

Global convergence of gradient descent



Global convergence of gradient descent

Theorem (Du et al. '18, Allen-Zhu et al. '18, Zou et al '19) If the width of each layer is $\text{poly}(n)$ where n is the number of data. Using random initialization with a particular scaling, gradient descent finds an approximate global minimum in polynomial time.

ϵ - close to global min

ϵ error

$\text{poly}(n, \frac{1}{\epsilon})$

Gradient Flow: a Kernel Point of View

$$\cdot \mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f(\theta, x_i), y_i)$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \frac{1}{n} \sum_{i=1}^n \ell'(f(\theta, x_i), y_i) \frac{\partial f(\theta, x_i)}{\partial \theta}$$

$$\text{GF: } \frac{d\theta(t)}{dt} = - \frac{\partial \mathcal{L}(\theta(t))}{\partial \theta}$$

If \mathcal{L} is *strongly convex*
 \Rightarrow unique θ^* , $\theta(t) \rightarrow \theta^*$

if *over-parameterized*
multiple θ^*

Gradient Flow: a Kernel Point of View

$$u_i(t) = f(\theta(t), x_i), \quad u(t) = \begin{pmatrix} u_1(t) \\ \vdots \\ u_n(t) \end{pmatrix} \in \mathbb{R}^n$$

$$\frac{du(t)}{dt} = -\frac{1}{n} H(t) \cdot l'(u(t), y)$$

$$H(t) \in \mathbb{R}^{n \times n}$$
$$[H(t)]_{ij} = \left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, \frac{\partial u_j(t)}{\partial \theta(t)} \right\rangle$$

$$l'(u(t), y) \in \mathbb{R}^n$$
$$[l'(u(t), y)]_i = l'(u_i(t), y_i)$$

$y \in \mathbb{R}^n$, labels

$$u(t) \rightarrow y$$

Gradient Flow: a Kernel Point of View

If l is quadratic, $l(u(t), y) = \frac{1}{2} \|u(t) - y\|_2^2$
 $l'(u(t), y) = u(t) - y$

$$\frac{d u(t)}{d t} = -\frac{1}{n} H(t) (u(t) - y)$$

(Claim:

If $H(t)$ is always positive definite:
 $\forall t, \exists \lambda_0 > 0, \lambda_{\min}(H(t)) \geq \lambda_0$

$$\rightarrow \frac{1}{2} \|u(t) - y\|_2^2 \rightarrow 0$$

(Pf:

$$\frac{d \left(\frac{1}{2} \|u(t) - y\|_2^2 \right)}{d t} = -\frac{1}{n} (u(t) - y)^T H(t) (u(t) - y)$$
$$\leq -\frac{1}{n} \lambda_0 \|u(t) - y\|_2^2$$

Gradient Flow: a Kernel Point of View

Consider $\frac{d}{dt} \left(\exp\left(\frac{\lambda_0 t}{n}\right) \frac{1}{2} \|u(t) - y\|_2^2 \right)$

$$= \frac{\lambda_0}{2n} \exp\left(\frac{\lambda_0 t}{n}\right) \|u(t) - y\|_2^2 + \frac{d\left(\frac{1}{2} \|u(t) - y\|_2^2\right)}{dt} \cdot \exp\left(\frac{\lambda_0 t}{n}\right)$$

$$\leq \exp\left(\frac{\lambda_0 t}{n}\right) \cdot \|u(t) - y\|_2^2 \left(\frac{\lambda_0}{2n} - \frac{\lambda_0}{n} \right) < 0$$

$\exp\left(\frac{\lambda_0 t}{n}\right) \cdot \frac{1}{2} \|u(t) - y\|_2^2$ is decreasing

when $t = 0$, assume $\frac{1}{2} \|u(0) - y\|_2^2 = C$, $\partial(C)$

$$\frac{1}{2} \|u(t) - y\|_2^2 \in \exp\left(-\frac{\lambda_0 t}{n}\right) \cdot C$$

$t \rightarrow \infty$, loss $\rightarrow 0$, $u(t) \rightarrow y$

$$t = O\left(\log\left(\frac{1}{\epsilon}\right)\right) \text{ verte}$$

Gradient Flow: a Kernel Point of View

$$f(\theta, x) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(w_r^T x)$$

m : width, $x \in \mathbb{R}^d$, $a_r \in \mathbb{R}$, $w_r \in \mathbb{R}^d$, $\sigma(\cdot)$: ReLU

