

Lecture 6: NTK continued.

$$f(x; W) := \frac{1}{\sqrt{m}} \sum_{j=1}^m a_j \sigma(w_j^T x)$$

$\uparrow \quad \uparrow$
 $\pm 1 \quad \text{Gaussian}$

$$f_0(x; W) := f(x; W_0) + \langle \nabla f(x; W_0), W - W_0 \rangle$$

NTK story:

near initialization, with large width,
 $f \approx f_0$, f_0 still universal approximator.

"kernel"

* In ML, whenever we have
"linear" (affine) predictor over some
feature mapping $x \mapsto \Phi(x)$,

we consider "kernel" $k(x, x') := \langle \Phi(x), \Phi(x') \rangle$.

(for us, $\Phi_m(x) \stackrel{\approx}{=} \nabla f(x; W_0)$.)
 \uparrow r.v.

$$k_m(x, x') := \langle \nabla f(x; W_0), \nabla f(x'; W_0) \rangle$$

$m \rightarrow \infty$ \rightarrow # problem 3 in hwk

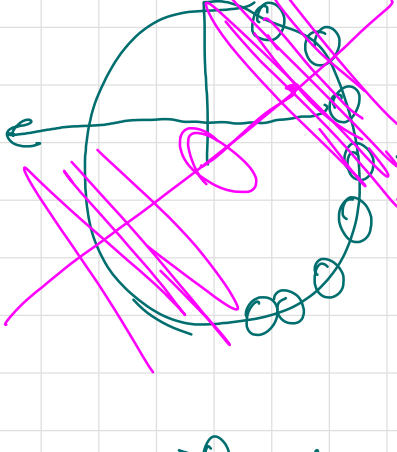
$f_x f_0$ near init for large width

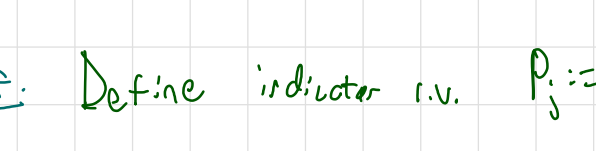
Prop. $\|\sigma''\| \leq \beta$, $\|k\| \leq 1$, $|a_j| \leq 1$, $\forall W$
 $|f(x; W) - f_0(x; W)| \leq \frac{\beta}{2\sqrt{m}} \|W - W_0\|_F^2$

Lemma. Let $B \geq 0$ be given,
 & fixed $\|k\| \leq 1$.
 With probability $\geq 1 - \delta$ (over W_0)
 for any $W \in \mathbb{R}^d$ with $\|W - W_0\|_F \leq B$,
 $|f(x; W) - f_0(x; W)| \leq \frac{2B^{4/3} + B \ln(1/\delta)^{1/4}}{m^{1/6}}$

Remark. Brake force: $(ReLU)$
 $\sigma(z) = z \sigma'(z)$
 $|f(x; W) - f_0(x; W)|$
 $= \left| \frac{1}{\sqrt{m}} \sum_j a_j (\sigma(w_j^T x) - \sigma(w_{0,j}^T x) - \sigma'(w_{0,j}^T x) x^T (w_j - w_{0,j})) \right|$
 $= \left| \frac{1}{\sqrt{m}} \sum_j a_j (\sigma'(w_j^T x) - \sigma'(w_{0,j}^T x)) w_j^T x \right|$
 $\leq \frac{1}{\sqrt{m}} \sum_j |a_j| \cdot |\sigma'(w_j^T x) - \sigma'(w_{0,j}^T x)| \cdot |w_j^T x|$
 $\leq \frac{1}{\sqrt{m}} \sum_j \|w_j\| \leq \frac{1}{\sqrt{m}} \sqrt{m} \cdot \sqrt{\sum_j \|w_j\|^2} = \|W\|_F$

Remark. First such bound was
 Allen-Zhu / Li / Song '18
 (multi-layers; this version one otherwise.)

Proof idea

 keep $\|W - W_0\|_F$ fixed,
 & increase m :
 $\sigma(w_j^T x) \approx \sigma(w_{0,j}^T x)$
 for "most" j

Lemma. For any $\tau > 0$, & any $x \neq 0$,
 with probability $\geq 1 - \delta$,
 $\sum_j \mathbb{1}[\|w_{0,j}^T x\| \leq \tau \|k\|] \leq \tau m + \sqrt{\frac{m}{2} \ln \frac{1}{\delta}}$


Proof. Define indicator r.v. $P_j := \mathbb{1}[\|w_{0,j}^T x\| \leq \tau \|k\|]$
 By rotational invariance of $w_{0,j}$, $P_j := \mathbb{1}[|g| \leq \tau]$
 where g is a univariate gaussian.

$$\mathbb{E} P_j = \int_{-\tau}^{\tau} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \leq \int_{-\tau}^{\tau} \frac{1}{\sqrt{2\pi}} dx = \frac{2\tau}{\sqrt{2\pi}} \leq \tau$$

By Hoeffding's inequality, w/ $\geq 1 - \delta$
 $\sum_j P_j \leq m \cdot \mathbb{E} P_j + \sqrt{\frac{m}{2} \ln \frac{1}{\delta}} \leq \tau m + \sqrt{\frac{m}{2} \ln \frac{1}{\delta}}$

Proof of $f_x f_0$ for ReLU.

