

## Module B: Algorithms for Optimization

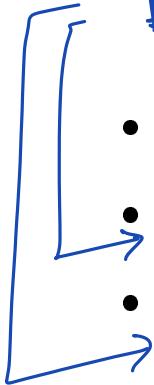
---

Recall that an optimization problem in standard form is given by

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{s.t.} \quad & g_i(x) \leq 0, i \in [m] := \{1, 2, \dots, m\}, \\ & h_j(x) = 0, j \in [p]. \end{aligned}$$

Most algorithms generate a sequence  $x_0, x_1, x_2, \dots$  by exploiting local information collected on the path.

- Zeroth Order: Only  $f(x_t), g_i(x_t), h_j(x_t)$  available.
- **First Order:** Gradients  $\nabla f(x_t), \nabla g_i(x_t), \nabla h_j(x_t)$  are used. Heavily used in ML.
- Second Order: Hessian information is used. Eg: Newton's Method, etc.
- Distributed Algorithms
- Stochastic/Randomized Algorithms



## Measure of progress

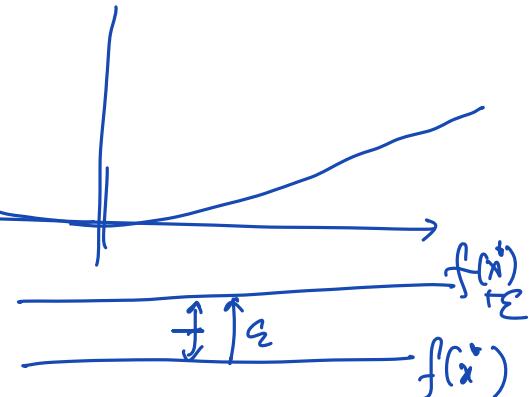
$$f(x) = 10^{-4} x^2$$

$$x_t = 10, x^* = 0$$

Let  $x^*$  be the optimal solution. The iterative algorithms continue till any of the following error metrics is sufficiently small.

- $\text{err}_t := \|x_t - x^*\|_2^2 = 100$
- $\text{err}_t := f(x_t) - f(x^*) = 0.01$
- A solution  $\bar{x}$  is  $\epsilon$ -optimal when

$$f(x^*) \leq f(\bar{x}) \leq f(x^*) + \epsilon.$$



We often run the algorithm till  $\text{err}_t$  is smaller than a sufficiently small  $\epsilon > 0$ .

- $\text{err}_t := \max(f(x_t) - f(x^*), g_1(x_t), g_2(x_t), \dots, g_m(x_t))$ ,  $(|h_1(x_t)|, |h_2(x_t)|, \dots)$



## First order methods: Gradient descent

Consider the unconstrained optimization problem:  $\min_{x \in \mathbb{R}^n} f(x)$

Gradient Descent (GD):  $x_{t+1} = x_t - \eta_t \nabla f(x_t)$ ,  $t \geq 0$  starting from an initial guess  $x_0 \in \mathbb{R}^n$ .

Convergence rate depends on choice of step size  $\eta_t$  and characteristic of the function.

• Bounded Gradient:  $\|\nabla f(x)\| \leq G$  for all  $x \in \mathbb{R}^n$ .

• Smooth: A differentiable convex  $f$  is  $\beta$ -smooth if for any  $x, y$ , we have

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

We can obtain a quadratic upper bound on the function from local information.

• Strongly Convex: A differentiable convex  $f$  is  $\alpha$ -strongly convex if for any  $x, y$ , we have

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

We can obtain a quadratic lower bound on the function from local information.

If  $f$  is twice differentiable, then

$$f \text{ is } \beta\text{-smooth} \Leftrightarrow \nabla^2 f(x) \preceq \beta I \quad \forall x \in \mathbb{R}^n$$

$$\Leftrightarrow \lambda_{\max}[\nabla^2 f(x)] \leq \beta \quad \forall x \in \mathbb{R}^n$$

$$f \text{ is } \alpha\text{-strongly convex} \Leftrightarrow \nabla^2 f(x) \succeq \alpha I \quad \forall x \in \mathbb{R}^n$$

$$\Leftrightarrow \lambda_{\min}[\nabla^2 f(x)] \geq \alpha \quad \forall x \in \mathbb{R}^n$$

$$\text{Ex: } f(x) = \frac{\|Ax - b\|_2^2}{2}$$

$$\nabla^2 f(x) = A^T A,$$

$$\therefore \beta = \lambda_{\max}(A^T A), \alpha = \lambda_{\min}(A^T A)$$

$$\|\nabla f(x)\| \leq G \quad \forall x \in \mathbb{R}^n.$$

## Gradient Descent with Bounded Gradient Assumption

Let  $x_0, x_1, \dots, x_{T-1}$  be the iterates generated by the GD algorithm.

For any  $t$ , we define  $\hat{x}_t = \frac{1}{t} \sum_{i=0}^{t-1} x_i$ . Let  $x^*$  be the optimal solution.

### Theorem 1: Convergence of Gradient Descent

Let the function  $f$  satisfy the bounded gradient property. Let  $\|x_0 - x^*\| \leq D$ . Then, for the choice of step size  $\eta_t = \frac{D}{G\sqrt{T}}$ , we have

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

To find an  $\epsilon$ -optimal solution, choose  $T \geq \left(\frac{DG}{\epsilon}\right)^2$  and  $\eta = \frac{\epsilon}{G^2}$ .

Possible Limitation: Need to know  $D$  and  $G$ .

