

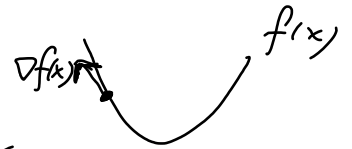
Optimization  $\min_{x \in \mathbb{R}^d} f(x)$  or  $\min_{x \in \mathcal{D} \subset \mathbb{R}^d} f(x)$

$x^* = \operatorname{argmin}_{x \in \mathcal{D} \subset \mathbb{R}^d} f(x)$

"Gradient-based optimization"

$$\checkmark \langle \nabla f(x), \vec{l} \rangle = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon \vec{l}) - f(x)}{\epsilon} \quad (*)$$

$$\checkmark \nabla f(x) = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right)$$



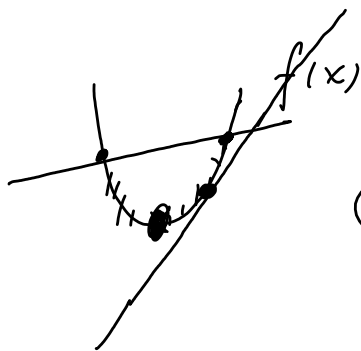
$$x^0 \xrightarrow{\checkmark} x^1 \xrightarrow{\checkmark} x^2 \xrightarrow{\checkmark} \dots \xrightarrow{\checkmark} x^{(N)} \approx x^*$$

$$x^k \xrightarrow{\nabla f(x)} x^{k+1}$$

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

$$\vdots$$

Convex function



Convex optimization

$$\min_{x \in \mathbb{R}^d \cap \mathcal{D}} f(x)$$

Non-Convex optimization



$f$  is convex

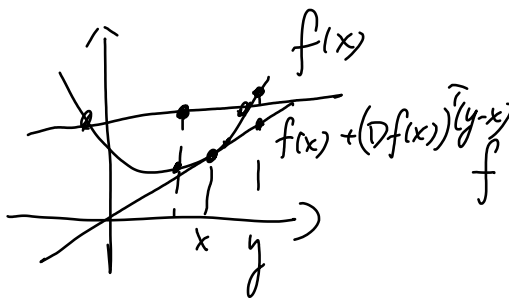
$$\forall \alpha \in [0, 1]$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\pm\infty\}$$

$$f((1-\alpha)x + \alpha y) \stackrel{(\text{convex})}{\leq} (1-\alpha)f(x) + \alpha f(y)$$

"="

say  $f$  is convex and  $\nabla f$  is continuous



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



Strongly Convex Functions  $\alpha \in [0, 1]$   $\underline{m > 0}$  Convexity Constant

$$f((1-\alpha)x + \alpha y) \leq (1-\alpha)f(x) + \alpha f(y) - \frac{1}{2} \underline{m \alpha (1-\alpha)} \underline{\|x-y\|_2^2}$$

$m = 0$  convex functions

Taylor's Theorems  $f: \mathbb{R}^n \rightarrow \mathbb{R} \quad x, p \in \mathbb{R}^n$

$$f(x+p) = f(x) + \int_0^1 \nabla f(x+\gamma p)^T p d\gamma \quad p = y - x$$

$$g(\gamma) = f(x + \gamma p) \quad g(1) - g(0) = \int_0^1 g'(\gamma) d\gamma = g'(\tilde{\gamma})$$

$\gamma \in [0, 1]$

$\gamma p = (\gamma p_1, \dots, \gamma p_n)$

$$= \int_0^1 \nabla f(x + \gamma p)^T p d\gamma$$

$$f(x+p) = f(x) + \nabla f(x + \tilde{\gamma} p)^T p \quad \text{for some } \gamma \in [0, 1]$$

$$\nabla f(x+p) = \nabla f(x) + \int_0^1 \nabla^2 f(x+\gamma p)^T p \, d\gamma$$

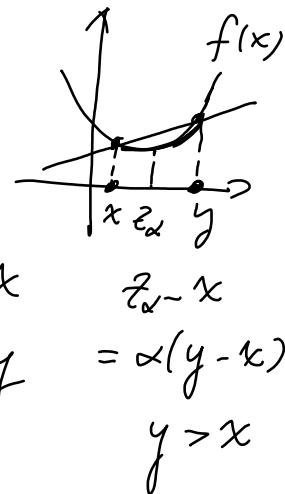
$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

$$\bullet \quad \boxed{f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x+\gamma p) p.}$$

for some  $\gamma \in [0, 1]$

Convex function

- $f((1-\alpha)x + \alpha y) \leq (1-\alpha)f(x) + \alpha f(y)$



$z_\alpha = (1-\alpha)x + \alpha y$

if  $\alpha \rightarrow 0$  then  $z_\alpha \rightarrow x$   
 if  $\alpha \rightarrow 1$  then  $z_\alpha \rightarrow y$

$$f(z_\alpha) \leq (1-\alpha)f(x) + \alpha f(y)$$

$$\Leftrightarrow \alpha f(y) - \alpha f(x) \geq f(z_\alpha) - f(x)$$

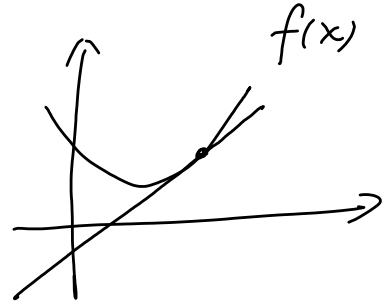
$$\Leftrightarrow \frac{\alpha f(y) - \alpha f(x)}{\alpha(y-x)} \geq \frac{f(z_\alpha) - f(x)}{z_\alpha - x} \quad \alpha \downarrow 0$$

$$\frac{f(y) - f(x)}{y - x} \geq \lim_{z \downarrow 0} \frac{f(x+z) - f(x)}{z} = f'(x)$$

$$f(y) - f(x) \geq f'(x)(y - x)$$

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

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





Strongly convex function

$$f(\underbrace{(1-\alpha)x + \alpha y}_{z_\alpha}) \leq (1-\alpha)f(x) + \alpha f(y) - \frac{m}{2} \alpha(1-\alpha) \|x-y\|_2^2$$

