# ORF 543
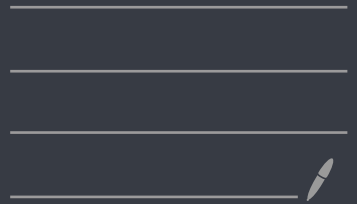
# Lecture 9/14 - NNs at Initialization

Consider a FCNN:

$$\text{input } \vec{x}_\alpha \in \mathbb{R}^{n_0} \to \vec{z}_\alpha^{(1)} = W^{(1)} x_\alpha + \vec{b}^{(1)} \in \mathbb{R}^{n_1}$$

$$\to \vec{z}_\alpha^{(2)} = W^{(2)} \sigma(\vec{z}_\alpha^{(1)}) + \vec{b}^{(2)} \in \mathbb{R}^{n_2}$$

$$\vdots$$

$$\to \vec{z}_\alpha^{(L+1)} \in \mathbb{R}^{n_{L+1}}$$

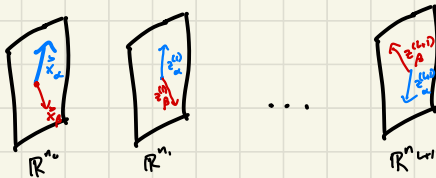So, $\quad z_{i\alpha}^{(l+1)} = b_i^{(l+1)} + \sum\limits_{j=1}^{n_l} W_{ij}^{(l+1)} \sigma(z_j^{(l)})$

and $\quad \vec{z}_\alpha^{(l)} = \langle z_{1\alpha}^{(l)}, \ldots, z_{n_l \alpha}^{(l)} \rangle$

We ask how to initialize $W_{ij}^{(l)}, b_i^{(l)}$ and learning rates for GD?

# Gaussian Initialization

Consider $\quad W_{ij}^{(l)} \sim N(0, V_w^{(l)}) \quad , \quad \vec{b}_i^{(l)} \sim N(0, V_b^{(l)})$

We use the **information propagation** framework where we want feature dot products to be conserved across layers.



$$\mathbb{R}^{n_0} \qquad \mathbb{R}^{n_1} \qquad \ldots \qquad \mathbb{R}^{n_{L+1}}$$

In math, we want to select $V_w^{(l)}(n_l, \sigma, L)$ and $V_b^{(l)}(n_l, \sigma, L)$ s.t.

$$\forall l \in \{0, \ldots, L\} \quad \frac{1}{n_l} \langle \vec{z}_\alpha^{(l)}, \vec{z}_\beta^{(l)} \rangle \approx \frac{1}{n_{l+1}} \langle \vec{z}_\alpha^{(l+1)}, \vec{z}_\beta^{(l+1)} \rangle \quad \text{preserved dot product averaged over dims}$$

There are two useful consequences of conservation of dot product

① We approximately preserve across $\ell \in \{0, \ldots, L-1\}$

$$\frac{1}{n_\ell} \| \vec{z}_\alpha^{(\ell)} \|^2, \qquad \left\langle \frac{\vec{z}_\alpha^{(\ell)}}{\| \vec{z}_\alpha^{(\ell)} \|}, \frac{\vec{z}_\beta^{(\ell)}}{\| \vec{z}_\alpha^{(\ell)} \|} \right\rangle$$

② The Law of Large Numbers suggests $\forall \ell$,

$$\frac{1}{n_\ell} \langle \vec{z}_\alpha^{(\ell)}, \vec{z}_\beta^{(\ell)} \rangle = \frac{1}{n_\ell} \sum_{j=1}^{n_\ell} z_{j i \alpha}^{(\ell)} z_{j i \beta}^{(\ell)} \approx \mathbb{E}\{ z_{i j \alpha}^{(\ell)} \cdot z_{i j \beta}^{(\ell)} \}$$

So, information propagation says we preserve the following across $\ell$:

$$\mathbb{E}\{ z_{i j \alpha}^{(\ell)} \}, \qquad \mathrm{Cov}\{ z_{i j \alpha}^{(\ell)}, z_{i j \beta}^{(\ell)} \}$$

Both conditions basically say mean and variance stay constant.

We can develop the following heuristic for $V_w^{(\ell)}, V_b^{(\ell)}$ with respect to $n_\ell$.



Suppose $|x_i| = O(1)$ ← low order, not unreasonably big or small

and $W_j \sim N(0, V_w), \quad b \sim N(0, V_b)$

Since $\vec{z} = \langle \vec{w}, \vec{x} \rangle + b$, we see $\vec{z} \sim N(0, V_b + \sum_{i=1}^{\hat{n}} V_w \| x_i \|^2)$

sum of gaussians

$$\Rightarrow \vec{z} \sim N(0, V_b + n V_w \cdot O(1))$$

We arrive at **fan-in scaling**, where

$$\boxed{\begin{array}{l} V_b^{(\ell)} = C_b, \quad C_b = O(1) \\[2mm] V_w^{(\ell)} = \dfrac{C_w}{n_{\ell-1}}, \quad C_w = O(1) \end{array}}$$

weight variance scales with width

**Def:** Let $T$ be a set. Then a <span style="color:red">Gaussian process</span> <span style="color:red">by $T$</span> is $\{X_t\}_{t \in T}$ such that $\langle X_{t_1}, ..., X_{t_k}\rangle \in \mathbb{R}^k$ is Gaussian $\forall \{t_1, ..., t_k\} \subseteq T$

ex/ Let $T = \{1, ..., n\}$. Let $X = \langle X_1, ..., X_n \rangle$ be jointly Gaussian.

ex/ Let $T = \mathbb{R}$. $X = X_t$ is gaussian process if $X$ is a random function on $\mathbb{R}$ with finite-dim distribution (fdd) $\langle X_{t_1}, ..., X_{t_n} \rangle \in \mathbb{R}^n$ Gaussian.

**Theorem:** <span style="color:blue">(Neal, Lee, ..., Hanin)</span>

Fix $n_0, n_{L+1}, \Theta$. Then as $n_1, ..., n_L \to \infty$    <span style="color:blue">input/output dims fixed, but else $\to \infty$</span>

$$\vec{z}^{(L+1)}_\alpha \to GP(0, k^{(L+1)})$$ <span style="color:blue">a Gaussian process</span>

i.e. $\mathbb{E}\{z_{i\alpha}^{(L+1)}\} = 0$ and $Cov(z_{i\alpha}^{(L+1)}, z_{j\beta}^{(L+1)}) = \delta_{ij} k_{\alpha\beta}^{(L+1)}$ <span style="color:red">Cov of two neurons when evaluated on $\vec{x}_\alpha, \vec{x}_\beta$</span>

This describes what happens when we send previous layers to infinite width. We can then recursively define

$$k_{\alpha\beta}^{(1)} = C_b^{(1)} + \frac{C_w^{(1)}}{n_0} \langle \vec{x}_\alpha, \vec{x}_\beta \rangle$$

$$K_{\alpha\beta}^{(l+1)} = C_b^{(l)} + C_w^{(l)} \underset{k_{\alpha\beta}^{(l)}}{\mathbb{E}} \{\Theta(\vec{z}_\alpha^{(l)}) \Theta(\vec{z}_\beta^{(l)})\}$$

<span style="color:blue">$$k_{\alpha\beta}^{(l)} = \lim_{n_0, ..., n_{l-1} \to \infty} Cov(z_{i\alpha}^{(l)}, z_{i\beta}^{(l)})$$

$$\approx \frac{1}{n_l}\langle \vec{z}_\alpha^{(l)}, \vec{z}_\beta^{(l)}\rangle$$</span>

<span style="color:blue">$\begin{pmatrix}\vec{z}_\alpha^{(l)} \\ \vec{z}_\beta^{(l)}\end{pmatrix} \sim N\left(0, \begin{pmatrix} k_{\alpha\alpha}^{(l)} & k_{\alpha\beta}^{(l)} \\ k_{\beta\alpha}^{(l)} & k_{\beta\beta}^{(l)} \end{pmatrix}\right)$</span>   call this $K^{(l)}$

<span style="color:blue">$$\int_{\mathbb{R}^2} \frac{\Theta(z_\alpha^{(l)})\Theta(z_\beta^{(l)}) e^{-\frac{1}{2}\langle (k^{(l)})^{-1}\begin{bmatrix} z_\alpha^{(l)} \\ z_\beta^{(l)} \end{bmatrix}, \begin{bmatrix} z_\alpha^{(l)} \\ z_\beta^{(l)} \end{bmatrix}\rangle}}{\det(2\pi k^{(l)})^{\frac{1}{2}}} dz_\alpha^{(l)} dz_\beta^{(l)}$$</span>

Info prop $\iff C_b^{(l)}, C_w^{(l)}$ are set such that $k_{\alpha\beta}^{(l)}$ is well-behaved at large $l$.