Proof: Define the following (potential) function:

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

We show that  $\Phi_t$  is decreasing in  $t$ . We compute  $\Phi_{t+1} - \Phi_t$  as:

$$\begin{aligned} \Phi_{t+1} - \Phi_t &= \frac{1}{2\eta} \left[ \|x_{t+1} - x^*\|^2 - \|x_t - x^*\|^2 \right] \\ &= \frac{1}{2\eta} \left[ \|x_{t+1} - x_t + x_t - x^*\|^2 - \|x_t - x^*\|^2 \right] \\ &= \frac{1}{2\eta} \left[ \|x_{t+1} - x_t\|^2 + 2\langle x_t - x^*, x_{t+1} - x_t \rangle + \|x_t - x^*\|^2 - \|x_t - x^*\|^2 \right] \\ &= \frac{1}{2\eta} \left[ \|\eta \nabla f(x_t)\|^2 - 2\eta \langle x_t - x^*, \nabla f(x_t) \rangle \right] \\ &= \frac{\eta}{2} \|\nabla f(x_t)\|^2 - \langle x_t - x^*, \nabla f(x_t) \rangle \end{aligned}$$

First pick  $\epsilon$ .  
Then set  $T$ ,  
then find  $\eta$ ,  
then run GD.

take avg. of solns.

$\hat{x}_T$  has  $\epsilon$ -optimality.

$$\frac{DG}{\sqrt{T}} \leq \epsilon$$

$$\Rightarrow T \geq \left(\frac{DG}{\epsilon}\right)^2$$

$$\eta = \frac{D}{G\sqrt{T}}$$

$$10^{-3}, \quad \frac{1}{\epsilon^2} \approx 10^6$$

$$10^{-3}, \quad \frac{1}{\epsilon^2} \approx 10^6$$

$$10^{-3}, \quad \frac{1}{\epsilon^2} \approx 10^6$$

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

Proof

$$\nabla f(x_t), x_t - x^* \geq [f(x_t) - f(x^*)]$$

$$\frac{\eta}{2} G^2 - [f(x_t) - f(x^*)]$$

$$\Rightarrow \phi_{t+1} - \phi_t + f(x_t) - f(x^*) \leq \frac{\eta}{2} G^2, \quad t=0, 1, 2, \dots, T$$

$$\sum_{t=0}^{T-1} \left[ \phi_{t+1} - \phi_t + f(x_t) - f(x^*) \right] \leq \frac{\eta T}{2} G^2$$

$$\Rightarrow \underbrace{\phi_T - \phi_0}_{\sum_{t=0}^{T-1} f(x_t) - T f(x^*)} + \frac{\eta T}{2} G^2 \leq \frac{\eta T}{2} G^2$$

$$\Rightarrow \left[ \frac{1}{T} \sum_{t=0}^{T-1} f(x_t) - f(x^*) \right] + \frac{\eta}{T} \phi_T \leq \frac{\eta}{2} G^2 + \frac{\phi_0}{T}$$

$$\Rightarrow \frac{1}{T} \sum_{t=0}^{T-1} f(x_t) - f(x^*) \leq \frac{\eta}{2} G^2 + \frac{1}{2\eta T} D^2$$

find  $\eta$  which minimizes RHS.

$$\frac{1}{2} G^2 = \frac{D^2}{2T} \cdot \frac{1}{\eta^2} \Rightarrow \eta^* = \frac{D}{G\sqrt{T}}$$

$$\Rightarrow \frac{1}{T} \sum_{t=0}^{T-1} f(x_t) - f(x^*) \leq \frac{G^2}{2} \cdot \frac{D}{G\sqrt{T}} + \frac{D^2}{2T} \cdot \frac{G\sqrt{T}}{D}$$

$$\underline{f(\hat{x}_T) - f(x^*)} = 5 \frac{GD}{2\sqrt{T}} + \frac{GD}{2\sqrt{T}} = \frac{GD}{\sqrt{T}}$$

From convexity of  $f$ ,  $f\left(\frac{1}{T} \sum_{t=0}^{T-1} x_t\right) \leq \frac{1}{T} \sum_{t=0}^{T-1} f(x_t)$

Proof Continues

---

Finally  $f(\hat{x}_T) - f(x^*) \leq \frac{GD}{T}$



## Gradient Descent with Smoothness Assumption

Recall that a differentiable convex  $f$  is  $\beta$ -smooth if for any  $x, y$ , we have

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

### Theorem 2

Let the function  $f$  be  $\beta$ -smooth. ~~Let  $\|x_0 - x^*\| \leq \frac{1}{\beta}$~~ . Then, for the choice of step size  $\eta_t = \frac{1}{\beta}$ , we have

$$f(x_T) - f(x^*) \leq \frac{\beta \|x_0 - x^*\|^2}{2T}.$$

Proof: Define the following (potential) function:

$$\Phi_t := t[f(x_t) - f(x^*)] + \frac{\beta}{2} \|x_t - x^*\|^2.$$

We show that  $\Phi_t$  is decreasing in  $t$ . We compute  $\Phi_{t+1} - \Phi_t$  as:

If we require  $\epsilon$ -optimal sol<sup>n</sup>, then  

$$\frac{\beta \|x_0 - x^*\|^2}{2T} \leq \epsilon$$

$$\phi_0 = \frac{\beta}{2} \|x_0 - x^*\|^2$$

$$\Rightarrow T \geq \frac{\beta \|x_0 - x^*\|^2}{2\epsilon} \approx O\left(\frac{1}{\epsilon}\right)$$

$$\phi_T = T[f(x_T) - f(x^*)] + \text{const}$$

$$\text{If } \phi_{t+1} \leq \phi_t + t, \text{ then } \phi_T \leq \phi_0$$

$$\Rightarrow [f(x_T) - f(x^*)] \leq \frac{\beta \|x_0 - x^*\|^2}{2T}$$

$$\phi_{t+1} - \phi_t = (t+1)[f(x_{t+1}) - f(x^*)] + \frac{\beta}{2} \|x_{t+1} - x^*\|^2$$

$$= t[f(x_t) - f(x^*)] + \frac{\beta}{2} \|x_t - x^*\|^2$$

$$= (t+1)[f(x_{t+1}) - f(x_t)] + f(x_t) - f(x^*)$$

$$+ \frac{\beta}{2} \left[ \|x_{t+1} - x^*\|^2 - \|x_t - x^*\|^2 \right]$$

Proof

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

$$\leq (t+1) \underbrace{[f(x_{t+1}) - f(x_t)]}_{\text{f(x_{t+1}) - f(x_t)}} + \frac{1}{2\beta} \|\nabla f(x_t)\|^2$$

$$\leq (t+1) \left[ \langle \nabla f(x_t), \underbrace{x_{t+1} - x_t}_{\text{f(x_{t+1}) - f(x_t)}} \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2 \right] + \frac{1}{2\beta} \|\nabla f(x_t)\|^2$$

