

• HW 1 Due Tomorrow 11:59 PM

• 2 Late Days

• Course Evaluation 2/3

Global convergence of gradient descent

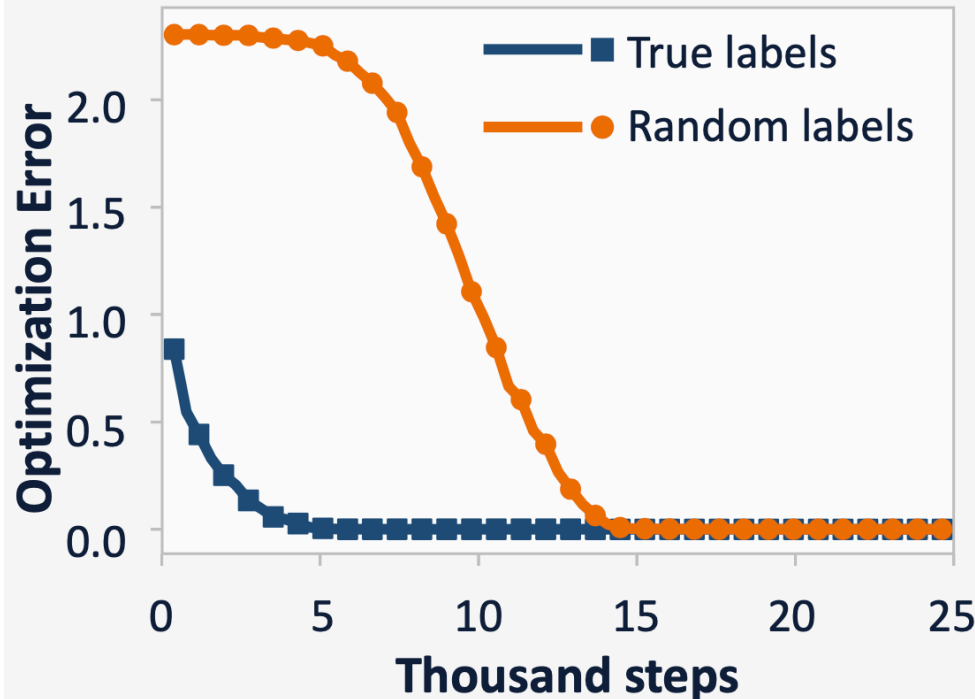


Gradient descent finds global minima

Practice: gradient descent

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

parameter
→ ↵



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

(convex loss)

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 - global min

$\text{poly}(n) \log(\frac{1}{\epsilon})$
for quadratic

Neural Tangent Kernel

proof for a two-layer NN

Gradient Flow: a Kernel Point of View

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

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

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

if $L(\theta)$ strongly convex, \exists unique θ^* , $\theta(t) \rightarrow \theta^*$
for NN, # of parameters, $\dim(\theta) > n$

we want show, $t \rightarrow \infty$, $f(\theta(t), x_i) \rightarrow y_i$

Gradient Flow: a Kernel Point of View

$$u_i(t) = f(\theta(t), x_i) ,$$

$$u(t) = \begin{pmatrix} u_1(t) \\ \vdots \\ u_n(t) \end{pmatrix}$$

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

$$l'(u(t), y) \in \mathbb{R}^n = \left\langle \frac{\partial u_i(t)}{\partial \theta(t)} , -\frac{1}{n} \sum_{j=1}^n l'(u_j(t), y_j) \cdot \frac{\partial u_j(t)}{\partial \theta(t)} \right\rangle$$

$$[l'(u(t), y)]_i = -\frac{1}{n} [l'(u_1(t), y_1), \dots, l'(u_n(t), y_n)]_i$$

$$= l'(u_i(t), y_i)$$

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

$$\left(\left\langle \frac{\partial u_1(t)}{\partial \theta(t)} , \frac{\partial u_1(t)}{\partial \theta(t)} \right\rangle , \dots , \left\langle \frac{\partial u_n(t)}{\partial \theta(t)} , \frac{\partial u_n(t)}{\partial \theta(t)} \right\rangle \right)$$

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

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

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

Gradient Flow: a Kernel Point of View

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

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

$$\frac{d(u(t), y)}{dt} = -\frac{1}{\eta} H(t) (u(t) - y)$$

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

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

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

H p.d.

