**Remarks** | Contraction principle:

1) For $\phi: \mathbb{R} \to \mathbb{R}$, $C_\phi$-Lip, $\mathbb{E}\left[\sup_{h \in \mathcal{H}} \frac{1}{n}\sum \sigma_i (\phi \circ h)(z_i) \mid z_1^n\right] \leq C_\phi \, \mathcal{R}_n \mathcal{H}$

2) For $\phi: \mathbb{R} \to \mathbb{R}$, $C_\phi$-Lip & $\phi(0)=0$,
$$\mathbb{E}\left[\sup_{h \in \mathcal{H}} \left|\frac{1}{n}\sum \sigma_i (\phi \circ h)(z_i)\right| \mid z_1^n\right] \leq 2C_\phi \, \mathbb{E}\left[\sup_{h \in \mathcal{H}} \left|\frac{1}{n}\sum \sigma_i \, h(z_i)\right| \mid z_1^n\right]$$

**Def** | A collection of zero mean RVs $\{V_h : h \in T\}$ is a sub-Gaussian process w.r.t. $d$ if
$$\mathbb{E}\, e^{\lambda(V_h - V_{h'})} \leq \exp\left(\frac{\lambda^2}{2} d(h,h')^2\right) \qquad \forall h,h' \in T, \quad \forall \lambda \in \mathbb{R}.$$
↳ tail of $V_h - V_{h'}$ is $d(h,h')^2$-sub G.

**Key Lemma** | Let $X_j$ be $\sigma_i^2$-sub-G RVs, $j=1,\cdots,N$. Then, $\mathbb{E}\max_{1 \leq j \leq N}|X_j| \leq \max_{1 \leq j \leq N}\sigma_j \cdot 2\sqrt{\log N}$, $N \geq 2$.

**Theorem** | (Dudley's entropy integral)  Let $\{V_h : h \in T\}$ be a sub-Gaussian process w.r.t. $d$ on $T$. For any $\delta \in [0, D]$,
$$\mathbb{E}\sup_{h \in T} V_h \leq \mathbb{E}\left[\sup_{h,h' \in T} V_h - V_{h'}\right] \leq 2\,\mathbb{E}\left[\sup_{\substack{r,r' \in T \\ d(r,r') \leq \delta}} V_r - V_{r'}\right] + 32 \int_{\delta/4}^{D} \sqrt{\log N(T,d,\varepsilon)}\; d\varepsilon$$

**Pf)** Let $N = N(T,d,\delta)$, and $\mathcal{U} := \{h_j\}_{j=1}^N$ be a $\delta$-cover of $T$. Fix an arbitrary $h \in T$. There exists $j$ s.t. $d(h,h_j) \leq \delta$. Then,
$$V_h - V_{h_1} = V_h - V_{h_j} + V_{h_j} - V_{h_1} \leq \sup_{\substack{r,r' \in T \\ d(r,r') \leq \delta}} (V_r - V_{r'}) + \max_{1 \leq j \leq N}|V_{h_j} - V_{h_1}|$$

Given another arbitrary $\tilde{h} \in T$, the same bound holds for $V_{h_1} - V_{\tilde{h}}$. Adding the two, and taking supremum over $h,\tilde{h} \in T$
$$\sup_{h,\tilde{h} \in T} V_h - V_{\tilde{h}} \leq 2 \sup_{\substack{r,r' \in T \\ d(r,r') \leq \delta}} (V_r - V_{r'}) + 2\max_{1 \leq j \leq N}|V_{h_j} - V_{h_1}|$$
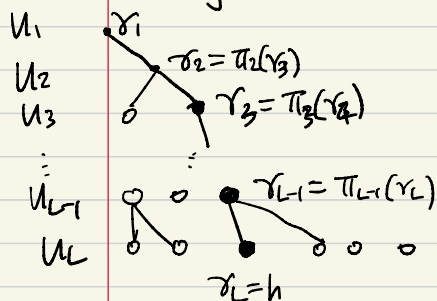$$\leq 2 \sup_{\substack{r,r' \in T \\ d(r,r') \leq \delta}} (V_r - V_{r'}) + 2 \sup_{h,\tilde{h} \in \mathcal{U}}|V_h - V_{\tilde{h}}|$$

Instead of bounding the last term via Lemma, we use a chaining argument. For $L$ s.t. $D2^{-L} \leq \delta$, think of $\mathcal{U}_L = \mathcal{U}$ as a $(D2^{-L})$-cover of $\mathcal{U}$. Now, for each $m = 1, \cdots, L-1$, define $\mathcal{U}_m :=$ minimal $(D \cdot 2^{-m})$-cover of $\mathcal{U}_{m+1}$ (we allow elements of $T$).

