Summary of previous lecture

Optimization: ① least-squares } condition number
② smooth optimization }
③ non-smooth optimization + SGD

---

Data: $(x_i, y_i) \in \mathcal{X} \times \mathbb{R}/\{-1,1\}$   learn a function $f: \mathcal{X} \to \mathbb{R}$

Expected risk: $R(f) = \mathbb{E}\,\ell(y, f(x))$ ← testing data

Empirical risk $\hat{R}(f)$
$$= \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(x_i))$$

ERM: class of models $\mathcal{F}$,  $\hat{f} \in \arg\min_{f \in \mathcal{F}} \hat{R}(f)$

<span style="color:red">goal: find $\hat{f}$ s.t. $\hat{R}(\hat{f}) - \inf_{f \in \mathcal{F}} \hat{R} \le \varepsilon$</span>

Estimation error:
$$R(\hat{f}) - \inf_{f \in \mathcal{F}} R(f) = R(\hat{f}) - \hat{R}(\hat{f}) + \underbrace{\hat{R}(\hat{f}) - \hat{R}(f_{\mathcal{F}}^{o})}_{} + \hat{R}(f_{\mathcal{F}}^{*}) - R(f_{\mathcal{F}}^{*})$$

$\underbrace{\quad}_{R(f_{\mathcal{F}}^{*})}$

$\le 0 \Rightarrow \varepsilon$

$$= \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| + \sup_{f \in \mathcal{F}} (R(f) - \hat{R}(f))$$

$\le \quad +\varepsilon$

Tool of choice: Rademacher complexity $= O\left(\frac{1}{\sqrt{n}}\right)$

Optimization = forget about ML

Goal = find a minimizer of $F: \mathbb{R}^d \to \mathbb{R}$

Motivation: ML $\quad F(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$

<span style="color:red">linear model</span>
<span style="color:red">$f_\theta(x_i) = \theta^T \varphi(x_i)$</span>

(1) Convex vs non convex

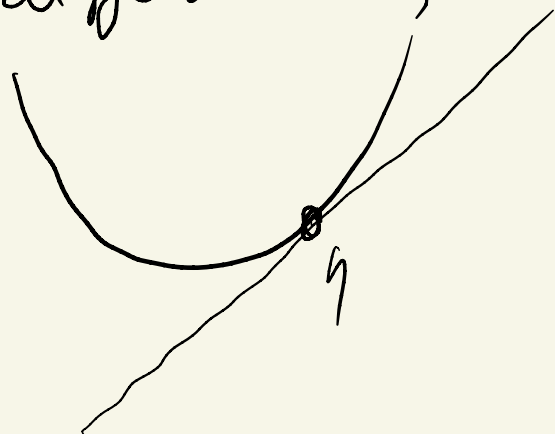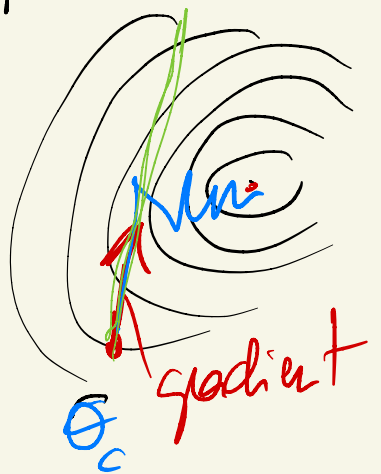<span style="color:red">Positive semi definite matrix</span>

F: convex
- if F twice differentiable, $F''(\theta) \succcurlyeq 0 \quad \forall \theta$
- if F differentiable, F above all of its tangents
$$F(\theta) \geq F(q) + F'(q)^T (\theta - q)$$

<span style="color:red">tangent at q</span>

# gradient descent (Cauchy, 1847)

$$\Theta_t = \Theta_{t-1} - \gamma F'(\Theta_{t-1})$$



$\Theta_c$

gradient

(step-size) $> 0$

constant

decaying with fixed schedule

line search

x Quadratic
x Smooth
o Non-smooth

quadratic function: $F(\theta) = \frac{1}{2}\theta^T H\theta - c^T\theta$, convex $H \succeq 0$

