

$$\min_x f(x)$$

① We know the exact form : $f(x) = \frac{1}{2} x^T A x - x^T b$

$$x^* = \underbrace{A^{-1}}_{} b$$

② $x \xrightarrow[\text{oracle}]{\text{Black box}} f(x)$] zeroth order oracle

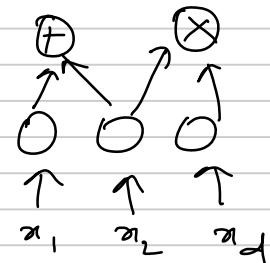
$x \longrightarrow f(x), \nabla f(x)$] First order oracle

$x \longrightarrow f(x), \nabla f(x), \nabla^2 f(x)$] Second order oracle.

Thm: If $f(\cdot)$ is represented by a circuit with simple nodes then
the $\nabla f(x)$ can be computed in time linear in the circuit size.

First order oracles can be implemented very efficiently.

We will focus only on first order algorithms.



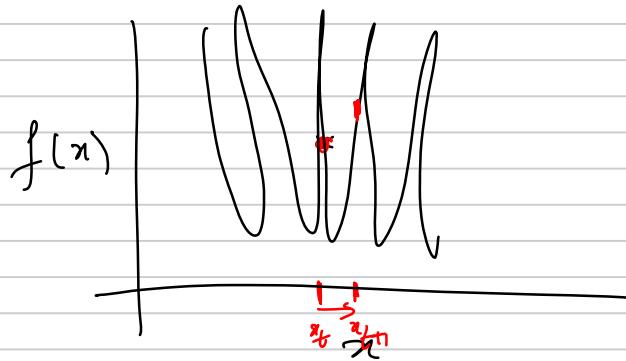
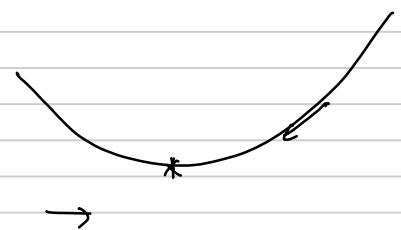
OUTLINE

- ① CONVEX OPTIMIZATION (DETERMINISTIC)
- ② STOCHASTIC CONVEX OPTIMIZATION
- ③ NONCONVEX OPT. (DET & STOCH.)
- ④ CONSTRAINED / MINIMAX OPT.

ACCESS TO $x \rightarrow f(x), \nabla f(x)$

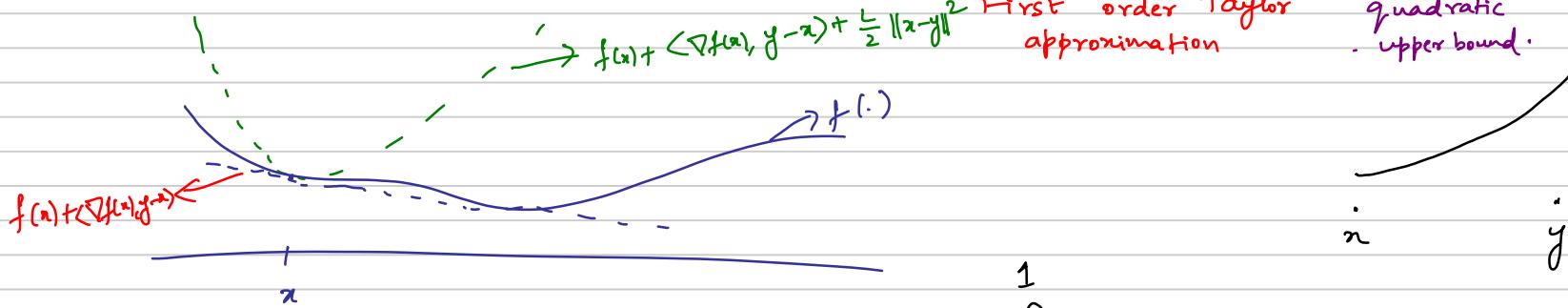
Gradient descent : Start with x_0

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



[DEFN.] SMOOTHNESS : A function $f(\cdot)$ is said to be L -Smooth if $\|\nabla f(x) - \nabla f(y)\| \leq L \cdot \|x-y\| \quad \forall x, y$.

LEMMA : If $f(\cdot)$ is L -Smooth then $f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|x-y\|^2 \quad \forall x, y$.



Proof : Given x and y ,
$$f(y) = f(x) + \int_{\tau=0}^1 \langle \nabla f((1-\tau)x + \tau y), y-x \rangle d\tau$$

$$= \boxed{f(x) + \langle \nabla f(x), y-x \rangle} \rightarrow \text{Linear Taylor expansion}$$

$$+ \int_{\tau=0}^1 \langle \nabla f((1-\tau)x + \tau y) - \nabla f(x), y-x \rangle d\tau$$

$$S \leq \int_{\tau=0}^1 \|\nabla f((1-\tau)x + \tau y) - \nabla f(x)\| \cdot \|y-x\| d\tau$$

S is the shaded region under the curve of the difference between the function value and its linear approximation.

$$\text{(Smoothness)} \leq \int_{\tau=0}^1 L \| (1-\tau)x + \tau y - x \| \cdot \|y - x\| d\tau$$

$$= L \|x - y\|^2 \int_{\tau=0}^1 \tau d\tau = \frac{L}{2} \|x - y\|^2 \quad \text{④}$$

$$GD: x_{t+1} = x_t - \eta \nabla f(x_t).$$

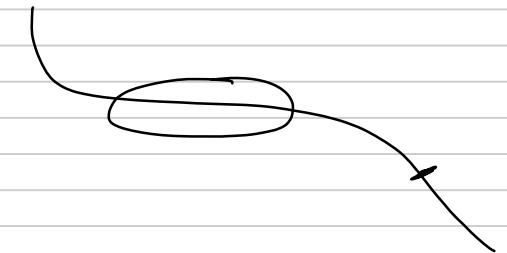
Thm: If $f(\cdot)$ is L -smooth and $\eta \leq \frac{1}{L}$ then

$$\min_{t=1,\dots,T} \|\nabla f(x_t)\|^2 \leq 2 \frac{(f(x_0) - f^*)}{\eta T} \xrightarrow{\eta \downarrow} \min_x f(x).$$

Proof:

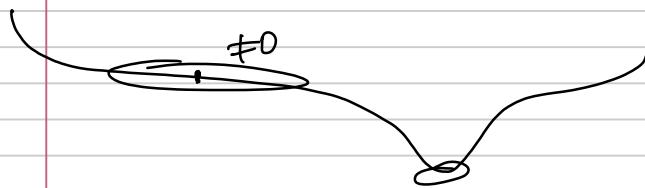
$$\begin{aligned} f(x_{t+1}) &\stackrel{\text{(Smoothness lemma)}}{\leq} f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{L}{2} \|x_t - x_{t+1}\|^2 \\ &= f(x_t) - \eta \|\nabla f(x_t)\|^2 + \frac{\eta^2}{2} \|\nabla f(x_t)\|^2 \\ &= f(x_t) - \eta \left(1 - \frac{\eta L}{2}\right) \|\nabla f(x_t)\|^2 \\ &\geq f(x_t) - \frac{\eta}{2} \|\nabla f(x_t)\|^2 \end{aligned}$$

$$f^* \leq f(x_{T+1}) \leq f(x_0) - \frac{1}{2} \sum_{t=0}^T \|\nabla f(x_t)\|^2$$



Reorganizing proves the theorem \square

Consider a special class of functions where $\nabla f(x) = 0 \Rightarrow x$ is global optimum.



\boxed{f} : $\left\{ \begin{array}{l} \text{① } \nabla f(x) = 0 \Rightarrow x \text{ is global opt. of } f(\cdot) \\ \text{② } \begin{cases} \text{if } g(x) = \underbrace{f(x)} + \langle \omega, x \rangle \text{ then} \\ \nabla g(x) = 0 \Rightarrow x \text{ is global opt. of } g(\cdot). \end{cases} \end{array} \right.$

Lemma: $\boxed{f} = \left\{ f : f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \forall x, y \right\}$.

Proof: Any $f \in f$ satisfies \nearrow

$$\text{Given } x \text{ and } y; \quad g(y) = \underbrace{f(y)} + \underbrace{\langle -\nabla f(x), y - x \rangle}$$

$$\nabla g(y) = \nabla f(y) - \nabla f(x) \Rightarrow \nabla g(y) = 0 \Rightarrow y \text{ is global opt. of } g(\cdot).$$

$$\Rightarrow g(x) \leq g(y) \forall y$$

$$\Rightarrow f(x) \leq f(y) - \langle \nabla f(x), y - x \rangle \quad \square$$

[Defn.]

CONVEX FNS: $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \quad \forall x, y$

Example: ① $f(x) = |x|^p \quad p > 1$

$$\frac{df}{dx} = p x^{p-1}$$

$$f(y) = f(x) + (y-x) \int_{\alpha=0}^1 \left. \frac{df}{dx} \right|_{(1-\alpha)x+\alpha y} \cdot d\alpha$$

$$\frac{df}{dx} \Big|_{(1-\alpha)x+\alpha y} = p [(1-\alpha)x + \alpha y]^{p-1} \geq p x^{p-1}$$

$$f(y) \geq f(x) + \underbrace{p x^{p-1}} \cdot \underbrace{(y-x)}.$$

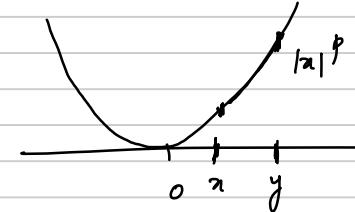
② If $f_1(\cdot)$ is convex and $f_2(\cdot)$ is convex then

$$f(x) = f_1(x) + f_2(x) \text{ is also convex.}$$

$$f(y) = f_1(y) + f_2(y)$$

$$\geq f_1(x) + \langle \nabla f_1(x), y-x \rangle + f_2(x) + \langle \nabla f_2(x), y-x \rangle$$

$$= f(x) + \langle \nabla f(x), y-x \rangle.$$



$$\textcircled{3} \quad x = (x_1, x_2)$$

$$f_1(x) = x_1^2$$

$$f_2(x) = x_2^2$$

so, $f(x) = \|x\|^2 = x_1^2 + x_2^2$ is also convex.

$$f(x) = \|x\|_p^p \text{ is also convex}$$

\textcircled{4} If $f(\cdot)$ is convex then $f(Ax+b)$ is also convex.

$$\text{let } g(x) \triangleq f(Ax+b)$$

$$\begin{aligned} g(y) &= \underbrace{f(Ay+b)}_{w_y} \stackrel{f(\cdot) \text{ is convex}}{\geq} \underbrace{f(Ax+b)}_{w_x} + \langle \nabla f \Big|_{Ax+b}, (Ay+b) - (Ax+b) \rangle \\ &= g(x) + \langle \nabla f \Big|_{Ax+b}, A(y-x) \rangle \end{aligned}$$

$$= g(x) + \underbrace{\langle A^T \nabla f \Big|_{Ax+b}, y-x \rangle}_{\text{Mean}}$$

$$= g(x) + \langle \nabla g(x), y-x \rangle.$$

$$\min_x \frac{1}{2} \|Ax-b\|^2.$$

$$f(x) = \underline{\|Ax-b\|^2} \text{ is convex.}$$

$$\text{Find } x \text{ s.t. } \underline{a^T x} \approx b$$

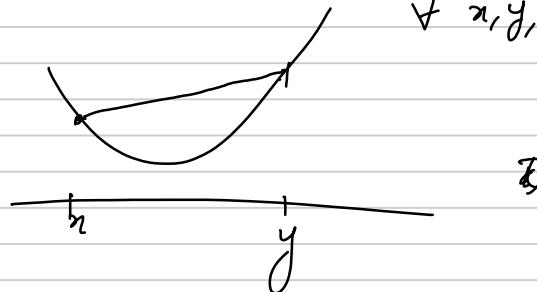
Temp	Humidity	...	Rainfall
a_1	a_2	\dots	b

LEC - 2

$$f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle$$

Lemma : If $f(\cdot)$ is a convex function then $f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$ $\forall x, y, \alpha \in [0, 1]$.

