\alpha}{2} \|x - v_s\|^2 \end{aligned}$$

$$\textcircled{1} \quad \phi_{s+1}(x) - f(x) \leq \left(1 - \frac{1}{\sqrt{k}}\right)^s (\phi_s(x) - f(x))$$

$$\begin{aligned} \text{pf)} \quad \phi_{s+1}(x) - f(x) &= \left(1 - \frac{1}{\sqrt{k}}\right) (\phi_s(x) - f(x)) + \frac{1}{\sqrt{k}} \left( f(x_s) + \nabla f(x_s)^T (x - x_s) + \frac{\alpha}{2} \|x - x_s\|^2 - f(x) \right) \\ &\leq \left(1 - \frac{1}{\sqrt{k}}\right) (\phi_s(x) - f(x)) \quad (\because \alpha\text{-strongly convex}) \\ &\leq \left(1 - \frac{1}{\sqrt{k}}\right)^s (\phi_s(x) - f(x)) \\ &= \left(1 - \frac{1}{\sqrt{k}}\right)^s \left( f(x_s) + \frac{\alpha}{2} \|x - x_s\|^2 - f(x) \right) \\ &\leq \left(1 - \frac{1}{\sqrt{k}}\right)^s \left( \nabla f(x_s)^T (x_s - x) + \frac{\alpha + \beta}{2} \|x - x_s\|^2 \right) \quad (\because \beta\text{-smooth}) \end{aligned}$$

$$\therefore \phi_{s+1}(x) - f(x) \leq \left(1 - \frac{1}{\sqrt{k}}\right)^s (\phi_s(x) - f(x))$$

$$\textcircled{2} \quad f(y_s) \leq \min_x \phi_s(x)$$

$$\bullet \phi_s(x) = \min_x (\phi_s^*(x) + \frac{\alpha}{2} \|x - v_s\|^2) = \phi_s^* + \frac{\alpha}{2} \|x - v_s\|^2$$

$$\bullet \phi_{s+1}(x) = \left(1 - \frac{1}{\sqrt{k}}\right) \phi_s(x) + \frac{1}{\sqrt{k}} (f(x_s) + \nabla f(x_s)^T (x - x_s) + \frac{\alpha}{2} \|x - x_s\|^2)$$

$$= \min_x (\phi_{s+1}(x)) + \frac{\alpha}{2} \|x - v_{s+1}\|^2 = \phi_{s+1}^* + \frac{\alpha}{2} \|x - v_{s+1}\|^2$$

$$\Rightarrow v_{s+1} = -\frac{1}{\alpha} \left( \left(1 - \frac{1}{\sqrt{k}}\right) \cdot \frac{\alpha}{2} \cdot (-2v_s) + \frac{1}{\sqrt{k}} \cdot \nabla f(x_s) + \frac{1}{\sqrt{k}} \cdot \frac{\alpha}{2} (-2x_s) \right)$$

$$= \left(1 - \frac{1}{\sqrt{k}}\right) \cdot v_s + \frac{1}{\sqrt{k}} x_s - \frac{1}{\alpha \sqrt{k}} \nabla f(x_s)$$

$$\Rightarrow \|x_s - v_{s+1}\|^2 = \left\| \left(1 - \frac{1}{\sqrt{k}}\right) (x_s - v_s) + \frac{1}{\alpha \sqrt{k}} \nabla f(x_s) \right\|^2$$

$$= \left(1 - \frac{1}{\sqrt{k}}\right)^2 \|x_s - v_s\|^2 + \frac{1}{\alpha^2 k} \|\nabla f(x_s)\|^2 + 2 \left(1 - \frac{1}{\sqrt{k}}\right) \cdot \frac{1}{\alpha \sqrt{k}} \nabla f(x_s)^T (x_s - v_s)$$

$$\bullet v_{s+1} - x_{s+1} = \left(1 - \frac{1}{\sqrt{k}}\right) v_s + \frac{1}{\sqrt{k}} x_s - \frac{1}{\alpha \sqrt{k}} \nabla f(x_s) - x_{s+1}$$