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

$$= (t+1) \left[ \langle \nabla f(x_t), -\frac{1}{\beta} \nabla f(x_t) \rangle + \frac{\beta}{2} \left\| -\frac{1}{\beta} \nabla f(x_t) \right\|^2 \right] + \frac{1}{2\beta} \|\nabla f(x_t)\|^2$$

$$= \|\nabla f(x_t)\|^2 \left[ -\frac{t+1}{\beta} + \frac{(t+1)}{2\beta} + \frac{1}{2\beta} \right]$$

$$= \|\nabla f(x_t)\|^2 \frac{1}{2\beta} (1 + t + 1 - 2t - 2)$$

$$= -\frac{t}{2\beta} \|\nabla f(x_t)\|^2 \leq 0 .$$

$\Rightarrow \underline{\phi_{t+1} \leq \phi_t \Rightarrow \phi_t \leq \phi_0}$  which proves the theorem.

# Gradient Descent with Smoothness and Strong Convexity

Recall that a differentiable convex  $f$  is  $\alpha$ -strongly convex if for any  $x, y$ , we have

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

### Theorem 3

Let the function  $f$  be  $\beta$ -smooth and  $\alpha$ -strongly convex with  $\alpha \leq \beta$ . Define condition number  $\kappa := \frac{\beta}{\alpha}$ . Then, for the choice of step size  $\eta_t = \frac{1}{\beta}$ , we have

$$f(x_T) - f(x^*) \leq e^{-\frac{T}{\kappa}} (f(x_0) - f(x^*)).$$

Note: To obtain  $\epsilon$ -optimal solution, choose  $T = \mathcal{O}(\log(\frac{1}{\epsilon}))$ .

Proof: Define the following (potential) function:

$$\Phi_t := (1 + \gamma)^t [f(x_t) - f(x^*)], \quad \text{where } \gamma = \frac{1}{\kappa - 1} = \frac{\alpha}{\beta - \alpha}.$$

We need to show that  $\Phi_{t+1} \leq \Phi_t$ .

$$\Rightarrow \Phi_{t+1} \leq \Phi_t$$

$$\Rightarrow (1 + \gamma)^{t+1} [f(x_{t+1}) - f(x^*)] \leq (1 + \gamma)^t [f(x_t) - f(x^*)]$$

$$\Rightarrow f(x_{t+1}) - f(x^*) \leq (1 + \gamma)^{t+1} [f(x_t) - f(x^*)]$$

$$\leq e^{-\frac{t+1}{\kappa}} [f(x_0) - f(x^*)]$$

$$\begin{aligned} \Phi_{t+1} - \Phi_t \\ = (1 + \gamma)^{t+1} [f(x_{t+1}) - f(x^*)] \\ - (1 + \gamma)^t [f(x_t) - f(x^*)] \end{aligned}$$

$$\Rightarrow (1 + \gamma)^t [\Phi_{t+1} - \Phi_t]$$

$$\begin{aligned} = (1 + \gamma) [f(x_{t+1}) - f(x^*)] - [f(x_t) - f(x^*)] \\ + (1 + \gamma) f(x_t) - \gamma f(x_t) \end{aligned}$$

$$= (1 + \gamma) [f(x_{t+1}) - f(x_t)] - \underbrace{(1 + \gamma) f(x^*) + f(x^*) + \gamma f(x_t)}_{9}$$

$$= (1 + \gamma) [f(x_{t+1}) - f(x_t)] + \gamma [f(x_t) - f(x^*)]$$

$$\leq (1+\gamma) \left( \frac{-1}{2\beta} \|\nabla f(x_t)\|_2^2 \right) + \frac{\gamma}{2\alpha} \|\nabla f(x_t)\|_2^2$$

1<sup>st</sup> term: use smoothness

Proof

Since  $f$  is  $\beta$ -smooth,  $f(x_{t+1}) \leq f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|_2^2$

$$\begin{aligned} \Rightarrow f(x_{t+1}) - f(x_t) &\leq \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|_2^2 \\ &= \langle \nabla f(x_t), -\frac{1}{\beta} \nabla f(x_t) \rangle + \frac{\beta}{2} \left\| -\frac{1}{\beta} \nabla f(x_t) \right\|_2^2 \\ &= -\frac{1}{2\beta} \|\nabla f(x_t)\|_2^2. \end{aligned}$$

2<sup>nd</sup> term: use strong convexity.

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

$$f(x^*) \geq f(x_t) + \langle \nabla f(x_t), x^* - x_t \rangle + \frac{\alpha}{2} \|x_t - x^*\|_2^2$$

$$\begin{aligned} \Rightarrow f(x_t) - f(x^*) &\leq \langle \nabla f(x_t), x_t - x^* \rangle - \frac{\alpha}{2} \|x_t - x^*\|_2^2 \\ &\leq \frac{1}{2\alpha} \|\nabla f(x_t)\|_2^2 \end{aligned}$$

$$\text{Note: } \|a-b\|_2^2 = a^T a + b^T b - 2a^T b \geq 0$$

$$\Rightarrow a^T a \geq 2a^T b - b^T b$$

let  $b = \sqrt{\frac{\alpha}{2}} (x_t - x^*)$ ,  $a = \frac{1}{\sqrt{2\alpha}} \nabla f(x_t)$

$f(x_t) - f(x^*) \leq \frac{1}{2\alpha} \|\nabla f(x_t)\|_2^2$  is called PL-inequality.

