

ECE 8101: Nonconvex Optimization for Machine Learning

Lecture Note 5-1: The Polyak-Łojasiewicz (PL) Condition
(feat. Neural Tangent Kernel)

Jia (Kevin) Liu

Assistant Professor
Department of Electrical and Computer Engineering
The Ohio State University, Columbus, OH, USA

Spring 2022

Outline

In this lecture:

- The Polyak-Łojasiewicz (PL) Condition
- Convergence of Various Methods under the PL Condition
- The PL Condition and the Over-parameterized Regime

Convergence Results of Methods We Learned Thus Far

- First-order and zeroth-order methods for nonconvex optimization in learning:
 - ▶ GD/SGD-style algorithms
 - ▶ Only focus on **stationarity gap**
 - ▶ Typically sublinear convergence rates: $O(1/K)$, $\frac{O(1/\sqrt{K})}{\text{SGD}}$, ... ($\frac{O(1/K^2)}{\text{Nesterov}}$ is order-optimal) (AGD, \dots)
- Meanwhile, it's well-known from convex optimization that:
 - ▶ GD achieves **linear convergence rate** under **strong convexity**
 - ▶ Convergence of **global optimality**

Can global linear convergence to optimality happen under weaker conditions?

The Polyak-Łojasiewicz Condition

Definition 1 ([Polyak, '63], [Łojasiewicz, '63])

A function $f(\mathbf{x})$ is said to satisfy the Polyak-Łojasiewicz (PL) condition if for all $\mathbf{x} \in \mathbb{R}^d$, there exists a constant $\mu > 0$ such that:

$$2\mu(f(\mathbf{x}) - f(\mathbf{x}^*)) \leq \|\nabla f(\mathbf{x})\|_2^2.$$

Remarks

- Aka “gradient dominated” condition (e.g., [Reddi et al., ICML'16])
- Implies any stationary point is a **global min**, although not necessarily unique
- Automatically holds for **strongly convex** functions
- Many **nonconvex** functions satisfy PL condition, especially in the over-parameterized regime
- **PL condition** means that the **optimality gap** $f(\mathbf{x}) - f^*$ is upper bounded by a quadratic function of the **stationarity gap**

Prop. SC \rightarrow PL.

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2, \forall x, y. \quad (\text{strong convexity})$$

Minimize both sides w.r.t. y :

LHS: $f(x^*)$

RHS: quad fn of y : take deriv. w.r.t. y , set it to 0:

$$\nabla f(x) + \mu(y-x) = 0 \Rightarrow y^* = x - \frac{1}{\mu} \nabla f(x).$$

plugging y^* back to RHS:

$$f(x) - \frac{1}{\mu} \|\nabla f(x)\|^2 + \frac{1}{2\mu} \|\nabla f(x)\|^2 = f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$$

So: $f(x^*) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$

$$\Rightarrow \|\nabla f(x)\|^2 \geq 2\mu (f(x) - f(x^*)).$$



Nice Features of the PL Condition

- Ease of verification compared to strong convexity (SC):
 - ▶ One only needs to access $\|\nabla f(\mathbf{x})\|$ and $f(\mathbf{x})$. In comparison, SC requires checking PD of the Hessian matrix \mathbf{H} (accurate estimation of $\lambda_{\min}(\mathbf{H})$)
- Robustness of the condition
 - ▶ $\|\nabla f(\mathbf{x})\|$ is more resilient to perturbation of the obj function than $\lambda_{\min}(\mathbf{H})$
- Allows multiple global minima:
 - ▶ Modern ML problems are over-parameterized and have manifolds of global minima, not compatible with SC in general but compatible with PL
- Invariance under transformation:
 - ▶ PL is invariant under a broad class of nonlinear coordinate transformations arising from feature extraction/transformation of many ML applications
- PL on manifolds:
 - ▶ PL allows for efficient optimization on manifolds, while being invariant under the choice of coordinates (see [Weber and Sra, arXiv:1710:10770])
- Linear convergence of GD and SGD: *Frank-Wolfe*
 - ▶ PL is sufficient not only for GD but also for SGD

Gradient Descent under the PL Condition

Theorem 2 (Linear Convergence Rate for GD)

Consider the unconstrained optimization problem $\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$, where f has an L -Lipschitz continuous gradient, a non-empty solution set \mathcal{X}^* , and satisfies the PL condition. Then, the *gradient descent* method with a step-size of $1/L$, i.e., $\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)$, has a *global linear convergence rate*:

$$f(\mathbf{x}_k) - f^* \leq \left(1 - \frac{\mu}{L}\right)^k (f(\mathbf{x}_0) - f^*).$$

