# Loss landscapes and optimization in over-parameterized non-linear systems and neural networks

[Liu, Zhu, Belkin 2022]

Presenter: Shirley
Xiaoqi Liu
27 Nov 2024

---

- Karimi, Nutini, Schmidt 16
- Liu, Zhu, Belkin 21
- Frei, Muthukumar, Yang 23 (Neurips tutorial)
- Frei 22 (Tutorial at Simons)

# Today's plan

# Today's plan

① Introduce Polyak–Lojasiewicz (PL) condition

# Today's plan

① Introduce Polyak–Lojasiewicz (PL) condition

② Can we establish PL for large non-linear systems?

- What do their loss landscapes look like?
- Whether PL is satisfied depends on the conditioning of the tangent kernel

# Today's plan

① Introduce Polyak–Lojasiewicz (PL) condition

② Can we establish PL for large non-linear systems?

- What do their loss landscapes look like?
- Whether PL is satisfied depends on the conditioning of the tangent kernel

③ Application to neural networks

④ Summary and extensions

# Today's plan

① Introduce Polyak-Lojasiewicz (PL) condition

② Can we establish PL for large non-linear systems?

- What do their loss landscapes look like?
- Whether PL is satisfied depends on the conditioning of the tangent kernel

③ Application to neural networks

④ Summary and extensions

Try to answer:

Why over-parameterized NNs are non-convex yet easy to optimize?

① Introduce Polyak–Lojasiewicz (PL) condition

- Strong convexity (SC)  $L: \mathbb{R}^m \to \mathbb{R}$

$$L(y) \geq L(x) + \langle \nabla L(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad \forall x, y \in \mathbb{R}^m$$

- Strong convexity (SC)  $\qquad L : \mathbb{R}^m \to \mathbb{R}$

$$L(y) \geq L(x) + \langle \nabla L(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad \forall x, y \in \mathbb{R}^m$$

- Polyak-Lojasiewicz (PL) condition

for some $\mu > 0$

$$\frac{1}{2} \|\nabla L(x)\|^2 \geq \mu (L(x) - \underbrace{L(x^*)}_{\text{global minimum } (L^*)}) \quad \forall x \in \mathbb{R}^m$$

- Strong convexity (SC)  $L: \mathbb{R}^m \to \mathbb{R}$

$$L(y) \geq L(x) + \langle \nabla L(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad \forall x, y \in \mathbb{R}^m$$

- Polyak-Lojasiewicz (PL) condition

  for some $\mu > 0$

$$\frac{1}{2} \|\nabla L(x)\|^2 \geq \mu(L(x) - \underbrace{L(x^*)}_{\text{global minimum } (L^*)}) \quad \forall x \in \mathbb{R}^m$$

- Strong convexity (SC)     $L : \mathbb{R}^m \to \mathbb{R}$

$$L(y) \geq L(x) + \langle \nabla L(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad \forall x, y \in \mathbb{R}^m$$

- Polyak-Lojasiewicz (PL) condition

for some $\mu > 0$

$$\frac{1}{2} \|\nabla L(x)\|^2 \geq \mu \left( L(x) - \underbrace{L(x^*)}_{\text{global minimum } (L^*)} \right) \quad \forall x \in \mathbb{R}^m$$

---

Remarks:

- $SC \implies PL$
- All stationary points are global minima
- PL is somewhat easier to verify than SC

# Smoothness + PL $\Rightarrow$ Linear convergence of gradient descent

Recall: $\mu$-PL gives $\frac{1}{2}\|\nabla L(x)\|^2 \geq \mu\left(L(x)-L(x^*)\right)$

GD    $x^{k+1} = x^k - \eta \nabla L(x^k)$

Recall: $\mu$-PL gives $\frac{1}{2}\|\nabla L(x)\|^2 \geq \mu\left(L(x) - L(x^*)\right)$

$\qquad$ GD $\quad x^{k+1} = x^k - \eta \nabla L(x^k)$

Proof :

# Smoothness + PL $\Rightarrow$ Linear convergence of gradient descent

Recall: $\mu$-PL gives $\frac{1}{2}\|\nabla L(x)\|^2 \geq \mu\left(L(x) - L(x^*)\right)$

GD $\quad x^{k+1} = x^k - \eta \nabla L(x^k)$

Proof:

$$L(x^{k+1}) - L(x^*)$$

Recall:    $\mu$-PL gives $\frac{1}{2}\|\nabla L(x)\|^2 \geq \mu\left(L(x) - L(x^*)\right)$

GD    $x^{k+1} = x^k - \eta \nabla L(x^k)$

Proof:

$L(x^{k+1}) - L(x^*)$

$\beta$-smoothness

$\leq L(x^k) + \langle \nabla L(x^k), x^{k+1} - x^k \rangle + \frac{1}{2}\beta\|x^{k+1} - x^k\|^2 - L(x^*)$

$= L(x^k) + \langle \nabla L(x^k), -\eta \nabla L(x^k)\rangle + \frac{1}{2}\beta\|\eta\nabla L(x^k)\|^2 - L(x^*)$

Recall: $\mu$-PL gives $\frac{1}{2}\|\nabla L(x)\|^2 \geq \mu\left(L(x) - L(x^*)\right)$

$\qquad$ GD $\quad x^{k+1} = x^k - \eta \nabla L(x^k)$

Proof:

$\qquad L(x^{k+1}) - L(x^*)$

$\beta$-smoothness

$\qquad \leq L(x^k) + \langle \nabla L(x^k), x^{k+1} - x^k \rangle + \frac{1}{2}\beta \|x^{k+1} - x^k\|^2 - L(x^*)$

$\qquad = L(x^k) + \langle \nabla L(x^k), -\eta \nabla L(x^k) \rangle + \frac{1}{2}\beta \| \eta \nabla L(x^k)\|^2 - L(x^*)$

$\eta = 1/\beta$

$\qquad = L(x^k) - \frac{1}{2}\eta \|\nabla L(x^k)\|^2 - L(x^*)$

Recall: $\mu$-PL gives $\frac{1}{2} \|\nabla L(x)\|^2 \geq \mu \left( L(x) - L(x^*) \right)$

$\qquad$ GD $\quad x^{k+1} = x^k - \eta \nabla L(x^k)$

Proof:

$\qquad L(x^{k+1}) - L(x^*)$

$\beta$-smoothness

$\qquad \leq L(x^k) + \langle \nabla L(x^k), x^{k+1} - x^k \rangle + \frac{1}{2} \beta \|x^{k+1} - x^k\|^2 - L(x^*)$

$\qquad = L(x^k) + \langle \nabla L(x^k), -\eta \nabla L(x^k) \rangle + \frac{1}{2} \beta \| \eta \nabla L(x^k) \|^2 - L(x^*)$

$\eta = 1/\beta$
$\qquad = L(x^k) - \frac{1}{2} \eta \|\nabla L(x^k)\|^2 - L(x^*)$

$\mu$-PL
$\qquad \leq L(x^k) - \eta \mu \left( L(x^k) - L(x^*) \right) - L(x^*)$

$\qquad = (1 - \eta \mu) \left( L(x^k) - L(x^*) \right) \quad \square$

Aside:

# Aside:

- Similar proofs for stochastic gradient descent, randomized coordinate descent, greedy coordinate descent...

# Aside:

- Similar proofs for stochastic gradient descent, randomized coordinate descent, greedy coordinate descent...

- Proximal gradient generalization    (via proximal-PL)

$$\min_{x} \tilde{L}(x) \equiv \underbrace{L(x)}_{\beta\text{-smooth}} + \underbrace{g(x)}_{\substack{\text{non-smooth} \\ \text{convex}}}$$

# Aside:

- Similar proofs for stochastic gradient descent, randomized coordinate descent, greedy coordinate descent...

- Proximal gradient generalization    (via proximal-PL)

$$\min_{x} \tilde{L}(x) \equiv \underbrace{L(x)}_{\beta\text{-smooth}} + \underbrace{g(x)}_{\substack{\text{non-smooth} \\ \text{convex}}}$$

- Other conditions for obtaining linear convergence:

    - Weak SC: $L(x^*) \geq L(x) + \langle \nabla L(x), x^* - x \rangle + \frac{\mu}{2} \| x^* - x \|^2$

    - Quadratic Growth: $L(x) - L(x^*) \geq \frac{\mu}{2} \| x^* - x \|^2$

        $\vdots$

② Can we establish PL for large non-linear systems?