$\eta_*$ minimizer unique

$H$ invertible

$F'(\eta_*) = c = H\eta_* - c \implies \eta_* = H^{-1}c$

Normal equations

gradient descent: $\theta_t = \theta_{t-1} - \gamma[H\theta_{t-1} - c] = \theta_{t-1} - \gamma[H\theta_{t-1} - H\eta_*]$
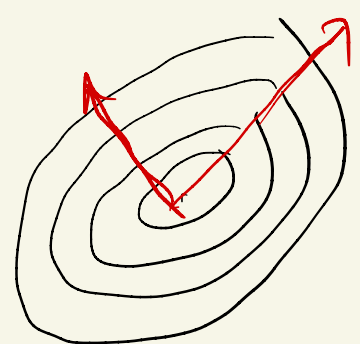
Linear iteration

$\theta_t - \eta_* = \theta_{t-1} - \eta_* - \gamma H(\theta_{t-1} - \eta_*)$

$\theta_t - \eta_* = [1 - \gamma H](\theta_{t-1} - \eta_*) = [1 - \gamma H]^t (\theta_0 - \eta_*)$

Eigenvalue decomposition: $H = \sum_{i=1}^{d} \lambda_i e_i e_i^T$

eigenvector

eigenvalue

$(1 - \gamma H)^t$ has same eigenvectors, and $(1 - \lambda_i \gamma)^t$ as eigenvalues

① when is it convergent? if $|1 - \lambda_i \gamma| < 1, \forall i$

② how fast?

$\implies 1 - \lambda_i \gamma > -1$

$\gamma < \frac{2}{\lambda_i} \forall i$

Def: $\mu$ = smallest eigenvalue

$L$ = largest eigenvalue

$0 < \mu < L$

$\boxed{\gamma < \frac{2}{L}}$

$\boxed{\gamma = \frac{1}{L}}$

$$\theta_t - q_* = (1 - \gamma H)^t (\theta_0 - q_*)$$

$$\|\theta_t - q_*\|^2 \leq \left(\text{largest eigenvalue of } (1-\gamma H)^t\right)^2 \|\theta_0 - q_*\|^2$$

with $\gamma = \frac{1}{L}$
$$(1-\gamma H)^t = (1 - \frac{1}{L} H)^t$$

$\textcolor{red}{\|M_3\|^2 \leq \|3\|^2}$
$\textcolor{red}{(\text{largest eig of})^2}$
$\textcolor{red}{\text{m}}$

with largest eig. $(1 - \frac{\mu}{L})^t$ because $\mu$ is the smallest eig. of $H$

$\textcolor{red}{\text{positive}}$

$$\|\theta_t - q_*\|^2 \leq (1 - \frac{\mu}{L})^{2t} \|\theta_0 - q_*\|^2$$

$$\leq \exp(-\frac{2\mu}{L} t) \|\theta_0 - q_*\|^2$$

$$\leq \boxed{\textcolor{red}{\exp(-2t/\kappa)}} \|\theta_0 - q_*\|^2 \quad \text{with } \kappa = \frac{L}{\mu} \overset{\geq 1}{} \text{ condition number}$$

$\textcolor{red}{\text{linear convergence}}$
$\textcolor{red}{\text{exponential convergence}}$

$\textcolor{red}{1 - d \leq e^{-d}}$

$\textcolor{red}{\varepsilon}$

$\textcolor{blue}{t = \frac{\kappa}{2} \log \frac{1}{\varepsilon}}$



$\textcolor{blue}{\text{small } \kappa = \frac{L}{\mu}}$

$\textcolor{blue}{\text{large } \kappa = \frac{L}{\mu}}$

$\textcolor{red}{\text{How small is } \mu?}$
$\textcolor{red}{\text{How big is } \kappa?}$
$\textcolor{red}{H = \frac{1}{m} \Phi^T \Phi}$
$\textcolor{red}{\text{for least squares}}$

$$\theta_t - \eta_r = (1 - \gamma H)^t (\theta_t - \eta_*) \qquad \gamma = 1/L$$

$$F(\theta_t) - F(\eta_*) = \frac{1}{2}(\theta_t - \eta_*)^T H (\theta_t - \eta_*)$$