By def, $|\mathcal{U}_m| \leq N(T, d, D \cdot 2^{-m})$.

Best approx of $h \in \mathcal{U}$ in $\mathcal{U}_m$.

For each $m = 1,\cdots, L-1$, define $\pi_m : \mathcal{U}_{m+1} \to \mathcal{U}_m$, $\pi_m(h) = \operatorname*{argmin}_{h' \in \mathcal{U}_m} d(h, h')$. Using this, we construct a chain from any $h \in \mathcal{U}$. $\gamma_{m-1} = \pi_{m-1}(\gamma_m)$



$\mathcal{U}_1 \quad \gamma_1$
$\mathcal{U}_2 \quad \gamma_2 = \pi_2(\gamma_3)$
$\mathcal{U}_3 \quad \gamma_3 = \pi_3(\gamma_4)$
$\vdots$
$\mathcal{U}_{L-1} \quad \gamma_{L-1} = \pi_{L-1}(\gamma_L)$
$\mathcal{U}_L$
$\gamma_L = h$

$$V_h - V_{\gamma_1} = \sum_{m=2}^{L} V_{\gamma_m} - V_{\gamma_{m-1}} \qquad \text{and}$$
$$\mathbb{E}|V_h - V_{\gamma_1}| \leq \sum_{m=2}^{L} \sup_{r \in \mathcal{U}_m}|V_r - V_{\pi_{m-1}(r)}| \qquad //$$

Similarly, for any other $\tilde{h} \in \mathcal{U}$, we have same bound with $\tilde{\gamma}_m$'s.

We arrive at $\quad |V_h - V_{\tilde{h}}| = |V_h - V_{r_1} + V_{r_1} - V_{\tilde{r}_1} + V_{\tilde{r}_1} - V_{\tilde{h}}|$

$$\leq |V_{r_1} - V_{\tilde{r}_1}| + |V_h - V_{r_1}| + |V_{\tilde{r}_1} - V_{\tilde{h}}|$$

$$\leq \max_{r_1, \tilde{r}_1 \in U_1} |V_{r_1} - V_{\tilde{r}_1}| + 2 \sum_{m=2}^{L} \max_{r \in U_m} |V_r - V_{\pi_{m-1}(r)}|$$

Now note that $\sup_{r_0, r_1 \in U_1} d(r_0, r_1) \leq D$, and $V_{r_1} - V_{\tilde{r}_1}$ is $d(r_1, \tilde{r}_1)^2$-subGaussian.

and $\max_{r \in U_m} d(r, \pi_{m-1}(r)) \leq D \cdot 2^{-(m+1)}$, and $|U_m| \leq N(T, d, D \cdot 2^{-m})$.

From Lemma, $\quad \mathbb{E} \max_{r_1, \tilde{r}_1 \in U_1} |V_{r_1} - V_{\tilde{r}_1}| \leq 2D \sqrt{\log N(T, d, \frac{D}{2})}$, and

$$\mathbb{E} \max_{r \in U_m} |V_r - V_{\pi_{m-1}(r)}| \leq 2D 2^{-(m-1)} \sqrt{\log N(T, d, D \cdot 2^{-m})}$$

Conclude that
$$\mathbb{E} \sup_{h, \tilde{h} \in U} |V_h - V_{\tilde{h}}| \leq 4 \sum_{m=1}^{L} D \cdot 2^{-(m-1)} \sqrt{\log(T, d, D 2^{-m})}$$

Since $\delta \mapsto \log N(T, d, \delta)$ is dec, $\quad D \cdot 2^{-m} \sqrt{\log N(T, d, D \cdot 2^{-m})} \leq 2 \int_{D 2^{-(m+1)}}^{D 2^{-m}} \sqrt{\log N(T, d, \varepsilon)} \, d\varepsilon$

$\Rightarrow \quad 2 \mathbb{E} \sup_{h, \tilde{h} \in T} |V_h - V_{\tilde{h}}| \leq 32 \int_{\delta/4}^{D} \sqrt{\log N(T, d, \varepsilon)} \, d\varepsilon$

Combining with $*$, we get the result. $\quad\quad\quad\quad \boxtimes$.