Proof is Exercise.



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

Lemma : If f is convex then $\forall x$ in the domain, \exists non-empty set $G(x)$ s.t. $f(y) \geq f(x) + \langle z, y-x \rangle \forall x, y$ and $z \in G(x)$.
↓
Subgradients

MORALLY : CONVEX = CONVEX.

↓
First defn.

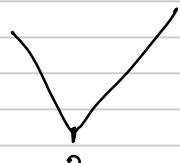
↓
Second defn.

Lasso :

$$\min_x \frac{1}{2} \|Ax-b\|^2 + \|x\|_1$$

$$\left\{ \begin{array}{l} f(x) = |x| \\ \frac{df}{dx} = \begin{cases} 1 & x > 0 \\ -1 & x < 0 \end{cases} \end{array} \right.$$

does not exist at $x=0$.



$$f(\cdot) \text{ is convex} \quad | \quad f(y) \geq f(x) + \underbrace{\langle \nabla f(x), y - x \rangle}_{\downarrow}$$

$$x_{t+1} = x_t - \eta \nabla f(x_t).$$

$$\text{Let } x^* = \arg \min_x f(x).$$

Subgradient

$$\eta_t = \frac{1}{\sqrt{T}}$$

Theorem: If $f(\cdot)$ is G -Lipschitz $\left[\|\nabla f(x)\| \leq G \forall x \right]$, then

$$\frac{1}{T} \sum_{t=1}^T \underbrace{\left[f(x_t) - f(x^*) \right]}_{\text{Suboptimality of } x_t} \leq \frac{\|x_0 - x^*\|^2}{2\eta T} + \frac{\eta G^2}{2}.$$

$$\text{Remarks: } \eta^* = \frac{\|x_0 - x^*\|}{G\sqrt{T}} ; \text{ RHS} = \frac{\|x_0 - x^*\| \cdot G}{\sqrt{T}}.$$

$$\begin{aligned} \text{Proof: } \|x_{t+1} - x^*\|^2 &= \|x_t - \eta \nabla f(x_t) - x^*\|^2 \\ &= \|x_t - x^*\|^2 - 2\eta \underbrace{\langle \nabla f(x_t), x_t - x^* \rangle}_{\leq G^2} + \eta^2 \underbrace{\|\nabla f(x_t)\|^2}_{\leq G^2} \end{aligned}$$

$$\begin{array}{c} y=x^* \\ x=x_t \end{array} \quad \left[f(x^*) \geq f(x_t) + \underbrace{\langle \nabla f(x_t), x^* - x_t \rangle}_{\leq G^2} \right]$$

$$\|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 - 2\eta [f(x_t) - f(x^*)] + \eta^2 G^2.$$

$$\text{Telescopic sum: } 0 \leq \|x_{T+1} - x^*\|^2 \leq \|x_0 - x^*\|^2 - 2\eta \sum_{t=0}^T [f(x_t) - f(x^*)] + \eta^2 G^2 (T+1).$$

Dividing by $2\eta(T+1)$ gives the theorem.

Remarks: ① The Convergence rate is independent of dimension

② In general we do not know G or $\|\mathbf{x}_0 - \mathbf{x}^*\|$.
Need to learn them on the go.

LOWER BOUNDS $\xrightarrow{\text{Gradient Span}}$
 $\xrightarrow{\text{Algorithms}}$ (GSA) : $\mathbf{x}_{t+1} \in \text{Span} \left\{ \mathbf{x}_0, \nabla f(\mathbf{x}_0), \nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_t) \right\}$.

GD \in Gradient Span algorithms.

Theorem: For any value of G, R and T if a function $f_{G,R,T}(\mathbf{x})$ s.t.

i) $\|\nabla f_{G,R,T}(\mathbf{x})\| \leq G$ o) f is convex

ii) $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$ and any vector in the span $\{\mathbf{x}_0, \dots, \mathbf{x}_T\}$

iii) for any GSA $f_{G,R,T}(\bar{\mathbf{x}}) - f_{G,R,T}(\mathbf{x}^*) \geq \frac{GR}{2(1 + \sqrt{T+1})}$.

Remark: Matches upper bound upto constant factors. [GD is optimal for this class of functions]

Proof:

$$f_{\mathbf{x}, \mu}(\mathbf{x}) = \underset{1 \leq i \leq T+1}{\text{min}} \mathbf{x}(i) + \frac{\mu}{2} \|\mathbf{x}\|^2$$

Ex: Show that $f_{\mathbf{x}, \mu}(\cdot)$ is convex.

i) ONLY PROPERTY WE NEED: $f(y) \geq f(x) + \langle z, y-x \rangle \forall x, y$

Any z that satisfies this is a Subgradient of f at x

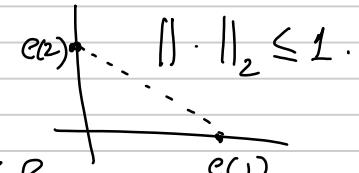
$$f(x) = \max_i f_i(x)$$

$$\text{let } I(x) = \{ i : f_i(x) = f(x) \}.$$

Any vector $z \in \text{Convex hull} \left(\underbrace{\{f_i(x) : i \in I(x)\}}_{\text{Convex hull}} \right)$ is a Subgradient of $f(\cdot)$ at x .

$$\nabla f(x) = \rho \text{Convex hull}(e(i) : i \in I(x)) + \mu z.$$

$$\begin{aligned} \|\nabla f(x)\| &\leq \| \rho \cdot 1 + \mu z \| \\ &\leq \rho \cdot 1 + \mu R \end{aligned}$$

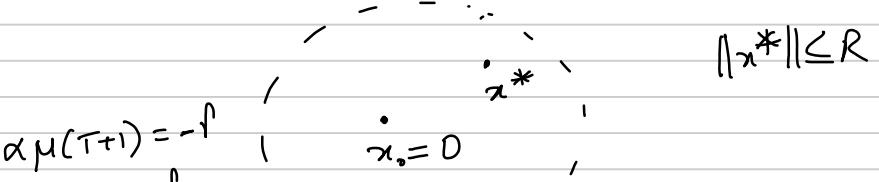


$$\text{i)} \quad f(x) = \rho \max_{1 \leq i \leq T+1} x(i) + \frac{\mu}{2} \|x\|^2$$

$$\textcircled{1} \quad \rho + \mu R \leq G$$

$$x(i) = \alpha$$

$$f(x) = \rho \alpha + \frac{\mu}{2} \alpha^2(T+1) \rightarrow \alpha \mu(T+1) = -\rho$$



$$\textcircled{1} \quad \rho + \mu R \leq G$$

$$x_0 = 0$$

$$\textcircled{2} \quad \|x_0 - x^*\|^2 = \|x^*\|^2$$

$$= \frac{\rho^2}{\mu^2(T+1)} \leq R^2.$$

$$\boxed{f(x^*) = \frac{-\rho}{2\mu(T+1)} \leq 0}$$

$$x^* = -\frac{\rho}{\mu(T+1)} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

Choice : $\rho = \frac{\sqrt{T+1} \cdot G}{1 + \sqrt{T+1}}$ and $\mu = \frac{G}{R(1 + \sqrt{T+1})}$

$$f(x) = \min_{1 \leq i \leq T+1} x(i) + \frac{\mu}{2} \|x\|^2$$

$$x_0 = 0$$

$$\nabla f(x_0) = \text{Conv. hull}(e_i : i \in I^{(0)}) + \mu \cdot 0$$

Resisting oracle: $\nabla f(x) = \sum e_{i(x)} + \mu x$ where $i(x) = \min I(x)$.

$$\begin{aligned} x_0 &= 0 \\ \nabla f(x_0) &= \sum e^{(1)} + \mu \cdot 0 = \sum e^{(1)} \end{aligned}$$

$$\left. \begin{aligned} x_1 &\in \text{Span} \left\{ \frac{x_0}{\|x_0\|}, \frac{\nabla f(x_0)}{\|\nabla f(x_0)\|} \right\} = \text{Span}(e^{(1)}) \\ \alpha e_1 &\in \text{Span}(e^{(1)}) \end{aligned} \right\}$$

Let $x_1 = \alpha e_1$. If $\alpha \geq 0$ then $\nabla f(x_1) = \sum e^{(1)} + \mu x_1$
 If $\alpha < 0$ then $\nabla f(x_1) = \sum e^{(2)} + \mu x_1$.

$$f_1(x) \xrightarrow{\substack{\nabla f(x) \\ \nabla f(x)}} \nabla f(x_1) \in \text{Span}\{e^{(1)}, e^{(2)}\}.$$

$$f(x) = \min \underbrace{\{x^{(1)}, x^{(2)}, \dots, x^{(T+1)}\}}_{=x} + \frac{\mu}{2} \|x\|^2.$$

$$x_1 = \alpha e_1 = \begin{bmatrix} \alpha \\ 0 \end{bmatrix}.$$

If $\alpha > 0$ then $I(x_1) = 1 \rightarrow \sum e_1$

If $\alpha = 0$ then $I(x_1) = 1, 2 \rightarrow \sum e_1$

If $\alpha < 0$ then $I(x_1) = 2 \rightarrow \sum e_2$

$$x_t \in \text{Span}\{e^{(1)}, \dots, e^{(t)}\}.$$

$$f(x_t) = \min_{1 \leq i \leq T+1} x_t^{(i)} + \frac{\mu}{2} \|x_t\|^2$$

$$x_t^{(T+1)} = 0$$

$$\text{So, } f(x_t) \geq 0.$$

$$\gamma = \frac{\sqrt{\tau+1} \cdot G}{1 + \sqrt{\tau+1}} ; \quad \mu = \frac{G}{R(1 + \sqrt{\tau+1})} ; \quad f(x_{\tau}) \geq 0.$$

$$f(x^*) = -\frac{\gamma^2}{2\mu(\tau+1)}$$

$$f(x_{\tau}) - f(x^*) \geq \frac{\gamma^2}{2\mu(\tau+1)} = \frac{(\tau+1)G^2}{(1+\sqrt{\tau+1})^2} \cdot \frac{R(1+\sqrt{\tau+1})}{2G} \cdot \frac{1}{\tau+1}$$

$$= \frac{GR}{2(1+\sqrt{\tau+1})}.$$

□

$$f(x) = \min \left\{ f_1(x), f_2(x) \right\}.$$

LECTURE - 3

Note Title

26-Jun-19

$$\begin{aligned} f \text{ is Convex} \\ \|\nabla f(x)\| \leq G \\ \|x_0 - x^*\| \leq R \end{aligned}$$

$$GD: \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{GR}{\sqrt{T}}$$

$$\text{Running example : } f(x) = \frac{1}{2} \|Ax - b\|^2.$$

$$x_0 = 0 \xrightarrow{\text{Sub.Gr.D.e.}}$$

$$R \triangleq \|x_0 - x^*\| = \|x^*\|$$

$$G \triangleq \underset{x: \|x\| \leq R}{\text{max}} A^T(Ax - b) \rightarrow G = \|A\|^2 \cdot R + \|A\| \cdot \|b\|$$

$$\frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{\epsilon R}{\sqrt{T}} = \frac{\|A\|^2 \cdot R^2 + \|A\| \cdot \|b\| \cdot R}{\sqrt{T}}.$$

Qn.: Can we do better?

↪ Say for Smooth functions.

Smoothness: $\|\nabla f(x) - \nabla f(y)\| \leq L \cdot \|x - y\|$.

