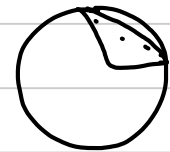


LECTURE - 6

MIRROR DESCENT

$$x \in \begin{cases} x_i \geq 0 \\ \sum_i x_i = 1 \end{cases} \quad \min \quad \underbrace{\frac{1}{2} \|Ax - b\|^2}_{= f(x)}$$



PGD: $x_{t+1} = P_x(x_t - \eta \nabla f(x_t)) \quad \left| \quad O\left(\frac{\|A^T A\| \sqrt{2}}{t}\right)\right.$

Smoothness $\rightarrow L = \|A^T A\| ; \quad \|x_0 - x^*\| \leq \text{Diam}(x) \leq \sqrt{2}$

each Column of A : $\|A_i\| \leq 1$

$$\|A^T A\| = \|A\|^2 \text{ Could be } \approx \underline{\underline{d.}}$$

$[| | | | |]$
↓
unit norm
columns.

When χ and $f(\cdot)$ have good properties w.r.t.

$\|\cdot\|$ which is not Euclidean, Can we do better?

$$\begin{array}{ccc}
 x & & \nabla f \\
 \downarrow & & \downarrow \\
 \|\cdot\| & & \|\cdot\|_* \\
 \text{e.g., for probabilities } \|\cdot\|_2 & \longrightarrow & \|\cdot\|_\infty
 \end{array}$$

Defn. of dual norm: Given a norm $\|\cdot\|$ on X , the dual norm on the space of linear operators $\mathcal{L}(X)$ is defined as:

$$\|\ell\|_* = \sup_{\substack{\|x\| \leq 1 \\ x \in X}} \ell(x) \quad \forall \ell \in \mathcal{L}(X).$$

$$\|x\| \triangleq \|x\|_1 = \sum_i |x_i| \quad \xrightarrow{\text{Dual norm}} \quad \|y\|_* = \|y\|_\infty = \max_i |y_i|$$

$$\|y\|_* = \sup_{\|x\|_1 \leq 1} \sum_i x_i y_i \quad x_{i^*} = \begin{cases} \text{sign}(y_{i^*}) & \text{on } i^* = \text{argmax}_i |y_i| \\ 0 & \text{o.w.} \end{cases}$$

$$= \|y\|_\infty$$

$\|\cdot\|$ on x , $\|\cdot\|_*$ on $\nabla f(x)$

$$\downarrow$$

$$\|\cdot\|_2$$

$$\downarrow$$

$$\|\cdot\|_\infty$$

$$f(y) \geq f(x) + \underbrace{\langle \nabla f(x), y-x \rangle}_{\nabla f(x)} \underbrace{\quad}_{\in x}$$

f -Convex & $\|\nabla f(x)\|_* \leq G$. $\|\cdot\|$
Lipschitz

$$\begin{aligned} \text{GD: } x_{t+1} &= \operatorname{argmin}_x f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|_2^2 \\ &= \operatorname{argmin}_x f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \underbrace{\|x - x_t\|_2^2}_{\text{red}} \end{aligned}$$

- ||
- ① Not clear how to solve this step $\left(\sum_i |x_i - x_{t,i}|\right)^2$
 - ② Not clear what the potential function is.

Defn. Bregman divergence: Given a Convex function $\phi: \mathcal{X} \rightarrow \mathbb{R}$,

$$D_{\phi}(x; y) \triangleq \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle \geq 0$$

\downarrow
 $\because \phi$ is convex.

Mirror descent:

$$x_{t+1} = \arg \min_x f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} D_{\phi}(x; x_t)$$

Requirement: ϕ is 1-strongly convex w.r.t. $\|\cdot\|$ on \mathcal{X} .

$$\phi(y) \geq \phi(x) + \langle \nabla \phi(x), y - x \rangle + \frac{1}{2} \|x - y\|^2.$$

ASSUMPTIONS

- ① f is Convex ③ $\phi(\cdot)$ is 1-Strongly Convex
in $\|\cdot\|$ on \mathcal{X} :
- ② $\|\nabla f(x)\|_* \leq G$

$$D_\phi(x^*; x_t)$$

$$x \in \mathcal{X} \text{ and } l \in \mathcal{L}(x) \\ = \nabla f(x).$$

$$\text{Hölder's. ineq.} \quad \langle l, x \rangle \leq \|l\|_* \cdot \|x\|$$

$$x_{t+1} = \operatorname{argmin}_x \underbrace{f(x_t)}_{\text{Const.}} + \underbrace{\langle \nabla f(x_t), x - x_t \rangle}_{\text{Linear in } x} + \underbrace{\frac{1}{2\eta} D_\phi(x; x_t)}_{\text{Convex in } x}.$$

$$\nabla f(x_t) + \frac{1}{2\eta} \nabla_x \underbrace{D_\phi(x; x_t)}_{\phi(x) - \phi(x_t) - \langle \nabla \phi(x_t), x - x_t \rangle} = 0.$$

$$\nabla f(x_t) + \frac{1}{2\eta} [\nabla \phi(x_{t+1}) - 0 - \nabla \phi(x_t)] = 0.$$

$$x_{t+1}: \underbrace{\nabla \phi(x_{t+1})}_{= x_{t+1}} = \underbrace{\nabla \phi(x_t)}_{x_t} - 2\eta \nabla f(x_t). \quad] \text{ MD Step.}$$