<u>Rmk</u> Measurability , asymptotic versions $\quad \sqrt{n} \left( \frac{1}{n} \sum \ell(\theta; z_i) - \mathbb{E}\ell(\theta; z) \right) \Rightarrow G(h)$ unif in $h$

# Asymptotics

We use the following notation.

**Def**  RVs  $X_n = O_p(1)$  if  $\sup_{n \geq 1} \mathbb{P}(\|X_n\| \geq M) \to 0$  as  $M \to \infty$

$X_n = o_p(1)$  if  $\mathbb{P}(\|X_n\| \geq M) \to 0$  $\forall M > 0$.

We write  $X_n = O_p(n)$  if  $n^{-1} X_n = O_p(1)$.

# ULLN

We want to show that  $\hat{\theta}_n \to \theta^*$, where  $\theta^* = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}\, \ell(\theta; z)$.

We use  uniform law of large numbers  $\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum \ell(\theta; z_i) - \mathbb{E}\ell(\theta; z) \right| \xrightarrow{p} 0$  to prove this.

To simplify notation,  let  $R(\theta) := \mathbb{E}\,\ell(\theta; z)$,  $\hat{R}_n(\theta) := \frac{1}{n} \sum \ell(\theta; z_i)$.

**Prop**  If ULLN holds, and  $\hat{\theta}_n$ is s.t.  $\hat{R}_n(\hat{\theta}_n) \leq \inf_{\theta \in \Theta} \hat{R}_n(\theta) + o_p(1)$,  $R(\hat{\theta}_n) \xrightarrow{p} \inf_{\theta \in \Theta} R(\theta)$.

**Pf)**  Let  $\theta^* \in \operatorname{argmin}_{\theta \in \Theta} R(\theta)$.  $R(\hat{\theta}_n) - R(\theta^*) = R(\hat{\theta}_n) - \hat{R}_n(\hat{\theta}_n) + \hat{R}_n(\hat{\theta}_n) - \hat{R}_n(\theta^*) + \hat{R}_n(\theta^*) - R(\theta^*)$

$\underbrace{\qquad}_{\text{by hypothesis}} \leq \sup_{\theta \in \Theta} |\hat{R}_n(\theta) - R(\theta)| + o_p(1) + o_p(1) \underbrace{\qquad}_{\text{by WLLN}} = o_p(1)$.  $\boxtimes$

**Cor**  Let  $R(\cdot)$  be s.t.  $\forall \varepsilon > 0 \; \exists \delta > 0$ s.t.  $R(\theta) \geq R(\theta^*) + \delta$  whenever  $\|\theta - \theta^*\| \geq \varepsilon$.
Under conditions of **Prop**,  $\hat{\theta}_n \xrightarrow{p} \theta^*$.

**Pf)**  $\mathbb{P}(\|\hat{\theta}_n - \theta\| > \varepsilon) \leq \mathbb{P}(R(\hat{\theta}_n) - R(\theta^*) \geq \delta) \to 0$  by **Prop**.  $\boxtimes$.

**Theorem**  Let  $H$  be an envelope function for  $\mathcal{H}$:  $\forall h \in \mathcal{H}$,  $|h| \leq H$.  Let  $\mathbb{E}|H(z)| < \infty$, and define truncated version of  $\mathcal{H}$:  $\mathcal{H}_M := \{ \underbrace{h \, \mathbb{1}\{|h| \leq M\}}_{=: h_M} : h \in H \}$.

If  $n^{-1} \log N(\mathcal{H}_M, \|\cdot\|_{L_1(\hat{P}_n)}, \varepsilon) \xrightarrow{p} 0$  for all fixed  $\varepsilon > 0$, $M < \infty$,
then  $\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum h(z_i) - \mathbb{E} h(z) \right| \xrightarrow{p} 0$.

