

INTRODUCTORY COURSE ON
OPTIMIZATION

PRANEETH NETRAPALLI.

CONVEX FUNCTIONS:

$$\min_x f(x)$$

Iterative methods: Start with $x_0 \rightarrow x_1 \rightarrow \dots$

How do we access f ?

Black-box models: Zeroth order $x \rightarrow f(x)$
First order $x \rightarrow \nabla f(x)$
Second order $x \rightarrow \nabla^2 f(x)$.

First-order hits a sweet spot. If f is given by a circuit with smooth components, can compute ∇f in essentially the same amount of time.

Classes of functions: Assume for now the function ' f ' is very smooth.

Prototypical first order method: GRADIENT DESCENT.

1. Start with x_0

2. Do $k=0, 1, \dots$

$$x_{k+1} = x_k - \eta \nabla f(x_k).$$

[LATER]

Lemma: Suppose $\|\nabla^2 f\| \leq L$ and $\eta = \frac{1}{L}$. Then

$$\min_{i=1, \dots, k} \|\nabla f(x_i)\|^2 \leq \frac{(f(x_0) - f^*)^2 L}{k}$$

Proof: $f(x_{k+1}) = f(x_k - \eta \nabla f(x_k))$

$$= f(x_k) + \langle \nabla f(x_k), -\eta \nabla f(x_k) \rangle + \eta^2 \nabla f(x_k)^T \nabla^2 f(\tilde{x}_k) \nabla f(x_k)$$

[by]

SMOOTHNESS: $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$

Lemma: If f is L -smooth, then

$$|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{L}{2} \|x - y\|^2$$

Proof: $f(y) = f(x) + \langle \nabla f(x), y - x \rangle + \int_0^1 \langle \nabla f(\tau y + (1-\tau)x) - \nabla f(x), y - x \rangle d\tau$

$$|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| = \left| \int_0^1 \langle \nabla f(\tau y + (1-\tau)x) - \nabla f(x), y - x \rangle d\tau \right|$$

$$\stackrel{(C-S)}{\leq} \int_0^1 \|\nabla f(\tau y + (1-\tau)x) - \nabla f(x)\| \|y - x\| d\tau$$

$$\stackrel{(Smoothness)}{\leq} \left[\int_0^1 \tau d\tau \right] \cdot L \|y - x\|^2$$

$$= \frac{L}{2} \|y - x\|^2 \quad \square$$

Proof of GD:

$$f(x_{k+1}) \leq f(x_k) - \eta \|\nabla f(x_k)\|^2 + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

$$= f(x_k) - \eta \left(1 - \frac{\eta L}{2}\right) \|\nabla f(x_k)\|^2$$

$$\Rightarrow \|\nabla f(x_k)\|^2 \leq 2L [f(x_k) - f(x_{k+1})]$$

Summing over k , we have

$$\frac{1}{K} \sum_{i=1}^K \|\nabla f(x_i)\|^2 \leq 2L [f(x_0) - f(x_{K+1})] \cdot \frac{1}{K}$$

$$\leq 2L [f(x_0) - f^*] \cdot \frac{1}{K}$$

$$\Rightarrow \min_i \|\nabla f(x_i)\|^2 \leq \frac{2L \cdot [f(x_0) - f^*]}{K} \quad \square$$

In general, this is all we can say using first order algorithms for general smooth functions.

Assumption I: $\nabla f(x) = 0 \Rightarrow$ Global optimum.

Assumption II: $f_1, f_2 \in \mathcal{F}$ then $\alpha f_1 + \beta f_2 \in \mathcal{F}$ for $\alpha, \beta \geq 0$.

Assumption III: Linear fns. $\in \mathcal{F}$.

Lemma: $\forall f \in \mathcal{F}$, we have $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$

\mathcal{F} is the set of convex functions

Proof: Let $g(y) = f(y) - f(x) - \langle \nabla f(x), y - x \rangle \in \mathcal{F}$

$\nabla g(y) = \nabla f(y) - \nabla f(x)$. $\nabla g(x) = 0 \Rightarrow x$ is global min of $g(\cdot)$.
 \Rightarrow Statement of lemma.

Examples of convex functions

① Linear regression: $f(x) = \frac{1}{2} \|Ax - b\|^2$.

$$\nabla f(x) = A^T(Ax - b)$$

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

$$= \frac{1}{2} \|Ay - b\|^2 - \frac{1}{2} \|Ax - b\|^2 - \langle A^T(Ax - b), y - x \rangle$$

$$= \frac{1}{2} \|A(x - y)\|^2 \geq 0.$$

② L_p norms: $f(x) = |x|^p$ Say p is even.
 $\nabla f(x) = p x^{p-1}$

$$g_x(y) = f(y) - f(x) - \langle \nabla f(x), y - x \rangle = |y|^p - |x|^p - p x^{p-1}(y - x)$$
$$= \cancel{y^p - x^p} - \cancel{p x^{p-1} y}$$

$$\nabla g_x(y) = p [y^{p-1} - x^{p-1}] = 0 \Rightarrow y = x.$$

③ Composition: If f is convex then $g(x) = f(Ax + b)$ is convex.

$$\nabla g = A^T \nabla f(Ax + b)$$

$$g(y) - g(x) - \langle \nabla g(x), y - x \rangle = f(Ay + b) - f(Ax + b)$$

$$- \langle A^T \nabla f(Ax + b), y - x \rangle$$

$$= f(Ay + b) - f(Ax + b) - \langle \nabla f(Ax + b), (Ay + b) - (Ax + b) \rangle$$

≥ 0

③ Max of Convex functions is also Convex.

④ L_p regression: $f(x) = \|Ax - b\|_p^p$.

⑤ LASSO: $f(x) = \frac{1}{2} \|Ax - b\|^2 + \|x\|_1$.

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

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

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

Theorem: If $\|\nabla f(x_t)\| \leq G$ then

$$\frac{1}{T} \sum_{t=1}^T f(x_t) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2\eta T} + \frac{\eta G^2}{2}$$

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

$$= \|x_t - x^*\|^2 - 2\eta \langle \nabla f(x_t), x_t - x^* \rangle + \eta^2 \|\nabla f(x_t)\|^2$$

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

$$\Rightarrow 2\eta \langle \nabla f(x_t), x_t - x^* \rangle \leq \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 + \eta^2 G^2$$

Adding and dividing by T ,

$$\frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq \frac{\|x_0 - x^*\|^2}{2\eta T} + \frac{\eta G^2}{2}$$

Use Convexity. ▮

Choosing $\eta = \frac{\|x_0 - x^*\|}{G\sqrt{T}}$ gives us,

$$\frac{1}{T} \sum_{t=1}^T [f(x_t) - f(x^*)] \leq \frac{\|x_0 - x^*\| G}{\sqrt{T}}$$

Remarks: ① The resulting rate is indep. of dimension.

① Need knowledge of $\|x_0 - x^*\|$ and G to set the optimal parameters. Designing parameter free algorithms is practically important. There are results but outside the scope of current course.

② All we needed from the function was

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

As long as we have a notion of $\nabla f(x)$ satisfying this, we are OK. We do not even need the function to be differentiable. The notion of Subgradient Suffices.

Subgradient: A vector g is said to be a Subgradient of f at x if

$$f(y) \geq f(x) + \langle g, y - x \rangle \quad \forall y.$$

Need access to Subgradient oracle.

Example: $f(x) = |x|$.

For $x \neq 0$, $\nabla f(x) = \text{sign}(x)$.

$x = 0$, $\nabla f(x) = [-1, 1]$. III

Lower bounds: Fix I_S the bound that we obtained,
i.e., $f(\bar{x}_T) - f(x^*) \leq \frac{\|x_0 - x^*\| \cdot G}{\sqrt{T}}$ optimal?

How do we define optimality? For any given problem,
 I can hardcode its output. We need to look at
function classes.

Let $\mathcal{F}_{G,R,x_0} = \left\{ f : \begin{array}{l} f \text{ is convex} \\ f \text{ is Lipschitz i.e., } |f(x) - f(y)| \leq G \|x - y\| \\ x^* \text{ exists and } x^* \in B(x_0, R) \end{array} \right\}$.

Note: f is Lipschitz \Leftrightarrow All subgradients are bounded.
 \rightarrow Use T instead of K

Theorem: For any class \mathcal{F}_{G,R,x_0} (i.e., for any values of
 G, R and x_0) and any $0 \leq k \leq n-1$ \exists a function $f \in \mathcal{F}_{G,R,x_0}$ s.t.

$$f(x_k) - f^* \geq \frac{GR}{2(1 + \sqrt{k+1})}$$

\forall optimization algorithm generating x_k s.t.

$$x_k \in x_0 + \text{Lin span} \left\{ \nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_{k-1}) \right\}$$

USE 'T' INSTEAD OF 'K'

Note: Gradient descent $x_{k+1} = x_k - \eta \nabla f(x_k)$ satisfies the hypothesis. Indeed in the black box model, this seems to be the most an algorithm can do.

Proof: WLOG $x_0 = 0$. Fix $0 \leq k \leq n-1$.

Pick $f(x) = \frac{\sqrt{k+2} \cdot G}{1 + \sqrt{k+2}} \max_{1 \leq i \leq k} x^{(i)} + \frac{G}{2(1 + \sqrt{k+2})R} \|x\|^2$.

Denote $\nu = \frac{\sqrt{k+2} \cdot G}{1 + \sqrt{k+2}}$ and $\mu = \frac{G}{(1 + \sqrt{k+2})R}$.

→ May be later.

- i) f is convex
- iii) $x^{*(i)} = \begin{cases} -\frac{\nu}{\mu(k+1)} & 1 \leq i \leq k \\ 0 & k+1 \leq i \leq n \end{cases}$

Aside: x^* optimizes $f \iff 0$ is a subgradient of f at x^* .

