

Neural network training.

For a given dataset \mathcal{D}_n , we aim to minimize the empirical risk

$$\hat{R}(\theta) = \hat{R}(f_\theta) = \frac{1}{n} \sum_{j=1}^n \ell(f_\theta(x_j), y_j) \quad (\text{training error})$$

For an algorithm A , we aim to characterize the optimization error

$$\hat{R}(A(\mathcal{D}_n)) - \inf_{f \in \mathcal{F}} \hat{R}(f)$$

While in some cases this can be done explicitly (e.g. linear regression) in general the problem of finding the minimum of a function \hat{R} is hard.

One method to (hopefully) solve this problem: move sequentially in the direction (in Θ) of steepest descent of \mathcal{R} by updating

$$\theta \leftarrow \theta - \gamma D_\theta R(\theta)$$



for a small timestep parameter γ . This method is called gradient descent:

$$\text{the update reads: } \underline{\theta_{k+1} = \theta_k - \gamma_k D_\theta R(\theta_k)} \quad (\text{GD})$$

Note: Why using this and not trying to solve $D_\theta R(\theta) = 0$?

Computation of $D_\theta R$ is cheap: consider

$$D_{w_{l+1}} f(\theta) = \sigma'_{l+1}(z_{l+1}(x)) W^l \sigma'_l(z_l(x)) \dots \sigma'_1(z_1(x)) \cdot x$$

In the above update, provided that we know σ'_e and z_e we are computing a complicated derivative by taking a product of known numbers (z_e were evaluated to find $f_\theta(x)$) so that

$$\begin{aligned}\theta_{k+1} &= \theta_k - \gamma_k D_\theta \hat{R}(\theta_k) \\ &= \theta_k - \gamma_k \frac{1}{n} \sum_{j=1}^n \underbrace{D_\theta f_{\theta_k}(x_j)}_{\text{"easy"}} \underbrace{D_f \ell(f_{\theta_k}(x_j), y_j)}_{\text{"easy": just evaluate } f_{\theta_k}}\end{aligned}$$

To study the dynamics of θ_k we write formally

$$\theta_{k+1} = \theta_k - \gamma D_\theta R(\theta_k) \rightarrow \frac{\theta_{k+1} - \theta_k}{\gamma} = -D_\theta R(\theta_k)$$

Based on the above, one expects that $\theta_k \approx \bar{\theta}_{k\gamma}$ where $\bar{\theta}: \mathbb{R}_+ \rightarrow \Theta$ solves

$$\begin{cases} \frac{d}{dt} \bar{\theta}_t = -D_\theta R(\bar{\theta}_t) \\ \bar{\theta}_0 = \theta_0 \end{cases} \quad (\text{gradient flow})$$

Def a differentiable function $F: \Theta \rightarrow \mathbb{R}$ is L -Lipschitz smooth if

$$\|F(\theta') - F(\theta) - D_\theta F(\theta)(\theta' - \theta)\| \leq \frac{L}{2} \|\theta' - \theta\|^2 \quad \text{for } L \geq 0$$

Lemma: F L -Lipschitz smooth $\iff D_\theta F$ is L -Lipschitz continuous

$$\begin{aligned}\text{Pf: } \Leftarrow \text{Taylor } F(\theta') &= F(\theta) + \int_\theta^{\theta'} D_\theta F(\eta) d\eta \\ &= F(\theta) + \int_\theta^{\theta'} D_\theta F(\theta) + (D_\theta F(\eta) - D_\theta F(\theta)) d\eta\end{aligned}$$

$$\begin{aligned}
&= F(\vartheta) + D_{\vartheta} F(\vartheta)(\vartheta' - \vartheta) + \int_{\vartheta}^{\vartheta'} D_{\vartheta} F(\eta) - D_{\vartheta} F(\vartheta) d\eta \\
&\leq F(\vartheta) + D_{\vartheta} F(\vartheta)(\vartheta' - \vartheta) + \frac{L}{2} \|\vartheta - \vartheta'\|^2
\end{aligned}$$

$$\stackrel{**}{=} (D_{\vartheta} F(\vartheta) - D_{\vartheta} F(\vartheta'))(\vartheta - \vartheta')$$

Remark: Let $F \in C^2(\Theta)$. Then if F is L -Lipschitz smooth we have

$$-L \mathbb{I} \leq \|D_{\vartheta}^2 F(\vartheta)\| \leq L \mathbb{I}$$

$$\text{Indeed for } \|\eta\|=1: \|\eta^T D_{\vartheta}^2 F(\vartheta)\| = \lim_{h \rightarrow 0} \left\| \frac{D_{\vartheta} F(\vartheta + h\eta) - D_{\vartheta} F(\vartheta)}{h} \right\| \leq \frac{L \|h\eta\|}{h} = L$$

Lemma: The solution to (GF) exists and is unique

Pf: Consequence of classical existence and uniqueness of solution to ODEs with Lipschitz vector fields.