Remarks

- For twice differentiable functions, L -smoothness means eigenvalues of $\nabla^2 f(\mathbf{x})$ are bounded from above by L (curvature upper bound)

Theorem 2 (Linear Convergence Rate for GD)

Consider the unconstrained optimization problem $\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$, where f has an L -Lipschitz continuous gradient, a non-empty solution set \mathcal{X}^* , and satisfies the PL condition. Then, the gradient descent method with a step-size of $1/L$, i.e., $\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)$, has a **global linear convergence rate**:

$$f(\mathbf{x}_k) - f^* \leq \left(1 - \frac{\mu}{L}\right)^k (f(\mathbf{x}_0) - f^*).$$

Proof: GD under PL, $\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)$.

$$\begin{aligned} f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) &\leq \nabla f(\mathbf{x}_k)^T \underbrace{(\mathbf{x}_{k+1} - \mathbf{x}_k)}_{-\frac{1}{L} \nabla f(\mathbf{x}_k)} + \frac{L}{2} \underbrace{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2}_{-\frac{1}{L} \nabla f(\mathbf{x}_k)} \\ &\stackrel{\text{using GD dynamic}}{\leq} -\frac{1}{L} \|\nabla f(\mathbf{x}_k)\|^2 + \frac{L}{2} \cdot \frac{1}{L^2} \|\nabla f(\mathbf{x}_k)\|^2 \end{aligned}$$

$$= -\frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|^2$$

$$\stackrel{\text{PL}}{\leq} -\frac{1}{2L} \cdot \mu (f(\mathbf{x}_k) - f^*) = -\frac{\mu}{L} (f(\mathbf{x}_k) - f^*).$$

Move $f(\mathbf{x}_k)$ to RHS, and subtract f^* on both sides:

$$f(\mathbf{x}_{k+1}) - f^* \leq f(\mathbf{x}_k) - f^* - \frac{\mu}{L} (f(\mathbf{x}_k) - f^*) = \left(1 - \frac{\mu}{L}\right) (f(\mathbf{x}_k) - f^*). \quad \square$$

Exact LS:

$$f(\mathbf{x}_{k+1}) = \min_s \left\{ f(\mathbf{x}_k) - s \nabla f(\mathbf{x}_k) \right\} \leq f(\mathbf{x}_k) - \frac{1}{L} \nabla f(\mathbf{x}_k).$$

Stochastic Gradient Descent under the PL Condition

- The finite-sum minimization problem: $\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x})$
- Consider the SGD method that uses the iteration: $\mathbf{x}_{k+1} = \mathbf{x}_k - s_k \nabla_{i_k} f(\mathbf{x}_k)$

Theorem 3 (Convergence Rate for SGD)

Assume that f has L -Lipschitz continuous gradients and a non-empty solution set \mathcal{X}^* , and it satisfies the PL condition, and f satisfies $\|\nabla f_{i_k}(\mathbf{x}_k)\| \leq C^2$ for all \mathbf{x}_k and some constant $C > 0$. Then, it holds that: $= O(\frac{1}{\mu k})$.

- SGD with diminishing step-size $s_k = \frac{2k+1}{2\mu(k+1)^2}$ has a convergence rate of:

$$\mathbb{E}[f(\mathbf{x}_k) - f^*] \leq \frac{LC^2}{2\mu^2 k} = O\left(\frac{1}{k}\right). \quad \left(\text{as opposed to } O\left(\frac{1}{k}\right) \right).$$

- SGD with constant step-size $s_k = s \leq \frac{1}{2\mu}$ has a convergence rate of:

$$\mathbb{E}[f(\mathbf{x}_k) - f^*] \leq (1 - 2s\mu)^k [f(\mathbf{x}_0) - f^*] + \frac{LC^2 s}{4\mu}.$$

Theorem 3 (Convergence Rate for SGD)

Assume that f has L -Lipschitz continuous gradients and a non-empty solution set \mathcal{X}^* , and it satisfies the PL condition, and f satisfies $\|\nabla f_{i_k}(\mathbf{x}_k)\| \leq C$ for all \mathbf{x}_k and some constant $C > 0$. Then, it holds that: $= O(\frac{1}{\mu k})$.