↓

$$\|A^T(Ax - b) - A^T(Ay - b)\| = \|A^T A(x - y)\|$$

$\therefore L \triangleq \|A^T A\|$.

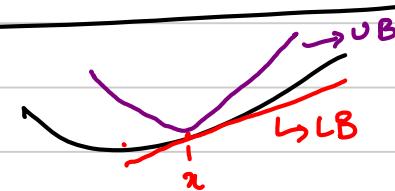
This lecture: ① Convergence rate of GD for Smooth Convex fun.

② Optimal method: Nesterov's accelerated gradient.

$$GD: \boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \eta \nabla f(\boldsymbol{x}_t)$$

$$\text{Convexity: } f(\boldsymbol{y}) \geq f(\boldsymbol{x}) + \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle$$

$$\text{Smoothness: } f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \langle \nabla f(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{\mu}{2} \|\boldsymbol{x} - \boldsymbol{y}\|^2$$



Thm: If f is an L -smooth convex fn. then GD with

$$\eta = \boxed{\frac{1}{2L}}, \text{ we have } f(x_T) - f(x^*) \leq \underbrace{\frac{C \cdot L \cdot \|x_0 - x^*\|^2}{T}}_{\rightarrow \text{Final iterate}}.$$

Lemma : If $f(\cdot)$ is L -Smooth then $\|\nabla f(x)\|^2 \leq 2L \cdot [f(x) - f(x^*)]$.

$$\text{Proof} : f(x^*) \leq f\left(x - \frac{1}{L} \nabla f(x)\right)$$

$$\begin{aligned} &\leq f(x) + \langle \nabla f(x), -\frac{1}{L} \nabla f(x) \rangle + \frac{L}{2} \left\| \frac{1}{L} \nabla f(x) \right\|^2 \\ &= f(x) - \frac{1}{2L} \left\| \nabla f(x) \right\|^2. \end{aligned}$$

□

Proof of thm :

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

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

$$= \|x_t - x^*\|^2 - 2\eta \underbrace{\langle \nabla f(x_t), x_t - x^* \rangle}_{\text{Convexity}} + \eta^2 \underbrace{\|\nabla f(x_t)\|^2}_{\substack{\text{Use Lemma} \\ \text{above}}}.$$

$$\leq \|x_t - x^*\|^2 - 2\eta \underbrace{\langle \nabla f(x_t), x_t - x^* \rangle}_{\text{Convexity}} + \eta^L \cdot 2L [f(x_t) - f(x^*)]$$

$$\leq \|x_t - x^*\|^2 - 2\eta [f(x_t) - f(x^*)] + 2\eta^2 L [f(x_t) - f(x^*)]$$

$$\|\boldsymbol{x}_{t+1} - \boldsymbol{x}^*\|^2 \leq \|\boldsymbol{x}_t - \boldsymbol{x}^*\|^2 - 2\eta(1-\eta L) [f(\boldsymbol{x}_t) - f(\boldsymbol{x}^*)].$$

Take telescopic Sum,

$$\begin{aligned} \|\boldsymbol{x}_{T+1} - \boldsymbol{x}^*\|^2 &\leq \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2 - 2\eta(1-\eta L) \sum_{t=0}^T [f(\boldsymbol{x}_t) - f(\boldsymbol{x}^*)] \\ \frac{1}{T+1} \sum_{t=0}^T [f(\boldsymbol{x}_t) - f(\boldsymbol{x}^*)] &\leq \frac{\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2 - \|\boldsymbol{x}_{T+1} - \boldsymbol{x}^*\|^2}{2\eta(1-\eta L) \cdot (T+1)} + \text{ve} \end{aligned}$$

$$\text{Choose } \eta = \frac{1}{2L} : \quad \leq \frac{2L\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2}{T+1} .$$

Exercise: For L -smooth, Convex $f(\cdot)$, Show that GD with

$$\eta \leq \frac{1}{2L} \text{ satisfies } f(x_{t+1}) \leq f(x_t) \quad \forall t.$$

Nesterov's accelerated gradient

$$x_{t+1} = x_t - \eta \nabla f(x_t).$$

$$x_{t+1} \triangleq \arg \min_x \left[f(x_t) + \underbrace{\langle \nabla f(x_t), x - x_t \rangle}_{\text{First order Taylor exp.}} + \underbrace{\frac{1}{2\eta} \|x - x_t\|^2}_{\text{Quadratic term}} \right]$$

Local upper bound when $\eta \leq \frac{1}{L}$.

$$f(x) \leq f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \underbrace{\frac{L}{2} \|x - x_t\|^2}_{\text{If } L \leq \frac{1}{\eta} \text{ then}} \leq \frac{1}{2\eta} \|x - x_t\|^2.$$

Estimate
Sequences

$$\phi_0(x) = \underbrace{f(v_0)}_{f(v_0)} + \frac{L}{2} \|x - v_0\|^2.$$

$$\phi_{t+1}(x) = (1 - \alpha_t) \phi_t(x) + \alpha_t \underbrace{\left[f(y_t) + \langle \nabla f(y_t), x - y_t \rangle \right]}_{\leq f(x)} \leq f(x)$$

Lemma: If $\exists x_t$ s.t. $f(x_t) \leq \min_x \phi_t(x)$ then
 $f(x_t) - f(x^*) \leq \underbrace{\left[\prod_{s=0}^{t-1} (1-\alpha_s) \right]}_{\text{Hypothesis}} [\phi_0(x^*) - f(x^*)].$

Proof : $f(x_t) - f(x^*) \stackrel{\text{Hypothesis}}{\leq} \phi_t^* - f(x^*)$

$$\leq \frac{\phi_t(x^*) - f(x^*)}{(1-\alpha_{t-1}) \phi_{t-1}(x^*) + \alpha_{t-1} [f(y_{t-1}) + \langle \nabla f(y_{t-1}), x^* - y_{t-1} \rangle] - (1-\alpha_{t-1}) f(x^*) - \alpha_{t-1} [f(x^*)]}$$

-ve

$$\begin{aligned}
 \phi_t(x^*) - f(x^*) &\leq (1-\alpha_{t-1}) [\phi_{t-1}(x^*) - f(x^*)] \\
 &\leq \prod_{s=0}^{t-1} (1-\alpha_s) [\phi_0(x^*) - f(x^*)].
 \end{aligned} \tag{11}$$

Given : x_t s.t. $f(x_t) \leq \min_x \phi_t(x)$.

Also given ϕ_t

Task : ① Choose y_t and query $\nabla f(y_t)$.

② choose $\alpha_t \longrightarrow \phi_{t+1}$

③ Find x_{t+1} s.t. $f(x_{t+1}) \leq \min_x \phi_{t+1}(x)$.

Lemma ; If $\phi_{t+1}(x) = (1-\alpha_t) \phi_t(x) + \alpha_t [f(y_t) + \langle \nabla f(y_t), x - y_t \rangle]$

then $\phi_{t+1}^* \triangleq \min_x \phi_{t+1}(x) =$ _____

and $v_{t+1} \triangleq \arg \min_x \phi_{t+1}(x) = v_t - \frac{\alpha_t}{\lambda_t(1-\alpha_t)L} \nabla f(y_t)$

Proof :

$$\phi_t^* \triangleq \min_x \phi_t(x)$$

$$v_t \triangleq \arg \min_x \phi_t(x)$$

$$\lambda_t \triangleq \prod_{s=0}^{t-1} (1-\alpha_s)$$

$$\begin{aligned} \phi_0(x) &= \phi_0^* + \frac{L}{2} \|x - v_0\|^2 \\ \phi_1(x) &= (1-\alpha_0) \phi_0(x) \\ &\quad + \alpha_0 [a + \langle b, x \rangle] \end{aligned}$$

$$\phi_t(x) = \phi_t^* + \frac{\alpha_t L}{2} \|x - v_t\|^2.$$

$$\phi_{t+1}(x) = (1-\alpha_t) \phi_t(x) + \alpha_t [f(y_t) + \langle \nabla f(y_t), x - y_t \rangle]$$

$$= (1-\alpha_t) \phi_t^* + \frac{(1-\alpha_t) \alpha_t L}{2} \|x - v_t\|^2$$

$$+ \alpha_t f(y_t) + \alpha_t \langle \nabla f(y_t), x - y_t \rangle.$$

$$= \frac{(1-\alpha_t) \alpha_t L}{2} \left[\|x\|^2 - 2 \langle v_t, x \rangle + \|v_t\|^2 \right]$$

Constant

Quadratic

Linear

$$+ \alpha_t \underbrace{\langle \nabla f(y_t), \alpha \rangle}$$

$$+ (1-\alpha_t) \phi_t^* + \alpha_t f(y_t)$$

$$- \alpha_t \langle \nabla f(y_t), y_t \rangle$$

$$= \frac{(1-\alpha_t)\alpha_t L}{2} \left[\|\alpha\|^2 - 2 \left\langle v_t - \frac{\alpha_t}{(1-\alpha_t)\alpha_t L} \nabla f(y_t), \alpha \right\rangle \right]$$

$$+ \left\| v_t - \frac{\alpha_t}{(1-\alpha_t)\alpha_t L} \nabla f(y_t) \right\|^2$$

Extra terms

$$\left\{ \begin{array}{l} + \frac{2\alpha_t}{(1-\alpha_t)\alpha_t L} \langle \nabla f(y_t), v_t \rangle \\ - \left(\frac{\alpha_t}{(1-\alpha_t)\alpha_t L} \right)^2 \|\nabla f(y_t)\|^2 \end{array} \right\} + \hat{\theta}$$

$$= \frac{(1-\alpha_t) \gamma_t L}{2} \| x - v_t + \underbrace{\frac{\alpha_t}{(1-\alpha_t)\gamma_t L} \nabla f(y_t)}_{v_{t+1}} \|^2 + \text{Extra terms} + \hat{\theta}$$

$$\phi_{t+1}^* = (1-\alpha_t) \phi_t^* + \alpha_t f(y_t) - \frac{\alpha_t^2}{2\gamma_t(1-\alpha_t)L} \|\nabla f(y_t)\|^2 + \alpha_t \langle \nabla f(y_t), v_t - y_t \rangle.$$

How to find v_{t+1} s.t. $f(x_{t+1}) \leq \phi_{t+1}^*$?

IV

By smoothness: $\underbrace{f(y_t - \eta \nabla f(y_t))}_{x_{t+1}} \leq f(y_t) - \eta \|\nabla f(y_t)\|^2 + \frac{\eta^2 L}{2} \|\nabla f(y_t)\|^2$

$$\leq f(y_t) - \frac{1}{2L} \|\nabla f(y_t)\|^2$$

$$f(x_{t+1}) - \phi_{t+1}^* \leq \underbrace{f(y_t)}_{\text{red}} - \frac{1}{2L} \|\nabla f(y_t)\|^2$$

$$- (1-\alpha_t) \phi_t^* - \alpha_t \underbrace{f(y_t)}_{\text{green}} + \frac{\alpha_t^2}{2\alpha_t(1-\alpha_t)L} \|\nabla f(y_t)\|^2$$

$$- \alpha_t \langle \nabla f(y_t), v_t - y_t \rangle$$

$$\leq (1-\alpha_t) f(y_t) - (1-\alpha_t) \underline{f(x_t)}$$

$$- \alpha_t \langle \nabla f(y_t), v_t - y_t \rangle$$

(Convexity: $x = y_t$
 $y = x_t$)

$$\leq \langle \nabla f(y_t), (1-\alpha_t)(y_t - x_t) \rangle - \alpha_t \langle \nabla f(y_t), v_t - y_t \rangle$$

$$= \langle \nabla f(y_t), y_t - (1-\alpha_t)x_t - \alpha_t v_t \rangle$$

$$\leq 0$$

Want: $\frac{\alpha_t^2}{\alpha_t(1-\alpha_t)} \leq 1$

Requirements

$$\frac{\alpha_t^2}{2\alpha_t(1-\alpha_t)L} \leq \frac{1}{2L}$$

Algorithm

$$y_t = (1-\alpha_t)x_t + \alpha_t v_t$$

$$x_{t+1} = y_t - \frac{1}{L} \nabla f(y_t)$$

v_{t+1} = Update from Lemma.

and minimize $\alpha_t = \prod_{s=0}^{t-1} (1-\alpha_s)$.

