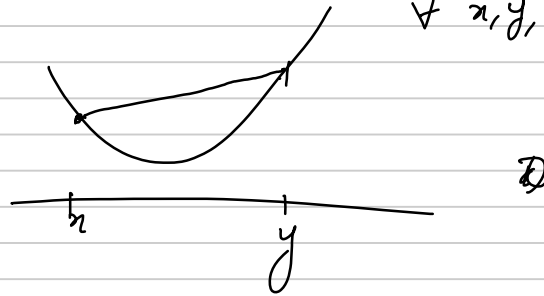


LEC-2

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

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: Exercise.



$$\widehat{\text{Convex}} : f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

Lemma: If f is $\widehat{\text{Convex}}$ 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

$$\text{MORALLY: CONVEX} = \widehat{\text{CONVEX}}$$

↓
First defn.

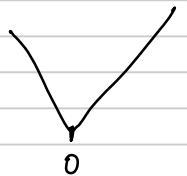
↓
Second defn.

Lasso:

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

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

does not exist at $x=0$.



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

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

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

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

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

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

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

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

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

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

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

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 $\|x_0 - x^*\|$.
Need to learn them on the go.

LOWER BOUNDS

Gradient span

Algorithms: $x_{t+1} \in \text{span} \{x_0, \nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_t)\}$.

(GSA)

GD \in Gradient span algorithms.

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

i) $\|\nabla f_{G,R,T}(x)\| \leq G$

o) f is convex

ii) $\|x_0 - x^*\| \leq R$

and

any vector in the span $\{x_0, \dots, x_T\}$

iii) for any GSA

$$f_{G,R,T}(\bar{x}) - f_{G,R,T}(x^*) \geq \frac{GR}{2(1+\sqrt{T+1})}$$

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

Proof: $f_{r,\mu}(x) = \min_{1 \leq i \leq T+1} x(i) + \frac{\mu}{2} \|x\|^2$

$x \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$ dimension = $T+1$

Ex: Show that $f_{r,\mu}(\cdot)$ 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

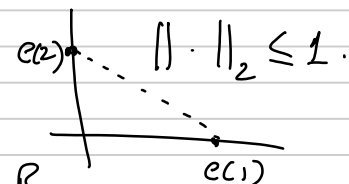
$$f(x) = \max_i f_i(x)$$

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

Any vector $z \in \text{Conv hull of } \{ \nabla f_i(x) : i \in \mathcal{I}(x) \}$ is a subgradient of $f(\cdot)$ at x .

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

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



$\forall \|x\| \leq R$

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

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

$$x(i) = \alpha$$

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

$$\alpha = \frac{-\rho}{\mu(T+1)}$$

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

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

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

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

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

$$x_0 = 0$$

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

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

$$x_0 = 0$$

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

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

$$\nabla f(x_0) = \min e_{(1)} + \mu \cdot 0 = \min e_{(1)}$$

$$x_1 \in \text{span} \left\{ \begin{array}{c} x_0 \\ \parallel \\ 0 \end{array}, \begin{array}{c} \nabla f(x_0) \\ \parallel \\ \text{span}(e_{(1)}) \end{array} \right\} = \text{span}(e_{(1)})$$

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

$$\nabla f(x_1) \in \text{span}\{e_{(1)}, e_{(2)}\}$$

$$f(x) = \min \left\{ \begin{array}{c} \uparrow f_1(x) \\ x(1) \\ = \alpha \end{array}, \begin{array}{c} \uparrow f_2(x) \\ x(2) \\ = 0 \end{array}, \dots \right\} + \frac{\mu}{2} \|x\|^2$$

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

If $\alpha > 0$ then $\mathcal{I}(x_1) = 1 \longrightarrow \min e_1$
 If $\alpha = 0$ then $\mathcal{I}(x_1) = 1, 2 \longrightarrow \min e_1$
 If $\alpha < 0$ then $\mathcal{I}(x_1) = 2 \longrightarrow \min e_2$

$x_t \in \text{span}\{e_{(1)}, \dots, e_{(t)}\}$.
 $f(x_t) = \min_{1 \leq i \leq T+1} x_t(i) + \frac{\mu}{2} \|x_t\|^2$
 $x_t(T+1) = 0$.
 So, $f(x_t) \geq 0$.

$$\sqrt{v} = \frac{\sqrt{T+1} \cdot G}{1 + \sqrt{T+1}} ; \mu = \frac{G}{R(1 + \sqrt{T+1})} ; f(x_T) \geq 0.$$

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

$$\begin{aligned} f(x_T) - f(x^*) &\geq \frac{v^2}{2\mu(T+1)} = \frac{(T+1)G^2}{(1 + \sqrt{T+1})^2} \cdot \frac{R(1 + \sqrt{T+1})}{2G} \cdot \frac{1}{T+1} \\ &= \frac{GR}{2(1 + \sqrt{T+1})} \quad \square \end{aligned}$$

$$f(x) = \max \{f_1(x), f_2(x)\}.$$