- SGD with diminishing step-size $s_k = \frac{2k+1}{2\mu(k+1)^2}$ has a convergence rate of:

$$\mathbb{E}[f(\mathbf{x}_k) - f^*] \leq \frac{LC^2}{2\mu^2 k} = O\left(\frac{1}{k}\right). \quad \left(\text{as opposed to } O\left(\frac{1}{k}\right)\right).$$

- SGD with constant step-size $s_k = s \leq \frac{1}{2\mu}$ has a convergence rate of:

$$\mathbb{E}[f(\mathbf{x}_k) - f^*] \leq (1 - 2s\mu)^k [f(\mathbf{x}_0) - f^*] + \frac{LC^2 s}{4\mu}.$$

Proof. $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x}_{k+1} - \mathbf{x}_k) + \frac{L}{2} \|\nabla f(\mathbf{x}_k)\|^2$.

$\stackrel{\text{SGD dynamic}}{=} f(\mathbf{x}_k) + s_k \nabla f(\mathbf{x}_k)^T \nabla f_{i_k}(\mathbf{x}_k) + \frac{L s_k^2}{2} \|\nabla f_{i_k}(\mathbf{x}_k)\|^2$.

Take full expectation w.r.t. $\{i_k\}$:

$$\mathbb{E}[f(\mathbf{x}_{k+1})] \leq \mathbb{E}[f(\mathbf{x}_k)] - s_k \mathbb{E}[\nabla f(\mathbf{x}_k)^T \nabla f_{i_k}(\mathbf{x}_k)] + \frac{L s_k^2}{2} \mathbb{E}[\|\nabla f_{i_k}(\mathbf{x}_k)\|^2]$$

$$\stackrel{\leq}{\leq} \mathbb{E}[f(\mathbf{x}_k)] - 2s_k \mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2] + \frac{LC^2 s_k^2}{2}$$

$$\stackrel{\text{PL}}{\leq} \mathbb{E}[f(\mathbf{x}_k)] - 2\mu s_k \mathbb{E}[f(\mathbf{x}_k) - f^*] + \frac{LC^2 s_k^2}{2}.$$

Subtracting f^* on both sides:

$$\mathbb{E}[f(\mathbf{x}_{k+1}) - f^*] \leq (1 - 2\mu s_k) \mathbb{E}[f(\mathbf{x}_k) - f^*] + \frac{LC^2 s_k^2}{2}. \quad (1)$$

• diminishing step-size: $s_k = \frac{2k+1}{2\mu(k+1)^2} = O\left(\frac{1}{\mu k}\right)$.

$$\mathbb{E}[f(\mathbf{x}_{k+1}) - f^*] \leq \frac{k^2}{(k+1)^2} \mathbb{E}[f(\mathbf{x}_k) - f^*] + \frac{LC^2 (2k+1)^2}{8\mu^2 (k+1)^4}$$

Multiplying both sides by $(k+1)^{-2}$ and letting $\delta_f(k) = k^2 \mathbb{E}[f(z_k) - f^*]$.

$$\delta_f(k+1) \leq \delta_f(k) + \frac{LC^2(2k+1)^2}{8\mu^2(k+1)^4} \leq \delta_f(k) + \frac{LC^2}{2\mu^2} \quad (2)$$

Summing (2) from $k=0$ to $k-1$ and $\delta_f(0) = 0^2 \mathbb{E}[f(z_0) - f^*] = 0$.

$$\delta_f(k) \leq \underbrace{\delta_f(0)}_{=0} + \frac{LC^2 k}{2\mu^2} \Rightarrow k^2 \mathbb{E}[f(z_k) - f^*] \leq \frac{LC^2 k}{2\mu^2}$$

$$\Rightarrow \mathbb{E}[f(z_k) - f^*] \leq \frac{LC^2}{2k\mu^2}$$

2° Constant step-size: $s_k = s$ for some $s < \frac{1}{2\mu}$. Applying in (1). (recursively)

$$\begin{aligned} \mathbb{E}[f(z_{k+1}) - f^*] &\leq (1-2\mu s)^k \mathbb{E}[f(z_0) - f^*] + \frac{LC^2 s^2}{2} \sum_{i=0}^k (1-2\mu s)^i \\ &\leq (1-2\mu s)^k \mathbb{E}[f(z_0) - f^*] + \frac{LC^2 s^2}{2} \sum_{i=0}^{\infty} (1-2\mu s)^i \\ &= (1-2\mu s)^k \mathbb{E}[f(z_0) - f^*] + \underbrace{\frac{LC^2 s^2}{2} \cdot \frac{1}{2\mu s}}_{= \frac{LC^2 s}{4\mu}} \end{aligned}$$