$$= \left(1 - \frac{1}{\sqrt{k}}\right) \left( x_s + \sqrt{k} (x_s - y_s) \right) + \frac{1}{\sqrt{k}} x_s - \frac{1}{\alpha \sqrt{k}} \nabla f(x_s) - x_{s+1} \quad (\because \text{mathematical induction})$$

$$= \sqrt{k} (x_s - y_s) + y_s - \frac{\sqrt{k}}{\beta} \nabla f(x_s) - x_{s+1}$$

$$= \sqrt{k} y_{s+1} - (\sqrt{k} - 1) y_s - x_{s+1} \quad (\because y_{s+1} = x_s - \frac{1}{\beta} \nabla f(x_s))$$

$$= \sqrt{k} y_{s+1} - (\sqrt{k} - 1) y_s - y_{s+1} - \frac{\sqrt{k} - 1}{\sqrt{k} + 1} (y_{s+1} - y_s)$$

$$= (\sqrt{k} - 1) (y_{s+1} - y_s) \frac{\sqrt{k}}{\sqrt{k} + 1} = \sqrt{k} (x_{s+1} - y_{s+1})$$

$$\bullet \phi_{s+1}(x_s) = \left(1 - \frac{1}{\sqrt{k}}\right) \phi_s(x_s) + \frac{1}{\sqrt{k}} f(x_s)$$

$$= \left(1 - \frac{1}{\sqrt{k}}\right) \phi_s^* + \frac{\alpha}{2} \left(1 - \frac{1}{\sqrt{k}}\right) \|x_s - v_s\|^2 + \frac{1}{\sqrt{k}} f(x_s) = \phi_{s+1}^* + \frac{\alpha}{2} \|x_s - v_{s+1}\|^2$$

$$\Rightarrow \phi_{s+1}^* = \left(1 - \frac{1}{\sqrt{k}}\right) \phi_s^* + \frac{\alpha}{2} \left(1 - \frac{1}{\sqrt{k}}\right) \|x_s - v_s\|^2 + \frac{1}{\sqrt{k}} f(x_s) - \frac{1}{2\beta} \|\nabla f(x_s)\|^2 - \frac{1}{\sqrt{k}} \left(1 - \frac{1}{\sqrt{k}}\right) \nabla f(x_s)^T (x_s - v_s)$$

$$= \left(1 - \frac{1}{\sqrt{k}}\right) \phi_s^* + \frac{\alpha \sqrt{k}}{2} \left(1 - \frac{1}{\sqrt{k}}\right) \|x_s - y_s\|^2 + \frac{1}{\sqrt{k}} f(x_s) - \frac{1}{2\beta} \|\nabla f(x_s)\|^2 + \left(1 - \frac{1}{\sqrt{k}}\right) \nabla f(x_s)^T (x_s - y_s)$$

$$\bullet f(y_{s+1}) \leq f(x_s) - \frac{1}{2\beta} \|\nabla f(x_s)\|^2 \quad (\because \beta-\text{smooth}, r := \frac{1}{\beta})$$

$$= \left(1 - \frac{1}{\sqrt{k}}\right) f(x_s) + \frac{1}{\sqrt{k}} f(x_s) - \frac{1}{2\beta} \|\nabla f(x_s)\|^2$$

$$= \left(1 - \frac{1}{\sqrt{k}}\right) (f(x_s) - f(y_s)) + \left(1 - \frac{1}{\sqrt{k}}\right) f(y_s) + \frac{1}{\sqrt{k}} f(x_s) - \frac{1}{2\beta} \|\nabla f(x_s)\|^2$$

$$\leq \left(1 - \frac{1}{\sqrt{k}}\right) \nabla f(x_s)^T (x_s - y_s) + \left(1 - \frac{1}{\sqrt{k}}\right) \phi_s^* + \frac{1}{\sqrt{k}} f(x_s) - \frac{1}{2\beta} \|\nabla f(x_s)\|^2 \quad (\because \text{mathematical induction})$$

