# Lecture 15: Linearized networks

## NTK analysis:

in fact we'll prove a generic result that doesn't even need the assumption that the student network is a one-hidden-layer MLP.

Consider a dataset $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$ and a student network $f_\theta: \mathbb{R}^d \to \mathbb{R}$, $\theta \in \mathbb{R}^p$, initialized to some $\theta_0$

We'll use shorthand $f_\theta(x) - y$ to denote

$$\left( f_\theta(x_1) - y_1, \ldots, f_\theta(x_n) - y_n \right).$$

Define: scaling param: $\alpha > 0$

empirical loss: $\hat{L}(g) \triangleq \frac{1}{2} \| g(x) - y \|_2^2$

$$\hat{L}_0 \triangleq \hat{L}(\gamma f_{\theta_0})$$

gradient flow:

$$d\theta_t \triangleq - \nabla_\theta \hat{L}(\gamma f_{\theta_t}) \, dt$$

$$= -\gamma J_t^T \nabla \hat{L}(\gamma f_{\theta_t}) \, dt, \quad \text{where}$$

Jacobian: $J_t \triangleq J_{\theta_t} \triangleq \begin{pmatrix} - \nabla_\theta f_{\theta_t}(x_1) - \\ \vdots \\ - \nabla_\theta f_{\theta_t}(x_n) - \end{pmatrix} \in \mathbb{R}^{n \times p}$

Will compare to linearized network / dynamics:

$$f_\Theta^{lin}(x) = f_{\Theta_0}(x) + J_0 \cdot (\Theta - \Theta_0)$$

$$d\breve{\Theta}_+ \overset{\triangle}{=} -\nabla_\Theta \hat{L}(\gamma f_{\breve{\Theta}_+}^{lin})$$

$$= -\gamma \boxed{J_0^T} \nabla \hat{L}(\gamma f_{\breve{\Theta}_+}^{lin})$$

Jacobian does not change for linearized network

Will assume that

1). $J_\Theta$ is Lipschitz in $\Theta$, i.e.

$$\|J_\Theta - J_{\Theta'}\|_{op} \leq \beta \|\Theta - \Theta'\|_2$$

2) $J_0 = J_{\Theta_0}$ is full rank (bounds will depend on $\sigma_{min}$, $\sigma_{max}$ of $J_0$)

Linearized dynamics very easy to analyze:

**Lemma 1:** If $Q(t) \succeq \lambda \cdot Id_n \; \forall \, t$, then for $(g_+)$ given by

$$dg_t = -Q(t) \nabla \hat{L}(g_+) \, dt,$$

we have

$$\hat{L}(g_+) \leq \hat{L}(g_0) \cdot \exp(-2\lambda t).$$

pf: $\frac{d}{dt} \hat{L}(g_t) = \langle -Q(t)(g_t(x)-y), g_t(x)-y \rangle$

$$\leq -\lambda \|g_t(x)-y\|_2^2$$

$$= -2\lambda \cdot \hat{L}(g_t),$$

So integrating this ( i.e. using Grönwall's inequality) completes the proof. □

Can apply this to $Q(t) = J_t J_t^T$ and $g_t = \gamma f_{\hat{\theta}_t}^{lin}$.

Then b/c

$$d\tilde{\Theta}_t = -\gamma J_0^T \nabla_\theta \hat{L}(\gamma f_{\hat{\theta}_t}^{lin}) dt,$$

by chain rule,

$$\frac{d}{dt}\left(\gamma f_{\hat{\theta}_t}^{lin}\right) = \gamma \underbrace{\nabla_\theta f_\theta^{lin}\Big|_{\theta=\tilde{\theta}_t}}_{J_0} \cdot \frac{d\tilde{\Theta}_t}{dt}$$

$$= -\gamma^2 \underbrace{J_0 J_0^T}_{Q(t) \geq \sigma_{min}^2(J_0) \cdot Id} \nabla_\theta \hat{L}(\gamma f_{\hat{\theta}_t}^{lin})$$

So by Lemma, $\hat{L}(\gamma f_{\hat{\theta}_t}^{lin}) \leq \exp\left(-2\gamma^2 t \sigma_{min}^2(J_0)\right)$

So (unsurprisingly), training loss for linearized network drops exponentially quickly.