<span style="color:red">by Taylor expansion around $\eta_*$</span>

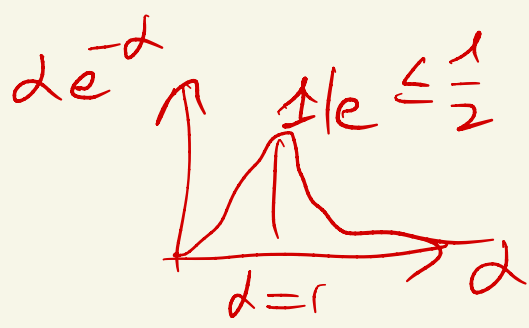<span style="color:red">$F'(\eta_*) = 0 \qquad F''(\eta_*) = H$</span>

$$= \frac{1}{2}(\theta_o - \eta_*)^T \underline{H}(1 - \gamma H)^{2t}(\theta_o - \eta_*)$$

$$\leq \frac{1}{2}\left(\text{largest eig. of } \underline{H}\left(1 - \frac{1}{L}H\right)^{2t}\right) \times \|\theta_o - \eta_*\|^2$$

$$\leq \frac{1}{2} \sup_{\lambda \in [\mu, L]} \lambda\left(1 - \frac{1}{L}\lambda\right)^{2t} \qquad \|\theta_o - \eta_*\|^2$$

$$\lambda\left(1 - \frac{\lambda}{L}\right)^{2t} \leq \lambda e^{-2t\lambda/L} = \frac{2t\lambda}{L} e^{-2t\lambda/L} \frac{L}{2t} \leq \frac{L}{4t}$$

<span style="color:red">$1 - \lambda \leq e^{-\lambda}$</span>

<span style="color:red">$\lambda e^{-\lambda}$</span>

<span style="color:red">$1/e \leq \frac{1}{2}$</span>

<span style="color:red">$\leq 1/2$</span>

<span style="color:red">$\lambda = 1$</span>   <span style="color:red">$\lambda$</span>

linear.

- Consequence $= F(\theta_t) - F(\eta_*) \leq \frac{L}{4t}\|\theta_o - \eta_*\|^2$
- Adaptivity
- Optimality with acceleration

$$k \to \sqrt{k} \qquad / \qquad \frac{1}{t} \to 1/t^2$$

[True as well for convex functions:   Spectrum of $H \subset [\mu, L]$

Assumption: F is smooth "$\Longleftrightarrow$" $\forall \theta$, $F''(\theta)$ has eigenvalues
less than L

$\boxed{\begin{array}{l} 10^{+8} \, 18 \\ 10^{+1} \, 40 \end{array}}$

F is strongly "$\Longleftarrow$" $\forall \theta$, $F''(\theta)$ ———————
convex                                        larger than $\mu$
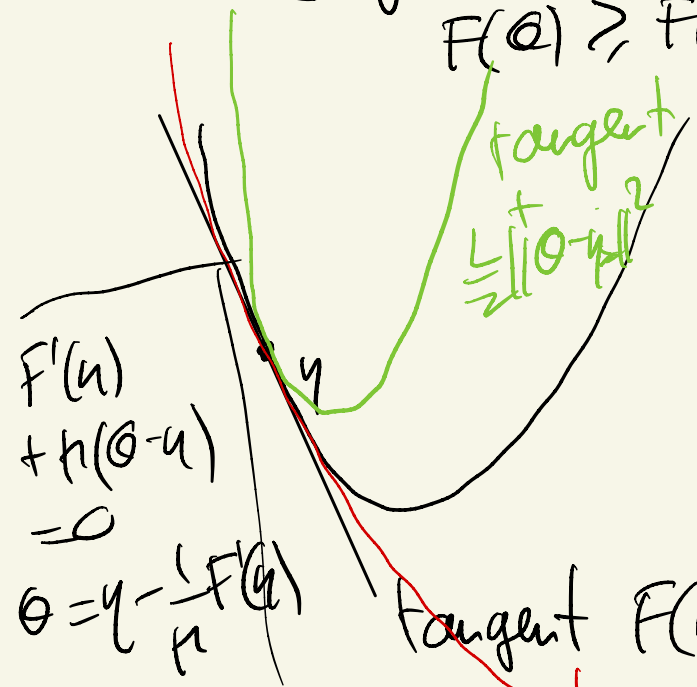
Only if F is twice differentiable

Lemmas: ① if F is smooth
$$F(\theta) \leq F(q) + F'(q)^T (\theta - q) + \boxed{\frac{L}{2} \|\theta - q\|^2}$$

② if F is strongly convex
$$F(\theta) \geq F(q) + F'(q)^T (\theta - q) + \boxed{\frac{\mu}{2} \|\theta - q\|^2}$$

proof:
Taylor expansion

tangent
$+ \frac{L}{2}\|\theta - q\|^2$

$$F(\theta) = F(q) + F'(q)^T (\theta - q) + \frac{1}{2}(\theta - q)^T F''(\theta)(\theta - q)$$