# Lecture 9/19 Tuning to Criticality

Note: at a particular layer $(\ell+1)$, we are given $\bar{z}_\alpha^{(\ell)}$ which is a random variable every neuron in the layer shares

$\Rightarrow z_{i\alpha}^{(\ell+1)}$ are i.i.d. Gaussian with variance

$$C_b + \frac{C_w}{n_\ell} \| \sigma(\bar{z}_\alpha^{(\ell)}) \|^2$$

$\hookleftarrow$ Gaussian with a random variance

Recall that the goal of info. prop. is to    converge

$$\frac{1}{n_\ell} \langle z_\alpha^{(\ell)}, z_\beta^{(\ell)} \rangle \approx \frac{1}{n_{\ell+1}} \langle z_\alpha^{(\ell+1)}, z_\beta^{(\ell+1)} \rangle$$

$\uparrow$ $K_{\alpha\beta}^{(\ell)}$            $\uparrow$ $K_{\alpha\beta}^{(\ell+1)}$

So, in the infinite limit, the goal is to find $C_b, C_w$ s.t. $K_{\alpha\beta}^{(\ell)}$ is as constant as possible across $\ell$.

## Ex/ $\sigma(t) = t$ (Deep Linear Networks)

$\alpha = \beta$: $K_{\alpha\alpha}^{(\ell+1)} = C_b + C_w \mathbb{E}_{K^{(\ell)}} \{ \sigma(z_\alpha)^2 \}$

This is like supposing the previous layer already went to infinite width $\Rightarrow z_{i\alpha}^{(\ell)} \sim N(0, K_{\alpha\alpha}^{(\ell)})$

$$= C_b + C_w \int_{-\infty}^{\infty} z_\alpha^2 e^{-\frac{z_\alpha^2}{2 K_{\alpha\alpha}^{(\ell)}}} \frac{dz_\alpha}{\sqrt{2\pi K_{\alpha\alpha}^{(\ell)}}} = C_b + C_w K_{\alpha\alpha}^{(\ell)}$$

$\alpha \neq \beta$: $K_{\alpha\beta}^{(\ell+1)} = C_b + C_w \mathbb{E}_{K^{(\ell)}} \{ \langle z_\alpha, z_\beta \rangle \} = C_b + C_w K_{\alpha\beta}^{(\ell)}$

So, if $\sigma(t) = t$ we want to choose $C_b = 0, C_w = 1$.

<u>Remark</u>: If $C_b = 0$ but $C_w \neq 1$, we have an initialization

$$K_{\alpha\beta}^{(\ell+1)} = (C_w)^\ell K_{\alpha\beta}^{(0)} \hookleftarrow \text{vanishes or explodes if } C_w \neq 1$$

Ex/ $\sigma(t) = \text{ReLU}(t) = \max\{0, t\}$

We have $k_{\alpha\alpha}^{(l+1)} = C_b + C_w \int_0^\infty \frac{z_\alpha^2 \, e^{-\frac{z_\alpha^2}{2k_{\alpha\alpha}^{(l)}}}}{\sqrt{2\pi k_{\alpha\alpha}^{(l)}}} \, dz_\alpha = C_b + \frac{C_w K_{\alpha\alpha}^{(l)}}{2}$

With $C_b = 0$, we require $1 = \frac{C_w^{(0)} \cdots C_w^{(l)}}{2^l} \; \forall_l \Rightarrow C_w = 2$ <span style="color:blue">"He-initialization"</span>

However, when $\sigma(t) \neq t$, $\mathbb{E}_{k^{(l)}}\{\sigma(z_\alpha)\sigma(z_\beta)\}$ is hard.

We can claim that the recursion
$$(\#) \quad K_{\alpha\beta}^{(l+1)} = C_b + C_w \mathbb{E}_{k^{(l)}}\{\sigma(z_\alpha), \sigma(z_\beta)\}$$
is a 3d dynamical system with variables $(k_{\alpha\alpha}^{(l)}, k_{\beta\beta}^{(l)}, K_{\alpha\beta}^{(l)})$ with time parameter $l$.

To solve such a system, we find fixed points, linearize about the fixed points, and ensure the points are stable & critical.

Fixed points at $(\circledast) \quad K_* = C_b + C_w \mathbb{E}_{k_*}\{\sigma^2(z)\}$
$$\left( k_{\alpha\alpha}^{(l)} = K_* \Rightarrow K_{\alpha\alpha}^{(l+1)} = k_* \right)$$

This condition will have that at deep layers, if $\overset{\circ}{x}_\alpha \sim \mathcal{N}(0, k_*)$, then at large $l$, $\frac{1}{n_l}\|z_\alpha^{(l)}\|^2 \approx \frac{1}{n_0}\|x_\alpha\|^2 = k_*$

The second condition is <span style="color:blue">parallel perturbation $d$ of $x_\alpha$ in direction $\frac{2}{x_\alpha}$</span>
$$(II) \quad \left.\frac{\partial k_{\alpha\alpha}^{(l+1)}}{\partial k_{\alpha\alpha}^{(l)}}\right|_{k_{\alpha\alpha}^{(l)} = K_*} = 1 \quad \color{blue}{\left( k_{\alpha\alpha}^{(l)} = k_* + \delta k \Rightarrow k_{\alpha\alpha}^{(l+1)} = k_* + \delta k + O(\delta k) \right)}$$
<span style="color:blue">linearized</span>

Thirdly,
$$(I) \quad \left.\frac{\partial K_{\alpha\beta}^{(l+1)}}{\partial k_{\alpha\beta}^{(l)}}\right|_{k_{\alpha\alpha}^l = k_{\beta\beta}^l = k_{\alpha\beta}^l = K_*} \quad \color{blue}{\left( k_{\alpha\beta}^{(l)} = k_* + \delta k \Rightarrow K_{\alpha\beta}^{(l+1)} = k_* + \delta k + O(\delta k) \right)}$$
<span style="color:blue">linearized</span>

These are the dynamical systems constraints for a fixed, stable, critical fixed points. Note that we treat this as small perturbation from a point to generate $x_\alpha, x_\beta$, which is why we use linear approx.

Now, 
$$\frac{\partial k_{\alpha\alpha}^{(l+1)}}{\partial k_{\alpha\alpha}^{(l)}} = \frac{\partial}{\partial k_{\alpha\alpha}^{(l)}}\left(C_b + C_w \,\mathbb{E}_{k^{(l)}}\left\{\theta(z_\alpha)^2\right\}\right)$$

$$= C_w \frac{\partial}{\partial k_{\alpha\alpha}^{(l)}} \int \theta(z_\alpha)^2 \frac{e^{-\frac{z_\alpha^2}{2k_{\alpha\alpha}^{(l)}}}}{\sqrt{2\pi k_{\alpha\alpha}^{(l)}}} dz_\alpha \qquad \left(\begin{array}{c}\text{Gaussians are} \\ \text{cuter in Fourier} \\ \text{Space}\end{array}\right)$$

$$\underset{\substack{\text{Fourier} \\ \text{Transform}}}{=\!=} C_w \frac{\partial}{\partial k_{\alpha\alpha}^{(l)}} \int \hat{\theta}^2(\gamma)\, e^{-k_{\alpha\alpha}^{(l)}\frac{\gamma^2}{2}} d\gamma$$

$$= C_w \int \hat{\theta}^2(\gamma)\left(-\tfrac{1}{2}\gamma^2\right) e^{-k_{\alpha\alpha}^{(l)}\frac{\gamma^2}{2}} d\gamma$$

