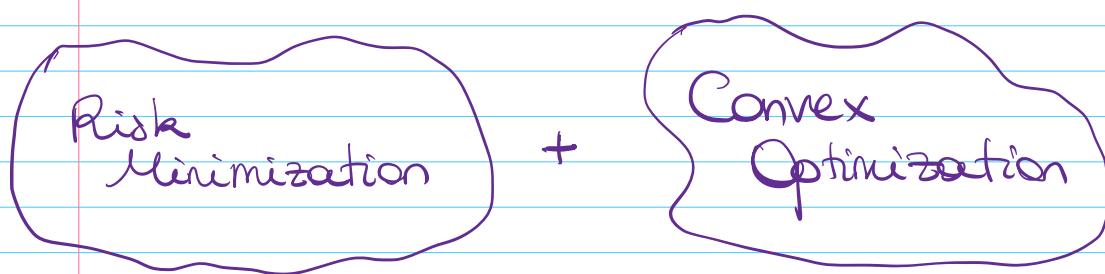


# STOCHASTIC CONVEX OPTIMIZATION



Back to the basics :

$$(x_1, y_1) \dots (x_n, y_n) \sim P$$

$\left\{ \begin{array}{l} X \text{ set of inputs} \\ Y \text{ set of outputs} \end{array} \right.$

$l(y, y')$  loss function.

Goal: Find a predictor  $f: X \rightarrow Y$

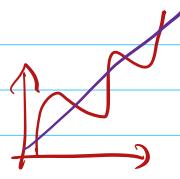
that minimizes the P-risk & test risk

$$R_p(f) = \mathbb{E}_p[\ell(f(x), y)]$$

To do so: introduce a class of predictors

$$f_\theta(x) = \langle x, \theta \rangle \quad \mathcal{F} = \{ f_\theta : X \rightarrow Y : \theta \in \mathbb{R}^d \}$$

$\Rightarrow$  Consider  $\theta^* = \underset{\theta}{\operatorname{arg\min}} R_p(f_\theta)$



How we did so far:

$$\frac{1}{n} \sum_i \ell(f_\theta(x_i), y_i)$$

① Approximate  $R_p(f_\theta)$  by  $R_n(f_\theta)$

② Minimize  $R_n(f_\theta)$  using GD.

$$\theta = (x^T x)^{-1} x^T y$$

(or closed form in certain cases)

$$R_n(f_\theta) = \frac{1}{n} \sum_{i=1}^n \ell_\theta(f_\theta(x_i), y_i)$$

$$\nabla R_n(f_\theta) = \underbrace{\frac{1}{n} \sum_{i=1}^n}_{n \text{ gradients to compute}} \underbrace{\nabla \ell_\theta(f_\theta(x_i), y_i)}$$

$$\text{GD: } \theta^{t+1} = \theta^t - s \nabla R_n(f_\theta)$$

$\Rightarrow T$  steps of GD =  $O(nT)$



Can we do better?

Idea: Apply GD directly on

$$\theta \mapsto R_p(f_\theta) = E_p[\ell(f_\theta(x), y)]$$

$$\nabla R_p(f_\theta) = E_p[\nabla_\theta \ell(f_\theta(x), y)]$$

... But we do not have access to  $\nabla R_p(f_\theta)$

$\Rightarrow$  Estimate it by  $\nabla_\theta \ell(f_\theta(x_i), y_i)$

$$[ \quad \theta_{t+1} = \theta_t - \gamma \nabla_\theta \ell(f_{\theta_t}(x_i), y_i) \quad ]$$

STOCHASTIC GRADIENT DESCENT

$\rightsquigarrow$  1 step of SGD =  $O(1)$

# ① Stochastic Gradient Descent

A more general setting

$$F: \mathbb{R}^d \rightarrow \mathbb{R}$$

SGD : .  $\theta^t$  initialization

For  $t = 1 \dots T-1$  :

- $v^t$  random vector such that

$$\mathbb{E}[v^t | \theta^t] = \nabla F(\theta^t)$$

Unbiased  
estimate  
of the gradient

$$\theta^{t+1} = \underbrace{\theta^t}_{\text{ }} - \underbrace{\gamma_t}_{\text{ }} \underbrace{v^t}_{\text{ }}$$

