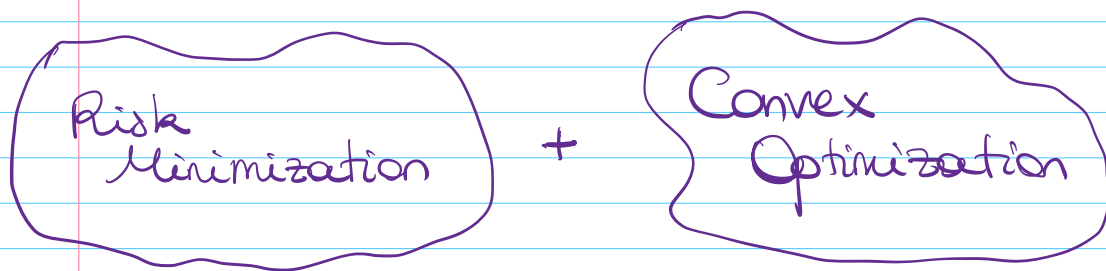


# STOCHASTIC CONVEX OPTIMIZATION



Back to the basics :

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

$X$  set of inputs  
 $Y$  set of outputs

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

Goal: Find a predictor  $f: X \rightarrow Y$   
that minimizes the P-risk  $\leftarrow$  test risk

$$\left[ R_P(f) = \mathbb{E}_P[l(f(x), Y)] \right]$$

To do so: introduce a class of predictors

$$\mathcal{F} = \{ f_\theta: X \rightarrow Y: \theta \in \mathbb{R}^d \}$$

$\Rightarrow$  Consider  $\theta^* = \underset{\theta}{\operatorname{argmin}} R_P(f_\theta)$

How we did so far:

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

② Minimize  $R_n(f_\theta)$  using GD.  
(or closed form in certain cases)

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

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

$n$  gradients to compute

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

⇒ T steps of GD =  $O(\ln T)$



Can we do better?

Idea: Apply GD directly on

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

$$\nabla R_p(f_\Theta) = \mathbb{E}_p[\nabla_\Theta \ell(f_\Theta(x), y)]$$

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

⇒ Estimate it by  $\nabla_\Theta \ell(f_{\Theta_t}(x_t), y_t)$

$$\left[ \Theta_{t+1} = \Theta_t - s \nabla_\Theta \ell(f_{\Theta_t}(x_t), y_t) \right]$$

STOCHASTIC GRADIENT DESCENT

⇒ 1 step of SGD =  $O(1)$

# ① Stochastic Gradient Descent

A more general setting

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

SGD: •  $\theta^1$  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

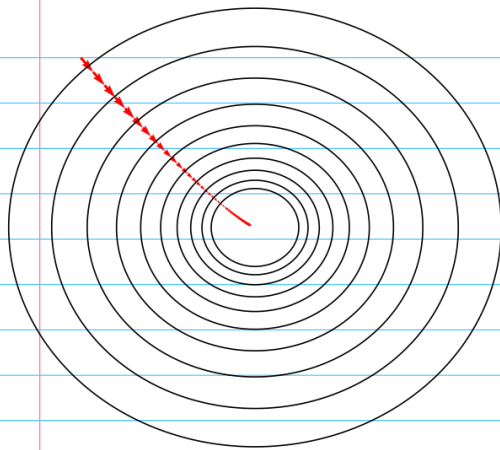
- $\theta^{t+1} = \theta^t - \eta_t V^t$

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

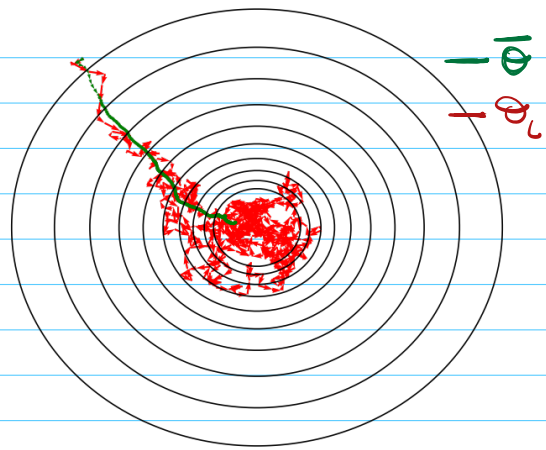
Toy Example:  $F(\theta_1, \theta_2) = a\theta_1^2 + b\theta_2^2$