$$\leq \phi_{s+1}^* \quad (\because \frac{\alpha \sqrt{k}}{2} \left(1 - \frac{1}{\sqrt{k}}\right) \|x_s - y_s\|^2 \geq 0)$$

$$\therefore f(y_{s+1}) \leq \phi_{s+1}^* = \min_x (\phi_{s+1}(x))$$

★ Convergence rate for Nesterov's accelerated gradient descent

\*  $\alpha$ -strongly convex &  $\beta$ -smooth

$$\bullet \quad \phi_s(x^*) - f(x^*) \leq \left(1 - \frac{1}{\sqrt{K}}\right)^{s-1} \left( \nabla f(x^*)^\top (x_1 - x^*) + \frac{\alpha + \beta}{2} \|x^* - x_1\|^2 \right)$$

$$= \frac{\alpha + \beta}{2} \|x^* - x_1\|^2 \left(1 - \frac{1}{\sqrt{K}}\right)^{s-1}$$

$$\leq \frac{\alpha + \beta}{2} \|x^* - x_1\|^2 \exp\left(-\frac{s-1}{\sqrt{K}}\right)$$

$$\bullet \quad f(y_t) - f(x^*) \leq \phi_t^* - f(x^*) \leq \phi_t(x^*) - f(x^*) \leq \frac{\alpha + \beta}{2} \|x_1 - x^*\|^2 \exp\left(-\frac{t-1}{\sqrt{K}}\right)$$

$\therefore O(\log \frac{1}{\epsilon})$  steps

# Opt Lec 7.

## • Stochastic Gradient Descent I

### 1. Stochastic oracle ( $\underline{g}$ )

$\Rightarrow$  fixed input point  $x \rightarrow g(x)$  s.t.  $E[g(x)] \in \partial f(x)$

$\Rightarrow$  random input point  $x \rightarrow g(x)$  s.t.  $E[g(x)|x] \in \partial f(x)$

### 2. Stochastic Gradient Descent

$\Rightarrow x_{t+1} = x_t - \gamma g_t$  where  $g_t = g(x_t)$  and  $E[g(x_t)] \in \partial f(x_t)$

### ★ Assumptions for simple analysis

①  $f$  is  $\alpha$ -strongly convex and  $\beta$ -smooth

$$(\alpha \|x-y\| \leq \|\nabla f(x) - \nabla f(y)\| \leq \beta \|x-y\|)$$

$$\textcircled{2} E[\|g_t - \nabla f(x_t)\|^2] = E[\|g_t\|^2] - \|\nabla f(x)\|^2 \leq \sigma^2 \quad \forall t$$