Output : average  $\bar{\theta} = \frac{1}{T} \sum_{t=1}^T \theta_t$

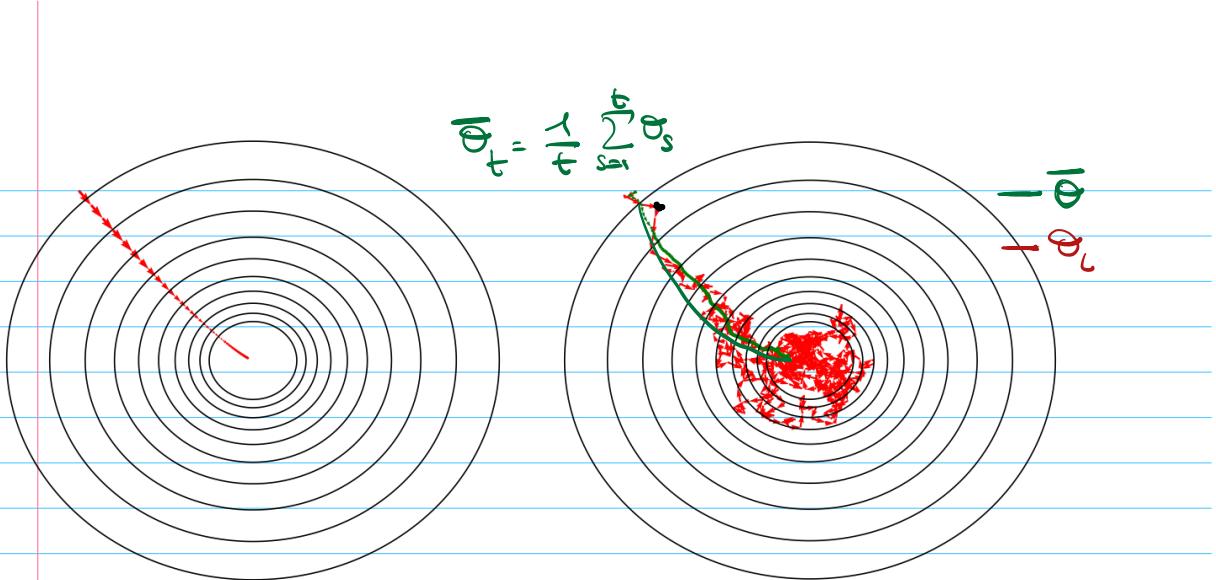
$$\alpha \cdot b \sim \| \theta \|^2$$

Toy Example :  $F(\theta_1, \theta_2) = a\theta_1^2 + b\theta_2^2$  ✓

$$v^t = \nabla F(\theta_1^t, \theta_2^t) + \epsilon$$

$\epsilon$

random noise



## GRADIENT DESCENT

## STOCHASTIC GRADIENT DESCENT

→ For risk minimization:

$$F(\theta) = R_p(f_\theta)$$

$T = n$  (number of samples)

$$v^i = \mathbb{E}[\nabla \ell(f_\theta(x_i), y_i)] = \nabla \mathbb{E}[\ell(f_\theta(x), y)]$$

$$\mathbb{E}[v^i | \theta^i] = \mathbb{E}[v^i] = \nabla F(\theta)$$

$\underbrace{v^i}_{f}$   
only depends  
on  $(x_j, y_j)_{j < i}$

## THEOREM:

$$\left\{ \begin{array}{l} F: \mathbb{R}^d \rightarrow \mathbb{R} \quad \text{convex + differentiable} \\ \text{minimizer } \theta^* \in \mathcal{B}(0; R) \quad R = \|\theta^*\| \end{array} \right.$$

