

The choice of metric in subgradient methods

Stephen Boyd & John Duchi
(with help from P. Giselsson and Y. Carmon)

EE364b, Stanford University

Dual averaging with general norms

distance generating function h , 1-strongly-convex w.r.t. $\|\cdot\|$:

$$h(y) \geq h(x) + \nabla h(x)^T (y - x) + \frac{1}{2} \|x - y\|^2$$

Fenchel conjugate

$$h^*(\theta) = \sup_{x \in C} \{\theta^T x - h(x)\}, \quad \nabla h^*(\theta) = \operatorname{argmax}_{x \in C} \{\theta^T x - h(x)\}$$

$\nabla h, \nabla h^*$ take us “through the mirror” and back

$$x \begin{array}{c} \xrightarrow{\nabla h} \\ \xleftarrow{\nabla h^*} \end{array} \theta$$

Dual averaging subgradient method

- (1) get subgradient $g^{(k)} \in \partial f(x^{(k)})$
- (2) $\theta^{(k+1)} = \theta^{(k)} - \alpha_k g^{(k)}$
- (3) $x^{(k+1)} = \nabla h^*(\theta^{(k+1)})$

$h(x) = \frac{1}{2} \|x\|_2^2$ recovers standard case

Convergence analysis

Dual averaging:

$$g^{(k)} \in \partial f(x^{(k)}), \theta^{(k+1)} = \theta^{(k)} - \alpha_k g^{(k)}, x^{(k+1)} = \nabla h^*(\theta^{(k+1)})$$

Bregman divergence

$$D_{h^*}(\theta' \leftarrow \theta) = h^*(\theta') - h^*(\theta) - \nabla h^*(\theta)^T(\theta' - \theta)$$

Let $\theta^* = \nabla h(x^*)$,

$$\begin{aligned} D_{h^*}(\theta^{(k+1)} \leftarrow \theta^*) &= D_{h^*}(\theta^{(k)} \leftarrow \theta^*) \\ &\quad + (\theta^{(k+1)} - \theta^{(k)})^T (\nabla h^*(\theta^{(k)}) - \nabla h^*(\theta^*)) \\ &\quad + D_{h^*}(\theta^{(k+1)} \leftarrow \theta^{(k)}) \end{aligned}$$

and

$$(\theta^{(k+1)} - \theta^{(k)})^T (\nabla h^*(\theta^{(k)}) - \nabla h^*(\theta^*)) = -\alpha_k g^{(k)T} (x^{(k)} - x^*)$$

Convergence analysis continued

From convexity and $g^{(k)} \in \partial f(x^{(k)})$,

$$f(x^{(k)}) - f(x^*) \leq g^{(k)T} (x^{(k)} - x^*)$$

Therefore

$$\begin{aligned} \alpha_k [f(x^{(k)}) - f(x^*)] &\leq D_{h^*}(\theta^{(k)} \leftarrow \theta^*) - D_{h^*}(\theta^{(k+1)} \leftarrow \theta^*) \\ &\quad + D_{h^*}(\theta^{(k+1)} \leftarrow \theta^{(k)}) \end{aligned}$$

Fact: h is 1-strongly-convex w.r.t. $\|\cdot\| \Leftrightarrow D_h(x' \leftarrow x) \geq \frac{1}{2} \|x' - x\|^2 \Leftrightarrow$
 h^* is 1-smooth w.r.t. $\|\cdot\|_* \Leftrightarrow D_{h^*}(\theta' \leftarrow \theta) \leq \frac{1}{2} \|\theta' - \theta\|_*^2$

Bounding the $D_{h^*}(\theta^{(k+1)} \leftarrow \theta^{(k)})$ terms and telescoping gives

$$\sum_{i=1}^k \alpha_i [f(x^{(i)}) - f(x^*)] \leq D_{h^*}(\theta^{(1)} \leftarrow \theta^*) + \frac{1}{2} \sum_{i=1}^k \alpha_i^2 \|g^{(i)}\|_*^2$$

Convergence guarantees

Note: $D_{h^*}(\theta^{(1)} \leftarrow \theta^*) = D_h(x^* \leftarrow x^{(1)})$

Most general guarantee,

$$\sum_{i=1}^k \alpha_i [f(x^{(i)}) - f(x^*)] \leq D_h(x^* \leftarrow x^{(1)}) + \frac{1}{2} \sum_{i=1}^k \alpha_i^2 \|g^{(i)}\|_*^2$$

Fixed step size $\alpha_k = \alpha$

$$\frac{1}{k} \sum_{i=1}^k f(x^{(i)}) - f(x^*) \leq \frac{1}{\alpha k} D_h(x^* \leftarrow x^{(1)}) + \frac{\alpha}{2} \max_i \|g^{(i)}\|_*^2$$

in general, converges if

- $D_h(x^* \leftarrow x^{(1)}) < \infty$
- $\sum_k \alpha_k = \infty$ and $\alpha_k \rightarrow 0$
- for all $g \in \partial f(x)$ and $x \in C$, $\|g\|_* \leq G$ for some $G < \infty$

Stochastic gradients are fine!