$\lambda_{\min}(H) \geq \lambda_0$
any vector \checkmark

$u^T H u$
 $\geq \lambda_0 \|u\|_2^2$

Gradient Flow: a Kernel Point of View

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

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

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

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

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

$$\log\left(\frac{1}{\epsilon}\right)$$

$$t \rightarrow \infty, \text{ loss} \rightarrow 0, u(t) - y$$

kernel: $w(t) = \phi(x_i)^T \theta(t)$

feature map for sample

$$[H(t)]_{ij} = \left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, \frac{\partial u_j(t)}{\partial \theta(t)} \right\rangle$$

$$= \langle \phi(x_i), \phi(x_j) \rangle$$

$$= k(x_i, x_j) \text{ does not depend on } t$$

$$H(t) = \begin{pmatrix} \ddots & \ddots & \ddots \\ \vdots & k(x_i, x_j) & \vdots \\ \ddots & \ddots & \ddots \end{pmatrix} \stackrel{\Delta}{=} K \in \mathbb{R}^{n \times n}$$

$$\lambda_{\min}(H(t)) > 0$$

for kernel $\lambda_{\min}(K) > 0 \Leftrightarrow K$ is full-rank

if kernel is universal: Gaussian, NTK
 $\Rightarrow K$ is full-rank

Gradient Flow: a Kernel Point of View

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

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

• Initialization: $a_r \sim \text{unit } \{1, -1\}$ for simplicity

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

• Training: only train w_1, \dots, w_m

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

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

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

H^* : NTK

Idea: $H(t)$ stays the same for $\forall t$

$$H(t) \approx H^*$$

$$H_{ij}^* = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{r=1}^m \left\langle \frac{\partial f(x_i, a, w_r)}{\partial \theta}, \frac{\partial f(x_j, a, w_r)}{\partial \theta} \right\rangle$$

Gradient Flow: a Kernel Point of View

$$H_{ij}'(t) = \left\langle \frac{\partial u_i(t)}{\partial w(t)}, \frac{\partial u_j(t)}{\partial w(t)} \right\rangle, \quad w \in \mathcal{D}^{m \times d}$$

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

$$\frac{\partial u_i(t)}{\partial w_r(t)} = \frac{1}{\sum_{k=1}^m a_k} a_r \cdot x_i = \frac{1}{\sum_{k=1}^m a_k} a_r x_i \mathbb{1}_{\{w_r^T x_i \geq 0\}}$$

$$H_{ij}'(t) = \sum_{r=1}^m \frac{1}{m} \left\langle a_r x_i \mathbb{1}_{\{w_r^T x_i \geq 0\}}, a_r x_j \mathbb{1}_{\{w_r^T x_j \geq 0\}} \right\rangle$$

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

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

Gradient Flow: a Kernel Point of View

Initialization!

Hoeffding inequality!

R.V. z_1, \dots, z_n i.i.d. \mathcal{U} , $|z_i| \leq 1$

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

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

$$H_{ij}(0) = X_i^T X_j \cdot \underbrace{\frac{1}{m} \sum_{l=1}^m}_{\text{average}} \underbrace{\mathbb{1}\{w_l(0)^T X_i \geq 0, w_l(0)^T X_j \geq 0\}}_{Z_l}$$

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

when m is large enough

$$|H_{ij}(0) - H_{ij}^*| \leq \epsilon$$

Gradient Flow: a Kernel Point of View

$$\begin{aligned} & \|H^* - H(\vartheta)\|_F \\ \leq & \sum_{i,j} |H_{ij}^* - H_{ij}(\vartheta)| \\ \leq & h^2 \cdot \Sigma \end{aligned}$$

$$\begin{aligned} a) \quad m & \rightarrow \infty \\ H(\vartheta) & \rightarrow H^* \end{aligned}$$

Gradient Flow: a Kernel Point of View

Want $H(t) \approx H(0)$

for simplicity

1) just train till time t

2) $y_i = 0$

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

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

Key idea every weight updates only m times
a little $O\left(\frac{1}{\sqrt{m}}\right)$: lazy training

