

①

Recap

$$\min_x f(x) \text{ such that } x \in C$$

→ GD + Projection  $x^{(k)} = P_C(x^{(k-1)} - t_k g^{(k-1)})$

Let  $I_C(x) = \begin{cases} \infty & x \notin C \\ 0 & x \in C \end{cases}$  then

$$\min_x f(x) + I_C(x) \text{ is equivalent}$$

Proximal Gradient Descent

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

↙      ↘

convex & differentiable      convex but non-smooth

e.g.  $\min_x f(x) \text{ such that } \|x\|_1 \leq \varepsilon$   
 $(\|x\|_1 - \varepsilon \leq 0)$

Lagrange Mult. →  $\min f(x) + \lambda (\|x\|_1 - \varepsilon)$

$$\equiv \min_{g(x)} f(x) + \lambda \|x\|_1$$

↙      ↘

g(x)      h(x)

Want to  $\min_z g(z) + h(z)$

(2)

$$g(z) \approx g(x) + \nabla g(x)^T (z - x) + \frac{1}{2t} \|z - x\|_2^2$$

$$(\lesssim) \quad (\nabla^2 g(x) \leq L \mathbf{I} \text{ & } t \leq \frac{1}{L})$$

$$\Rightarrow \min_z g(z) + h(z)$$

$$\cong \min_z g(x) + \nabla g(x)^T (z - x) + \frac{1}{2t} \|z - x\|_2^2 + h(z)$$

$$\cong \min_z \frac{1}{2t} \|z - (x - t \nabla g(x))\|_2^2 + h(z)$$

↖ close to gradient  
at  $g$

} minimize  $h$

Definition

Proximal Operator

$$\text{prox}_{h,t}(x) = \underset{z}{\operatorname{argmin}} \|x - z\|_2^2 + h(z)$$

Proximal Update:  $x^{(k)} = \text{prox}_{h,t_k}(x^{(k-1)} - t_k \nabla g(x^{(k-1)}))$

$$\text{Let } G_t(x) = \frac{x - \text{prox}_{h,t}(x - t \nabla g(x))}{t} \quad (3)$$

then  $x^{(k)} = x^{(k-1)} - t_k G_{t_k}(x^{(k-1)})$   
 is proximal update

Going back to  $\min_x f(x) \text{ s.t. } x \in C$

$$\equiv \min_x f(x) + I_C(x)$$

$$\text{prox}_{I_C, t} = \arg \min_z \|x - z\|_2^2 + I_C(z)$$

↳ which is projection to  $C$

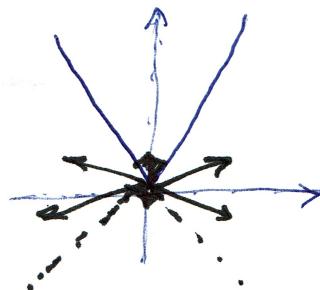
Example: Iterative Soft Thresholding Algorithm

LASSO:  $\min_{\beta} \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{g(\beta)} + \underbrace{\lambda \|\beta\|_1}_{h(\beta)}$

(4)

$$\text{prox}_{h,+}(\beta) = \underset{z}{\operatorname{argmin}} \frac{1}{2t} \|x - z\|_2^2 + \lambda \|z\|_1$$

$$\partial \|z\|_1 = \begin{cases} \text{sgn}(z) & z \neq 0 \\ \alpha \in [-1, 1] & z = 0 \end{cases}$$



Rewrite :  $\underset{z}{\operatorname{argmin}} \frac{1}{2} \|x - z\|_2^2 + \lambda t \|\bar{z}\|_1$

$$0 = \frac{\partial}{\partial z} \Rightarrow 0 = z - x + \lambda t \partial \|z\|_1$$