Can also show that, relative to drop in loss, movement of parameters is negligible:

**Lemma 2:** Suppose process $(\hat{\Theta}_t)$ satisfies

$$d\hat{\Theta}_t = - S(t)^T \nabla \hat{L}(g_{\hat{\Theta}_t})$$

for some network $g$, and $\underline{\lambda} \cdot \text{Id} \preceq S(t)S(t)^T \preceq \bar{\lambda} \cdot \text{Id} \ \forall t$. Then

$$\| \hat{\Theta}_t - \hat{\Theta}_0 \| \leq \frac{\sqrt{\bar{\lambda}}}{\underline{\lambda}} \| g_{\hat{\Theta}_0}(x) - y \|$$

**Pf:** $\| \hat{\Theta}_t - \hat{\Theta}_0 \| = \left\| \int_0^t \left( - S(s)^T \nabla \hat{L}(g_{\hat{\Theta}_s}) \right) ds \right\|$

$$\leq \int_0^t \underbrace{\| S(s) \|_{op}}_{\leq \sqrt{\bar{\lambda}}} \cdot \underbrace{\| \nabla \hat{L}(g_{\hat{\Theta}_s}) \|}_{g_{\hat{\Theta}_s}(x) - y} ds$$

$$\leq \sqrt{\bar{\lambda}} \cdot \underbrace{\int_0^t \| g_{\hat{\Theta}_s}(x) - y \| ds}_{\substack{\leq \exp(-\underline{\lambda}s) \cdot \| g_{\hat{\Theta}_0}(x) - y \| \\ \text{(by prev. lemma)}}}$$

$$\leq \sqrt{\bar{\lambda}} \cdot \| g_{\hat{\Theta}_0}(x) - y \| \cdot \underbrace{\int_0^s \exp(-\underline{\lambda}s) ds}_{1/\underline{\lambda}} \qquad \square$$

Applying this to $\hat{\Theta}_t = \breve{\Theta}_t$, $g = f_\theta^{lin}$, $S(t) = \gamma J_0$,

$$\Rightarrow \quad \left\| \breve{\Theta}_t - \Theta_c \right\| \leq \frac{\sqrt{\sigma_{max}^2(J_0)}}{\sigma_{min}^2(J_0)} \cdot \left\| f_{\theta_0}(x) - y \right\|$$

$$\leq \frac{\sqrt{2\delta^2 \sigma_{max}^2(J_0)}}{\gamma^2 \sigma_{min}^2(J_0)} \cdot \sqrt{\hat{L}_0}$$

$$< \frac{\sigma_{max}(J_0)}{\gamma \, \sigma_{min}^2(J_0)} \cdot \sqrt{\hat{L}_0}$$

Remains to show can apply Lemmas 1+2 to $(\Theta_t)$. Complication is that $J_t$ is changing over time. We will show it does not change that much, provided $\gamma$ sufficiently large and $\Theta_t$ remains close to initialization.