# Setting

# Setting

- Training data $\{x_i, y_i\}_{i=1}^{n}$ $\quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$

# Setting

- Training data $\{x_i, y_i\}_{i=1}^{n}$    $x_i \in \mathbb{R}^d$    $y_i \in \mathbb{R}$
- A parametric model $f(w; x_i)$    e.g. a neural network

# Setting

- Training data $\{x_i, y_i\}_{i=1}^{n}$  $x_i \in \mathbb{R}^d$  $y_i \in \mathbb{R}$

- A parametric model $f(w; x_i)$  e.g. a neural network

  Define shorthand:  $F(w) = \begin{bmatrix} f(w; x_1) \\ f(w; x_2) \\ \vdots \\ f(w; x_n) \end{bmatrix} \in \mathbb{R}^n$  $F: \mathbb{R}^m \to \mathbb{R}^n$

  $$\nabla F = \frac{\partial F}{\partial w} \in \mathbb{R}^{n \times m}$$

  $$H_{F_i} = \frac{\partial^2 F_i}{\partial w^2} \in \mathbb{R}^{m \times m}$$

# Setting

- Training data $\{x_i, y_i\}_{i=1}^{n}$ $\quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$

- A parametric model $f(w; x_i)$ $\quad$ e.g. a neural network

  Define shorthand: $F(w) = \begin{bmatrix} f(w; x_1) \\ f(w; x_2) \\ \vdots \\ f(w; x_n) \end{bmatrix} \in \mathbb{R}^n \qquad F: \mathbb{R}^m \rightarrow \mathbb{R}^n$

  $$\nabla F = \frac{\partial F}{\partial w} \in \mathbb{R}^{n \times m}$$

  $$H_{F_i} = \frac{\partial^2 F_i}{\partial w^2} \in \mathbb{R}^{m \times m}$$

  $\beta_F$-smooth, $L_F$-Lipschitz

# Setting

- Training data $\{x_i, y_i\}_{i=1}^n$    $x_i \in \mathbb{R}^d$    $y_i \in \mathbb{R}$
- A parametric model $f(w; x_i)$    e.g. a neural network

Define shorthand: $F(w) = \begin{bmatrix} f(w; x_1) \\ f(w; x_2) \\ \vdots \\ f(w; x_n) \end{bmatrix} \in \mathbb{R}^n$    $F: \mathbb{R}^m \to \mathbb{R}^n$

$$\nabla F = \frac{\partial F}{\partial w} \in \mathbb{R}^{n \times m}$$

$$H_{F_i} = \frac{\partial^2 F_i}{\partial w^2} \in \mathbb{R}^{m \times m}$$

$\beta_F$ - smooth, $L_F$ - Lipschitz

- Loss function for optimization

$$L(w) \overset{e.g.}{=} \frac{1}{2} \| F(w) - y \|^2$$

$$H_L = \frac{\partial^2 L}{\partial w^2} \in \mathbb{R}^{m \times m}$$

# Loss landscapes

# Loss landscapes
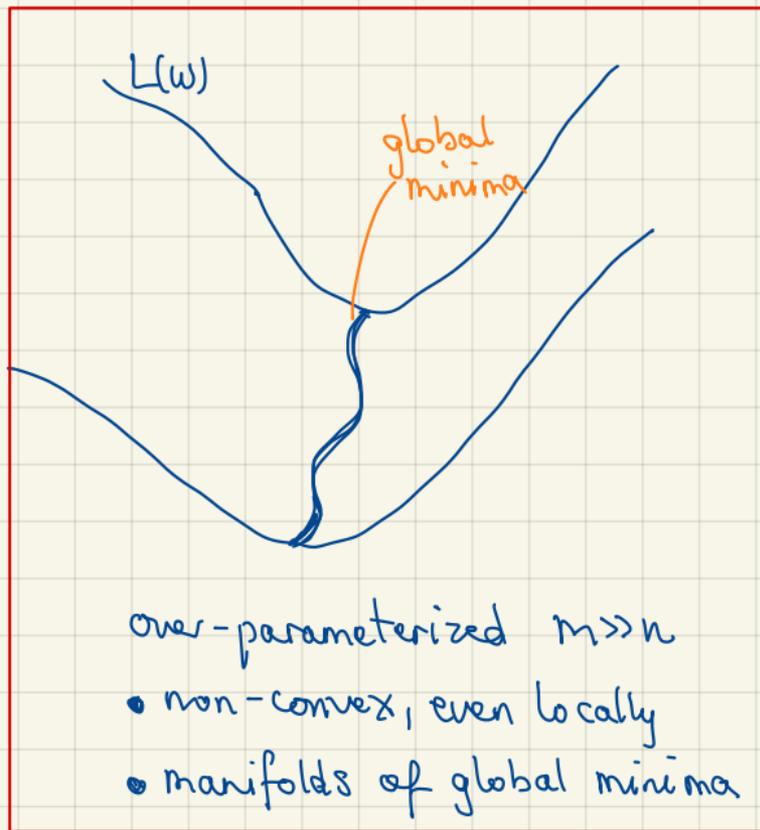


$L(\omega)$

under – parameterized $m \subset n$

isolated local minima

# Loss landscapes

$L(w)$

$L(w)$

global minima

under-parameterized $m < n$
isolated local minima

over-parameterized $m \gg n$
- non-convex, even locally
- manifolds of global minima

# Loss landscape in overparameterized case.

# Loss landscape in overparameterized case.

Prop 1. (Local non-convexity)  — can be generalized

Consider $L(w) = \frac{1}{2} \| F(w) - y \|^2$. Suppose $\nabla F(w^*) \neq 0$ and a mild technical assumption on $\{ H_{F_i} \}_{i=1}^n$, then $L(w)$ isn't convex in any neighbourhood of $w^*$

# Loss landscape in overparameterized case.

Prop 1. (Local non-convexity).    — can be generalized
Consider $L(w) = \frac{1}{2} \| F(w) - y \|^2$. Suppose $\nabla F(w^*) \neq 0$ and
a mild technical assumption on $\{H_{F_i}\}_{i=1}^n$, then
$L(w)$ isn't convex in any neighbourhood of $w^*$

Proof idea:

# Loss landscape in overparameterized case.

Prop 1. (Local non-convexity)  — can be generalized
Consider $L(w) = \frac{1}{2} \| F(w) - y \|^2$. Suppose $\nabla F(w^*) \neq 0$ and
a mild technical assumption on $\{ H_{F_i} \}_{i=1}^n$, then
$L(w)$ isn't convex in any neighbourhood of $w^*$

Proof idea:

- $H_L(w) = \nabla F(w)^T \dfrac{\partial^2 L}{\partial F^2} \nabla F(w) + \displaystyle\sum_{i=1}^n \left( F(w) - y \right)_i H_{F_i}(w)$

# Loss landscape in overparameterized case.

Prop 1. (Local non-convexity) — can be generalized
Consider $L(w) = \frac{1}{2}\|F(w) - y\|^2$. Suppose $\nabla F(w^*) \neq 0$ and
a mild technical assumption on $\{H_{F_i}\}_{i=1}^n$, then
$L(w)$ isn't convex in any neighbourhood of $w^*$

Proof idea:

- $\underbrace{H_L(w)}_{m \times m} = \underbrace{\nabla F(w)^T}_{m \times n} \; \underbrace{\frac{\partial^2 L}{\partial F^2}}_{n \times n} \; \underbrace{\nabla F(w)}_{n \times m} + \sum_{i=1}^n \underbrace{\left(F(w) - y\right)_i}_{\text{scalar}} \underbrace{H_{F_i}(w)}_{m \times m}$

for $m \gg n$, first term is generally low-rank, but second term isn't

# Loss landscape in overparameterized case.

Prop 1. (Local non-convexity) ← can be generalized
Consider $L(w) = \frac{1}{2} \| F(w) - y \|^2$. Suppose $\nabla F(w^*) \neq 0$ and
a mild technical assumption on $\{H_{F_i}\}_{i=1}^{n}$, then
$L(w)$ isn't convex in any neighbourhood of $w^*$

Proof idea:

- $\underbrace{H_L(w)}_{m \times m} = \underbrace{\nabla F(w)^T}_{m \times n} \underbrace{\frac{\partial^2 L}{\partial F^2}}_{n \times n} \underbrace{\nabla F(w)}_{n \times m} + \sum_{i=1}^{n} \underbrace{\left( F(w) - y \right)_i}_{\text{scalar}} \underbrace{H_{F_i}(w)}_{m \times m}$

  for $m \gg n$, first term is generally low-rank, but second term isn't

- Consider the special case where $n=1$, we have
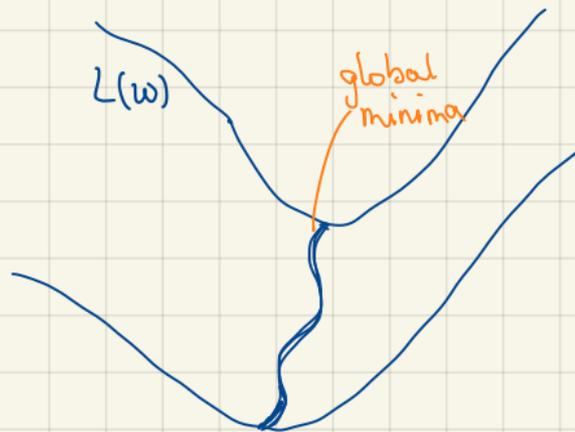
  $H_L(w) = \underbrace{\nabla F(w)^T \nabla F(w)}_{\text{rank-1}} + \underbrace{\left( F(w) - y \right)}_{=0 \text{ at } w = w^*} H_F(w)$   Assume > rank-1

# Loss landscape in overparameterized case.

Prop 1. (Local non-convexity)

— can be generalized

Consider $L(w) = \frac{1}{2}\|F(w) - y\|^2$. Suppose $\nabla F(w^*) \neq 0$ and a mild technical assumption on $\{H_{F_i}\}_{i=1}^n$, then $L(w)$ isn't convex in any neighbourhood of $w^*$

Proof idea:

- $\underbrace{H_L(w)}_{m \times m} = \underbrace{\nabla F(w)^T}_{m \times n} \underbrace{\frac{\partial^2 L}{\partial F^2}}_{n \times n} \underbrace{\nabla F(w)}_{n \times m} + \sum_{i=1}^n \underbrace{\left(F(w) - y\right)_i}_{\text{scalar}} \underbrace{H_{F_i}(w)}_{m \times m}$

for $m \gg n$, first term is generally low-rank, but second term isn't

- Consider the special case where $n=1$, we have

$H_L(w) = \underbrace{\nabla F(w)^T \nabla F(w)}_{\text{rank-1}} + \underbrace{\left(F(w) - y\right)}_{=0 \text{ at } w = w^*} \underbrace{H_F(w)}_{\text{Assume} > \text{rank-1}}$