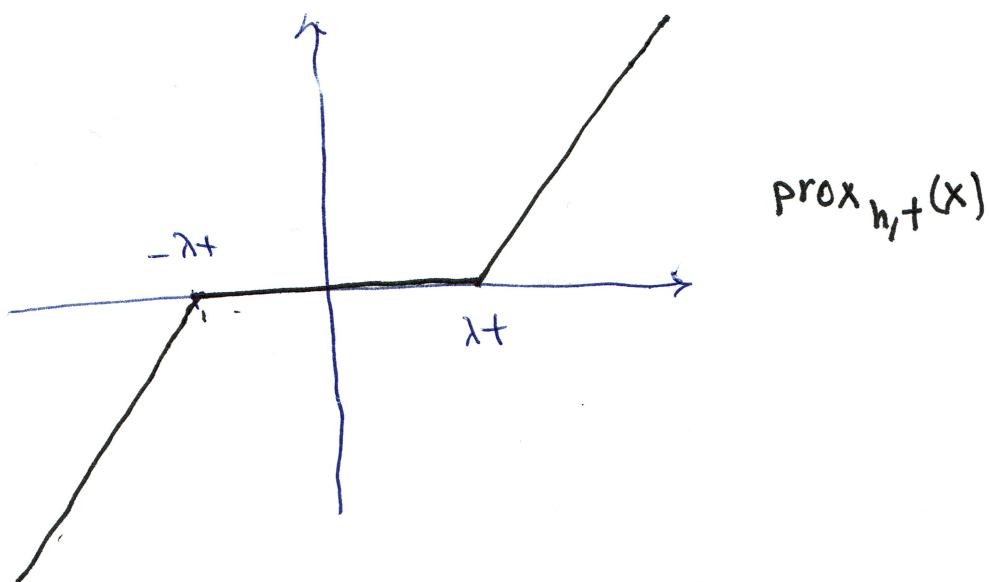
$$\left\{ \begin{array}{l} \text{if } \cancel{|x_i|} \leq \lambda t \quad \boxed{\text{Case I}} \\ \text{if } |x_i| > \lambda t \quad \boxed{\text{Case II}} \end{array} \right.$$

Case I  $\Rightarrow z_i = 0 ; \partial \|z\|_{1,i} = \frac{x}{\lambda t}$

Case II  $\Rightarrow z_i = x_i - \lambda t ; \partial \|z\|_{1,i} = \text{sgn}(x_i)$

$$\text{S}_{\lambda t}(\beta) = \text{prox}_{h,+}(\beta) = \begin{cases} \beta_i - \lambda t & \beta_i > \lambda t \\ 0 & -\lambda t \leq \beta_i \leq \lambda t \\ \beta_i + \lambda t & \beta_i < -\lambda t \end{cases}$$

(5)



$$\Rightarrow \beta^{(k)} = S_{\lambda+}(\beta^{(k-1)} - X^T(y - X\beta))$$

Convergence :  $f(x^{(k)}) - f^* \leq \frac{\|x^{(0)} - x^*\|_2^2}{2+k}$

Accelerated Proximal Gradient Descent

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

$$\left\{ \begin{array}{l} v = x^{(k-1)} + \frac{k-2}{k+1} (x^{(k-1)} - x^{(k-2)}) \\ x^{(k)} = \text{prox}_{h,t_k}(v - t_k \nabla g(v)) \end{array} \right.$$

(6)

Convergence:

$$f(x^{(k)}) - f^* \leq \frac{2 \|x^{(0)} - x^*\|_2^2}{t(k+1)^2}$$

↗ acceleration improvement



### practical) Notes

Arcelaration VS. Warin Start

$$\lambda_1 > \lambda_2 > \dots > \lambda_r$$

Let  $x^{(0)} = x^*(\lambda_{j-1})$

Often warm start does just as well

notice  $\|x^{(0)} - x^*\|_2^2$  is always

a factor in convergence with  
limited # of steps

# Stochastic Gradient Descent

7

Let's go back to  $\min_x f(x)$  problem.

Often,  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$

↳ observations

example: Regression  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$f(\beta) = \frac{1}{n} \sum_{i=1}^n \|y_i - x_i \beta\|_2^2$$

↳ mean squared error

Gradient Descent:  $x^{(k)} = x^{(k-1)} - t_k \frac{1}{n} \sum_{i=1}^n \nabla f_i$

Stochastic " pick  $i_k \in \{1, \dots, n\}$  at random

$$x^{(k)} = x^{(k-1)} - t_k \nabla f_{i_k}(x^{(k-1)})$$

" ↳ only a single term "

(8)