Strongly convex functions satisfy this.

Some non-convex functions also satisfy this.

$$\text{Recall: } \gamma = \frac{\alpha}{\beta - \alpha} \Rightarrow 1 + \gamma = \frac{\beta}{\beta - \alpha} \Rightarrow \frac{1 + \gamma}{2\beta} = \frac{1}{2(\beta - \alpha)}$$

### Proof Continues

---

We have

$$(1 + \gamma)^{-t} [\phi_{t+1} - \phi_t] \leq -\frac{1 + \gamma}{2\beta} \|\nabla f(x_t)\|_2^2 + \frac{\gamma}{2\alpha} \|\nabla f(x_t)\|_2^2$$

$$= \|\nabla f(x_t)\|_2^2 \left[ \frac{1}{2(\beta - \alpha)} - \frac{1 + \gamma}{2\beta} \right]$$

$$\Rightarrow \phi_{t+1} \leq \phi_t + t. \quad = \textcircled{1}.$$

## Summary of gradient descent convergence rates

Consider the unconstrained optimization problem:  $\min_{x \in \mathbb{R}^n} f(x)$

Gradient Descent (GD):  $x_{t+1} = x_t - \eta_t \nabla f(x_t)$ ,  $t \geq 0$  starting from an initial guess  $x_0 \in \mathbb{R}^n$ .

### Theorem 4: GD Convergence rates

Let  $\|x_0 - x^*\| \leq D$ .

- If  $\|\nabla f(x)\| \leq G$  for all  $x \in \mathbb{R}^n$ , then with  $\eta_t = \frac{D}{G\sqrt{T}}$ ,  $f(\hat{x}_T) - f(x^*) \leq \frac{DG}{\sqrt{T}}$ .
- If  $f$  is  $\beta$ -smooth, for  $\eta_t = \frac{1}{\beta}$ ,  $f(x_T) - f(x^*) \leq \frac{\beta\|x_0 - x^*\|^2}{2T}$ .
- If  $f$  is  $\beta$ -smooth and  $\alpha$ -strongly convex, for  $\eta_t = \frac{1}{\beta}$ ,  $f(x_T) - f(x^*) \leq e^{-\frac{T}{\kappa}}(f(x_0) - f(x^*))$  where  $\kappa := \frac{\beta}{\alpha}$  is the condition number.

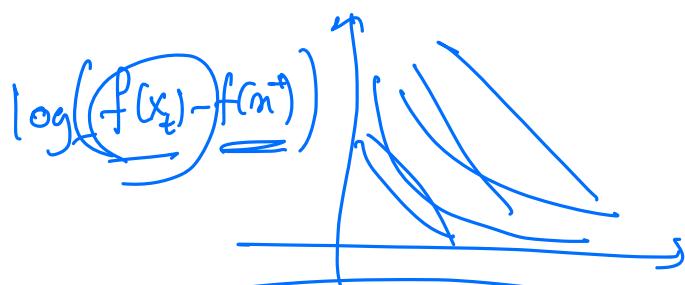
If  $K$  is large,  $e^{-T/K}$  decreases slowly  $\Rightarrow$  ill-conditioned

If  $f(x) = x^T x$ ,  $K = 1$  since  $\nabla^2 f(x) = I$  problems

$$f(x) = 10^6 x_1^2 + \sum_{j=2}^n x_j^2, \quad K = 10^6$$

for a desired  $\epsilon$ ,

$$e^{-T/K} (f(x_0) - f(x^*)) \leq \epsilon$$

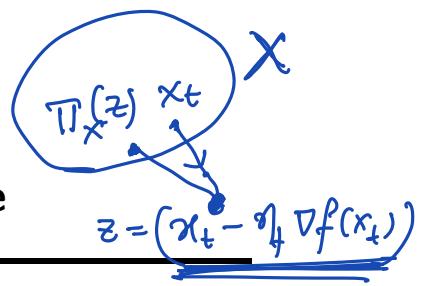


$$\Rightarrow -\frac{T}{K} + \log(\text{const}) \leq \log(\epsilon)$$

$$\Rightarrow \frac{T}{K} \geq -\log(\epsilon) + \text{const.}$$

$$\Rightarrow T \geq K \left[ \log\left(\frac{1}{\epsilon}\right) + \text{const.} \right]$$

## Gradient descent: Constrained Case



Consider the unconstrained optimization problem:  $\min_{x \in X} f(x)$  where  $X \subseteq \mathbb{R}^n$  is a convex feasibility set.

Projected Gradient Descent (PGD):  $x_{t+1} = \Pi_X[x_t - \eta_t \nabla f(x_t)]$ ,  $t \geq 0$  starting from an initial guess  $x_0 \in \mathbb{R}^n$  where  $\Pi_X(y)$  is the projection of  $y$  on the set  $X$ .

### Theorem 5

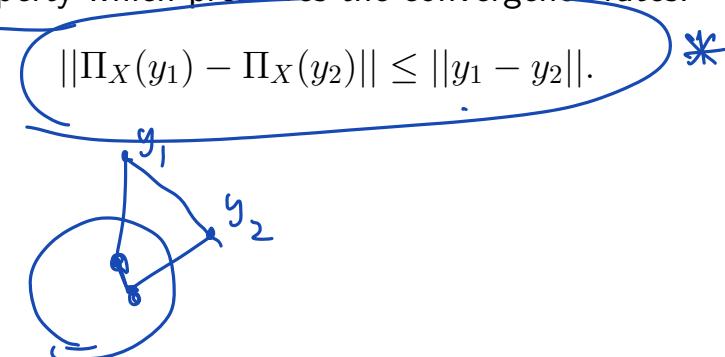
Let  $\|x_0 - x^*\| \leq D$ .