$$\mathbb{E}[\|\nabla^t\|^2] \leq p^2$$

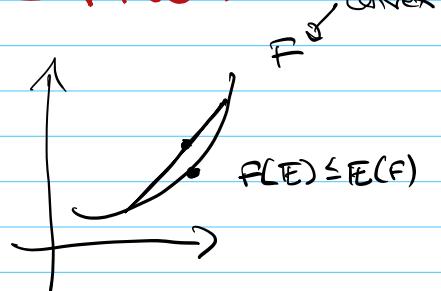
$$\left\{ \begin{array}{l} \text{Step size } \delta = \sqrt{\frac{4R^2}{p^2 T}} \\ \text{Initialization } \theta^0 \in \mathcal{B}(0; R) \end{array} \right.$$

$$\Rightarrow \mathbb{E}[F(\bar{\theta})] - F(\theta^*) \leq 2 \frac{Rp}{\sqrt{T}}$$

proof 1: special case  $\nabla^t = \nabla F(\theta^t)$

proof 2: General case

$$\bar{\theta} = \frac{1}{T} \sum_{t=1}^T \theta_t$$

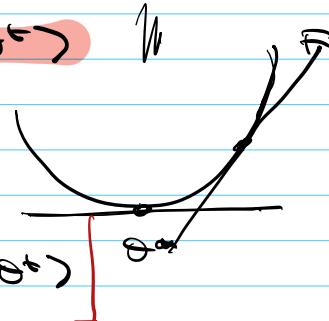


$$F(\bar{\theta}) - F(\theta^*) \leq \frac{1}{T} \sum_{t=1}^T F(\theta_t) - F(\theta^*)$$

$$F(\theta^*) \geq F(\theta^t) + \langle \nabla F(\theta^t), \theta^* - \theta^t \rangle$$

$$\langle \nabla F(\theta^t), \theta^* - \theta^t \rangle \geq F(\theta^t) - F(\theta^*)$$

$$(+) \left[ F(\bar{\theta}) - F(\theta^*) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla F(\theta^t), \theta^* - \theta^t \rangle \right]$$



$$\mathbb{E}[v^t | \theta^t] = \nabla F(\theta^t)$$

$$\Rightarrow \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\langle v^t, \theta^* - \theta^t \rangle | \theta^t]$$

$$\begin{aligned} \mathbb{E}[F(\theta) - F(\theta^*)] &\leq \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\langle v^t, \theta^* - \theta^t \rangle] \\ &= \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \langle v^t, \theta^* - \theta^t \rangle\right] \end{aligned}$$

Claim:  $\langle v^t, \theta^t - \theta^* \rangle = \frac{1}{2s} (||\theta^t - \theta^*||^2 - ||\theta^{t+1} - \theta^*||^2 + \frac{s}{2} ||v^t||^2)$

$$\theta^{t+1} = \theta^t - s v^t$$

$$||\theta^t - \theta^* - s v^t||^2$$

$$\theta^{t+1} = \theta^t - s \nabla F(\theta^t)$$

(\*) | Claim:  $\langle \nabla F(\theta^t), \theta^* - \theta^t \rangle$

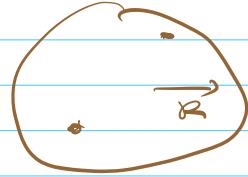
$$= \frac{1}{2s} \left( \|\theta^t - \theta^*\|^2 - \|\theta^{t+1} - \theta^*\|^2 \right) + \frac{s}{2} \|\nabla F(\theta^t)\|^2$$

$$\frac{1}{T} \sum_{t=1}^T \langle \nabla F(\theta^t), \theta^* - \theta^t \rangle = \frac{1}{T} \sum_{t=1}^T \frac{1}{2s} \left( \|\theta^t - \theta^*\|^2 - \|\theta^{t+1} - \theta^*\|^2 \right) + \frac{1}{T} \frac{s}{2} \sum_{t=1}^T \|\nabla F(\theta^t)\|^2$$

$$\leq \frac{1}{T2s} \left( \underbrace{\|\theta^t - \theta^*\|^2}_{\leq 4R^2} - \underbrace{\|\theta^T - \theta^*\|^2}_{\leq 0} \right) + \frac{s}{2} \rho^2$$