Can choose α_t : $\boxed{\alpha_{t+1} \leq \frac{4}{t^2}}$

$$\begin{aligned}
 f(x_T) - f(x^*) &\leq \alpha_T [\phi_0(x^*) - f(x^*)] \\
 &\leq \frac{4}{(T-1)^2} \left[f(v_0) + \frac{L}{2} \|x^* - v_0\|^2 - f(x^*) \right] \\
 &\leq \frac{4}{(T-1)^2} [L \|x^* - v_0\|^2].
 \end{aligned}$$

Comparing with GD: $\frac{1}{T^2}$ instead of $\frac{1}{T}$.

LECTURE - 4

$$f(\alpha) = \frac{1}{2} \|A\alpha - b\|^2 ; \quad L = \|A\|^2$$

$$\text{AGD} \quad f(\alpha_T) - f(\alpha^*) \leq \frac{6 \|A\|^2 \cdot \|\alpha_0 - \alpha^*\|^2}{T^2}$$

Can we do better?

\hookrightarrow Strong convexity.

Strong Convexity : f is μ -Strongly convex

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

Linear regression : $\hat{y} = \text{argmin } (AA^T)$.

$$\frac{1}{2} \|Ax - b\|^2$$

$$GD: \quad x_{t+1} = x_t - \eta \nabla f(x_t).$$

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

$$= \|x_t - x^*\|^2 - 2\eta \underbrace{\langle \nabla f(x_t), x_t - x^* \rangle}_{\text{Strong Convexity}} + \eta^2 \|\nabla f(x_t)\|^2$$

$$f(x^*) \geq f(x_t) + \underbrace{\langle \nabla f(x_t), x^* - x_t \rangle}_{\text{Smoothness}} + \frac{\eta}{2} \|x_t - x^*\|^2$$

$$\begin{aligned}
 \|\gamma_{t+1} - x^*\|^2 &\leq \|x_t - x^*\|^2 - 2\eta \left[f(x_t) - f(x^*) + \frac{\mu}{2} \|x_t - x^*\|^2 \right] \\
 &\quad + \eta^2 \cdot 2L [f(x_t) - f(x^*)] \\
 &\leq (1 - \eta\mu) \|x_t - x^*\|^2 - \underline{2\eta(1-\eta\mu)[f(x_t) - f(x^*)]}
 \end{aligned}$$

If we choose $\boxed{\eta \leq \frac{1}{L}}$ then ≤ 0

$$\begin{aligned}
 \|x_{t+1} - x^*\|^2 &\leq (1 - \eta\mu) \|x_t - x^*\|^2 \\
 &\leq (1 - \eta\mu)^{t+1} \|x_0 - x^*\|^2.
 \end{aligned}$$

$$\begin{aligned}
 f(x_t) - f(x^*) &\leq \frac{\mu}{2} \|x_t - x^*\|^2 \\
 &\leq \frac{\mu}{2} \cdot (1-\eta\mu)^t \|x_0 - x^*\|^2. \\
 &\leq \frac{\mu \cdot \|x_0 - x^*\|^2}{2} \cdot e^{-\eta\mu t}.
 \end{aligned}$$

HWER

FACT: $1 - \omega \leq e^{-\omega}$

Choose η as large as possible

$$\text{so, } \eta = \frac{1}{L}.$$

$$f(x_t) - f(x^*) \leq \frac{\mu \|x_0 - x^*\|^2}{2} \cdot e^{-\frac{\mu}{L} t}.$$

condition number
 $n = \frac{L}{\mu} \geq 1$

$$f(x_t) - f(x^*) \leq \frac{\mu \|x_0 - x^*\|^2}{2} e^{-\frac{t}{\mu}}$$

Acceleration for Strongly Convex functions

Smooth and

Is this the best possible rate?

Nesterov's AGD : $f(x_t) - f(x^*) \leq \frac{\mu \|x_0 - x^*\|^2}{2t^2}$

$$\leq \frac{\mu}{\mu} \cdot [f(x_0) - f(x^*)] \cdot \frac{1}{t^2} + \frac{\mu}{2} \|x_0 - x^*\|^2$$

$$= \frac{1}{t^2} [f(x_0) - f(x^*)].$$

$$\frac{\|x_0 - x^*\|^2}{t^2} \leq \frac{2}{\mu} [f(x_0) - f(x^*)]$$

Choosing $t = \sqrt{2K}$ iterations,

$$f(x_{\sqrt{2K}}) - f(x^*) \leq \frac{1}{2} [f(x_0) - f(x^*)]$$

$$f(x_{2\sqrt{2K}}) - f(x^*) \leq \frac{1}{2} [f(x_{\sqrt{2K}}) - f(x^*)]$$

$$f(x_{i\sqrt{2K}}) - f(x^*) \leq \frac{1}{2} [f(x_{(i-1)\sqrt{2K}}) - f(x^*)]$$

$$f(x_t) - f(x^*) \leq \left(\frac{1}{2}\right)^{(t/\sqrt{2k})} [f(x_0) - f(x^*)].$$

$$\text{GD: } e^{-t/k} \quad ; \quad \text{AGD: } e^{-t/\sqrt{k}}$$

CLASS	Algorithm	Guarantee
$\ x_0 - x^*\ \leq R$ $\ \nabla f\ \leq L$	SUBGRADIENT DESCENT	$\frac{LR}{\sqrt{T}}$
L-SMOOTH $\ \nabla f(x) - \nabla f(y)\ \leq L\ x-y\ $	G.D. A.G.D.	$\frac{LR^2}{T}$ $\frac{LR^2}{T^2}$
L-SMOOTH μ -Str. Conv.	G.D. A.G.D.	$LR^2 \cdot \exp(-\frac{\mu T}{L})$ $LR^2 \cdot \exp(-\sqrt{\frac{\mu}{L}} \cdot T)$

\mathcal{L} -Lipschitz
 H -Str. conv.

Sub. $\mathcal{L} \cdot D$

$$\frac{\mathcal{L}^2}{\mu T}$$

CONSTRAINED OPTIM.

$$\min_{x \in X} f(x)$$

X is a simple, convex set.
closed

We have access to projection operator

$$\text{given } x, \quad P_X(x) = \operatorname{argmin}_{y \in X} \|x - y\|^2.$$

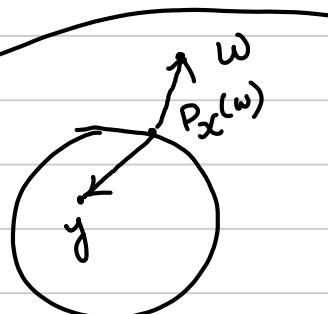
$X = B_1(\omega)$ is a simple set

Pythagoras theorem: $\forall \omega$, we have

EXERCISE

i) $\langle \omega - P_X(\omega), y - P_X(\omega) \rangle \leq 0$

ii) $\|P_X(\omega) - y\|^2 \leq \|\omega - y\|^2 \quad \boxed{\forall y \in X}$



Projected GD: $x_{t+1} = P_X(x_t - \eta \nabla f(x_t))$

$\xrightarrow{\text{Equivalent}}$

$$x_{t+1} = \arg \min_{x \in X} \left[f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|^2 \right].$$

Linear

$\xrightarrow{\text{Str. Conv.}}$ $\xrightarrow{\text{quadric}}$

Exercise: $\hat{f}_\eta(x_t; x)$ is $\frac{1}{\eta}$ -Strongly Convex

Doing now
Conv. rate of PGD
for L-Smooth fun.

① $\rightarrow \hat{f}_\eta(x_t; x) \geq \hat{f}_\eta(x_t; x_{t+1}) + \frac{1}{2\eta} \|x - x_{t+1}\|^2$

$\xrightarrow{\text{use Pythag. thm.}} (\because x_{t+1} \text{ is argmin } \hat{f}_\eta(x_t; x))$

For $\eta \leq \frac{1}{L}$; $f(x) \leq \hat{f}_\eta(x_t; x)$ [by L-smoothness]
 $\text{g } f(\cdot)$

$$\begin{aligned}
 f(x) &\stackrel{\text{Conv.}}{\geq} f(x_t) + \langle \nabla f(x_t), x - x_t \rangle \\
 &= \hat{f}_\eta(x_t; x) - \frac{1}{2\eta} \|x - x_t\|^2 \\
 &\stackrel{(1)}{\geq} \hat{f}_\eta(x_t; x_{t+1}) + \frac{1}{2\eta} \|x - x_{t+1}\|^2 - \frac{1}{2\eta} \|x - x_t\|^2 \\
 &\stackrel{\text{Smoothness}}{\geq} f(x_{t+1}) + \frac{1}{2\eta} \|x - x_{t+1}\|^2 - \frac{1}{2\eta} \|x - x_t\|^2
 \end{aligned}$$

$$\|\boldsymbol{x} - \boldsymbol{x}_{t+1}\|^2 \leq \|\boldsymbol{x} - \boldsymbol{x}_t\|^2 - 2\eta [f(\boldsymbol{x}_{t+1}) - f(\boldsymbol{x})]$$

$\boldsymbol{x} = \boldsymbol{x}^*$ and iterate

$$\begin{aligned} \|\boldsymbol{x}^* - \boldsymbol{x}_{t+1}\|^2 &\leq \|\boldsymbol{x}^* - \boldsymbol{x}_0\|^2 - 2\eta \underbrace{\sum_{s=1}^{t+1} [f(\boldsymbol{x}_s) - f(\boldsymbol{x}^*)]}_{\frac{1}{2\eta(t+1)} [\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2 - \|\boldsymbol{x}^* - \boldsymbol{x}_{t+1}\|^2]}. \\ \frac{1}{t+1} \sum_{s=1}^{t+1} [f(\boldsymbol{x}_s) - f(\boldsymbol{x}^*)] &\leq \frac{\|\boldsymbol{x}^* - \boldsymbol{x}_0\|^2 - \|\boldsymbol{x}^* - \boldsymbol{x}_{t+1}\|^2}{2\eta(t+1)}. \\ [\eta = \frac{1}{L}] &\leq \frac{\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2}{2(t+1)}. \quad \text{④} \end{aligned}$$

Ref:

INTRODUCTORY LECTURES ON CONVEX OPT.
— YURII NESTEROV

LECTURE-5

Note Title

27-Jun-19

Faster algorithms (better than the black-box bounds)

for structured non-smooth problems

$$f(x) = \underbrace{g(x)}_{\text{Convex \& Smooth}} + \underbrace{h(x)}_{\text{Convex \& non-smooth}}$$

perhaps
non-smooth

Access to : $\arg \min_x \langle \omega, x \rangle + h(x) + \frac{1}{2\eta} \|x - y\|^2$
 $\forall \omega, \eta, y.$

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

GD
with
Prox

Algorithm:

$$\begin{aligned} x_{t+1} = \arg \min_x & g(x_t) + \langle \nabla g(x_t), x - x_t \rangle \\ & + h(x) + \frac{1}{2\eta} \|x - x_t\|^2. \end{aligned}$$

$$\begin{aligned} \text{GD : } x_{t+1} = \arg \min_x & g(x_t) + \langle \nabla g(x_t), x - x_t \rangle \\ & + h(x_t) + \langle \nabla h(x_t), x - x_t \rangle \\ & + \frac{1}{2\eta} \|x - x_t\|^2. \end{aligned}$$

(Compressed
Sensing)

