

CONVERGENCE OF GRADIENT FLOW UNDER OVERPARAMETERIZATION

Let $f(x; \omega, c) := \sum_{i=1}^m c_i \sigma(\omega^T x)$, $\forall x \in \mathbb{R}^d$.

Given $(x_j, y_j) \in \mathbb{R}^d \times \mathbb{R}$, $j = 1, 2, \dots, n$, let

$$f(\omega) := \begin{bmatrix} f(x_1; \omega, c) \\ f(x_2; \omega, c) \\ \vdots \\ f(x_n; \omega, c) \end{bmatrix} \in \mathbb{R}^n.$$

Our objective is $\frac{1}{\alpha^2} g(\alpha f(\omega))$ for a scale factor $\alpha > 0$.
Recall that $\alpha = \frac{1}{\sqrt{m}}$ yields the usual NTK scaling.

Gradient Flow Given a symmetric Xavier initialization $(\omega(0), c)$,

$$\dot{\omega}(t) = -\frac{1}{\alpha} J_{\omega(t)}^T \nabla g(\alpha f(\omega(t))), \quad t \geq 0,$$

where

$$J_{\omega} = \begin{bmatrix} \nabla^T f(x_1; \omega, c) \\ \vdots \\ \nabla^T f(x_n; \omega, c) \end{bmatrix} \in \mathbb{R}^{n \times md}.$$

Target Gradient Flow

$$T_f(u) = J_{\omega(0)} [u - \omega(0)] = [\nabla^T f(x_j; \omega(0), c) [\omega - \omega(0)]]_j;$$

Note that this is the linearized model with random features $\nabla f(x; \omega(0), c)$.

For the linear model, the loss is: $\frac{1}{\alpha^2} g(\alpha T_f(u))$.

Starting from $u(0) = \omega(0)$,

$$\dot{u}(t) = -\frac{1}{\alpha} J_{\omega(0)}^T \nabla g(\alpha T_f(u(t))), \quad t \geq 0.$$

Note: GF is the continuous-time version of GD.

Assumptions:

$$\text{rank}(J_{\omega(0)}) = n$$

$$\sigma_{\min} = \sqrt{\lambda_{\min}(J_{\omega(0)} J_{\omega(0)}^T)} > 0,$$

$$\|J_{\omega} - J_{\nu}\|_2 \leq L \cdot \|\omega - \nu\|, \quad \forall \omega, \nu.$$

Moreover, f is μ -sc. and ν -smooth.

Exercise: Verify that $\nu \mapsto \|\nu - z\|_2^2$ satisfies sc and smoothness assumptions, for any $z \in \mathbb{R}^n$.

THEOREM $\kappa = \frac{\nu}{\mu}$. Let f^* be the minimizer of

If $\kappa > \frac{\|f^*\|}{M}$ for $M := \frac{\sigma_{\min}^2}{32\kappa^{3/2} \|J_{\omega(0)}\|_2 L}$, then, for $t \geq 0$,

$$\|\kappa f(\omega(t)) - f^*\| \leq \sqrt{\kappa} \|\kappa f(\omega(0)) - f^*\|_2 \cdot e^{-\mu \sigma_{\min}^2 t/4}.$$

Furthermore, as $\kappa \rightarrow \infty$,

$$\sup_{t \geq 0} \|\omega(t) - \omega(0)\|_2 = O(1/\kappa),$$

Notes: We will relax the full-rank assumption later in **Lemma 3**, and show that overparameterization ($m = \tilde{O}(n^2)$) is required to ensure $\sigma_{\min} > 0$.

Above, $\|\omega(t) - \omega(0)\| \leq O(1/\kappa)$ is satisfied without projection.

Previously, we needed projection to keep $\|\omega(t) - \omega(0)\|$ bounded.

But, we just required $m = O(\log(n))$ neurons rather than $\tilde{O}(n^2)$.

Proof | Trajectory in the function space:

$$\frac{d}{dt} \kappa f(\omega(t)) = -J_{\omega(t)} J_{\omega(t)}^T \nabla_g(\kappa f(\omega(t))).$$

Let $\tau_0 := \frac{\sigma_{\min}}{2L}$. Then,

$$\|\omega - \omega(0)\|_2 < \tau_0 \Rightarrow J_{\omega} J_{\omega}^T \geq \frac{\sigma_{\min}^2}{4} I. \quad \begin{array}{l} \text{(Lemma 1)} \\ \text{(next page)} \end{array}$$

Let $T = \inf \{t \geq 0 : \|\omega(t) - \omega(0)\|_2 > \tau_0\}$.

We want to show that $T = \infty$.