SGD under PL Condition in Over-parameterized Regime

- Consider ERM in **over-parameterized regime**: $\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x})$
- Applied by Lipschitz cont. and smoothness of f_i .*
- $f(\mathbf{x})$ is L -smooth: $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y}$
 - $f_i(\mathbf{x})$ satisfies: $\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq \tilde{L}|f_i(\mathbf{x}) - f_i(\mathbf{y})|$ for some $\tilde{L} > 0$
 - In ML problems, w.l.o.g., we can assume that $\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = 0$ and so the PL condition can be modified as μ -PL*: $2\mu f(\mathbf{x}) \leq \|\nabla f(\mathbf{x})\|_2^2$

$$\|\nabla f(\mathbf{x})\|_2^2 \geq 2\mu(f(\mathbf{x}) - f^*)$$

- Over-parameterized regime**: $d \gg N$

- The **interpolation** effect: for every sequence $\mathbf{x}_1, \mathbf{x}_2, \dots$ such that $\lim_{k \rightarrow \infty} f(\mathbf{x}_k) = 0$, we have

$$\lim_{k \rightarrow \infty} f_i(\mathbf{x}_k) = 0, \quad 1 \leq i \leq N.$$

- Meaning**: In the over-parameterized regime, the richness of the model is so high such that **fit all** training samples

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathbb{R}^{N \times d}, \quad \mathbf{x} \in \mathbb{R}^d, \quad f = \min_{\mathbf{z} \in \mathbb{R}^d} \frac{1}{2} \|A\mathbf{z} - \mathbf{b}\|^2.$$

$N > d$. Null space: $N-d$ dimension.

SGD under PL Condition in Over-parameterized Regime

- Consider the general mini-batched version of SGD with constant step-size s :

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{s}{B} \sum_{j=1}^B \nabla f_{i_k^j}(\mathbf{x}_k),$$

- B : the mini-batch size; the sample indices $\{i_k^1, \dots, i_k^B\}$ in the mini-batch are drawn uniformly with replacement in each iteration k from $\{1, \dots, N\}$

Theorem 4 ([Bassily et al., arXiv:1811.02564])

Consider the mini-batch SGD with smooth losses as stated. Suppose the *interpolation* condition holds. Suppose that the ERM function $f(\mathbf{x})$ is μ -PL* for some $\mu > 0$. For any mini-batch size $B \in \mathbb{N}$, the mini-batch SGD with *constant* step-size $s^*(B) \triangleq \frac{2\mu B}{L(\bar{L} + L(B-1))}$ guarantees that:

$$\mathbb{E}[f(\mathbf{x}_k)] \leq (1 - \mu s^*(B))^k f(\mathbf{x}_0)$$

Theorem 4 ([Bassily et al., arXiv:1811.02564])

Consider the mini-batch SGD with smooth losses as stated. Suppose the interpolation condition holds. Suppose that the ERM function $f(\mathbf{x})$ is μ -PL* for some $\mu > 0$. For any mini-batch size $B \in \mathbb{N}$, the mini-batch SGD with constant step-size $s^*(B) \triangleq \frac{2\mu B}{L(\tilde{L} + L(B-1))}$ guarantees that:

$$\mathbb{E}[f(\mathbf{x}_k)] \leq (1 - \mu s^*(B))^k f(\mathbf{x}_0)$$

Proof: From descent lemma:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_{k+1} - \mathbf{x}_k) + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2.$$

Using SGD dynamics:

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -s \nabla f(\mathbf{x}_k)^\top \left[\frac{1}{B} \sum_{j=1}^B \nabla f_{i_k^j}(\mathbf{x}_k) \right] + \frac{Ls^2}{2} \left\| \frac{1}{B} \sum_{j=1}^B \nabla f_{i_k^j}(\mathbf{x}_k) \right\|^2.$$

Fix \mathbf{x}_k , and take expectation over choice $\{i_k^1, \dots, i_k^B\}$.