$z_\alpha = (1-\alpha)x + \alpha y \quad z_\alpha - x = \alpha(y-x)$

$$\alpha f(y) - \alpha f(x) \geq f(z_\alpha) - f(x) + \frac{1}{2} m \alpha(1-\alpha) \|x-y\|_2^2$$

$$= (\nabla f(x))^T \alpha(y-x) + \cancel{O(\alpha^2 \|y-x\|^2)}$$

$$\alpha \rightarrow 0 \quad \underline{f(y) - f(x) \geq (\nabla f(x))^T (y-x) + \frac{1}{2} m \|x-y\|_2^2}$$

$$\bullet f(y) \geq f(x) + (\nabla f(x))^T (y-x) + \frac{1}{2} \underline{m} \|x-y\|_2^2$$

$$\nabla^2 f(x) \geq m \mathbf{I}$$

$$\lambda(\nabla^2 f(x)) \geq m$$

$$f(x) \geq f(z) + (\nabla f(z))^T (x-z) + \frac{m}{2} \|x-z\|_2^2 \quad \textcircled{1}$$

$$f(y) \geq f(z) + (\nabla f(z))^T (y-z) + \frac{m}{2} \|y-z\|_2^2 \quad \textcircled{2}$$

$$(1-\alpha) \times \textcircled{1} + \alpha \times \textcircled{2}$$

$$(1-\alpha) f(x) + \alpha f(y) \geq f(z) + \underbrace{(\nabla f(z))^T}_{\text{cancel}} \left[ \underbrace{(1-\alpha)(x-z) + \alpha(y-z)}_{\text{cancel}} \right] + \frac{m}{2} \left[ (1-\alpha) \|x-z\|^2 + \alpha \|y-z\|^2 \right]$$

$$\underline{z} = z_\alpha = \underline{(1-\alpha)x + \alpha y}$$

$$z - x = -\alpha(x-y)$$

$$\begin{aligned} z - y &= (1-\alpha)x + (\alpha-1)y \\ &= (1-\alpha)(x-y) \end{aligned}$$

$$(1-\alpha)(x-z) + \alpha(y-z) = (1-\alpha)\alpha(x-y) + \alpha(1-\alpha)(y-x) = 0$$

$$(1-\alpha) \|x-z\|^2 + \alpha \|y-z\|^2 = \alpha(1-\alpha) \|x-y\|^2$$

$$(1-\alpha)f(x) + \alpha f(y) \geq \underbrace{f(\alpha x + (1-\alpha)y)} + \frac{\mu}{2} \alpha(1-\alpha) \|x-y\|^2$$

$$f(\alpha x + (1-\alpha)y) \leq (1-\alpha)f(x) + \alpha f(y) - \frac{\mu}{2} \alpha(1-\alpha) \|x-y\|^2$$

### Optimization

#### Assumptions

- ①  $\nabla f$  is  $L$ -Lipschitz  
 $\|\nabla f(x) - \nabla f(y)\| \leq L\|x-y\|$
- ②  $f$  is  $\mu$ -strongly convex  
 $(1-\alpha)f(x) + \alpha f(y) - \frac{\mu}{2} \alpha(1-\alpha) \|x-y\|^2 \geq f(\alpha x + (1-\alpha)y)$ .

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

Wednesday, July 22, 2020 9:31 AM

$$L I \geq \nabla^2 f(x) \geq m I$$

$$\lambda_{\max}(\nabla^2 f(x)) \leq L$$

$$\lambda_{\min}(\nabla^2 f(x)) \geq m$$

$$f(y) = f(x) + \int_0^1 \nabla f(x + \gamma(y-x))^T (y-x) d\gamma$$

$$f(y) - f(x) - (\nabla f(x))^T (y-x) = \int_0^1 [\nabla f(x + \gamma(y-x))^T - (\nabla f(x))^T] (y-x) d\gamma$$

$$a^T b \leq \|a\| \cdot \|b\|$$

$$\langle a, b \rangle \leq \|a\| \cdot \|b\|$$

$$\leq \int_0^1 \|\nabla f(x + \gamma(y-x)) - \nabla f(x)\| \cdot \|y-x\| d\gamma$$

$$\leq \int_0^1 L \gamma \|y-x\|^2 d\gamma = \frac{L}{2} \|y-x\|^2$$

$$\nabla f(x) = Qx + b \quad \nabla^2 f(x) = Q$$

$$f(x) = \frac{1}{2} x^T Q x + b^T x + c$$

$$\frac{1}{2} (x_1 \dots x_d) \begin{pmatrix} q_{11} & \dots & q_{1d} \\ \dots & \dots & \dots \\ q_{d1} & \dots & q_{dd} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$$

$$= \frac{1}{2} \sum_{i,j} q_{ij} x_i x_j$$

$$\langle \nabla f(x), \vec{l} \rangle$$

$$= \lim_{\epsilon \rightarrow 0} \frac{\frac{1}{2} (x + \epsilon \vec{l})^T Q (x + \epsilon \vec{l}) + b^T (x + \epsilon \vec{l}) + c - (\frac{1}{2} x^T Q x + b^T x + c)}{\epsilon}$$



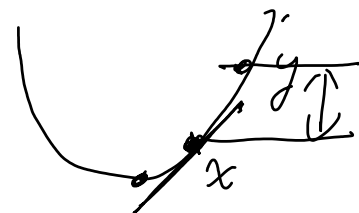
$$= \lim_{\epsilon \rightarrow 0} \frac{\frac{1}{2} \epsilon \vec{l}^T Q x + \frac{1}{2} x^T Q \vec{l} + \epsilon b^T \vec{l}}{\epsilon} = \langle Qx + b, \vec{l} \rangle$$

$\nabla^2 f \gg 1$        $\nabla^2 f < 1$

If  $f$  is strongly convex with convexity constant  $\underline{m} > 0$

$$f(y) - f(x) \geq -\frac{1}{2\underline{m}} \|\nabla f(x)\|^2$$

$$\|\nabla f(x)\|^2 \geq 2\underline{m} [f(x) - f(y)]$$



$$f(x^k) \rightarrow f(x^*) \iff \nabla f(x^k) \rightarrow 0$$

$$\|\nabla f(x^k)\|^2 \geq 2\underline{m} [f(x^k) - f(x^*)]$$

$$\begin{aligned}
 f(y) &\geq f(x) + (\nabla f(x))^T (y-x) + \frac{\mu}{2} \|y-x\|^2 \\
 &= f(x) + \frac{\mu}{2} \left[ \|y-x\|^2 + \frac{2}{\mu} (\nabla f(x))^T (y-x) \right] \\
 &= f(x) + \frac{\mu}{2} \left[ \left\| y-x + \frac{1}{\mu} \nabla f(x) \right\|^2 - \frac{1}{\mu^2} \|\nabla f(x)\|^2 \right] \\
 &\geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2
 \end{aligned}$$