We hope that

$$\mathbb{E} [\nabla f_{i_k}(x^{(k-1)})] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^{(k-1)})$$

mini batch : pick  $1 \leq b \leq n$

and  $\underbrace{I_k \subseteq \{1, \dots, n\}}_{\text{random}}$  s.t.  $|I_k| = b$

$$x^{(k)} = x^{(k-1)} - \tau_k \frac{1}{b} \sum_{i \in I_k} \nabla f_i(x^{(k-1)})$$

there is higher chance that

$$\mathbb{E} \left[ \frac{1}{b} \sum_{i \in I_k} \nabla f_i(x^{(k-1)}) \right] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^{(k-1)})$$

convex

GD

SGD

$$O(\frac{1}{\sqrt{k}})$$

$$O(\frac{1}{\sqrt{k}})$$

Lipschitz

$$O(\frac{1}{k})$$

$$O(\frac{1}{\sqrt{k}})$$

Strongly Convex

$$O(\gamma^k)$$

$$O(1/k)$$

Theorem

$f$  is convex and M-Lipschitz

(9)

$$(f(y) \leq f(x) + \nabla f(x)(y-x) + \frac{M}{2} \|y-x\|_2^2)$$

OR:  $\nabla^2 f \leq M I$

$$\text{Let } t_k \leq \frac{1}{M + \sigma \sqrt{k}} = \frac{\|x^{(0)} - x^*\|_2}{M \|x^{(0)} - x^*\|_2 + \sigma \sqrt{k}}$$

$$\text{with } \sigma^2 = \max_x \|\hat{\nabla} f(x) - \nabla f(x)\|_2^2$$

how good our estimate of gradient is

notice if  $\sigma=0$  then

$$t_k \leq \frac{1}{M} \text{ (exactly GD)}$$

Then

$$\mathbb{E}\left[\frac{1}{k} \sum f(x^{(k)})\right] - f(x^*) \leq$$

$$\frac{M \|x^{(0)} - x^*\|_2^2}{k} + \frac{\sigma \|x^{(0)} - x^*\|_2}{\sqrt{k}}$$

notice if  $\sigma=0$  then

this is exactly GD error bound

Proof

$$f(x^{(k+1)}) \leq f(x^{(k)}) + \nabla f(x^{(k)})(x^{(k+1)} - x^{(k)}) + \frac{M}{2} \|x^{(k+1)} - x^{(k)}\|_2^2$$

by convexity:

$$f(x^*) \geq f(x^{(k)}) + \nabla f(x^{(k)})(x^* - x^{(k)})$$

Hence:

$$f(x^{(k+1)}) - f(x^*) \leq \nabla f(x^{(k)})(x^{(k)} - x^*) + \nabla f(x^{(k)})(x^{(k+1)} - x^{(k)}) + \frac{M}{2} \|x^{(k+1)} - x^{(k)}\|_2^2$$

By SGD step,  $x^{(k+1)} = x^{(k)} - t_k \hat{\nabla} f(x^{(k)})$ ; hence

$$\begin{aligned} f(x^{(k+1)}) - f(x^*) &\leq \nabla f(x^{(k)})(x^{(k)} - x^*) + \nabla f(x^{(k)})(-t_k \hat{\nabla} f(x^{(k)})) \\ &\quad + \frac{M t_k^2}{2} \|\hat{\nabla} f(x^{(k)})\|_2^2 \\ &= \cancel{(\nabla f(x^{(k)}) - \hat{\nabla} f(x^{(k)}))}(x^{(k)} - x^*) \\ &\quad + \hat{\nabla} f(x^{(k)})(x^{(k)} - x^*) \\ &\quad + (\nabla f(x^{(k)}) - \hat{\nabla} f(x^{(k)}))(-t_k \hat{\nabla} f(x^{(k)})) \\ &\quad + \frac{M t_k^2}{2} \|\hat{\nabla} f(x^{(k)})\|_2^2 \end{aligned}$$

$$(\nabla - \hat{\nabla})(x^{(k)} - x^*) + (\nabla - \hat{\nabla})(-\tau_k \hat{\nabla})$$

$$= (\nabla - \hat{\nabla})(x^{(k)} - x^* - \underbrace{\tau_k \nabla}_{\alpha^2}) + \tau_k (\nabla - \hat{\nabla})(\underbrace{\nabla - \hat{\nabla}}_{\alpha^2})$$

$$\Rightarrow E[(\nabla - \hat{\nabla})(x^{(k)} - x^*) + (\nabla - \hat{\nabla})(-\tau_k \hat{\nabla})] = \tau_k \alpha^2$$

Notice :

$$\begin{aligned} \frac{1}{2\tau_k} \|x^{(k+1)} - x^*\| &= \frac{1}{2\tau_k} \|x^{(k)} - x^* - \tau_k \hat{\nabla}\|_2^2 \\ &= \frac{1}{2\tau_k} \|x^{(k)} - x^*\|_2^2 + \frac{\tau_k}{2} \|\hat{\nabla}\|_2^2 \\ &\quad - \hat{\nabla}(x^{(k)} - x^*) \end{aligned}$$