Rewrite: $f(x) = \nu \max_{1 \leq i \leq k+1} x^{(i)} + \frac{\mu}{2} \|x\|^2$

RULES FOR COMPUTING SUBGRADIENT.

$$\begin{aligned} \partial f(x) &= \frac{\nu}{k+1} \cdot (e_1 + \dots + e_{k+1}) + \mu x \\ &= 0. \quad \text{So } x^* \text{ is opt of } f. \end{aligned}$$

$$\|x^*\| = \frac{\gamma}{\mu\sqrt{k+1}} \Rightarrow x^* \in B(x_0, \frac{\gamma}{\mu\sqrt{k+1}})$$

ii) SUBgradient at x :

$$\partial f(x) = \mu x + \gamma \cdot \text{Conv} \{ e_i : i \in I(x) \},$$

$$\text{where } I(x) = \left\{ j : x^{(j)} = \max_{1 \leq i \leq k+1} x^{(i)} \right\}.$$

$$\begin{aligned} \|\partial f(x)\| &\leq \mu \|x\| + \gamma \cdot \max_{v \in \text{Conv}\{\dots\}} \|v\| \\ &\leq \mu \cdot R + \gamma \cdot \end{aligned}$$

$$\begin{aligned} \text{Want: } \mu R + \gamma &\leq G \\ \frac{\gamma}{\mu\sqrt{k+1}} &\leq R \end{aligned} \left. \begin{array}{l} \\ \end{array} \right\} \rightarrow \text{Choose } \begin{aligned} \gamma &= \frac{\sqrt{k+1} \cdot G}{1 + \sqrt{k+1}} \\ \mu &= \frac{G}{R(1 + \sqrt{k+1})} \end{aligned}$$

Now: ~~$\nabla f(x)$~~ $\partial f(x) = \mu x + \gamma \cdot \text{Conv} \{ e_i : i \in I(x) \}.$

The subgradient oracle could return any item of the subgradient set. Let us say, we are dealing with the subgradient oracle which gives $\nabla f(x) = \mu x + \gamma e_l$ where

$$l = \min \left\{ j : x_j = \max_{1 \leq i \leq k+1} x_i \right\}.$$

$$x_0 = 0$$

$$\nabla f(0) = \nabla f(x_0) = \mu \cdot 0 + \nu e_1 = \nu e_1$$

$$x_1 \in 0 + \text{span}(e_1).$$

$$\nabla f(x_1) = \mu x_1 + \nu \{e_1 \text{ or } e_2\}$$

$$\text{So, } x_2 \in 0 + \text{span}\{\nabla f(x_0), \nabla f(x_1)\}$$

$$\subseteq \text{span}\{e_1, e_2\}.$$

Similarly if $x_{k-1} \in \text{span}\{e_1, \dots, e_{k-1}\}$,

then $\nabla f(x_k) \in \text{span}\{e_1, \dots, e_k\}$

and so $x_k \in \text{span}\{e_1, \dots, e_k\}$.

On the other hand,

$$\min_{x: x \in \text{span}\{e_1, \dots, e_k\}} f(x) \geq \nu \cdot 0 + \mu \cdot 0 = 0.$$

as long as $l \leq k$.

$$\text{So, } f(x_k) \geq 0 \text{ while } f(x^*) = \frac{-\nu^2}{\mu(k+1)} + \frac{\mu}{2} \cdot \frac{\nu^2}{\mu^2(k+1)}$$

$$= \frac{-\nu^2}{2\mu(k+1)}$$

$$= \frac{-(k+2)GR}{2(k+1)[1+\sqrt{k+2}]}$$

$$= \frac{-GR}{2[1+\sqrt{k+2}]}$$

$$\text{So, } f(x_k) - f(x^*) \geq \frac{GR}{2[1+\sqrt{k+4}]} \quad \square$$

Comparing with the performance bound of GD, we see that this is tight.

$$\text{GD: } f(\bar{x}_T) - f(x^*) \leq \frac{GR}{\sqrt{T}}$$

Lower bound: $f(x_T) - f(x^*) \geq \frac{GR}{2[1+\sqrt{T+1}]}$
for any "span" algorithm.

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

Linear model for (a_i, b_i) : $\langle a_i, x \rangle \approx b_i$.

~~So~~ Given $A = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}$ and $b = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$

we would like $x^* = \underset{x}{\operatorname{argmin}} f(x)$

Let us suppose $\|x^*\| \leq R$.

$$\|\nabla f(x)\| = \|Ax - b\| \leq \|A\| \cdot R + \|b\| \triangleq G.$$

Then GD starting at $x_0 = 0$ will give

$$f(\bar{x}_T) - f(x_0^*) \leq \frac{R[\|A\| \cdot R + \|b\|]}{\sqrt{T}}$$

Can we do better for this problem?

Perhaps. The gradients exist everywhere.
In fact they are Lipschitz continuous.

$$\begin{aligned}\|\nabla f(x) - \nabla f(y)\| &= \|Ax - b - Ay + b\| \\ &\leq \|A\| \cdot \|x - y\|.\end{aligned}$$

The lower bound example we constructed does not have this property. ~~So, may~~

Theorem: Suppose f has Lipschitz gradients (also known as smoothness): $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$.
Suppose x^* exists. ~~and $x^* \in \mathbb{R}^n$ (not)~~. Then,

$$f(\bar{x}_T) - f(x^*) \leq \frac{6L \cdot \|x_0 - x^*\|^2}{T}$$

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

$$\begin{aligned}\|x_{t+1} - x^*\|^2 &= \|x_t - \eta \nabla f(x_t) - x^*\|^2 \\ &= \|x_t - x^*\|^2 - 2\eta \langle \nabla f(x_t), x_t - x^* \rangle + \eta^2 \|\nabla f(x_t)\|^2.\end{aligned}$$

On the other hand,

$$\begin{aligned}f(x^*) &\leq f(x_t - \eta \nabla f(x_t)) \leq f(x_t) - \eta \|\nabla f(x_t)\|^2 + L\eta^2 \|\nabla f(x_t)\|^2 \\ \Rightarrow \eta(1 - \eta L) \|\nabla f(x_t)\|^2 &\leq f(x_t) - f(x^*).\end{aligned}$$

By convexity, $f(x^*) \geq f(x_t) + \langle \nabla f(x_t), x^* - x_t \rangle$.

This gives, $\|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 - \eta \left[2 - \frac{1}{\eta L}\right] [f(x_t) - f(x^*)]$

$$\Rightarrow f(x_t) - f(x^*) \leq \frac{1}{\eta \left[2 - \frac{1}{\eta L}\right]} \left\{ \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 \right\}.$$

Taking a telescopic sum and using $f(x_{t+1}) \leq f(x_t)$, we have

$$f(x_T) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{T \cdot \eta \left[2 - \frac{1}{\eta L}\right]}.$$

Choosing $\eta = \frac{1}{3L}$, we have

$$f(x_T) - f(x^*) \leq \frac{6 \|x_0 - x^*\|^2 \cdot L}{T}. \quad \square$$

Compare this to the non-smooth setting where we obtained

$$f(\bar{x}_T) - f(x^*) \leq \frac{GR}{\sqrt{T}}.$$

For the linear least squares $f(x) = \frac{1}{2} \|Ax - b\|^2$,
Smoothness parameter = $\|A\|^2$.

$$\text{So, } f(x_T) - f(x^*) \leq \frac{6 \|A\|^2 \|x_0 - x^*\|^2}{T}.$$

Is this tight? Best possible?

Lower bound for Smooth functions

$$f(x) = \frac{1}{2} x^T A x - x^T b$$

where $A = \begin{bmatrix} 2 & -1 & & \\ -1 & \ddots & & \\ & & \ddots & \\ & & & 2 \end{bmatrix}$

* $\|A\| \leq 4$. [by Gershgorin disk theorem]

$$\text{So, } \|\nabla f(x) - \nabla f(y)\| \leq 4 \cdot \|x - y\|.$$

* Let $x^{(0)} = 0$.

$$f(x) = \frac{1}{2} \left[x_1^2 + (x_1 - x_2)^2 + (x_2 - x_3)^2 + \dots \right] - x_1.$$

If we have $x_{k+1} = x_{k+2} = \dots = 0$ then the

minimizer over only x_1, \dots, x_k will satisfy

$$\left. \begin{aligned} 2x_1 - x_2 &= 1 \\ x_3 - 2x_2 + x_1 &= 0 \\ x_4 - 2x_3 + x_2 &= 0 \\ &\vdots \\ x_{k-1} - 2x_k &= 0 \end{aligned} \right\}$$

$$x_i = \frac{(i+1)2^{i-k}}{k+1}$$

$$f_k^* = \frac{1}{2} \left[\left(\frac{k}{k+1} \right)^2 + k \cdot \left(\frac{1}{k+1} \right)^2 \right] - \frac{k}{k+1}$$

$$= \frac{1}{2} \cdot \frac{k}{k+1} - \frac{k}{k+1} = \frac{-k}{2(k+1)} = -\frac{1}{2} \left\{ \frac{1}{k+1} \right\}$$

$$\|x^*\|^2 = \frac{\sum_{i=0}^{k-1} (i+1)^2}{(k+1)^2} = \frac{(k-1)k(2k-1)}{6(k+1)^2} \rightarrow \text{Choose } k=2k$$

$$\frac{f_k^* - f^*}{L \cdot \|x^* - x^0\|^2} \geq \frac{\frac{1}{4(k+1)}}{k^1} \approx \frac{1}{k^2}$$

The lower complexity bound is $\frac{1}{k^2}$.

GD upper bound $\rightarrow \frac{1}{k}$.

Can we actually achieve $\frac{1}{k^2}$?

YES! NESTEROV'S AGD. Quad. U.B. minimization viewpoint of GD

In GD, we use gradient information from only the current point. The main idea is to make use of gradients from previous time steps as well.

Estimate Sequences: Start with a quadratic fn.

$$\Phi_0(x) = \Phi_0^* + \frac{L}{2} \|x - x_0\|^2, \text{ for some } x_0.$$

*USE τ instead of k

When we query the gradient at a point y_k ,
we get a lower bound $f(x) \geq f(y_k) + \langle \nabla f(y_k), x - y_k \rangle$.

Update $\phi_{k+1}(x) = (1 - \alpha_k) \phi_k(x) + \alpha_k [f(y_k) + \langle \nabla f(y_k), x - y_k \rangle]$

Choosing $d_0 = 1$ and $d_{k+1} = (1 - \alpha_k) d_k$, we have

$$\phi_{k+1}(x) \leq d_{k+1} \phi_0(x) + (1 - d_{k+1}) f(x).$$

Let $\phi_k(x) = \phi_k^* + \frac{d_k}{2} L \|x - v_k\|^2$.

Lemma: If $\exists x_k$ s.t. $f(x_k) \leq \phi_k^*$ then

$$\min_x f(x) \leq f(x_k) - f^* \leq d_k [\phi_0(x^*) - f^*]$$

Proof: $f(x_k) - f^* \leq \phi_k^* - f^*$

$$\leq \phi_k(x^*) - f(x^*)$$

$$\leq d_k [\phi_0(x^*) - f(x^*)] \quad \square$$

Two requirements: ① d_k should decrease rapidly.

② Can always find x_k s.t. $f(x_k) \leq \phi_k^*$.