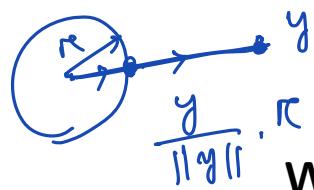
- If  $\|\nabla f(x)\| \leq G$  for all  $x \in \mathbb{R}^n$ , then with  $\eta_t = \frac{D}{G\sqrt{T}}$ ,  $f(\hat{x}_T) - f(x^*) \leq \frac{DG}{\sqrt{T}}$ .
- If  $f$  is  $\beta$ -smooth, for  $\eta_t = \frac{1}{\beta}$ ,  $f(x_T) - f(x^*) \leq \frac{\beta\|x_0 - x^*\|^2}{2T}$ .
- If  $f$  is  $\beta$ -smooth and  $\alpha$ -strongly convex, for  $\eta_t = \frac{1}{\beta}$ ,  $f(x_T) - f(x^*) \leq e^{-\frac{T}{\kappa}}(f(0) - f(x^*))$  where  $\kappa := \frac{\beta}{\alpha}$  is the condition number.

Note: Convergence rates remain unchanged.

Note: Projection itself is another optimization problem!

Non-expansive Property which preserves the convergence rates:





$$\underset{x \in X}{\operatorname{argmin}} \|x - y\|_2^2$$

to cross check, one should verify  $\Pi_X(y) \in X$

When is Projection easy to find?

Note that  $\Pi_X(y) = \underset{x \in X}{\operatorname{argmin}} \|y - x\|^2$ . Find closed form expression of the projection for the following cases.

- $X_r = \{x \in \mathbb{R}^n \mid \|x\|_2 \leq r\}$ .

$$x_l, x_u \in \mathbb{R}^n$$

$$\Pi_{X_r}(y) = \frac{r y}{\|y\|_2},$$

$$\begin{aligned} \|\Pi_{X_r}(y)\|_2 &= \left\| \frac{r y}{\|y\|_2} \right\|_2 \\ &= \frac{r}{\|y\|_2} \cdot \|y\|_2 = r \end{aligned}$$

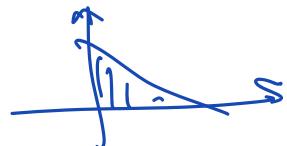
- $X_2 = \{x \in \mathbb{R}^n \mid x_l \leq x \leq x_u\}$ .

$$\Pi_{X_2}(y) = \begin{cases} y_i & \text{if } (x_e)_i \leq y_i \leq (x_u)_i \\ (x_u)_i & \text{if } y > (x_u)_i \\ (x_e)_i & \text{if } y < (x_e)_i \end{cases} \Rightarrow \Pi_{X_2}(y) \in X_2$$

~~•  $X_3 = \{x \in \mathbb{R}^n \mid Ax = b\}$ .~~

- $X_y = \{x \in \mathbb{R}^n \mid x \geq 0, \sum_{i=1}^n x_i \leq 1\}$ .

$$\Pi_{X_y}(y) = \begin{cases} y & \text{if } y \in X_y \\ \max(y - \mu^*, 0) & \text{where } \mu^* \text{ is the solution of} \end{cases}$$



~~• HW~~

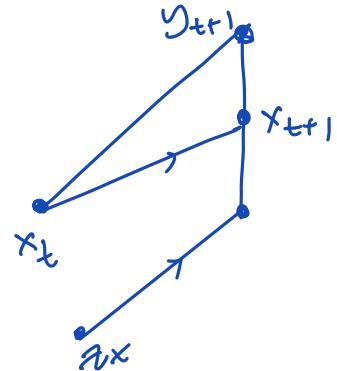
$$\Pi_{X_y}(y) = \max(y_i - \mu^*, 0)$$

$$\sum_{i=1}^n \max(y_i - \mu^*, 0) = 1 \quad (\Pi_{X_y}(y))_i$$

## Accelerated Gradient Descent

Start with  $x_0 = y_0 = z_0 \in \mathbb{R}^n$ . At every time-step  $t$ ,

$$\begin{aligned} y_{t+1} &= x_t - \frac{1}{\beta} \nabla f(x_t) \\ z_{t+1} &= z_t - \eta_t \nabla f(x_t) \\ x_{t+1} &= (1 - \tau_{t+1}) y_{t+1} + \tau_{t+1} z_{t+1} \end{aligned}$$



### Theorem 6

Let  $f$  be  $\beta$ -smooth,  $\eta_t = \frac{t+1}{2\beta}$  and  $\tau_t = \frac{2}{t+2}$ . Then, we have

$$f(y_T) - f(x^*) \leq \frac{2\beta \|x_0 - x^*\|^2}{T(T+1)}.$$

Proof: Define  $\phi_t = t(t+1)(f(y_t) - f(x^*)) + 2\beta \|z_t - x^*\|^2$  and show that  $\phi_{t+1} \leq \phi_t$ .

$$\phi_0 = 2\beta \|x_0 - x^*\|^2$$

$$\phi_T = T(T+1) [f(y_T) - f(x^*)] + \text{const.}$$

If  $\phi_{t+1} \leq \phi_t$  for  $t$ , then  $\phi_T \leq \phi_0$

$$\Rightarrow f(y_T) - f(x^*) \leq \frac{2\beta \|x_0 - x^*\|^2}{T(T+1)} \quad \checkmark$$

HW

## Accelerated Gradient Descent 2

---

Start with  $x_0 = y_0$ . At every state  $t$ ,

$$\begin{aligned} y_{t+1} &= x_t - \frac{1}{\beta} \nabla f(x_t) \\ \underline{x_{t+1}} &= (1 + \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}) \underline{y_{t+1}} - \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \underline{y_t} \end{aligned}$$

### Theorem 7

Let  $f$  be  $\beta$ -smooth,  $\alpha$ -strongly convex with  $\kappa = \frac{\beta}{\alpha}$  and let  $\gamma = \frac{1}{\sqrt{\kappa} - 1}$ . Then, we have

$$f(y_T) - f(x^*) \leq \underbrace{(1 + \gamma)^{-T}}_{\gamma} \left( \frac{\alpha + \beta}{2} \|x_0 - x^*\|^2 \right).$$

Improvement upon the previous rate where we had  $\gamma = \frac{1}{\kappa - 1}$ .