From now on we set $F = R$

Lemma: Let $D_{\vartheta} R$ be Lipschitz. For every $T > 0$ there exists $C > 0$ s.t.

$$\|\bar{\vartheta}_{k\delta} - \vartheta_k\| \leq C\delta \quad \forall k \in \{0, \dots, \lfloor \frac{T}{\delta} \rfloor\}$$

$$\text{Proof: } \bar{\vartheta}_{t+\delta} = \bar{\vartheta}_t - \delta D_{\vartheta} R(\bar{\vartheta}_t) + \underbrace{\int_t^{t+\delta} D_{\vartheta} R(\bar{\vartheta}_s) - D_{\vartheta} R(\bar{\vartheta}_t) ds}_{\delta^2}$$

And we have that

$$\left\| \int_t^{t+\delta} D_{\vartheta} R(\bar{\vartheta}_s) - D_{\vartheta} R(\bar{\vartheta}_t) ds \right\| \leq \delta \sup_{s \in (t, t+\delta)} \|D_{\vartheta} R(\bar{\vartheta}_s) - D_{\vartheta} R(\bar{\vartheta}_t)\| \leq \delta \sup_{s \in (t, t+\delta)} \|\bar{\vartheta}_s - \bar{\vartheta}_t\|$$

$$\leq \gamma^2 \sup_{s \in [t, t+\gamma]} \|D_g R(\bar{\theta}_s)\|$$

By Lipschitz smoothness we have that $\sup_{t \in [0, T]} D_g R(\bar{\theta}_t) \leq C$

$$\begin{aligned} e_{k+1} &= \|\bar{\theta}_{(k+1)\gamma} - \theta_{k+1}\| = \|\bar{\theta}_{k\gamma} - \gamma D_g R(\bar{\theta}_{k\gamma}) - \theta_k + \gamma D_g R(\theta_k)\| + C\gamma^2 \\ &\leq (1 + \gamma\lambda) e_k + C\gamma^2 \end{aligned}$$

$$\Rightarrow e_{k+1} = C\gamma^2 \sum_{j=1}^k (1 + \gamma\lambda)^j = \frac{C\gamma^2}{\lambda\gamma} ((1 + \gamma\lambda)^{k+1} - 1) \leq \frac{C\gamma}{\lambda} (e^{\lambda\gamma(k+1)} - 1) \leq C'\gamma \quad \square$$

This result justifies the use of GF instead of GD.

Q: when does (GF) converge, and when can we bound

$$\hat{R}(\theta_0) - \inf_{\theta \in \mathcal{T}} \hat{R}(\theta) \quad ?$$

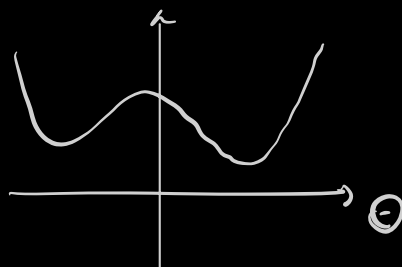
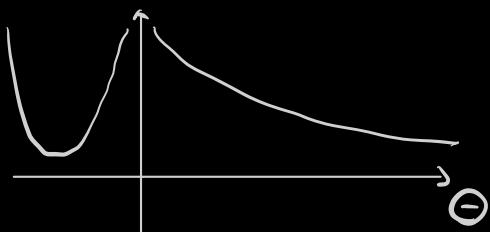
Remark: Since

$$\frac{d}{dt} R(\theta_t) = D_g R(\theta_t) \cdot (-D_g R(\theta_t)) = -\|D_g R(\theta_t)\| \leq 0$$

we have $R(\theta_t) \leq R(\theta_0)$.

So if R is bounded from below, as $t \rightarrow \infty$ by monotone convergence theorem $R(\theta_t)$ converges.

⚠ This does NOT show that θ_t converges!



However, we can't say much about optimality of limit points.

