

§3. Neural Tangent Kernel (NTK)

2. Themes of DL Theory.

• "Traditional" supervised learning:

- Why GD + NN can (over)fit the training set?

(Global convergence)

- What solutions can GD find? Why they can generalize?

(Algorithmic regularization.)

-

• Not so supervised. (Pre-training + finetuning).

- Different pre-training tasks. (contrastive learning - reconstruction - ...)

- What representations do they learn?

- Why they can be used in downstream tasks.

-

• In-context learning. (Prompting)

- We don't know what questions to ask yet - - -

• Generic template.

- Why X works?

- Why X is better than Y?

2. Background of NTK

- (Zhang et al., 2016). (experimental)

GD + NN can often globally minimize the loss, even when the labels are random.

- Over-parameterized networks.

- Original meaning: # params $> n$
 \uparrow # training samples

- What people mean in DL theory:

neurons (per layer) = poly(r) \leftarrow latent dim.
or poly(d) \leftarrow input dim.

or poly(d, n)

or exp(d) or ∞

3. Setup

- Training set: $\{(x_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^d \times \mathbb{R}$.

- Learner network.

$$f(x; W, a) = \frac{1}{\sqrt{m}} a^T \sigma(Wx) = \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k G(W_k \cdot x)$$

where $x \in \mathbb{R}^d$ input.

$m \in \mathbb{N}$. # neurons.

$W \in \mathbb{R}^{m \times d}$ first layer weights

$\sigma: \mathbb{R} \rightarrow \mathbb{R}$. ReLU.

$a \in \mathbb{R}^m$ output weights

- Initialization. $a_k \sim \text{Unif}(\{\pm 1\})$

$W_k \sim N(0, I_d)$.

- Training algorithm: Fix a . Train W with

GD + MSE. $L(W) = \frac{1}{2} \sum_{i=1}^n (y_i - f(x_i; W))^2$

4. Derivatives.

$$\begin{aligned}\nabla_{w_k} L(w) &= \sum_{i=1}^n (y_i - f(x_i; w)) \nabla_{w_k} f(x_i; w) \\ &= - \sum_{i=1}^n (y_i - f(x_i; w)) \frac{a_k}{\sqrt{m}} \nabla_{w_k} \sigma(w_k \cdot x_i) \\ &= - \frac{a_k}{\sqrt{m}} \sum_{i=1}^n (y_i - f(x_i; w)) \sigma'(w_k \cdot x_i) x_i\end{aligned}$$

Recall GF: $\frac{d}{dt} w_k = - \nabla_{w_k} L(w)$.

$$\begin{aligned}\frac{d}{dt} f(x_j) &= \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k \frac{d}{dt} \sigma(w_k \cdot x_j) \\ &= \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k \langle \sigma'(w_k \cdot x_j) x_j, \frac{d}{dt} w_k \rangle \\ &= \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k \langle \sigma'(w_k \cdot x_j) x_j, \frac{a_k}{\sqrt{m}} \sum_{i=1}^n (y_i - f(x_i; w)) \sigma'(w_k \cdot x_i) x_i \rangle \\ &= \frac{1}{m} \sum_{i=1}^n (y_i - f(x_i; w)) \sum_{k=1}^m a_k^2 \sigma'(w_k \cdot x_i) \sigma'(w_k \cdot x_j) \langle x_i, x_j \rangle\end{aligned}$$

Define the NTK,

$$H(x, x') = \frac{1}{m} \sum_{k=1}^m \sigma'(w_k \cdot x) \sigma'(w_k \cdot x') \langle x_i, x_j \rangle$$

and $H_{ij} = H(x_i, x_j)$, $i, j \in [n]$.

Define $y = (y_i)_{i=1}^n$, $f_t = (f(x_i; w_t))_{i=1}^n$.

Then $\frac{d}{dt} f_t = H_t (y - f_t)$.



Remarks.

a) H_t depends on t .

b) NTK and these formula themselves are quite generic, while the NTK technique is not.

c) If $H_t \succeq \lambda_0 \text{Id}$ for some $\lambda_0 > 0$, $\forall t$, then $f_t \rightarrow y$ linearly.

5. The NTK technique.

Choose m and the initialization scale (im)properly to ensure $H_t \approx H_0$ throughout training, and reduce the problem to a convex one.