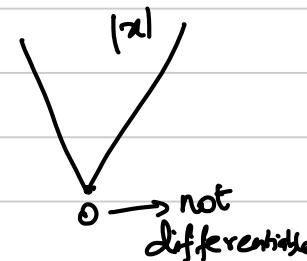
$$\text{LASSO: } f(x) = \frac{1}{2} \|Ax - b\|^2 + \alpha \|x\|_1$$

$\downarrow g(x)$

Linear regression

$\downarrow h(x)$

encourage sparsity
of x .



$$x_{t+1} = \arg \min_x \cancel{g(x_t)} + \langle \nabla g(x_t), x - x_t \rangle + h(x) + \frac{1}{2\eta} \|x - x_t\|^2$$

$$= \arg \min_x \langle \nabla g(x_t), x \rangle + h(x) + \frac{1}{2\eta} \|x - x_t\|^2$$

$\downarrow \frac{1}{2} \|x\|_1$

$$= \arg \min_{\boldsymbol{x}} \sum_{i=1}^d x_i \nabla g(\boldsymbol{x}_t)_i + |x_i| + \frac{1}{2\eta} (x_i - (\boldsymbol{x}_t)_i)^2$$

$$(\boldsymbol{x}_{t+1})_i = \arg \min_{x_i} x_i \nabla g(\boldsymbol{x}_t)_i + |x_i| + \frac{1}{2\eta} (x_i - (\boldsymbol{x}_t)_i)^2$$

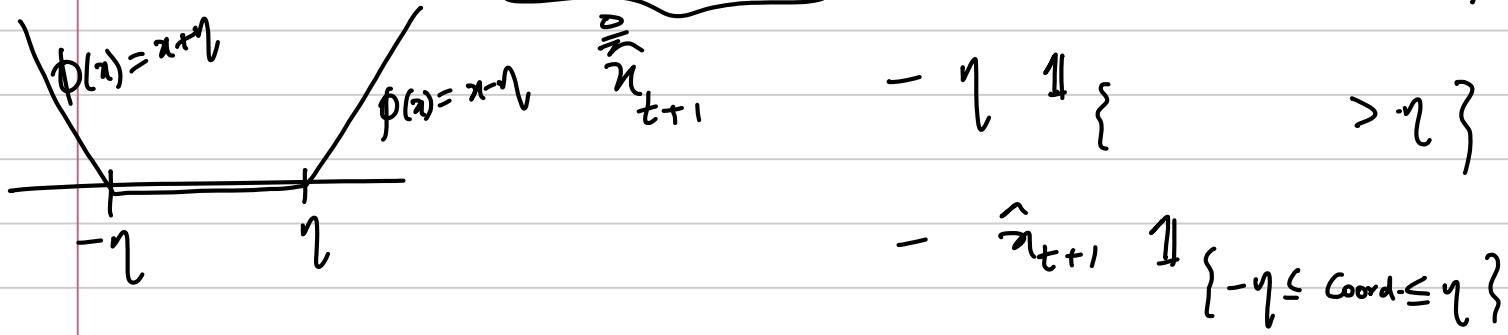
$$= \arg \min_{x_i} \frac{1}{2\eta} \underbrace{(x_i - (\boldsymbol{x}_t)_i + \eta \nabla g(\boldsymbol{x}_t)_i)^2}_{=} + |x_i|.$$

Exercise:

$$\begin{aligned}
 (\boldsymbol{x}_{t+1})_i &= (\boldsymbol{x}_t)_i - \eta \nabla g(\boldsymbol{x}_t)_i - \eta && \text{if } (\boldsymbol{x}_t)_i - \eta \nabla g(\boldsymbol{x}_t)_i \geq \eta \\
 &= \text{,,} && \text{if } \text{,,} \leq -\eta \\
 &= 0 && \text{otherwise.}
 \end{aligned}$$

$= \phi(x_t - \eta \nabla g(x_t))$, where ϕ is applied to each coordinate.

$$x_{t+1} = x_t - \eta \nabla g(x_t) + \eta \mathbb{1}_{\{\text{Coordinates} < -\eta\}}$$



$$- \eta \mathbb{1}_{\{>\eta\}} - \hat{x}_{t+1} \mathbb{1}_{\{-\eta \leq \text{coord.} \leq \eta\}}$$

$$\hat{f}_\eta(x_t; x) \triangleq \underbrace{g(x_t)}_{\text{ }} + \underbrace{\langle \nabla g(x_t), x - x_t \rangle}_{\text{ }} + \underbrace{h(x)}_{\text{ }} + \underbrace{\frac{1}{2\eta} \|x - x_t\|^2}_{\text{ }}$$

$$\begin{aligned}
 & \boxed{x^* = \underset{x}{\operatorname{argmin}} f(x)} \quad x_{t+1} = \underset{x}{\operatorname{argmin}} \hat{f}_\eta(x_t; x) . \\
 & \boxed{\begin{array}{l} g \text{ is } L\text{-Smooth} \\ \eta \leq \frac{1}{L} \end{array}}
 \end{aligned}$$

$$\underline{f(x^*)} = g(x^*) + h(x^*)$$

$$\begin{aligned}
 & \text{Conv. } \underline{f(g(\cdot))} \geq g(x_t) + \langle \nabla g(x_t), x^* - x_t \rangle + h(x^*) \\
 & \quad + \frac{1}{2\eta} \|x^* - x_t\|^2 - \frac{1}{2\eta} \|x^* - x_t\|^2
 \end{aligned}$$

$$\begin{aligned}
 & \left[\begin{array}{l} x_{t+1} = \underset{x}{\operatorname{argmin}} \hat{f}_\eta \\ \text{& } \hat{f}_\eta - \frac{1}{\eta} \text{ str. conv.} \end{array} \right] \Rightarrow \\
 & \quad = \hat{f}_\eta(x_t; x^*) - \frac{1}{2\eta} \|x^* - x_t\|^2 \\
 & \quad \geq \hat{f}_\eta(x_t; x_{t+1}) + \frac{1}{2\eta} \|x^* - x_{t+1}\|^2 - \frac{1}{2\eta} \|x^* - x_t\|^2
 \end{aligned}$$

$$\hat{f}_\eta(x_t; x) = \underbrace{g(x_t) + \langle \nabla g(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|^2}_{\text{L-smoothness of } g(\cdot)} + h(x)$$

$\left[\begin{array}{l} \text{L-smoothness} \\ \text{& } \eta \leq \frac{1}{L} \end{array} \right] \Rightarrow g(x) + h(x) = f(x) \cdot \quad \forall x.$

$x_t \quad x$
 $\uparrow \quad \uparrow$

smoothness lemma: $g(y) \leq g(x) + \langle \nabla g(x), y - x \rangle + \frac{L}{2} \|x - y\|^2$

$$f(x^*) \geq f(x_{t+1}) + \frac{1}{2\eta} \|x^* - x_{t+1}\|^2 - \frac{1}{2\eta} \|x^* - x_t\|^2.$$

Rearranging,

$$\|\boldsymbol{x}_{t+1} - \boldsymbol{x}^*\|^2 \leq \|\boldsymbol{x}_t - \boldsymbol{x}^*\|^2 - 2\eta [f(\boldsymbol{x}_{t+1}) - f(\boldsymbol{x}^*)].$$

$$[\text{Telescope}] \leftarrow \leq \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2 - 2\eta \sum_{s=1}^{t+1} [f(\boldsymbol{x}_s) - f(\boldsymbol{x}^*)].$$

$$\frac{1}{t+1} \sum_{s=1}^{t+1} [f(\boldsymbol{x}_s) - f(\boldsymbol{x}^*)] \leq \frac{\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2 - \|\boldsymbol{x}_{t+1} - \boldsymbol{x}^*\|^2}{2\eta(t+1)}.$$

$\eta = \frac{L}{2}$

$$\leq \frac{L\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2}{2(t+1)}.$$

When the non-smooth Component is Simple, we can get $O\left(\frac{1}{T}\right)$ Convergence rate. [ISTA]

Iterative Shrinkage
Thresholding Algorithm

We can in fact improve this to $O\left(\frac{1}{T^2}\right)$ using

Nesterov's AGD. [FISTA]

STOCHASTIC ALGORITHMS

$$1) \quad f(x) = \mathbb{E}_{(a,b)} [(a^T x - b)^2]$$

$$\nabla f(x) = \mathbb{E}_{(a,b)} [(a^T x - b)a].$$

Oracle : (a_i, b_i)

$$\hat{\nabla} f(x) \triangleq (a_i^T x - b_i) a_i \Rightarrow \mathbb{E}[\hat{\nabla} f(x)] = \underline{\nabla f(x)}.$$

② Empirical Risk Minimization (ERM)

$$f(x) = \frac{1}{2n} \sum_{i=1}^n (a_i^T x - b_i)^2$$

If $a_i \in \mathbb{R}^d$

$O(nd)$ time $\leftarrow \nabla f(x) = \frac{1}{n} \sum_{i=1}^n (a_i^T x - b_i) a_i$

Old time $\leftarrow \hat{\nabla} f(x) \triangleq (a_i^T x - b_i) a_i$ where $i \sim \text{Unif}[1, \dots, n]$.

Setting:

$$x \rightarrow \boxed{\text{oracle}} \rightarrow \hat{\nabla} f(x)$$

$\hat{\nabla} f(x)$ independent of everything else $\mathbb{E}[\hat{\nabla} f(x)] = \nabla f(x)$; $\mathbb{E}[\|\hat{\nabla} f(x) - \nabla f(x)\|^2] \leq \sigma^2$.

SGD [Robbins & Monro]

$$\|\nabla f(\mathbf{x})\| \leq L$$

f : Convex & L -Lipschitz
Random noise has bounded variance

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \hat{\nabla} f(\mathbf{x}_t)$$

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2\right] &= \mathbb{E}\left[\|\mathbf{x}_t - \mathbf{x}^* - \eta \hat{\nabla} f(\mathbf{x}_t)\|^2\right] \\ &= \mathbb{E}\left[\|\mathbf{x}_t - \mathbf{x}^*\|^2 - 2\eta \langle \hat{\nabla} f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle + \eta^2 \|\hat{\nabla} f(\mathbf{x}_t)\|^2\right] \\ &= \mathbb{E}\left[\|\mathbf{x}_t - \mathbf{x}^*\|^2\right] - 2\eta \underbrace{\mathbb{E}\left[\langle \hat{\nabla} f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle\right]}_{+ \eta^2 \underbrace{\mathbb{E}\left[\|\hat{\nabla} f(\mathbf{x}_t)\|^2\right]}} \end{aligned}$$

$$1) \mathbb{E} \left[\langle \hat{\nabla} f(\pi_t), \pi_t - \pi^* \rangle \middle| \pi_t \right] = \langle \mathbb{E} \left[\hat{\nabla} f(\pi_t) \middle| \pi_t \right], \pi_t - \pi^* \rangle$$

$$= \langle \nabla f(\pi_t), \pi_t - \pi^* \rangle$$

[Convexity] $\Leftrightarrow f(\pi_t) - f(\pi^*) \leq \hat{f}(\pi_t) - \mathbb{E}[f(\pi_t)]$

$$2) \mathbb{E} \left[\langle \hat{\nabla} f(\pi_t), \hat{\nabla} f(\pi_t) \rangle \middle| \pi_t \right] = \mathbb{E} \left[\langle \nabla f(\pi_t) + \eta_t, \nabla f(\pi_t) + \eta_t \rangle \right]$$

$$= \underbrace{\mathbb{E} [\|\nabla f(\pi_t)\|^2]}_{\text{Since } f(\cdot) \text{ is } L\text{-Lipschitz}} + \mathbb{E} [\|\eta_t\|^2]$$

$$\leq \frac{\sigma^2}{\text{Since noise variance } \leq \sigma^2}$$

$$\text{So, } \mathbb{E}[|\boldsymbol{x}_{t+1} - \boldsymbol{x}^*|^2] \leq \mathbb{E}[|\boldsymbol{x}_t - \boldsymbol{x}^*|^2] \\ - 2\eta [\mathbb{E}[f(\boldsymbol{x}_t)] - f(\boldsymbol{x}^*)] \\ + \eta^2 (\zeta^2 + \sigma^2).$$