Example (linear regression) $n \geq d$ $X^T X \succeq \lambda \mathbb{1}$

$$\hat{R}(\theta) = \frac{1}{n} \|X\theta - y\|_2^2 \implies \theta_* = (X^T X)^{-1} X^T y$$

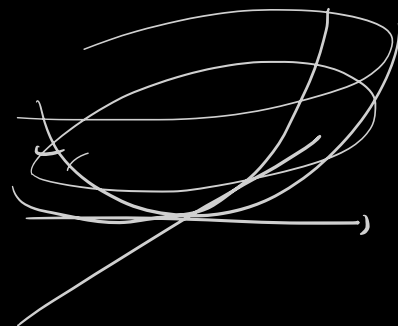
$$\implies R(\theta_*) = \frac{1}{n} \|X(X^T X)^{-1} X^T y\|_2^2$$

$$\begin{aligned} \frac{d}{dt} \|\theta_t - \theta_*\|_2^2 &= \langle \theta_t - \theta_*, -D_\theta \hat{R}(\theta_t) \rangle = -\langle \theta_t - \theta_*, X^T (X\theta_t - y) \rangle \\ &= -\langle \theta_t - \theta_*, X^T X (\theta_t - \theta_*) \rangle \leq -\lambda \|\theta_t - \theta_*\|_2^2 \end{aligned}$$

$$\stackrel{G}{\implies} \|\theta_t - \theta_*\|_2^2 \leq e^{-\lambda t} \|\theta_0 - \theta_*\|_2^2 \quad \text{☺}$$

This happens because \hat{R} is strongly convex:

$$\begin{aligned} \hat{R}(\theta) - \hat{R}(\theta_*) &= D_\theta R(\theta_*) \cdot (\theta - \theta_*) + \frac{1}{2} (\theta - \theta_*)^T D_\theta^2 R(\theta_*) (\theta - \theta_*) \\ &= \frac{1}{2} (\theta - \theta_*)^T X^T X (\theta - \theta_*) \end{aligned}$$



The convergence properties of GF are guaranteed when R is convex:

Def: a differentiable $\hat{R}: \Theta \rightarrow \mathbb{R}$ is λ -strongly convex for $\lambda \in \mathbb{R}$ if.

$$(*) \quad \hat{R}(\theta') \geq \hat{R}(\theta) + \langle D_{\theta} \hat{R}(\theta), \theta' - \theta \rangle + \frac{\lambda}{2} \|\theta' - \theta\|^2 \quad \text{for any } \theta, \theta' \in \Theta$$

Lemma: If R is convex (λ -strongly convex with $\lambda=0$) if $DR(\theta_*)=0$, $\theta_* \in \min R$

Pf: trivial from definition.

Lemma let \hat{R} be λ -convex, then the following holds:

$$\langle D_{\theta} \hat{R}(\theta') - D_{\theta} \hat{R}(\theta), \theta' - \theta \rangle \geq \lambda \|\theta - \theta'\|^2$$

Proof: Consider:

$$\hat{R}(\theta') \geq \hat{R}(\theta) + \langle D_{\theta} \hat{R}(\theta), \theta' - \theta \rangle + \frac{\lambda}{2} \|\theta - \theta'\|^2$$

adding the above to the definition of λ strongly convex we are done \square

Thm: Let R be λ -strongly convex and θ_* be the minimizer, then

$$\|D_{\theta} R(\theta)\|_2^2 \geq 2\lambda (R(\theta) - R(\theta_*))$$

Proof: Setting $\theta' = \theta - \frac{1}{\lambda} D_{\theta} R(\theta)$ in $(*)$ we get

$$R(\theta_*) \geq R(\theta) - \frac{1}{\lambda} \|D_{\theta} R(\theta)\|^2 + \frac{1}{2\lambda} \|D_{\theta} R(\theta)\|^2 = R(\theta) - \frac{1}{2\lambda} \|D_{\theta} R(\theta)\|^2$$

Thm: Let $R(\theta)$ be λ -convex, ^{L -smooth.} then there exists a unique θ^* and

$$\|\bar{\theta}_t - \theta_*\|^2 \leq \|\bar{\theta}_0 - \theta_*\|^2 e^{-2\lambda t}$$

$$R(\bar{\theta}_t) - R(\theta_*) \leq (R(\bar{\theta}_0) - R(\theta_*)) e^{-2\lambda t}$$