Gradient Flow: a Kernel Point of View

$$\begin{aligned} & \|w_r(t) - w_r(0)\|_2 \\ &= \left\| \int_0^t \frac{dw_r(\tau)}{d\tau} d\tau \right\|_2 \\ &\leq \int_0^t \left\| \frac{dw_r(\tau)}{d\tau} \right\| d\tau \\ &= \int_0^t \left\| \frac{1}{\sqrt{m}} \frac{1}{n} \sum_{i=1}^n (y_i(\tau) - y_i) \cdot \underbrace{a_r X_i}_{\mathcal{O}(1)} \cdot \left\{ w_{r(0)}^\top X_i \geq \tau_0 \right\} \right\|_2 d\tau \\ &\leq C \cdot \int_0^t \frac{1}{\sqrt{m}} d\tau \\ &\leq \frac{C \cdot t}{\sqrt{m}} \end{aligned}$$

Gradient Flow: a Kernel Point of View

ReLU Smoothness

smoothness: small movement in parameter \Rightarrow small deviation in function's derivative

$$H_{ij}(t) = X_i^T X_j \frac{1}{m} \sum_{v=1}^m \mathbb{1} \{ w_v(t)^T X_i > 0, w_v(t)^T X_j > 0 \}$$

$$H_{ij}(0) = X_i^T X_j \frac{1}{m} \sum_{v=1}^m \mathbb{1} \{ w_v(0)^T X_i > 0, w_v(0)^T X_j > 0 \}$$

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

\Rightarrow want to bound # of pattern changes

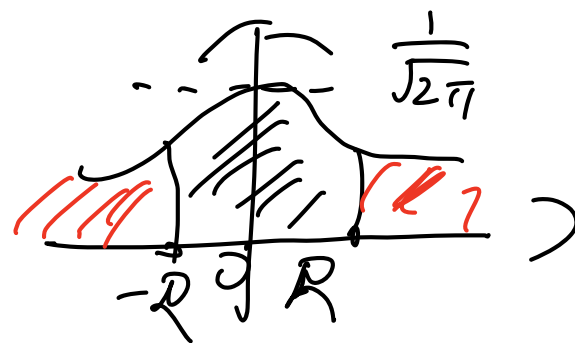
$$\frac{1}{m} \sum_{v=1}^m \mathbb{1} \{ \text{sgn}(w_v(t)^T X_i) \neq \text{sgn}(w_v(0)^T X_i) \}$$

Gradient Flow: a Kernel Point of View $\mathcal{N}(0, 1)$

Gaussian Anti-concentration

$$P(|Z| \leq R) \leq \frac{2R}{\sqrt{2\pi}}$$

$Z \sim \mathcal{N}(0, 1)$



$$W_v(0) \sim \mathcal{N}(0, I) \Rightarrow W_v(0)^T X_2 \sim \mathcal{N}(0, 1) \quad (\text{when } \|X_i\|_2 = 1)$$

If $\forall t, \|W_v(t) - W_v(0)\|_2 \leq \Delta W$
 ($\Delta W \rightarrow 0$ as $m \rightarrow \infty$, $\Delta W = O(\frac{1}{\sqrt{m}})$)

Let's choose $R > \Delta W$

Suppose $|W_v(0)^T X_i| \geq R$

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

$$\begin{aligned} & |W_v(0)^T X_i - W_v(t)^T X_i| \\ & \leq \|X_i\|_2 \cdot \|W_v(0) - W_v(t)\|_2 \\ & \leq R \leq \|W_v(0)^T X_i\| \end{aligned}$$

Gradient Flow: a Kernel Point of View

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

We know $\Delta u \rightarrow 0$ as $m \rightarrow \infty$

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1} \left\{ \text{sgn}(w_r(t)^T x_i) \neq \text{sgn}(w_r(0)^T x_i) \right\} \rightarrow 0$$

$$\|H(t) - H(0)\|_F \leq \epsilon, \quad m = \frac{u^6}{\epsilon^2}$$

$$(t(t) -) H^*$$

$$\frac{dU(t)}{dt} = -\frac{1}{n} H(t) (U(t) - y)$$

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

universal \rightarrow full rank

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

For predicting

$m \rightarrow \infty$, x_{te} test point

$f(\theta, x_{te}) \rightarrow$ kernel predictor

$$k_{te} (H^*)^T y, \quad y \in \mathcal{R}^y$$

$$H^*: \mathcal{N} \times \mathcal{X}$$

$$k_{te} \in \mathcal{R}^y, \quad \begin{pmatrix} K(x_1, x_{te}) \\ \vdots \\ K(x_n, x_{te}) \end{pmatrix}$$