Using the above
inequality \leq
iteratively

$$\mathbb{E}[|\boldsymbol{x}_0 - \boldsymbol{x}^*|^2] - 2\eta \sum_{s=0}^t [\mathbb{E}[f(\boldsymbol{x}_s)] - f(\boldsymbol{x}^*)] \\ + \eta^2 (t+1) (\zeta^2 + \sigma^2).$$

$$\frac{1}{t+1} \sum_{s=0}^t [\mathbb{E}[f(\boldsymbol{x}_s)] - f(\boldsymbol{x}^*)] \leq \frac{1}{2\eta(t+1)} [|\boldsymbol{x}_0 - \boldsymbol{x}^*|^2 - \mathbb{E}[|\boldsymbol{x}_{t+1} - \boldsymbol{x}^*|^2]]$$

$$+ \frac{1}{2} (G^2 + \sigma^2).$$

Choose $\eta = \sqrt{\frac{\|x_0 - x^*\|^2}{(t+1)(G^2 + \sigma^2)}}$

Avg. expected Suboptimality \leq

Stochastic oracle

function parameter

$$\frac{\sqrt{(G^2 + \sigma^2)} \cdot \|x_0 - x^*\|}{\sqrt{t+1}}.$$

For nonSmooth : this is the best possible rate ($\frac{1}{\sqrt{t}}$)

For Smooth : Cannot improve on $\frac{1}{\sqrt{t}}$; but can obtain

$$\frac{L \|x_0 - x^*\|^2}{t^2} + \frac{\sigma \|x_0 - x^*\|}{\sqrt{t}}.$$

Exercise: SGD for L -Smooth $f(\cdot)$ gets rate of

$$O\left(\frac{L\|x_0 - x^*\|^2}{t} + \sigma \frac{\|x_0 - x^*\|}{\sqrt{t}}\right).$$

Exercise: If f is g -Lipschitz and μ -strongly convex,
then SGD [with diff. η] gets rate of

$$O\left(\frac{g^2 + \sigma^2}{\mu t}\right).$$

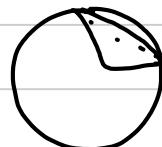
LECTURE - 6

Note Title

27-Jun-19

MIRROR DESCENT

$$\min_{\substack{x_i \geq 0 \\ \sum x_i = 1}} \frac{1}{2} \|Ax - b\|^2.$$



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

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

Smoothness

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

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

[| | | |]
↓
unit norm
columns.

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

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

$$\begin{array}{ccc}
 x & \downarrow & \nabla f \\
 \| \cdot \| & & \downarrow \\
 & & \| \cdot \|_*
 \end{array}$$

e.g., for probabilities $\| \cdot \|_2$ $\xrightarrow{\quad}$ $\| \cdot \|_\infty$

Defn. of dual norm: Given a norm $\| \cdot \|$ on X , the dual norm on the space of linear operators $L(X)$ is defined as : $\| l \|_* = \sup_{\substack{\| x \| \leq 1 \\ x \in X}} |l(x)| \quad \forall l \in L(X)$.

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

$$\begin{aligned} \left\| \boldsymbol{y} \right\|_* &= \sup_{\left\| \boldsymbol{x} \right\|_1 \leq 1} \sum_i x_i y_i \\ &= \left\| \boldsymbol{y} \right\|_\infty. \end{aligned}$$

$x_{i^*} = \text{sign}(y_i)$ On $i^* = \arg \max_i |y_i|$
 $= 0 \quad 0 \cdot \omega.$

$$\|\cdot\| \text{ on } x,$$

$$\downarrow$$

$$\|\cdot\|_2$$

$$\|\cdot\|_* \text{ on } \nabla f(x)$$

$$\downarrow$$

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

$$f(y) \geq f(x) + \underbrace{\langle \nabla f(x), y - x \rangle}_{\nabla f(x)^T e_x}$$

f -Convex & $\|\nabla f(x)\|_* \leq L$ $\| \cdot \|$

Lipschitz

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

$$= \operatorname{argmin}_x f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \underbrace{\|x - x_t\|_2^2}$$

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

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

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

$\therefore \phi$ 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

② $\|\nabla f(x)\|_* \leq \gamma$

③ $\phi(\cdot)$ is 1-strongly convex
in $\|\cdot\|$ on X .

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

$x \in X$ and $\lambda \in L(x)$
 $= \nabla f(x)$.
Hölder's ineq.: $\langle \lambda, x \rangle \leq \|\lambda\|_* \cdot \|x\|$

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

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

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

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

①

By optimality of x_{t+1} ,

$$\begin{aligned} f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{1}{2\eta} D_\phi(x_{t+1}; x_t) &\leq f(x_t) + \langle \nabla f(x_t), x_t - x_t \rangle \\ &\quad + \frac{1}{2\eta} D_\phi(x_t; x_t) \\ &= f(x_t) \end{aligned}$$

$$\begin{aligned} \frac{1}{4\eta} \|x_t - x_{t+1}\|^2 &\stackrel{\text{Str. Conv. of } \phi}{\leq} \frac{1}{2\eta} D_\phi(x_{t+1}; x_t) \leq - \langle \nabla f(x_t), x_{t+1} - x_t \rangle. \\ &\stackrel{\text{H\"older's}}{\leq} \|\nabla f(x_t)\|_* \cdot \|x_{t+1} - x_t\|. \end{aligned}$$

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

$$\begin{aligned}
 D_\phi(x; x_{t+1}) &= \phi(x) - \phi(x_{t+1}) - \langle \nabla \phi(x_{t+1}), x - x_{t+1} \rangle \\
 &\stackrel{\textcircled{1}}{=} \phi(x) - \phi(x_{t+1}) - \langle \nabla \phi(x_t) - 2\eta \nabla f(x_t), x - x_{t+1} \rangle \\
 &\stackrel{\text{Conv. of } \phi}{\leq} \phi(x) - \phi(x_t) - \langle \nabla \phi(x_t), x_{t+1} - x_t \rangle \\
 &\quad - \langle \nabla \phi(x_t) - 2\eta \nabla f(x_t), x - x_{t+1} \rangle \\
 &= \boxed{\phi(x) - \phi(x_t) - \langle \nabla \phi(x_t), x - x_t \rangle} \\
 &\quad - 2\eta \langle \nabla f(x_t), x_{t+1} - x \rangle \\
 &= D_\phi(x; x_t) - 2\eta \langle \nabla f(x_t), x_t - x \rangle
 \end{aligned}$$

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

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

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

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

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

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

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 - Str. $\xrightarrow{C \propto n}$
 w.r.t. $\|\cdot\|_2$

① If $\|\cdot\| = \|\cdot\|_2$ $\xrightarrow{\text{Euclidean}}$, then choose $\phi = \frac{1}{2} \|x\|_2^2$

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

In this case

- a) Mirror descent = GD
- b) Obtained bounds are the same.

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

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

$$\begin{aligned} D_\phi(x; y) &= \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle \\ &= \sum_i x_i \log x_i - \sum_i y_i \log y_i - \sum_i (1 + \log y_i)(x_i - y_i) \\ &= \sum_i x_i \log \frac{x_i}{y_i} + \underbrace{\sum_i (x_i - y_i)}_{=0} \end{aligned}$$

$$= KL(x||y).$$

ASIDE

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

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

$$D_\phi(x||y) = KL(x||y).$$

① ϕ [or eq. D_ϕ] is 1-St strongly convex

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

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

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

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

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

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

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

$$= \underset{\pi \in \mathcal{X}}{\operatorname{argmin}} \quad KL(\pi \| z) = z.$$

Ind. discovered

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

MULTIPLICATIVE
WEIGHTS
UPDATE
EXPONENTIATED
GRADIENT ALGORITHM

$$\min_{x: x_i \geq 0} \frac{1}{2} \|Ax - b\|^2.$$

$$\sum x_i = 1$$

$$\|\cdot\|_1 \xrightarrow{\text{on } X} \|\cdot\|_\infty \xrightarrow{\text{on } \nabla f(x)}.$$

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

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

$$\nabla f(x) = A^T(Ax - b) \rightarrow \|A^T(Ax - b)\|_\infty \leq \|A^TAx\|_\infty + \|A^Tb\|_\infty$$

Holder's inequality \leftarrow i^{th} row of A^TA \leftarrow $\max_i \|(A^TA)_i\|_\infty$
 $+ \|A^Tb\|_\infty$

Non-Euclidean

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

\downarrow
largest element
in $A^T A$

Euclidean

$$\|A^T A\|_2 \text{ Could be}$$

d times larger.

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

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

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

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

Could be
 $\propto (\alpha)$

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

LECTURE - 7

Note Title

28-Jun-19

Stochastic Variance Reduced Gradient (SVRG)

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}). \quad \text{E.g., } f_i(\mathbf{x}) = \frac{1}{2} (\mathbf{a}_i^\top \mathbf{x} - b_i)^2$$

Each f_i is L -Smooth & Convex

$$L = \max_i \|q_i\|^2$$

f is μ -Strongly convex.

$$\mu = \tau_{\min} \left(\frac{1}{n} \sum_i q_i q_i^\top \right)$$

GD: $O\left(n \cdot d \cdot \frac{L}{\mu} \log \frac{1}{\epsilon}\right)$ Computations.

One ∇f Computation

iterations.

$$\text{SGD: } O\left(d \cdot \frac{\sigma \cdot \|x_0 - x^*\|}{\epsilon^2}\right)$$

→ Std. dev. in the Stochastic gradients

↓ ↗

One ∇f_i computation. # iterations

In most cases]: $\sigma \ll n \cdot \frac{L}{\mu}$

	Dep. on param.	Desired acc.
SGD	Usually better	Usually better
GD		

→ Best of both worlds?

the slow convergence rate of SGD ($\frac{1}{\sqrt{t}}$) is due to large variance in the stochastic gradients even close to the optimal point.

$$\boxed{x^* = \arg \min_x f(x) \text{ then } \nabla f(x^*) = 0}$$

$$= \frac{1}{n} \sum_i \nabla f_i(x^*)$$

but $\nabla f_i(x^*) \neq 0$ for most i .

Variance reduction step

Pick x_0
 Compute $\nabla f(x_0) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_0)$
 Do $t = 1, \dots, T$

Pick 'i' uniformly at random from $1, \dots, n$
 $x_t = x_{t-1} - \eta [\nabla f_i(x_{t-1}) - \nabla f_i(x_0) + \nabla f(x_0)]$.