Mirror descent

$$x^{(k+1)} = \operatorname{argmin}_{x \in C} \left\{ \alpha_k g^{(k)T} x + D_h(x \leftarrow x^{(k)}) \right\} = \nabla h^* \left(\nabla h(x^{(k)}) - \alpha_k g^{(k)} \right)$$

Dual averaging

$$\theta^{(k+1)} = \theta^{(k)} - \alpha_k g^{(k)}$$

$$x^{(k+1)} = \nabla h^*(\theta^{(k+1)})$$

Mirror descent

$$\tilde{\theta}^{(k+1)} = \theta^{(k)} - \alpha_k g^{(k)}$$

$$x^{(k+1)} = \nabla h^*(\tilde{\theta}^{(k+1)})$$

$$\theta^{(k+1)} = \nabla h(x^{(k+1)})$$

Analysis essentially the same:

$$D_{h^*}(\tilde{\theta}^{(k+1)} \leftarrow \theta^*) \leq D_{h^*}(\theta^{(k)} \leftarrow \theta^*) - \alpha_k (f(x^{(k)}) - f(x^*)) + \frac{\alpha_k^2}{2} \|g^{(k)}\|_*^2$$

To telescope, we note that

$$D_{h^*}(\theta^{(k+1)} \leftarrow \theta^*) \leq D_{h^*}(\tilde{\theta}^{(k+1)} \leftarrow \theta^*)$$

Hence, the guarantee we derived holds for mirror descent as well.

Mirror descent examples

- Usual (projected) subgradient descent: $h(x) = \frac{1}{2} \|x\|_2^2$
- With constraints of simplex, $C = \{x \in \mathbf{R}_+^n \mid \mathbf{1}^T x = 1\}$, use negative entropy

$$h(x) = \sum_{i=1}^n x_i \log x_i$$

- (1) Strongly convex with respect to ℓ_1 -norm
- (2) With $x^{(1)} = \mathbf{1}/n$, have $D_h(x^* \leftarrow x^{(1)}) \leq \log n$ for $x^* \in C$
- (3) If $\|g\|_\infty \leq G_\infty$ for $g \in \partial f(x)$ for $x \in C$,

$$f_{\text{best}}^{(k)} - f^* \leq \frac{\log n}{\alpha k} + \frac{\alpha}{2} G_\infty^2$$

- (4) Can be much better than regular subgradient decent...

Example

Robust regression problem (an LP):

$$\text{minimize } f(x) = \|Ax - b\|_1 = \sum_{i=1}^m |a_i^T x - b_i|$$

$$\text{subject to } x \in C = \{x \in \mathbf{R}_+^n \mid \mathbf{1}^T x = 1\}$$

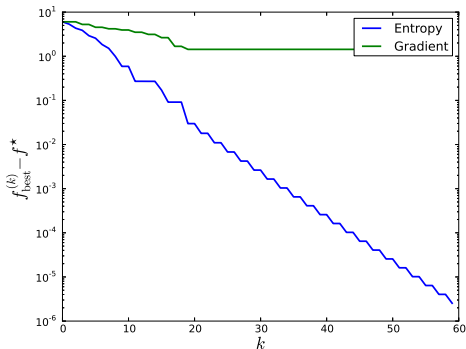
subgradient of objective is $g = \sum_{i=1}^m \text{sign}(a_i^T x - b_i) a_i$

- Projected subgradient update ($h(x) = (1/2) \|x\|_2^2$): homework
- Mirror descent update ($h(x) = \sum_{i=1}^n x_i \log x_i$):

$$x_i^{(k+1)} = \frac{x_i^{(k)} \exp(-\alpha g_i^{(k)})}{\sum_{j=1}^n x_j^{(k)} \exp(-\alpha g_j^{(k)})}$$

Example

Robust regression problem with $a_i \sim N(0, I_{n \times n})$ and $b_i = (a_{i,1} + a_{i,2})/2 + \varepsilon_i$ where $\varepsilon_i \sim N(0, 10^{-2})$, $m = 20$, $n = 3000$



stepsizes chosen according to best bounds (but still sensitive to stepsize choice)

Variable metric subgradient methods

subgradient method with variable metric $H_k \succ 0$:

- (1) get subgradient $g^{(k)} \in \partial f(x^{(k)})$
 - (2) update (diagonal) metric H_k
 - (3) update $x^{(k+1)} = x^{(k)} - H_k^{-1} g^{(k)}$
- matrix H_k generalizes step-length α_k

there are many such methods (Ellipsoid method, AdaGrad, ...)