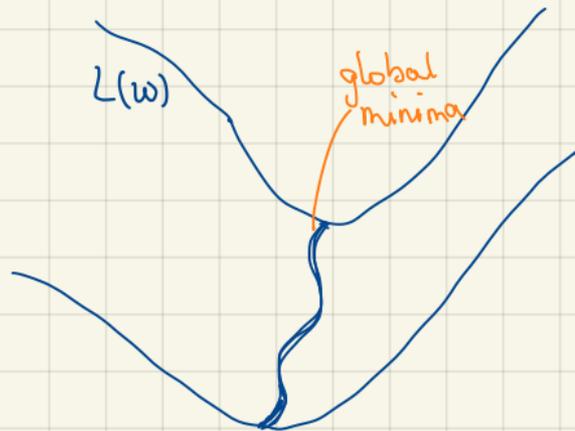$H_L(w+\delta)$ and/or $H_L(w-\delta)$ are not PSD $\therefore$ no local convexity $\square$

# Convexity → PL

$L(w)$

global minima

# Convexity → PL



$L(w)$

global minima

- Recall: PL requires
$$\frac{1}{2}\|\nabla L(w)\|^2 \geq \mu\left(L(w) - L(w^*)\right) \qquad \forall w \in \mathbb{R}^m$$

global optimum

# Convexity → PL → Local PL



$L(w)$

global minima

- Recall: PL requires **global optimum**

$$\frac{1}{2}\|\nabla L(w)\|^2 \geq \mu\left(L(w) - L(w^*)\right) \qquad \forall w \in \mathbb{R}^m$$

- Local PL / PL$^*$:

  A non-negative function $L$ satisfies $\mu$-PL$^*$ on a set $S \subset \mathbb{R}^m$ if

  $$\|\nabla L(w)\|^2 \geq \mu L(w) \qquad \forall w \in S$$

Can we get PL* for our overparam. system ? 1/23

# Can we get PL* for our overparam. system?

Theorem 1 (Uniform conditioning $\Rightarrow$ PL*)

    The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

    on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

**Theorem 1** (Uniform conditioning $\Rightarrow$ PL*)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

Theorem 1 (Uniform conditioning $\Rightarrow$ $PL^*$)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu - PL^*$

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left((F(w)-y)\,\nabla F(w)\,\nabla F(w)^T\,(F(w)-y)^T\right)$$

# Can we get PL* for our overparam. system? 

Theorem 1 (Uniform conditioning $\Rightarrow$ PL*)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left(\underbrace{(F(w) - y)}_{\text{len}-n \text{ vector}} \underbrace{\nabla F(w)}_{n \times m} \underbrace{\nabla F(w)^T}_{m \times n} (F(w) - y)^T\right)$$

**Theorem 1** (Uniform conditioning ⇒ PL*)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$K(w)$: tangent kernel

$$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left( \underbrace{(F(w)-y)}_{len-n\ vector} \underbrace{\nabla F(w)}_{n \times m} \underbrace{\nabla F(w)^\top}_{m \times n} (F(w)-y)^\top \right)$$

Theorem 1 (Uniform conditioning $\Rightarrow$ PL*)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left( (F(w) - y) \underbrace{\nabla F(w)}_{n \times m} \underbrace{\nabla F(w)^T}_{m \times n} (F(w) - y)^T \right)$

$\underbrace{(F(w) - y)}_{\text{len-}n \text{ vector}}$

$K(w)$: tangent kernel

$\geq \frac{1}{2} \lambda_{min}(K(w)) \|F(w) - y\|^2$

# Can we get PL* for our overparam. system?

> **Theorem 1** (Uniform conditioning $\Rightarrow$ PL*)
> The square loss $L(w) = \frac{1}{2}\|F(w)-y\|^2$ satisfies $\mu$-PL*
> on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left((F(w)-y)\underbrace{\nabla F(w)\,\nabla F(w)^T}_{K(w):\ \text{tangent kernel}}(F(w)-y)^T\right)$$

$\underbrace{(F(w)-y)}_{\text{len-}n\ \text{vector}}$  $\underbrace{\nabla F(w)}_{n\times m}$  $\underbrace{\nabla F(w)^T}_{m\times n}$

$$\geq \frac{1}{2}\lambda_{min}(K(w))\|F(w)-y\|^2 \geq \mu L(w) \quad \square$$

Theorem 1 (Uniform conditioning ⇒ PL*)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$K(w)$ : tangent kernel

$$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left((F(w) - y)\underbrace{\nabla F(w) \nabla F(w)^T}_{}(F(w) - y)^T\right)$$

$(F(w) - y)$ : len-n vector $\quad \nabla F(w)$ : n×m $\quad \nabla F(w)^T$ : m×n

$$\geq \frac{1}{2}\lambda_{min}(K(w))\|F(w) - y\|^2 \geq \mu L(w) \quad \square$$

Theorem 1 (Uniform conditioning $\Rightarrow$ PL*)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$K(w)$: tangent kernel

$$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left(\underbrace{(F(w) - y)}_{\text{len-}n \text{ vector}} \underbrace{\nabla F(w)}_{n \times m} \underbrace{\nabla F(w)^T}_{m \times n} (F(w) - y]^T\right)$$

$$\geq \frac{1}{2} \lambda_{min}(K(w)) \|F(w) - y\|^2 \geq \mu L(w) \quad \square$$

① For Gaussian random matrix $M \in \mathbb{R}^{n \times m}$, the condition number $\kappa$ of $MM^T$ satisfies $\mathbb{E}[\log \kappa] \sim \log \frac{m}{|m - n| + 1}$ [Chen, Dongarra 08]

# Can we get PL* for our overparam. system?

**Theorem 1** (Uniform conditioning $\Rightarrow$ PL*)

The square loss $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL*

on $S \subset \mathbb{R}^m$ if $\lambda_{min}(K(w)) \geq \mu$ on $S$

Proof:

$$\frac{1}{2}\|\nabla L(w)\|^2 = \frac{1}{2}\left(\underbrace{(F(w) - y)}_{\text{len-}n\text{ vector}}\ \overbrace{\underbrace{\nabla F(w)}_{n\times m}\ \underbrace{\nabla F(w)^T}_{m\times n}}^{K(w):\text{ tangent kernel}}\ (F(w) - y)^T\right)$$

$$\geq \frac{1}{2}\ \lambda_{min}(K(w))\ \|F(w) - y\|^2 \geq \mu L(w) \quad \square$$

① For Gaussian random matrix $M \in \mathbb{R}^{n\times m}$, the condition number $\kappa$
of $MM^T$ satisfies $\mathbb{E}[\log \kappa] \sim \log \frac{m}{|m-n|+1}$ [Chen, Dongarra 05]