• Initialization: $a_r \sim \text{unif}\{-1, 1\}$ *only for simplicity of proof*

$$\Rightarrow f(\theta(0), x) = \mathcal{O}(1)$$

$$w_r \sim \mathcal{N}(0, I)$$

NTK

• Training: only train w_1, \dots, w_m

$$\min_{w_1, \dots, w_m} \frac{1}{n} \sum_{i=1}^n (f(x_i, a, w) - y_i)^2$$

$$u_i(t) = f(x_i, a, w(t))$$

$$\frac{du(t)}{dt} = -\frac{1}{n} H(t) (u(t) - y)$$

Idea:
 $H(t) \approx H^*$
 for all t
 $(H^*)_{ij} = \lim_{n \rightarrow \infty} \frac{\partial u_i(t)}{\partial w_j}$
 (or $\frac{\partial u_i(t)}{\partial \theta}$)

Gradient Flow: a Kernel Point of View

$\lambda_{\min}(H^*) \geq \lambda_0$
 \Leftrightarrow Kernel is universal

$$H_{ij}(t) = \left\langle \frac{\partial U_i(t)}{\partial w(t)}, \frac{\partial U_j(t)}{\partial w(t)} \right\rangle$$

$$= \sum_{r=1}^m \left\langle \frac{\partial U_i(t)}{\partial w_r(t)}, \frac{\partial U_j(t)}{\partial w_r(t)} \right\rangle$$

$$\frac{\partial U_i(t)}{\partial w_r(t)} = \frac{1}{\sqrt{m}} a_r \cdot x_i \cdot \mathbb{1} \{ w_r(t)^T x_i \geq 0 \}$$

$$H_{ij}(t) = \frac{1}{m} \sum_{r=1}^m \left\langle a_r x_i \cdot \mathbb{1} \{ w_r(t)^T x_i \geq 0 \}, a_r x_j \cdot \mathbb{1} \{ w_r(t)^T x_j \geq 0 \} \right\rangle$$

$$= \frac{1}{m} x_i^T x_j - \sum_{r=1}^m \mathbb{1} \{ w_r(t)^T x_i \geq 0, w_r(t)^T x_j < 0 \}$$

To show: $H(t) \approx H^*$,
 1) $H(0) \approx H^*$
 2) $H(t) \approx H(0), \forall t$

Hoeffding inequality:

R.V. $z_1, \dots, z_n \stackrel{i.i.d.}{\sim} D, |z_i| \leq 1$

if $n = \Omega\left(\frac{\log(1/\delta)}{\epsilon^2}\right), 0 < \delta < 1$

w.p. $1 - \delta, \left| \frac{1}{n} \sum_{i=1}^n z_i - \mathbb{E}[z_i] \right| \leq \epsilon$

$$H_{ij}(0) = X_i^T X_j \underbrace{\frac{1}{m} \sum_{r=1}^m \mathbb{1}_{\{w_r(0)^T X_i \geq 0, w_r(0)^T X_j \geq 0\}}}_{\text{average}} \underbrace{z_r}_{z_r}$$

$$H_{ij}^* = \mathbb{E}_{w \sim \mathcal{N}(0, I)} (X_i^T X_j \cdot \mathbb{1}_{\{w^T X_i \geq 0, w^T X_j \geq 0\}})$$

when m is sufficiently large

$$H_{ij}(0) \rightarrow H_{ij}^*$$

$$\Rightarrow H(0) \rightarrow H^*$$

Want to show $H(t) \approx H(0)$
for simplicity:

1) just train till time t

2) $y_i = 0$ (1)

3) $\|X_i\|_2 = 1$

$$|h_{ij}^*| = X_i^T X_j \cdot \frac{\pi - \arccos(X_i^T X_j)}{2\pi}$$

★ every weight w_{ij} only moves

a little: $O\left(\frac{1}{\sqrt{m}}\right)$

lazy training

$$\begin{aligned}
& \|w_V(t) - w_P(0)\|_2 \\
&= \left\| \int_0^t \frac{dw_V(\tau)}{d\tau} d\tau \right\|_2 \\
&\leq \int_0^t \left\| \frac{dw_V(\tau)}{d\tau} \right\|_2 d\tau \\
&= \int_0^t \left\| -\frac{1}{\sqrt{m}} \frac{1}{n} \sum_{i=1}^n (u_i(\tau) - y_i) u_V(x_i) \right\| \left\{ w_V(\tau) x_i^T \right\|_2 d\tau \\
&\qquad\qquad\qquad \underbrace{\hspace{15em}}_{O(1)}
\end{aligned}$$

$$\leq C \cdot \int_0^t \frac{1}{\sqrt{m}} d\tau$$

$$\leq \frac{C \cdot t}{\sqrt{m}}$$

ReLU Smoothness

smooth: small deviation in parameters

\Rightarrow small deviation in function derivative

$$\begin{cases} H_{ij}(t) = X_i^T X_j \frac{1}{m} \sum_{v=1}^m \mathbb{1} \{ w_v(t)^T X_i \geq 0, w_v(t)^T X_j \geq 0 \} \\ H_{ij}(0) = X_i^T X_j - \dots - \dots \cdot w_v(0) - \dots - \dots \cdot w_v(0) - \dots - \dots \end{cases}$$