F.T. diagonalizes the derivative

$$= C_w \int \tfrac{1}{2}\partial_{z_\alpha}^2 (\theta(z_\alpha)) \frac{e^{-\frac{z_\alpha^2}{2k_{\alpha\alpha}^{(l)}}}}{\sqrt{2\pi k_{\alpha\alpha}^{(l)}}} dz_\alpha$$

$$\chi_{\parallel}(k_*) = \frac{C_w}{2} \mathbb{E}_{k_*}\left\{\partial^2(\theta^2(z))\right\} = 1$$

We can do the same thing to find
$$\chi_{\perp}(k_*) = C_w \mathbb{E}_{k_*}\left\{(\partial\theta(z))^2\right\} = 1$$

So, the constraints of "tuning to criticality" result with

$$\boxed{\begin{array}{l}(*)\ \ k_* = C_b + C_w\, \mathbb{E}_{k_*}\left\{\theta^2(z)\right\} \\[2mm] (\parallel)\ \ \chi_{\parallel}(k_*) = \frac{C_w}{2} \mathbb{E}_{k_*}\left\{\partial^2(\theta^2(z))\right\} = 1 \\[2mm] (\perp)\ \ \chi_{\perp}(k_*) = C_w\, \mathbb{E}_{k_*}\left\{(\partial\theta(z))^2\right\} = 1\end{array}}$$

These conditions confirm that if you have two inputs $x_\alpha, x_\beta$ "close" with $Cov(x_\alpha, x_\beta) = 1-\epsilon$, things don't exponentially explode or vanish ($k_*$ is fixed point).

We can return to $\theta(t) = \text{ReLU}(t)$
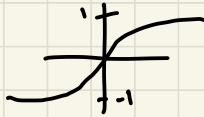
$(*)$ $K_* = C_b + C_w \frac{K_*}{2}$

$\Longrightarrow$

$(\parallel)$ $1 = \frac{C_w}{2} \mathbb{E}_{k_*}\left\{ \frac{\partial^2}{\partial z^2}\left(z^2 \mathbb{1}_{z>0}\right)\right\} = \frac{C_w}{2} \mathbb{E}_{k_*}\left\{2\, \mathbb{1}_{z>0}\right\} = \frac{C_w}{2}$

Gaussian integral

$(\perp)$ $1 = C_w \mathbb{E}_{k_*}\left\{ (\partial z \, \mathbb{1}_{z>0})^2\right\} = C_w \mathbb{E}_{k_*}\left\{\mathbb{1}_{z>0}\right\} = \frac{C_w}{2}$

So, $C_w = 2$, $C_b = 0$, $K_* \geq 0$ arbitrary

EX/ $\theta(t) = \tanh(t)$



Note: the only fixed point is $K_* = 0$.

$\chi_\parallel(K_*) = \frac{C_w}{2} \mathbb{E}_{k_*}\left\{\partial^2(\theta^2(z))\right\} = C_w \mathbb{E}\left\{\partial(\theta(z)\theta'(z))\right\}$

$= C_w \mathbb{E}_{k_*}\left\{\theta(z)\theta''(z)\right\} + \chi_\perp(K_*)$

So, if you want $\chi_\parallel(K_*) = \chi_\perp(K_*) = 1$, we require

$C_w \mathbb{E}_{k_*}\left\{\theta(z)\theta''(z)\right\} = 0 \quad \Longleftrightarrow \quad K_* = 0$

$\theta\theta''$ is even and 0 at origin

So what happens is that, at criticality

$C_b = 0, \quad C_w = 1, \quad K_{aa}^{(\ell)} = \frac{(C_w)^\ell}{2\ell} = \frac{1}{2\ell}$ at large $\ell$.

Covariances approach fixed point, don't do exponential stuff.

# Lecture 9/21 - NNGP

Theorem: For $L \geq 1$, $n_0, n_L \geq 1$, $\theta: \mathbb{R} \to \mathbb{R}$. Define

$$z_{i\alpha}^{(l+1)} = \begin{cases} b_i^{(l+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(l+1)} \theta(z_{j\alpha}^{(l)}) & \ell \geq 1 \\ b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_{j\alpha} & \ell = 0 \end{cases}$$

with $W_{ij}^{(l+1)} \sim N(0, \frac{C_W}{n_\ell})$, $b_i^{(l+1)} \sim N(0, C_b)$

If $\theta$ is poly bounded $\left(\text{i.e. } \exists n \geq 1, C > 0 \text{ s.t. } \sup_{x \in \mathbb{R}} \frac{|\theta(x)|}{1 + x^{2n}} \leq C\right)$

then for any $\vec{X}_A = (x_{\alpha_1}, \dots, x_{\alpha_k})$, $x_{\alpha_j} \in \mathbb{R}^{n_0}$, the output vector $\vec{z}_A^{(l+1)} = \langle z_{\alpha_1}^{(l+1)}, \dots, z_{\alpha_k}^{(l+1)} \rangle \in \mathbb{R}^{k \times n_{l+1}}$ converge in distribution as $n_1, \dots, n_L \to \infty$ to a mean 0 Gaussian with

$$\lim_{n_1, \dots, n_L \to \infty} \text{Cov}(z_{i\alpha}^{(l+1)}, z_{j\beta}^{(l+1)}) = \delta_{ij} k_{\alpha\beta}^{(l+1)}$$

where

$$\begin{cases} k_{\alpha\beta}^{(l+1)} = C_b + C_W \, \mathbb{E}_{k^{(l)}} \{ \theta(z_\alpha) \theta(z_\beta) \} & \ell \geq 1 \\ k_{\alpha\beta}^{(l+1)} = C_b + \frac{C_W}{n_0} \vec{x}_\alpha \cdot \vec{x}_\beta & \ell = 0 \end{cases}$$

Recall:

(1) Suppose $\vec{X}_n \in \mathbb{R}^k$ is a random variable with
$$\mathbb{E}\{ e^{-i \vec{X}_n \cdot \vec{\jmath}} \} \xrightarrow{n \to \infty} \mathbb{E}\{ e^{-i \vec{X} \cdot \vec{\jmath}} \} \quad \forall \vec{\jmath} \in \mathbb{R}^k.$$

Then, $X_n \xrightarrow{d} X$ ← in distribution

(2) Suppose $X \sim N(\vec{\mu}, \Sigma) \in \mathbb{R}^k$. Then,
$$\mathbb{E}\{ e^{-i \vec{X}_n \cdot \vec{\jmath}} \} = e^{-i \vec{\mu} \cdot \vec{\jmath} - \frac{1}{2} \vec{\jmath}^T \Sigma \vec{\jmath}}$$

Proof: We WTS that for any $\vec{\jmath} = \langle \vec{\jmath}_1, \dots, \vec{\jmath}_{n_{l+1}} \rangle$, $\vec{\jmath}_j \in \mathbb{R}^k$,
$$\lim_{n_1, \dots, n_L \to \infty} \mathbb{E}\{ e^{-i \vec{z}_A^{(l+1)} \cdot \vec{\jmath}} \} = e^{-\frac{1}{2} \sum_{j=1}^{n_{l+1}} \vec{\jmath}_j^T k_A^{(l+1)} \vec{\jmath}_j} \quad (*)$$

where
$$k_A^{(l+1)} = \begin{pmatrix} k_{\alpha_1 \alpha_1}^{(l+1)} & \cdots & \\ \vdots & \ddots & \\ & & k_{\alpha_k \alpha_k}^{(l+1)} \end{pmatrix}$$

<u>**Step 1**</u>: Vibes: we can think of the layers moving through the network as a Markov chain.