**Pf)**  From symmetrization,  $\mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum h(z_i) - \mathbb{E}h(z) \leq 2 \mathbb{E} \mathcal{R}_n \mathcal{H}$

$\leq 2 \mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum \sigma_i \left( h(z_i) - h_M(z_i) \right) + 2 \mathbb{E} \mathcal{R}_n \mathcal{H}_M$

$\leq 2 \mathbb{E}\, H(z) \mathbb{1}\{H(z) > M\} + 2 \mathbb{E} \mathcal{R}_n \mathcal{H}_M$

Take a  $\varepsilon$-cover  $\mathcal{H}_{M, \varepsilon}$  of  $\mathcal{H}_M$  in  $\|\cdot\|_{L_1(\hat{P}_n)}$.  $\mathcal{R}_n \mathcal{H}_M \leq \mathcal{R}_n \mathcal{H}_{M, \varepsilon} + \varepsilon$
Now, note that since  $\sup_{h \in \mathcal{H}} \|h\|_{L_2(\hat{P}_n)} \leq M$,  Lemma gives

$\sqrt{n}\, \mathcal{R}_n \mathcal{H}_{M, \varepsilon} \leq 2M \sqrt{\log N(\mathcal{H}_M, \|\cdot\|_{L_1(\hat{P}_n)}, \varepsilon)}$.  $\Rightarrow$  $\mathcal{R}_n \mathcal{H}_{M, \varepsilon} \xrightarrow{p} 0$

Same bound holds for  $\mathcal{R}_n(\mathcal{H}_{M, \varepsilon})$.

So  $\mathbb{E} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum h(z_i) - \mathbb{E}h(z) \right| \leq 4 \mathbb{E}\, H(z) \mathbb{1}\{H(z) > M\} + 4 \mathbb{E} \mathcal{R}_n \mathcal{H}_{M, \varepsilon} + \varepsilon$

$\equiv \limsup_{n \to \infty} \mathbb{E} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum h(z_i) - \mathbb{E}h(z) \right| \leq 4 \mathbb{E}\, H(z) \mathbb{1}\{H(z) > M\} + \varepsilon$  $\forall \varepsilon > 0, M < \infty$.
Taking  $M \to \infty$, $\varepsilon \downarrow 0$,  we get the result.  $\boxtimes$.

# Rates of convergence

We now characterize the rate of convergence for $\hat{\theta}_n \to \theta^*$.

**Intuition**



/ $R(\theta)$    If curvature of $R$ is higher than perturbations $\hat{R}_n - R$, then we're good.

$$\hat{R}_n(\theta) - \hat{R}_n(\theta^*) = \underbrace{\hat{R}_n(\theta) - R(\theta) - (\hat{R}_n(\theta^*) - R(\theta^*))}_{=: \Delta_n(\theta) \text{ Fluctuation}} + \underbrace{R(\theta) - R(\theta^*)}_{\text{Growth}}$$

We call $\Delta_n$ the localized process.

**Def**  The modulus of continuity around $\theta^*$ is $W_n(\delta) := \sup_{\|\theta - \theta^*\| \leq \delta} |\Delta_n(\theta)|$.

**Theorem**  Let $\hat{\theta}_n \in \arg\min_{\theta \in \Theta} \hat{R}_n(\theta)$, and assume $\hat{\theta}_n \xrightarrow{p} \theta^*$.   _We assume $W_n$ is small compared to curvature._

**Fluctuation**
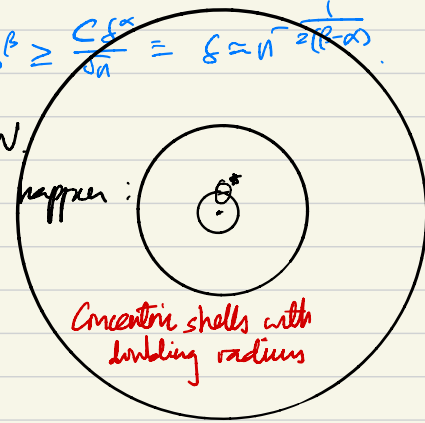$$\mathbb{E} W_n(\delta) \leq \frac{C}{\sqrt{n}} \delta^\alpha \quad \text{for some } M < \infty, \alpha > 0.$$

**Growth**
$$\exists \beta \geq 1, \lambda > 0, \varepsilon > 0 \quad \text{s.t.} \quad R(\theta) \geq R(\theta^*) + \lambda \|\theta - \theta^*\|^\beta \quad \forall \theta \text{ s.t. } \|\theta - \theta^*\| \leq \varepsilon.$$

Then, $\|\hat{\theta}_n - \theta^*\| = O_p(n^{-\frac{1}{2(\beta - \alpha)}})$ if $\beta > \alpha$    _Intuition:_ $\lambda \delta^\beta \geq \frac{C}{\sqrt{n}} \delta^\alpha \Rightarrow \delta \simeq n^{-\frac{1}{2(\beta - \alpha)}}$.