② Standard initialization $W_0 \overset{iid}{\sim} N(0,1)$, for wide NN $F(w)$,
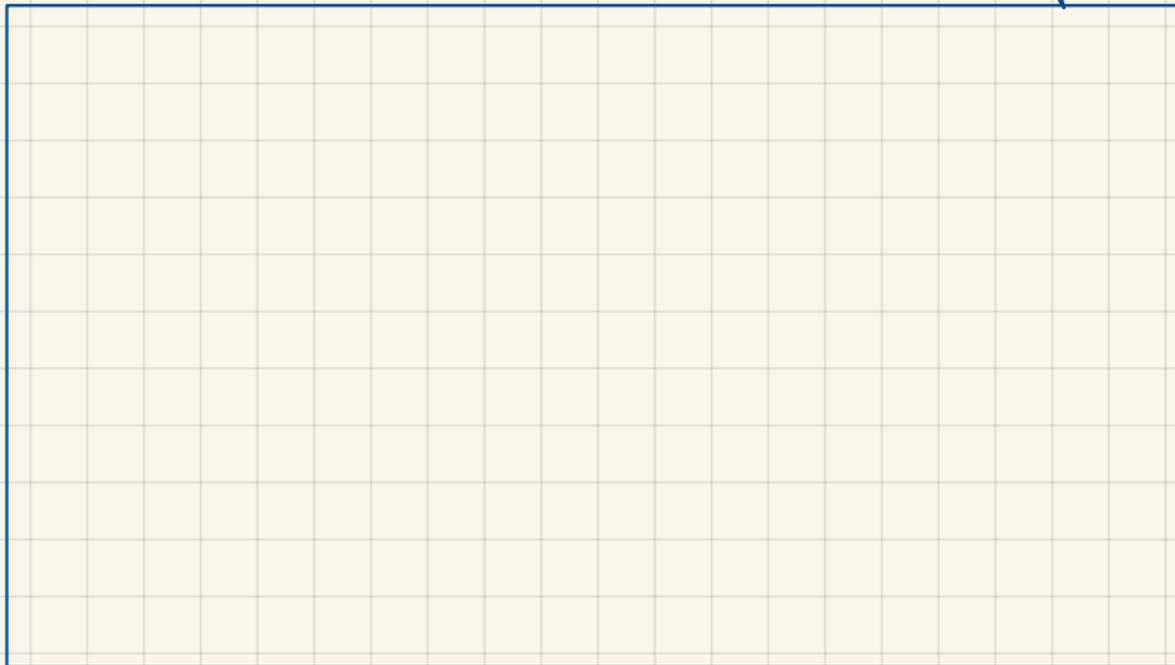$\lambda_{min}(K(w_0)) = O(1)$ w.h.p. assuming data is not degenerate
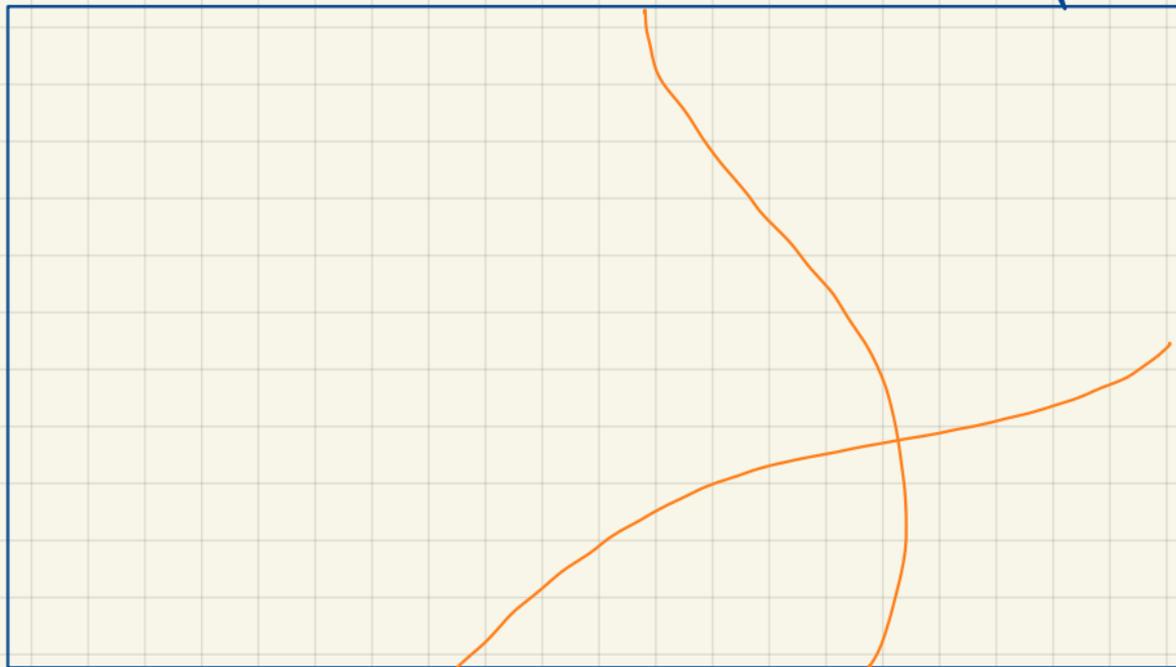
# Picture so far

# Picture so far
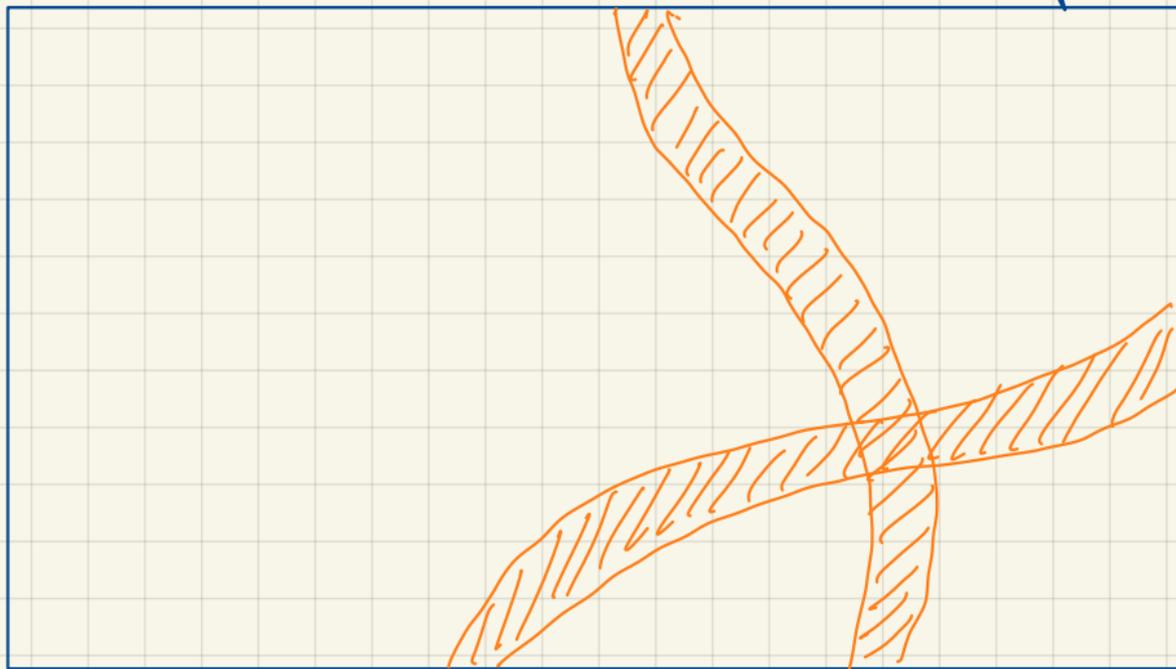
Parameter space $W \in \mathbb{R}^m$

# Picture so far

$\lambda_{min}(K(w)) = 0$ Parameter space $w \in \mathbb{R}^m$

# Picture so far



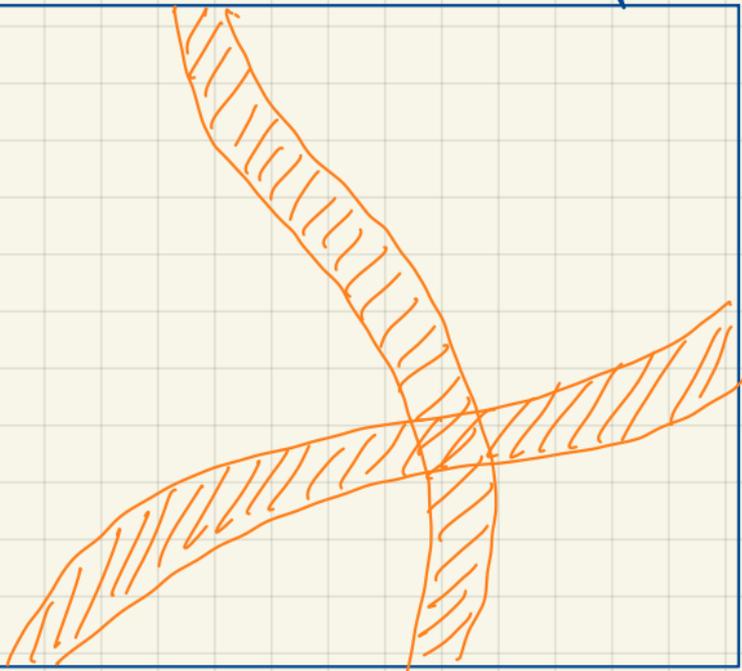$\lambda_{min}(K(w)) \leq \mu$  Parameter space $w \in \mathbb{R}^m$

# Picture so far
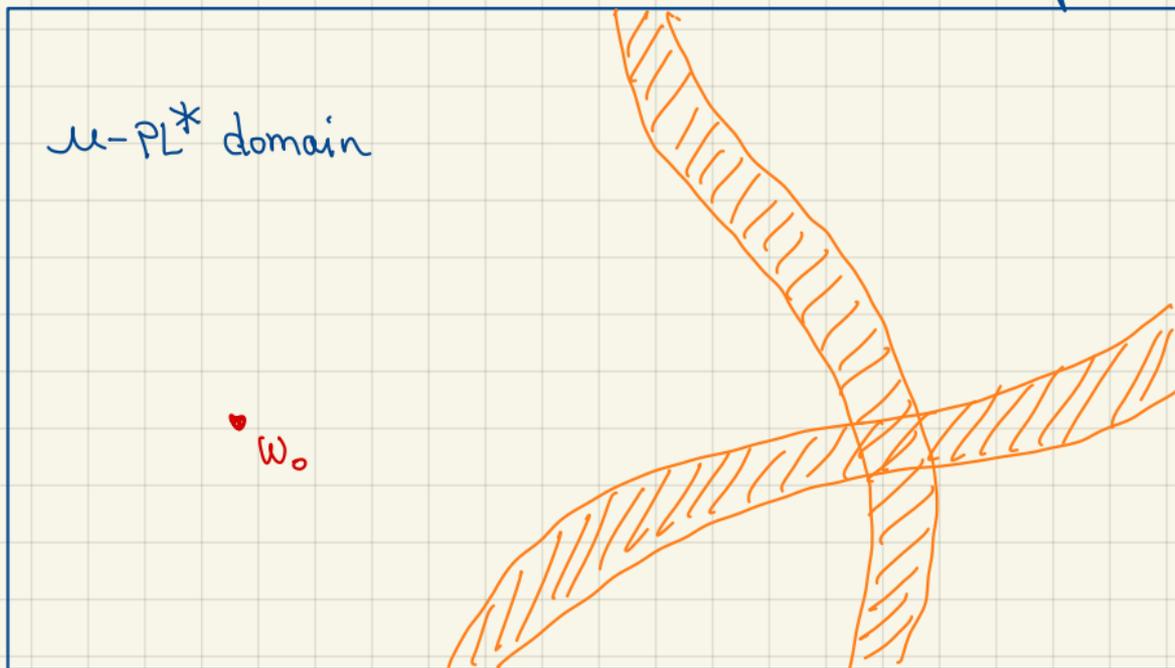
$\lambda_{min}(K(w)) \leq \mu$ Parameter space $w \in \mathbb{R}^m$

$\mu$-PL$^*$ domain

# Picture so far

$\lambda_{min}(K(w)) < \mu$ Parameter space $w \in \mathbb{R}^m$

$\mu$-PL* domain

$w_0$

w.h.p. the initializer $w_0$ is good (i.e. $\lambda_{min}(K(w_0)) \geq \mu$)

# Picture so far



$\lambda_{min}(K(w)) < \mu$ Parameter space $w \in \mathbb{R}^m$

$\mu$-PL$^*$ domain

$w_t$

$w_o$

w.h.p. the initializer $w_o$ is good (i.e. $\lambda_{min}(K(w_o)) \geq \mu$)

- Next steps: want to show $w \in B(w_o, R)$ is good too, and adjust $R$ to cover training trajectory in the ball.