Given $z_A^{(l)}$, we find $\quad$ (men 0 independent Gaussians)

$$z_{jA}^{(l+1)} \equiv \langle z_{j\alpha_1}^{(l+1)}, \ldots, z_{j\alpha_k}^{(l+1)} \rangle \quad \text{and}$$

$$\text{Cov}\left( z_{j\alpha}^{(l+1)}, z_{j\beta}^{(l+1)} \mid z_A^{(l)} \right)$$

<u>Recall</u>: If $\vec{X} \sim N(0, \Sigma) \in \mathbb{R}^K$ and $\vec{u}, \vec{v} \in \mathbb{R}^K$, $\langle \vec{X} \cdot \vec{u}, \vec{X} \cdot \vec{v} \rangle$ is Gaussian with mean 0 and $\text{Cov}(\vec{X} \cdot \vec{u}, \vec{X} \cdot \vec{v}) = \vec{u}^T \Sigma \vec{v}$

Note that $z_{i\alpha}^{(l+1)} = \langle b_i^{(l+1)}, W_{i_1}^{(l+1)}, \ldots, W_{in_l}^{(l+1)} \rangle \cdot \langle 1, \theta(z_{i\alpha}^{(l)}), \ldots, \theta(z_{A_l \alpha}^{(l)}) \rangle$

$$\Rightarrow \text{Cov}\left( z_{j\alpha}^{(l+1)}, z_{j\beta}^{(l+1)} \right) = \begin{bmatrix} 1 \\ \theta(\vec{z}_\alpha^{(l)}) \end{bmatrix}^T \begin{bmatrix} C_b & \frac{C_w}{n_l} & 0 \\ 0 & \ddots & \frac{C_w}{n_l} \end{bmatrix} \begin{bmatrix} 1 \\ \theta(\vec{z}_\beta^{(l)}) \end{bmatrix}$$

$$= C_b + \frac{C_w}{n_l} \sum_{j=1}^{n_l} \theta(z_{j\alpha}^{(l)}) \theta(z_{j\beta}^{(l)}) = \hat{K}_{\alpha\beta}^{(l+1)}$$

Thus, $\mathbb{E}\left\{ e^{-i \vec{z}_A^{(l+1)} \cdot \vec{y}} \right\} = \mathbb{E}\left\{ \mathbb{E}\left\{ e^{-i \vec{z}_A^{(l+1)} \cdot \vec{y}} \mid \vec{z}_A^{(l)} \right\} \right\}$

$$= \mathbb{E}\left\{ e^{-\frac{1}{2} \sum_{j=1}^{n_{l+1}} \vec{y}_j^T \hat{K}_A^{(l+1)} \vec{y}_j} \right\} \quad (\#)$$

<span style="color:gray">↑ we want this to approach constant $\Sigma$</span>

$\vec{z}_{iA}^{(l+1)}$ are i.i.d. Gaussian with mean 0 but with some covariance $\hat{K}_A^{(l)}$.

<u>**Step 2**</u>: Vibes: Each transition between layers is symmetric to permutation of the neurons. So, only averages can matter.

Each entry of $\hat{K}_A^{(l+1)}$ has form $O_\rho^{(l)} = \frac{1}{n_l} \sum_{j=1}^{n_l} f(z_{jA}^{(l)})$

$$= \frac{1}{n_l} \sum_{j=1}^{n_l} \left( b + C_w \theta(z_{j\alpha}^{(l)}) \theta(z_{j\beta}^{(l)}) \right)$$

We can use the following Proposition:

**Prop**: If $f$ is poly bounded, $\sup\limits_{n_1, \ldots, n_l \geq 1} \left| \mathbb{E}\left\{ O_\rho^{(l)} \right\} \right| < \infty$ $\quad$ (always bounded)

and $\lim\limits_{n_1, \ldots, n_l \to \infty} \text{Var}\left( O_\rho^{(l)} \right) = 0$ $\quad$ (goes to constant)

**Corollary:** If we define $K_{\alpha\beta}^{(L+1)} = \lim_{n_1,\ldots,n_L \to \infty} \mathbb{E}\{\hat{K}_{\alpha\beta}^{(L+1)}\}$, then (#) $\Rightarrow$ (*).

**Proof of corollary:** The proposition gives $\hat{K}_{\alpha\beta}^{(L+1)} \xrightarrow{d} K_{\alpha\beta}^{(L+1)}$.

continuous

Also, the map
$$K \mapsto e^{-\frac{1}{2}\frac{3}{2}^T K \frac{3}{2}} \quad \text{is bounded \& } \overset{\vee}{C^0}.$$

So, all the network outputs' variances converge to the same shared deterministic covariance $K_{\alpha\beta}^{(L+1)}$. $\square$

We now know that the output vectors converge in distributions to mean 0 Gaussians with $\lim_{n_1,\ldots,n_L \to \infty} \mathrm{Cov}\left(z_{i\alpha}^{(L+1)}, z_{j\beta}^{(L+1)}\right) = \delta_{ij} K_{\alpha\beta}^{(L+1)}$

We complete the proof by deriving a recurrence relation for $k_{\alpha\beta}^{(L+1)}$. We know

$$K_{\alpha\beta}^{(L+1)} = \lim_{n_1,\ldots,n_L \to \infty} \mathrm{Cov}\left(z_{i\alpha}^{(L+1)}, z_{i\beta}^{(L+1)}\right)$$

$\mathrm{Cov}(X,Y) = \mathbb{E}\{\mathrm{Cov}(X,Y|Z)\}$
$+ \mathrm{Cov}(\mathbb{E}\{X|Z\}, \mathbb{E}\{Y|Z\})$ = Law of Total Covariance

$$= \lim_{n_1,\ldots,n_L \to \infty} \mathbb{E}\left\{\mathrm{Cov}\left(z_{i\alpha}^{(L+1)}, z_{i\beta}^{(L+1)} \mid z_A^{(L)}\right)\right\}$$
$$+ \mathrm{Cov}\left(\mathbb{E}\{z_{i\alpha}^{(L+1)} \mid z_A^{(L)}\}, \mathbb{E}\{z_{i\beta}^{(L+1)} \mid z_A^{(L)}\}\right)$$

mean 0 Gaussian

$$= \lim_{n_1,\ldots,n_L \to \infty} \mathbb{E}\left\{C_b + \frac{C_w}{n_L} \sum_{j=1}^{n_L} \mathcal{O}\left(z_{j\alpha}^{(L)}\right)\mathcal{O}\left(z_{j\beta}^{(L)}\right)\right\}$$

all these have same expectation because of symmetry

$$= C_b + C_w \underset{h^{(L)}}{\mathbb{E}}\left\{\mathcal{O}(z_\alpha)\mathcal{O}(z_\beta)\right\}$$

limiting outputs. So limit of expectation is expectation of limits via Continuous Mapping Theorem

We can repeat this logic via induction to get the recurrence relation.

We finish by going back and proving the proposition:

**Prop:** If $f$ is poly bounded, $\sup_{n_1,\ldots,n_L \geq 1} \left|\mathbb{E}\{\mathcal{O}_\beta^{(L)}\}\right| < \infty$ (always bounded)

and $\lim_{n_1,\ldots,n_L \to \infty} \mathrm{Var}(\mathcal{O}_\beta^{(L)}) = 0$ (goes to content)

**Proof:** We induct on $L$. When $L=1$,
$$z_{iA}^{(1)} = \langle z_{i\alpha_1}^{(1)}, \ldots, z_{i\alpha_K}^{(1)}\rangle \text{ are i.i.d Gaussian with}$$
mean 0 and $\mathrm{Cov}(z_{i\alpha}^{(1)}, z_{i\beta}^{(1)}) = C_b + \frac{C_w}{n_0} \vec{x}_\alpha \cdot \vec{x}_\beta$

Thus, $\mathbb{E}\{O_f^{(1)}\} = \mathbb{E}\{f(z_{jA}^{(1)})\}$ is finite because $f$ is poly
bounded independently of $n_1$. Furthermore,

$$\text{Var}\left(O_f^{(1)}\right) = \text{Var}\left(\frac{1}{n_1} \sum_{j=1}^{n} f(z_{jA}^{(1)})\right) = \frac{1}{n_1}\text{Var}\left(f(z_{jA}^{(1)})\right)$$

$$\leq \frac{1}{n_1}\mathbb{E}\{f(z_{jA}^{(1)})^2\}$$

$$\to 0 \text{ as } L \to \infty.$$

The inductive step happens because $f$ is poly bounded.