6. Define $H^\infty \in \mathbb{R}^{d \times d}$ by

$$\begin{aligned} H_{i,j}^\infty &= \lim_{m \rightarrow \infty} H_{i,j}^{(m)} \\ &= \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{k=1}^m \sigma'(W_k \cdot x_i) \sigma'(W_k \cdot x_j) \langle x_i, x_j \rangle \\ &= \mathbb{E}_{W \sim \mathcal{N}(0, 2d)} \left[\sigma'(W \cdot x_i) \sigma'(W \cdot x_j) \langle x_i, x_j \rangle \right]. \end{aligned}$$

Define $\lambda_0 := \lambda_{\min}(H^\infty)$.

Fact. If $x_i \not\propto x_j$, $\forall i \neq j$, then $\lambda_0 > 0$.

Assume, $\lambda_0 > 0$.

{ Du et al., 2019. Gradient descent
{ provably optimizes over-parametrized
{ neural networks
{ -----

7. Lemma. Choose $m = \Omega\left(\frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)$.

With proba. $\geq 1 - \delta$, we have

$$\|H_0 - H_0^w\|_2 \leq \lambda_0/4, \text{ whence } \lambda_{\min}(H_0) \geq \frac{3}{4}\lambda_0.$$

Lemma. Initialization w_1, \dots, w_m .

Any $\hat{w}_1, \dots, \hat{w}_m \in \mathbb{R}^d$ with $\|w_k - \hat{w}_k\| \leq \mathcal{O}\left(\frac{\delta \lambda_0}{n^2}\right) =: R$.

Let H and \hat{H} be the corresponding NTKs. With proba. $\geq 1 - \delta$ over the random init., we have

$$\|\hat{H} - H\|_2 \leq \lambda_0/4, \text{ whence } \lambda_{\min}(\hat{H}) \geq \lambda_0/2.$$

Deterministic.

The infinite-width limit.

Depends only on the initialization scheme

H^∞

Depends only on the actual $\rightarrow H_0$ random initialization.

H_{tr} \leftarrow Also depends on the training procedure.

8. Thm. Assume $\|x_i\|, |y_i| \leq C$. Choose $m = \Omega\left(\frac{n^6}{\lambda_0^4}\right)$.
 With prob. $\geq 1 - O(\delta)$, GF converges to a point with
 $\|y - f_+\| \leq \varepsilon$ within $O\left(\frac{1}{\lambda_0} \log\left(\frac{\|y - f_+\|}{\varepsilon}\right)\right)$ amount of time.

In words, GF will fit the training set before
 W_k moves too far away from the initialization.

Pf. Define $T_* := \min\{T_1, T_2\}$ where

$$T_1 := \inf\{t \geq 0 : \|f_t - y\| \leq \varepsilon\}$$

$$T_2 := \inf\{t \geq 0 : \exists k, \|W_k(t) - W_k(0)\| \geq R\}$$

By the previous lemmas, $\forall t \leq T_* \leq T_2$, we have

$$\begin{aligned} \frac{d}{dt} \|y - f_t\|^2 &= -2 \langle y - f_t, H_t(y - f_t) \rangle \\ &\leq -\lambda_0 \|y - f_t\|^2. \end{aligned}$$

$$\Rightarrow \|y - f_t\|^2 \leq \|y - f_0\|^2 \exp(-\lambda_0 t).$$

$$\Rightarrow \|y - f_t\| \leq \|y - f_0\| \exp(-\lambda_0 t/2)$$

Meanwhile, we have

$$\| \nabla_{W_k} L(W) \| = \left\| \frac{g_k}{\sqrt{m}} \sum_{i=1}^n (y_i - f(x_i)) \sigma'(W_k \cdot x) x \right\|$$

$$\leq \frac{C}{\sqrt{m}} \sum_{i=1}^n |y_i - f(x_i)|$$

$$\leq C \sqrt{\frac{y}{m}} \|y - f_+\|$$

$$\Rightarrow \|W_k(t) - W_k(0)\| \leq C \sqrt{\frac{y}{m}} \int_0^t \|y - f_s\| ds$$

$$\leq C \sqrt{\frac{y}{m}} \|y - f_0\| \int_0^t \exp(-\lambda_0 s/2) ds$$

$$\leq \underbrace{2C \sqrt{\frac{y}{m}} \|y - f_0\|}_{\leq R}$$

$\Rightarrow T_*$ cannot be attained by T_2 .

$$\Rightarrow T_* = T_1 \leq O\left(\frac{1}{\lambda_0} \log\left(\frac{\|y - f_0\|}{\varepsilon}\right)\right).$$

Question. Where did we use $a_k \sim \text{Unif}\{\pm 1\}$?