Note: indices are i.i.d., then, we have:

$$\mathbb{E} \left[f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \mid \mathbf{x}_k \right] \leq -s \|\nabla f(\mathbf{x}_k)\|^2 + \frac{s^2 L}{2} \left(\frac{1}{B} \mathbb{E}_{i_k} \left[\|\nabla f_{i_k}(\mathbf{x}_k)\|^2 \right] + \frac{B-1}{B} \|\nabla f(\mathbf{x}_k)\|^2 \right).$$

Since $\forall i \in \{1, \dots, N\}$, $f_i(\cdot)$ satisfies:

$$\|\nabla f_{i_k}(\mathbf{x}_k) - \underbrace{\nabla f_{i_k}(\mathbf{x}^*)}_{=0}\| \leq \tilde{L} \|f_{i_k}(\mathbf{x}_k) - \underbrace{f_{i_k}(\mathbf{x}^*)}_{=0}\| = \tilde{L} \underbrace{\|f_{i_k}(\mathbf{x}_k)\|}_{\geq 0} \leq 2\tilde{L} f_{i_k}(\mathbf{x}_k)$$

Thus,

$$\mathbb{E} \left[f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \mid \mathbf{x}_k \right] \leq -s \left(1 - \frac{sL}{2} \cdot \frac{B-1}{B} \right) \|\nabla f(\mathbf{x}_k)\|^2 + \frac{s^2 L \tilde{L}}{B} f(\mathbf{x}_k).$$

$$\begin{aligned} \text{non-neg. } f(z) & \stackrel{PL}{\leq} 2\mu s \left(1 - \frac{sL(B-1)}{B} \right) f(z_k) + \frac{s^2 L \tilde{L}}{B} f(z_k) \\ & = s \left[-2\mu - \frac{sL}{B} (\mu(B-1) + \tilde{L}) \right] f(z_k). \end{aligned}$$

Finally, rearranging & taking full expectation:

$$\mathbb{E}[f(z_{k+1})] \leq \left[1 - 2s\mu + \frac{s^2 L}{B} (\mu(B-1) + \tilde{L}) \right] \mathbb{E}[f(z_k)]. \quad (3)$$

Optimizing the quad term in (3) w.r.t. s yield $s^*(B)$,

$$\mathbb{E}[f(z_{k+1})] \leq (1 - \mu s^*(B)) \mathbb{E}[f(z_k)]. \quad \square$$

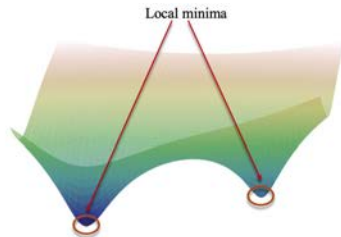
Other Methods under the PL Condition

Similar linear convergence rate results can be shown for other methods under the μ -PL, L -smoothness, and uniform variance bound conditions, which implies the following sample complexity results:

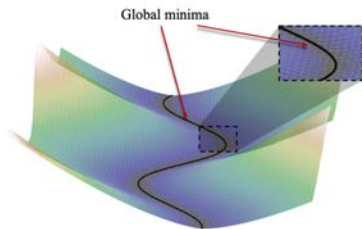
- GD [Polyak, '63]: $\frac{L}{\mu} \log \frac{\Delta_0}{\epsilon}$
- SGD [Karimi et al., ECML-KDD'16]: $\frac{L}{\mu} \left(\frac{\max_i L_i}{\mu} \log\left(\frac{\Delta_0}{\epsilon}\right) + \frac{\max_i L_i \Delta_*}{\mu \epsilon} \right)$
- SVRG [Reddi et al., NeurIPS'16]: $\left(N + \frac{N^{2/3} \max_i L_i}{\mu} \right) \log\left(\frac{\Delta_0}{\epsilon}\right)$
- SAGA [Reddi et al., NeurIPS'16]: $\left(N + \frac{N^{2/3} \max_i L_i}{\mu} \right) \log\left(\frac{\Delta_0}{\epsilon}\right)$
- PAGE [Li et al., ICML'21]: $\left(b + \sqrt{b} \frac{L_{\text{avg}}}{\mu} \right) \log\left(\frac{\Delta_0}{\epsilon}\right)$, where $b = \min\left\{\frac{\sigma^2}{\mu \epsilon}, N\right\}$

PL Condition and Over-parameterized Regime

- Landscape of under-parameterized and over-parameterized models (figure from [Liu et al., arXiv:2003:00307])



(a) Loss landscape of under-parameterized models



(b) Loss landscape of over-parameterized models

$d \gg N.$

- Key Insight:
 - ▶ Convexity is not the right framework for analyzing the loss landscape of over-parameterized systems, even locally
 - ▶ Instead, the μ -PL* condition (i.e., $\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu f(\mathbf{w}), \forall \mathbf{w}$) is a more appropriate framework

PL Condition and Over-parameterized Regime

The essence of supervised learning:

- Given a dataset of size N , $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$, $\mathbf{x}_i \in \mathbb{R}^d$, $y \in \mathbb{R}$
- A parametric family of models $f(\mathbf{w}, \mathbf{x})$ (e.g., a neural network)
- **Goal:** To find a model with parameter \mathbf{w}^* that **fits** the training data:

$$f(\mathbf{w}^*, \mathbf{x}_i) \approx y_i, \quad i = 1, 2, \dots, N$$

- **Mathematically:** Equivalent to solving (exactly or approximately) a system of N nonlinear equations:

$$\mathcal{F}(\mathbf{w}) = \mathbf{y},$$

vector-valued

where $\mathbf{w} \in \mathbb{R}^d$, $\mathbf{y} \in \mathbb{R}^N$, and $\mathcal{F}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^N$ with $(\mathcal{F}(\mathbf{w}))_i = f(\mathbf{w}, \mathbf{x}_i)$.

- The system of equations is solved by minimizing a certain loss function $\mathcal{L}(\mathbf{w})$
 - ▶ E.g., the square loss: $\mathcal{L}(\mathbf{w}) = \frac{1}{2} \|\mathcal{F}(\mathbf{w}) - \mathbf{y}\|^2 = \frac{1}{2} \sum_{i=1}^N (f(\mathbf{w}, \mathbf{x}_i) - y_i)^2$

PL Condition and Over-parameterized Regime

μ -PL* condition emerges through the spectrum of the tangent kernel

- Let $D\mathcal{F}(\mathbf{w}) \in \mathbb{R}^{N \times d}$ be the differential of the mapping \mathcal{F} at \mathbf{w}
- The **tangent kernel** of \mathcal{F} is defined as an $N \times N$ matrix:

$$\mathbf{K}(\mathbf{w}) \triangleq D\mathcal{F}(\mathbf{w})D\mathcal{F}^\top(\mathbf{w})$$

- ▶ It follows from the definition that $\mathbf{K}(\mathbf{w})$ is PSD

- The square loss \mathcal{L} is μ -PL* at \mathbf{w} [Liu, et al., arXiv:2003:00307], where

$$\| \nabla \mathcal{L}(\mathbf{z}) - \nabla \mathcal{L}(\mathbf{z}^*) \| \geq \mu \| \mathbf{z} - \mathbf{z}^* \|$$

$\mu = \lambda_{\min}(\mathbf{K}(\mathbf{w}))$

is the smallest eigenvalue of the kernel matrix

Thus, the PL* condition is inherently tied to the spectrum of the tangent kernel matrix associated with \mathcal{F}

PL Condition and Over-parameterized Regime

Wide (hence over-parameterized) neural networks satisfy PL* condition:

- A powerful tool: the **neural tangent kernel** (NTK)
 - ▶ First appeared in a landmark paper [Jacot et al., NeurIPS'18]
 - ▶ Tangent kernel of a ^{hidden}single-layer wide neural networks with linear output layer ($f(\mathbf{x}) = \sum_{i=1}^d \sigma(\mathbf{w}^\top \mathbf{x})$) are nearly constant in a ball \mathcal{B} of a certain radius around the ball with a random center (note: d is also the width of the NN):

$$\|\mathbf{H}_{\mathcal{F}}(\mathbf{w})\| = O^*(1/\sqrt{d}),$$

where $\mathbf{H}_{\mathcal{F}}(\mathbf{w})$ is a $N \times d \times d$ tensor with $(\mathbf{H}_{\mathcal{F}})_{ijk} = \frac{\partial^2 \mathcal{F}_i}{\partial w_j \partial w_k}$

- ▶ **Constancy** of NTK implies training dynamic of wide NNs is approximately a **linear** model \Rightarrow **linear convergence** of gradient flow (hence GD)
- ▶ It can be shown that [Liu, et al., arXiv:2003:00307]:

$$|\lambda_{\min}(\mathbf{K}(\mathbf{w})) - \lambda_{\min}(\mathbf{K}(\mathbf{w}_0))| < O\left(\sup_{\mathbf{w} \in \mathcal{B}} \|\mathbf{H}_{\mathcal{F}}(\mathbf{w})\|\right) = O(1/\sqrt{d})$$

Thus, the PL* condition holds for single-layer wide NN

NTK:

High-level intuition:

- 1° If loss is convex, then GD (SGD) converges to global min.
 - 2° Linear/kernel model exhibit convex loss landscape.
 - 3° We'll prove wide NN, landscape looks like a kernel model.
- $3^\circ \rightarrow 2^\circ \rightarrow 1^\circ$

1) Gradient dynamic for linear models:

Dataset $\{(x_i, y_i)\}_{i=1}^N$ $u_i = \underline{w}^T x_i$ $x_i \in \mathbb{R}^d$

$$\underline{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_N \end{bmatrix} \in \mathbb{R}^N \quad \underline{X} = \begin{bmatrix} \cdots & x_1^T & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & x_N^T & \cdots \end{bmatrix}_{N \times d} \quad \underline{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} \quad d \gg N$$

Square loss: $L(\underline{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - x_i^T \underline{w})^2 = \frac{1}{2} \|\underline{y} - \underline{X} \underline{w}\|^2$

Optimize via GD: $\underline{w}^T = \underline{w} - s \nabla L(\underline{w})$

$$\nabla_{\underline{w}} L(\underline{w}) = - \sum_{i=1}^N x_i (y_i - u_i) = - \underline{X}^T (\underline{y} - \underline{u})$$

$s \rightarrow 0$: $\frac{d\underline{w}}{dt} = -\nabla L(\underline{w}) = \underline{X}^T (\underline{y} - \underline{u})$ (ODE for \underline{w} evolution)

Gradient flow.

GD: a finite-time discretization of this ODE

$$\frac{d\underline{u}}{dt} = \frac{d(\underline{X} \underline{w})}{dt} = \underline{X} \frac{d\underline{w}}{dt} = \underbrace{\underline{X} \underline{X}^T}_{\underline{K}} (\underline{y} - \underline{u}) = \underline{K} (\underline{y} - \underline{u}).$$

Remarks:

1° Linear ODE can be solved in closed form. Let $\underline{r} = y - \underline{u}$

$$\frac{d\underline{r}}{dt} = \frac{d(y - \underline{u})}{dt} = -\frac{d\underline{u}}{dt} = -\underline{K}\underline{r}$$

$$\underline{r}(t) = \exp(-\underline{K}t) \underline{r}(0)$$

If \underline{K} is full rank, $\lambda_{\min}(\underline{K}) > 0$.

GD converges exp. fast $\rightarrow 0$ loss.

2° $\underline{K} = \underline{X}\underline{X}^T$ const. Configuration of data pts with $\lambda_{\min}(\underline{K})$ allows GD to converge fast.

3° All that matters is set of pair wise product $[\underline{K}]_{ij} = \underline{x}_i^T \underline{x}_j$
"kernel trick". kernel fn:

$$[\underline{K}]_{ij} = \langle \phi(\underline{x}_i), \phi(\underline{x}_j) \rangle, \text{ where } \phi \text{ is "feature map"}$$

2). General dynamics for non-linear model: $f(\underline{w})$.

For \underline{x}_i , $u_i = f(\underline{w}, \underline{x}_i)$.

$$\text{Square loss: } L(\underline{w}) = \frac{1}{2} \sum_{j=1}^N (y_j - \underbrace{f(\underline{w}, \underline{x}_j)}_{u_j})^2$$

The grad w.r.t. any one weight parameter:

$$\nabla_{w_i} L(\underline{w}) = - \sum_{j=1}^N \frac{\partial f(\underline{w}, \underline{x}_j)}{\partial w_i} (y_j - u_j)$$

$$\frac{du_i}{dt} = \sum_{k=1}^d \frac{\partial u_i}{\partial w_k} \frac{dw_k}{dt} = \sum_{k=1}^d \frac{\partial f(\underline{w}, \underline{x}_i)}{\partial w_k} \frac{dw_k}{dt}$$

$$= \sum_{k=1}^d \frac{\partial f(\underline{w}, \underline{x}_i)}{\partial w_k} \left[\sum_{j=1}^N \frac{\partial f(\underline{w}, \underline{x}_j)}{\partial w_k} (y_j - u_j) \right]$$

$$= \sum_{j=1}^N \left\langle \frac{\partial f(\underline{w}, \underline{x}_i)}{\partial \underline{w}}, \frac{\partial f(\underline{w}, \underline{x}_j)}{\partial \underline{w}} \right\rangle (y_j - u_j)$$

$$= \sum_{j=1}^N [\underline{K}]_{ij} (y_j - u_j), \quad \text{where } [\underline{K}]_{ij} = \left\langle \frac{\partial f(\underline{w}, \underline{x}_i)}{\partial \underline{w}}, \frac{\partial f(\underline{w}, \underline{x}_j)}{\partial \underline{w}} \right\rangle$$

$$= \sum_{k=1}^d \frac{\partial f(\underline{w}, \underline{x}_i)}{\partial w_k} \cdot \frac{\partial f(\underline{w}, \underline{x}_j)}{\partial w_k}$$

$$\underline{K}_t = DF(\underline{w}_t) DF^T(\underline{w}_t). \quad \text{kernel matrix.}$$

$$\frac{d\underline{u}}{dt} = -\underline{K}_t (\underline{y} - \underline{u}). \quad \leftarrow \text{nonlinear ODE.}$$