Proof: $\frac{d}{dt} \frac{1}{2} \|\bar{\theta}_t - \theta_*\|^2 = \langle \bar{\theta}_t - \theta_*, \frac{d}{dt} \bar{\theta}_t \rangle$

$$= -\langle \bar{\theta}_t - \theta_*, D_{\theta} R(\bar{\theta}_t) - D_{\theta} R(\theta_*) \rangle$$

$$\leq -\lambda \|\bar{\theta}_t - \theta_*\|^2$$

→ Gronwall

Similarly $\frac{d}{dt} R(\theta_t) - R(\theta_*) = \langle D_{\theta} R(\theta_t), \dot{\theta}_t \rangle$

Thm $\quad \quad \quad = -\|D_{\theta} R(\theta_t)\|^2$

$$\leq -2\lambda (R(\theta_t) - R(\theta_*))$$

→ Gronwall

Lemma: Assume R is L -smooth, λ strongly convex.

Choose $\eta = \frac{1}{L}$. Then $R(\theta_t) - R(\theta_*) \leq e^{-\frac{t}{\kappa}} (0)$

PF: $R(\theta_t) = R(\theta_{t+1} - \frac{1}{L} D_{\theta} R(\theta_{t+1})) \leq R(\theta_{t+1}) - \frac{1}{L} \|D_{\theta} R(\theta_{t+1})\|^2$

$$\Rightarrow R(\theta_t) - R(\theta_*) \leq R(\theta_{t+1}) - R(\theta_*) - \frac{1}{2L} \|D_{\theta} R(\theta_{t+1})\|^2$$

$$\leq (1 - \frac{\lambda}{2L}) (R(\theta_{t+1}) - R(\theta_*))$$

□

Gradient flows for convex functionals.

Example: Consider now the case where some eigenvalues of $X^T X$ are 0.

$$\frac{d}{dt} \vartheta_t - \vartheta_* = -X^T X (\vartheta_t - \vartheta_*) \implies \vartheta_t - \vartheta_* = e^{-X^T X t} (\vartheta_0 - \vartheta_*)$$

$$R(\vartheta_t) - R(\vartheta_*) = \left(\frac{1}{2} (\vartheta_t - \vartheta_*)^T X^T X (\vartheta_t - \vartheta_*) \right) = \frac{1}{2} (\vartheta_0 - \vartheta_*)^T \underbrace{e^{-X^T X t} X^T X e^{-X^T X t}}_{A(t)} (\vartheta_0 - \vartheta_*)$$

$A(t)$

$$\text{Eig}(A(t)) = e^{-\lambda t} \cdot \lambda \cdot e^{-\lambda t} = \frac{1}{t} (\lambda t) e^{-\lambda t} \leq \frac{1}{t}$$

Theorem: Let $R(\vartheta)$ be convex and Lipschitz smooth. Then



$$R(\bar{\vartheta}_t) - R(\bar{\vartheta}_*) \leq C \frac{1}{t}.$$

Proof: $R(\vartheta') \geq R(\vartheta) - D_\vartheta R(\vartheta)(\vartheta - \vartheta')$

$$\begin{aligned} \implies \frac{d}{dt} \frac{1}{2} \|\vartheta_t - \vartheta_*\|^2 &= (\vartheta_t - \vartheta_*)^T (-D_\vartheta R(\vartheta_t)) \\ &\leq R(\vartheta_t) - R(\vartheta_t) \end{aligned}$$

$$\implies \int_0^t R(\vartheta_s) - R(\vartheta_*) ds \leq -\frac{1}{2} \|\vartheta_t - \vartheta_*\|^2 + \frac{1}{2} \|\vartheta_0 - \vartheta_*\|^2 \leq \frac{1}{2} \|\vartheta_0 - \vartheta_*\|^2$$

$$\implies \frac{1}{t} \int_0^t R(\vartheta_s) ds - R(\vartheta_*) \leq \frac{1}{2t} C_0$$

$$\text{but } \frac{1}{t} \int_0^t R(\vartheta_s) ds \geq R(\vartheta_t)$$

$$\frac{d}{dt} R(\vartheta_t) \leq 0$$

$$\implies R(\vartheta_t) - R(\vartheta_*) \leq \frac{1}{2t} C_0$$

□

Lemma: Let $\bar{\theta}_t$ converge to θ_* and there exists $N(\theta_*)$ s.t.
 $D_\theta R(\theta) \succeq \varepsilon^2 I \quad \forall \theta \in N$ then for γ sufficiently small.

$$\sup_t \|\bar{\theta}_t - \theta_t\| \leq C\gamma$$

Pf: Let T be such that $\bar{\theta}_t \in N$.