$\theta^*, \theta^t \in B(\theta; R)$

$$\leq \frac{4R^2}{2Ts} + \frac{s}{2} \rho^2.$$



$$\begin{aligned} & \sum_{t=1}^T \left( \|\theta^t - \theta^*\|^2 - \|\theta^{t+1} - \theta^*\|^2 \right) \\ &= \|\theta^1 - \theta^*\|^2 - \cancel{\|\theta^2 - \theta^*\|^2} \\ &\quad + \cancel{\|\theta^2 - \theta^*\|^2} - \cancel{\|\theta^3 - \theta^*\|^2} \\ &\quad + \cancel{\|\theta^3 - \theta^*\|^2} - \cancel{\|\theta^4 - \theta^*\|^2} \\ &\quad + \cancel{\|\theta^T - \theta^*\|^2} - \|\theta^T - \theta^*\|^2 \end{aligned}$$

Remark: in practice, a good idea is

to discard the first iterates when

computing  $\bar{\theta}$ : *forget about initialization*

$$\bar{\theta} = \frac{1}{T-T_0} \sum_{t=T_0}^T \theta_t$$

Another Example: *arbitrary functions*

$$F(\theta) = \frac{1}{n} \sum_{i=1}^n F_i(\theta)$$

Let  $I \sim \text{Unif}\{1 \dots n\}$

$$\Rightarrow F(\theta) = E[F_I(\theta)] = \sum_{i=1}^n P(I=i) F_i(\theta)$$

To get estimates of  $\nabla F(\theta)$ :

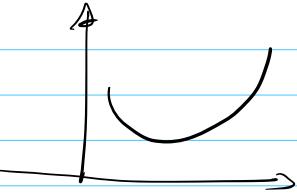
$$V^t = \nabla F_{I_t}(\theta)$$

$I \dots I_T$   
iid samples.

$d=1$        $\beta$  smooth:  $F''(\theta) \leq \beta$   
 $\alpha$  strongly conv:  $F''(\theta) \geq \alpha$

Comparison with GD:

Recall:



①  $F$   $\beta$ -smooth +  $\alpha$ -strongly convex

$T$  steps of GD ;  $s = 1/\beta$

$$F(\theta^t) - F(\theta^*) \leq \exp\left(-\frac{\alpha}{\beta} T\right)$$

$T = O(\log(1/\epsilon))$  steps to reach precision  $\epsilon$ .

②  $F$   $\beta$ -smooth + conv

$T$  steps of GD ;  $s = 1/\beta$

$$F(\theta^t) - F(\theta^*) \leq \frac{\beta}{T}$$

$T = O(1/\epsilon)$  steps to reach precision  $\epsilon$ .

③  $F$  conv + diff

$T$  steps of SGD

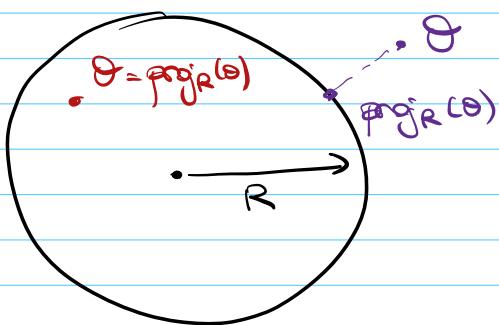
$$E[F(\bar{\theta})] - F(\theta^*) \leq \frac{1}{\sqrt{T}}$$

**Question :** Can we improve the rate  
for SGD with more assumptions?

→  $\beta$ -smoothness? no no

→  $\alpha$ -strongly conv? no yes! 😊

$$\text{proj}_R(\theta) = \begin{cases} \theta & \text{if } \|\theta\| \leq R \\ \frac{R}{\|\theta\|}\theta & \text{else.} \end{cases}$$



## SGD with projection:

For  $t = 1, \dots, T-1$ :

- $v^t$  random vector such that

$$\mathbb{E}[v^t | \theta^t] = \nabla F(\theta^t)$$

- $\theta^{t+1} = \text{proj}_Q(\theta^t - \gamma_t v^t)$

Output: average  $\bar{\theta} = \frac{1}{T} \sum_{t=1}^T \theta_t$ .