$F'(q)$
$+ h(\theta - q)$
$= 0$
$\theta = q - \frac{1}{h} F'(q)$

③ Losajevitch inequality
$$F(q) - F(q_*) \leq \frac{1}{2\mu} \|F'(q)\|^2$$

tangent $F(q) + F'(q)^T (\theta - q)$

tangent $+ \frac{\mu}{2}\|\theta - q\|^2$

proof: $F(q_*) \geq F(q) - \frac{1}{2\mu}\|F'(q)\|^2$

# Proof of convergence of gradient descent

$$\Theta_t = \Theta_{t-1} - \gamma F'(\Theta_{t-1})$$ <span style="color:red">(smoothness)</span>

$$F(\Theta_t) \leq F(\Theta_{t-1}) + F'(\Theta_{t-1})^T(\Theta_t - \Theta_{t-1}) + \frac{L}{2}\|\Theta_t - \Theta_{t-1}\|^2$$

$$\leq F(\Theta_{t-1}) + F'(\Theta_{t-1})^T(-\gamma F'(\Theta_{t-1})) + \frac{L}{2}\|\gamma F'(\Theta_{t-1})\|^2 \quad \text{<span style=\"color:red\">using GD iteration</span>}$$

$$= F(\Theta_{t-1}) - \left(\gamma - \frac{L}{2}\gamma^2\right)\|F'(\Theta_{t-1})\|^2 \qquad \text{<span style=\"color:red\">}\gamma = 1/L\text{</span>}$$

<span style="color:red">$\gamma/2$</span>

$$F(\Theta_t) \leq F(\Theta_{t-1}) - \frac{\gamma}{2}\underbrace{\|F'(\Theta_{t-1})\|^2}_{\color{red}{\geq 2\mu(F(\Theta_{t-1}) - F(q_\infty))}}$$

<span style="color:red">Losojewitz ineq.</span>

$$F(\Theta_t) - F(q_\infty) \leq \underbrace{(1 - \gamma\mu)}_{\color{red}{1 - \mu/L}}(F(\Theta_{t-1}) - F(q_\infty))$$

**other result**

$$\frac{L}{t}\|\Theta_c - q_\infty\|^2 \qquad \leq \left(1 - \frac{\mu}{L}\right)^t (F(\Theta_c) - F(q_*))$$
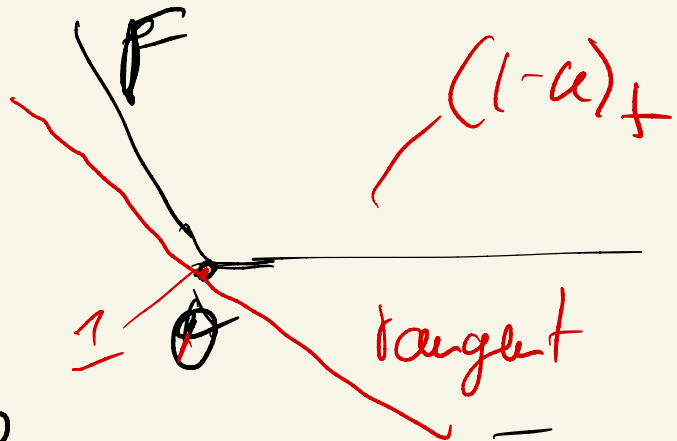
# Non-smooth optimization

F convex, not necessarily differentiable
gradient may not exist
use "subgradient" at $\Theta$ : any vector $z$
such that $F(\eta) \geq F(\Theta) + z^T(\eta - \Theta)$

Example : $|\Theta|$

at $\Theta > c$, gradient $= 1$
$\Theta < c$, gradient $= -1$
$\Theta = c$, subgradient : any element of $[-1, 1]$

Example from me
$$F(\Theta) = \frac{1}{m} \sum_{i=1}^{m} (1 - y_i \Theta^T \varphi(x_i))_+ \quad \text{SVM}$$

$$B = \max_i \| \varphi(x_i) \|$$

Assumption : F is convex and Lipschitz-continuous
$$\forall \Theta, \eta \quad |F(\Theta) - F(\eta)| \leq B \|\Theta - \eta\|_2$$

Proposition : $\forall s$ subgradient of F, $\|s\|_2 \leq B$

$(1-\alpha)_+$

tangent

Subgradient method : $\theta_t = \theta_{t-1} - \boxed{\gamma_t} F'(\theta_{t-1})$