Theorem 2 ( Small Hessian $\Rightarrow$ PL$^*$ in $B(w_0, R)$ )

Theorem 2 ( Small Hessian $\Rightarrow$ PL$^*$ in $B(w_0, R)$)

Suppose $\lambda_{min}(K(w_0)) = \lambda_0 > 0$. If $\boxed{\|H_{F_i}(w)\|_2 \leq \dfrac{\lambda_0 - \mu}{2 L_F \sqrt{n} R}}$

for $w \in B(w_0, R)$, then $K(w)$ is $\mu$-uniformly conditioned in the ball.

Hence, $L(w) = \frac{1}{2}\| F(w) - y\|^2$ satisfies $\mu$-PL$^*$ for $w \in B(w_0, R)$.

Theorem 2 ( Small Hessian $\Rightarrow$ PL$^*$ in $B(w_0, R)$ )

Suppose $\lambda_{min}(K(w_0)) = \lambda_0 > 0$. If $\boxed{\|H_{F_i}(w)\|_2 \leq \dfrac{\lambda_0 - \mu}{2 L_F \sqrt{n} R}}$

for $w \in B(w_0, R)$, then $K(w)$ is $\mu$-uniformly conditioned in the ball.

Hence, $L(w) = \frac{1}{2} \| F(w) - y \|^2$ satisfies $\mu$-PL$^*$ for $w \in B(w_0, R)$.

Proof idea:

Theorem 2 ( Small Hessian $\Rightarrow$ PL$^*$ in $B(w_0, R)$)

Suppose $\lambda_{min}(K(w_0)) = \lambda_0 > 0$. If $\boxed{\|H_{F_i}(w)\|_2 \leq \dfrac{\lambda_0 - \mu}{2L_F \sqrt{n} R}}$

for $w \in B(w_0, R)$, then $K(w)$ is $\mu$-uniformly conditioned in the ball.

Hence, $L(w) = \frac{1}{2}\|F(w) - y\|^2$ satisfies $\mu$-PL$^*$ for $w \in B(w_0, R)$.

Proof idea: fundamental theorem of calculus

$$\nabla F_i(w) = \nabla F_i(w_0) + \int_0^1 \underbrace{H_{F_i}(w_0 + \tau(w - w_0))(w - w_0)}_{\text{Cauchy - Schwartz} \cdots} d\tau$$

Theorem 2 ( Small Hessian $\Rightarrow$ PL$^*$ in $B(w_0, R)$ )

Suppose $\lambda_{min}(K(w_0)) = \lambda_0 > 0$. If $\boxed{\| H_{F_i}(w) \|_2 \leq \dfrac{\lambda_0 - \mu}{2 L_F \sqrt{n} R}}$

for $w \in B(w_0, R)$, then $K(w)$ is $\mu$-uniformly conditioned in the ball.

Hence, $L(w) = \frac{1}{2} \| F(w) - y \|^2$ satisfies $\mu$-PL$^*$ for $w \in B(w_0, R)$.

Proof idea: fundamental theorem of calculus

$$\nabla F_i(w) = \nabla F_i(w_0) + \int_0^1 \underbrace{H_{F_i}(w_0 + \tau(w - w_0))(w - w_0)}_{\text{Cauchy-Schwartz} \cdots} d\tau$$

- How do results so far apply to neural nets?

- When do we have small $\| H_{\bar{F}_i}(w) \|_2$?

③ Application to neural networks

# A shallow neural network.

# A shallow neural network.

- We'll show that for large enough width $m$, $\|H_{F_i}(w)\|_2$ is small

# A shallow neural network.

- We'll show that for large enough width $m$, $\|H_{F_i}(\omega)\|_2$ is small

- For simplicity, let's consider data $x_i \in \mathbb{R}$, $|x_i| \leq 1$ $\forall i \in [n]$

# A shallow neural network.

- We'll show that for large enough width m, $\|H_{F_i}(w)\|_2$ is small

- For simplicity, let's consider data $x_i \in \mathbb{R}$, $|x_i| \leq 1$ $\forall i \in [n]$

$$\bar{F_i}(w) = f(w; x_i) = \frac{1}{\sqrt{m}} \sum_{j=1}^{m} v_j \, \sigma(w_j x_i)$$

# A shallow neural network.

- We'll show that for large enough width $m$, $\|H_{F_i}(w)\|_2$ is small

- For simplicity, let's consider data $x_i \in \mathbb{R}$, $|x_i| \leq 1$ $\forall i \in [n]$

$$F_i(w) = f(w; x_i) = \frac{1}{\sqrt{m}} \sum_{j=1}^{m} v_j \, \sigma(w_j x_i)$$

- Fix $v_j \in \{1, -1\}$
- Initialize $w_j \sim N(0,1)$ then train
- $\sigma$: $\beta_0$-smooth, $L_0$-Lipschitz
  e.g. tanh, sigmoid

# A shallow neural network.

- We'll show that for large enough width m, $\| H_{F_i}(w) \|_2$ is small

- For simplicity, let's consider data $x_i \in \mathbb{R}$, $|x_i| \le 1$ $\forall i \in [n]$

$$F_i(w) = f(w; x_i) = \frac{1}{\sqrt{m}} \sum_{j=1}^{m} v_j \, \sigma(w_j x_i)$$

- Fix $v_j \in \{1, -1\}$
- Initialize $w_j \sim N(0,1)$ then train
- $\sigma$: $\beta_\sigma$-smooth, $L_\sigma$-Lipschitz
  e.g. tanh, sigmoid

① $[H_{F_i}(w)]_{jk} = \dfrac{\partial^2 f(w; x_i)}{\partial w_j \, \partial w_k} = \dfrac{1}{\sqrt{m}} v_j \, \sigma''(w_j x_i) \, x_i^2 \, \mathbb{1}\{j=k\}$

$\Rightarrow \| H_{F_i}(w) \|_2 \le \dfrac{1}{\sqrt{m}} \beta_\sigma = O\left(\dfrac{1}{\sqrt{m}}\right)$

② $\| \nabla F_i(w) \|^2 = \dfrac{1}{m} \sum_{j=1}^{m} x_i^2 \left( \sigma'(w_j x_i) \right)^2 \le L_\sigma^2 = \Theta(1)$

# Deeper neural nets

# Deeper neural nets

L-layer NN:
$$a^{(0)} = x$$
$$a^{(\ell)} = \sigma_\ell \left( \frac{1}{\sqrt{m_{\ell-1}}} W^{(\ell)} a^{(\ell-1)} \right) \qquad \forall \ \ell = 1, \cdots L+1$$
$$f(W; x) = a^{(L+1)}$$

# Deeper neural nets

$L$-layer NN :
$$a^{(0)} = x$$
$$a^{(\ell)} = \sigma_\ell \left( \frac{1}{\sqrt{m_{\ell-1}}} W^{(\ell)} a^{(\ell-1)} \right) \quad \forall\, \ell = 1, \cdots L+1$$
$$f(W; x) = a^{(L+1)}$$

Theorem 3 (Wide neural nets w. linear output layer satisfy PL*)

# Deeper neural nets

L-layer NN:
$$a^{(0)} = x$$
$$a^{(\ell)} = \sigma_\ell \left( \frac{1}{\sqrt{m_{\ell-1}}} W^{(\ell)} a^{(\ell-1)} \right) \quad \forall \, \ell = 1, \cdots L+1$$
$$f(W; x) = a^{(L+1)}$$

---

**Theorem 3** (Wide neural nets w. linear output layer satisfy PL*)

Consider $W_0^{(\ell)} \sim N(0, I) \quad \forall \, \ell \in [L+1]$, and $\sigma_{\ell+1}(z) = z$. Suppose $\lambda_0 = \lambda_{\min}(K(W_0)) > 0$. Then $L(w) = \frac{1}{2} \| F(w) - y \|^2$ satisfies $\mu$-PL* on $B(W_0, R)$ if

$$m = \tilde{\Omega} \left( \frac{n R^{6L+2}}{(\lambda_0 - \mu P^{-2})^2} \right)$$