Given x_k and $\phi_k(\cdot)$, let us see how to find x_{k+1} and $\phi_{k+1}(\cdot)$ satisfying

i) $f(x_{k+1}) \leq \phi_{k+1}^*$ and

ii) α_k as large as possible.

$$\phi_k(x) = \phi_k^* + \frac{d_k \cdot L}{2} \|x - v_k\|^2.$$

Let us say we update ϕ_{k+1} as

$$\begin{aligned} \phi_{k+1}(x) = & (1 - \alpha_k) \left[\phi_k^* + \frac{d_k \cdot L}{2} \|x - v_k\|^2 \right] \\ & + \alpha_k \left[f(y_k) + \langle \nabla f(y_k), x - y_k \rangle \right] \end{aligned}$$

Lemma: $\phi_{k+1}^* =$ see next page.

$$v_{k+1}^* = v_k - \frac{\alpha_k}{d_k(1-\alpha_k) \cdot L} \nabla f(y_k).$$

Proof: The terms involving x can be written as

$$\begin{aligned} & \frac{d_{k+1} \cdot L}{2} \left[\|x\|^2 - 2 \langle x, v_k \rangle + \frac{2\alpha_k}{d_{k+1} \cdot L} \langle x, \nabla f(y_k) \rangle \right] \\ & = \frac{d_{k+1} \cdot L}{2} \left[\|x\|^2 - 2 \langle x, v_k - \frac{\alpha_k}{d_k(1-\alpha_k) \cdot L} \nabla f(y_k) \rangle \right]. \end{aligned}$$

$$\text{So, } v_{k+1} = v_k - \frac{\alpha_k}{d_k(1-\alpha_k) \cdot L} \nabla f(y_k).$$

$$\begin{aligned}
 \text{And } \Phi_{k+1}^* &= (1-\alpha_k) \Phi_k^* + \frac{\alpha_k(1-\alpha_k) \cdot L}{2} \|v_k\|^2 \\
 &\quad + \alpha_k f(y_k) - \alpha_k \langle \nabla f(y_k), y_k \rangle \\
 &\quad - \frac{\alpha_k(1-\alpha_k) \cdot L}{2} \|v_k - \frac{\alpha_k}{\alpha_k(1-\alpha_k) \cdot L} \nabla f(y_k)\|^2 \\
 &= (1-\alpha_k) \Phi_k^* + \alpha_k f(y_k) - \frac{\alpha_k^2}{2\alpha_k(1-\alpha_k) \cdot L} \|\nabla f(y_k)\|^2 \\
 &\quad + \alpha_k \langle \nabla f(y_k), v_k - y_k \rangle. \quad \square
 \end{aligned}$$

Given x_k, y_k and v_k , one way to get x_{k+1} s.t. $f(x_{k+1}) \leq \Phi_{k+1}^*$, is to try GD step at one of these points.

In fact if $x_{k+1} = y_k - \eta \nabla f(y_k)$ with $\eta \leq \frac{1}{L}$,

$$\text{then } f(x_{k+1}) \leq f(y_k) - \frac{\eta^2}{2} \|\nabla f(y_k)\|^2.$$

$$\text{So, } f(x_{k+1}) - \Phi_{k+1}^*$$

$$\leq f(y_k) - \frac{\eta^2}{2} \|\nabla f(y_k)\|^2 - (1-\alpha_k) f(x_k)$$

$$- \alpha_k f(y_k) + \frac{\alpha_k^2}{2\alpha_k(1-\alpha_k)L} \|\nabla f(y_k)\|^2 - \alpha_k \langle \nabla f(y_k), v_k - y_k \rangle$$

$$\leq (1-\alpha_k) [f(y_k) - f(x_k)]$$

$$- \alpha_k \langle \nabla f(y_k), v_k - y_k \rangle.$$

$$\eta \geq \frac{\alpha_k^2}{L \cdot \alpha_k(1-\alpha_k)}$$

$$\leq (1-\alpha_k) \langle \nabla f(y_k), y_k - x_k \rangle - \alpha_k \langle \nabla f(y_k), v_k - y_k \rangle$$

$$= \langle \nabla f(y_k), y_k - (1-\alpha_k)x_k - \alpha_k v_k \rangle.$$

Remember, we have a choice in choosing y_k .

So, choose $y_k = (1-\alpha_k)x_k + \alpha_k v_k$.

$$\Rightarrow f(x_{k+1}) \leq \underline{\underline{\phi_{k+1}^*}}.$$

What are the conditions? $\frac{\alpha_k^2}{d_k(1-\alpha_k)} \leq \frac{1}{2}$.

Eventually, we want the smallest d_k .

$\Rightarrow \alpha_k$ should be as large as possible.

So, choose α_k as the largest root of

$$\frac{\alpha_k^2}{d_k(1-\alpha_k)} = \frac{1}{2} \quad \text{and} \quad \alpha_k \leq \frac{1}{2}.$$

Let $a_k = \frac{1}{\sqrt{d_k}}$. Then

$$a_{k+1} - a_k = \frac{\sqrt{d_k} - \sqrt{d_{k+1}}}{(\sqrt{d_k} \cdot \sqrt{d_{k+1}})} = \frac{d_k - d_{k+1}}{\sqrt{d_k d_{k+1}} (\sqrt{d_k} + \sqrt{d_{k+1}})}$$

$$\Rightarrow \frac{a_k}{2\sqrt{d_{k+1}}} = \frac{1}{2} \cdot \frac{1}{2}$$

$$\Rightarrow \frac{1}{\sqrt{d_{k+1}}} \geq \frac{k}{2} \frac{1}{\sqrt{k}} \Rightarrow d_{k+1} \leq \frac{4\epsilon^2}{k^2}$$

Plugging this in the previous lemma gives us

$$\begin{aligned} f(x_k) - f^* &\leq \frac{4\epsilon^2}{k^2} [\Phi_0(x^*) - f(x^*)] \\ &= \frac{4}{k^2} \left[f(x_0) - f(x^*) + \frac{L}{2} \|x_0 - x^*\|^2 \right]. \end{aligned}$$

□

For linear least squares problem $f(x) = \frac{1}{2} \|Ax - b\|^2$,

$$\text{we get } f(x_k) - f^* \leq \frac{4}{k^2} \cdot \|A\|^2 \cdot \|x_0 - x^*\|^2$$

$$\text{Compare to GD: } \frac{6 \|A\|^2 \|x_0 - x^*\|^2}{k}$$

$$\frac{1}{k} \quad \text{vs} \quad \frac{1}{k^2}$$

By the lower bound argument, we also know this is optimal for gradient methods.

STRONG CONVEXITY

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

$\mu \rightarrow$ Strong convexity parameters.

Smooth + Strongly Convex

~~GD~~ ~~GD~~

By defn.: $\langle \nabla f(x), \frac{y-x}{\|y-x\|} \rangle \leq f(x) - f(y) - \frac{\mu}{2} \|x-y\|^2$

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

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

Smoothness $\Rightarrow \eta(1-\eta L) \|\nabla f(x_t)\|^2 \leq f(x_t) - f(x^*)$

Substituting, we have,

$$\|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 - \eta \left[\frac{2 - \eta L}{1 - \eta L} \right] [f(x_t) - f(x^*)] - \eta \mu \|x_t - x^*\|^2$$

For $\eta \leq \frac{1}{2L}$, we have

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

So, $f(x_{t+1}) - f(x^*) \leq \frac{L}{2} \|x_{t+1} - x^*\|^2$

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

$$\text{So, } f(x_t) - f(x^*) \leq K \left(1 - \frac{1}{K}\right)^t [f(x_0) - f(x^*)].$$

$$\boxed{O\left(K \log \frac{K}{\epsilon}\right)} \leftarrow \text{Important factors.}$$

Can we improve? YES using AGD.

Reduction: Apply AGD

$$f(x_t) - f(x^*) \leq \frac{4L \|x_0 - x^*\|^2}{t^2}$$

$$\leq \frac{4K [f(x_0) - f(x^*)]}{t^2}$$

After $t = 4\sqrt{K}$ steps, we have

ALSO TALK

LOWER

BOUND

$$f(x_t) - f(x^*) \leq \frac{1}{4} [f(x_0) - f(x^*)]$$

For ϵ -accuracy, we need $O\left(\sqrt{K} \log \frac{1}{\epsilon}\right)$ iterations.

Improvement from K to \sqrt{K} .

Can also design a direct algorithm (without restarts).

~~NON SMOOTH + STRONG CONV~~

For least squares $f(x) = \frac{1}{2} \|Ax - b\|^2$.

$$K = \frac{\sigma_{\max}(AA^T)}{\sigma_{\min}(AA^T)} = K(A)^2.$$

$$\text{AGD: } f(x_t) - f(x^*) \leq \exp\left(-\frac{t}{\sqrt{K}}\right) \cdot [f(x_0) - f(x^*)].$$

NON-SMOOTH + STRONGLY CONVEX

Recall + LIPSCHITZ [$\|\nabla f\| \leq G$].

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

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

Recall by strong convexity

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

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

$$\Rightarrow f(x_t) - f(x^*) \leq \frac{1}{2} \left[\frac{(1 - \eta\mu) \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2}{\eta} + \eta G^2 \right]$$

Choose $\eta_t = \frac{1}{\mu t}$.

$$\Rightarrow f(x_t) - f(x^*) \leq \frac{1}{2\mu} \left[\mu(t-1) \|x_t - x^*\|^2 - \mu t \|x_{t+1} - x^*\|^2 \right] + \frac{G^2}{2\mu t}.$$

$$\Rightarrow \sum_{t=1}^T f(x_t) - f(x^*) \leq \frac{G^2}{2\mu} \sum_{t=1}^T \frac{1}{t} \leq \frac{G^2 \log T}{2\mu}.$$

$$\Rightarrow f(\bar{x}_T) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x^*)$$

$$\leq \frac{G^2 \log T}{2 \mu T}$$

Lower bounds: Consider a non smooth fn $f(x)$

Consider $f_\epsilon(x) = f(x) + \frac{\epsilon}{2} \|x - x_0\|^2$.

If f has G -Lipschitz then in $B_D(x_0)$,

f_ϵ has $G + \epsilon D$ -Lipschitz.

Furthermore, $f(\hat{x}) \leq f(\hat{x}) + \frac{\epsilon}{2} \|\hat{x} - x_0\|^2$

$$\geq f(x^*) + \frac{\epsilon}{2} \|x^* - x_0\|^2$$

$$\geq f^* + \frac{\epsilon D^2}{2}$$

for $\hat{x} \triangleq \arg \min_x f_\epsilon(x)$

and $f(\hat{x}) \leq f(x^*) + \frac{\epsilon}{2} \|x^* - x_0\|^2$

$$\leq f^* + \frac{\epsilon}{2} D^2$$

If for strongly convex, we have suboptimality $\frac{1}{T^\alpha}$ then, by choosing $\epsilon \propto T^{-\alpha/2}$, we have

Suboptimality $\frac{\epsilon D^2}{2} + \frac{G^2 \log T}{2 \cdot \epsilon \cdot T^\alpha}$

$$= \left[\frac{D^2}{2} + \frac{G^2 \log T}{2} \right] \cdot T^{-\alpha/2}$$

So, $O(\frac{1}{\gamma})$ is the optimal rate achievable for non-smooth + strongly convex functions.

PROJECTION & GRADIENT MAPPING

$\min_{x \in X} f(x)$ X is a simple ^{convex} set.
E.g., $\mathbb{R}_{\geq 0}^n$.

Simple means that we have access to a projection oracle.

Given a point w , $P_x(w) = \arg \min_{x \in X} \|x - w\|^2$.