Contrast with SGD step: $x_t = x_{t-1} - \eta \nabla f_i(x_{t-1})$.

$$E_i[\nabla f_i(x_0) - \nabla f(x_0)] = E_i[\nabla f_i(x_0)] - \nabla f(x_0) = \nabla f(x_0) - \nabla f(x_0) = 0.$$

Suppose $x_0 = x^*$

$$\begin{aligned} x_1 &= x_0 - \eta [\nabla f_i(x_0) \\ &\quad - \nabla f_i(x_0) \\ &\quad + \nabla f(x_0)] \\ &= x_0 = x^* \end{aligned}$$

Recall that the convergence guarantee of SGD depends on the variance in stochastic gradient.

$$\sigma^2 \triangleq \mathbb{E} \left[\left\| \underbrace{\nabla f_i(x_t) - \nabla f_i(x_0)}_{\text{Stochastic gradient}} + \nabla f(x_0) - \nabla f(x_t) \right\|^2 \right]$$

$\overbrace{\quad\quad\quad}$ True gradient
 $= \mathbb{E}[\text{Stoch. grad.}]$

$$\begin{aligned} & \because (\alpha + b)^2 \leq 2\alpha^2 + 2b^2 \\ & \leq 2 \mathbb{E} \left[\left\| \nabla f_i(x_t) - \nabla f_i(x_0) \right\|^2 \right] + 2 \mathbb{E} \left[\left\| \nabla f(x_0) - \nabla f(x_t) \right\|^2 \right] \\ & \stackrel{\text{Jensen's ineq.}}{\geq} 4 \mathbb{E} \left[\left\| \nabla f_i(x_t) - \nabla f_i(x_0) \right\|^2 \right] \quad \left\{ \begin{array}{l} \mathbb{E} [\|\omega\|^2] \geq \|\mathbb{E} [\omega]\|^2 \\ \text{Jensen's inequality} \end{array} \right\} \end{aligned}$$

$$\xrightarrow{\quad} x^* = \underset{x}{\operatorname{argmin}} f(x)$$

$$(a+b)^2 \leq 2a^2 + 2b^2$$

$$\mathbb{E} \left[\| \nabla f_i(x_t) - \nabla f_i(x_0) \|^2 \right] \leq 2 \left[\mathbb{E} \left[\| \nabla f_i(x_t) - \nabla f_i(x^*) \|^2 \right] + \mathbb{E} \left[\| \nabla f_i(x_0) - \nabla f_i(x^*) \|^2 \right] \right].$$

$$\mathbb{E}_i \left[\| \nabla f_i(x) - \nabla f_i(x^*) \|^2 \right] \leq \frac{2L}{3} [f(x) - f(x^*)].$$

Goal: Bound ↑

Approach:

$$g_i(x) = \underbrace{f_i(x)}_{\substack{\text{L-smooth} \\ \text{Convex}}} - \underbrace{\langle \nabla f_i(x^*), x \rangle}_{\substack{\text{L-smooth} \\ \text{Convex}}} \quad \downarrow \text{Linear}$$

$$\nabla g_i(x^*) = \nabla f_i(x^*) - \nabla t_i(x^*) = 0$$

$$\Rightarrow g_i(x^*) \stackrel{\text{Conv. of } g_i}{\leq} g_i(\tilde{x} - \eta \nabla g_i(\tilde{x}))$$

$$g_i \stackrel{L-\text{smooth}}{\leq} g_i(\tilde{x}) - \underbrace{\eta(2 - \frac{1}{2L}) \|\nabla g_i(\tilde{x})\|^2}_{\substack{= \arg \min_x g_i(x) \\ \downarrow \\ \tilde{x}^*}} + \eta \nabla g_i(\tilde{x})^\top \tilde{x}$$

$$\eta = \frac{1}{L} \longrightarrow \leq g_i(\tilde{x}) - \frac{3}{2L} \|\nabla g_i(\tilde{x})\|^2$$

$$\|\nabla g_i(\tilde{x})\|^2 \leq [g_i(\tilde{x}) - g_i(x^*)]^{\frac{2L}{3}}.$$

$$\|\nabla f_i(\tilde{x}) - \nabla f_i(x^*)\|^2 \leq [g_i(\tilde{x}) - g_i(x^*)]^{\frac{2L}{3}}$$

Taking an average,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \| \nabla f_i(\tilde{x}) - \nabla f_i(x^*) \|^2 &\leq \frac{1}{n} \cdot \frac{2L}{3} \underbrace{\sum_{i=1}^n [g_i(\tilde{x}) - g_i(x^*)]}_{f_i(\tilde{x}) - \langle \nabla f_i(x^*), \tilde{x} \rangle} \\ &= \frac{2L}{3n} \sum_{i=1}^n f_i(\tilde{x}) - f_i(x^*) - \langle \nabla f_i(x^*), \tilde{x} - x^* \rangle \\ &= \frac{2L}{3} \left\{ \frac{1}{n} \sum_{i=1}^n f_i(\tilde{x}) - \frac{1}{n} \sum_{i=1}^n f_i(x^*) \right. \\ &\quad \left. - \langle \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^*), \tilde{x} - x^* \rangle \right\} \end{aligned}$$

$$= \frac{2L}{3} \left\{ f(\tilde{x}) - f(x^*) - \underbrace{\langle \nabla f(x^*), \tilde{x} - x^* \rangle}_0 \right\}$$

Lemma: $\mathbb{E} \left[\| \nabla f_i(\tilde{x}) - \nabla f_i(x^*) \|^2 \right] \leq \frac{2L}{3} \left[f(\tilde{x}) - f(x^*) \right].$

Furthermore,

$$\begin{aligned} \sigma^2 &\triangleq \mathbb{E} \left[\| \nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) - \nabla f(x_t) \|^2 \right] \\ &\stackrel[\text{2nd page computation}]{\text{3rd}}{\leq} 8 \left\{ \mathbb{E} \left[\| \nabla f_i(x_t) - \nabla f_i(x^*) \|^2 \right] + \mathbb{E} \left[\| \nabla f_i(x_0) - \nabla f_i(x^*) \|^2 \right] \right\} \\ &\leq \frac{16L}{3} \left[f(x_t) - f(x^*) + f(x_0) - f(x^*) \right]. \end{aligned}$$

$$SVRG : \quad x_{t+1} = x_t - \eta \left[\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) \right]$$

$$\begin{aligned} \mathbb{E} [\|x_{t+1} - x^*\|^2] &= \mathbb{E} [\|x_t - x^*\|^2] - 2\eta \underbrace{\langle \mathbb{E} [\quad], x_t - x^* \rangle}_{+ \eta^2 \mathbb{E} [\|\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0)\|^2]} \\ &\quad + \eta^2 \underbrace{\mathbb{E} [\|\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0)\|^2]}_{+ \eta^2 \|\nabla f(x_t)\|^2 + \eta^2 \sigma^2}. \end{aligned}$$

$$\begin{aligned} &= \mathbb{E} [\|x_t - x^*\|^2] - 2\eta \underbrace{\langle \nabla f(x_t), x_t - x^* \rangle}_{+ \eta^2 \|\nabla f(x_t)\|^2 + \eta^2 \sigma^2} \\ &\quad + \eta^2 \underbrace{\|\nabla f(x_t)\|^2}_{+ \eta^2 \|\nabla f(x_t)\|^2 + \eta^2 \sigma^2} + \eta^2 \sigma^2 \end{aligned}$$

[By L -smoothness of $f(\cdot)$,
 $\|\nabla f(x_t)\|^2 \leq L \left[f(x_t) - f(x^*) \right]$]

$$\leq \mathbb{E} [\|x_t - x^*\|^2] - 2\eta \underbrace{[f(x_t) - f(x^*)]}_{+ \eta^2 \|\nabla f(x_t)\|^2 + \eta^2 \sigma^2} + \eta^2 L (f(x_t) - f(x^*))$$

$$\begin{aligned}
& + \eta^2 \cdot \frac{16L}{3} \left[f(x_t) - f(x^*) + f(x_0) - f(x^*) \right] \\
E[\|x_{t+1} - x^*\|^2] & \leq E[\|x_t - x^*\|^2] - \eta \left(2 - \frac{19\eta L}{3} \right) (f(x_t) - f(x^*)) \\
& + \frac{16\eta^2 L}{3} (f(x_0) - f(x^*)).
\end{aligned}$$

Telescoping and averaging,

$$\frac{1}{T+1} \sum_{t=0}^T \overbrace{\left[E[f(x_t)] - f(x^*) \right]}^{\|x_0 - x^*\|^2} \leq \frac{\|x_0 - x^*\|^2}{\eta \left(2 - \frac{19\eta L}{3} \right) (T+1)} + \frac{16\eta L}{3 \left(2 - \frac{19\eta L}{3} \right)} [f(x_0) - f(x^*)].$$

f is μ -Str. CVx. $\therefore f(x_0) - f(x^*) \geq \frac{\mu}{2} \|x_0 - x^*\|^2$.

$$\text{Av. exp. Suboptimality} \leq \left[\underbrace{\frac{2}{\mu \gamma \left(2 - \frac{19\eta L}{3}\right) T_{f+1}}}_{\frac{L}{\mu} = K} + \underbrace{\frac{16\eta L}{3 \left(2 - \frac{19\eta L}{3}\right)}}_{\leq \frac{1}{10}} \right] \underbrace{\{f(x_0) - f(x^*)\}}_{\text{Initial Suboptimality}}$$

$$\eta = \frac{1}{100 \cdot L}$$

$$\approx \left[\frac{2K}{T_{f+1}} + \frac{1}{10} \right] [\text{Initial Suboptimality}] .$$

$$T > 6 \cdot K \Rightarrow \underline{\text{Av. exp. Suboptimality} \leq \frac{1}{2} \text{ Initial Subopt.}}$$

Pick x_0 .

For epochs = 1, ..., $\log(\frac{1}{\epsilon})$.

SVRG epoch

Compute $\nabla f(x_0) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_0)$ $\leftarrow O(nd)$

For $t = 0, \dots, T = 5K$

 Pick i u.a.r. from 1, ..., n .

$x_{t+1} = x_t - \eta [\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0)]$ $\leftarrow O(Kd)$

Reset $x_0 = x_1, \dots, x_{T+1}$ uniformly at random.

$O(\underbrace{(nd + Kd)}_{\text{Time}} \underbrace{\log \frac{1}{\epsilon}}_{\text{Space}})$

$$f_i(x) = (a_i^T x - b_i)^2$$

$$\mathcal{O} \left(\underset{\downarrow}{(nd + kd)} \cdot \log \frac{1}{\epsilon} \right)$$

$$\frac{\max_i \|a_i\|^2}{\mu}$$

$$\text{Smoothness} \left(\frac{1}{n\mu} (a_i^T x - b_i)^2 \right) = \frac{\sum_{j=1}^n \|a_j\|^2}{n}$$

Importance Sampling

$$i \text{ w.p. } p_i \propto \|a_i\|^2.$$

$$f(x) = \sum_{i=1}^n p_i \underbrace{\left(\frac{1}{n p_i} (a_i^T x - b_i)^2 \right)}$$

$$p_i = \frac{\|a_i\|^2}{\sum_{j=1}^n \|a_j\|^2}$$

Uniform Sampling

$$K_{us} = \frac{\max_i \|a_i\|^2}{\mu}$$

Weighted Sampling

$$K_{ws} = \frac{\sum_{j=1}^n \|a_j\|^2}{n \mu}$$

$$\underline{K_{ws} \leq K_{us}}$$

LECTURE-8

Note Title

28-Jun-19

CONVEX - CONCAVE MINIMAX OPTIMIZATION

$$\min_{x \in X} \max_{y \in Y} f(x, y)$$

$f(\cdot, y)$ is convex

$f(x, \cdot)$ is concave

$$\left\{ \begin{array}{l} g \text{ is concave} \\ \triangleq -g \text{ is convex} \\ 1) g(\alpha y_1 + (1-\alpha)y_2) \geq \alpha g(y_1) + (1-\alpha)g(y_2) \\ 2) g(z) \leq g(y) + \langle \nabla g(y), z - y \rangle \end{array} \right.$$

i) Constrained Optimization :

$$\min_{\boldsymbol{x}} \max_{i=1, \dots, m} f_i(\boldsymbol{x}) \quad \leftarrow$$

$$\begin{aligned} & \text{Find } \boldsymbol{x} \\ & \text{s.t. } f_i(\boldsymbol{x}) \leq 0 \quad i=1, \dots, m. \end{aligned}$$

Each f_i is convex.

