

Lecture 5:

Ann.

- * TA OH
- * hwl Q?
- * typed notes

Lec 1-4: constructive approx
 Lec 5-8: near initialization / over parameterization

Lec 1-4: care about # nodes to approx

Remark: training one ReLU has negative results & many papers.

GD seems to care about two function classes:

- ① Near initialization
- ② low norm.

Contradictory: ① is close (large) Gaussian initialization
 ② measured wrt origin.

Next few lectures: ① (aka NTK) in these steps:

① Near initialization \Rightarrow near Taylor expansion around initialization

[critical point: scaling of error with m .]

② Characterize Taylor expansion as $m \rightarrow \infty$.
 Remark: This will reveal a "signal/noise" separation; this lets us choose our scaling constants

③ kernel perspective ("neural tangent kernel").

④ Refined complexity estimates.

Remark: On one hand, Taylor expansion not a good approximator of network in practice; on the other hand, this perspective has good predictive power.

Definition (NTK lectures)

$$F_0(x; w) := F(x; w_0) + \left\langle \frac{\partial}{\partial w} F(x; w_0), w - w_0 \right\rangle$$

minimum norm element of Clarke differential

Moreover we will typically only vary V (where $w = (a, V)$), so that optimization problem is "convex".

E.g. $F(x; (a, V)) = \sum a_j \sigma(v_j^T x)$; $F_0(x; (a, V)) := \sum a_{0,j} \sigma(v_{0,j}^T x) + \sum_j a_{0,j} \sigma'(v_{0,j}^T x) x^T (v_j - v_{0,j})$

Remark (Initialization) * These lectures: $a_j \sim \text{Unif}(\{-1, +1\})$, $v_j \sim \mathcal{N}(0, I_d)$

* Literature, (NTK) ① $a_j \sim \text{Unif}(\{\frac{-1}{\sqrt{m}}, \frac{+1}{\sqrt{m}}\})$, $v_j \sim \mathcal{N}(0, I_d)$

② $a_j \sim \text{Unif}(\{\frac{-\epsilon}{\sqrt{m}}, \frac{\epsilon}{\sqrt{m}}\})$, $v_j \sim \mathcal{N}(0, I_d)$

③ $a_j =$ some deterministic choices so that $F(x; w_0) = 0$.
 "symmetric"

* "Pytorch": ① $a_j \sim \text{Unif}(\{-\frac{1}{\sqrt{m}}, \frac{+1}{\sqrt{m}}\})$, $v_{j,i} \sim \text{Unif}(\{\frac{-1}{\sqrt{2}}, \frac{+1}{\sqrt{2}}\})$

(TensorFlow different).

Open: Uniform vs Gaussian

Another scaling: mean-field

(A) Near initialization \Rightarrow Near Taylor expansion

Warm-up: "smooth" activations.

Proposition. Suppose σ is β -smooth ($|\sigma'(r) - \sigma'(s)| \leq \beta|r-s|$).

Then, given any $V \in \mathbb{R}^{m \times d}$,

$$|F(x; V) - F_0(x; V)| \leq \frac{\beta \|a\|_{\infty} \|x\|^2 \|V-V_0\|^2}{2}$$

Remark. (a) open to prove matching for the ReLU. (c) no m in rhs.
 (b) this theorem has no randomness.

Proof. Note

$$\begin{aligned} & \left| \sigma(r) - (\sigma(s) + \sigma'(s)(r-s)) \right| \\ &= \left| \int_s^r \sigma'(t) dt - \int_s^r \sigma'(s) dt \right| \\ &\leq \int_s^r |\sigma'(t) - \sigma'(s)| dt \leq \frac{\beta(r-s)^2}{2}. \end{aligned}$$

Therefore

$$\begin{aligned} |F(x; V) - F_0(x; V)| &= \left| \sum_{j=1}^m a_j (\sigma(v_j^T x) - (\sigma(v_{0,j}^T x) + \sigma'(v_{0,j}^T x) \dots \langle x, v_j - v_{0,j} \rangle)) \right| \\ &\leq \|a\|_{\infty} \sum_{j=1}^m |\sigma(v_j^T x) - \dots| \\ &\leq \|a\|_{\infty} \sum_{j=1}^m \frac{\beta}{2} (v_j^T x - v_{0,j}^T x)^2 \leq \|a\|_{\infty} \|x\|^2 \cdot \frac{\beta}{2} \sum_j \|v_j - v_{0,j}\|^2 \end{aligned}$$

Let's try this for ReLU.

$$|F(x; V) - F_0(x; V)| =$$

$$\left| \sum_j a_j \left(\underbrace{\sigma(v_j^T x)}_{v_j^T x \sigma'(v_j^T x)} - \underbrace{(\sigma(v_{0,j}^T x))}_3 + \underbrace{\sigma'(v_{0,j}^T x)}_2 (v_j^T x - \underbrace{v_{0,j}^T x}_1) \right) \right|$$

$$= \left| \sum_j a_j v_j^T x \left(\sigma'(v_j^T x) - \sigma'(v_{0,j}^T x) \right) \right|$$

Handled via
 Gaussian concentration

$\leq ???$

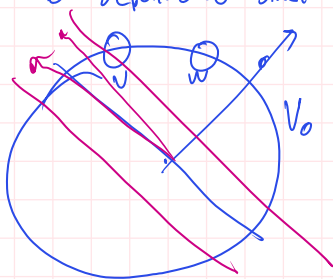
$$\begin{aligned} &\leq ?? \leq \|a\|_{\infty} \sum_j |v_j^T x| \cdot \underbrace{|\sigma'(v_j^T x) - \sigma'(v_{0,j}^T x)|}_{\leq 1} \\ &\leq \|a\|_{\infty} \sum_j \|w_j\| \leq \|a\|_{\infty} \sqrt{m} \|V\|_F \end{aligned}$$

Lemma (ReLU linearization). For any $B > 0$ and any $\|x\| \leq 1$,
 with probability $\geq 1 - \delta$ (over V_0), for any $\|W - V_0\| \leq B$, $\|V - V_0\| \leq B$,

$$|F(x; V) - (F(x; W) + \langle \bar{\sigma}' F(x; W), V - W \rangle)| \leq \|x\|_{\infty} m^{1/3} (4B^{4/3} + 2B \ln(1/\delta)^{1/4}).$$

Remark. (i) Maybe don't need $m^{1/3}$, perhaps if do two layers (or other change).
 (ii) vs smooth theorem: B dependence smaller (B^2), less probability.

Why proof works



Lemma. For any $\tau > 0$, $x \in \mathbb{R}^d$, $\|x\| \geq 0$,
 with $pr \geq 1 - \delta$

$$\sum_j \mathbb{1} [|v_j^T x| \leq \tau \|x\|] \leq m \tau + \sqrt{\frac{m}{2} \ln(1/\delta)}.$$

Proof. rotational invariance, explicit integral of Gaussian density
 & Hoeffding