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

↑  
random noise



GRADIENT  
DESCENT



STOCHASTIC GRADIENT  
DESCENT

→ For risk minimization:

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

$$T = n \quad (\text{number of samples})$$

$$V^i = \nabla \ell(f_{\theta^i}(x_i), y_i)$$

$$E[V^i | \theta^i] = E[V^i] = \nabla F(\theta)$$

↑  
only depends  
on  $(x_j, y_j)_{j < i}$

### THEOREM :

$F: \mathbb{R}^d \rightarrow \mathbb{R}$  convex + differentiable  
minimizer  $\theta^* \in B(0; R)$

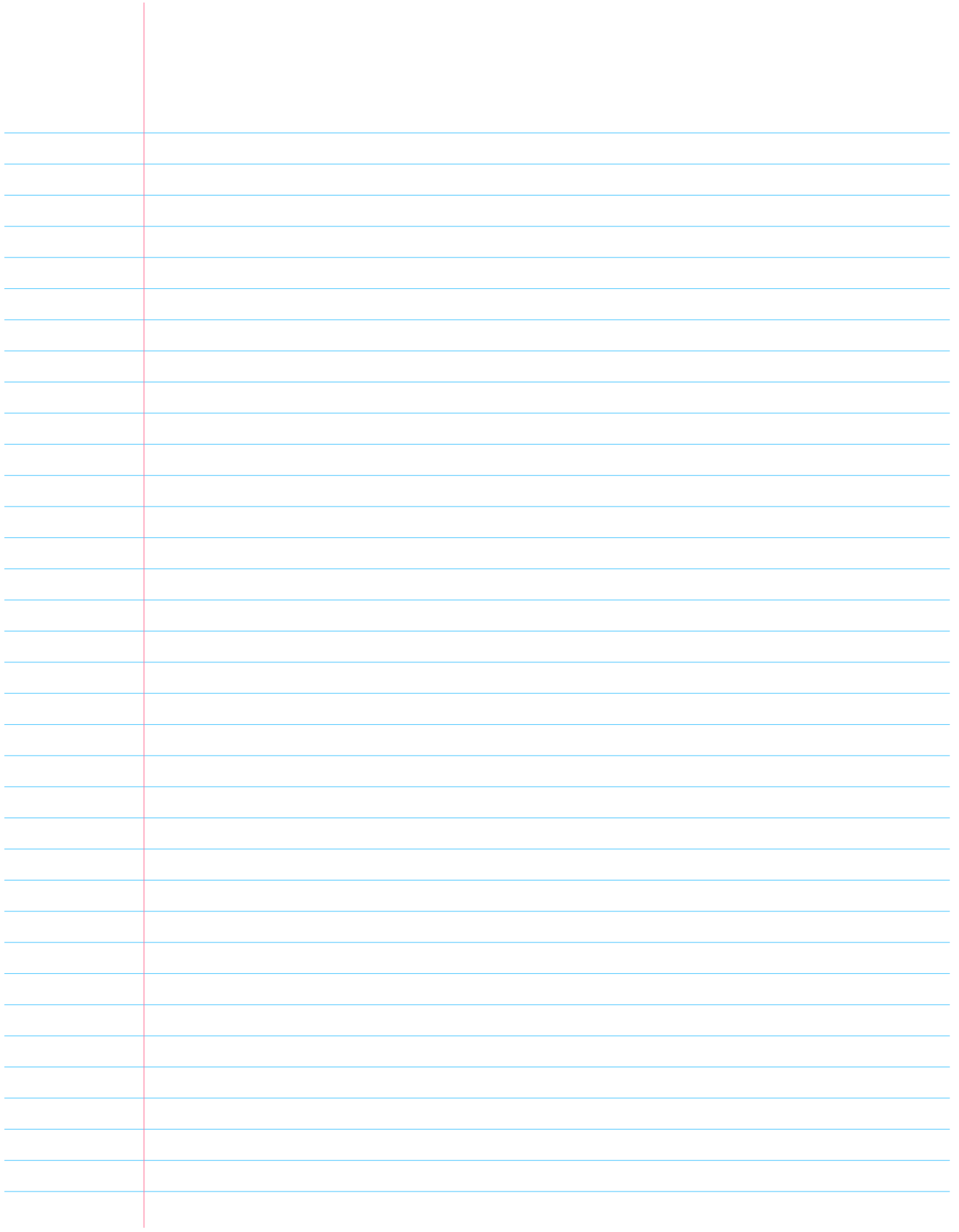
$$E[\|V^t\|^2] \leq \rho^2$$

Step size  $s = \sqrt{\frac{4R^2}{\rho^2 T}}$   
Initialization  $\theta^0 \in B(0; R)$

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

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

proof 2: General case



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

$$\Rightarrow F(\theta) = \mathbb{E}[F_I(\theta)]$$