1° $\underline{K}_t \geq 0, \forall t.$ kernel mapping: $\phi: \underline{x} \mapsto \frac{\partial f(\underline{w}, \underline{x})}{\partial \underline{w}} \in \mathbb{R}^d$

If f is NN, then ϕ is "NTK".

3) Wide NN exhibits linear model dynamics. [Du, ICLR'19]

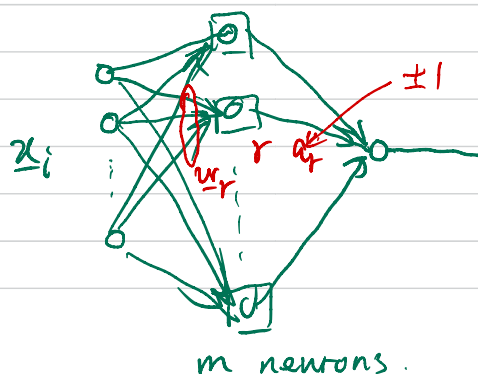
1° Randomly initialize \underline{w} at $t=0$.

2° At $t=0$, we'll show NTK \underline{K}_0 is full rank.

3° For wide NNs, $\underline{K}_t \approx \underline{K}_0$, hence \underline{K}_t is full rank $\forall t$.

Consider 2-layer NN w/ m hidden neurons, with twice diff'ble ψ activation fn.

Fix 2nd layer, only train 1st layer.



$$f(\underline{w}, \underline{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \psi(\langle \underline{w}_r, \underline{x} \rangle), \quad a_r = \pm 1$$

Initialize $[\underline{w}_1(0) \dots \underline{w}_m(0)]^T$ stand normal distr.

$$\frac{\partial f(\underline{w}(0), \underline{x}_i)}{\partial \underline{w}_r} = \frac{1}{\sqrt{m}} a_r \underline{x}_i \psi'(\langle \underline{w}_r(0), \underline{x}_i \rangle).$$

$$\begin{aligned} \text{So, NTK at } t=0: [\underline{K}]_{ij} &= \left\langle \frac{\partial f(\underline{w}, \underline{x}_i)}{\partial \underline{w}_r}, \frac{\partial f(\underline{w}, \underline{x}_j)}{\partial \underline{w}_r} \right\rangle \\ &= \underline{x}_i^T \underline{x}_j \left[\frac{1}{m} \sum_{r=1}^m a_r^2 \psi'(\langle \underline{w}_r(0), \underline{x}_i \rangle) \psi'(\langle \underline{w}_r(0), \underline{x}_j \rangle) \right] \end{aligned}$$

Each entry of $[\underline{K}]_{ij}$ is a r.v. with mean being equal to:

$$\underline{x}_i^T \underline{x}_j \mathbb{E}_{\underline{w} \sim N(0, \mathbb{I})} \psi'(\underline{x}_i^T \underline{w}) \psi'(\underline{x}_j^T \underline{w}) \triangleq [\underline{K}^*]_{ij}$$

As $m \rightarrow \infty$, NTK at $t=0$ is equal to \underline{K}^*

1° for $\varepsilon > 0$ if $m > \tilde{O}\left(\frac{N^4}{\varepsilon^2}\right)$, then $\|\underline{K}(0) - \underline{K}^*\| \leq \varepsilon$ w.h.p.

2° Suppose $y_i = \pm 1$, $w_i(\tau)$ bnded throughout training, $0 \leq \tau \leq t$.

for $\varepsilon > 0$, if $m \geq \tilde{O}\left(\frac{N^6 t^2}{\varepsilon^2}\right)$, then $\|\underline{K}(\tau) - \underline{K}^*\| \leq \varepsilon$, w.h.p.

Remarks:

(1). width scales play w.r.t. N . [Song et al. NeurIPS'21] $\tilde{O}(N^{\frac{3}{2}})$

(2). Dependence on data: [Nguyen et al.]: $O(Nd)$

(3). For L -layer NN, widths need to scale as $\text{poly}(N, L)$.

Next Class

First-Order Methods under Additional Assumptions