

Lecture 15: NTK opt continued

Goal for these lectures:

* Idealized/abstract (smooth) network
 Minimizes training error of a strongly convex loss by staying close to initialization.

Recall

$$f(w) = \begin{bmatrix} f(x_1; w) \\ \vdots \\ f(x_n; w) \end{bmatrix} \in \mathbb{R}^n, \quad f \text{ smooth}$$

$$J_w^T = \begin{bmatrix} -\nabla f(x_1; w)^T \\ \vdots \\ -\nabla f(x_n; w)^T \end{bmatrix}^T \quad \|J_w - J_v\| \leq \beta \|w - v\|$$

Rem: smooth?

if $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ smooth, is $\sum_{j=1}^n a_j \sigma(w_j^T x)$?

xt $w_2 \sigma(w_{2-1} \sigma(\dots \sigma(w_{1,k}) \dots))$
 $\dots w_k w_{k-1} w_{k-2}$

Two gradient flows:

$$\dot{w}(s) = -\nabla_w L(\alpha f(w(s)))$$

$$= -\alpha J_s \nabla L(\alpha f(w(s))) = -\alpha J(\alpha f(w(s)) - y)$$

"tangent flow"

$$f_0(u) = f(w_0) + J_0^T (u - w_0)$$

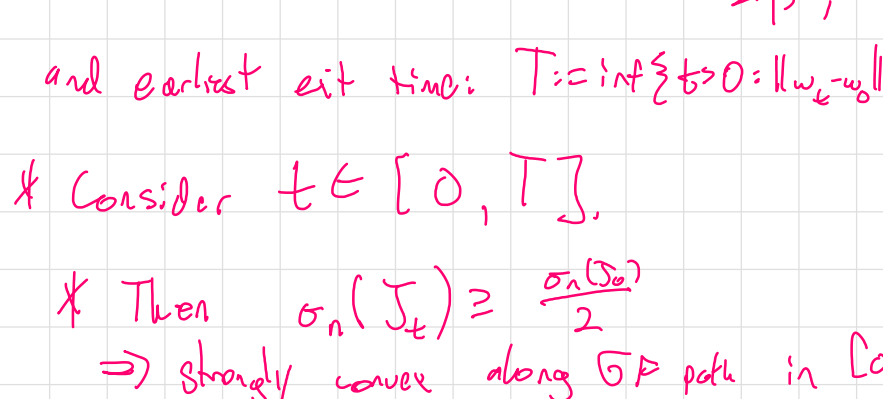
$$\dot{u}(s) = -\nabla_u L(\alpha f_0(u(s)))$$

$$= -\alpha J_0 \nabla L(\alpha f_0(u(s)))$$

Theorem. Assume $\|J_w - J_v\| \leq \beta \|w - v\|$, $\sigma_n(J_0) > 0$
 $\alpha \geq \frac{\beta \sqrt{106} \sigma_{\max}^3 L(\alpha f(w_0))}{\sigma_{\min}^3}$

Then: $\max \{ L(\alpha f(w(t))), L(\alpha f(u(t))) \} \leq L(\alpha f(w_0)) \exp\left(\frac{-t \alpha^2 \sigma_n^2}{2}\right)$

$$\max \{ \|w(t) - w_0\|, \|u(t) - u_0\| \} \leq \frac{3 \sqrt{8} \sigma_n^2(J_0) L(\alpha f(w_0))}{\alpha \sigma_n^2(J_0)}$$



Rem: $\sigma_n > 0$?

$$\frac{1}{2} \|f_0(u) - y\|^2 = \frac{1}{2} \|f(w_0) + J_0^T (u - w_0) - y\|^2$$

$$= \frac{1}{2} \|J_0^T (u - w_0) - (y - f(w_0))\|^2$$

$$= \frac{1}{2} \|J_0^T u - y_0\|^2$$

$$\Downarrow$$

$$J_0^T J_0 u = J_0^T y_0$$

$$\sigma_n(J_0) > 0 \rightarrow u = (J_0^T J_0)^{-1} J_0^T y_0 //$$

Proof plan

Define "convenient" radius $B := \frac{\sigma_n}{2\beta}$, and earliest exit time: $T := \inf \{ t > 0 : \|w_t - w_0\| > B \}$.

* Consider $t \in [0, T]$.