"first-order algorithms"

GD

AGD

Optimization Algorithm

$$\min_{x \in \mathbb{R}^n} f(x) = f(x^*)$$

$$x^0 \rightarrow x^1 \rightarrow x^2 \rightarrow \dots \rightarrow x^k \rightarrow \dots$$

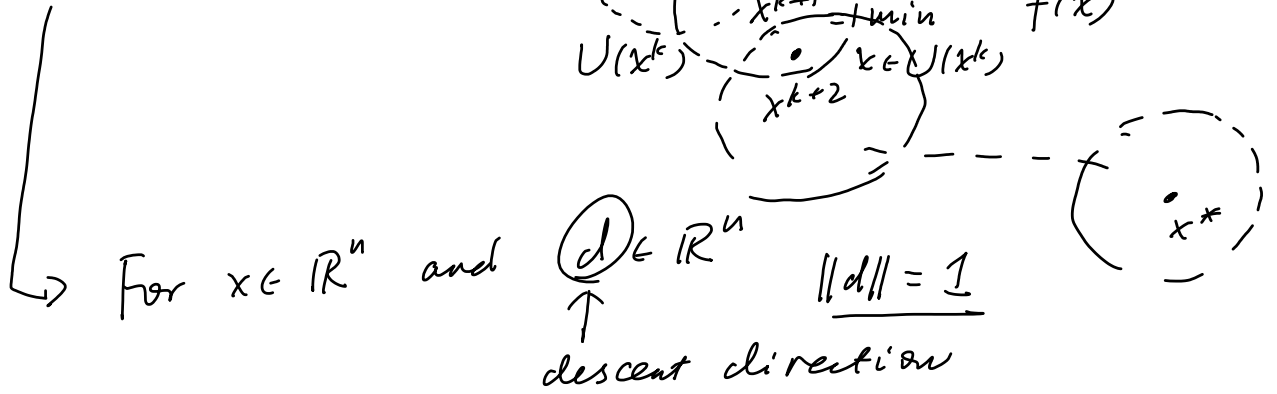
$$x^k \rightarrow x^* \quad \text{as } k \rightarrow \infty$$

$$x^k \xrightarrow{?} \underline{x^{k+1}}$$

$$f(x^k) + \nabla f(x^k)^T (x - x^k) + \frac{1}{2} (x - x^k)^T \nabla^2 f(x^k) (x - x^k)$$

Local neighborhood  
"trust region"

## Line Search Methods



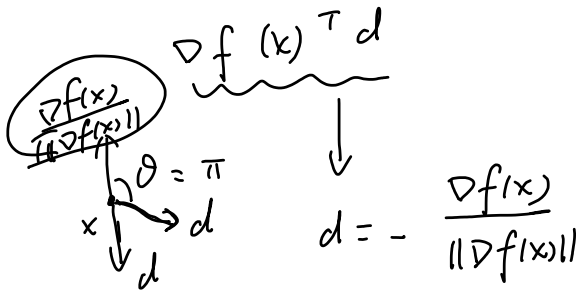
$$f(x + t d) < f(x) \text{ for some } t \text{ small}$$

$$f(x + td) = f(x) + t \frac{\nabla f(x + \gamma td)^T d}{\text{for some } \gamma \in (0, 1)}$$

$$< f(x)$$

$$\nabla f(x + \gamma td)^T d < 0$$

ss



$$\inf_{\|d\|=1} (\nabla f(x))^T d$$

$$= \inf_{\|d\|=1} \|\nabla f(x)\| \cdot \left\langle \frac{\nabla f(x)}{\|\nabla f(x)\|}, d \right\rangle$$

$$= \|\nabla f(x)\| \inf_{\|d\|=1} \left\langle \frac{\nabla f(x)}{\|\nabla f(x)\|}, d \right\rangle$$

$$x \rightarrow x + t \cdot \left( \frac{-\nabla f(x)}{\|\nabla f(x)\|} \right)$$
$$= x - \underbrace{\left( \frac{t}{\|\nabla f(x)\|} \right)}_{\alpha(x)} \nabla f(x)$$

$$x^{k+1} = x^k - \underbrace{\alpha_k}_{\text{learning rate}} \nabla f(x^k)$$

Gradient Descent

# Proof of Convergence

Assume  $f$  is  $m$ -strongly convex  $\nabla f$  is  $L$ -Lipschitz

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

$$m \leq \lambda(\nabla^2 f(x)) \leq L$$

## ① Spectral Method

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

$$= \Phi_\alpha(x^k)$$



$$\Phi_\alpha(x) = x - \alpha \nabla f(x) \rightsquigarrow \begin{matrix} x = \phi(x) \\ \Leftrightarrow \nabla f(x) = 0 \end{matrix}$$

$$x^{k+1} = \Phi_{\alpha_k}(x^k) \quad x^k \rightarrow \underline{x^*} \quad \underbrace{x^* = \Phi_{\alpha}(x^*)}_{\text{fixed pt.}}$$

$$\begin{aligned} x^k &= \Phi_{\alpha}(x^{k-1}) = \Phi_{\alpha} \circ \Phi_{\alpha}(x^{k-2}) \\ &= \dots = \underbrace{(\Phi_{\alpha} \circ \Phi_{\alpha} \circ \dots \circ \Phi_{\alpha})}_k(x^0) \end{aligned}$$