□

# Lecture ? - LR in NTK/GP Regime

Last time: We saw

① How to set $C_B, C_W$ in a random FCNN at large width of the form $z_{i;\alpha}^{(l+1)} = b_i^{(l+1)} + \sum_{j=1}^{n_l} W_{ij}^{(l+1)} \sigma\left(z_{j;\alpha}^{(l)}\right)$

with $W_{ij}^{(l+1)} \sim N\left(0, \frac{C_W}{n_l}\right)$ and $b_i^{(l+1)} \sim N(0, C_B)$

② That as $n_1, ..., n_L \to \infty$ $\quad z_\alpha^{(l+1)} \to GP(0, K^{(l+1)})$
with
$$\lim_{n_1, ..., n_L \to \infty} Cov\left(z_{i;\alpha}^{(l+1)}, z_{j;\beta}^{(l+1)}\right) = \delta_{ij} K_{\alpha\beta}^{(l+1)} \quad \text{and the relation}$$

$$K_{\alpha\beta}^{(l+1)} = C_b + C_W \, \mathbb{E}_{K^{(l)}}\left[\sigma(z_\alpha) \sigma(z_\beta)\right]$$

with $\chi_{||} = \frac{C_W}{2} \mathbb{E}_{K_*}\left[\sigma^2 \sigma^2(z)\right] = 1$, $\quad \chi_\perp = C_W \mathbb{E}_{K_*}\left[\sigma'(z)^2\right] = 1$

Today: We ask how to set LR for GD to be "well-behaved"!
$$\theta(t+1) = \theta(t) - \eta_t \vec\nabla_\theta \mathcal{L}(\theta(t))$$

## Intuition 1: Denote that $z(\vec{x}, \theta) = \theta \vec{x}$, $Y = \theta_* X$, $\mathcal{L}(\theta) = \frac{1}{2} \|\theta X - Y\|_2^2$
This yields
$$\vec\nabla_\theta \mathcal{L}(\theta) = (\theta X - Y) X^T = (\theta - \theta_*) X X^T$$

So, the GD update step becomes
$$\theta(t+1) - \theta_* = \theta(t) - \theta_* - \eta(\theta(t) - \theta_*) X X^T$$
$$= (\theta(t) - \theta_*)(I - \eta X X^T)$$

$$\Rightarrow \boxed{\eta < \frac{2}{\lambda_{max}(XX^T)} = \frac{2}{\|Hess(\mathcal{L})\|_{op}} = \frac{2}{\lambda_{max}(\vec\nabla_\theta z (\vec\nabla_\theta z)^T)}}$$

Under this condition,
$$\|\theta(t+1) - \theta_*\|_2 \leq \|\theta(t) - \theta_*\|_2 (1 - \eta \lambda_{min}(XX^T))$$
$$\leq \|\theta(t) - \theta_*\|_2 \, e^{-\eta t \lambda_{min}(XX^T)}$$

So, the best convergence rate is $e^{\frac{-2t}{\kappa(XX^T)}}$, $\quad \kappa(A) = \frac{\lambda_{max}(A)}{\lambda_{min}(A)}$

**Intuition 2:** Suppose we have noisy gradients

$$\theta(t+1) = \theta(t) - \eta_t \left( \vec{\nabla}_\theta \mathcal{L}(\theta(t)) + \xi_t \right), \quad \xi_t \sim \mathcal{N}(0, \sigma^2)$$

$$\Rightarrow \theta(t+1) - \theta_* = (\theta(t) - \theta_*)(I - \eta_t x x^T) + \eta_t \xi_t$$

$$\Rightarrow \|\theta(t+1) - \theta_*\|_2 \leq \|\theta(t) - \theta_*\|_2 \, e^{-\eta_t \lambda_{min}} + \eta_t \|\xi_t\|$$

$$\leq \|\theta(0) - \theta_*\|_2 \, e^{-\sum_{s=0}^{t} \eta_t \lambda_{min}} + \sum_{s=0}^{t} \eta_s \|\xi_s\| \, e^{-\sum_{s'=s+1}^{t} \eta_{s'} \lambda_{min}}$$

So, we need $\boxed{\displaystyle\sum_{s=0}^{\infty} \eta_s = \infty \quad \text{and} \quad \eta_s \to 0}$ $\left( \text{also } \displaystyle\sum_{s=0}^{\infty} \eta_s^2 < \infty \right)$

Now, returning to wide NNs with scalar output ($n_{L+1} = 1$), the **effective Jacobian** is

$$\vec{\eta} \odot \overset{\text{elementwise}}{\vec{\nabla}_\theta z_{1\alpha}^{(L+1)}} = \left( \eta_{\theta_j} \partial_{\theta_j} z_{1\alpha}^{(L+1)}, \quad j \in \{1, \ldots, \# \text{ params}\} \right)$$

($\eta_{\theta_j}$ for each param)

$$\Rightarrow \lambda_{max} = \| \vec{\eta} \odot \vec{\nabla}_\theta z_{1\alpha}^{(L+1)} \|^2$$

$$\Rightarrow \quad \cdots$$

$$\Rightarrow \boxed{\begin{array}{l} \eta_b^{(\ell)} = O(1) \quad \left( \text{or } O\left(\tfrac{1}{L}\right) \right) \\[3mm] \eta_w^{(\ell)} = O\left(\tfrac{1}{\sqrt{n_{\ell-1}}}\right) \quad \left( \text{or } O\left(\tfrac{1}{L\sqrt{n_{\ell-1}}}\right) \right) \end{array}}$$

# Lecture — Pathologies of NTK/GP Regime

**Pathologies:** ① As $n_1, \ldots, n_L \to \infty$, GD on MSE equivalent to
linear/kernel method $\boxed{z_\alpha^{(L+1)}(\theta) \to \tilde{z}_\alpha^{(L+1)}(\theta) \equiv z_\alpha^{(L+1)}(\theta(0)) + \eta \odot \vec{\nabla}_\theta z_\alpha^{(L+1)}(\theta(0))(\theta - \theta(0))}$

learning happens in last layer
and to first order in
hidden layers ② **No feature learning!**

**Fix:** Mean-field init $\quad W_{ij}^{(\ell)} \sim \begin{cases} \mathcal{N}\left(0, \frac{C_w}{n_{\ell-1}}\right) & \ell \leq L \\ \mathcal{N}\left(0, \frac{1}{n_L^2}\right) & \ell = L+1 \end{cases} \quad$ and $\quad \eta_w^{(\ell)} = \eta_b^{(\ell)} = O(1)$

# Lecture 10/3- Loss Hessian

We can summarize the optimization of our network via

- Loss Hessian $\text{Hess}_\theta \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} & \cdots \\ \vdots & \ddots \\ \vdots & & \ddots \end{bmatrix} \approx$ negative inverse of Fisher Information

- NNGP (NN Gaussian Process) $\quad z^{(L)}(D)^T z^{(L)}(D) \in \mathbb{R}^{n_c \times n_c}$
  - infinite-width limit of Bayesian network, such as a randomly initialized NN like Lecture 9/21

- NTK (Neural Tangent Kernel) $\quad \vec{\nabla}_\theta z(D; \dot\theta)(\vec{\nabla}_\theta z(D; \dot\theta))^T$
  - Kernel methods replace learning feature vectorizations with weighting the training inputs and interacting via Kernel $K(\vec{x}, \vec{x}'): \mathcal{X} \times \mathcal{X} \to \mathbb{R}$
  - Kernels are great when $K(\vec{x}, \vec{x}') = \langle \varphi(\vec{x}), \varphi(\vec{x}') \rangle_V$ for some vector space $V$ and some $\varphi: \mathcal{X} \to V$
  - The NTK is a kernel $\varphi: \mathbb{R}^{n_{in}} \times \mathbb{R}^{n_{in}} \to \mathbb{R}^{n_{out} \times n_{out}}$ with
    $$\varphi_{jk}(\vec{x}, \vec{y}; \theta) = \sum_{\theta_i} \partial\theta_i z_j(\vec{x}; \theta) \partial\theta_i z_k(\vec{y}; \theta)$$

  - The NTK represents the influence of the loss gradient $\partial_w \mathcal{L}(w, \vec{y}_i)\big|_{w=z(\vec{x}_i; \theta)}$ w.r.t. example $(\vec{x}_i, \vec{y}_i)$ on the evolution of the NN $z(\cdot, \theta)$ through GD step.
  - In large width (large parameter) limit, NTK is constant & deterministic!