Examples: i) $X = B_1(0)$: $P_x(w) = w$ if $\|w\| \leq 1$
 $= \frac{w}{\|w\|}$ if $\|w\| > 1$.

ii) $X = \mathbb{R}_{\geq 0}^n$: $P_x(w) = w \odot \mathbb{1}_{\{w \geq 0\}}$.

Pythagoras theorem: $\langle w - P_x(w), y - P_x(w) \rangle \leq 0$
and $\|w - y\|^2 \geq \|P_x(w) - y\|^2$ $\forall y \in X$.

Proof: ~~$\|y - w\|^2 \geq \|P_x(w) - w\|^2$~~

$\Rightarrow \|y - P_x(w)\|^2 = 2$ Since X is a convex set,
 $P_x(w) + \alpha(y - P_x(w)) \in X$ for $\alpha \in [0, 1]$.

$$\|w - P_X(w)\|^2 \leq \|w - P_X(w) - \alpha(y - P_X(w))\|^2$$

$$= \|w - P_X(w)\|^2 + \alpha^2 \|y - P_X(w)\|^2$$

$$- 2\alpha \langle w - P_X(w), y - P_X(w) \rangle$$

$$\Rightarrow \langle w - P_X(w), y - P_X(w) \rangle \leq \frac{\alpha}{2} \|y - P_X(w)\| \quad \forall \alpha > 0.$$

Taking $\alpha \rightarrow 0$ finishes the proof. \square

~~NON SMOOTH GD~~

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

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

$$= \|x_t - x^*\|^2 - \eta \langle \nabla f(x_t), x_t - x^* \rangle$$

$$+ \eta^2 \|\nabla f(x_t)\|^2$$

Same argument: $\frac{1}{T} \sum_{t=1}^T f(x_t) - f(x^*)$

$$\leq \frac{\|x_0 - x^*\|^2}{\eta T} + \eta G^2.$$

SMOOTH GD [Gradient mapping]

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

$$= \arg \min_{x \in X} f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|^2$$

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

Use $\hat{f}_\eta(\cdot; x)$

$$\text{So, we have } x_{t+1} = \arg \min_{x \in X} \hat{f}_\eta(x_t; x)$$

Since $\hat{f}_\eta(x_t; \cdot)$ is $\frac{1}{\eta}$ strongly convex, we have

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

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

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

~~$$= \hat{f}_\eta(x_t; x) - \frac{1}{2\eta} \|x - x_t\|^2$$~~

~~$$\geq f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle$$~~

~~$$+ \frac{1}{2\eta} \|x_{t+1} - x_t\|^2 + \frac{1}{2\eta} \|x - x_{t+1}\|^2 - \frac{1}{2\eta} \|x - x_t\|^2$$~~

$$\begin{aligned} \hat{f}_\eta(x_t; x) &\geq \hat{f}_\eta(x_t; x_{t+1}) + \frac{1}{2\eta} \|x - x_{t+1}\|^2 \\ &\geq f(x_{t+1}) + \frac{1}{2\eta} \|x - x_{t+1}\|^2. \end{aligned}$$

Also,

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

$$= \hat{f}_\eta(x_t; x) - \frac{1}{2\eta} \|x - x_t\|^2$$

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

$$\text{Firstly, } f(x_t) \geq f(x_{t+1}) + \frac{1}{2\eta} \|x_t - x_{t+1}\|^2$$

Secondly,

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

Taking a telescopic sum and averaging, we have

$$\frac{1}{T} \sum_{t=1}^T f(x_t) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2\eta T}.$$

Since $\eta = \frac{1}{L}$, we have

$$f(\bar{x}_T) - f(x^*) \leq \frac{L \|x_0 - x^*\|^2}{2T}. \quad \square$$

We can similarly extend the accelerated gradient descent algorithm also to the constrained setting with simple constraints. This gives convergence rate of $O(\frac{1}{T^2})$.

NON-EUCLIDEAN SETTINGS

Can we obtain similar results when the Lipschitz constant, Smoothness, Strong Convexity etc. are all w.r.t. to a non-Euclidean norm, say $\|\cdot\|_1$?

E.g.,
$$\min_{x: \|x\|_1 \leq 1} \frac{1}{2} x^T A x - x^T b$$

where $\|A\|_\infty$ and $\|b\|_\infty$ are bounded.

If we use the previous results, the

Lipschitz parameter = $\max_x \|Ax - b\|$ Could be \sqrt{d}

Smoothness parameter = $\|A\|$ Could be d .

Can we avoid such dependence?

This is not just looseness of the results. The actual behavior of GD is this.

How can we change the algorithm depending on the structure of the problem?

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

↓

$$x_{t+1} = \arg \min_x f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta} \|x - x_t\|_2^2$$

- ① Computing such projections could be hard in general since the squared term couples diff. coordinates.
- ② What would be the potential function that we can use to track the progress of GD?

BREGMAN DIVERGENCES

Given a fn. ϕ which is μ -strongly convex w.r.t. $\|\cdot\|$ i.e., $\phi(y) \geq \phi(x) + \langle \nabla \phi(x), y - x \rangle + \frac{\mu}{2} \|x - y\|^2$,

Bregman divergence

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

Since ϕ is 1-strongly convex, $\mathcal{D}_\phi(\cdot, y)$ is also 1-strongly convex.

$$\text{MD: } x_{t+1} = \arg \min_{x \in \mathcal{X}} f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{\eta} \mathcal{D}_\phi(x; x_t)$$

So, we have We will first consider the non-smooth setting: $\|\nabla f(x)\|_* \leq G$ Dual norm.

~~$$\begin{aligned} \mathcal{D}_\phi(x; x_{t+1}) &= \phi(x) - \langle \nabla \phi(x_{t+1}), x - x_{t+1} \rangle - \phi(x_{t+1}) \\ &= \phi(x) - \phi(x_{t+1}) - \langle \nabla \phi(x_{t+1}), x - x_{t+1} \rangle \\ &= \phi(x) - \phi(x_t) + \frac{1}{2} \|x_t - x_{t+1}\|^2 - \langle \nabla \phi(x_{t+1}), x - x_t \rangle. \end{aligned}$$~~

By the update rule for x_{t+1} , we have

$$\langle \nabla f(x_t), x \rangle + \frac{1}{\eta} \nabla \phi(x_{t+1}) - \frac{1}{\eta} \nabla \phi(x_t), x - x_{t+1} \rangle \geq 0 \quad \forall x \in \mathcal{X}$$

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

$$\text{So, } \mathcal{D}_\phi(x; x_{t+1}) = \phi(x) - \phi(x_{t+1}) - \langle \nabla \phi(x_{t+1}), x - x_{t+1} \rangle$$

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

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

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

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

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

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

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

$$= \mathcal{D}_\phi(x; x_t) + \eta \langle \nabla f(x_t), x - x_t \rangle$$

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

$$\leq \mathcal{D}_\phi(x; x_t) - \eta [f(x_t) - f(x)]$$

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

Since $x_{t+1} = \underset{x \in \mathcal{X}}{\operatorname{argmin}} f(x_t) + \langle \nabla f(x_t), x - x_t \rangle$

$$+ \frac{1}{2\eta} \mathcal{D}_\phi(x; x_t)$$

we have

$$\langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{1}{2\eta} D_\phi(x_{t+1}; x_t) \leq 0.$$

$$\begin{aligned} \Rightarrow D_\phi(x_{t+1}; x_t) &\leq -2\eta \langle \nabla f(x_t), x_{t+1} - x_t \rangle \\ &\leq 2\eta \|\nabla f(x_t)\|_* \cdot \|x_{t+1} - x_t\|. \end{aligned}$$

since $D_\phi(\cdot; x_t)$ is $\frac{1}{2}$ -strongly convex, we have

$$\frac{1}{2} \|x_{t+1} - x_t\|^2 \leq \eta \cdot G \cdot \|x_{t+1} - x_t\|.$$

$$\text{So, } \|x_{t+1} - x_t\| \leq 2\eta G.$$

Plugging this back, we obtain

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

Taking a telescopic sum, we get.

$$\frac{1}{T} \sum_{t=1}^T f(x_t) - f(x) \leq \frac{D_\phi(x; x_0)}{\eta T} + 2\eta G^2$$

Taking $\eta = \sqrt{\frac{D_\phi(x; x_0)}{2GT}}$ and by convexity, we have

$$f(\bar{x}_T) - f(x) \leq G \sqrt{D_\phi(x; x_0)} \cdot \sqrt{\frac{2}{T}}. \quad \square$$

Example:- let $f(x) = \frac{1}{2} x^T A x - x^T b$. over ~~simplex~~ ^{simplex}.

$$\|A\|_\infty \leq 1 \text{ and } \|b\|_\infty \leq 1$$

$$\begin{aligned} x_i &\geq 0 \\ \sum x_i &= 1 \end{aligned}$$

$$\Rightarrow \|\nabla f(x)\| = \|Ax - b\|_\infty \leq 2.$$

Take $\phi(x) = \sum_i x_i \log \frac{x_i}{y_i} = -H(x)$
 \downarrow
 entropy.

$$\begin{aligned} D_\phi(x; y) &= \sum_i x_i \log \frac{x_i}{y_i} - \sum y_i \log y_i \\ &\quad - \sum (\log \frac{x_i}{y_i} - 1) (x_i - y_i) \end{aligned}$$

$$= \sum_i x_i \log \frac{x_i}{y_i} = KL(x \| y).$$

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

Gradient step:

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

$$= \operatorname{argmin}_{x \in \mathcal{X}} \langle w, x \rangle + \sum_i x_i \log \frac{x_i}{x_{t,i}}$$

\downarrow
 $\cong \eta \nabla f(x_t)$

$$= \operatorname{argmin}_{x \in \mathcal{X}} \sum_i x_i \log \frac{x_i}{\sum_i x_i \exp(-w_i)}$$

$$\text{let } z = \frac{x_t \odot \exp(-w)}{\sum_i x_{t,i} \exp(-w_i)} \rightarrow \alpha$$

$$x_{t+1} = \operatorname{argmin}_{x \in \mathcal{X}} \sum_i x_i \log \frac{x_i}{z_i} - x_i \log \alpha$$

$$= \operatorname{argmin}_{x \in \mathcal{X}} \text{KL}(x \| z) - \log \alpha \quad [\because \sum x_i = 1]$$

$$= z. \quad [\because z \in \mathcal{X}].$$

So, the gradient update is

$$x_{t+1} = \frac{x_t \odot \exp(-\eta \nabla f(x_t))}{\sum_i x_{t,i} \exp(-\eta \nabla f(x_{t,i}))}$$

Our result says,

$$f(\bar{x}_T) - f(x^*) \leq \frac{\sqrt{8D_\Phi(x^*; x_0)}}{\sqrt{T}}$$

If we choose $x_0 = \begin{bmatrix} 1/d \\ \vdots \\ 1/d \end{bmatrix}$ then $D_\Phi(x^*; x_0)$

$$= \text{KL}(x^*; \begin{bmatrix} 1/d \\ \vdots \\ 1/d \end{bmatrix})$$

$$= \log d - H(x^*) \leq \log d.$$

$$\text{So, } f(\bar{x}_T) - f(x^*) \leq \sqrt{\frac{8 \log d}{T}}$$

Compare with ~~gr~~ Euclidean gradient descent.

$$f(\bar{x}_T) - f(x^*) \leq \frac{[\|A\|_{1 \rightarrow 2} \|b\|]}{\sqrt{T}}$$

If $\|A\|_{\infty} \leq 1$, it is still possible that

$\|A\|_{1 \rightarrow 2}$ could still be \sqrt{d}

and $\|b\|_{\infty} \leq 1$, it is still possible that $\|b\| \approx \sqrt{d}$.