"contraction mapping theorem"

$$\boxed{\|\Phi(x) - \Phi(y)\| \leq \beta \|x - y\|} \quad \underline{0 < \beta < 1}$$

$$\begin{aligned} \|x^{k+1} - x^*\| &= \|\Phi(x^k) - \Phi(x^*)\| \leq \beta \|x^k - x^*\| \\ x^{k+1} &\rightarrow x^* \quad \leq \beta^2 \|x^{k-1} - x^*\| \leq \dots \leq \beta^k \|x^0 - x^*\| \\ &\rightarrow 0 \quad (k \rightarrow \infty) \end{aligned}$$

$$\begin{aligned}
 \|\phi(x) - \phi(y)\| &= \|x - \alpha \nabla f(x) - (y - \alpha \nabla f(y))\| \\
 &= \|(x-y) - \alpha (\nabla f(x) - \nabla f(y))\| \\
 &= \left\| (x-y) - \alpha \int_0^1 \underbrace{\nabla^2 f(y + \gamma(x-y))}_{\in [m, L]} \underbrace{(x-y)}_{\gamma} d\gamma \right\| \\
 &= \left\| \underbrace{\left( I - \alpha \int_0^1 \nabla^2 f(y + \gamma(x-y)) d\gamma \right)}_{\substack{1 - \alpha L \\ \leq \lambda(A(x, y)) \leq 1 - \alpha m}} \right\| \underline{\underline{(x-y)}} \left\| \right.
 \end{aligned}$$

$$\begin{aligned} \|\phi(x) - \phi(y)\| &= \|A(x, y)(x-y)\| \\ &\leq \beta \|x-y\| \quad 0 < \beta < 1 \end{aligned}$$

$|\lambda(A)| \leq \beta$  "spectral method"

$$-\beta \leq 1 - \alpha L \leq \lambda(A) \leq 1 - \alpha m \leq \beta$$

$\alpha \leq \frac{1+\beta}{L}$ 
 $\alpha \geq \frac{1-\beta}{m}$

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

$$\alpha \in \left[ \frac{1-\beta}{m}, \frac{1+\beta}{L} \right]$$

$$\frac{1-\beta}{m} = \frac{1+\beta}{L} \Leftrightarrow \beta = \frac{L-m}{L+m}$$

$$\alpha = \frac{1 - \frac{L-m}{L+m}}{m} = \frac{2}{L+m}$$



② Taylor Expansion Method

"Generally"

$$f(x + \alpha d) = f(x) + \alpha (\nabla f(x))^T d + \alpha \int_0^1 \frac{\langle \nabla f(x + \gamma \alpha d) - \nabla f(x), d \rangle}{\|d\|} d\gamma$$

$$\leq f(x) + \alpha (\nabla f(x))^T d + \alpha \int_0^1 \|\nabla f(x + \gamma \alpha d) - \nabla f(x)\| \cdot \|d\| d\gamma$$

" $\nabla f$  is  $L$ -Lipschitz"

$$\leq f(x) + \alpha (\nabla f(x))^T d + \alpha \int_0^1 L \gamma \alpha \|d\|^2 d\gamma$$

$$= f(x) + \alpha (\nabla f(x))^T d + \frac{L \alpha^2 \|d\|^2}{2}$$

"if  $\alpha = \frac{1}{L}$  then  $\alpha - \frac{\alpha^2 L}{2} = \frac{1}{L} - \frac{\frac{1}{L^2} L}{2} = \frac{1}{2L}$ "

Monday, July 27, 2020 9:19 AM

$$f(x + \alpha d) \leq f(x) + \underbrace{\alpha (\nabla f(x))^T d + \frac{\alpha^2 L}{2} \|d\|^2}_{< f(x)}$$

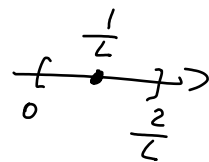
GD  $x^{k+1} = x^k - \alpha \nabla f(x^k) \quad d = -\nabla f(x)$

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

$$= f(x^k) - \underbrace{\left(\alpha - \frac{\alpha^2 L}{2}\right)}_{\alpha \in (0, \frac{2}{L})} \|\nabla f(x^k)\|^2$$

$$0 < \alpha \left[1 - \frac{\alpha L}{2}\right] \Leftrightarrow \alpha < \frac{2}{L}$$

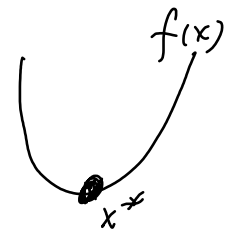
$$\alpha \in (0, \frac{2}{L}) \implies f(x^{k+1}) < f(x^k)$$



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

$$x^{k+1} = x^k - \frac{1}{L} \nabla f(x^k)$$

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



If  $f$  is  $m$ -strongly convex  $\|\nabla f(x^k)\|^2 \stackrel{\textcircled{2}}{\geq} 2m [f(x^k) - f(x^*)]$

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

$$\stackrel{\textcircled{2}}{\leq} f(x^k) - \frac{1}{L} m [f(x^k) - f(x^*)]$$

$$f(x^{k+1}) - f(x^*) \leq \left(1 - \frac{m}{L}\right) f(x^k) + \frac{m}{L} f(x^*) = \left(1 - \frac{m}{L}\right) (f(x^k) - f(x^*))$$

$$f(x^{k+1}) - f(x^*) \leq \left(1 - \frac{m}{L}\right) (f(x^k) - f(x^*))$$

$$\leq \dots$$

$$\leq \underbrace{\left(1 - \frac{m}{L}\right)^k}_{\rightarrow 0} [f(x^0) - f(x^*)]$$

$$\rightarrow 0 \quad f(x^k) \rightarrow f(x^*)$$

$$1 - \frac{m}{L} \in (0, 1)$$

$$\frac{m}{L} \approx 0.5 \quad 1 - \frac{m}{L} \approx 0.5$$

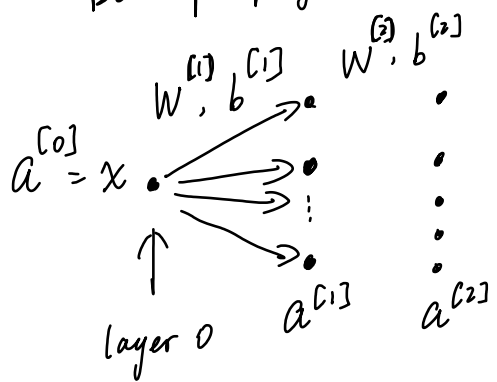
$$\frac{m}{L} = 0.0000001 \quad 1 - \frac{m}{L} \approx 1$$

$k \gg 1$

$$mI \leq D^2 f(x) \leq LI \quad \kappa(D^2 f) = \frac{L}{m} \quad \text{"ill conditioned"}$$

"Neural Networks"