**Lemma 3:** If $\left\| \theta - \theta_0 \right\| \leq \frac{\sigma_{min}(J_0)}{2\beta} \triangleq B$, then

$$\frac{\sigma_{min}(J_0)}{2} \cdot Id \preceq J_\theta \preceq \frac{3}{2} \frac{\sigma_{max}(J_0)}{2} \cdot Id$$

PF:

$$\|J_\theta\|_{op} \leq \|J_o\|_{op} + \|J_\theta - J_o\|_{op}$$

$$\leq \sigma_{max}(J_o) + \underbrace{\beta\|\theta - \theta_o\|}_{\leq \frac{\sigma_{min}(J_o)}{2} \leq \frac{\sigma_{max}(J_o)}{2}}$$

$$\leq \frac{3\sigma_{max}(J_o)}{2}.$$

for lower bound, for any $\|v\| = 1$,

$$\|J_\theta v\| \geq \|J_o v\| - \|(J_\theta - J_o)v\|$$

$$\geq \|J_o v\| - \beta\|\theta - \theta_o\| \geq \frac{\sigma_{min}(J_o)}{2}. \ \square$$

So we can safely apply Lemma 1+2 to get

$$\hat{L}(\gamma f_{\theta_+}) \leq \exp\left(-\tfrac{1}{2}\gamma^2 \hat{\sigma}_{min}^2(J_o) t\right)$$

$$\|\theta_+ - \theta_o\| \leq \frac{\sigma_{max}}{\gamma \sigma_{min}^2} \cdot \sqrt{\hat{L}_o} \qquad (\maltese)$$

as long as $\|\theta_s - \theta_o\| \leq \beta \quad \forall \ s \in [0, t]$

note that bound in $(\maltese)$ $<<$ $\beta$ as long as

$$\gamma >> \frac{\beta \sigma_{max}}{\sigma_{min}^3}\sqrt{\hat{L}_o}.$$

We conclude

**Thm :** Linearized network $f_{\hat{\Theta}_t}^{lin}$ and

true network $f_{\Theta_t}$ stay $\frac{\sigma_{max}}{\gamma \sigma_{min}^2} \sqrt{\hat{L}_0}$ -close

for all $t \geq 0$, and training loss for $f_{\Theta_t}$

drops exponentially quickly.

**Example :** Consider $\gamma f_\Theta = \gamma \sum_{i=1}^{N} a_i \sigma(\langle w_i, x \rangle)$. For

simplicity, suppose $a_i$'s are random $\{\pm 1\}$'s that

are not subsequently trained, so $\Theta = \{w_i\}_{i=1}^{N}$.

$$\boxed{\beta} \quad J_\Theta = \begin{pmatrix} x_1^T \cdot \{a_i \sigma'(\langle w_i, x \rangle)\}_i \\ \vdots \\ x_n^T \cdot \{a_i \sigma'(\langle w_i, x \rangle)\}_i \end{pmatrix}$$

So $\| J_\Theta - J_{\Theta'} \|_{op}^2$

$$\leq \sum_{i,j} \| x_i \|_2^2 \cdot \underbrace{\left( \sigma'(\langle w_j, x \rangle) - \sigma'(\langle w'_j, x \rangle) \right)^2}_{\lesssim \| w_j - w'_j \|^2}$$

$$= \| X \|_F^4 \cdot \| \Theta - \Theta' \|_2^2$$

So can take $\beta \doteq \|X\|_F^2 \approx \underline{nd}$ (e.g. if $X \sim S^{d-1} \cdot \sqrt{d}$)

$\boxed{\hat{L}_0}$: Can initialize in such a way that

$f_{\theta_0}(x)$ is dominated by $y$. So

$$\hat{L}_0 \approx \|y\|^2 \approx n$$

$\boxed{\sigma_{min}(J_0), \sigma_{max}(J_0)}$: entries of $J_0$ are $O(1)$,

So because $Nd \gg n$, singular values are of order

$\sqrt{Nd}$

Putting everything together, $\gamma \gg \dfrac{\beta \, \sigma_{min}(J_0)}{\sigma_{max}^2(J_0)} \sqrt{\hat{L}_0}$

yields

$$\gamma \gg \frac{nd \cdot \sqrt{Nd}}{(\sqrt{Nd})^3} \cdot \sqrt{n}$$

$$= \frac{n^{3/2}}{N},$$

So provided we are in this regime, gradient flow
well-approx'd by linearized dynamics.