_Concentric shells with doubling radius_

**Pf)**  We use a peeling argument. Let $r_n := n^{\frac{1}{2(\beta - \alpha)}}$, and fix $M \in \mathbb{N}$.
If $r_n \|\hat{\theta}_n - \theta^*\| \geq 2^M$ then at least one of the following must happen:
- $\|\hat{\theta}_n - \theta^*\| > \varepsilon$ so growth condition doesn't apply
- $\exists j$ s.t. $2^{j-1} < r_n \|\hat{\theta}_n - \theta^*\| \leq 2^j$    cf. $2^j \leq r_n \varepsilon$

  $\hookrightarrow$ In this case, $\Delta_n(\hat{\theta}_n) = \underbrace{\hat{R}_n(\hat{\theta}_n) - \hat{R}_n(\theta^*)}_{\leq 0} - \underbrace{(R(\hat{\theta}_n) - R(\theta^*))}_{\geq 0}$    satisfies

$$W_n(r_n^{-1} 2^j) \geq |\Delta_n(\hat{\theta}_n)| \geq R(\hat{\theta}_n) - R(\theta^*) \geq \lambda \|\hat{\theta}_n - \theta^*\|^\beta \geq \lambda \cdot (r_n^{-1} \cdot 2^{j-1})^\beta.$$

So taking a union bound gives

$$\mathbb{P}(r_n \|\hat{\theta}_n - \theta^*\| \geq 2^M) \leq \sum_{j \geq M, 2^j \leq r_n \varepsilon} \mathbb{P}(r_n \|\hat{\theta}_n - \theta^*\| \in [2^{j-1}, 2^j]) + \mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq \varepsilon).$$

$$\leq \sum_{j \geq M, 2^j \leq r_n \varepsilon} \mathbb{P}(W_n(r_n^{-1} 2^j) \geq \lambda \cdot (r_n^{-1} 2^{j-1})^\beta) + \mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq \varepsilon)$$

$$\leq \sum_{j \geq M, 2^j \leq r_n \varepsilon} \frac{1}{\lambda (r_n^{-1} 2^{j-1})^\beta} \mathbb{E} W_n(r_n^{-1} 2^j) + \mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq \varepsilon)$$

$$\leq \frac{C}{\lambda \sqrt{n}} 2^\beta \cdot r_n^{\beta - \alpha} \sum_{j \geq M, 2^j \leq r_n \varepsilon} \frac{1}{2^{j \cdot (\beta - \alpha)}} + \mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq \varepsilon)$$

$$= \underbrace{\frac{C}{\lambda} \cdot 2^\beta \sum_{j \geq M} \frac{1}{2^{j(\beta - \alpha)}}}_{\to 0 \text{ as } M \to \infty} + \underbrace{\mathbb{P}(\|\hat{\theta}_n - \theta^*\| \geq \varepsilon)}_{\to 0 \text{ as } n \to \infty}$$

We've shown that $r_n \|\hat{\theta}_n - \theta^*\| = O_p(1)$.    ∎

**Example**  Let $\theta \mapsto \ell(\theta; z)$ be $C^3$ $\forall z$, $\theta \mapsto \ell(\theta; z)$ is $L(z)$-Lipschitz $\forall z$, with $\mathbb{E} L(z)^2 < \infty$.
Assume $\nabla^2 R(\theta^*) \succ 0$. From Taylor expansion,

$$R(\theta) = R(\theta^*) + \nabla R(\theta^*)^T (\theta - \theta^*) + \tfrac{1}{2} (\theta - \theta^*)^T \nabla^2 R(\theta^*) (\theta - \theta^*) + O(\|\theta - \theta^*\|^3)$$

$$\geq R(\theta^*) + \tfrac{1}{2} (\theta - \theta^*)^T \nabla^2 R(\theta^*) (\theta - \theta^*) + O(\|\theta - \theta^*\|^3)$$

$$\geq R(\theta^*) + \tfrac{1}{4} \lambda_{min}(\nabla^2 R(\theta^*)) \|\theta - \theta^*\|^2.    \quad \text{for } \|\theta - \theta^*\| \text{ small enough.}$$

So $\theta \mapsto R(\theta)$ satisfies growth condition with $\beta = 2$ and $\lambda = \frac{1}{4}\lambda_{min}(\nabla^2 R(\theta^*))$.