III

$$\begin{aligned} & \min_{\boldsymbol{x}} \max_{\substack{\boldsymbol{y} \\ \boldsymbol{y} \geq 0 \\ \sum_{i=1}^m y_i = 1}} \sum_{i=1}^m y_i f_i(\boldsymbol{x}) \end{aligned}$$

$$f(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i=1}^m y_i f_i(\boldsymbol{x})$$

$f(\cdot, \boldsymbol{y})$ is convex if \boldsymbol{y}
 $f(\boldsymbol{x}, \cdot)$ is linear in \cdot if \boldsymbol{x} concave

2) Many non-smooth functions can be written as smooth minimax problems.

$$a) \quad \|x\|_1 = \max_{-1 \leq y_i \leq 1} \sum_i y_i x_i$$

$$\min_x \frac{1}{2} \|Ax - b\|^2 + d \|x\|_1 = \min_x \max_{-1 \leq y_i \leq 1} \underbrace{\frac{1}{2} \|Ax - b\|^2 + d \sum_i y_i x_i}_{\text{---}}$$

b)

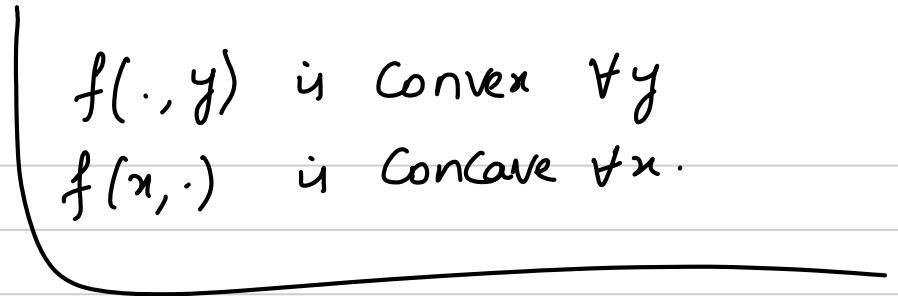
Given a_1, \dots, a_m find the
smallest ball covering a_1, \dots, a_m .

→ Again can be written as smooth
minimax problem.

3) Minimax ≡ zero sum games.

$$\min_x \underbrace{\max_y f(x, y)}_{\equiv g(x)}$$

$f(\cdot, y)$ is convex w.r.t y
 $f(x, \cdot)$ is concave w.r.t x .



Exercise: $g(x)$ is convex

How to evaluate $\nabla g(x)$?

(Informal) Danskin's thm: Under smoothness assumptions on $f(\cdot, \cdot)$,

$$\nabla g(x) = \text{Conv. hull} \left\{ \nabla_x f(x, y) : y \in \arg \underline{\max}_z f(x, z) \right\}.$$

$$\left. \begin{array}{l} 1) \text{ Find } y_t \in \arg \max_y f(x_t, y) \\ 2) \quad x_{t+1} = x_t - \eta \nabla_x f(x_t, y_t) \end{array} \right\} \begin{array}{l} \text{Convex opt.} \\ \text{Subgradient descent} \\ \text{on } g(x). \end{array}$$

Issues : ① Finding y_t takes time
 ② g could be non smooth \Rightarrow slow convergence rate.

Non smooth : Gradient descent ascent $\rightarrow O(\frac{1}{\sqrt{T}})$

Smooth : Mirror-Prox $\rightarrow O(\frac{1}{T})$.

$g(x)$ could still be non smooth

[f could be arbitrary] maximin \leq minimax

Weak duality: $\underbrace{\max_{y \in Y} \min_{x \in X} f(x, y)}_{LHS} \leq \underbrace{\min_{x \in X} \max_{y \in Y} f(x, y)}_{RHS}.$

$$y^* = \arg \max_{y \in Y} \left[\min_{x \in X} f(x, y) \right]$$

$$x^* = \arg \min_{x \in X} \left[\max_{y \in Y} f(x, y) \right]$$

$$LHS = \min_{x \in X} f(x, y^*) \leq f(x^*, y^*) \leq \max_{y \in Y} f(x^*, y) = RHS.$$

[Sion's minimax thm.]

Strong duality: If $f(x, y)$ is Convex-Concave and X and Y

are Compact then $\min_{x \in X} \max_{y \in Y} f(x, y) = \max_{y \in Y} \min_{x \in X} f(x, y)$.

Defn. ϵ -primal dual pair: (\bar{x}, \bar{y}) is Said to be an ϵ -primal

dual pair if $\max_y f(\bar{x}, y) - \min_x f(x, \bar{y}) \leq \epsilon$.

Note: $\max_y f(\bar{x}, y) \geq f(\bar{x}, \bar{y}) \geq \min_x f(x, \bar{y})$.

Let (\bar{x}, \bar{y}) be an ϵ -primal dual pair

Exercise: Show that \bar{x} is an ϵ -optimal soln. for

$$\min_x g(x) \triangleq \max_y f(x, y)$$

Exercise: \bar{y} is an ϵ -optimal soln. for $\max_y h(y) \triangleq \min_x f(x, y)$.

$$\text{diam}(X), \text{diam}(Y) \leq R$$

Gradient descent ascent :

$$x_{t+1} = x_t - \eta \cdot \nabla_x f(x_t, y_t) \rightarrow \text{GD for min}$$

$$y_{t+1} = y_t + \eta \cdot \nabla_y f(x_t, y_t) \rightarrow \text{GA for max}$$

$f(\cdot, y)$ is convex

$f(x, \cdot)$ is concave

$$\|\nabla_x f(x, y)\| \leq G$$

$$\|\nabla_y f(x, y)\| \leq G$$

Thm: GDA has Conv. rate $O\left(\frac{\epsilon R}{\sqrt{T}}\right)$.

Proof: $\|x_{t+1} - x\|^2 = \|x_t - x\|^2 - 2\eta \underbrace{\langle \nabla_x f(x_t, y_t), x_t - x \rangle}_{\text{Conv. of } f(\cdot, y_t)} + \eta^2 \underbrace{\|\nabla_x f(x_t, y_t)\|^2}_{G\text{-Lipschitz}}$

x is arb.

$$\leq \|x_t - x\|^2 - 2\eta [f(x_t, y_t) - f(x, y_t)] + \eta^2 G^2$$

$$\boxed{y \text{ in arb.}} \quad \|y_{t+1} - y\|^2 = \|y_t - y\|^2 - 2\eta \underbrace{\langle \nabla_y f(x_t, y_t), y - y_t \rangle}_{\text{Concavity of } f(x_t, \cdot)} + \eta^2 \underbrace{\|\nabla_y f(x_t, y_t)\|^2}_{L\text{-Lipschitz}}$$

$$\leq \|y_t - y\|^2 - 2\eta [f(x_t, y) - f(x_t, y_t)] + \eta^2 \varsigma^2$$

$$\|x_{t+1} - x\|^2 + \|y_{t+1} - y\|^2 \leq \|x_t - x\|^2 + \|y_t - y\|^2 - 2\eta [f(x_t, y) - f(x, y_t)] + 2\eta^2 \varsigma^2.$$

$$\frac{1}{T+1} \sum_{t=0}^T [f(x_t, y) - f(x, y_t)] \leq \frac{\|x_0 - x\|^2 + \|y_0 - y\|^2}{2\eta(T+1)} + \eta \varsigma^2$$

$$f\left(\underbrace{\frac{1}{T+1} \sum_{t=0}^T x_t}_{\bar{x}_T}, y\right) \leq \frac{1}{T+1} \sum_{t=0}^T f(x_t, y) \quad \text{and} \quad f\left(x, \underbrace{\frac{1}{T+1} \sum_{t=0}^T y_t}_{\bar{y}_T}\right) \geq \frac{1}{T+1} \sum_{t=0}^T f(x, y_t)$$

$$f(\bar{x}_T, y) - f(x, \bar{y}_T) \leq \frac{\|\bar{x}_0 - x\|^2 + \|y_0 - y\|^2}{2\eta(T+1)} + \eta\zeta^2.$$

$$\begin{aligned} \max_y f(\bar{x}_T, y) - \min_x f(x, \bar{y}_T) &\leq \frac{\max_x \|\bar{x}_0 - x\|^2 + \max_y \|y_0 - y\|^2}{2\eta(T+1)} + \eta\zeta^2 \\ &\leq \frac{R^2}{\eta(T+1)} + \eta\zeta^2 \\ \eta = \frac{R}{6\sqrt{T+1}} &\leq \frac{2GR}{\sqrt{T+1}}. \end{aligned}$$

(\bar{x}_T, \bar{y}_T) is an $\frac{2GR}{\sqrt{T+1}}$ - primal dual pair.

⑩

[Prox] Proximal algorithm: $x_{t+1} = \arg \min_{x \in X} \left[f(x) + \frac{1}{2\eta} \|x - x_t\|^2 \right]$

GD Step: $x_{t+1} = \arg \min_{x \in X} \left[f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|^2 \right]$

Exercise: Prox has convergence rate $O(\frac{1}{T})$ for
[possibly non-smooth] convex functions.

Exercise: Prox step is efficiently implementable for

L -Smooth functions and $\eta \leq \frac{1}{2L}$.

Efficient: Can find ϵ -approximate Soln. to the prox step
in $O(\log \frac{1}{\epsilon})$ iterations.

Prox Step: Given x , find w^* s.t.

$$w^* = \underset{z \in X}{\operatorname{arg\,min}} \quad f(z) + \frac{1}{2\eta} \|x - z\|^2$$

$$\equiv \quad w^* = \quad x - \eta \nabla f(w^*) \rightarrow \text{Prox.}$$

$$\nabla f(x) \rightarrow \text{GD}$$

Algorithm to find]. Pick $w_0 = x$
 From Step $w_{t+1} = x - \eta \nabla f(w_t) \rightarrow ① //$

Ⓐ Claim: w^* is the unique fixed point of ①

Ⓑ Claim: ① is a $\frac{1}{2}$ Contraction.

$$\begin{aligned} w \rightarrow w^+ &= x - \eta \nabla f(w) \\ \tilde{w} \rightarrow \tilde{w}^+ &= x - \eta \nabla f(\tilde{w}) \end{aligned} \quad \left. \begin{aligned} \|w^+ - \tilde{w}^+\| &= \eta \|\nabla f(w) - \nabla f(\tilde{w})\| \\ &\leq \eta L \|w - \tilde{w}\| \end{aligned} \right\}$$

If $\eta \leq \frac{1}{2L}$ $\leq \frac{1}{2} \|w - \tilde{w}\|$.

Ⓐ + Ⓑ \Rightarrow ① finds ϵ -approximate w^* in $\log \frac{1}{\epsilon}$ iterations.

Conceptual Mirror-Prox

$$x_{t+1} = \underset{x \in X}{\operatorname{argmin}} f(x, y_{t+1}) + \frac{1}{2\eta} \|x - x_t\|^2$$

$$y_{t+1} = \underset{y \in Y}{\operatorname{argmax}} f(x_{t+1}, y) - \frac{1}{2\eta} \|y - y_t\|^2$$

Implementation of each step

$$\begin{aligned} x_t^0 &= x_t \\ y_t^0 &= y_t \end{aligned}$$

$$i=1, \dots, \log \frac{1}{\epsilon}, \quad x_t^i = \underset{x \in X}{\operatorname{argmin}} f(x, y_t^{i-1}) + \frac{1}{2\eta} \|x - x_t\|^2$$

2' in the actual mirror prox alg.

$$y_t^i = \underset{y \in Y}{\operatorname{argmax}} f(x_t^{i-1}, y) - \frac{1}{2\eta} \|y - y_t\|^2.$$

Exercise: For Conceptual Mirror-Prox, obtain $O\left(\frac{1}{\tau}\right)$ convergence rate for L-Smooth minimax opt.

Exercise: Show that the algorithm for implementing the Prox Step Converges in $O\left(\log \frac{1}{\epsilon}\right)$ iterations.

Exercise: Careful accounting of all the terms to show that 2 inner steps suffice for $O\left(\frac{1}{\tau}\right)$ convergence rate.