Backpropagation,



Loss function

$$C(W, b) = \frac{1}{2} (y - a^{[L+1]})^2$$

$$W = (W^{[1]} \dots W^{[L+1]}) \quad b = (b^{[1]} \dots b^{[L+1]})$$

$$a^{[L+1]} = g(x, y; W, b)$$

$$a^{[0]} = x$$

$$\begin{cases} z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} = \sigma(z^{[l]}) \end{cases} \quad l = 1, 2, \dots, L+1$$

$$l = 1, 2, \dots, L+1$$

GD on NN  $(W^{(k+1)}, b^{(k+1)}) = (W^{(k)}, b^{(k)}) - \alpha \nabla C(\underline{W}^{(k)}, \underline{b}^{(k)})$

$$C = C(W, b) = \frac{1}{2} (y - a^{[L+1]})^2$$

$$\frac{\partial C}{\partial W_{ij}^{[k]}}$$

$$\frac{\partial C}{\partial b_i^{[k]}}$$

"Backpropagation"

$$\left\{ \begin{aligned} z^{[l+1]} &= W^{[l]} a^{[l]} + b^{[l]} \\ a^{[l+1]} &= \sigma(z^{[l+1]}) \end{aligned} \right\} z^{[l+1]} = W^{[l]} \sigma(z^{[l]}) + b^{[l]}$$

$$C = \frac{1}{2} (a^{[L+1]} - y)^2$$

$$= f_1 \circ f_2 \circ \dots \circ f_n$$

$$a^{[L+1]} = a^{[L+1]}(W, b)$$

$$= \sigma(W^{[L]} a^{[L]} + b^{[L]})$$

$$= \sigma(W^{[L]} (\sigma(W^{[L-1]} a^{[L-1]} + b^{[L-1]})))$$

$$\dots$$

"error"

$$\delta^{[l]} = \frac{\partial C}{\partial z^{[l]}}$$

$$\frac{\partial C}{\partial a^{[L+1]}} = a^{[L+1]} - y$$

$$\frac{\partial C}{\partial z^{[L+1]}} = \frac{\partial C}{\partial a^{[L+1]}} \frac{\partial a^{[L+1]}}{\partial z^{[L+1]}} = (a^{[L+1]} - y) \cdot \sigma'(z^{[L+1]})$$

$$\frac{\partial C}{\partial W^{[L]}} = \frac{\partial C}{\partial z^{[L+1]}} \frac{\partial z^{[L+1]}}{\partial W^{[L]}} = (a^{[L+1]} - y) \sigma'(z^{[L+1]}) a^{[L]}$$

$$\frac{\partial C}{\partial b^{[L]}} = \frac{\partial C}{\partial z^{[L+1]}} \frac{\partial z^{[L+1]}}{\partial b^{[L]}} = (a^{[L+1]} - y) \sigma'(z^{[L+1]})$$

$$\frac{\partial C}{\partial z^{[L]}} = \frac{\partial C}{\partial z^{[L+1]}} \frac{\partial z^{[L+1]}}{\partial z^{[L]}} = \frac{\partial C}{\partial z^{[L+1]}} \frac{\partial z^{[L+1]}}{\partial a^{[L]}} \frac{\partial a^{[L]}}{\partial z^{[L]}}$$

$$\delta^{[L]} = \delta^{[L+1]} W^{[L]} \sigma'(z^{[L]})$$

"Back propagation"



$$f^{[L+1]} \rightarrow f^{[L]} \rightarrow f^{[L-1]} \rightarrow \dots \rightarrow f^{[0]}$$

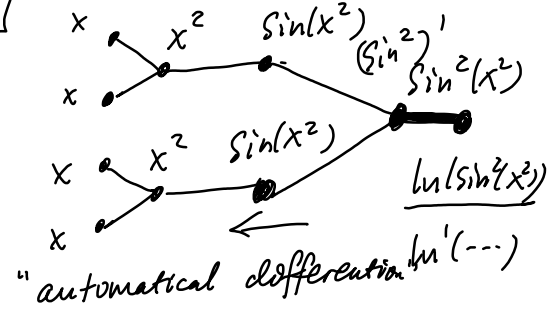
$$\frac{\partial C}{\partial W^{[l]}} = \frac{\partial C}{\partial z^{[l+1]}} \frac{\partial z^{[l+1]}}{\partial W^{[l]}} = \delta^{[l+1]} a^{[l]}$$

$$\frac{\partial C}{\partial b^{[l]}} = \frac{\partial C}{\partial z^{[l+1]}} \frac{\partial z^{[l+1]}}{\partial b^{[l]}} = \delta^{[l+1]}$$

"computational graph"

$$f(x) = \ln(\sin^2(x^2))$$

$$f'(x) = \ln'(\sin^2)' \sin'$$



# Accelerated Gradient Descent

$f$  is  $m$ -strongly convex     $\nabla f$  is  $L$ -Lipschitz     $\Leftrightarrow mI \leq \nabla^2 f \leq LI$