To show fluctuation, we use Dudley's entropy integral. From symmetrization,

$$\mathbb{E}[W_n(\delta) \mid Z_1^n] \leq 2\,\mathbb{E}\left[\sup_{\|\theta - \theta^*\| \leq \delta} \left| \frac{1}{n}\sum (\ell(\theta; z_i) - \ell(\theta^*; z_i))\sigma_i \right| \mid Z_1^n\right]$$

Recall that $\varepsilon$-covering number of $\{z \mapsto \ell(\theta; z) - \ell(\theta^*; z) : \|\theta - \theta^*\| \leq \delta\}$ is held by

$$N\left(\{\theta : \|\theta - \theta^*\| \leq \delta\}, \|\cdot\|, \frac{\varepsilon}{\|L\|_{L_2(\hat{P}_n)}}\right) \leq \left(1 + \frac{\delta \cdot \|L\|_{L_2(\hat{P}_n)}}{\varepsilon}\right)^d.$$

$$\mathbb{E}[W_n(\delta) \mid Z_1^n] \leq \frac{1}{\sqrt{n}} \int_0^{\delta \|L\|_{L_2(\hat{P}_n)}} \sqrt{d\,\log\left(1 + \frac{\delta\|L\|_{L_2(\hat{P}_n)}}{\varepsilon}\right)}\,d\varepsilon$$

$$\leq \sqrt{\frac{d}{n}} \cdot \delta \cdot \|L\|_{L_2(\hat{P}_n)}$$

Noting that $\mathbb{E}\|L\|_{L_2(\hat{P}_n)} \leq \sqrt{\frac{1}{n}\sum \mathbb{E}L(z_i)^2} = \sqrt{\mathbb{E}L(z)^2}$,  $\mathbb{E}W_n(\delta) \leq \sqrt{\frac{d}{n}}\,\delta \cdot \sqrt{\mathbb{E}L(z)^2}$.

$(\alpha = 1)$

We conclude that $\sqrt{n}\,\|\hat{\theta}_n - \theta^*\| = O_p(1)$.

# SGD

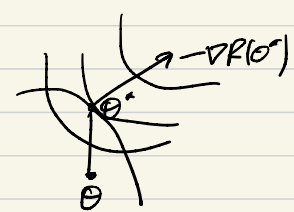**Def**  A function $R: \mathbb{R}^d \to \mathbb{R}$ is convex if $\forall \theta, \theta' \in \mathbb{R}^d$, $R(t\theta + (1-t)\theta') \leq tR(\theta) + (1-t)R(\theta')$ $\forall t \in [0,1]$.

**Lemma**  Let $R: \mathbb{R}^d \to \mathbb{R}$ be differentiable on the interior of its domain. $R$ is convex iff
$\forall \theta, \theta' \in \mathbb{R}^d$   $R(\theta') \geq R(\theta) + \nabla R(\theta)^T (\theta' - \theta)$. ← 1st order approx is a global minorization

**Pf)**  '⇒' From def of convexity, $R(\theta + t(\theta' - \theta)) \leq R(\theta) + t(R(\theta') - R(\theta))$. $\equiv$ $R(\theta') - R(\theta) \geq \frac{1}{t}(R(\theta + t(\theta'-\theta)) - R(\theta))$. Send $t \downarrow 0$.
'⇐' Define $\theta_t = t\theta + (1-t)\theta'$. Combining $R(\theta) \geq R(\theta_t) + \nabla R(\theta_t)^T(\theta - \theta_t)$, $R(\theta') \geq R(\theta_t) + \nabla R(\theta_t)^T(\theta' - \theta_t)$,
$tR(\theta) + (1-t)R(\theta') \geq R(\theta_t) + \nabla R(\theta_t)^T(t\theta + (1-t)\theta' - \theta_t)$   $\forall t \in [0,1]$.  ☒

**Rmk**  The latter def of convexity motivates generalization of gradients to nonsmooth, convex functions.

## Optimality

Consider $\min_{\theta \in \Theta} R(\theta)$, for $R: \mathbb{R}^d \to \mathbb{R}$ diff; convex.