From gradient flow (i.e., $\dot{\omega}(t) = -\frac{1}{\alpha} J_{\omega(t)}^T \nabla_g(\kappa f(\omega(t)))$),

$$\begin{aligned} \|\dot{\omega}(t)\|_2 &\leq \frac{1}{\alpha} \|J_{\omega(t)}\|_2 \cdot \|\nabla_g(\kappa f(\omega(t))) - \underbrace{\nabla_g(f^*)}_{=0}\|_2 \\ &\leq \frac{2\nu}{\alpha} \|J_{\omega(0)}\|_2 \cdot \|\kappa f(\omega(t)) - f^*\|_2 \quad \text{since } g \text{ is } \nu\text{-sm.} \end{aligned}$$

By Lemma 2, for $t \in [0, T]$,

$$\begin{aligned} \|\omega(t) - \omega(0)\| &= \left\| \int_0^t \dot{\omega}(s) ds \right\| \leq \int_0^t \|\dot{\omega}(s)\| ds \leq \int_0^t \|\kappa f(\omega(s)) - f^*\| ds \frac{2\nu}{\alpha} \|J_{\omega(0)}\| \\ &\leq \frac{2\nu}{\alpha} \|J_{\omega(0)}\| \cdot \int_0^t \sqrt{\frac{\nu}{\mu}} \cdot \|\kappa f(\omega(0)) - f^*\| \cdot e^{-\frac{\mu \sigma_{\min}^2}{4}s} ds \\ &\leq \frac{8\kappa^{3/2}}{\alpha \sigma_{\min}^2} \|J_{\omega(0)}\| \cdot \|\kappa f(\omega(0)) - f^*\|_2 = \frac{8\kappa^{3/2}}{\alpha \sigma_{\min}^2} \|J_{\omega(0)}\| \cdot \|f^*\| \end{aligned}$$

$$\text{Sufficiently large } \alpha \Rightarrow \frac{8\kappa^{3/2}}{\alpha \sigma_{\min}^2} \|J_{\omega(0)}\| \cdot \|f^*\| < \tau_0$$

$$\Rightarrow T = \infty. \quad \blacksquare$$

Lemma 1 and Lemma 2 are on the following pages.

Lemma 1

Suppose $\|\omega - \omega(0)\| \leq r_0 = \frac{\sigma_{\min}}{2L}$. Then,

$$\sigma_{\min}(J_\omega) \geq \sigma_{\min}(J_0) - L \cdot \|\omega - \omega(0)\| \geq \frac{\sigma_{\min}}{2}.$$

Proof

Given u , let $Au = J_{\omega(0)}^T u$ and $Bu = (J_\omega - J_{\omega(0)})^T u$.

Then,

$$\sigma_{\min}^2(J_\omega) = \min_{\|u\|=1} u^T J_\omega J_\omega^T u$$

$$= \min_u \left[(J_{\omega(0)} + J_\omega - J_{\omega(0)})^T u \right]^T \left[(J_{\omega(0)} + J_\omega - J_{\omega(0)})^T u \right]$$

$$= \min_u \|Au\|^2 + 2Au^T Bu + \|Bu\|_2^2$$

$$\stackrel{\text{C.S.}}{\geq} \min_u \|Au\|^2 - 2\|Au\| \cdot \|Bu\| + \|Bu\|^2$$

$$\begin{aligned} &= \min_u (\|Au\| - \|Bu\|)^2 \geq \min_u (\sigma_{\min} - \beta L)^2 \|u\|^2 \\ &= \frac{\sigma_{\min}^2}{4}. \quad \blacksquare \end{aligned}$$

Lemma 2

Suppose that g is μ -sc and ν -smooth, z^* is the global minimizer of g . Let $Q(t)$ be a continuous linear operator s.t. \rightarrow auto adjoint

$$\inf_{t \in [0, \tau]} \lambda_{\min}(Q(t)) \geq \lambda > 0.$$

Then, solutions on $[0, \tau]$ to

$$\dot{z}(t) = -Q(t) \nabla g(z(t))$$

satisfy, for $t \in [0, \tau]$,

$$\|z(t) - z^*\| \leq \frac{\nu}{\mu} \cdot \|z(0) - z^*\| \cdot e^{-\mu \lambda t}$$

Pf: From μ -sc. of g , (Łojasiewicz Lemma)

$$\bar{g}(z) := g(z) - g(z^*) \leq \frac{1}{2\mu} \cdot \|\nabla g(z)\|^2$$

Then,

$$\begin{aligned} \frac{d}{dt} \bar{g}(z(t)) &\stackrel{\text{chain rule}}{=} -\nabla g^T(z(t)) Q(t) \nabla g(z(t)) \\ &\leq -\lambda \|\nabla g(z(t))\|^2 \\ &\leq -2\mu\lambda \cdot [g(z(t)) - g(z^*)] = -\bar{g}(z(t)) \end{aligned}$$

By Grönwall's lemma,

$$\bar{g}(z(t)) \leq \bar{g}(z(0)) \cdot e^{-2\mu\lambda t}$$

Using μ -sc and ν -smooth properties of g above,

$$\frac{\mu}{2} \|z - z^*\|_2^2 \leq \bar{g}(z) \leq \frac{\nu}{2} \|z - z^*\|_2^2$$

$$\frac{\mu}{2} \|z(t) - z^*\|_2^2 \leq e^{-2\mu\lambda t} \cdot \frac{\nu}{2} \|z(0) - z^*\|_2^2$$

$$\Rightarrow \|z(t) - z^*\| \leq e^{-\mu\lambda t} \cdot \frac{\nu}{\mu} \cdot \|z(0) - z^*\|. \quad \blacksquare$$