# Summary so far

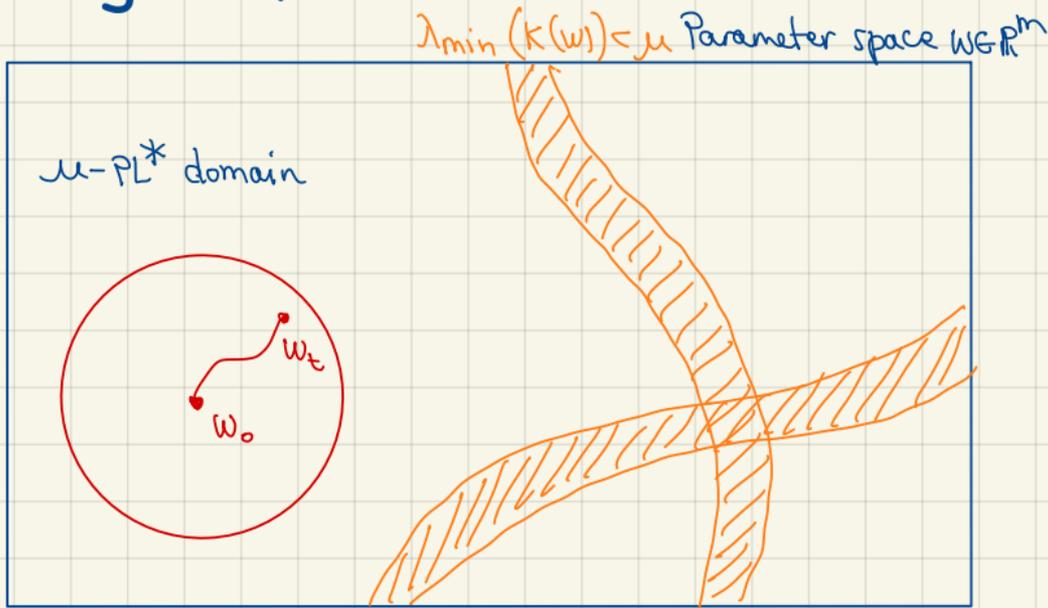# Summary so far



$\lambda_{min}(K(w)) < \mu$ Parameter space $w \in \mathbb{R}^m$

$\mu$-PL$^*$ domain

$w_t$

$w_0$

# Summary so far



$\lambda_{min}(K(w)) < \mu$ Parameter space $w \in \mathbb{R}^m$
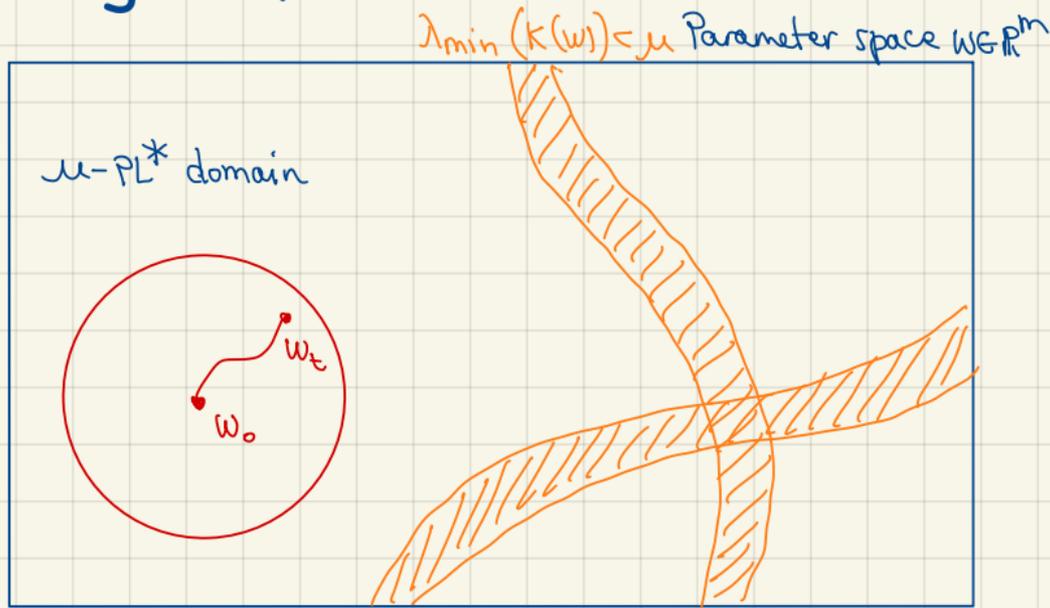
$\mu$-PL* domain

$w_t$

$w_0$

Neural net sufficiently wide

$\Rightarrow K(w) = \nabla F(w) \nabla F(w)^T$ well-conditioned in $B(w_0, R)$

$\Rightarrow L(w) = \frac{1}{2} \| F(w) - y \|^2$ satisfies PL* in $B(w_0, R)$

# Summary so far



$\lambda_{min}(K(w)) < \mu$ Parameter space $w \in \mathbb{R}^m$

$\mu$-PL$^*$ domain

$w_t$

$w_0$

Neural net sufficiently wide

$\Rightarrow K(w) = \nabla F(w) \nabla F(w)^\top$ well-conditioned in $B(w_0, R)$

$\Rightarrow L(w) = \frac{1}{2} \| F(w) - y \|^2$ satisfies PL$^*$ in $B(w_0, R)$

$\Rightarrow$ Linear convergence of (S) GD

# Theorem 4 ($PL^*$ ⇒ existence of solution + linear convergence)

**Theorem 4** (PL* $\Rightarrow$ existence of solution + linear convergence)

Consider the neural net and random initialization $W_0$ as described in Theorem 3. For large enough $m$, then with an appropriate step size $\eta$,

Theorem 4 (PL* ⇒ existence of solution + linear convergence)

Consider the neural net and random initialization $W_0$ as described in Theorem 3. For large enough $m$, then with an appropriate step size $\eta$,

① GD converges to a global minimizer in $B(W_0, R)$ with $R = O(\frac{1}{\mu})$

$w^*$ s.t. $F(w^*) = y$

**Theorem 4** (PL* ⇒ existence of solution + linear convergence)

Consider the neural net and random initialization $W_0$ as described in Theorem 3. For large enough $m$, then with an appropriate step size $\eta$,

① GD converges to a global minimizer in $B(W_0, R)$ with $R = O(\frac{1}{\mu})$

$$W^* \text{ s.t. } F(W^*) = y$$

② with exponential rate: $\boxed{L(W_t) \leq (1 - \eta \mu)^t L(W_0)}$

**Theorem 4** (PL* ⇒ existence of solution + linear convergence)

Consider the neural net and random initialization $W_0$ as described in Theorem 3. For large enough $m$, then with an appropriate step size $\eta$,

① GD converges to <u>a global minimizer</u> in $B(W_0, R)$ with $R = O(\frac{1}{\mu})$

$w^*$ s.t. $F(w^*) = y$

② with exponential rate: $\boxed{L(w_t) \leq (1 - \eta\mu)^t L(w_0)}$

Remarks:
- Note $R = O(\frac{1}{\mu})$

**Theorem 4** (PL* $\Rightarrow$ existence of solution + linear convergence)

Consider the neural net and random initialization $W_0$ as described in Theorem 3. For large enough $m$, then with an appropriate step size $\eta$,

① GD converges to a global minimizer in $B(W_0, R)$ with $R = O(\frac{1}{\mu})$

$\quad\quad\quad\quad\quad\quad$ $W^*$ s.t. $F(W^*) = y$

② with exponential rate: $\boxed{L(W_t) \leq (1 - \eta\mu)^t L(W_0)}$

Remarks:
- Note $R = O(\frac{1}{\mu})$
- Condition number $\kappa := \frac{1}{\eta\mu} = \sup_{W \in B(W_0, R)} \frac{\lambda_{max}(H_L(W))}{\mu}$

**Theorem 4** ( PL* $\Rightarrow$ existence of solution + linear convergence)

Consider the neural net and random initialization $W_0$ as described in Theorem 3. For large enough $m$, then with an appropriate step size $\eta$,

① GD converges to <u>a global minimizer</u> in $B(W_0, R)$ with $R = O(\frac{1}{\mu})$
   $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad w^* \text{ s.t. } F(w^*) = y$

② with exponential rate: $\boxed{L(W_t) \leq (1 - \eta \mu)^t \, L(W_0)}$

Remarks:
- Note $R = O(\frac{1}{\mu})$
- Condition number $\kappa := \frac{1}{\eta \mu} = \sup_{w \in B(W_0, R)} \dfrac{\lambda_{max}(H_L(w))}{\mu}$
- In the special case of a linear system $F(w) = Aw$ with
  $L(w) = \frac{1}{2} \| Aw - y \|^2$, then $\kappa = \dfrac{\lambda_{max}(AA^T)}{\lambda_{min}(AA^T)}$

Aside:

# Aside:

This theory covers:

- Wide NN with linear output layer $\left(\text{ie. } \sigma_{L+1}(z) = z\right)$
- CNN, resnet
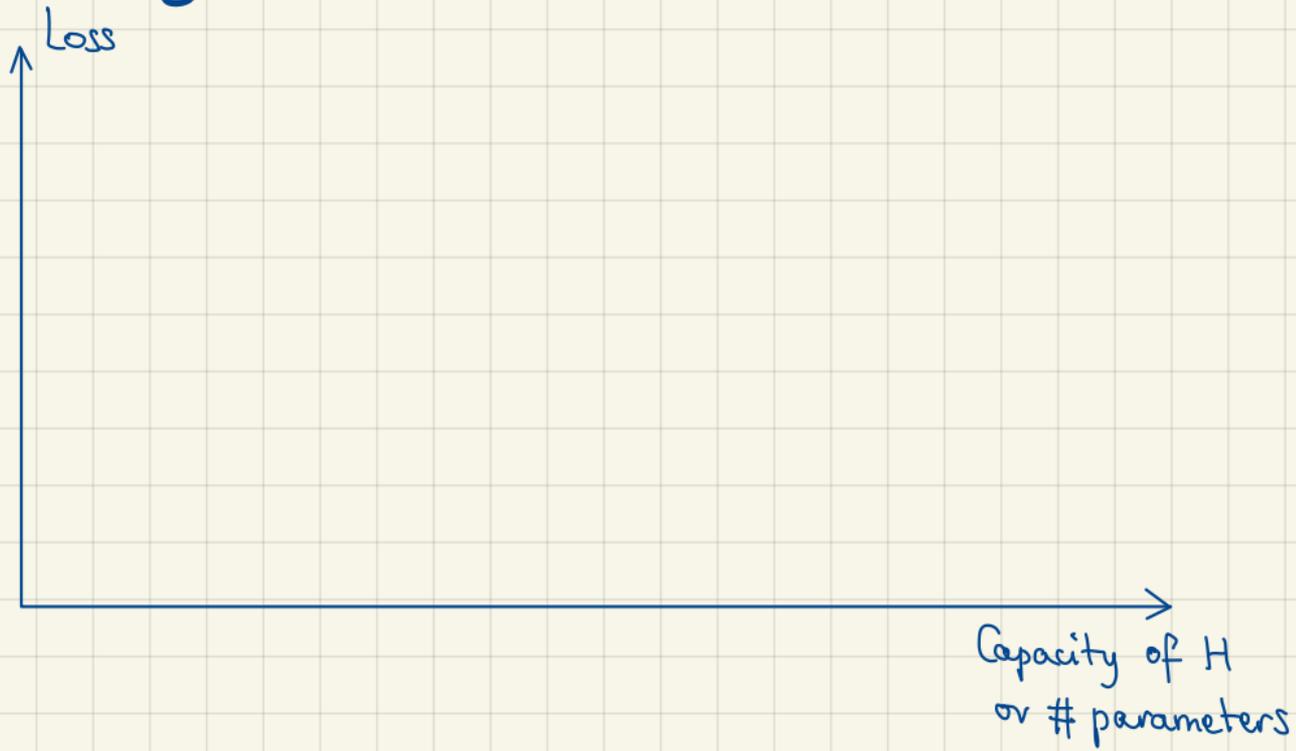
Doesnt cover

- NN with non-linear output layer
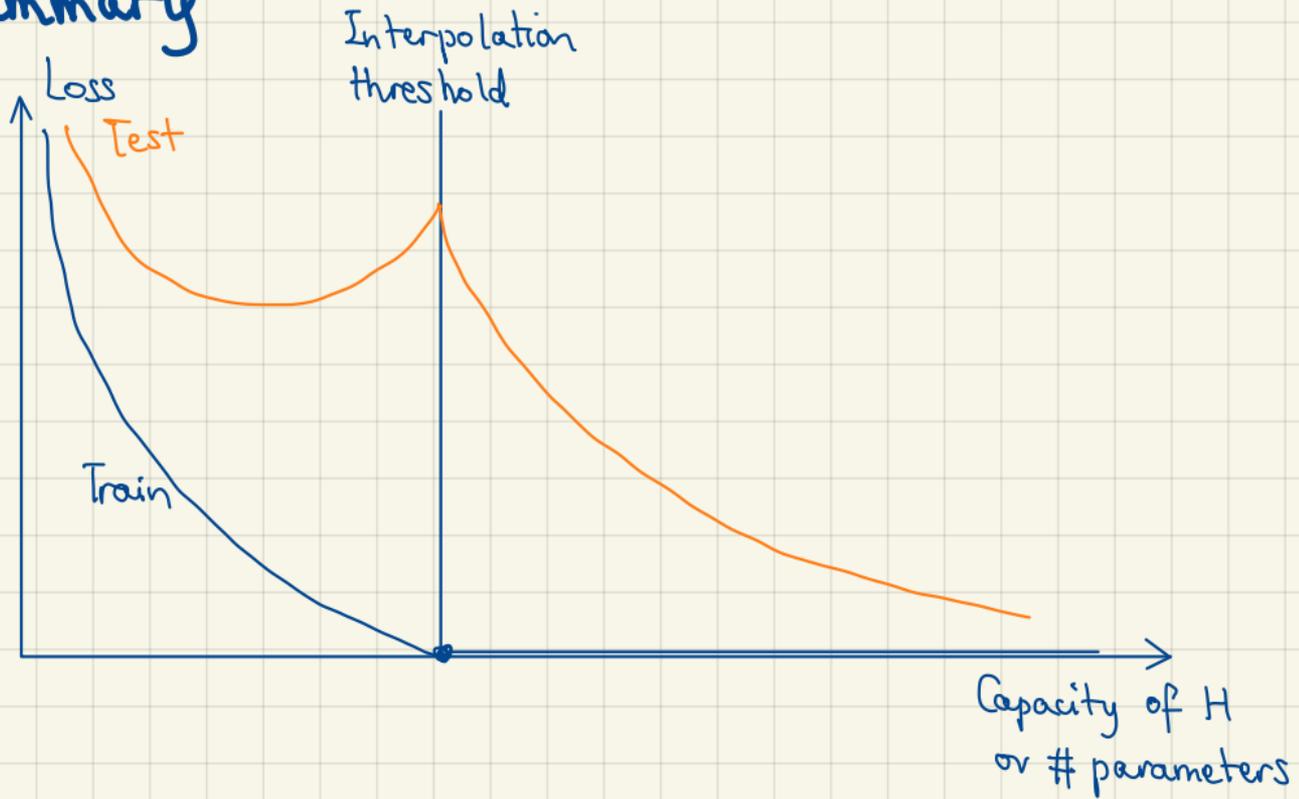- Bottleneck layers.
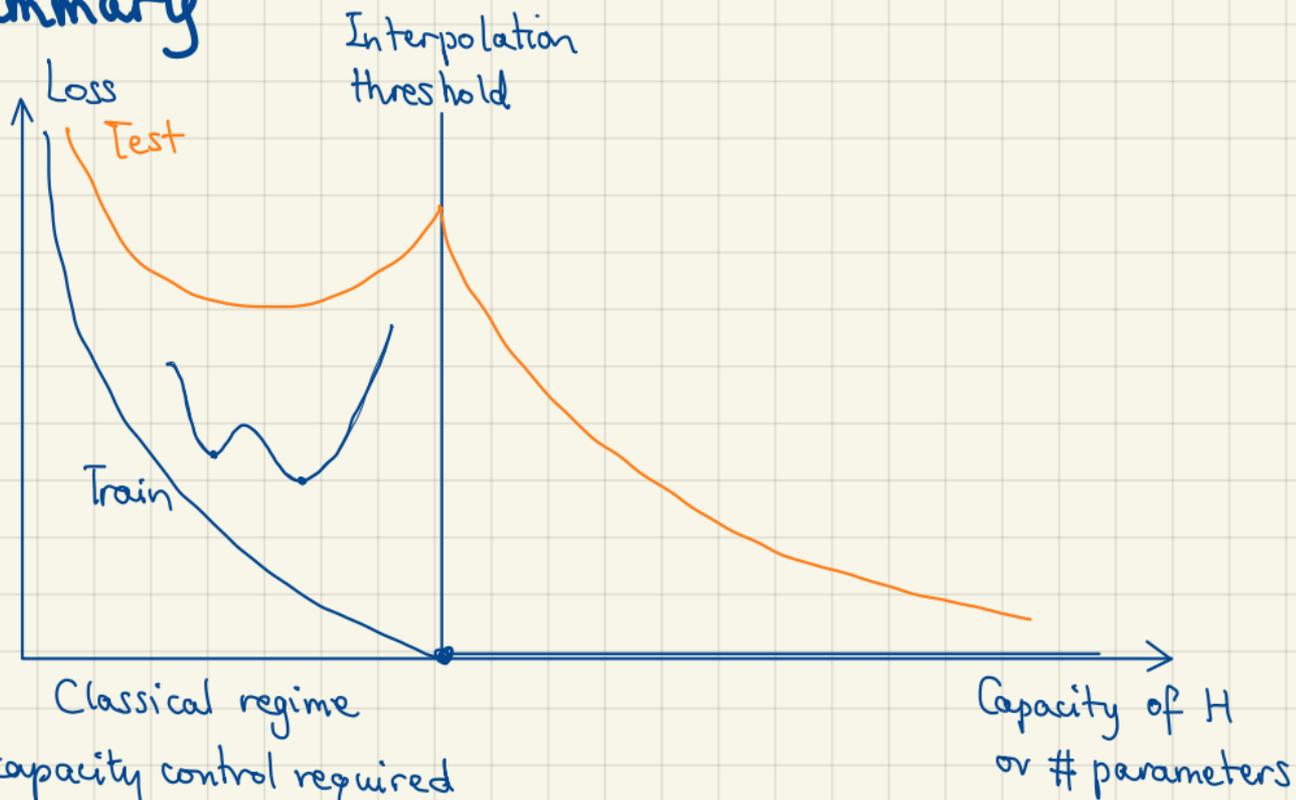
④ Summary and extensions

# Summary

# Summary

Loss (y-axis)

Capacity of H or # parameters (x-axis)

# Summary



Loss

Test

Interpolation threshold

Train

Capacity of H or # parameters

# Summary



- Classical regime
  - capacity control required
  - many non-global minima
  - local convexity, GD oscillates

# Summary



Loss

Test

Interpolation threshold

Train

Today's talk

Classical regime

Capacity of H or # parameters

- capacity control required
- many non-global minima
- local convexity, GD oscillates

# Summary



Loss

Interpolation threshold

Test

Train

- not locally convex
- satisfies PL*
  (i.e. all stationary points are global minima)
- fast convergence under (S)GD

Today's talk

Classical regime

Capacity of H or # parameters

- capacity control required
- many non-global minima
- local convexity, GD oscillates

# Limitations

# Limitations

- m unrealistically large. $m = \tilde{\Omega}\left(n R^{6L+2}\right)$

# Limitations

- m unrealistically large. $m = \tilde{\Omega}\left(n R^{6L+2}\right)$
- still in the NTK regime $\nabla F(w_0) \approx \nabla F(w_t)$

  this ignores the network's ability to learn features.