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

$$V^t = \nabla F_{I^t}(\theta) \quad I^1, \dots, I^T \\ \text{iid samples.}$$



## 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 + cvx

$T$  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  $F$  cvx + diff

$T$  steps of SGD

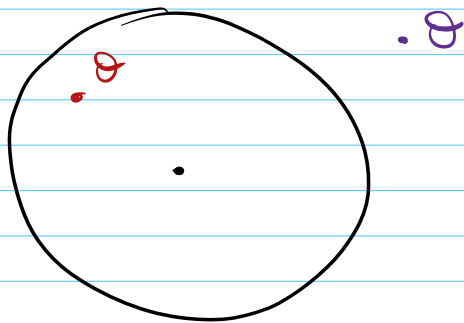
$$E[F(\theta)] - F(\theta^*) \approx \frac{1}{\sqrt{T}}$$

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

→  $\beta$ -smoothness? no no

→  $\alpha$ -strongly convex? 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

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

- $\theta^{t+1} = \text{proj}_R(\theta^t - \delta_t V^t)$

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

## THEOREM:

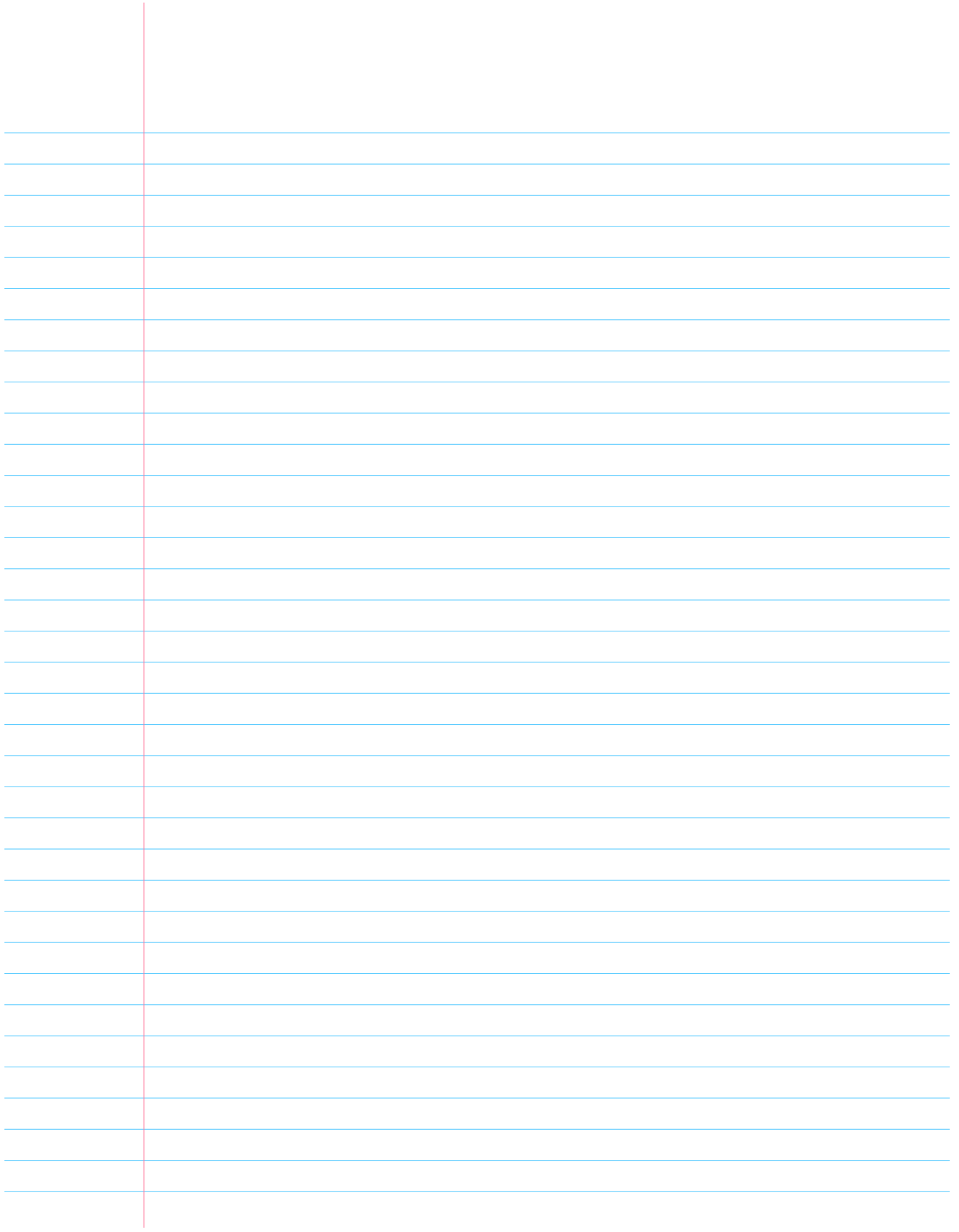
$F: \mathbb{R}^d \rightarrow \mathbb{R}$   $\alpha$ -strongly convex  
minimizer  $\theta^* \in B(0; R)$