$$Q \geq 0 \quad f(x) = \frac{1}{2} x^T Q x - b^T x + c$$

$$\nabla^2 f(x) = Q$$

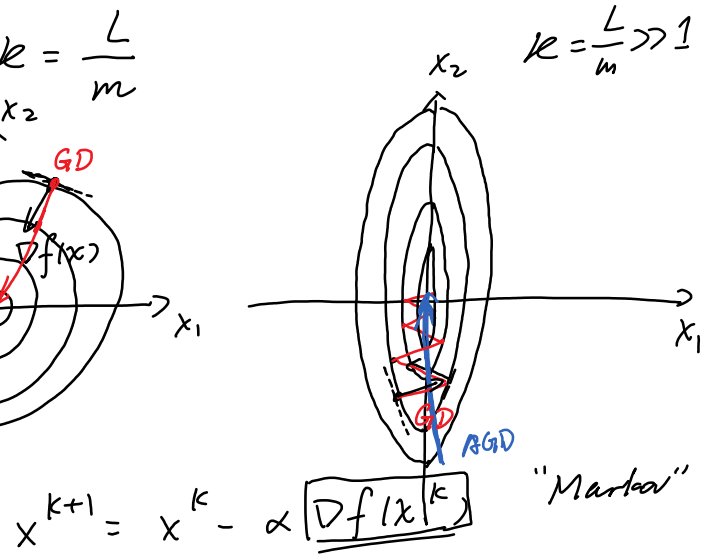
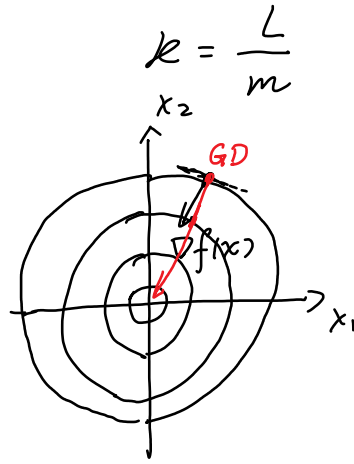
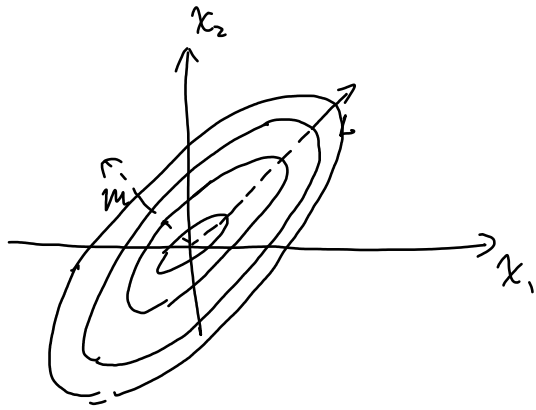
$$\lambda_{\min}(Q) = m$$

$$\lambda_{\max}(Q) = L$$

$$x = (x_1, x_2)^T$$

$$f(x) = \frac{1}{2} x^T Q x = \frac{1}{2} [q_{11} x_1^2 + q_{22} x_2^2 + 2 q_{12} x_1 x_2]$$

"level set"



"inertia"  $\leftrightarrow$  "acceleration"

$$F = ma = m \frac{d^2 x}{dt^2}$$

GD  
 $x^0$

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

$$\frac{x^{k+1} - x^k}{\Delta t} = - \nabla f(x^k)$$

↓

$$\frac{dx}{dt} = - \nabla f(x)$$

Gradient Flow

AGD  $\mu$  — Heavy Ball

$$\mu \frac{d^2 x}{dt^2} = - \nabla f(x) - \mu b \frac{dx}{dt}$$

$$\mu \frac{x(t+\Delta t) - 2x(t) + x(t-\Delta t)}{(\Delta t)^2}$$

$$= - \nabla f(x) - \mu b \frac{x(t+\Delta t) - x(t)}{\Delta t}$$

$$\mu \frac{x(t+\Delta t) - 2x(t) + x(t-\Delta t))}{(\Delta t)^2} = -\nabla f(x)^{(t)} - \mu b \frac{x(t+\Delta t) - x(t)}{\Delta t}$$

$$\rightsquigarrow x(t+\Delta t) = \dots x(t) + \dots x(t-\Delta t) \quad (\star)$$

$$\mu \left[ \frac{1}{(\Delta t)^2} + \frac{b}{\Delta t} \right] x(t+\Delta t) = -\nabla f(x) + \left[ \frac{2\mu}{(\Delta t)^2} + \frac{\mu b}{\Delta t} \right] x(t)$$

$$(1 + b\Delta t)\mu x(t+\Delta t) = -(\Delta t)^2 \nabla f(x) + \mu \left[ 2 + b\Delta t \right] x(t) - \frac{\mu}{(\Delta t)^2} x(t-\Delta t)$$

$$x(t+\Delta t) = - \frac{(\Delta t)^2}{(1+b\Delta t)\mu} \nabla f(x(t)) + \frac{1+b\Delta t+1}{1+b\Delta t} x(t) - \frac{1}{1+b\Delta t} x(t-\Delta t)$$

$$x(t+\Delta t) = -\alpha \nabla f(x(t)) + x(t) + \frac{\beta(x(t) - x(t-\Delta t))}{(AGD)}$$

$$x(t+\Delta t) = -\alpha \nabla f(x(t)) + x(t) \quad (GD)$$

# Polyak's Heavy Ball method (GD with momentum)

$$x^{(k+1)} = x^{(k)} - \alpha \nabla f(x^{(k)}) + \beta (x^{(k)} - x^{(k-1)}), \quad x^{(0)} = x^{(-1)}$$

Annotations:  
-  $\alpha$ : learning rate  
-  $\beta$ : momentum constant  
-  $(x^{(k)} - x^{(k-1)})$ : momentum

$x$

$$m \frac{dx}{dt}$$

"momentum"  $p^k = x^{k+1} - x^k$

$$\begin{cases} p^k = -\alpha \nabla f(x^k) + \beta p^{k-1} \\ x^{k+1} = x^k + p^k \end{cases}$$

GD }  $x^{k+1} = x^k + p^k$  }  $\begin{cases} \text{GD} \rightarrow p^k = -\alpha \nabla f(x^k) \\ \text{AGD} \rightarrow p^k = -\alpha \nabla f(x^k) + \beta p^{k-1} \\ \phantom{\text{AGD}} = -\alpha \nabla f(x^k) - \beta \alpha \nabla f(x^{k-1}) + \beta^2 p^{k-2} \\ \phantom{\text{AGD}} = -\alpha \nabla f(x^k) - \beta \alpha \nabla f(x^{k-1}) - \beta^2 \alpha \nabla f(x^{k-2}) \\ \phantom{\text{AGD}} \phantom{=} + \beta^3 p^{k-3} \\ \phantom{\text{AGD}} \phantom{=} \phantom{=} = -\alpha \nabla f(x^k) - \beta \alpha \nabla f(x^{k-1}) - \beta^2 \alpha \nabla f(x^{k-2}) - \dots \end{cases}$

$\beta^m \alpha \nabla f(x^{k-m})$   
 "moving-average"



# Polyak's Heavy Ball

$$x^{k+1} = x^k - \alpha \nabla f(x^k) + \beta(x^k - x^{k-1})$$
"Explicit Scheme"

$$x^{k+1} = x^k - \alpha \nabla f(x^k + \beta(x^k - x^{k-1})) + \beta(x^k - x^{k-1})$$
 $x^{k+1} = F(x^k, x^{k-1}, \dots)$

$$x^{k+1} = x^k - \alpha \nabla f(x^{k+1}) + \beta(x^k - x^{k-1})$$
"Implicit Scheme"

$$x^k - \alpha \nabla f(x^k + \beta(x^k - x^{k-1})) = \alpha \nabla f(x^{k+1}) + \beta(x^k - x^{k-1})$$
 $x^{k+1} = F(x^{k+1}, x^k, x^{k-1}, \dots)$