$$\begin{aligned} \text{Answer. } \mathbb{E}_{a, W} \|f_0\|^2 &= \mathbb{E}_{a, W} \left\| \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k \sigma(W_k \cdot X) \right\|^2 \\ &= \frac{1}{m} \sum_{i, j=1}^m \mathbb{E}_{a, W} \left\{ a_i a_j \sigma(W_i \cdot X) \sigma(W_j \cdot X) \right\} \end{aligned}$$

1) If $a_k \sim \text{Unif}\{\pm 1\}$, then

$$= \frac{1}{m} \sum_{k=1}^m \mathbb{E}_W \sigma^2(W_k \cdot X) = \mathbb{E}_W \sigma^2(W \cdot X)$$

2) If $a_k = 1$, then

$$\begin{aligned} &= \frac{1}{m} \sum_{i, j=1}^m \mathbb{E}_W \sigma(W_i \cdot X) \sigma(W_j \cdot X) \\ &\geq \Omega \left(m \left(\mathbb{E}_W \sigma(W_k \cdot X) \right)^2 \right) \end{aligned}$$

As a result -

$$\begin{aligned} \|W_k(t) - W_k(s)\| &\leq 2 \left(\sqrt{\frac{d}{m}} \|y - f_0\| \right) \\ &\rightarrow 0 \text{ as } m \rightarrow \infty. \end{aligned}$$

9. NTK \Leftrightarrow random feature.

If the movement is small.

$$f(\pi; W(t)) = \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k \sigma(W_k \cdot X)$$

$$\approx \underbrace{f(\pi; W(0))} + \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k \langle W_k(t) - W_k(0), \pi \rangle \sigma'(W_k \cdot X)$$

$\hookrightarrow \approx 0$ when m is large.

Define $\alpha_k = (W_k(t) - W_k(0)) / \sqrt{m}$. Then

$$f(\pi; W(t)) \approx \sum_{k=1}^m \langle \alpha_k, a_k \sigma'(W_k(0) \cdot X) X \rangle.$$

\hookrightarrow a linear model over the
fixed random feature mapping

$$\pi \mapsto \left(a_k \sigma'(W_k(0) \cdot X) \pi \right)_{k=1}^m$$

10. Random linear features const
learn a single linear model

• Input distr. $x \sim N(0, I_d)$

• Target function, $f_*(x) = \langle W_*, x \rangle$,

for some fixed unit vector W_* .

• Learner: $f(x; u) = \sum_{k=1}^m u_k \langle W_k, x \rangle$

$$= \left\langle \underbrace{\sum_{k=1}^m u_k W_k}_{=: W_u}, x \right\rangle$$

where $W_k \sim N(0, I_d/d)$

Claim. If $m \leq d/2$, then w.h.p.

$$\min_u \text{MSE} = \min_u \mathbb{E}_x \left\{ (f(x; u) - f_*(x))^2 \right\} \geq \frac{1}{4}$$

Pf. $\text{MSE} = \mathbb{E}_x \langle W_u - W_*, x \rangle^2 = \|W_u - W_*\|^2$

$$\Rightarrow \min_u \text{MSE} = \text{dist. from } W_* \text{ to } \text{span}(W_1, \dots, W_m)$$

By symmetry, we may assume w.l.o.g. that (W_1, \dots, W_m) are fixed and W_* is a random unit vector. Moreover assume $W_k = e_k$.

$$\Rightarrow L(W_*) := \min_u \text{MSE} = 1 - \sum_{k=1}^m [W_*]_k^2$$
$$\geq 1 - \sum_{k=1}^{d/2} [W_*]_k^2$$

\hookrightarrow (GC) -Hölder with median = 1/2

\implies with proba. $\geq 1 - \exp(-\text{GC}(d))$,

$$L(W_*) \geq 1/4.$$

Lemma. (Lévy's inequality). Let $V \sim \text{Unif}(S^{d-1})$.

$g: \mathbb{R}^d \rightarrow \mathbb{R}$. L -Lipschitz. $\exists C, c > 0$ st. $\forall \varepsilon$.

$$\mathbb{P}[|g(V) - \text{median}(g)| \geq \varepsilon] \leq C \exp\left(-\frac{c\varepsilon^2 d}{L^2}\right).$$

• Informed version of Thm. 1.2 of Thm 4.2 of Lehtonen and Shamir (2019).

$$\mathcal{F} = \left\{ x \mapsto f(Wx) : W = \begin{bmatrix} W_1 \\ \vdots \\ W_m \end{bmatrix}, W_k \sim \text{Unif}(S^{d-1}) \right\}, \quad N(x) = \sum_{k=1}^m u_k f(W_k x)$$

