



SGD

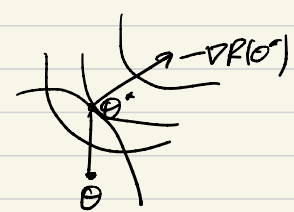
Def A function $R: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if $\forall \theta, \theta' \in \mathbb{R}^d, R(t\theta + (1-t)\theta') \leq tR(\theta) + (1-t)R(\theta') \quad \forall t \in [0,1]$.

Lemma Let $R: \mathbb{R}^d \rightarrow \mathbb{R}$ 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, $R(\theta + t(\theta' - \theta)) \leq R(\theta) + t(R(\theta') - R(\theta)) \Leftrightarrow R(\theta') - R(\theta) \geq \frac{1}{t}(R(\theta + t(\theta' - \theta)) - R(\theta))$. Send $t \rightarrow 0$
 '⇐' Define $\theta_t = t\theta + (1-t)\theta'$. Combining $R(\theta) \geq R(\theta_t) + \nabla R(\theta_t)^T(\theta - \theta_t)$, $R(\theta') \geq R(\theta_t) + \nabla R(\theta_t)^T(\theta' - \theta_t)$,
 $tR(\theta) + (1-t)R(\theta') \geq R(\theta_t) + \nabla R(\theta_t)^T(t\theta + (1-t)\theta' - \theta_t) \quad \forall t \in [0,1]$. □

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

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



Lemma $\theta^* = \operatorname{argmin}_{\theta \in \Theta} R(\theta)$ iff $\nabla R(\theta^*)^T(\theta - \theta^*) \geq 0 \quad \forall \theta \in \Theta$

Pf '⇐' From prev lemma, $R(\theta) - R(\theta^*) \geq \nabla R(\theta^*)^T(\theta - \theta^*) \geq 0 \quad \forall \theta \in \Theta$.

'⇒' $\nabla R(\theta^*)^T(\theta - \theta^*) = \lim_{t \rightarrow 0} \frac{1}{t}(R(\theta^* + t(\theta - \theta^*)) - R(\theta^*)) \geq 0 \quad \forall \theta \in \Theta$. □

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

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, $\|\theta' - \Pi_{\Theta}(\theta)\|_2^2 \leq \|\theta - \theta'\|_2 \|\theta' - \Pi_{\Theta}(\theta)\|_2 \quad \forall \theta' \in \Theta$. □

Stochastic gradients A stochastic gradient $G(\theta)$ is a RV st. $\mathbb{E}G(\theta) = \nabla R(\theta)$.

We study first-order optimization methods based on stoch. gradients.

(Canonical Problem)

minimize $\theta \in \Theta \quad \{\mathbb{E} \ell(\theta; z) =: R(\theta)\}$

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 Θ .

• 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.

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

Rank 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 (Θ) be compact. Assume $\exists R > 0$ s.t. $\forall \theta \in \Theta \|\theta - \theta^*\|_2 \leq R$, $\exists M > 0$ s.t. $\mathbb{E} \|G(\theta)\|_2^2 \leq M^2 \forall \theta \in \Theta$.
Let α_k be dec. pos. step sizes, and $\bar{\theta}_k = \frac{1}{k} \sum_{i=1}^k \theta_i$. Then,

$$\mathbb{E}[R(\bar{\theta}_k) - R(\theta^*)] \leq \frac{D^2}{2k\alpha_k} + \frac{1}{2k} \sum_{i=1}^k \alpha_i M^2.$$

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

$$\begin{aligned} \frac{1}{2} \|\theta_{k+1} - \theta^*\|_2^2 &= \frac{1}{2} \|\Pi_{\Theta}(\theta_k - \alpha_k G(\theta_k)) - \theta^*\|_2^2 \\ &\leq \frac{1}{2} \|\theta_k - \alpha_k G(\theta_k) - \theta^*\|_2^2 \text{ by non-expansiveness of } \Pi_{\Theta} \\ &= \frac{1}{2} \|\theta_k - \theta^*\|_2^2 - \alpha_k \langle G(\theta_k), \theta_k - \theta^* \rangle + \frac{\alpha_k^2}{2} \|G(\theta_k)\|_2^2. \end{aligned}$$

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

$$\begin{aligned} &= \frac{1}{2} \|\theta_k - \theta^*\|_2^2 - \alpha_k \langle \nabla R(\theta_k), \theta_k - \theta^* \rangle + \frac{\alpha_k^2}{2} \|G(\theta_k)\|_2^2 - \alpha_k \langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^* \rangle \\ &\leq \frac{1}{2} \|\theta_k - \theta^*\|_2^2 - \alpha_k (R(\theta_k) - R(\theta^*)) + \frac{\alpha_k^2}{2} \|G(\theta_k)\|_2^2 - \alpha_k \langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^* \rangle \text{ by convexity} \end{aligned}$$