## ★ Convergence rate for Stochastic Gradient Descent

\*  $\alpha$ -strongly convex &  $\beta$ -smooth

- $$\begin{aligned} \mathbb{E}[f(x_t) - f(x_{t+1}) | x_t] &\geq \mathbb{E}\left[\nabla f(x_t)^T (x_t - x_{t+1}) - \frac{\beta}{2} \|x_t - x_{t+1}\|^2 | x_t\right] \quad (\because \text{smooth}) \\ &= \mathbb{E}\left[\gamma \nabla f(x_t)^T g_t - \frac{\beta \gamma^2}{2} \|g_t\|^2 | x_t\right] \\ &= \gamma \|\nabla f(x_t)\|^2 - \frac{\beta \gamma^2}{2} \mathbb{E}[\|g_t\|^2 | x_t] \quad (\because \mathbb{E}[g_t | x_t] = \nabla f(x_t)) \\ &\geq \gamma \left(1 - \frac{\beta \gamma}{2}\right) \|\nabla f(x_t)\|^2 - \frac{\beta \gamma^2}{2} \sigma^2 \quad (\because \text{assumption}) \end{aligned}$$
- $$\begin{aligned} \mathbb{E}[f(x_{t+1}) - f(x^*) | x_t] &\leq \mathbb{E}[f(x_t) - f(x^*) | x_t] - \gamma \left(1 - \frac{\beta \gamma}{2}\right) \|\nabla f(x_t)\|^2 + \frac{\beta \gamma^2}{2} \sigma^2 \\ &\leq f(x_t) - f(x^*) - \gamma \left(1 - \frac{\beta \gamma}{2}\right) \cdot 2\alpha \cdot (f(x_t) - f(x^*)) + \frac{\beta \gamma^2}{2} \sigma^2 \\ &\quad (\because \alpha\text{-strongly convex} \Rightarrow \text{Polyak-Lojasiewicz inequality}) \\ &\leq \left(1 - \alpha \gamma (2 - \beta \gamma)\right) (f(x_t) - f(x^*)) + \frac{\beta \gamma^2}{2} \sigma^2 \end{aligned}$$
- $$\begin{aligned} \mathbb{E}[f(x_T) - f(x^*)] &\leq \left(1 - \alpha \gamma (2 - \beta \gamma)\right)^T (f(x_0) - f(x^*)) + \frac{\beta \gamma^2}{2} \sigma^2 \left(\sum_{t=0}^{T-1} (1 - \alpha \gamma (2 - \beta \gamma))^t\right) \\ &\leq \left(1 - \alpha \gamma (2 - \beta \gamma)\right)^T (f(x_0) - f(x^*)) + \frac{\beta \gamma^2}{2} \sigma^2 \cdot \frac{1 - (1 - \alpha \gamma (2 - \beta \gamma))^T}{\alpha \gamma (2 - \beta \gamma)} \end{aligned}$$

(To boost convergence speed, we can increase batch size ( $\sigma^2 \downarrow$ ) or decrease learning rate ( $\gamma \downarrow$ ))

- When  $\gamma = \frac{1}{\beta}$

$$\begin{aligned} \Rightarrow \mathbb{E}[f(x_T) - f(x^*)] &\leq \left(1 - \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*)) + \frac{\sigma^2}{2} \cdot \frac{1 - (1 - \frac{\alpha}{\beta})^T}{\alpha} \\ &\leq \left(1 - \frac{\alpha}{\beta}\right)^T (f(x_0) - f(x^*)) + \frac{\sigma^2}{2\alpha} \end{aligned}$$

$\therefore O(\log \frac{1}{\varepsilon})$  steps

- When  $\gamma = \frac{1}{\sqrt{T}}$

$$\begin{aligned} \Rightarrow \mathbb{E}[f(x_T) - f(x^*)] &\leq \left(1 - \frac{\alpha}{\sqrt{T}} \left(2 - \frac{\beta}{\sqrt{T}}\right)\right)^T (f(x_0) - f(x^*)) + \frac{\beta}{2} \sigma^2 \frac{1 - (1 - \frac{\alpha}{\sqrt{T}} \left(2 - \frac{\beta}{\sqrt{T}}\right))^T}{\frac{\alpha}{\sqrt{T}} \left(2 - \frac{\beta}{\sqrt{T}}\right)} \\ &\leq \exp(-2\alpha\sqrt{T} + \alpha\beta) (f(x_0) - f(x^*)) + \frac{\beta}{2T} \sigma^2 \cdot \frac{1}{\frac{\alpha}{\sqrt{T}} \left(2 - \frac{\beta}{\sqrt{T}}\right)} \\ &\leq \exp(-2\alpha\sqrt{T} + \alpha\beta) (f(x_0) - f(x^*)) + \frac{\beta}{2} \sigma^2 \cdot \frac{1}{\alpha(2\sqrt{T} - \beta)} \end{aligned}$$

$\therefore O(\frac{1}{\varepsilon^2})$  steps

$\Rightarrow O(\exp(-\sqrt{T}))$

$\Rightarrow O(\frac{1}{\varepsilon})$

$\Rightarrow O((\log \frac{1}{\varepsilon})^2)$  steps

$\Rightarrow O(\frac{1}{\varepsilon^2})$  steps  
(dominant)