$\forall W^* \in \mathbb{R}^d$ with $\|W^*\| = d^2$. \exists bias $b_x \in \mathbb{R}$ with $|b_x| = O(d^3)$ st. u.h.p.

$$\mathbb{E} \left\{ (N(x) - \sigma(W^* x + b_x))^2 \right\} \leq 1/50.$$

$$\Rightarrow m \cdot \max_k |u_k| \geq \exp(\Omega(d)).$$

11. Continuity argument etc. induction/minimal counterexample in continuous time

• Example $f \in C^0(\mathbb{R})$. $\dot{x}_t = f(x_t)$. $f(0) > 0$. $x_0 > 0 \Rightarrow x_t > 0 \forall t$.

• Example. Comparative Grönwall.

"Progress" X_t . "error" Y_t . Suppose when $Y_t \leq \delta = \frac{1}{\text{poly}(d)}$, we have

$$\begin{cases} \dot{X}_t \leq -A_t X_t \\ \dot{Y}_t \leq K A_t Y_t \end{cases} \quad \sim \quad \begin{cases} \dot{X}_t \leq -A_t X_t \\ \dot{Y}_t \leq K A_t Y_t + B \end{cases}$$

• Example. Suppose when $Y_t \leq \delta \leq 1/\text{poly}(d)$, we have $\dot{X}_t \leq -A_t X_t$. $\dot{Y}_t \leq X_t Y_t$.

*. Present this before the NTK results.

§4. Mean-Field Networks.

1. Recall $f(x; W, a) = \frac{1}{m} a^T \sigma(Wx) = \frac{1}{m} \sum_{k=1}^m a_k \sigma(w_k \cdot x)$.

Let $\mu = \frac{1}{m} \sum_{k=1}^m \delta_{(a_k, w_k)} \leftarrow$ the empirical distr. of the neurons.

$$f(x; W, a) = \int a \sigma(w \cdot x) d\mu(a, w) =: f(x; \mu)$$

Allow μ to be any (reasonably regular) distribution on \mathbb{R}^{1+d}

\Rightarrow mean-field networks.

MSE: $L(\mu) = \frac{1}{2} \mathbb{E}_{x, y} \left\{ (y - f(x; \mu))^2 \right\}$

\hookrightarrow a functional over probability measures.

GF of L ?

2. The vector space structure

Notations. $\mathcal{M}_{\pm}(\mathbb{R}^d) = \{ \text{signed measures over } \mathbb{R}^d \}$

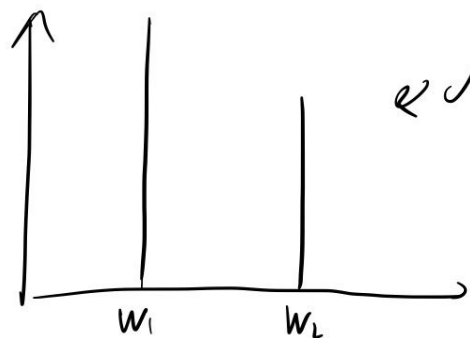
$$\mathcal{P}(\mathbb{R}^d) = \{ \text{probability measures } \sim \}, \quad \mathcal{P}_2(\mathbb{R}^d) = \{ \mu \in \mathcal{P}(\mathbb{R}^d) : \int \|x\|^2 d\mu < \infty \}$$

The vector space structure: $\mu, \nu \in \mathcal{M}_{\pm}(\mathbb{R}^d), a, b \in \mathbb{R}$

measurable E $(a\mu \oplus b\nu)(E) = a\mu(E) + b\nu(E)$.

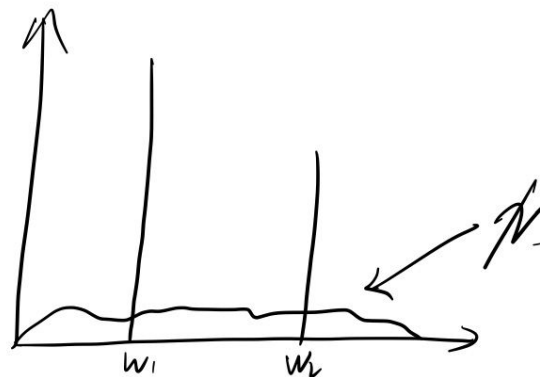
Let $\mathcal{P}(\mathbb{R}^d)$ inherit this structure.

Q: Is this the "correct" structure?