# Hessian Eigenvalues (Sagun et. al.)

Spectrum of $\text{Hess}_\theta \, z(x, \theta(\infty))$   [converged params]
- decomposes into bulk + outliers
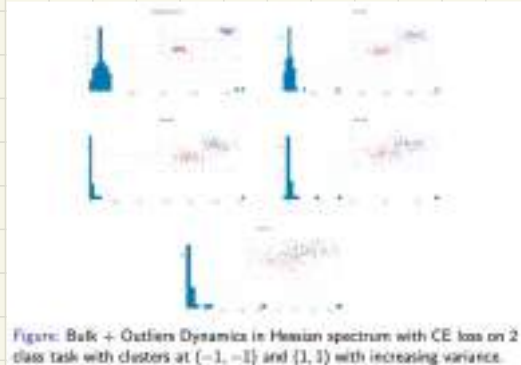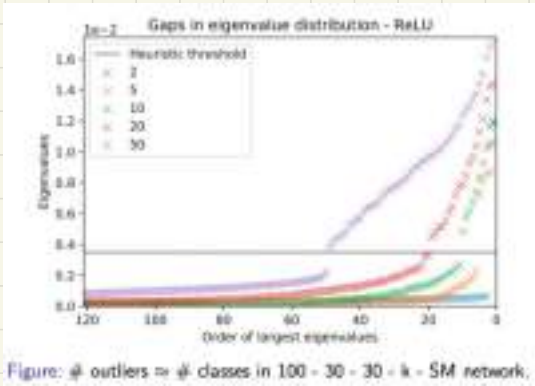- bulk has small eigenvalues (some negative)



Figure: Bulk + Outliers Dynamics in Hessian spectrum with CE loss on 2 class task with clusters at $(-1, -1)$ and $(1, 1)$ with increasing variance.

- # outliers $\approx$ # of classes
- outlier size depends on batch size
- left edge of spectrum gets negative!

# Properties in the Wild

- Hessian has rank at most min{# data, #params}
- Larger eigenvalue ⟹ sharper loss surface, faster optimization



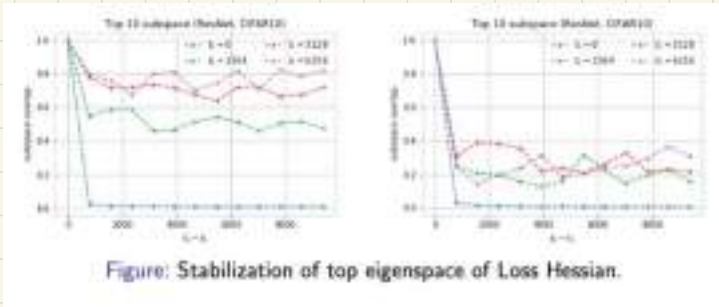Figure: # outliers ≈ # classes in 100 - 30 - 30 - k - SM network.
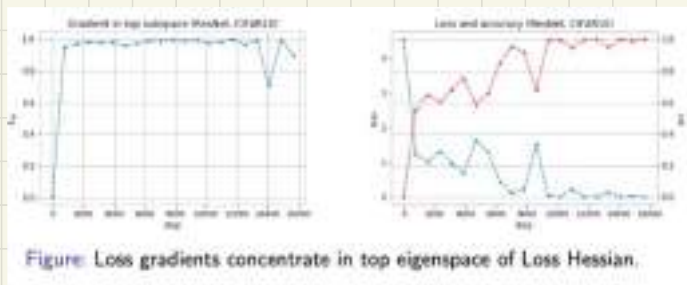
Outlier eigenvalues correspond to class means??

# Hessian Eigenvectors (Gur-Ari, Roberts, Dyer, "Gradient Descent... Happens")

2 results that are robust!

① Top eigenvectors stabilize as training converges.



Figure: Stabilization of top eigenspace of Loss Hessian.

② Loss gradients are in span of top eigenvectors.



Figure: Loss gradients concentrate in top eigenspace of Loss Hessian.

Let $x_{ic}$ be input, $c \in \{1, \ldots, C\}$ is class and $i \in \{1, \ldots, n\}$ a index.
Model output is $f(x_i, c) \in \mathbb{R}^C$ with softmax as $p(x_i, c) = \begin{bmatrix} e^{f(x_i, c)_1} / \sum\limits_{c=1}^{C} e^{f(x_i, c)_{c'}} \\ \vdots \end{bmatrix}$
Let $g_{icc'} := \vec{\nabla}_\theta \ell$ for an example
$x_{ic}$ if assigned label was $c'$.

## Structure of loss Hessian

Gauss-Newton decomposition — where interesting eigenvalues happen

$$H_\theta[\ell(f(\theta))] = \vec{\nabla}_\theta f \, H_f \ell \, \vec{\nabla}_\theta f^T + \vec{\nabla}_\theta \ell \, H_\theta f := G + E$$

$\approx 0$ near converged

For cross-entropy, $G$ is $2^{nd}$ moment matrix:

$$g_{icc'} = \partial_\theta \ell(f(x_{ic}; \theta), y_{c'}) \qquad G = \underset{i, c, c'}{Avg} \{g_{icc'} \, g_{icc'}^T\}$$

We decompose $G$ into $\quad G = G_{class} + G_{cross} + G_{within} + G_{c=c'}$

Covariance in a class $\quad G_{class} = \sum\limits_c w_c g_c g_c^T$

covariance within class group $\quad G_{within} = \sum\limits_{i,c,c'} w_{icc'} (g_{icc'} - g_{cc'})(g_{icc'} - g_{cc'})^T$

covariance between class groups and global incorrect average $\quad G_{cross} = \sum\limits_{c,c'} w_{cc'} (g_{cc'} - g_c)(g_{cc'} - g_c)^T$

where

$$g_{cc'} = \underset{i}{Avg} \{g_{icc'}\} \qquad\qquad g_c = \underset{c \neq c'}{Avg} \{g_{cc'}\}$$

avg. gradient for class $c$ if label was $c'$ $\qquad$ Avg. incorrect gradient

*3 level structure in unsupervised learning? What is $C$, and is the structure there — cool research question*

We see the contributions of different parts to the 3-level structure.
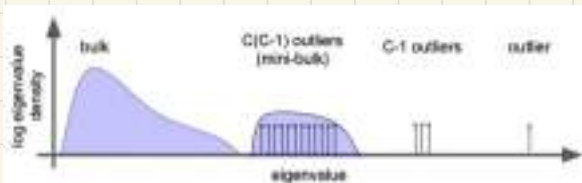(Bulk, $C^2$ outliers with higher eigenvalue of $H$, $C$ outliers with even higher).



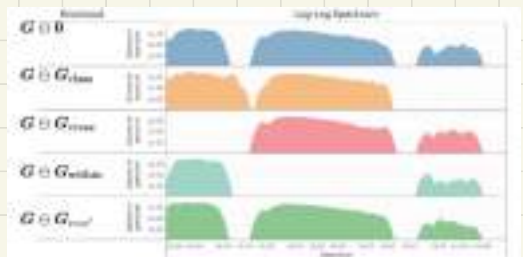Figure: Cartoon of bulk + outlier structure in Hessian spectrum.