NOW : the rest of the terms

$$\begin{aligned} &\hat{\nabla}(x^{(k)} - x^*) - \underbrace{\tau_k \|\hat{\nabla}\|_2^2}_{\text{assumption on } \tau_k} + \frac{M\tau_k^2}{2} \|\hat{\nabla}\|_2^2 \\ &\leq \hat{\nabla}(x^{(k)} - x^*) - \frac{\tau_k}{2} \|\hat{\nabla}\|_2^2 \\ &\stackrel{(I)}{=} \frac{1}{2\tau_k} \left[ \|x^{(k)} - x^*\|_2^2 - \|x^{(k+1)} - x^*\|_2^2 \right] \end{aligned}$$

Finally

(12)

$$f(x^{(k+1)}) - f(x^*) \leq \frac{1}{2t_k} \left[ \|x^{(k)} - x^*\|_2^2 - \|x^{(k+1)} - x^*\|_2^2 \right] + t_k \sigma^2$$

$$\Rightarrow \mathbb{E} \left[ \frac{1}{k} \sum f(x^{(k)}) \right] - f(x^*) \leq \frac{\|x^{(0)} - x^*\|_2^2}{2t_k k} + t_k \sigma^2$$

plugging  $t_k \leq \frac{1}{M + \frac{\sigma\sqrt{k}}{\|x^{(0)} - x^*\|_2}}$  we get  
the bound

— 0 — 0 —

## Variance Reduction

— We want " $\sigma$ " to be small !

## Stochastic Average Gradient (SAG)

→ Initialize  $g_i^{(0)} = \nabla f_i(x^{(0)})$

→ Randomly pick  $i_k \in \{1, \dots, n\}$  and  $\begin{cases} g_{i_k}^{(k)} = \nabla f_{i_k}(x^{(k-1)}) \\ g_{\neq i_k}^{(k)} = g_{\neq i_k}^{(k-1)} \end{cases}$

→  $\hat{V} = \frac{1}{n} \sum_{i=1}^n g_i^{(k)}$

$$\text{SAG convergence: } t_k \leq \frac{1}{16M}; g_i^{(0)} = \nabla f_i(x^{(0)}) - \frac{\cancel{13}}{\nabla f(x^{(0)})}$$

$$\mathbb{E} \left[ \frac{1}{k} \sum f(x^{(k)}) \right] - f(x^*) \leq \frac{48n}{k} (f(x^{(0)}) - f(x^*))$$

$$+ \frac{128M}{k} \|x^{(0)} - x^*\|_2^2$$

Does NOT depend on "σ"

SAG, strong convexity convergence:

$$\mathbb{E} \left[ \frac{1}{k} \sum f(x^{(k)}) \right] - f(x^*) \leq \gamma(m, M, n) \sum_{i=1}^k \left( \frac{3}{2} f(x^{(0)}) - f(x^*) \right)$$

exponential

BUT

$$+ \frac{4M}{n} \|x^{(0)} - x^*\|_2^2$$

↓  
indep. of "k"

SAGA

$$\hat{\nabla} = \overset{\wedge}{g}_{ik}^{(k)} - \overset{\wedge}{g}_{ik}^{(k-1)} + \underbrace{\frac{1}{n} \sum_{i=1}^n \overset{\wedge}{g}_i^{(k-1)}}_{\text{SAG}}$$

AdaGrad (Adam)

$$\hat{\nabla}_i = \frac{\overset{\wedge}{g}_i^{(k)}}{\sqrt{\sum_{j=1}^k (\overset{\wedge}{g}_j^{(0)})^2}}$$

(8)

We hope that

$$\mathbb{E} [\nabla f_{i_k}(x^{(k-1)})] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^{(k-1)})$$

mini batch : pick  $1 \leq b \leq n$

and  $\underbrace{I_k \subseteq \{1, \dots, n\}}_{\text{random}}$  s.t.  $|I_k| = b$

$$x^{(k)} = x^{(k-1)} - \tau_k \frac{1}{b} \sum_{i \in I_k} \nabla f_i(x^{(k-1)})$$

there is higher chance that

$$\mathbb{E} \left[ \frac{1}{b} \sum_{i \in I_k} \nabla f_i(x^{(k-1)}) \right] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^{(k-1)})$$

convex

GD

SGD

$O(\frac{1}{\sqrt{k}})$

$O(\frac{1}{\sqrt{k}})$

Lipschitz

$O(\frac{1}{k})$

$O(\frac{1}{\sqrt{k}})$

Strongly Convex

$O(\gamma^k)$

$O(1/k)$