Divide each side by α_k , and rearrange

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

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

So summing both sides of (*),

$$\mathbb{E} \sum_{k=1}^K R(\theta_k) - R(\theta^*) \leq \frac{D^2}{2\alpha_K} + \frac{1}{2} \sum_{k=1}^K \alpha_k M^2 - \sum_{k=1}^K \mathbb{E} \langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^* \rangle.$$

Taking expectations on both sides and noting

$$\begin{aligned} \mathbb{E} \langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^* \rangle &= \mathbb{E} \left[\mathbb{E} \left[\langle G(\theta_k) - \nabla R(\theta_k), \theta_k - \theta^* \rangle \mid \theta_k \right] \right] \\ &= \mathbb{E} \left[\langle \mathbb{E}[G(\theta_k) \mid \theta_k] - \nabla R(\theta_k), \theta_k - \theta^* \rangle \right] = 0, \end{aligned}$$

we get $\sum_{k=1}^K R(\theta_k) - R(\theta^*) \leq \frac{D^2}{2\alpha_K} + \frac{1}{2} \sum_{k=1}^K \alpha_k M^2$. Noting $R(\bar{\theta}_k) \leq \frac{1}{k} \sum_{i=1}^k R(\theta_i)$, we get the result. \square

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

Pf) Noting $\sum_{j=1}^k \frac{1}{j^2} \leq \int_0^k \frac{1}{j^2} dt = 2\sqrt{k}$, RHS $\leq \frac{DM}{2\sqrt{k}} + \frac{DM}{\sqrt{k}}$. \square

Remark Think of K as # access to gradient oracle. If $G(\theta) = \nabla_{\theta} l(\theta; z_i)$, then $K = \#$ samples.

Remark 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.

SGD



Def A function $R: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if $\forall \theta, \theta' \in \mathbb{R}^d$,

$$R(t\theta + (1-t)\theta') \leq tR(\theta) + (1-t)R(\theta') \quad \forall t \in [0, 1].$$

Lemma Let $R: \mathbb{R}^d \rightarrow \mathbb{R}$ be differentiable on the interior of its domain.

R is convex iff $R(\theta') \geq R(\theta) + \nabla R(\theta)^T(\theta' - \theta)$, $\forall \theta, \theta' \in \mathbb{R}^d$

↑ 1st order approx is a global minorization

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

Lemma $\theta^* = \operatorname{argmin}_{\theta \in \Theta} R(\theta)$ iff $\nabla R(\theta^*)^T(\theta - \theta^*) \geq 0 \quad \forall \theta \in \Theta$

Cor. Let Θ be a closed convex set in \mathbb{R}^d ,

Define the projection operator $\pi_{\Theta}(\theta) := \underset{\theta' \in \Theta}{\operatorname{argmin}} \|\theta - \theta'\|_2$.

Then, $\|\pi_{\Theta}(\theta) - \theta'\|_2 \leq \|\theta - \theta'\|_2 \quad \forall \theta' \in \Theta \quad \forall \theta \in \mathbb{R}^d$.

Stochastic gradients

A stochastic gradient $G(\theta)$ is a RV st. $\mathbb{E}G(\theta) = \nabla R(\theta)$.

We study first-order optimization methods based on stoch. gradients.

(Canonical Problem)

$$\text{minimize}_{\theta \in \mathcal{H}} \{ \mathbb{E} l(\theta; z) =: R(\theta) \}$$

- SGD Idea: Go in the direction of stoch. gradient, then project to \mathcal{H} .
- Algo: let $G_k(\theta)$ be a stoch. gradient of $R(\theta)$.

At each iteration k ,

$$\theta_{k+1} = \Pi_{\mathcal{H}}(\theta_k - \alpha_k G_k(\theta_k)) \quad \text{for some stepsize } \alpha_k > 0.$$

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

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

Theorem Let Θ be compact.

Assume $\exists D > 0$ s.t. $\sup_{\theta \in \Theta} \|\theta - \theta^*\|_2 \leq D$, $\exists M > 0$ s.t. $\mathbb{E} \|G(\theta)\|_2^2 \leq M^2 \forall \theta \in \Theta$.

Let α_k be dec, pos. step sizes, and $\bar{\theta}_K = \frac{1}{K} \sum_{i=1}^K \theta_k$

$$\mathbb{E}[R(\bar{\theta}_K) - R(\theta^*)] \leq \frac{D^2}{2K\alpha_K} + \frac{1}{2K} \sum_{i=1}^K \alpha_k M^2.$$

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

Pf) Noting $\sum_{i=1}^K \frac{1}{\sqrt{i}} \leq \int_0^K \frac{1}{\sqrt{t}} dt = 2\sqrt{K}$, RHS $\leq \frac{DM}{2\sqrt{K}} + \frac{DM}{\sqrt{K}}$.