## THEOREM:

$F: \mathbb{R}^d \rightarrow \mathbb{R}$   $\alpha$ -strongly convex  
 minimizer  $\theta^* \in \text{B}(0; R)$

$$\mathbb{E}[\|v^t\|^2] \leq \rho^2$$

Step size  $\gamma_t = 1/\alpha t$

Initialization  $\theta^0 \in \text{B}(0; R)$

$$\Rightarrow \mathbb{E}[F(\bar{\theta})] - F(\theta^*) \leq \frac{\rho^2}{2\alpha T} (1 + \log T)$$

$$\text{proof : } F(\theta) - F(\theta^*) \leq \frac{1}{T} \sum (F(\theta^t) - F(\theta^*))$$

$$\langle \nabla F(\theta^t), \theta^t - \theta^* \rangle = \frac{1}{2\alpha_t} \left( \|\theta^t - \theta^*\|^2 - \|\theta^{t+1} - \theta^*\|^2 \right) + \frac{\alpha_t}{2} \|v^t\|^2$$

$$E[F(\theta^t) - F(\theta^*)] \leq E[\langle \nabla F(\theta^t), \theta^t - \theta^* \rangle - \frac{\alpha_t}{2} \|\theta^t - \theta^*\|^2]$$

$$E \left( \frac{1}{T} \sum (F(\theta^t) - F(\theta^*)) \right) \leq$$

$$E[v^t | \theta^t] = \nabla F(\theta^t)$$

$$E \left( \frac{1}{T} \sum_t \left( \langle \nabla F(\theta^t), \theta^t - \theta^* \rangle - \frac{\alpha_t}{2} \|\theta^t - \theta^*\|^2 \right) \right)$$

$$= E \left[ \frac{1}{T} \sum_t \left( \langle v^t, \theta^t - \theta^* \rangle - \frac{\alpha_t}{2} \|\theta^t - \theta^*\|^2 \right) \right]$$

$$\frac{\alpha_t}{2} (\|\theta^t - \theta^*\|^2 - \|\theta^{t+1} - \theta^*\|^2)$$

$$\xi_t = \frac{1}{\alpha_t}$$

$$+ \frac{1}{2\alpha_t} \|v^t\|^2$$

$$E \left( \frac{1}{T} \sum_t \left( \frac{\alpha_t}{2} \|\theta^t - \theta^*\|^2 - \frac{\alpha_t}{2} \|\theta^{t+1} - \theta^*\|^2 + \frac{1}{2\alpha_t} \|v^t\|^2 \right) \right)$$

$$- \frac{\alpha_t}{2} \|\theta^t - \theta^*\|^2$$

$$\underbrace{\frac{\alpha_t(t-1)}{2} \|\theta^t - \theta^*\|^2}_{\frac{\alpha_t(t-1)}{2} \|\theta^t - \theta^*\|^2}$$

$$\frac{1}{T} \sum (v_t - v_{t+1}) = \frac{1}{T} (v_1 - v_T) = \frac{1}{T} (0 - 0) \leq 0$$

$$\leq \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \frac{1}{2\alpha t} \|v_t^\epsilon\|^2 \right]$$

$$\int_1^T \frac{1}{t} dt = \log T$$

$$\leq \frac{1}{T} \sum_{t=1}^T \frac{\rho^2}{2\alpha t} \leq \frac{\rho^2}{2\alpha T} \underbrace{\left( \sum_{t=1}^T \frac{1}{t} \right)}_{\approx \log T}$$

□

## ② Application to Risk Minimization

Back to risk minimization

$$F(\theta) = R_p(f_\theta) \quad \theta \in \mathbb{R}^k$$

minimizer  $\theta^*$

Sample  $(x_1, y_1), \dots, (x_n, y_n)$

Two methods to approximate  $\theta^*$ :

