

# Global convergence of gradient descent

---

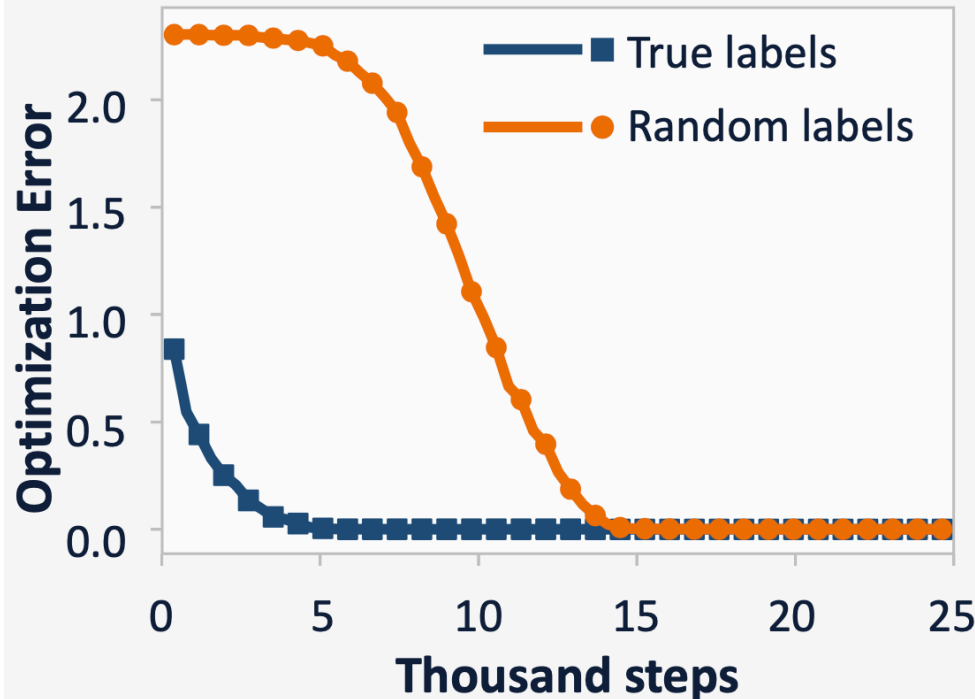


# Gradient descent finds global minima

over-parameterized

Practice: gradient descent

$$\theta(t+1) \leftarrow \theta(t) - \eta \frac{\partial L(\theta(t))}{\partial \theta(t)}$$



Optimization error  $\rightarrow 0$  for both *true labels* and *random labels* !

Zhang Bengio Hardt Recht Vinyals 2017

Understanding DL Requires Rethinking Generalization

# 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 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.

in  $\epsilon$ -optimal  
 $\text{poly}\left(n, \frac{1}{\epsilon}\right)$  time

# Gradient Flow: a Kernel Point of View

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

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

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

- If  $L$  is strongly convex  $\Rightarrow$  unique  $\theta^*$ ,  $\theta(t) \rightarrow \theta^*$
- If over-parameterized  $\Rightarrow$  multiple  $\theta^*$

# Gradient Flow: a Kernel Point of View

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

0-loss:  $u(t) \rightarrow y$   
 $y$ : labels,  $\mathbb{R}^n$

$$\frac{dM(t)}{dt} = - \frac{1}{n} H(t) \cdot \lambda'(u(t), y) \in \mathbb{R}^n$$

$$H(t) \in \mathbb{R}^{n \times n}$$

$$[H(t)]_{ij}$$

$$= \left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, \frac{\partial u_j(t)}{\partial \theta(t)} \right\rangle$$

$$\lambda'(u(t), y) \in \mathbb{R}^n$$

$$[\lambda'(u(t), y)]_i$$

$$= \lambda'(u_i(t), y_i)$$

# Gradient Flow: a Kernel Point of View

If  $l$  is quadratic,  $l(u(t), y) = \frac{1}{2} (u(t) - y)^2$

$$l'(u(t), y) = u(t) - y$$

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

Claim: If  $H(t)$  is always positive definite  
 $\forall t, \exists \lambda_0 > 0, \lambda_{\min}(H(t)) > 0$   
 $\rightarrow \frac{1}{2} \|u(t) - y\|_2^2 \rightarrow 0$