## Further details

---

- AGD invented by Nesterov in a series of papers in the 80s and early 2000s, later popularized by ML researchers
- The convergence rates in the previous two theorems are the best possible ones.
- Book by Nesterov:  
<https://link.springer.com/book/10.1007/978-1-4419-8853-9>
- <https://francisbach.com/continuized-acceleration/>
- <https://www.nowpublishers.com/article/Details/0PT-036>

weights  $w$  : s.t.  $w^T \phi(x) \approx y$  . Given  $(\hat{x}_i, \hat{y}_i), i=1, 2, \dots, N$

$$\min_w \sum_{i=1}^N (\hat{y}_i - w^T \phi(\hat{x}_i))^2$$

**Finite Sum Setting**

- A large number of problems that arise in (supervised) ML can be written as

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{1}{N} \sum_{i=1}^N f_i(x) = \frac{1}{N} \sum_{i=1}^N l(x, \xi_i).$$

weights  
data

- Example: Regression/Least Squares, SVM, NN Training

- The above problem can also be viewed as *sample average approximation* of a stochastic optimization problem

$$f(x) = \mathbb{E}[l(x, \xi)]$$

involving uncertain parameter or random variable  $\xi$ .

- Challenge:  $N$  (number of samples) or  $n$  (dimension of decision variable) both may be large. Samples may be located in different servers.

$$\min_{x \in \mathbb{R}^n} f(x), \quad f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x_i)$$

## Gradient Descent vs. Stochastic Gradient Descent

$$\nabla f(x) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_i)$$

Gradient Descent (GD)  $x_{t+1} = x_t - \eta_t \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_t)$ ,  $t \geq 0$  starting from an initial guess  $x_0 \in \mathbb{R}^n$ .

Each step requires  $N$  gradient computations.

Stochastic Gradient Descent (SGD) At every time step  $t$ ,

- Pick an index (sample)  $i_t$  uniformly at random from the set  $\{1, 2, \dots, N\}$ .
- Set  $x_{t+1} = x_t - \eta_t \nabla f_{i_t}(x_t)$ .

Each step requires 1 gradient computation, which is a noisy version of the true gradient of the cost function at  $x_t$ .

$$f(x) = \|Ax - b\|_2^2 = \sum_{i=1}^n (a_i^T x - b_i)^2$$

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

## Key result for SGD convergence

$$f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x)$$

Under the following assumptions

- Convexity: each  $f_i$  is convex,
- Bounded variance:  $\mathbb{E}[\|\nabla f_{i_t}(x)\|^2] \leq \sigma^2$  for some  $\sigma$  for all  $x$ ,
- Unbiased gradient estimate:  $\mathbb{E}[\nabla f_{i_t}(x)] = \nabla f(x)$  for all  $x$ ,

the solutions generated by SGD algorithm satisfies

$$\begin{aligned} x_{t+1} - x_t &= -\eta_t \nabla f_{i_t}(x_t) \\ &= -\eta_t \nabla f_{i_t}(x_t) \\ &\quad \sum_{t=0}^{T-1} \eta_t [\mathbb{E}[f(x_t)] - f(x^*)] \leq \frac{1}{2} \|x_0 - x^*\|^2 + \frac{\sigma^2}{2} \sum_{t=0}^{T-1} \eta_t^2 \\ &\quad \mathbb{E}[f(\hat{x}_T)] - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2 \sum_{t=0}^{T-1} \eta_t} + \frac{\sigma^2}{2} \frac{\sum_{t=0}^{T-1} \eta_t^2}{\sum_{t=0}^{T-1} \eta_t}, \end{aligned}$$

where  $\hat{x}_T = \frac{1}{\sum_{t=0}^{T-1} \eta_t} \sum_{t=0}^{T-1} \eta_t x_t$ .

We compute

$$\begin{aligned} \|x_{t+1} - x^*\|_2^2 &= \|x_{t+1} - x_t + x_t - x^*\|_2^2 \\ &= \|x_{t+1} - x_t\|_2^2 + \|x_t - x^*\|_2^2 + 2(x_{t+1} - x_t)^T (x_t - x^*) \end{aligned}$$

$$\Rightarrow \|x_{t+1} - x^*\|_2^2 - \|x_t - x^*\|_2^2$$

$$= \eta_t^2 \|\nabla f_{i_t}(x_t)\|_2^2 - 2\eta_t \nabla f_{i_t}(x_t)^T (x_t - x^*)$$

taking expectation on both sides:

$$\Delta \phi_t = \mathbb{E} [\|x_{t+1} - x^*\|_2^2] - \mathbb{E} [\|x_t - x^*\|_2^2]$$

$$\begin{aligned} &= \eta_t^2 \mathbb{E} [\|\nabla f_{i_t}(x_t)\|_2^2] - 2\eta_t \mathbb{E} [\nabla f_{i_t}(x_t)^T (x_t - x^*)] \\ &\quad \text{20} \\ &\quad \mathbb{E} [i_t | x_t] = \mathbb{E} [i_t] \end{aligned}$$