← Nesterov's Accelerated GD

Poljak  $x^{k+1} = x^k - \alpha \nabla f(x^k) + \beta(x^k - x^{k-1})$

Nesterov  $x^{k+1} = x^k - \alpha \nabla f(x^k + \beta(x^k - x^{k-1})) + \beta(x^k - x^{k-1})$

Meta-Theorem  $\alpha, \beta$  tuned appropriately

$f$  is  $m$ -strongly convex  $\nabla f$  is  $L$ -Lipschitz

$$f(x^{k+1}) - f(x^*) \leq \left(1 - \sqrt{\frac{m}{L}}\right)^k [f(x^0) - f(x^*)] + \frac{m}{2} \|x^0 - x^*\|^2$$

	$k = \frac{L}{m}$
GD	$1 - \frac{1}{k}$
AGD	$1 - \frac{1}{\sqrt{k}}$

$$k = 10^8$$
$$1 - 10^{-8} \approx 1$$
$$1 - 10^{-4} \approx 1 - 10^{-8}$$

# Gradient-Based Methods

GD

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

AGD

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

ERM

$(a_i, b_i)$   
↑  
input    label

$$l_i(x) = l(g(a_i; x), b_i) \stackrel{eg}{=} |g(a_i; x) - b_i|^2$$

$$L_n(x) = \frac{1}{n} \sum_{i=1}^n l_i(x)$$

Empirical Risk

$$x_n^* = \underset{x}{\operatorname{argmin}} L_n(x)$$

# Gradient-Based Optimization applied to ERM

$$\nabla L_n(x) = \frac{1}{n} \sum_{i=1}^n \nabla l_i(x)$$

GD  $x^{k+1} = x^k - \frac{\alpha}{n} \sum_{i=1}^n \nabla l_i(x^k) \quad k=1 \dots T \gg 1$

Large-Scale ML problems  $x \in \mathbb{R}^D \quad D \gg 1$   
training data  $(a_i, b_i) \quad i=1 \dots n \quad \underline{n \gg 1}$

IFO (Incremental First-Order Oracle)  $\text{IFO}[\text{GD}] \rightarrow \underline{n \times T}$

# Stochastic Gradient Descent (SGD)

GD

$$x^{k+1} = x^k - \frac{\alpha}{n} \sum_{i=1}^n \nabla l_i(x^k) \quad B=n$$

minibatch-SGD

minibatch  $\boxed{1 \leq B \leq n}$  take randomly uniformly size  $B$  subsets from  $\{1, \dots, n\}$

full-batch  $\{1, 2, 3, \dots, n-1, n\}$

e.g.  $n=3$   $\{1, 2, 3\}$

$B=2$   $\frac{\{1, 2\}}{\uparrow p=1/3}$   $\frac{\{1, 3\}}{\uparrow p=1/3}$   $\frac{\{2, 3\}}{\uparrow p=1/3}$

$$\underline{B=1} \quad \{1 \dots n\}$$

mini-batch  $i$  taken randomly uniformly from  $\{1 \dots n\}$

$B \subset \{1 \dots n\}$  s.t.  $|B| = B$  and  $B$  taken randomly uniformly

$B_1, B_2, \dots, B_K, \dots$  i.i.d

$$\underline{B \ll n}$$

mini-batch SGD  $\frac{1}{B} \sum_{i \in \mathcal{B}} \nabla l_i(x)$  "gradient estimator"

$$x^{k+1} = x^k - \frac{\alpha}{B} \sum_{i \in \mathcal{B}_k} \nabla l_i(x^k) = x^k - \alpha g_{\mathcal{B}_k}(x^k)$$

B = 1

$$x^{k+1} = x^k - \alpha \nabla l_{i_k}(x^k)$$

$$L_n(x) = \frac{1}{n} \sum_{i=1}^n l_i(x)$$

$i_1, i_2, \dots, i_k, \dots$  is drawn randomly uniformly from  $\{1, \dots, n\}$  in an i.i.d way

$$\mathbb{P}(i=k) = \frac{1}{n} \quad k=1 \dots n$$

$$L_n(x) = \mathbb{E} l_i(x)$$



SGD

$n=3$

$$\begin{aligned} x^1 &= x^0 - \alpha \nabla \ell_2(x^0) \\ x^2 &= x^1 - \alpha \nabla \ell_1(x^1) \\ x^3 &= x^2 - \alpha \nabla \ell_3(x^2) \\ &\vdots \end{aligned}$$

GD

$$\begin{aligned} x^1 &= x^0 - \alpha \frac{\nabla \ell_1(x^0) + \nabla \ell_2(x^0) + \nabla \ell_3(x^0)}{3} \\ x^2 &= x^1 - \alpha \frac{\nabla \ell_1(x^1) + \nabla \ell_2(x^1) + \nabla \ell_3(x^1)}{3} \\ x^3 &= x^2 - \alpha \frac{\nabla \ell_1(x^2) + \nabla \ell_2(x^2) + \nabla \ell_3(x^2)}{3} \\ &\vdots \end{aligned}$$

In general if  $\mathcal{B}$  is taken randomly uniformly from size  $B$  subsets of  $\{1, \dots, n\}$

$$\binom{n}{B} = \frac{n!}{B!(n-B)!}$$

then  $\frac{1}{n} \sum_{i=1}^n \nabla \ell_i(x) = \mathbb{E}_{\mathcal{B}} \left( \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla \ell_i(x) \right)$

"unbiasedness"  $\nabla L_n(x) = \mathbb{E}_{\mathcal{B}} (g_{\mathcal{B}}(x))$   $g_{\mathcal{B}} = \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla \ell_i(x)$

GD

$$x^{k+1} = x^k - \alpha \nabla L_n(x^k)$$

SGD

$$x^{k+1} = x^k - \alpha \underbrace{g_{B_k}(x^k)}$$

$$\mathbb{E} g_{B_k}(x^k) = \nabla L_n(x^k)$$

$$\mathbb{E} [g_{B_k}(x^k) - \mathbb{E} g_{B_k}(x^k)]^2 \uparrow \text{ as } B \downarrow$$