Prf:

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

# Gradient Flow: a Kernel Point of View

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

$$= \frac{\lambda_0}{2n} \exp\left(-\frac{\lambda_0 t}{n}\right) \|u(t) - y\|_2^2 + \frac{d\left(\frac{1}{2} \|u(t) - y\|_2^2\right)}{dt} \exp\left(-\frac{\lambda_0 t}{n}\right)$$
$$\leq \exp\left(-\frac{\lambda_0 t}{n}\right) \cdot \|u(t) - y\|_2^2 \left( \frac{\lambda_0}{2n} - \frac{\lambda_0}{n} \right) < 0$$

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

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

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

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

$$\Rightarrow u(t) \rightarrow y \text{ as } t \rightarrow \infty$$

□

# Gradient Flow: a Kernel Point of View

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

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

Initialization:  $a_r \sim \text{unif}\{-1, 1\}$   
 $w_r \sim \mathcal{N}(0, I)$

just for simplicity

$$\Rightarrow f(\theta(0), x) = 0$$

Training: only train  $w_1, \dots, w_m$

$$\min_{w_1, \dots, w_m}$$

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

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

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

$H^{\#}$ :  $n \times n$   
 Kernel:  $\sqrt{TK}$

Idea:

- $H(t)$  almost constant for  $t \ll \frac{1}{\lambda}$
- $H(t) \approx H^{\#}$
- $[H^{\#}]_{ij}$
- $\lim_{m \rightarrow \infty} \mathbb{E}_{\text{init}} \left[ \frac{\partial u_i}{\partial \theta_j} \right]$



# Gradient Flow: a Kernel Point of View

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

use functional analysis  
to show  $\lambda_{min}(H^*) = \lambda_0 > 0$

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

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

$$= \frac{1}{\sqrt{m}} \sum_{r=1}^m \left\langle a_r X_i \mathbb{1}_{\{w_r^T X_i \geq 0\}}, a_r X_j \mathbb{1}_{\{w_r^T X_j \geq 0\}} \right\rangle$$

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

$$(a_r^2 = 1)$$

$$= \frac{1}{m} \sum_{r=1}^m X_i^T X_j \mathbb{1}_{\{w_r^T X_i \geq 0, w_r^T X_j \geq 0\}}$$

$$1) H(0) \approx H^*$$

To show:  $H(t) \approx H^*$

2)  $H(t) \approx H(0), \forall t$

Hoeffding inequality

r.v.  $z_1, \dots, z_n \stackrel{\text{i.i.d.}}{\sim} D$ ,  $|z_i| \leq 1$

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

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

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

$z_r$

$m$  r.v. average

when  $m$  is sufficiently large

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

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

|

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

1) just train till some time  $t$

2)  $y_i = 0$  (1)

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

(note  $H_{ij}^* = x_i^T x_j$  .  $\frac{\pi - \arccos(x_i^T x_j)}{2\pi}$ )

will show: every  $w_r$  only moves

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

$$\begin{aligned}
& \|w_r(t) - w_r(0)\|_2 \\
&= \left\| \int_0^t \frac{dw_r(\tau)}{d\tau} d\tau \right\|_2 \\
&\leq \int_0^t \left\| \frac{dw_r(\tau)}{d\tau} \right\|_2 d\tau \\
&= \int_0^t \left\| -\frac{1}{\sqrt{m}} \frac{1}{N} \sum_{i=1}^N (u_i(\tau) - y_i) \text{cov}(x_i) \right\| \{w_r^T(\tau) x_i\}_2 d\tau \\
&\qquad\qquad\qquad O(1) \\
&\leq \frac{C \cdot t}{\sqrt{m}}
\end{aligned}$$