$$E[\|V^t\|^2] \leq \rho^2$$

Step size  $\delta_t = 1/\alpha t$   
Initialization  $\theta^0 \in B(0; R)$

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

prof :



## ② Application to Risk Minimization

Back to risk minimization

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

minimizer  $\theta^*$

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

Two methods to approximate  $\theta^*$ :

(SGD)  $\hat{\theta}_{SGD}$ :  $n$  steps of SGD  
with gradients  $\nabla_{\theta} l(f_{\theta}(x_i), y_i)$

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

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

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

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

## Two Questions:

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

② What is the associated time complexity?

⇒ We answer ① and ② in the

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

$\theta \mapsto \ell(f_\theta(x), y)$  is  $\left\{ \begin{array}{l} \alpha\text{-strongly convex} \\ \beta\text{-smooth} \end{array} \right.$

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

$$\text{no } \|v^i\| \leq \beta \|\theta^i - \theta^*\| \leq 2\beta R \quad \epsilon$$

$$\Rightarrow \mathbb{E}[\mathcal{R}_p(f_{\theta_{SGD}}) - \mathcal{R}_p(f_{\theta^*})] \lesssim \frac{1}{n}$$

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

② Time complexity:  $\tilde{O}(kn)$

(GD-ER) Let  $\theta_\infty = \arg\min_{\theta} R_n(f_\theta)$ .

$$R_p(f_{\theta_T}) - R_p(f_{\theta^*}) \leq$$

$$\leq 2 \cdot \underbrace{\sup_{\theta} |R_n(f_\theta) - R_p(f_\theta)|}_{\text{estimation error}} + \underbrace{R_n(f_{\theta_T}) - R_n(f_{\theta_\infty})}_{\lesssim \exp(-\frac{\alpha}{\beta} T)}$$
$$\gtrsim \frac{1}{\sqrt{n}} \qquad \lesssim \exp(-\frac{\alpha}{\beta} T)$$

To get

$$R_p(f_{\theta_T}) - R_p(f_{\theta^*}) \leq \varepsilon$$

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  $\left\{ \begin{array}{l} \alpha\text{-strongly convex} \\ \beta\text{-smooth} \end{array} \right.$  setting

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

To go further . . . .

Variance Reduction technique :

SVRG / SAGA

↳ see lecture notes for references.