MIRROR DESCENT + SMOOTHNESS

Smoothness: $\|\nabla f(x) - \nabla f(y)\|_* \leq L \|x - y\|$

Dual norm ← Original norm.

~~$$x_{t+1} = \underset{x \in X}{\operatorname{argmin}} f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{\eta} D_{\Phi}(x; x_t)$$~~

Firstly ~~$f(x_{t+1}) \leq f(x_t) +$~~

Smoothness again gives us a quadratic upper bound.

$$\begin{aligned}
 f(y) &= f(x) + \int_{\alpha=0}^1 \langle \nabla f(x + \alpha(y-x)), y-x \rangle d\alpha \\
 &= f(x) + \langle \nabla f(x), y-x \rangle + \int_{\alpha=0}^1 \langle \nabla f(x + \alpha(y-x)) - \nabla f(x), y-x \rangle d\alpha \\
 &\leq f(x) + \langle \nabla f(x), y-x \rangle + \int_{\alpha=0}^1 \|\nabla f(x + \alpha(y-x)) - \nabla f(x)\|_* \|y-x\| d\alpha \\
 &\stackrel{[\text{SMOOTHNESS}]}{\leq} f(x) + \langle \nabla f(x), y-x \rangle + \int_{\alpha=0}^1 L \cdot \alpha \|y-x\| \cdot \|y-x\| d\alpha \\
 &= f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|x-y\|^2.
 \end{aligned}$$

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

we have

$$f(x_t) \stackrel{[\text{S.C.}]}{\leq} f(x_{t+1}) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{1}{2\eta} D_{\phi}(x_{t+1}; x_t)$$

$$\stackrel{[\text{S.C.}]}{\leq} f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{1}{2\eta} \|x_t - x_{t+1}\|^2$$

$$\geq f(x_{t+1}) + \left(\frac{1}{2\eta} - L\right) \|x_t - x_{t+1}\|^2$$

So, mirror descent is a descent algorithm.

Furthermore, from previous computations, we have

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

By optimality defn. of x_{t+1} , and s.c. of D_{ϕ} and smoothness, we have

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

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

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

$$\leq \mathcal{D}_\phi(x; x_t) - \eta [f(x_t) - f(x)]$$

$$- \eta [f(x_t) - f(x_{t+1})]$$

$$\Phi = \mathcal{D}_\phi(x; x_t) + \eta \langle \nabla f(x_t), x_t - x_{t+1} \rangle$$

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

$$\leq \mathcal{D}_\phi(x; x_t) - \eta [f(x_t) - f(x)]$$

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

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

$$\leq \mathcal{D}_\phi(x; x_t) - \eta [f(x_t) - f(x)]$$

$$- \eta [f(x_t) - f(x_{t+1})]$$

By ~~defn of x_{t+1}~~ and Strong Convexity of Φ , we have

$$\langle \nabla \Phi(x_t), x_{t+1} - x_t \rangle \geq \Phi(x_{t+1}) - \Phi(x_t) + \langle \nabla \Phi(x_t), x_{t+1} - x_t \rangle + \frac{1}{2} \|x_t - x_{t+1}\|^2$$

So, from the previous computations, we have

$$\begin{aligned} \mathcal{D}_\phi(x; x_{t+1}) &= \phi(x) - \phi(x_{t+1}) - \langle \nabla \phi(x_{t+1}), x - x_{t+1} \rangle \\ &\leq \phi(x) - \phi(x_t) - \langle \nabla \phi(x_t), x_{t+1} - x_t \rangle - \frac{1}{2} \|x_t - x_{t+1}\|^2 \\ &\quad - \langle \nabla \phi(x_{t+1}), x - x_{t+1} \rangle \end{aligned}$$

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

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

$$= \mathcal{D}_\phi(x; x_t) + \eta \langle \nabla f(x_t), x - x_{t+1} \rangle - \frac{1}{2} \|x_t - x_{t+1}\|^2.$$

$$\begin{aligned} &= \mathcal{D}_\phi(x; x_t) + \eta \langle \nabla f(x_{t+1}), x - x_{t+1} \rangle - \frac{1}{2} \|x_t - x_{t+1}\|^2 \\ &\quad + \eta \langle \nabla f(x_t) - \nabla f(x_{t+1}), x_t - x_{t+1} \rangle. \end{aligned}$$

$$\text{Now, } \langle \nabla f(x_t), x_t - x_{t+1} \rangle$$

$$\begin{aligned} &= \underbrace{\langle \nabla f(x_{t+1}), x_t - x_{t+1} \rangle}_{\leq f(x_{t+1}) - f(x_t)} + \underbrace{\langle \nabla f(x_t) - \nabla f(x_{t+1}), x_t - x_{t+1} \rangle}_{\leq \frac{L}{2} \|x_t - x_{t+1}\|^2} \\ &\leq 0 \end{aligned}$$

$$\text{So, } D_{\phi}(x; x_{t+1}) \leq D_{\phi}(x; x_t) - \eta [f(x_t) - f(x)] - \left(\frac{1}{2} - \eta L\right) \|x_t - x_{t+1}\|^2.$$

~~If we~~ Since $\eta < \frac{1}{2L}$, we have

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

Telescoping and averaging gives us

$$\frac{1}{T} \sum_{t=1}^T f(x_t) - f(x) \leq \frac{D_{\phi}(x; x_0)}{\eta T}. \quad \square$$

In fact we can extend all the Euclidean algorithms to the non-Euclidean setting

- Gradient descent
- Accelerated gradient descent
- Faster rates with Strong Convexity, Smoothness etc.

FASTER ALGORITHMS FOR STRUCTURED NON-SMOOTH PROBLEMS

$$\min_x \frac{1}{2} x^T A x + x^T b + \|x\|_1$$

$$\text{LASSO: } \min_x \frac{1}{2} \|Ax - b\|^2 + \underbrace{\lambda}_{=1} \|x\|_1$$

The function $f(x) = \frac{1}{2} \|Ax - b\|^2 + \|x\|_1$ is not smooth. So, the best rate we get using GD is $O(\frac{1}{\sqrt{T}})$. Can we do better?

Main idea: The nonsmooth part $\|x\|_1$ is very simple.

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

Notice $x \leftarrow$
not x_t

$$\text{Implementability: } x_{t+1} = \arg \min_x \langle \nabla g(x_t), x \rangle + \|x\|_1 + \frac{1}{2\eta} \|x - x_t\|^2$$

Decouples into 'd' independent problems, one for each direction.

$$i: \arg \min_{x_i} \nabla g(x_t)_i \cdot x_i + |x_i| + \frac{1}{2\eta} (x_i - x_{t,i})^2$$

$$= \arg \min_{x_i} \frac{1}{2\eta} (x_i - x_{t,i} + \eta \nabla g(x_t)_i)^2 + |x_i|$$

$$\Rightarrow x_{t+1,i}^* = \frac{1}{1+\eta} [x_{t,i} - \eta \nabla g(x_t)_i]$$

$$\Rightarrow x_{t+1} = \frac{1}{1+\eta} (x_t - \eta \nabla g(x_t))$$

Convergence rate:

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

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

$$\geq g(x_t) + \langle \nabla g(x_t), x - x_t \rangle + h(x)$$

$$= \hat{f}_\eta(x; x_t) - \frac{1}{2\eta} \|x - x_t\|^2$$

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

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

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

Reorganizing, we have

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

$$\text{So, } \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x) \leq \frac{\|x - x_0\|^2}{2\eta T}.$$

$$\Rightarrow f(\bar{x}_T) - f(x) \leq \frac{2L\|x - x_0\|^2}{T}.$$

$$[\text{Since } \eta = \frac{1}{2L}]$$

Can again extend Nesterov's AGD to get a rate of $O(\frac{1}{\sqrt{T}})$.

STOCHASTIC ALGORITHMS

Oracle models: So far we assumed we had access to EXACT gradients of the function.
 $x \rightarrow \nabla f(x)$.

In many cases this may not be reasonable

For e.g., in machine learning

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

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

↓

Cannot compute expectations.

But if we sample (a_i, b_i) and compute

$$\hat{\nabla} f(x) = (a_i^T x - b_i) a_i \quad \text{then} \quad \mathbb{E}[\hat{\nabla} f(x)] = \mathbb{E}[(a_i^T x - b_i) a_i] \\ = \nabla f(x).$$

So, we may assume we have access to $\hat{\nabla} f(x)$ [Stochastic] s.t. $\mathbb{E}[\hat{\nabla} f(x)] = \nabla f(x)$.

The simplest setting, from an analytical point of view is if we assume independent and bounded variance noise.

$$\mathbb{E}[\|\hat{\nabla} f(x) - \nabla f(x)\|^2] \leq \sigma^2 \quad \forall x.$$

Stochastic gradient descent [Robbins & Monro 1952]

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

$$\mathbb{E}[\|x_{t+1} - x\|^2] = \mathbb{E}[\|x_t - x\|^2] - 2\eta_t \mathbb{E}[\langle \hat{\nabla} f(x_t), x_t - x \rangle] \\ + \eta_t^2 \mathbb{E}[\|\hat{\nabla} f(x_t)\|^2] \\ \leq \mathbb{E}[\|x_t - x^*\|^2] - 2\eta_t \mathbb{E}[\langle \nabla f(x_t), x_t - x \rangle] \\ + \eta_t^2 [G^2 + \sigma^2].$$

$$\leq \mathbb{E}[\|x_t - x\|^2] - 2\eta_t \mathbb{E}[f(x_t) - f(x)] + \eta_t^2 [G^2 + \sigma^2].$$

Rearranging and telescoping,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[f(x_t)] - f(x) \leq \frac{\mathbb{E}[\|x_0 - x\|^2]}{\eta T} + \eta [G^2 + \sigma^2]$$

$$\exists \eta = \frac{\|x_0 - x\|}{\sqrt{T} \sqrt{G^2 + \sigma^2}}, \text{ then}$$

$$\mathbb{E}[f(\bar{x}_T)] - f(x) \leq \frac{2 \cdot \sqrt{G^2 + \sigma^2} \cdot \|x_0 - x\|}{\sqrt{T}}.$$

STRONG CONVEXITY

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

$$\Rightarrow \mathbb{E}[\|x_{t+1} - x\|^2] \leq (1 - \eta_t \mu) \mathbb{E}[\|x_t - x\|^2] - 2\eta_t \mathbb{E}[f(x_t) - f(x)] + \eta_t^2 [G^2 + \sigma^2]$$

$$\text{Choosing } \eta_t = \frac{1}{\mu t}.$$

$$2 \mathbb{E}[f(x_t) - f(x)] \leq \mu(t-1) \mathbb{E}[\|x_t - x\|^2] - \mu t \mathbb{E}[\|x_{t+1} - x\|^2] + \eta_t [G^2 + \sigma^2]$$

$$\Rightarrow \frac{1}{T} \sum_{t=1}^T \mathbb{E}[f(x_t) - f(x)] \leq \frac{(G^2 + \sigma^2) \log T}{\mu T}$$

With Smoothness: (No S.C.)