**Lemma**  $\theta^* = \arg\min_{\theta \in \Theta} R(\theta)$   iff   $\nabla R(\theta^*)^T(\theta - \theta^*) \geq 0$ $\forall \theta \in \Theta$

**Pf)**  '⇐' From prev lemma, $R(\theta) - R(\theta^*) \geq \nabla R(\theta^*)^T(\theta - \theta^*) \geq 0$ $\forall \theta \in \Theta$.
'⇒' $\nabla R(\theta^*)^T(\theta - \theta^*) = \lim_{t \downarrow 0} \frac{1}{t}(R(\theta^* + t(\theta - \theta^*)) - R(\theta^*)) \geq 0$ $\forall \theta \in \Theta$.  ☒

**Cor**  Let $\Theta$ be a closed convex set in $\mathbb{R}^d$. Define the projection operator $\Pi_\Theta(\theta) := \arg\min_{\theta \in \Theta} \|\theta - \theta'\|_2$.
Then, $\|\Pi_\Theta(\theta) - \theta'\|_2 \leq \|\theta - \theta'\|_2$ $\forall \theta' \in \Theta$ $\forall \theta \in \mathbb{R}^d$.

**Pf)**  From first order conditions for $\min_{\theta \in \Theta} \|\theta - \theta'\|_2^2$,
$0 \leq (\Pi_\Theta(\theta) - \theta)^T(\theta' - \Pi_\Theta(\theta)) = (\Pi_\Theta(\theta) - \theta' + \theta' - \theta)^T(\theta' - \Pi_\Theta(\theta)) = -\|\theta' - \Pi_\Theta(\theta)\|_2^2 + (\theta' - \theta)^T(\theta' - \Pi_\Theta(\theta))$.
From Cauchy-Schwarz, $\|\theta' - \Pi_\Theta(\theta)\|_2^2 \leq \|\theta - \theta'\| \|\theta' - \Pi_\Theta(\theta)\|_2$. $\forall \theta' \in \Theta$  ☒.

## Stochastic gradients

A stochastic gradient $G(\theta)$ is a RV s.t. $\mathbb{E}G(\theta) = \nabla R(\theta)$.
We study first-order optimization methods based on stoch. gradients.

**(Canonical Problem)**  $\text{minimize}_{\theta \in \Theta} \{\mathbb{E}\, \ell(\theta; Z) =: R(\theta)\}$
If $\theta \mapsto \ell(\theta; z)$ is differentiable, then $\nabla_\theta \ell(\theta; z)$ is a stochastic gradient if $\mathbb{E}, \nabla_\theta$ can be interchanged.

- **SGD Idea:** Go in the direction of stoch. gradient, then project to $\Theta$.
- **Algo:** Let $G_k(\theta)$ be a stoch. gradient of $R(\theta)$.
At each iteration $k$, $\theta_{k+1} = \Pi_\Theta(\theta_k - \alpha_k G_k(\theta_k))$ for some stepsize $\alpha_k > 0$.

We're implicitly assuming that projections are efficient to compute.

**Rmk**  We can't even evaluate $\mathbb{E}\ell(\theta; Z)$. So SGD takes samples. In its simplest form, draw $Z_k \sim \mathbb{P}$, then take $G(\theta_k) := \nabla_\theta \ell(\theta_k; Z_k)$. We could take multiple samples and average over them.

**Rmk 2**  We could consider ERM $\min_{\theta \in \Theta} \frac{1}{n}\sum \ell(\theta; Z_i)$, and think of $\nabla_\theta \ell(\theta; Z_i)$ as a stoch. gradient of the empirical loss. Our following convergence results still apply in this case. The rationale for SGD w.r.t. empirical loss is purely computational: instead of incurring $O(n)$ to evaluate each gradient, I want to compute an approximate gradient in $O(1)$.

Convergence     Assume   $\theta^* \in \arg\min_{\theta \in \Theta} R(\theta) > -\infty$   exists.

Theorem   Let $\Theta$ (+) be compact. Assume $\exists R > 0$ s.t. $\sup_{\theta \in \Theta} \|\theta - \theta^*\|_2 \le D$, $\exists M > 0$ s.t. $\mathbb{E}\|G(\theta)\|_2^2 \le M^2 \; \forall \theta \in \Theta$.
Let $\alpha_k$ be dec. pos. step sizes, and $\bar{\theta}_k = \frac{1}{k}\sum^k \theta_k$  Then,

$$\mathbb{E}[R(\bar{\theta}_k) - R(\theta^*)] \le \frac{D^2}{2k\alpha_k} + \frac{1}{2k}\sum_1^k \alpha_k M^2.$$

Pf)  We expand on the error $\|\theta_{k+1} - \theta^*\|_2^2$.