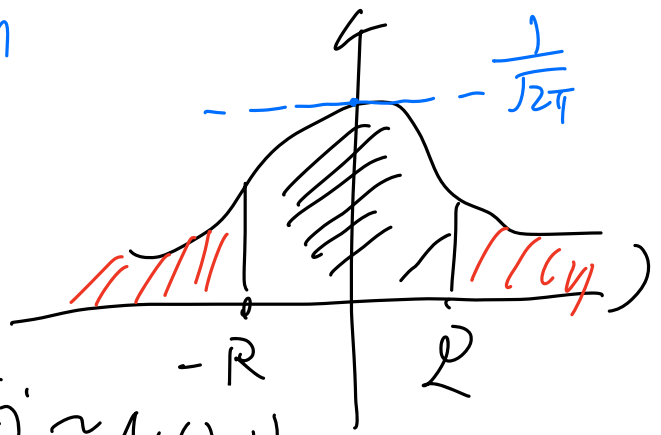
$$|H_{ij}(t) - H_{ij}(0)| \leq \frac{X_i^T X_j}{m} \left[\sum_{v=1}^m \mathbb{1} \{ \text{sgn}(w_v(t)^T X_i) \neq \text{sgn}(w_v(0)^T X_i) \} + \sum_{v=1}^m \mathbb{1} \{ \text{sgn}(w_v(t)^T X_j) \neq \text{sgn}(w_v(0)^T X_j) \} \right]$$

$\Rightarrow O\left(\frac{1}{\sqrt{m}}\right)$

\Rightarrow want to show # of pattern changes

Cauchy's Anti-Concentration

$$P_{Z \sim \mathcal{N}(0,1)} (|Z| \leq R) \leq \frac{2R}{\sqrt{2\pi}}$$



$$w_v(0) \sim \mathcal{N}(0, I) \Rightarrow w_v(0)^T X_i \sim \mathcal{N}(0, 1)$$

$$\text{If } \forall \gamma, \|w_v(t) - w_v(0)\| \leq \Delta w$$

($\Delta w \rightarrow 0$, as $m \rightarrow \infty$)

Let's choose $R > \Delta w$

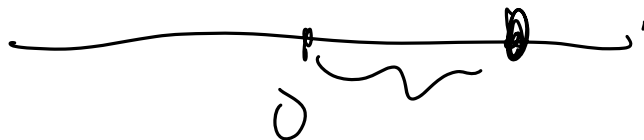
$$\text{Suppose } |w_v(0)^T X_i| \geq R, \text{ u.p. } 1 - \frac{2R}{\sqrt{2\pi}}$$

$$\Rightarrow \text{sgn}(w_v(t)^T X_i) = \text{sgn}(w_v(0)^T X_i)$$

$$|w_v(0)^T X_i - w_v(t)^T X_i|$$

$$\leq \|X_i\|_2 \cdot \|w_v(0) - w_v(t)\|_2$$

$$< R \leq |w_v(0)^T X_i|$$



$$P_V (|W_V(0)^T X_i| \leq \delta w) \leq \frac{2\delta w}{\sqrt{2\pi}}$$

We know $\delta w \rightarrow 0$, as $m \rightarrow \infty$

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1} \{ \text{sgn}(W_V^T(t) X_i) \neq \text{sgn}(W_V(0)^T X_i) \} \rightarrow 0$$

$$\|H(t) - H(0)\|_F \xrightarrow{m \rightarrow \infty} 0$$

$$\|H(0) - H^*\|_F \xrightarrow{m \rightarrow \infty} 0$$

$$\frac{du(t)}{dt} = -\frac{1}{h} H(t) (u(t) - y)$$

$$\approx -\frac{1}{h} \underbrace{H^*}_{\text{universal}} (u(t) - y)$$

universal, \rightarrow $\dim(H^*) > 0$

$$\Rightarrow u(t) \rightarrow y$$

if $m = \text{poly}(n, \frac{1}{\epsilon})$

$$\|H(t) - H^*\| \leq \epsilon$$