Variable metric projected subgradient method

same, with projection carried out in the H_k metric:

- (1) get subgradient $g^{(k)} \in \partial f(x^{(k)})$
- (2) update (diagonal) metric H_k
- (3) update $x^{(k+1)} = P_{\mathcal{X}}^{H_k} (x^{(k)} - H_k^{-1} g^{(k)})$

where

$$\Pi_{\mathcal{X}}^H(y) = \operatorname{argmin}_{x \in \mathcal{X}} \|x - y\|_H^2$$

and $\|x\|_H = \sqrt{x^T H x}$.

Convergence analysis

since $\Pi_{\mathcal{X}}^{H_k}$ is non-expansive in the $\|\cdot\|_{H_k}$ norm, we get

$$\begin{aligned}\|x^{(k+1)} - x^*\|_{H_k}^2 &= \left\| P_{\mathcal{X}}^{H_k} \left(x^{(k)} - H_k^{-1} g^{(k)} \right) - P_{\mathcal{X}}^{H_k} (x^*) \right\|_{H_k}^2 \\ &\leq \|x^{(k)} - H_k^{-1} g^{(k)} - x^*\|_{H_k}^2 \\ &= \|x^{(k)} - x^*\|_{H_k}^2 - 2(g^{(k)})^T (x^{(k)} - x^*) + \|g^{(k)}\|_{H_k^{-1}}^2 \\ &\leq \|x^{(k)} - x^*\|_{H_k}^2 - 2(f(x^{(k)}) - f^*) + \|g^{(k)}\|_{H_k^{-1}}^2.\end{aligned}$$

using $f^* = f(x^*) \geq f(x^{(k)}) + g^{(k)T}(x^* - x^{(k)})$

apply recursively, use

$$\sum_{i=1}^k \left(f(x^{(i)}) - f^* \right) \geq k \left(f_{\text{best}}^{(k)} - f^* \right)$$

and rearrange to get

$$f_{\text{best}}^{(k)} - f^* \leq \frac{\|x^{(1)} - x^*\|_{H_1}^2 + \sum_{i=1}^k \|g^{(i)}\|_{H_i^{-1}}^2}{2k} + \frac{\sum_{i=2}^k \left(\|x^{(i)} - x^*\|_{H_i}^2 - \|x^{(i)} - x^*\|_{H_{i-1}}^2 \right)}{2k}$$

numerator of additional term can be bounded to get estimates

- for general $H_k = \mathbf{diag}(h_k)$

$$f_{\text{best}}^k - f^* \leq \frac{R_\infty^2 \|H_1\|_1 + \sum_{i=1}^k \|g^{(i)}\|_{H_i^{-1}}^2}{2k} + \frac{R_\infty^2 \sum_{i=2}^k \|H_i - H_{i-1}\|_1}{2k}$$

- for $H_k = \mathbf{diag}(h_k)$ with $h_i \geq h_{i-1}$ for all i

$$f_{\text{best}}^k - f^* \leq \frac{\sum_{i=1}^k \|g^{(i)}\|_{H_i^{-1}}^2}{2k} + \frac{R_\infty^2 \|h_k\|_1}{2k}$$

where $\max_{1 \leq i \leq k} \|x^{(i)} - x^*\|_\infty \leq R_\infty$

converges if

- $R_\infty < \infty$ (e.g. if \mathcal{X} is compact)
- $\sum_{i=1}^k \|g^{(i)}\|_{H_i^{-1}}^2$ grows slower than k
- $\sum_{i=2}^k \|H_i - H_{i-1}\|_1$ grows slower than k **or**
 $h_i \geq h_{i-1}$ for all i and $\|h_k\|_1$ grows slower than k

AdaGrad

AdaGrad — adaptive subgradient method

- (1) get subgradient $g^{(k)} \in \partial f(x^{(k)})$
- (2) choose metric H_k :
 - set $S_k = \sum_{i=1}^k \mathbf{diag}(g^{(i)})^2$
 - set $H_k = \frac{1}{\alpha} S_k^{\frac{1}{2}}$
- (3) update $x^{(k+1)} = P_{\mathcal{X}}^{H_k} (x^{(k)} - H_k^{-1} g^{(k)})$

where $\alpha > 0$ is step-size

AdaGrad – motivation

- for fixed $H_k = H$ we have estimate:

$$f_{\text{best}}^{(k)} - f^* \leq \frac{1}{2k} (x^{(1)} - x^*)^T H (x^{(1)} - x^*) + \frac{1}{2k} \sum_{i=1}^k \|g^{(i)}\|_{H^{-1}}^2$$