We already have

$$\mathbb{E}[\|x_{t+1} - x\|^2] \leq \mathbb{E}[\|x_t - x\|^2] - 2\eta \mathbb{E}[f(x_t) - f(x)] + \eta^2 [\mathbb{E}[\|\nabla f(x_t)\|^2] + \sigma^2]$$

By Smoothness, $\|\nabla f(x_t)\|^2 \leq L \cdot [f(x_t) - f(x)]$.

$$\text{So, } \mathbb{E}[\|x_{t+1} - x\|^2] \leq \mathbb{E}[\|x_t - x\|^2] - \eta(2 - \eta L) \mathbb{E}[f(x_t) - f(x)] + \eta^2 \sigma^2.$$

Telescoping and averaging,

$$\frac{1}{T} \mathbb{E}\left[\sum_{t=1}^T f(x_t)\right] - f(x) \leq \frac{\|x_0 - x\|^2}{\eta(2 - \eta L)T} + \frac{\eta \sigma^2}{2 - \eta L}$$

$$\eta \leq \min\left\{\frac{1}{L}, \frac{1}{\sigma \|x_0 - x^*\| \sqrt{T}}\right\} \Rightarrow \frac{2\sigma \|x_0 - x\|}{(\cancel{L}) \sqrt{T}}$$

STOCHASTIC MIRROR DESCENT [TALK ABOUT IT]

Is additive independent noise ~~a good~~ with bounded variance a good assumption?

* LOT of work in the Stochastic Approximation literature for more general noise i.e., Markov structure and noise not growing fast.

E.g.,
$$\mathbb{E}[\|\nabla f(x) - \hat{\nabla} f(x)\|^2] \leq L \|x\|^2.$$

* However the bounds are usually not very fully non-asymptotic [dependence on all the parameters not clear].

Primary reason for the resurgence of Stochastic algorithms in large scale machine learning

FAST COMPUTATION.

Empirical risk minimization

$$f(x) = \frac{1}{2n} \sum_{i=1}^n (a_i^T x - b_i)^2. \quad a_i, x \in \mathbb{R}^d.$$

We would like to $\min_x f(x)$.

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

One gradient computation : $O(nd)$ time.

Total number of GD steps : $O(K \log(\frac{1}{\epsilon}))$



$$\text{Condition number} = \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)}$$

Total time : $O(ndK \log \frac{1}{\epsilon})$.

Can we do better?

$$\text{SGD: } \hat{\nabla} f(x) = (a_i^T x - b_i) a_i$$

$$\text{Variance: } \mathbb{E}_i [\|\hat{\nabla} f(x) - \nabla f(x)\|^2] = \mathbb{E}_i \left[\left\| (a_i^T x - b_i) a_i - \frac{1}{n} \sum_{i=1}^n (a_i^T x - b_i) a_i \right\|^2 \right]$$

$$\sigma^2 \approx \left\{ \max_i \|a_i\|^2 \right\} \|x - x^*\|^2 + f(x^*) \cdot \max_i \|a_i\|^2$$

SGD: One gradient computation : $O(d)$ time

$$\# \text{ SGD steps : } \frac{\sigma^2}{\epsilon^2}$$

$$ndK \log \frac{1}{\epsilon}$$

$$\frac{nd}{\sqrt{\epsilon}}$$

vs

$$\frac{d\sigma^2}{\epsilon^2}$$

↳ n does not appear.

They have tradeoffs between various parameters.

SVRG: Use SGD but with decreasing variance. \rightarrow How?

$$\text{Let } f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$$

Let each $f_i(x)$ be L -smooth: $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L \|x - y\|$.

Suppose further that $f(\cdot)$ is μ -strongly convex.

$$\text{Compute } \nabla f(x_0) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_0).$$

Instead of $x_{t+1} = x_t - \eta \nabla f_i(x_t)$, do

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

$$\mathbb{E}[\cdot] = \nabla f(x_t).$$

$$\text{Var} = \mathbb{E}_i \left[\left\| \nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) - \nabla f(x_t) \right\|^2 \right]$$

$$\# = \mathbb{E}_i \left[\left\| \nabla f_i(x_t) - \nabla f_i(x_0) \right\|^2 \right] + \mathbb{E} \left[\left\| \nabla f(x_0) - \nabla f(x_t) \right\|^2 \right]$$

$$\leq 2 \cdot L \cdot \|\bar{x}_0 - x_t\|^2$$

$$\begin{aligned} \mathbb{E}[\|x_{t+1} - x^*\|^2] &= \mathbb{E}\left[\|x_t - \gamma \begin{pmatrix} -\nabla f_i(x_0) \\ \nabla f_i(x_t) \\ + \nabla f(x_0) \end{pmatrix} - x^*\|^2\right] \\ &= \mathbb{E}[\|x_t - x^*\|^2] - 2\gamma \mathbb{E}[\langle \nabla f_i(x_t), x_t - x^* \rangle] \\ &\quad + \gamma^2 \mathbb{E}[\|\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0)\|^2] \end{aligned}$$

$$\begin{aligned} &\leq \mathbb{E}[\|x_t - x^*\|^2] - 2\gamma \left\{ \mathbb{E}[f(x_t)] - f(x^*) \right\} \\ &\quad + \gamma^2 L \left\{ \mathbb{E}[f(x_0)] - f(x^*) \right\} \\ &\quad + 2\gamma^2 L \|\bar{x}_0 - x_t\|^2 \cdot \boxed{?} \end{aligned}$$

$$\mathbb{E}[\|\nabla f_i(x_t) - \nabla f_i(x_0)\|^2]$$

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

$$\mathbb{E}[\|\nabla f_i(x_t) - \nabla f_i(x^*)\|^2] \leq ?$$

Consider $g_i(x) = f_i(x) - \langle \nabla f_i(x^*), x \rangle$

$$\Rightarrow g_i(x^*) \leq g_i(x_t - \gamma \nabla f_i(x_t) + \gamma \nabla f_i(x^*))$$

$$\begin{aligned} &\leq \cancel{f_i(x_t - \gamma \nabla f_i(x_t)) - \langle \nabla f_i(x^*), x_t - \gamma \nabla f_i(x_t) \rangle} \\ &\leq g_i(x_t) - \gamma(1 - \gamma L) \|\nabla f_i(x_t) - \nabla f_i(x^*)\|^2 \end{aligned}$$

~~$$g_i(x^*) \leq f_i(x_t) - \eta \|\nabla f_i(x_t)\|^2 + \eta^2 L \|\nabla f_i(x_t)\|^2 - \langle \nabla f_i(x^*), x_t - \eta \nabla f_i(x_t) \rangle.$$~~

$$\frac{1}{n} \sum_i g_i(x^*) = f(x^*) \leq f(x_t) - \frac{\eta(1-\eta L)}{n} \sum_i \|\nabla f_i(x_t)\|^2$$

$$\Rightarrow \sum_i \|\nabla f_i(x_t)\|^2 \leq \frac{1}{\eta(1-\eta L)} [f(x_t) - f(x^*)].$$

$$\text{So, } \mathbb{E}[\|x_{t+1} - x^*\|^2] \leq \mathbb{E}[\|x_t - x^*\|^2] - \eta(2-8L\eta) [f(x_t) - f(x^*)] + \eta^2 L [f(x_0) - f(x^*)].$$

Telescoping, we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[f(x_t)] - f(x^*) &\leq \frac{\|x_0 - x^*\|^2}{\eta T} + \eta L [f(x_0) - f(x^*)] \\ &\leq \left[\frac{1}{\eta T} + \eta L \right] [f(x_0) - f(x^*)]. \end{aligned}$$

Choose $\eta = \frac{1}{10L}$ and $T > 20K$.

$$\Rightarrow \mathbb{E}[f(\bar{x}_T)] - f(x^*) \leq \frac{2}{3} \{ f(x_0) - f(x^*) \}.$$

Get run time complexity of $O\left[(nd + dK) \log \frac{1}{\epsilon}\right]$

$$GD : O\left(nd \log \frac{1}{\epsilon}\right)$$

$$SGD : O\left(\frac{d\sigma^2}{\epsilon}\right).$$

IMPORTANCE SAMPLING :

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

Smoothness parameter $L = \max_i \|a_i\|^2$.

Instead let us sample each i w.p. $\propto \|a_i\|^2$.

$$\text{then } f(x) = \sum_{i=1}^n p_i \left\{ \frac{\sum_j \|a_j\|^2}{2n \|a_i\|^2} (a_i^T x - b_i)^2 \right\}.$$

$$p_i = \frac{\|a_i\|^2}{\sum_j \|a_j\|^2}.$$

Here the Smoothness parameter = $\frac{\sum_j \|a_j\|^2}{n}$.

Remaining calculations in SVRG go through.

$$\text{Total time } O\left[(nd + d \cdot \left(\frac{\sum_j \|a_j\|^2}{n}\right)) \log \frac{1}{\epsilon}\right].$$

NONCONVEX OPTIMIZATION

NP-hard in general. Can only find ~~locally~~ OFF STATIONARY PTS.

$$\boxed{f(x)}$$

FOSP: $\nabla f = 0$

SOSP: $\nabla f = 0$ and $\nabla^2 f \succeq 0$.

Approximate versions: $\|\nabla f\| \leq \epsilon_g$
and $\nabla^2 f \succeq -\epsilon_H \text{Id}$.

Let us first focus on approximate FOSP.

Lemma: If $f(\cdot)$ is L -Smooth, then GD finds an ϵ -FOSP in $O\left(\frac{L\{f(x_0) - f^*\}}{\epsilon^2}\right)$ iterations.

Proof: Already done in the first lecture.

Finite sum setting: $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$

Can we get SVRG-style complexity?

Algorithm: Start with \tilde{x}_0 .