<span style="color:red">dodging stepsize</span>

<span style="color:red">avg subgradient of $F'$ at $\theta_{t-1}$</span>
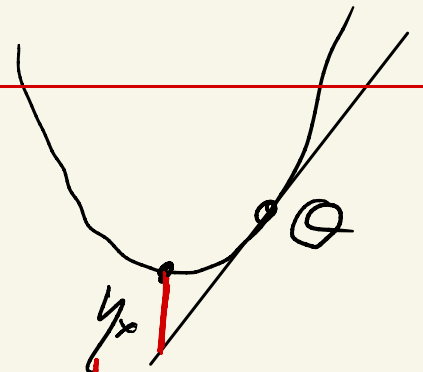
Convergence proof:

Assumption: $\eta_*$ is a global minimizer

$$\|\theta_t - \eta_*\|^2 = \|\theta_{t-1} - \eta_* - \gamma_t F'(\theta_{t-1})\|^2$$

$$= \|\theta_{t-1} - \eta_*\|^2 - 2\gamma_t \underbrace{F'(\theta_{t-1})^\top (\theta_{t-1} - \eta_*)}_{\color{red}{\geq F(\theta_{t-1}) - F(\theta_*)}} + \gamma_t^2 \underbrace{\|F'(\theta_{t-1})\|^2}_{\color{red}{\leq \gamma_t^2 B^2}}$$

$$\|\theta_t - \eta_*\|^2 \leq \|\theta_{t-1} - \eta_*\|^2 - 2\gamma_t \left[ F(\theta_{t-1}) - F(\eta_*) \right] + \gamma_t^2 B^2$$

<hr>

Lemma : $\forall \theta; \; F(\theta) - F(\eta_*) \leq F'(\theta)^\top (\theta - \eta_*)$

Proof: $F(\eta_*) \geq F(\theta) + F'(\theta)^\top (\eta_* - \theta)$



<span style="color:red">function above tangent</span>

$$\|\theta_t - \theta_*\|^2 \leq \|\theta_{t-1} - \theta_*\|^2 - 2\gamma_t\left(F(\theta_{t-1}) - F(q_*)\right) + \gamma_t^2 B^2$$

$$2\gamma_t\left(F(\theta_{t-1}) - F(\theta_*)\right) \leq \|\theta_{t-1} - q_*\|^2 - \|\theta_t - q_*\|^2 + \gamma_t^2 B^2$$

$$2\sum_{t=1}^{S} \gamma_t\left(F(\theta_{t-1}) - F(\theta_*)\right) \leq \|\theta_c - q_*\|^2 - \|\theta_S - q_*\|^2 + \sum_{t=1}^{S}\gamma_t^2 B^2$$

$\geq \min\limits_{t \in \{1, \ldots, S\}} F(\theta_{t-1})$

telescoping sum

$$\left(2\sum_{t=1}^{S}\gamma_t\right)\left[\min_t F(\theta_{t-1}) - F(q_*)\right] \leq$$

$$\min_{t \in \{1, \ldots, S\}} F(\theta_{t-1}) - F(q_*) \leq \frac{\|\theta_c - q_*\|^2}{2\sum_{t=1}^{S}\gamma_t} + \frac{\sum_{t=1}^{S}\gamma_t^2}{\sum_{t=1}^{S}\gamma_t} B^2$$

$\gamma_t = \gamma, \forall t$

$\gamma_t = 1/t^\alpha \rightarrow \alpha = 1 \quad \sum\limits_{t=1}^{S}\frac{1}{t} \sim \log(S)$

$\frac{\cancel{B}}{S\gamma} \quad + \quad \gamma B^2$

$\alpha = 1/2$

$$\text{if } \gamma_t = \frac{\gamma}{\sqrt{t}}; \quad \sum_{t=1}^{S} \gamma_t = \gamma \sum_{t=1}^{S} \frac{1}{\sqrt{t}} \geq \gamma S / \sqrt{S} = \gamma\sqrt{S}$$

<span style="color:red">$$\geq \frac{1}{\sqrt{S}}$$</span>

$$\sum_{t=1}^{S} \gamma_t^2 = \gamma^2 \sum_{t=1}^{S} \frac{1}{t}$$

$$\leq \gamma^2 (1 + \log S)$$

<span style="color:red">Comparison with integral</span>

<span style="color:red">$$\sum_{t=1}^{S} \frac{1}{t} \leq 1 + \sum_{t=2}^{S} \frac{1}{t} \leq 1 + \int_{1}^{S-1} \frac{du}{u}$$</span>