Figure: Bulk + outliers in FIM spectrum of VGG11 on CIFAR10

## Structure of activations

Consider $\vec{h}_{ic}^{\ell} = \sigma\left(W^{(\ell)} \vec{h}_{ic}^{\ell-1}\right)$. Let $H^{\ell} \equiv \underset{i,c}{\text{Avg}}\left(h_{ic}^{\ell} \left(h_{ic}^{\ell}\right)^{T}\right)$

↑ post-activations    feature covariance

We decompose $H^{\ell} = H_{class}^{\ell} + H_{within}^{\ell}$

between-class moment
$$H_{class}^{\ell} = \underset{c}{\text{Avg}}\left\{\vec{h}_{c}^{\ell} \vec{h}_{c}^{\ell T}\right\} \quad \text{(mean)}$$

within class 2nd moment
$$H_{within}^{\ell} = \underset{i,c}{\text{Avg}}\left\{\left(\vec{h}_{ic}^{\ell} - \vec{h}_{c}^{\ell}\right)\left(\vec{h}_{ic}^{\ell} - \vec{h}_{c}^{\ell}\right)^{T}\right\} \quad \text{(variance)}$$

where $\vec{h}_{c}^{\ell} \equiv \underset{i}{\text{Avg}}\left\{\vec{h}_{ic}^{\ell}\right\} \qquad \vec{h}_{G}^{\ell} \equiv \underset{c}{\text{Avg}}\left\{\vec{h}_{c}^{\ell}\right\}$

feature class means

We find larger eigenvalues and interesting outlier stuff happening for $H_{class}^{\ell}$. The largest eigenvalue is class-agnostic.



Figure: Eigenvalues of $H^{\ell} \perp H_{class}^{\ell}$ (x axis) vs $H_{class}^{\ell}$ (y axis). Outliers come from $H_{class}^{\ell}$, especially in later layers.



Figure: Eigenvalues of $H^{\ell}$ (blue) vs $H_{class}^{\ell}$ (orange).

## Structure of backprop. grads

Let $\delta_{icc'}^{\ell} = $ layer $\ell$ grads, $\qquad \Delta^{\ell} = \underset{icc'}{\text{Avg}}\left\{\delta_{icc'}^{\ell}\left(\delta_{icc'}^{\ell}\right)^{T}\right\}$

We decompose $\Delta^{\ell} = \Delta_{class}^{\ell} + \Delta_{cross}^{\ell} + \Delta_{within}^{\ell} + \delta_{c \cdot c'}$
where

$$\Delta_{class}^{\ell} = \underset{c}{\text{Avg}}\left\{\delta_{c}^{\ell}\left(\delta_{c}^{\ell}\right)^{T}\right\}$$

$$\Delta_{within}^{\ell} = \underset{i,c}{\text{Avg}}\left\{\left(\delta_{icc'}^{\ell} - \delta_{cc'}^{\ell}\right)\left(\delta_{icc'}^{\ell} - \delta_{cc'}^{\ell}\right)^{T}\right\}$$

$$\Delta_{cross}^{\ell} = \underset{c \neq c'}{\text{Avg}}\left\{\left(\delta_{cc'}^{\ell} - \delta_{c}^{\ell}\right)\left(\delta_{cc'}^{\ell} - \delta_{c}^{\ell}\right)^{T}\right\}$$

with

$\delta_{cc'}^{\ell} = \underset{i}{\text{Avg}}\left\{\delta_{icc'}^{\ell}\right\}$ grad cross-class means

$\delta_{c}^{\ell} = \underset{c \neq c'}{\text{Avg}}\left\{\delta_{cc'}^{\ell}\right\}$ grad class means

# Neural collapse (Papyan, "Prevalence of Neural Collapse")

Call the layer $\ell$ output $\vec{h}_{ic}^{\ell}$.

We want to understand the late-time dynamics of $\vec{h}_{ic}^{\ell}$ via means

$$\vec{\mu}_G = \underset{ic}{Avg}\ \vec{h}_{ic} \qquad\qquad \vec{\mu}_c = \underset{i}{Avg}\ \vec{h}_{ic}$$

<span style="color:blue">global mean</span>             <span style="color:blue">class mean</span>

and covariances

$$\Sigma_B = \underset{c}{Avg}\left\{(\vec{\mu}_c - \vec{\mu}_G)(\vec{\mu}_c - \vec{\mu}_G)^T\right\}$$

$$\Sigma_v = \underset{ic}{Avg}\left\{(\vec{h}_{ic} - \vec{\mu}_c)(\vec{h}_{ic} - \vec{\mu}_c)^T\right\}$$

## Phenomena of Neural Collapse

(1) Variability collapse $\Sigma_w \to 0$ (predictions approach class mean)

(2) $\{\mu_c \mid c \in 1,...,C\}$ approaches simplex vertices (class means are *maximally* orthogonal and same magnitude)

(3) Classification becomes nearest neighbor
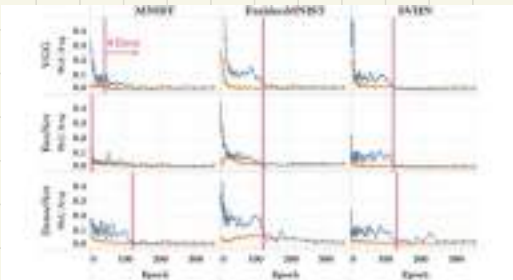
(4) $W_c \approx \vec{\mu}_c - \vec{\mu}_G$



Figure: Coefficients of variation for $\|\mu_c\|$ (orange) and $\|\mu_c - \mu_G\|$ (blue). Simple datasets.



Figure: Mismatch between nearest-neighbor and NN classifiers. Complex datasets.

<span style="color:blue">"Var($\|\vec{\mu}_c\|$) and Var($\|\vec{\mu}_c - \vec{\mu}_G\|$) $\to 0$"
Sorta holds for complex datasets
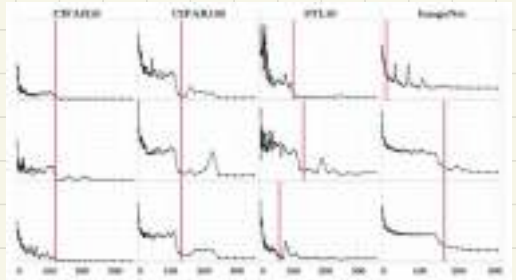**(1)** and **(2)**</span>

<span style="color:blue">"Nearest-neighbor and NN behave similarly"
Doesn't hold for complex datasets
**(3)**</span>

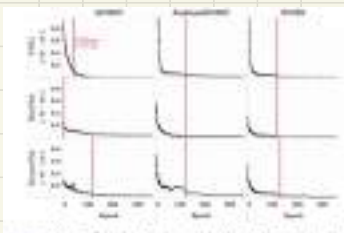<span style="color:blue">$\emptyset$ Pic of $\|W - \{\vec{\mu}_c - \vec{\mu}_G \mid c \in 1,...,C\}\|_F$ for (4)</span>



Figure: $\|M - W\|_F^2$ where $M = \|\mu_c - \mu_G\|$, $W = \|W_c\|$ are final layer weights. Simple datasets.

<span style="color:blue">"final layer approaches class means" **(4)**</span>
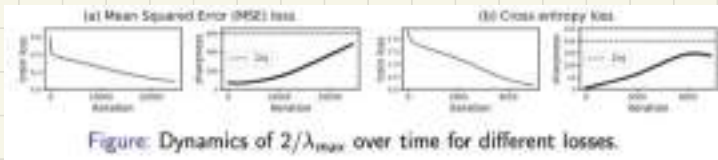
# Sharpness (Cohen et al, "Edge of Stability")

Train a NN with a fixed learning rate $\eta$.
We track "sharpness", or $\lambda_{max} = \lambda_{max}(\text{Hessian})$ over time
← Hessian of loss

The theoretical expectation is that $\eta$ should not be much larger than $2/\lambda_{max}$.

The empirical observation is that $\lambda_{max}$ grows until $\lambda_{max} \approx 2/\eta$

They interpret this that the model finds "sharpest" parts during training so that steps are most meaningful.



Figure: Dynamics of $2/\lambda_{max}$ over time for different losses.