Compute $\tilde{x}_{j+1} =$ and $\nabla f(\tilde{x}_j)$

Inner loop $\left\{ \begin{array}{l} \text{For } k = 1, \dots, m \\ x_{k+1} = x_k - \eta \nabla f_i(x_k) + \eta \nabla f_i(\tilde{x}_j) \\ \quad \quad \quad - \eta \nabla f(\tilde{x}_j). \end{array} \right.$

We have $E[f(x_{k+1})] = E[f(x_k - \eta \nabla f_i(x_k) + \eta \nabla f_i(\tilde{x}_j) - \eta \nabla f(\tilde{x}_j))]$

$$\leq f(x_k) - \eta \|\nabla f(x_k)\|^2 + \eta^2 L E[\|\nabla f_i(x_k) - \nabla f_i(\tilde{x}_j) + \nabla f(\tilde{x}_j)\|^2]$$

What is $E[\|\nabla f_i(x_k) - \nabla f_i(\tilde{x}_j) + \nabla f(\tilde{x}_j) - \nabla f(x_k)\|^2]$?

$$\leq 4L^2 E[\|x_k - \tilde{x}_j\|^2].$$

$$x_k - x_{k+1} = \eta \nabla f_i(x_k) - \eta \nabla f_i(\tilde{x}_j) + \eta \nabla f(\tilde{x}_j).$$

$$\mathbb{E}[\|x_k - x_{k+1} - \eta \nabla f(x_k)\|^2] \leq 4\eta^2 L^2 \mathbb{E}[\|x_k - \tilde{x}\|^2].$$

and

$$\mathbb{E}[f(x_{k+1})] \leq f(x_k) - \frac{\eta}{2} \mathbb{E}[\|\nabla f(x_k)\|^2] + 4\eta^2 L^2 \mathbb{E}[\|x_k - \tilde{x}\|^2]$$

$$\mathbb{E}[\|\nabla f(x_k)\|^2] \geq \mathbb{E}[\|\nabla f(x_k)\|^2]$$

$$\geq \frac{1}{\eta^2} \mathbb{E}[\|x_k - x_{k+1} - \eta \nabla f(x_k)\|^2]$$

$$\text{Let } v_k = \nabla f_i(x_k) - \nabla f_i(\tilde{x}) + \nabla f(\tilde{x}).$$

$$x_{k+1} = x_k - \eta v_k.$$

$$\mathbb{E}[f(x_{k+1})] = \mathbb{E}[f(x_k - \eta v_k)]$$

$$\leq \mathbb{E}\left[f(x_k) - \eta \langle \nabla f(x_k), v_k \rangle + \frac{\eta^2 L}{2} \|v_k\|^2\right]$$

$$= f(x_k) - \eta \|\nabla f(x_k)\|^2 + \frac{\eta^2 L}{2} \|\nabla f(x_k)\|^2$$

$$+ \frac{\eta^2 L}{2} \cdot 4L^2 \|x_k - \tilde{x}\|^2.$$

$$\leq f(x_k) - \frac{\eta}{2} \|\nabla f(x_k)\|^2 + 2\eta^2 L^3 \|x_k - \tilde{x}\|^2.$$

$$\begin{aligned}
\mathbb{E}[\|x_{k+1} - \bar{x}\|^2] &= \mathbb{E}[\|x_k - \eta v_k - \bar{x}\|^2] \\
&= \mathbb{E}[\|x_k - \bar{x}\|^2] - 2\eta \langle \nabla f(x_k), x_k - \bar{x} \rangle \\
&\quad + \eta^2 \mathbb{E}[\|v_k\|^2] \\
&\leq \mathbb{E}[\|x_k - \bar{x}\|^2] + \eta \left[\frac{1}{\beta} \|\nabla f(x_k)\|^2 + \frac{\beta}{\beta} \|x_k - \bar{x}\|^2 \right] \\
&\quad + \eta^2 \|\nabla f(x_k)\|^2 + 2\eta^2 L^2 \|x_k - \bar{x}\|^2 \\
&= \left(1 + \frac{\eta\beta}{\beta} + 2\eta^2 L^2\right) \|x_k - \bar{x}\|^2 + \frac{\eta^2 + \eta}{\beta} \|\nabla f(x_k)\|^2 \\
\mathbb{E}[f(x_{k+1})] + c_{k+1} \mathbb{E}[\|x_{k+1} - \bar{x}\|^2] \\
&\leq f(x_k) - \left[\frac{\eta}{2} - \frac{\eta}{\beta} c_{k+1} - \eta^2 c_{k+1} \right] \|\nabla f(x_k)\|^2 \\
&\quad + \left\{ 2\eta^2 L^3 + c_{k+1} \left(1 + \frac{\eta}{\beta} + 2\eta^2 L^2\right) \right\} \|x_k - \bar{x}\|^2.
\end{aligned}$$

~~Pick $\alpha_{k+1} = \frac{1}{4c_{k+1}}$. Then,~~

~~$$\begin{aligned}
\mathbb{E}[f(x_{k+1})] + c_{k+1} \mathbb{E}[\|x_{k+1} - \bar{x}\|^2] \\
&\leq f(x_k) - \frac{\eta}{4} [1 - 4\eta c_{k+1}] \|\nabla f(x_k)\|^2 \\
&\quad + \left\{ 2\eta^2 L^3 + c_{k+1} (1 + 4\eta c_{k+1} + 2\eta^2 L^2) \right\} \|x_k - \bar{x}\|^2.
\end{aligned}$$~~

If $\eta < \frac{1}{2L}$, then

$$\begin{aligned} & \mathbb{E}[f(x_{k+1})] + C_{k+1} \mathbb{E}[\|x_{k+1} - \tilde{x}\|^2] \\ & \leq f(x_k) + C_k \|x_k - \tilde{x}\|^2 - \left[\frac{\eta}{2} - \frac{\eta}{\beta} C_{k+1} - \eta^2 C_{k+1} \right] \|\nabla f(x_k)\|^2. \end{aligned}$$

where $C_{k+1} = C_k (1 + \eta\beta + 2\eta^2 L^2) + 2\eta^2 L^3$.

If we let $C_{m+1} = 0$ then

$$C_0 = 2\eta^2 L^3 \cdot \frac{(1 + \eta\beta + 2\eta^2 L^2)^{m+1} - 1}{\eta\beta + 2\eta^2 L^2}$$

$$\left. \begin{aligned} \eta &= \frac{M_0}{L\eta^\alpha} \\ \beta &= \frac{L}{\eta^{\alpha/2}} \\ m &= \frac{\eta^{3\alpha/2}}{3M_0} \end{aligned} \right\} \leq \frac{2M_0^2 L}{\eta^{2\alpha}} \cdot \frac{10}{\frac{M_0}{\eta^{3\alpha/2}}} = \frac{20M_0 L}{\eta^{\alpha/2}}$$

On the other hand,

$$\begin{aligned} \min_k \left[\frac{\eta}{2} - \frac{\eta}{\beta} C_{k+1} - \eta^2 C_{k+1} \right] &= \frac{\eta}{2} - \frac{\eta}{\beta} C_0 - \eta^2 C_0 \\ &= \frac{\eta}{2} \left[1 - 20M_0 - \frac{20M_0^2}{\eta^{3\alpha/2}} \right] \\ &\geq \frac{\eta}{4}. \end{aligned}$$

Taking a telescopic sum, we have

$$\frac{1}{m} \sum_{k=0}^m \|\nabla f(x_k)\|^2 \leq \frac{4[f(x_0) - f(x_m)]}{m \eta}$$

~~$$= \frac{4}{m \eta} [f(x_0) - f(x_m)]$$~~

Taking further sum over outer loops, we have

$$\frac{1}{mT} \sum_{k=0}^m \sum_{t=1}^T \|\nabla f(x_k^t)\|^2 \leq \frac{4[f(x_0) - f^*]}{(mT) \cdot \eta}$$

With n calls, we get for $\alpha = 2/3$ and $T=1$,

$$\text{that } \epsilon = \frac{1}{n^{1/3}}$$

If we continue further, for larger T , we get

$$nT \text{ calls } \rightarrow \epsilon = \frac{1}{Tn^{1/3}}$$

$$\hookrightarrow O\left(n + \frac{n^{2/3}}{\epsilon}\right)$$

Focus on finding SOSPs

Importance: Many problems such as matrix completion, robust PCA, ... all SOSPs are good. Give example of PCA.

Want: $\|\nabla f(x)\| \leq \epsilon$
and $\nabla^2 f(x) \succeq -\epsilon_H \text{Id}$.

Smoothness = Gradient Lipschitz. \Leftarrow

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

Hessian Lipschitz: $\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq \rho \|x - y\|.$

Hessian Lipschitz is necessary for SOSPs.

Gradient descent can get stuck at FOSPs that are not SOSPs.

However this can happen only for measure zero set of bad initialization.

However there are still situations where escaping saddle points can take exponential time.

Perturbation to the rescue.

If $\|\nabla f(x_t)\| \leq \epsilon$ and perturbation not added for
Some time, then

$$x_t \leftarrow x_t - \eta \xi_t$$

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

Descent lemma: $f(x_{t+1}) \leq f(x_t) - \frac{\eta}{2} \|\nabla f(x_t)\|^2$.

Improve or localize: $\|x_t - x_0\| \leq \sqrt{2\eta t [f(x_0) - f(x_t)]}$

Lemma: If \tilde{x} is a saddle pt. i.e., $\|\nabla f(\tilde{x})\| \leq \epsilon$
and $d_{\min}(\nabla^2 f(\tilde{x})) \leq -\sqrt{\epsilon}$.

Then letting $x_0 = \tilde{x} + \eta \xi$ ($\xi \sim \text{Unit}(B_0(r))$) and
running GD for τ iterations then w.p. $\geq 1 - \delta$,

$$f(x_\tau) - f(\tilde{x}) \leq -\frac{\mathcal{F}}{2} \text{ where}$$

$$\eta = \frac{1}{L}, \quad \gamma = \frac{\epsilon}{400}, \quad \alpha = C \log\left(\frac{dL\Delta_f}{\beta\epsilon\delta}\right)$$

$$\tau = \frac{\beta}{\sqrt{\beta\epsilon}}, \quad \mathcal{F} = \frac{1}{50\alpha^3} \sqrt{\frac{\epsilon}{\beta}}, \quad \mathcal{S} = \frac{1}{4\alpha} \sqrt{\frac{\epsilon}{\beta}}$$

$$f(x_t) - f(\tilde{x}) = f(x_t) - f(x_0) + f(x_0) - f(\tilde{x})$$

w.h.p. $f(x_t) - f(x_0)$ is small.
 $\leq -\mathcal{F}$.

~~max~~
 (1) Set of bad points from where GD does not escape quickly.

(2) Volume of that set is small
 \iff Any two points which are separated in the -ve direction do not belong to the set.

(3) x_0 \cdot
 $\downarrow = \eta \gamma_0 \epsilon_1$
 $x'_0 \cdot$ then $\min \left\{ f(x_t) - f(x_0), f(x'_t) - f(x'_0) \right\} \leq -\mathcal{F}$.