## Proof Continues

$$\begin{aligned}
 \text{Now: } \Delta \phi_t &\leq \eta_t^2 \sigma^2 - 2\eta_t \underbrace{\left( \mathbb{E}_{i_t|x_t} [\nabla f_{i_t}(x_t)] \right)^T (x_t - \bar{x}^*) } \\
 &= \eta_t^2 \sigma^2 - 2\eta_t \underbrace{\nabla f(x_t)^T (x_t - \bar{x}^*) } \\
 &\leq \eta_t^2 \sigma^2 - 2\eta_t [f(x_t) - f(\bar{x}^*)] \quad \left| \begin{array}{l} \text{Since function } f \text{ is convex,} \\ f(\bar{x}^*) \geq f(x_t) + \nabla f(x_t)^T (\bar{x}^* - x_t) \\ \Rightarrow \nabla f(x_t)^T (x_t - \bar{x}^*) \geq f(x_t) - f(\bar{x}^*) \end{array} \right. \\
 \Rightarrow 2\eta_t [f(x_t) - f(\bar{x}^*)] &\leq \eta_t^2 \sigma^2 - \Delta \phi_t \\
 &= - \mathbb{E} \left[ \|x_{t+1} - \bar{x}^*\|_2^2 \right] + \mathbb{E} \left[ \|x_t - \bar{x}^*\|_2^2 \right]
 \end{aligned}$$

$$\begin{aligned}
 2 \sum_{t=0}^{T-1} \eta_t f(x_t) - \underbrace{\left( 2 \sum_{t=0}^{T-1} \eta_t \right) f(\bar{x}^*)} &\leq \sigma^2 \sum_{t=0}^{T-1} \eta_t^2 + \|x_0 - \bar{x}^*\|_2^2 + \eta_t^2 \sigma^2 \\
 &\quad - \mathbb{E} \left[ \|x_T - \bar{x}^*\|_2^2 \right] \\
 &\leq \sigma^2 \sum_{t=0}^{T-1} \eta_t^2 + \|x_0 - \bar{x}^*\|_2^2 \\
 \Rightarrow \frac{1}{\sum_{t=0}^{T-1} \eta_t} \sum_{t=0}^{T-1} \eta_t f(x_t) - f(\bar{x}^*) &\leq \frac{\sigma^2}{2} + \frac{\|x_0 - \bar{x}^*\|_2^2}{2 \sum_{t=0}^{T-1} \eta_t} \\
 f(\hat{x}_T) &\leq \left( \frac{\sum_{t=0}^{T-1} \eta_t^2}{\sum_{t=0}^{T-1} \eta_t} \right) f(\bar{x}^*) + \frac{\|x_0 - \bar{x}^*\|_2^2}{2 \sum_{t=0}^{T-1} \eta_t} \\
 \Rightarrow f(\hat{x}_T) - f(\bar{x}^*) &\leq \left( \frac{\sum_{t=0}^{T-1} \eta_t^2}{\sum_{t=0}^{T-1} \eta_t} \right) f(\bar{x}^*) - f(\bar{x}^*)
 \end{aligned}$$

## **Proof Continues**

---

## **Proof Continues**

---

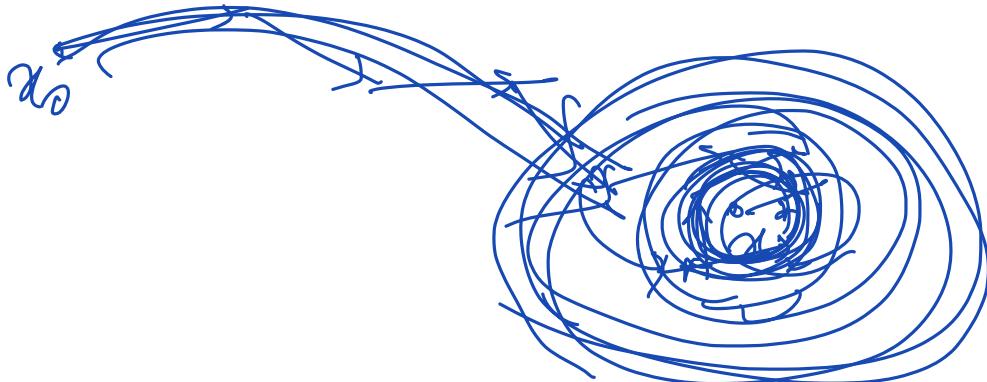
## Choice of stepsize

Constant step-size will not give us convergence. For convergence, we need to choose step sizes that are diminishing and square-summable, i.e.,

$$\lim_{T \rightarrow \infty} \sum_{t=0}^{T-1} \eta_t = \infty, \quad \lim_{T \rightarrow \infty} \sum_{t=0}^{T-1} \eta_t^2 < \infty.$$

- If  $\eta_t := \frac{1}{c\sqrt{t+1}}$ , then  $\mathbb{E}[f(\hat{x}_T)] - f(x^*) \leq \mathcal{O}\left(\frac{\log T}{\sqrt{T}}\right)$ . This rate does not improve when the function is smooth.
- When the function is smooth, then for  $\eta_t := \eta$  chosen appropriately, then R.H.S. will be of order  $\mathcal{O}\left(\frac{1}{\eta T}\right) + \mathcal{O}(\eta)$ .

Recall that for vanilla GD under bounded gradient assumption, we had  $f(\hat{x}_T) - f(x^*) \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$