If $x = 0$, $f(x; W) = 0 = f_0(x; W) \forall W$,
 proof is complete; henceforth, $\|k\| > 0$.

Let W be given with $\|W - W_0\|_F \leq B$,
 $S_1 := \{j \in \{1, \dots, m\} : |w_{0,j}^T x| \leq r \cdot \|k\|\}$
 $S_2 := \{j \in \{1, \dots, m\} : \|w_j - w_{0,j}\| \geq r\}$
 $S := S_1 \cup S_2$

Where r is a free parameter optimized shortly.

By the Hoeffding proof above,
 $|S_1| \leq r m + \sqrt{m \ln \frac{1}{\delta}}$

For $|S_2|$,
 $B^2 \geq \|W - W_0\|_F^2 = \sum_{j=1}^m \|w_j - w_{0,j}\|^2 \geq \sum_{j \in S_2} \|w_j - w_{0,j}\|^2 \geq r^2 |S_2|$
 $\Rightarrow |S_2| \leq \frac{B^2}{r^2}$

For any $j \notin S$, if $w_{0,j}^T x \geq 0$ (is analogous)
 $w_j^T x = w_{0,j}^T x + (w_j - w_{0,j})^T x$
 $\geq r \cdot \|k\| - \|w_j - w_{0,j}\| \cdot \|k\| > \|k\| (r - r) = 0$
 $\Rightarrow \text{sgn}(w_{0,j}^T x) = \text{sgn}(w_j^T x)$

Summarizing, S satisfies
 $|S| \leq r m + \sqrt{m \ln \frac{1}{\delta}} + \frac{B^2}{r^2}$ and $j \notin S \Rightarrow \begin{cases} \text{sgn}(w_j^T x) \\ = \text{sgn}(w_{0,j}^T x) \end{cases}$
 $r := \frac{B^{2/3}}{m^{1/3}} \Rightarrow \frac{B^2}{r^2} = (Bm)^{2/3} + \sqrt{m \ln \frac{1}{\delta}} + (Bm)^{2/3} \Rightarrow 2(Bm)^{2/3} + \sqrt{m \ln \frac{1}{\delta}}$

$$|f(x; W) - f_0(x; W)| \leq \frac{1}{\sqrt{m}} \sum_j |a_j| \cdot |\sigma'(w_j^T x) - \sigma'(w_{0,j}^T x)| \cdot |w_j^T x|$$

$$\leq \frac{1}{\sqrt{m}} \sum_{j \in S} |w_j^T x| \mathbb{1}[\text{sgn}(w_j^T x) \neq \text{sgn}(w_{0,j}^T x)]$$

$$\leq \frac{1}{\sqrt{m}} \sum_{j \in S} |w_j^T x - w_{0,j}^T x| \cdot \mathbb{1}[\text{sgn}(w_j^T x) \neq \text{sgn}(w_{0,j}^T x)]$$

$$\leq \frac{1}{\sqrt{m}} \sum_{j \in S} \|w_j - w_{0,j}\| \cdot \|k\|$$

$$\leq \frac{1}{\sqrt{m}} \sqrt{|S|} \cdot \|W - W_0\|_F$$

○ f file rows

if $\text{sgn}(w_j^T x) \neq \text{sgn}(w_{0,j}^T x)$

$$\Rightarrow |w_j^T x - w_{0,j}^T x| = |w_j^T x| + |w_{0,j}^T x| \geq |w_j^T x|$$

$$\Rightarrow \sum_{j \in S} |w_j^T x| \cdot \mathbb{1}(\text{sgn}(w_j^T x) \neq \text{sgn}(w_{0,j}^T x)) \leq \sum_{j \in S} |w_j^T x - w_{0,j}^T x| \cdot \mathbb{1}(\dots)$$

$$\sum_{j \in S} \|w_j\| \cdot \|x\|$$

$$\leq \sqrt{|S|} \cdot \|W\|_F$$

$$\mathbb{E} \|W_0\|_F^2 = \sum_{j=1}^m \sum_{i=1}^d w_{0,j,i}^2 = md$$

$$\frac{1}{\sqrt{m}} \cdot \sqrt{md} \cdot md$$

$\frac{1}{m^{1/2}}$