↳ ①

By optimality of x_{t+1} ,

$$f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{1}{2\eta} D\phi(x_{t+1}; x_t) \leq f(x_t) + \langle \nabla f(x_t), x_t - x_t \rangle + \frac{1}{2\eta} D\phi(x_t; x_t) = f(x_t)$$

$\frac{1}{4\eta} \|x_t - x_{t+1}\|^2 \stackrel{1\text{-Str. conv. of } \phi}{\leq} \frac{1}{2\eta} D\phi(x_{t+1}; x_t) \leq -\langle \nabla f(x_t), x_{t+1} - x_t \rangle$

$\stackrel{\text{Hölder's}}{\leq} \|\nabla f(x_t)\|_* \cdot \|x_{t+1} - x_t\|$

$$\|x_t - x_{t+1}\| \leq 4\eta \|\nabla f(x_t)\|_* \leq 4\eta G \rightarrow \textcircled{2}$$

$$D_{\phi}(x; x_{t+1}) = \phi(x) - \phi(x_{t+1}) - \langle \nabla \phi(x_{t+1}), x - x_{t+1} \rangle$$

$$\stackrel{\textcircled{1}}{=} \phi(x) - \phi(x_{t+1}) - \langle \nabla \phi(x_t) - 2\eta \nabla f(x_t), x - x_{t+1} \rangle$$

Conv. of ϕ

$$\leq \phi(x) - \phi(x_t) - \langle \nabla \phi(x_t), x_{t+1} - x_t \rangle$$

$$- \langle \nabla \phi(x_t) - 2\eta \nabla f(x_t), x - x_{t+1} \rangle$$

$$= \phi(x) - \phi(x_t) - \langle \nabla \phi(x_t), x - x_t \rangle$$

$$- 2\eta \langle \nabla f(x_t), x_{t+1} - x \rangle$$

$$= D_{\phi}(x; x_t) - 2\eta \langle \nabla f(x_t), x_t - x \rangle$$

$$-2\eta \langle \nabla f(x_t), x_{t+1} - x_t \rangle.$$

$$\stackrel{\text{Hölder's}}{\leq} D_\phi(x; x_t) - 2\eta \langle \nabla f(x_t), x_t - x \rangle + 2\eta \|\nabla f(x_t)\|_* \cdot \|x_{t+1} - x_t\|$$

→ Convexity

②

$$\leq D_\phi(x; x_t) - 2\eta [f(x_t) - f(x)]$$

$$+ 2\eta \underbrace{\|\nabla f(x_t)\|_*}_{\leq G} \cdot 4\eta G.$$

$$\leq D_\phi(x; x_t) - 2\eta [f(x_t) - f(x)] + 8\eta^2 G^2.$$

$$D_{\phi}(x; x_{t+1}) \leq D_{\phi}(x; x_t) - 2\eta[f(x_t) - f(x)] + 8\eta^2 G^2.$$

Telescoping and averaging gives us

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [f(x_t) - f(x)] &\leq \frac{D_{\phi}(x; x_0) - D_{\phi}(x; x_T)}{2\eta T} + 4\eta G^2. \\ &\leq \frac{D_{\phi}(x; x_0)}{2\eta T} + 4\eta G^2. \end{aligned}$$

$$\eta = \frac{1}{G} \sqrt{\frac{D_{\phi}(x; x_0)}{8T}}$$

$$\frac{1}{T} \sum_{t=0}^{T-1} [f(x_t) - f(x)] \leq C \cdot \frac{G \sqrt{D_\phi(x; x_0)}}{\sqrt{T}}$$

$$\leq C \cdot \frac{G \|x - x_0\|}{\sqrt{T}}$$

① If $\|\cdot\| = \|\cdot\|_2 \rightarrow$ Euclidean

, then choose $\phi = \frac{1}{2} \|x\|_2^2$
 \rightarrow 1-str. convex w.r.t. $\|\cdot\|_2$

$$D_\phi(x; x_t) = \frac{1}{2} \|x - x_t\|^2$$

In this case

a) Mirror descent = GD

b) Obtained bounds are the same.