* Then $\sigma_n(J_t) \geq \frac{\sigma_n(J_0)}{2} \Rightarrow$ strongly convex along GP path in $[0, T]$.

* This implies quickly minimize $L(\alpha f(w(t)))$

* Because $\nabla L(\alpha f(w_t)) = (\alpha f(w_t) - y)$, small loss \Rightarrow small gradient norm \Rightarrow iterates don't move much

* We'll show $T = \infty$ (stay in NTK; rare setting).

Lemma. If $\|w - w_0\| \leq B = \frac{\sigma_n}{2\beta}$:

① $\sigma_n(J_w) \geq \frac{\sigma_n}{2}$ & ② $\sigma_1(J_w) \leq \frac{3}{2} \sigma_1$

Pf. ② $\sigma_1(J_w) = \|J_w\|_2 = \|J_w - J_0 + J_0\|_2 \leq \|J_w - J_0\| + \|J_0\| \leq \beta \cdot B + \sigma_1 = \frac{3}{2} \sigma_1$

① $\sigma_n(J_w) = \lambda_{\min}(J_w^T J_w) \in \mathbb{R}^{n \times n}$

$$= \min_{\|u\|=1} u^T J_w^T J_w u = \min_{\|u\|=1} u^T (J_w - J_0 + J_0)^T (J_w - J_0 + J_0) u$$

$$= \min_{\|u\|=1} \|u^T (J_w - J_0)\|^2 + 2u^T (J_w - J_0) J_0 u + \|J_0 u\|^2$$

$$\geq \min_{\|u\|=1} \|u^T (J_w - J_0)\|^2 - 2\|u^T (J_w - J_0)\| \cdot \|J_0\| + \|J_0 u\|^2$$

$$= \min_{\|u\|=1} \left(\|u^T (J_w - J_0)\| - \|J_0 u\| \right)^2$$

$$\geq \left(\sigma_n - \frac{\sigma_n}{2} \right)^2 = \left(\frac{\sigma_n}{2} \right)^2 //$$

$$t \in [0, T] \Rightarrow \|w(t) - w_0\| \leq B \Rightarrow \sigma_n(J_t) \geq \frac{\sigma_n}{2}, \sigma_1(J_t) \leq \frac{3}{2} \sigma_1$$

Next: Loss quickly $\downarrow 0$ for $u(t)$ & $u(t)$ via abstract lemma.

Lemma Consider $\dot{z}(t) = -Q(t) \nabla L(z(t))$ where $\inf_{t \in [0, T]} \lambda_{\min}(Q(t)) \geq \lambda > 0$.

$\forall t \in [0, T] \quad L(z(t)) \leq L(z(0)) \exp(-2\lambda t)$

Example. $z(t) = \alpha f(w(t)) \Rightarrow \dot{z}(t) \approx \frac{d}{dt} \alpha f(w(t)) = \alpha J_t^T \dot{w}(t) = -\alpha J_t^T J_t \nabla L(\alpha f(w(t)))$

$\Rightarrow L(\alpha f(w(t))) \leq L(\alpha f(w_0)) \exp(-\alpha^2 \frac{\sigma_n^2}{2} t)$

Pf. $\frac{d}{dt} L(z(t)) = \langle \nabla L(z(t)), \dot{z}(t) \rangle = \langle z(t) - y, -Q(t) \nabla L(z(t)) \rangle = -(z(t) - y)^T Q(t) (z(t) - y) \leq -\lambda \|z(t) - y\|^2 = -2\lambda L(z(t))$

Grönwall. $\Rightarrow L(z(t)) \leq L(z(0)) \exp(-2\lambda t)$

Still need to show T large.

Lemma. $\dot{v}(t) = -S(t) \nabla L(g(v(t)))$ Suppose $A \geq \sigma_1(S(t)) \geq \sigma_n(S(t)) \geq B$. Then

$$\|v(t) - v(0)\| \leq \frac{\sigma_n}{\sigma_1} \sqrt{2 L(g(v(0)))}$$

Why true: v inc shown $L(\alpha f(w(t))) \downarrow 0$ exponentially fast,

$$\| \dot{w}(t) \| = \| J_t \| \cdot \| \nabla L(\alpha f(w(t))) \| \geq \| J_t \| \cdot \| \alpha f(w(t)) - y \| = \sqrt{2 L(\alpha f(w(t)))}$$

\Rightarrow must also go down exponentially fast.