Proof: If not, by improve or localize,

$$\|x_t - \tilde{x}_0\| \leq \sqrt{2\eta^2 \mathcal{F} + \eta \gamma} \leq \mathcal{F}$$

$$\text{and } \|x'_t - \tilde{x}_0\| \leq \mathcal{F}.$$

Denote $\hat{x}_t \triangleq x_t - x'_t$.

$$\hat{x}_{t+1} = x_{t+1} - x'_{t+1}$$

$$= \hat{x}_t - \eta [\nabla f(x_t) - \nabla f(x'_t)]$$

$$= \underbrace{(\mathbb{I} - \eta H)}_1 \hat{x}_t - \eta \Delta_t \hat{x}_t$$

~~$$\nabla f(x_t) = \nabla f(x'_t) + \int_{\theta=0}^1 \nabla^2 f(\theta x_t + (1-\theta)x'_t) (x_t - x'_t) d\theta$$~~

$$\nabla f(x_t) = \nabla f(x'_t) + \int_{\theta=0}^1 \nabla^2 f(x'_t + \theta(x_t - x'_t)) (x_t - x'_t) d\theta$$

$$\hat{x}_{t+1} = (\mathbb{I} - \eta H) \hat{x}_t - \eta \Delta_t \hat{x}_t, \text{ where}$$

$$\Delta_t \triangleq \int_{\theta=0}^1 (\nabla^2 f(x'_t + \theta(x_t - x'_t)) - H) d\theta$$

$$= \underbrace{(\mathbb{I} - \eta H) \hat{x}_0}_{\triangleq p(t+1)} - \eta \underbrace{\sum_{\tau=0}^t (\mathbb{I} - \eta H) \Delta_\tau \hat{x}_\tau}_{\triangleq q(t+1)}$$

Claim: $\|q(t)\| \leq \frac{\|p(t)\|}{2} \quad \forall t \in \mathcal{T}$.

Proof: By induction.

Base case: $q(0) = 0$ so, true.

By Herian Lipschitz, $\|\Delta z\| \leq \rho \max \{ \|x_t - \tilde{x}\|, \|x'_t - \tilde{x}'\| \}$

$$\leq \rho \delta$$

and $\|\hat{x}_\tau\| \leq \|\rho(\tau)\| + \|\zeta(\tau)\|$

(By induction hyp.) $\leq 2\|\rho(\tau)\| = 2(1+\eta)^{\tau} \eta r_0$

$$\begin{aligned} \|q(t+1)\| &\leq \eta \sum \|(\mathbb{I} - \eta H)\|^{t-\tau} \|\Delta z\| \cdot \|\hat{x}_\tau\| \\ &\leq 2\eta \rho \delta \sum_{\tau=0}^t (1+\eta)^{\tau} \eta r_0 \\ &\leq 2\eta \rho \delta \tau \|\rho(t+1)\|. \end{aligned}$$

By choice, $2\eta \rho \delta \tau \leq \frac{1}{2}$. \square

This implies $\max \{ \|x_\tau - \tilde{x}\|, \|x'_\tau - \tilde{x}'\| \}$

$$\begin{aligned} &\geq \frac{1}{2} \|\hat{x}_\tau\| \geq \frac{1}{4} \|\rho(\tau)\| \\ &= \frac{(1+\eta)^{\tau} \eta r_0}{4} > \delta. \end{aligned}$$

Contradiction. \square

CONVEX - CONCAVE MINIMAX OPT.

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

$f(\cdot, y)$ is Convex $\forall y$

$f(x, \cdot)$ is Concave $\forall x$.

Example applications:

① Constrained optimization.

Find x

$$\text{s.t. } f_i(x) \leq 0 \\ i = 1, \dots, m$$

$$\longrightarrow \min_x \max_{i=1, \dots, m} f_i(x)$$

$$\parallel \quad \min_x \max_{\substack{d_i \geq 0 \\ \sum_i d_i = 1}} \sum_{i=1}^m d_i f_i(x)$$

② Non Smooth functions can be sometimes be written as smooth minimax problems.

$$\|x\|_1 \equiv \max_{-1 \leq t_i \leq 1} \sum_i d_i x_i$$

③ Zero Sum games.

Exercise 1: $g(x) \triangleq \max_y f(x, y)$ is a convex function since $f(\cdot, y)$ is convex $\forall y$.

Danskin's theorem: If $f(x, y)$ is differentiable with respect to x $\forall y \in \mathcal{Y}$ and if $\frac{\partial f}{\partial x}$ is continuous w.r.t. y $\forall x$, then

$$\nabla g(x) = \text{Conv. hull} \left\{ \frac{\partial f(x, y)}{\partial x} : y \in \underset{z \in \mathcal{Y}}{\text{argmax}} f(x, z) \right\}$$

Subgradient

One way to solve the minimax problem is to do Subgradient descent on $g(x)$.

① Find $y_t \in \underset{y}{\text{argmax}} f(x_t, y)$

② $\nabla g(x_t) = \nabla_x f(x_t, y_t)$

③ $x_{t+1} = x_t - \eta \nabla g(x_t)$.

ISSUES: ① very time consuming to find y_t

② Requires f to be smooth.

Non-smooth setting: Gradient descent ascent
 $\hookrightarrow \frac{1}{\sqrt{t}}$ Conv. rate.

Smooth setting: Mirror-Prox
 $\hookrightarrow \frac{1}{t}$ Conv. rate.

Weak duality: $\min_x f(\hat{x}, y) \leq f(\hat{x}, \hat{y})$
 $\leq \max_x f(x, \hat{y})$

Weak duality: $\min_x f(x, \hat{y}) \leq f(\hat{x}, \hat{y})$
 $\leq \max_y f(\hat{x}, y)$

$\hat{y} \quad y^* \in \operatorname{argmax}_{\hat{y} \in Y} \left[\min_x f(x, \hat{y}) \right]$

and $x^* \in \operatorname{argmin}_{\hat{x} \in X} \left[\max_y f(\hat{x}, y) \right]$

then $\max_{\hat{y}} \min_x f(x, \hat{y}) = \min_x f(x, y^*)$
 $\leq f(x^*, y^*)$
 $\leq \max_y f(x^*, y)$
 $= \min_{\hat{x}} \max_y f(\hat{x}, y)$ \square

So, $\max_y \min_x f(x, y) \leq \min_x \max_y f(x, y)$.

Strong duality: If X and Y are compact and $f(\cdot, \cdot)$ is convex-concave, then

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

A point (\bar{x}, \bar{y}) is said to be an ϵ -primal dual pair if

$$\max_{y \in Y} f(\bar{x}, y) - \min_{x \in X} f(x, \bar{y}) < \epsilon.$$

Exercise: If (\bar{x}, \bar{y}) is an ϵ -primal dual pair of $f(x, y)$, then

\bar{x} is an ϵ -approx ~~soln~~ of minimizer of $g(x) \triangleq \max_y f(x, y)$

and \bar{y} is an ϵ - " maximizes of $h(y) \triangleq \min_x f(x, y)$.

Gradient descent ascent:

$$x_{t+1} = x_t - \eta \nabla_x f(x_t, y_t)$$

$$y_{t+1} = y_t + \eta \nabla_y f(x_t, y_t)$$

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

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

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

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

Fix x .

$$\begin{aligned} \|x_{t+1} - x\|^2 &= \|x_t - x\|^2 - 2\eta \langle \nabla_x f(x_t, y_t), x_t - x \rangle \\ &\quad + \eta^2 \|\nabla_x f(x_t, y_t)\|^2 \end{aligned}$$

$$\begin{aligned} &\leq \|x_t - x\|^2 - 2\eta [f(x_t, y_t) - f(x, y_t)] \\ &\quad + \eta^2 G^2 \end{aligned}$$

$$\begin{aligned} \|y_{t+1} - y\|^2 &= \|y_t - y\|^2 + 2\eta \langle \nabla_y f(x_t, y_t), y_t - y \rangle \\ &\quad + \eta^2 \|\nabla_y f(x_t, y_t)\|^2 \end{aligned}$$

$$\begin{aligned} &\leq \|y_t - y\|^2 - 2\eta [f(x_t, y) - f(x_t, y_t)] \\ &\quad + \eta^2 G^2 \end{aligned}$$

Adding,

$$\begin{aligned} \|x_{t+1} - x\|^2 + \|y_{t+1} - y\|^2 &\leq \|x_t - x\|^2 + \|y_t - y\|^2 \\ &\quad - 2\eta [f(x_t, y) - f(x, y_t)] \\ &\quad + \eta^2 G^2 \end{aligned}$$

Telescoping and reorganizing,

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

$$+ \eta G^2$$

By Convexity

$$f(\bar{x}_T, y) \leq \frac{1}{T+1} \sum_{t=0}^T f(x_t, y)$$

$$\text{where } \bar{x}_T \triangleq \frac{1}{T+1} \sum_{t=0}^T x_t$$

$$\text{Similarly } f(x, \bar{y}_T) \geq \frac{1}{T+1} \sum_{t=0}^T f(x, y_t)$$

$$\text{So, } f(\bar{x}_T, y) - f(x, \bar{y}_T)$$

$$\leq \frac{\|x_0 - x\|^2 + \|y_0 - y\|^2}{2\eta(T+1)} + \eta G^2$$

$$\max_y f(\bar{x}_T, y) - \min_x f(x, \bar{y}_T) \leq \frac{R^2}{\eta(T+1)} + \eta G^2$$

$$\eta = \frac{R}{G\sqrt{T+1}} \rightarrow \leq \frac{2RG}{\sqrt{T+1}}$$

Gradient ^{descent} algorithm: $x_{t+1} = x_t - \eta \nabla f(x_t)$

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

Exercise: Proximal algorithm has convergence rate $\frac{1}{T}$ for [even non-smooth] convex functions.

Exercise: Given a point x , and an L -smooth function, $f(\cdot)$, and $\eta \leq \frac{1}{L}$, we can ~~implement this~~ find z satisfying $z = x - \eta \nabla f(z)$ up to an accuracy of ϵ in $O(\log \frac{1}{\epsilon})$ steps.

(Conceptual)

Mirror-Prox:

$$x_{t+1} = x_t - \eta \nabla_x f(x_{t+1}, y_{t+1})$$

$$y_{t+1} = y_t + \eta \nabla_y f(x_{t+1}, y_{t+1})$$

Let $\bar{x}^+ = x - \eta \nabla_x f(\bar{x}, \bar{y})$

$$\bar{y}^+ = y - \eta \nabla_y f(\bar{x}, \bar{y})$$

Similarly, \tilde{x}^+ and \tilde{y}^+ .

$$\|\bar{x}^+ - \tilde{x}^+\| + \|\bar{y}^+ - \tilde{y}^+\| \leq 2\eta L [\|\bar{x} - \tilde{x}\| + \|\bar{y} - \tilde{y}\|].$$

These iterations can quickly find the CMP step.