$$\textcircled{2} \quad \|\cdot\| = \|\cdot\|_1 \quad \text{and} \quad \mathcal{X} = \{x : x_i \geq 0 ; \sum x_i = 1\}.$$

$$\phi(x) = \sum_i x_i \log x_i$$

$$D_\phi(x; y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle$$

$$= \sum_i x_i \log x_i - \sum_i y_i \log y_i - \sum_i (1 + \log y_i) (x_i - y_i)$$

$$= \sum_i x_i \log \frac{x_i}{y_i} + \underbrace{\sum_i (x_i - y_i)}_{=0}$$

$$= \text{KL}(x \parallel y).$$

ASIDÉ

$$x_i \geq 0 : \phi(x) = \sum_i x_i \log x_i - \sum_i x_i \quad \left[\sum_i x_i \neq 1 \right].$$

$$\sum_i x_i = 1, \quad x_i \geq 0. \rightarrow \phi(x) = \sum_i x_i \log x_i$$

$$D_\phi(x \parallel y) = \text{KL}(x \parallel y).$$

① ϕ [or eq. D_ϕ] is 1-Strongly convex

Pinsker's inequality: $KL(x \parallel y) \geq 2 \|x - y\|_1^2$

$$\textcircled{2} \quad x_{t+1} = \underset{x \in \mathcal{X}}{\operatorname{argmin}} \underbrace{f(x_t)}_{\text{Const.}} + \underbrace{\langle \nabla f(x_t), x - x_t \rangle}_{\text{Const.}} + \frac{1}{2\eta} D_\phi(x; x_t)$$

$$= \underset{x \in \mathcal{X}}{\operatorname{argmin}} \langle \nabla f(x_t), x \rangle + \frac{1}{2\eta} \sum_i x_i \log \frac{x_i}{(x_t)_i}$$

$$= \underset{x \in \mathcal{X}}{\operatorname{argmin}} \sum_i x_i \left[\underbrace{\nabla f(x_t)_i}_{\text{Const.}} + \frac{1}{2\eta} \log \frac{x_i}{(x_t)_i} \right].$$

$$= \operatorname{argmin}_{x \in X} \sum_i x_i \log \frac{x_i}{(x_t)_i \exp(-2\eta(\nabla f(x_t))_i)}$$

$$Z_i \stackrel{\Delta}{=} \frac{(x_t)_i \exp(-2\eta(\nabla f(x_t))_i)}{\sum_{j=1}^d (x_t)_j \exp(-2\eta(\nabla f(x_t))_j)}$$

$$= \operatorname{argmin}_{x \in X} \sum_i x_i \log \frac{x_i}{Z_i} - \underbrace{\sum_i x_i}_{=1} \underbrace{\log \left(\sum_{j=1}^d (x_t)_j \exp(-2\eta(\nabla f(x_t))_j) \right)}_{\text{Constant}}$$

Constant

$$= \operatorname{argmin}_{x \in \mathcal{X}} \text{KL}(x \| z) = z.$$

Ind. discovered

$$x_{t+1} = z = \frac{x_t \odot \exp(-2\eta \nabla f(x_t))}{\|x_t \odot \exp(-2\eta \nabla f(x_t))\|_1}.$$

MULTIPLICATIVE
WEIGHTS
UPDATE

EX PONENTIATED

GRADIENT ALGORITHM

$$\begin{aligned} \text{Min} \quad & \frac{1}{2} \|Ax - b\|^2 \\ \text{on } \mathcal{X} \quad & x_i \geq 0 \\ & \sum x_i = 1 \end{aligned}$$

$$\|\cdot\|_1 \longrightarrow \|\cdot\|_\infty$$

on \mathcal{X} on $\nabla f(\mathcal{X})$.

$$\phi(x) = \sum_i x_i \log x_i$$

$$\text{Av. Suboptimality} \leq C \cdot \frac{\|\nabla f(x)\|_\infty \sqrt{D_\phi(x; x_0)}}{\sqrt{J}}$$

$$\begin{aligned} \nabla f(x) = A^T(Ax - b) &\rightarrow \|A^T(Ax - b)\|_\infty \leq \|A^T A x\|_\infty + \|A^T b\|_\infty \\ &\leq \max_i \|(A^T A)_i\|_\infty + \|A^T b\|_\infty \end{aligned}$$

Hölder's inequality
ith row of $A^T A$

Non-Euclidean

$$\|A^T A\|_\infty \leq$$

↓
largest element
in $A^T A$

Euclidean

$\|A^T A\|_2$ Could be
 d times larger.

$$G = \|A^T A\|_\infty + \|A^T b\|_\infty \leq \|A^T A\|_2 + \|A^T b\|_2$$

$$D_\phi(x; x_0) : x_0 = \frac{1}{d} \cdot \mathbf{1}$$

$$= \sum_i x_i \log \frac{x_i}{(x_0)_i} = \underbrace{\sum_i x_i \log x_i}_{\leq 0} + \underbrace{\sum_i x_i \log d}_{= \log d}$$

$$\text{Av. Suboptimality} \leq C \cdot \frac{(\|A^T A\|_\infty + \|A^T b\|_\infty) \cdot \sqrt{\log d}}{\sqrt{T}} \cdot \sqrt{\log d}$$

Could be $\Omega(d) \downarrow$

$$C \cdot \frac{(\|A^T A\|_2 + \|A^T b\|_2) \cdot 1}{\sqrt{T}} \cdot \sqrt{\log d}$$