# ReLU smoothness

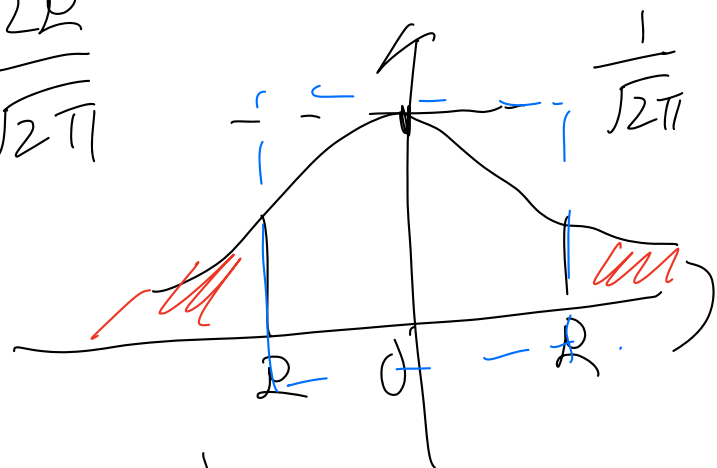
$$H_{ij}(t) = X_i^T X_j \frac{1}{m} \sum_{r=1}^m \mathbb{1} \{ w_r(t)^T X_i \geq 0, w_r(t)^T X_j \geq 0 \}$$

$$H_{ij}(0) = \dots \dots \dots w_r(0) \dots, w_r(0) \dots$$

$$\begin{aligned} & |H_{ij}(t) - H_{ij}(0)| \\ & \leq \frac{X_i^T X_j}{m} \sum_{r=1}^m \left[ \mathbb{1} \{ \text{sgn}(w_r(t)^T X_i) \neq \text{sgn}(w_r(0)^T X_i) \} \right. \\ & \quad \left. + \mathbb{1} \{ \text{sgn}(w_r(t)^T X_j) \neq \text{sgn}(w_r(0)^T X_j) \} \right] \end{aligned}$$

# Gaussian Anti-concentration

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



$$w_r(0) \sim N(0, I)$$

$$\Rightarrow w_r(0)^T X_i \sim N(0, 1)$$

Let's choose  $\rho > \Delta W$

$$\text{we know } |w_r(0)^T X_i| \geq \rho, \text{ w.p. } 1 - \frac{2\rho}{\sqrt{2\pi}}$$

$$\text{if } \|w_r(t) - w_r(0)\| \leq \Delta W < \rho$$

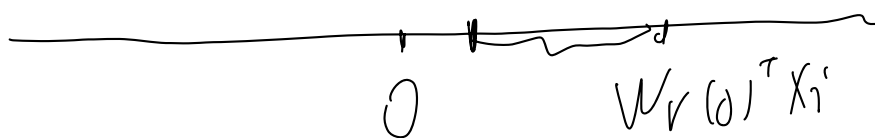
$$\text{Claim: } \text{sgn}(w_r(t)^T X_i) = \text{sgn}(w_r(0)^T X_i)$$

$$\|w_r(0)^T X_i - w_r(t)^T X_i\|$$

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

$$\leq \rho < |w_r(0)^T X_i|$$

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



$$P_r ( |w_r(0)^T x_i| \leq \Delta W ) \leq \frac{2\Delta W}{\sqrt{2\pi}}$$

We know  $\Delta W \rightarrow 0$ , as  $m \rightarrow \infty$

$$\Rightarrow \frac{1}{m} \sum_{r=1}^m \mathbb{1} \{ \text{sgn}(w_r^T(t) x_i) \neq \text{sgn}(w_r^T(0) x_i) \}$$

$\rightarrow 0$

$$H_{ij}(t) \rightarrow H(0) \quad \text{as } m \rightarrow \infty$$

$\Downarrow$

$$\| H(t) - H(0) \|_F \rightarrow 0 \quad \text{as } m \rightarrow \infty$$

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

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

$$(m \rightarrow \infty) \rightarrow -\frac{1}{n} H^* (u(t) - y)$$

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

$$m = \text{poly}(n), \quad u(t) \rightarrow y \quad \square$$