(SGD)  $\hat{\theta}_{\text{SGD}}$ :  $n$  steps of SGD

with gradients  $\nabla l(f_\theta(x_i), y_i)$

(GD-ER)  $\hat{\theta}_T$ : Gradient Descent for  $T$  steps

applied to

$$R_n(f_\theta) = \frac{1}{n} \sum_{i=1}^n l(f_\theta(x_i), y_i)$$

$$\mathbb{E}[R_p(f_\theta) - R_p^*] = \underbrace{\mathbb{E}[R_p(f_\theta) - R_p(f_{\theta^*})]}_{\text{optimization error}}$$

$$+ \underbrace{R_p(f_{\theta^*}) - R_p^*}_{\text{approximation error}}$$

Two Questions:

(1) What is the minimal number  $n$  of samples required to get an optimization error smaller than  $\epsilon^2$

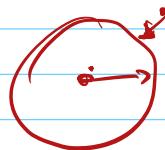
(2) What is the associated time complexity?

→ We answer (1) and (2) in the

setting where:  $\forall x \in \mathcal{X}, y \in \mathcal{Y}$

$\theta \mapsto l(f_\theta(x), y)$  is  $\begin{cases} \alpha\text{-strongly convex} \\ \beta\text{-smooth} \end{cases}$

$$(\text{SGD}) \quad v^i = \nabla l_{f_\theta}(x_i, y_i)$$



$$\text{so } \|v^i\| \leq \beta \|\theta^i - \theta^*\| \leq 2\beta R \quad \pi \epsilon$$

$$\Rightarrow E[R_p(f_{\theta_{\text{SGD}}}^*) - R_p(f_{\theta^*})] \leq \frac{1}{n}$$

(1)  $n = \tilde{\mathcal{O}}(1/\epsilon)$  up to log factors

(2) Time complexity:  $\tilde{\mathcal{O}}(kn)$

$$\theta_{\infty} = \arg\min \mathcal{R}_p(f_{\theta}) \quad \dagger$$

(GD-ER) Let  $\theta_{\infty} = \arg\min \mathcal{R}_n(f_{\theta})$ .

$$\mathcal{R}_p(f_{\theta_{\infty}}) - \mathcal{R}_p(f_{\theta^*}) \leq \mathcal{R}_p(f_{\theta_{\infty}}) - \mathcal{R}_n(f_{\theta_{\infty}})$$

$$+ \mathcal{R}_n(f_{\theta_{\infty}}) - \mathcal{R}_n(f_{\theta^*})$$

$$\mathcal{R}_n(f_{\theta_{\infty}}) \leq \mathcal{R}_n(f_{\theta^*}) \forall \theta$$

$$+ \mathcal{R}_n(f_{\theta_{\infty}}) - \mathcal{R}_n(f_{\theta^*}) \leq 0$$

$$- \mathcal{R}_p(f_{\theta^*}) + \mathcal{R}_n(f_{\theta^*})$$

$$\leq 2 \cdot \sup_{\theta} |\mathcal{R}_n(f_{\theta}) - \mathcal{R}_p(f_{\theta})| + \mathcal{R}_n(f_{\theta_{\infty}}) - \mathcal{R}_n(f_{\theta^*})$$

estimation error

$$\approx \frac{1}{\sqrt{n}}$$

$$\lesssim \exp(-\frac{\alpha}{\beta} T)$$

To get

$$\mathcal{R}_p(f_{\theta_{\infty}}) - \mathcal{R}_p(f_{\theta^*}) \leq \varepsilon$$

①  $n = O(1/\varepsilon^2)$  samples

②  $T = O(\log(\varepsilon^{-1}))$

Complexity  $O(kTn) = \tilde{O}\left(\frac{k}{\varepsilon^2}\right)$

SUMMARY: in the  $\begin{cases} \alpha\text{-strongly convex} \\ \beta\text{-smooth} \end{cases}$  setting

Algo	Num. of samples	Complexity
SGD	$n = \tilde{\mathcal{O}}(1/\varepsilon)$	$\tilde{\mathcal{O}}(k/\varepsilon)$
GD	$n = \mathcal{O}(1/\varepsilon^2)$	$\tilde{\mathcal{O}}(k/\varepsilon^2)$

To go further...

Variance Reduction technique :

SVRG / SAGA

→ see lecture notes for references.