$\frac{1}{2}\|\theta_{k+1} - \theta^*\|_2^2 = \frac{1}{2}\|\Pi_\Theta(\theta_k - \alpha_k G(\theta_k)) - \theta^*\|_2^2$

$\qquad \le \frac{1}{2}\|\theta_k - \alpha_k G(\theta_k) - \theta^*\|_2^2$   by non-expansiveness of $\Pi_\Theta$

$\qquad = \frac{1}{2}\|\theta_k - \theta^*\|_2^2 - \alpha_k \langle G(\theta_k), \theta_k - \theta^*\rangle + \frac{\alpha_k^2}{2}\|G(\theta_k)\|_2^2$.

Add & subtract $\alpha_k \langle \nabla R(\theta_k), \theta_k - \theta^*\rangle$ to get

$\qquad = \frac{1}{2}\|\theta_k - \theta^*\|_2^2 - \alpha_k \langle \nabla R(\theta_k), \theta_k - \theta^*\rangle + \frac{\alpha_k^2}{2}\|G(\theta_k)\|_2^2 - \alpha_k \langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^*\rangle$

$\qquad \le \frac{1}{2}\|\theta_k - \theta^*\|_2^2 - \alpha_k(R(\theta_k) - R(\theta^*)) + \frac{\alpha_k^2}{2}\|G(\theta_k)\|_2^2 - \alpha_k\langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^*\rangle$   by convexity

Divide each side by $\alpha_k$, and rearrange

$\qquad R(\theta_k) - R(\theta^*) \le \frac{1}{2\alpha_k}\left(\|\theta_k - \theta^*\|_2^2 - \|\theta_{k+1} - \theta^*\|_2^2\right) + \frac{\alpha_k}{2}\|G(\theta_k)\|_2^2 - \langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^*\rangle$   $\cdots$ (*)

Now, note that $\sum_{h=1}^k \frac{1}{2\alpha_k}\left(\|\theta_k - \theta^*\|_2^2 - \|\theta_{k+1} - \theta^*\|_2^2\right) = \frac{1}{2\alpha_1}\|\theta_1 - \theta^*\|_2^2 - \frac{1}{2\alpha_k}\|\theta_k - \theta^*\|_2^2 + \sum_{k=2}^k \left(\frac{1}{2\alpha_k} - \frac{1}{2\alpha_{k-1}}\right)\|\theta_k - \theta^*\|_2^2$

$\qquad\qquad \le \frac{D^2}{2\alpha_1} + \frac{D^2}{2}\sum_{k=2}^k \left(\frac{1}{\alpha_k} - \frac{1}{\alpha_{k-1}}\right) = \frac{D^2}{2\alpha_k}$.

So summing both sides of (*),

$\mathbb{E}\sum R(\theta_h) - R(\theta^*) \le \frac{D^2}{2\alpha_k} + \frac{1}{2}\sum \alpha_k M^2 - \sum_{h=1}^k \mathbb{E}\langle G(\theta_h) - \nabla R(\theta_h), \theta_h - \theta^*\rangle$.

Taking expectations on both sides and noting

$\mathbb{E}[\langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^*\rangle] = \mathbb{E}\left[\mathbb{E}[\langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^*\rangle | \theta_k]\right]$

$\qquad\qquad = \mathbb{E}[\langle \mathbb{E}[G(\theta_k)|\theta_k] - \nabla R(\theta_k), \theta_k - \theta^*\rangle] = 0,$

we get $\boxed{\sum\left(R(\theta_h) - R(\theta^*)\right) \le \frac{D^2}{2\alpha_k} + \frac{1}{2}\sum \alpha_h M^2}$.  Noting $R(\bar\theta_k) \le \frac{1}{k}\sum R(\theta_h)$,  we get the result. $\boxtimes$

Cor   For $\alpha_k = \frac{D}{M\sqrt{k}}$,    $\mathbb{E}R(\bar\theta_k) - R(\theta^*) \le \frac{3DM}{2\sqrt{k}}$.

Pf)  Noting $\sum_1^k \frac{1}{\sqrt{k}} \le \int_0^k \frac{1}{\sqrt{t}}dt = 2\sqrt{k}$, RHS $\le \frac{DM}{2\sqrt{k}} + \frac{DM}{\sqrt{k}}$.         $\boxtimes$.

Rmk   Think of $k$ as # access to gradient oracle. If $G(\theta) = \nabla_\theta \ell(\theta; z_i)$, then $k =$ # samples.

Rmk   Often, we iterate through data $C$ times. This gives gains on empirical loss. But population loss-wise, theory doesn't give gains as $C$ grows. In fact, we can't do better. We show this next class.