## Analysis for Smooth and Strongly Convex Functions

---

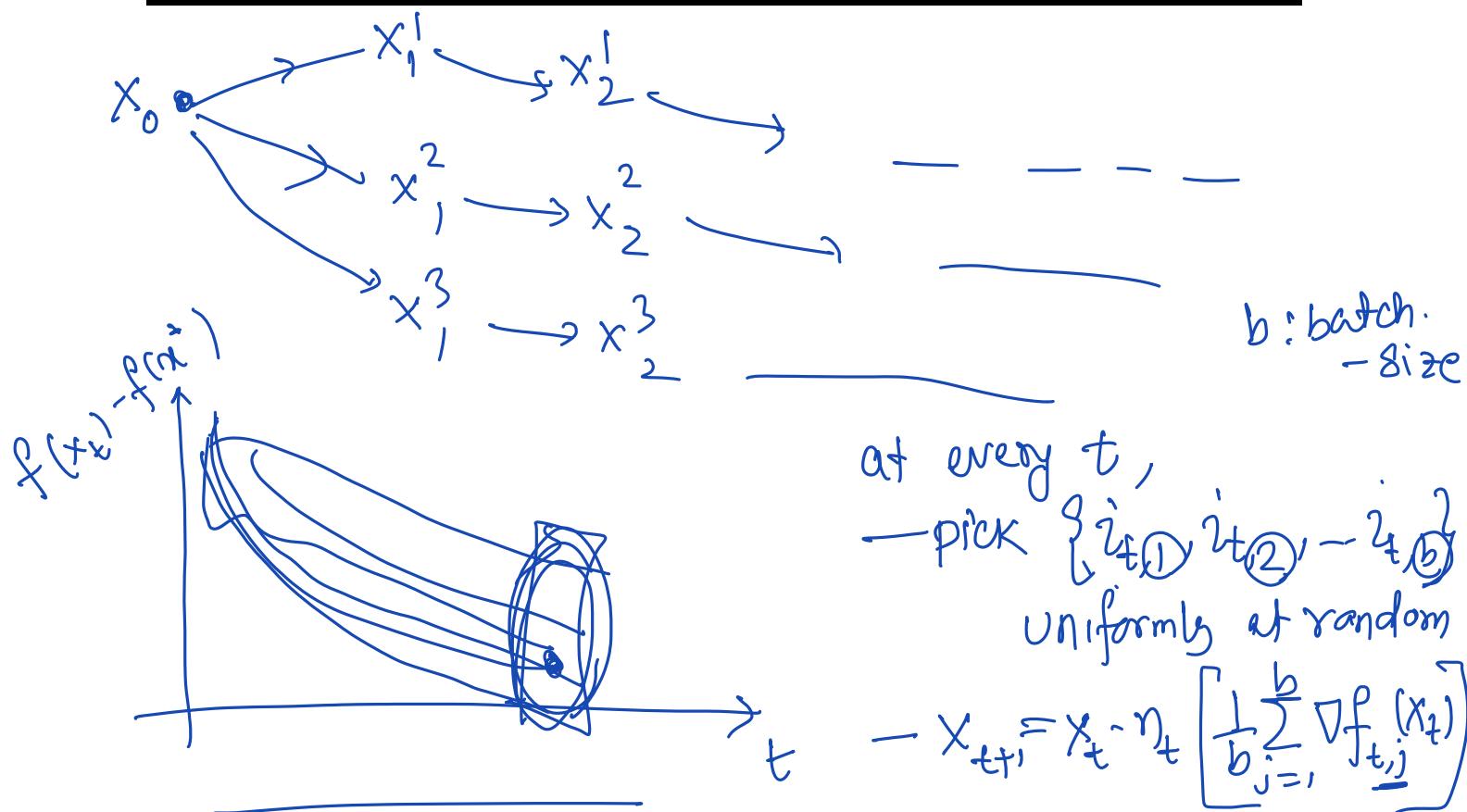
When the function  $f$  is  $\beta$ -smooth and  $\alpha$ -strongly convex, we have the following guarantees for SGD after  $T$  iterations.

- If  $\eta_t := \frac{1}{ct}$  for a suitable constant  $c$ , then error bound is  $\mathcal{O}\left(\frac{\log T}{T}\right)$ . Can be improved to  $\mathcal{O}\left(\frac{1}{T}\right)$ .
- If  $\eta_t := \eta$ , then error bound

$$\mathbb{E}[||x_T - x^*||^2] \leq (1 - \eta\alpha)^T ||x_0 - x^*||^2 + \frac{\eta\beta\sigma^2}{2\alpha}.$$

With constant step-size  $\eta < \frac{1}{\alpha}$ , convergence is quick to a neighborhood of the optimal solution.

## Extension: Mini-Batch



When  $b=1$ : SGD

$b=N$ : GD

The variance of error is smaller by a factor of  $b$  compared to vanilla SGD

## Extension: Stochastic Averaging

---

$$x_{t+1} = x_t - \eta_t g_t, \quad g_t = \frac{1}{N} \sum_{i=1}^N g_t^i,$$
$$g_t^i = \begin{cases} \overline{g_{t-1}^i} & \text{if } i \neq i_t \\ \nabla f_i(x_t) & \text{if } i = i_t \end{cases}$$

## Further Reading

---

~~SAG~~: Schmidt, Mark, Nicolas Le Roux, and Francis Bach. "Minimizing finite sums with the stochastic average gradient." *Mathematical Programming* 162 (2017): 83-112.

~~SAGA~~: Defazio, Aaron, Francis Bach, and Simon Lacoste-Julien. "SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives." *Advances in neural information processing systems* 27 (2014).

~~Recent Review~~: Gower, Robert M., Mark Schmidt, Francis Bach, and Peter Richtárik. "Variance-reduced methods for machine learning." *Proceedings of the IEEE* 108, no. 11 (2020): 1968-1983.

~~Allen-Zhu~~, Zeyuan. "Katyusha: The First Direct Acceleration of Stochastic Gradient Methods." *Journal of Machine Learning Research* 18 (2018): 1-51.

~~Varre~~, Aditya, and Nicolas Flammarion. "Accelerated SGD for non-strongly-convex least squares." In *Conference on Learning Theory*, pp. 2062-2126. PMLR, 2022.

~~Hanzely~~, Filip, Konstantin Mishchenko, and Peter Richtárik. "SEGA: Variance reduction via gradient sketching." *Advances in Neural Information Processing Systems* 31 (2018).

## Extension: Adaptive Step-sizes

---

AdaGrad Duchi, John, Elad Hazan, and Yoram Singer. "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of machine learning research* 12, no. 7 (2011).

Adam Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." *arXiv preprint arXiv:1412.6980* (2014).