Remark (An implication of the full-rank assumption)

$$T_f(\omega) = \begin{bmatrix} \nabla_{\omega}^T f(x_1; \omega(0), c) [\omega - \omega(0)] \\ \vdots \\ \nabla_{\omega}^T f(x_n; \omega(0), c) [\omega - \omega(0)] \end{bmatrix} \quad \text{is the linearized regressor.}$$

Consider

$$\min_{\omega \in \mathbb{R}^m} \frac{1}{2} \|T_f(\omega) - y\|_2^2 = \min_{\omega} \frac{1}{2} \|J_{\omega(0)} \omega - y_0\|_2^2$$

where

$$y_0 = y + \underbrace{J_{\omega(0)} \omega(0)}_{=0 \text{ for ReLU}}$$

Then, the normal equations:

$$J_{\omega(0)}^T J_{\omega(0)} \omega = J_{\omega(0)}^T y_0. \quad (+)$$

$$J_{\omega(0)} \text{ is full-rank} \Rightarrow J_{\omega(0)} = \sum_{i=1}^n s_i u_i v_i^T \quad \text{SVD.}$$

$$J_{\omega(0)}^+ = \sum_{i=1}^n s_i^{-1} v_i u_i^T$$

Then, from (+):

$$(J_{\omega(0)}^+)^T J_{\omega(0)}^T J_{\omega(0)} \omega = (J_{\omega(0)}^+)^T J_{\omega(0)}^T y_0$$

$$\Rightarrow J_{\omega(0)} \omega = \underbrace{\left[\sum_{i=1}^n u_i u_i^T \right]}_{\text{idempotent and full rank}} y_0 = y_0.$$

$$\text{Choose } \hat{\omega} = J_{\omega(0)}^+ y_0. \text{ Then, } J_{\omega(0)} \hat{\omega} = \sum_i u_i u_i^T y_0 = y_0.$$

$$\Rightarrow \frac{1}{2} \|T_f(\hat{\omega}) - y\|_2^2 = \frac{1}{2} \|J_{\omega(0)} \hat{\omega} - y_0\|_2^2 = 0.$$

No assumptions on y_0 !

Remark

How to ensure that $\lambda_{\min}(\mathcal{J}_{\omega(0)} \mathcal{J}_{\omega(0)}^T) > 0$?
→ overparameterization.

$$\begin{aligned} \frac{1}{n} [\mathcal{J}_{\omega(0)} \mathcal{J}_{\omega(0)}^T]_{i,j} &= \frac{1}{n} \nabla^T f(x_i; \omega(0), c) \nabla f(x_j; \omega(0), c) \\ &= \frac{1}{n} x_i^T x_j \sum_{k=1}^m \sigma'(\omega_k^T(0) x_i) \sigma'(\omega_k^T(0) x_j) \\ &\xrightarrow{n \rightarrow \infty} x_i^T x_j \mathbb{E}_{\omega_0 \sim N(0, I_d)} [\sigma'(\omega_0^T x_i) \sigma'(\omega_0^T x_j)]. \end{aligned}$$

$$= K_{ij}.$$

Lemma 3

Let $\|x_i\|_2 \leq 1$, $|\sigma'| \leq 1$.

Assume $\lambda_{\min}(K) = \lambda_0 > 0$. Then,

$$\text{if } m = \Omega\left(\frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right), \quad \text{w.p.} \geq 1 - \delta,$$

$$\left\| \frac{1}{n} \mathcal{J}_{\omega(0)} \mathcal{J}_{\omega(0)}^T - K \right\|_2 \leq \frac{\lambda_0}{4}, \text{ and}$$

$$\lambda_{\min}\left(\frac{1}{n} \mathcal{J}_{\omega(0)} \mathcal{J}_{\omega(0)}^T\right) \geq \frac{3\lambda_0}{4}.$$

Proof: $\forall i, j \in \{1, \dots, n\}$,

$$\left| \frac{1}{n} [\mathcal{J}_{\omega(0)} \mathcal{J}_{\omega(0)}^T]_{i,j} - K_{ij} \right| \leq \frac{2\sqrt{\log(n^2/\delta)}}{\sqrt{m}}$$

simultaneously w.p. $\geq 1 - \delta$. Then,

$$\left\| \frac{1}{n} \mathcal{J}_{\omega(0)} \mathcal{J}_{\omega(0)}^T - K \right\|_2^2 \leq \left\| \frac{1}{n} \mathcal{J}_{\omega(0)} \mathcal{J}_{\omega(0)}^T - K \right\|_F^2$$

$$\leq \frac{8 \log(n/\delta) n^2}{m} \quad \blacksquare$$