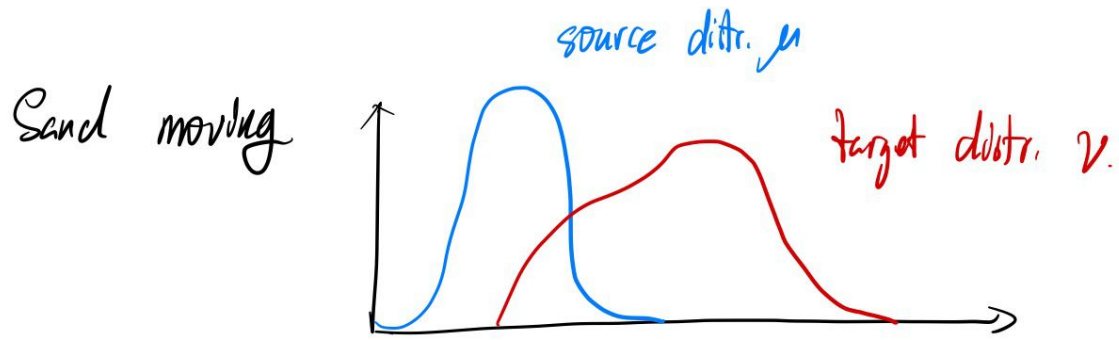


$$\leftarrow \mu = \frac{1}{2} \delta_{w_1} + \frac{1}{2} \delta_{w_2}$$

add a "small" perturbation χ



3. Wasserstein-2 space W_2



cost of moving 1 unit of sand from x to y : $\|x - y\|^2$

Goal: minimize the total cost.

$$\underset{T}{\text{minimize}} \int \|x - T(x)\|^2 d\mu(x) \quad \text{s.t.} \quad T\#\mu = \nu.$$

Def. For two (sufficiently regular) probability measure $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$, the Wasserstein-2 distance between μ and ν is defined as.

$$W_2^2(\mu, \nu) = \inf \left\{ \int \|x - T(x)\|^2 d\mu(x) : T\#\mu = \nu \right\}$$

Remark. Use couplings instead of transport maps for less regular distributions.

transport map

Def. Given $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$

$$T\#\mu(E) = \mu(T^{-1}(E)).$$

Thm. $W_2(\mathbb{R}^d) := (P_2(\mathbb{R}^d), W_2)$ is a metric space

the constant speed geodesic between μ and ν is given by

$$t \mapsto ((1-t)I_d + tT) \# \mu, \quad t \in [0, 1],$$

where T is the optimal transport map.

4. Wasserstein gradient flow

Remark. It's possible to define GF in general metric spaces. but we'll focus on W_2 .

Def. $F: P(\mathbb{R}^d) \rightarrow \mathbb{R}$, $\mu \in P(\mathbb{R}^d)$. We say $G: \mathbb{R}^d \rightarrow \mathbb{R}$ is the first variation of F at μ if $\forall \nu \in \mathcal{M}_+(\mathbb{R}^d)$ with $\mu + \varepsilon\nu \in P(\mathbb{R}^d) \quad \forall \varepsilon > 0$, we have

$$\left. \frac{d}{d\varepsilon} \right|_{\varepsilon=0} F(\mu + \varepsilon\nu) = \int G(x) d\nu(x) \rightarrow =: \frac{\delta F}{\delta \mu}[\mu]$$

Remark. 1) directional derivative. 2) local linear approximation.

Example, 1) $F(\mu) = \int f(x) d\mu(x) \Rightarrow \frac{\delta F}{\delta \mu[\mu]}(x) = f(x)$

2) $F(\mu) = \iint W(x,y) d\mu(x) d\mu(y)$. W symmetric $\Rightarrow \frac{\delta F}{\delta \mu[\mu]}(x) = 2 \int W(x,y) d\mu(y)$.

Q: How to minimize $F(\mu) = \int G(x) d\mu(x)$
 \searrow cost of placing 1 unit of particles at x .

A: At each x , move the particles along $-\nabla G(x)$.

$$\Rightarrow \frac{d}{dt} \mu_t = -\nabla G(\mu_t) = -\nabla \frac{\delta F}{\delta \mu[\mu]}(\mu_t) \quad \forall \mu_t \in \text{supp}(\mu_t) \quad (**)$$

$$\Leftrightarrow \partial_t \mu_t - \nabla \cdot (\mu_t \nabla \frac{\delta F}{\delta \mu[\mu]}) = 0.$$

Def. $F: \mathcal{M}_2 \rightarrow \mathbb{R}$. We say μ_t is the Wasserstein gradient flow of F if it satisfies the continuity equation $(*)$ or $(**)$.