- **idea:** Choose *diagonal* $H_k \succ 0$ that minimizes this estimate in hindsight:

$$H_k = \operatorname{argmin}_h \max_{x, y \in C} (x - y)^T \mathbf{diag}(h) (x - y) + \sum_{i=1}^k \|g^{(i)}\|_{\mathbf{diag}(h)^{-1}}^2$$

- optimal $H_k = \frac{1}{R_\infty} \mathbf{diag} \left(\sqrt{\sum_{i=1}^k (g_1^{(i)})^2}, \dots, \sqrt{\sum_{i=1}^k (g_n^{(i)})^2} \right)$
- **intuition:** adapt step-length based on historical step lengths

AdaGrad – convergence

by construction, $H_i = \frac{1}{\alpha} \mathbf{diag}(h_i)$ and $h_i \geq h_{i-1}$, so

$$\begin{aligned} f_{\text{best}}^{(k)} - f^* &\leq \frac{1}{2k} \sum_{i=1}^k \|g^{(i)}\|_{H_i^{-1}}^2 + \frac{1}{2k\alpha} R_\infty^2 \|h_k\|_1 \\ &\leq \frac{\alpha}{k} \|h_k\|_1 + \frac{1}{2k\alpha} R_\infty^2 \|h_k\|_1 \end{aligned}$$

(second line is a theorem)

also have (with $\alpha = R_\infty^2$) and for compact sets C

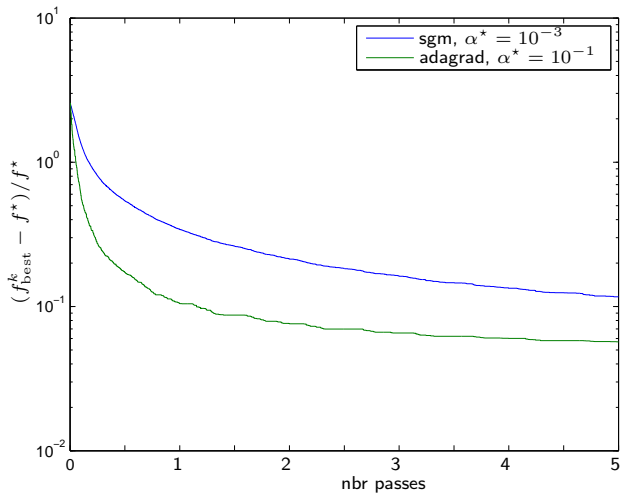
$$f_{\text{best}}^{(k)} - f^* \leq \frac{2}{k} \inf_{h \geq 0} \left\{ \sup_{x, y \in C} (x - y)^T \mathbf{diag}(h)(x - y) + \sum_{i=1}^k \|g^{(i)}\|_{\mathbf{diag}(h)^{-1}}^2 \right\}$$

Example

Classification problem:

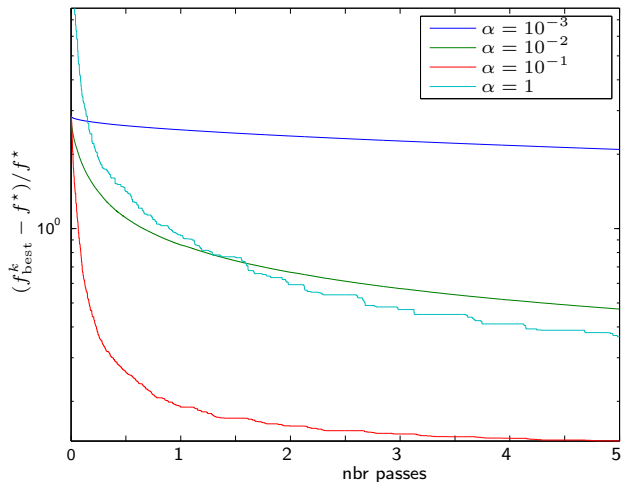
- **Data:** $\{a_i, b_i\}$, $i = 1, \dots, 50000$
 - $a_i \in \mathbf{R}^{1000}$
 - $b \in \{-1, 1\}$
 - Data created with 5% mis-classifications w.r.t. $w = \mathbf{1}$, $v = 0$
- **Objective:** find classifiers $w \in \mathbf{R}^{1000}$ and $v \in \mathbf{R}$ such that
 - $a_i^T w + v > 1$ if $b = 1$
 - $a_i^T w + v < 1$ if $b = -1$
- **Optimization method:**
 - Minimize hinge-loss: $\sum_i \max(0, 1 - b_i(a_i^T w + v))$
 - Choose example uniformly at random, take sub-gradient step w.r.t. that example

Best subgradient method vs best AdaGrad



Often best AdaGrad performs better than best subgradient method

AdaGrad with different step-sizes α :



Sensitive to step-size selection (like standard subgradient method)