<span style="color:red">$$\leq 1 + \log(S-1)$$</span>

$$\min_{t \in \{1, \ldots, S\}} F(\theta_{t-1}) - F(\eta_\delta) \leq \frac{1}{2} \frac{\|\theta_0 - \eta_\delta\|^2}{\gamma\sqrt{S}} + \frac{1}{2} B^2 \gamma \frac{1 + \log S}{\sqrt{S}}$$

<span style="color:green">$$\gamma_t = 1/t^\alpha$$</span>

<span style="color:green">$$\sum \gamma_t = S^{1-\alpha}$$</span>

<span style="color:green">$$\sum \gamma_A^2 \sim$$</span>

<span style="color:red">$$\longrightarrow 0 \quad \text{as } S \longrightarrow +\infty$$</span>

<span style="color:red">$$\text{as} \quad \frac{\log S}{\sqrt{S}}$$</span>

# Stochastic gradient descent

**Version 1**: Minimize $\frac{1}{n} \sum_{i=1}^{n} F_i(\theta)$    $= \mathbb{E}\, G(\theta, \mathfrak{z})$ using empirical dist.

"finite sum"

ex: empirical risk

Algo:
$$\theta_t = \theta_{t-1} - \gamma_t \, F'_{i(t)}(\theta_{t-1})$$

$i(t)$: index chosen at random in $\{1, \dots, n\}$ independently

will converge to
$$\eta_\star = \arg\min \frac{1}{n} \sum F_i$$

---

**Version 2**: $F(\theta) = \mathbb{E}_{\mathfrak{z}}\, G(\theta, \mathfrak{z})$   "expectation"

ex: expected risk

Algo: $\theta_t = \theta_{t-1} - \gamma_t \, G'(\theta_{t-1}, \mathfrak{z}_t)$

$\boxed{t = M}$

will converge to $\theta_\star = \arg\min F$

independent observation

$$\theta_t = \theta_{t-1} - \gamma_A G'(\theta_{t-1}, \mathfrak{z}_t)$$

$$\|\theta_t - \theta_*\|^2 = \|\theta_{t-1} - \theta_*\|^2 - 2\gamma_A (\theta_{t-1} - \theta_*)^T G'(\theta_{t-1}, \mathfrak{z}_t)$$

$$+ \gamma_t^2 \|G'(\theta_{t-1}, \mathfrak{z}_t)\|^2$$

$$\leq \gamma_A^2 B^2$$

$$E(\quad | \theta_{t-1})$$

$$E(\|\theta_t - \theta_*\|^2 | \theta_{t-1}) = \|\theta_{t-1} - \theta_*\|^2 + \gamma_t^2 B^2$$

$$- 2\gamma_A (\theta_{t-1} - \theta_*)^T F(\theta_{t-1})$$

$$\geq F(\theta_{t-1}) - F(\theta_*)$$

$$E \|\theta_t - \theta_*\|^2 \leq E \|\theta_{t-1} - \theta_*\|^2 + \gamma_t^2 B^2$$

$$- 2\gamma_A \left( E(F(\theta_{t-1})) - F(\theta_*) \right)$$

$$\sum_{t=1}^{S} \gamma_t E\left( F(\theta_t) \right) - F(\theta_*) \leq \underline{\quad\quad}$$

Jensen's inequality

$$\sum_{t=1}^{S} \gamma_t F(\theta_t) \geq \left( \sum_{t=1}^{S} \gamma_t \right) F\left( \frac{\sum_{t=1}^{S} \gamma_t \theta_t}{\sum_{t=1}^{S} \gamma_t} \right) \quad \overline{\theta_S}$$

$$\boxed{\mathbb{E} \, F(\theta) \geq F(\mathbb{E}\theta)}$$

min $F(\theta_{t-1})$

$t \in \{1, S\}$

$$\gamma_t = \frac{\gamma}{\sqrt{t}}$$

$$\mathbb{E}\left( F(\overline{\theta_S}) \right) F(\theta_*) \leq \frac{1}{2\gamma\sqrt{S}} \|\theta_0 - \theta_*\|^2 + \frac{\gamma B^2 (1 + \log S)}{2\sqrt{S}}$$

$$S = n \qquad \frac{1}{\sqrt{n}} \qquad \boxed{\text{Complexity } O(n)}$$