# Limitations

- $m$ unrealistically large. $m = \widetilde{\Omega}\left(n R^{6L+2}\right)$
- still in the NTK regime $\nabla F(w_0) \approx \nabla F(w_t)$

  this ignores the network's ability to learn features.

- beyond PL ?

# Limitations

- $m$ unrealistically large. $m = \tilde{\Omega}\left(n R^{6L+2}\right)$
- still in the NTK regime $\nabla F(w_0) \approx \nabla F(w_t)$

  this ignores the network's ability to learn features.

- beyond PL?      to incorporate (good) non-global minima?

# Limitations

- $m$ unrealistically large. $m = \tilde{\Omega}\left(n R^{6L+2}\right)$
- still in the NTK regime $\nabla F(w_0) \approx \nabla F(w_t)$

  this ignores the network's ability to learn features.

- beyond PL?    to incorporate (good) non-global minima?



- square loss $L(w) = \frac{1}{2} \| F(w) - y \|^2$. isn't the most commonly used
  loss function. What about other loss functions?

# Limitations

- m unrealistically large. $m = \tilde{\Omega}(n R^{6L+2})$
- still in the NTK regime $\nabla F(w_0) \approx \nabla F(w_t)$

  this ignores the network's ability to learn features.

- beyond PL?   to incorporate (good) non-global minima?



- square loss $L(w) = \frac{1}{2} \| F(w) - y \|^2$. isn't the most commonly used
  loss function. What about other loss functions?
- model $F(w)$ may not be smooth, e.g. ReLU

  [Oymak, Soltanolkotabi 20]

# Remedies.

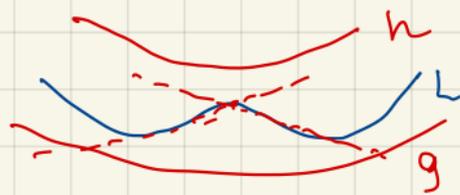# Remedies. [Frei, Gu 21]

# Remedies. [Frei, Gu 21]

- Proxy convexity

# Remedies. [Frei, Gu 21]

- Proxy convexity

  $L: \mathbb{R}^m \to \mathbb{R}$ is $(g, h)$-proxy-PL if $\exists\, g, h : \mathbb{R}^m \to \mathbb{R}$ s.t. $\forall\, \omega, v$,

  $$h(v) \geq g(\omega) + \langle \nabla L(\omega), v - \omega \rangle$$

# Remedies. [Frei, Gu 21]

- Proxy convexity

  $L: \mathbb{R}^m \to \mathbb{R}$ is $(g,h)$-proxy-PL if $\exists\, g, h: \mathbb{R}^m \to \mathbb{R}$ s.t. $\forall\, w, v,$

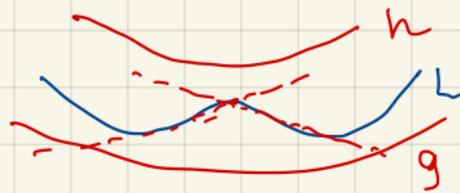  $$h(v) \geq g(w) + \langle \nabla L(w), v - w \rangle$$



- Proxy PL
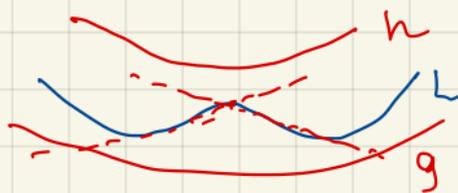
# Remedies. [Frei, Gu 21]

• Proxy convexity

$L: \mathbb{R}^m \to \mathbb{R}$ is $(g, h)$-proxy-PL if $\exists\, g, h: \mathbb{R}^m \to \mathbb{R}$ s.t. $\forall\, w, v,$

$$h(v) \geq g(w) + \langle \nabla L(w), v - w \rangle$$



• Proxy PL

$L: \mathbb{R}^m \to \mathbb{R}$ satisfies $(h, \xi)$-proxy PL inequality if

$\exists\, h: \mathbb{R}^m \to \mathbb{R}$ and $\xi \geq 0$, $\alpha, \mu > 0$ s.t. $\forall\, w$

$$\| \nabla L(w) \|^{\alpha} \geq \mu \, (h(w) - \xi)$$

# Remedies. [Frei, Gu 21]

- Proxy convexity

  $L: \mathbb{R}^m \to \mathbb{R}$ is $(g, h)$-proxy-PL if $\exists g, h : \mathbb{R}^m \to \mathbb{R}$ s.t. $\forall w, v,$

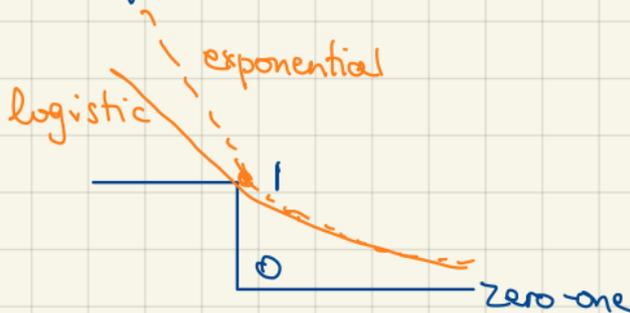  $$h(v) \geq g(w) + \langle \nabla L(w), v - w \rangle$$



- Proxy PL

  $L: \mathbb{R}^m \to \mathbb{R}$ satisfies $(h, \xi)$-proxy PL inequality if

  $\exists h: \mathbb{R}^m \to \mathbb{R}$ and $\xi \geq 0$, $\alpha, \mu > 0$ s.t. $\forall w$

  $$\| \nabla L(w) \|^{\alpha} \geq \mu \, (h(w) - \xi)$$
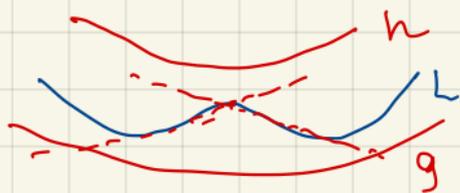
  — think of $h(w) \geq L(w)$

# Remedies. [Frei, Gu 21]

- Proxy convexity

  $L: \mathbb{R}^m \to \mathbb{R}$ is $(g,h)$-proxy-PL if $\exists \, g, h : \mathbb{R}^m \to \mathbb{R}$ s.t. $\forall \, w, v,$

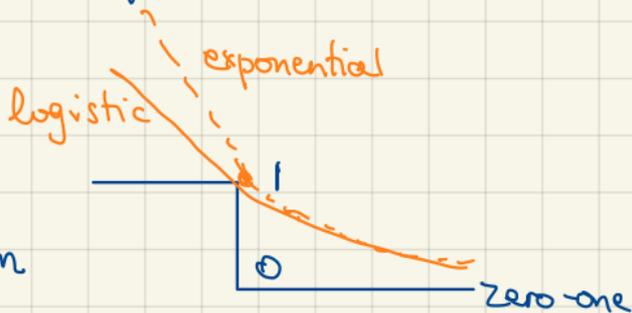  $$h(v) \geq g(w) + \langle \nabla L(w), v-w \rangle$$

- Proxy PL

  $L: \mathbb{R}^m \to \mathbb{R}$ satisfies $(h, \xi)$-proxy PL inequality if

  $\exists \, h: \mathbb{R}^m \to \mathbb{R}$ and $\xi \geq 0$, $\alpha, \mu > 0$ s.t. $\forall \, w$

  $$\| \nabla L(w) \|^\alpha \geq \mu \, (h(w) - \xi)$$

  - think of $h(w) \geq L(w)$
  - $\xi \geq 0$ allows good local minimum

Thanks.

# Non-linear vs linear output layer

last layer activation

$$g(w; x) = \phi(f(w; x))$$

$$(\phi(f))'' = \underbrace{\phi''(f)}_{} \underbrace{f'}_{O(1)} + \underbrace{\phi'(f)}_{O(1)} \underbrace{f''}_{O(\frac{1}{\sqrt{m}})}$$

$$\begin{cases} = 0 \text{ For linear output layer} \\ \neq 0 \text{ For non-linear output layer.} \end{cases}$$

# Classical statistical learning theory

Expected risk      $R(f) = \mathbb{E}_{P(x,y)}\, \ell\big(f(x), y\big)$

Bayes-optimal      $f^* = \underset{f: \mathbb{R}^d \to \mathbb{R}}{\text{argmin}}\ R(f)$

Empirical risk minimization

$$R_{emp}(f) = \frac{1}{n} \sum_{i=1}^{n} \ell\big(f(x_i), y_i\big)$$

$$f_{emp} = \underset{f \in H}{\text{argmin}}\ R_{emp}(f)$$

$\Rightarrow$ w.h.p. over the randomness in data, for any $f \in H$,

$$\underbrace{R(f)}_{\approx \text{Test}} - \underbrace{R_{emp}(f)}_{\text{Train}} < \tilde{O}\left( \sqrt{\frac{cap(H)}{n}} \right)\ \text{e.g. VC dim, covering number}$$