# "Convergence"

Monday, August 3, 2020 9:13 AM

SGD  $x^{k+1} = x^k - \alpha \underset{\text{Gradient Estimator}}{g_{\zeta_k}(x^k)}$

$$\mathbb{E} g_{\zeta_k}(x^k) = \nabla f(x^k)$$

"unbiased"

$\zeta_1, \dots, \zeta_k, \dots$  is i.i.d

$f$  is  $m$ -strongly convex

$x^*$

$\mathbb{E} = \mathbb{E}^k$

$g_{\zeta_k}$

$$\|x^{k+1} - x^*\|^2 = \|x^k - \alpha g_{\zeta_k}(x^k) - x^*\|^2$$

$$= \|x^k - x^*\|^2 - 2\alpha \langle x^k - x^*, g_{\zeta_k}(x^k) \rangle + \alpha^2 \|g_{\zeta_k}(x^k)\|^2$$

$$\mathbb{E} \left( \|x^{k+1} - x^*\|^2 \middle| x^k \right) = \|x^k - x^*\|^2 - 2\alpha \langle x^k - x^*, \nabla f(x^k) \rangle + \alpha^2 \mathbb{E} \left( \|g_{\zeta_k}(x^k)\|^2 \right)$$

$$\langle \nabla f(x), x - x^* \rangle \geq m \|x - x^*\|^2 \quad (\text{Ex.})$$

Hint  $\nabla^2 f \geq mI$ ,

$$\mathbb{E} \left[ \|x^{k+1} - x^*\|^2 \mid x^k \right] \leq \|x^k - x^*\|^2 - 2\alpha m \|x^k - x^*\|^2 + \alpha^2 \mathbb{E} \|g_\zeta(x^k)\|^2$$

$$\mathbb{E} \left[ \mathbb{E}[A \mid B] \right] = \mathbb{E}A \quad (\text{"telescoping"})$$

$$\mathbb{E} \|x^{k+1} - x^*\|^2 \leq (1 - 2\alpha m)^k \mathbb{E} \|x^0 - x^*\|^2 + 2\alpha^2 \left[ \text{Var}(g_\zeta) + \mathbb{E} \|\nabla f\|^2 \right]$$

$$\mathbb{E} g_\zeta = \nabla f \quad \mathbb{E} \|g_\zeta\|^2 = \mathbb{E} \|(g_\zeta - \nabla f) + \nabla f\|^2 \leq 2 \text{Var}(g_\zeta) + 2 \mathbb{E} \|\nabla f\|^2$$

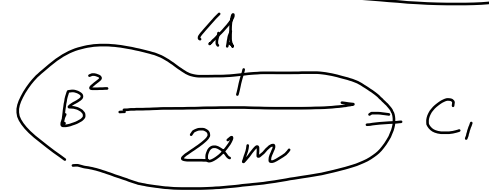
$$\begin{aligned} \mathbb{E} \|x^{k+1} - x^*\|^2 &\leq (1 - 2\alpha m) \mathbb{E} \|x^k - x^*\|^2 + B^2 \\ &\leq (1 - 2\alpha m)^2 \mathbb{E} \|x^{k-1} - x^*\|^2 + (1 - 2\alpha m)B^2 + B^2 \\ &\leq \dots \end{aligned}$$

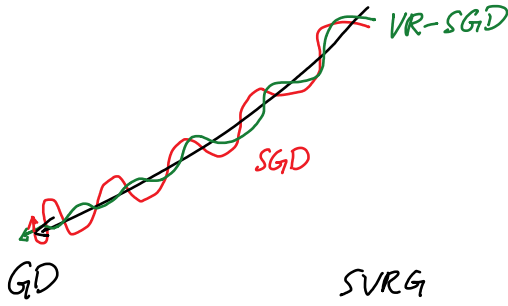
$$\leq (1 - 2\alpha m)^k \mathbb{E} \|x^1 - x^*\|^2 + \underbrace{(1 - 2\alpha m)^{k-1} B^2 + \dots + (1 - 2\alpha m) B^2 + B^2}_{\leq C_1}$$

$$\mathbb{E} \|x^k - x^*\|^2 \leq \lambda^k \cdot C_0 + C_1$$

$$0 < \lambda < 1$$

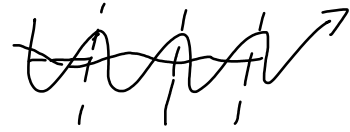
$C_1^k \rightarrow 0$   
(VR) as  $k \rightarrow \infty$





• Variance-Reduced SGD

GD SVRG  
 2014? Johnson-Zhang

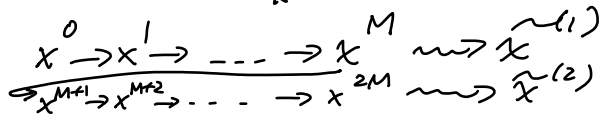


SGD

$$x^{k+1} = x^k - \alpha \nabla l_{i_k}(x^k)$$

VR-SGD

"Snapshots"



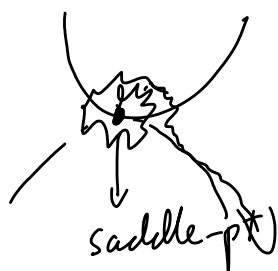
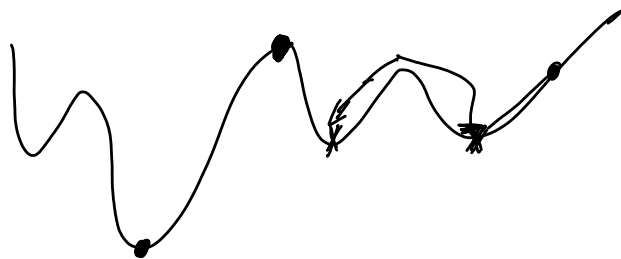
$$x^{k+1} = x^k - \alpha \left( \nabla l_{i_k}(x^k) - \nabla l_{i_k}(\tilde{x}) + \frac{1}{n} \sum_{i=1}^n \nabla l_i(\tilde{x}) \right)$$

VR gradient estimator

- Non-Convex Optimization

"Implicit Regularization"

SGD



"escape from saddles"

Practical Issues

GD  
gradient

Adaptive Methods

Momentum

Stochastic Gradients

$\nabla \rightarrow \hat{\nabla}$