Sharpness $\lambda_{max}$ approaches $2/\eta$

# Large Learning Rates - (Lewkowycz et al, "Catapult Phase")

We ask about fixing the NN and varying large $\eta$.
The finding is three phases

- lazy phase $0 < \eta < 2/\lambda_{max}$ (NTK)
  output sensitive to parameters
- catapult phase $2/\lambda_{max}$ (NTK) $< \eta < C_*/\lambda_{max}$ (NTK) ← mystery threshold
- divergent phase $C_*/\lambda_{max}$ (NTK) $< \eta$

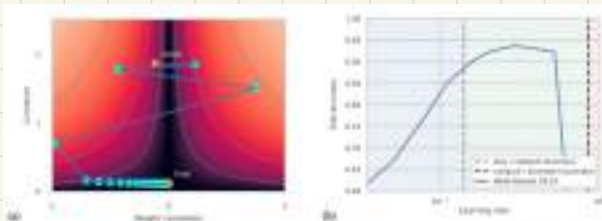here things look like losses grow divergently, but then they walk back



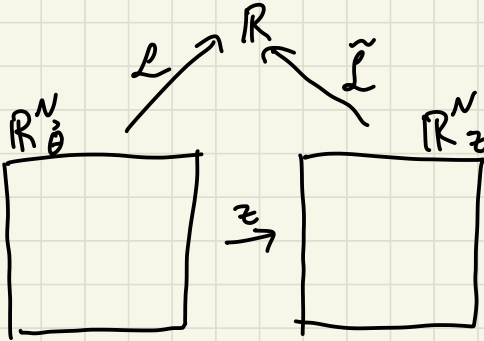Figure: Weight correlation over training and test accuracy on CIFAR10 with fixed number of training steps.

Best $\eta$ is in catapult region.

# Lecture 10/10 - Intro to NTK

Consider GD with

$$\vec{\theta}(t+1) = \vec{\theta}(t) - \gamma \vec{\nabla}_\theta \mathcal{L}(\vec{\theta}(t))$$

Suppose that $\mathcal{L}(\vec{\theta}) = \tilde{\mathcal{L}}(z(\vec{\theta}))$ for some $z$: ← change of coordinates



In $z(\vec{\theta})$ variables,

$$z(t+1) = z(\vec{\theta}(t+1))$$
$$= z(\vec{\theta}(t)) - \gamma \vec{\nabla}_\theta \mathcal{L}(\vec{\theta}(t))$$
$$= z(\vec{\theta}(t)) - \gamma \vec{\nabla}_\theta \tilde{\mathcal{L}}(z(\vec{\theta}(t)))$$
$$= z(\vec{\theta}(t)) - \gamma \vec{\nabla}_\theta z|_{\theta(t)} \cdot \vec{\nabla}_z \tilde{\mathcal{L}}|_{z=z(\vec{\theta}(t))} + O(\gamma^2)$$
$$= z(t) - \gamma (\vec{\nabla}_\theta z)^T (\vec{\nabla}_\theta z)|_{\theta(t)} \cdot \vec{\nabla}_z \tilde{\mathcal{L}}(z(t))$$

Jacobian $\in \mathbb{R}^{N \times n}$

So,

$$\boxed{z(t+1) = z(t) - \gamma K_{\vec{\theta}(t)} \vec{\nabla}_z \tilde{\mathcal{L}}(z(t))}$$

where $K_{\vec{\theta}(t)} \equiv (\nabla_{\vec{\theta}} z)^T (\nabla_{\vec{\theta}} z) \in \mathbb{R}^{n \times n}$

Neural Tangent Kernel

The picture looks like



$K_\theta$ makes update steps move along the manifold allowed by $\text{im}(z)$

$K_\theta \vec{\nabla}_z \tilde{\mathcal{L}}$

$\nabla_z \tilde{\mathcal{L}}$

# Ex/ MSE

Consider a NN $z(\vec{x}; \vec{\theta})$ and $m$ training data points
$D = \{(\vec{x}_i, y_i), i = 1, \ldots, m\}$ and loss
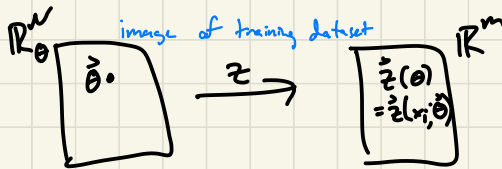$$\mathcal{L}(\vec{\theta}) = \tilde{\mathcal{L}}(z(\vec{\theta})) = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2}\left(y_i - z(\vec{x}_i, \vec{\theta})\right)^2$$

We can use the change of coordinates induced by $z$ to get that
$$\vec{z}(t) = \{z(\vec{x}_i; \vec{\theta}(t))\}_{i=1}^{m} \quad \text{and} \quad \vec{Y} = \{y_i\}_{i=1}^{m}$$


image of training dataset

We have $\mathcal{L}(\vec{\theta}(t)) = \frac{1}{2}\|\hat{Y} - \vec{z}(t)\|^2$ and $\vec{z}(t+1) = \vec{z}(t) - \eta K_{\theta(t)}(\vec{z}(t) - \vec{Y})$

where
$$\boxed{(K_{\theta(t)})_{ij} = \left(\vec{\nabla}_\theta z(\vec{x}_i; \vec{\theta}(t))\right)^T \vec{\nabla}_\theta z(\vec{x}_j; \vec{\theta}(t)) \qquad i,j \in \{1, \ldots, m\}}$$
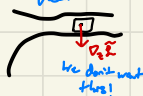Gram Matrix

## Key points!
- If $K_{\theta(t)} = K$ is independent of $\vec{\theta}$, this is "Kernel methods"
  on $\tilde{\mathcal{L}}(z) = \frac{1}{2}\|\hat{Y} - \vec{z}\|^2$. This is a time-varying Kernel
- Suppose $\exists \lambda_0 > 0$ s.t. $\forall t \geq 0$, $\lambda_{min}(K_{\theta(t)}) \geq \lambda_0 \Longleftrightarrow \lambda_\infty > K_{\theta(t)} \geq \lambda_0 I$ $\boxed{(\#)}$

  z-image manifold

  if $K_0$ PSD
  $\Rightarrow \langle K_0 \vec{\nabla}_z \tilde{\mathcal{L}}, \vec{\nabla}_z \tilde{\mathcal{L}} \rangle > 0$,
  we can always move
  in the direction given by
  $\vec{\nabla}_z \tilde{\mathcal{L}}$

  $K$ is PD and bounded
  This condition promises successful optimization.

## Proof:
$\mathcal{L}(\vec{\theta}(t+1)) = \frac{1}{2}\|\hat{Y} - \vec{z}(t+1)\|^2$ and $\vec{z}(t+1) - \vec{Y} = (I - \eta K_{\theta(t)})(\vec{z}(t) - \vec{Y})$

Thus, if $\eta < \frac{1}{\lambda_{max}}$, $\mathcal{L}(\vec{\theta}(t+1)) = \frac{1}{2}\|(I - \eta K_{\theta(t)})(\vec{z}(t) - \vec{Y})\|^2$

$\leq \frac{1}{2}\|\vec{z}(t) - \vec{Y}\|^2 (1 - \eta \lambda_0)^2 \leq \mathcal{L}(\theta(t)) e^{-2\eta\lambda_0}$

$\Rightarrow \mathcal{L}(\vec{\theta}(t+1)) \leq e^{-2\eta\lambda_0} \mathcal{L}(\theta(0))$

$\square$

<u>The goal is as follows:</u>

For wide NNs w/ NTK init, MSE loss, small $\eta$, and $L < \infty$ fixed, the "Meta Theorem" is that $K_{\theta(t)}$ satisfies **(#)** (bounded PD).
$$\Rightarrow \mathcal{L} \to 0.$$

The intuition is that if $K_{\theta(t)} \succeq \lambda_0 I$, we can move $z(x_i; \hat{\theta})$ at will and always make progress. So, the data points cannot fight each other.

> To show **(#)** typically,
> (i) Show $K_{\theta(0)} \succeq \lambda I$    and   (ii) Show   $\sup\limits_{t > 0} \| K_{\theta(t)} - K_{\theta(0)} \| \leq \dfrac{\lambda_0}{2}$

**Ex/ Simple NN**   $z_\alpha^{(2)} = \sum\limits_{j=1}^{n} \frac{1}{\sqrt{n}} W_j^{(2)} \sigma(W_j^{(1)} \cdot \vec{x}_\alpha)$   with 1D output and 1 hidden layer.

Suppose that $|\sigma|, |\sigma'|, |\sigma''| \leq 1$ and $\|\vec{x}_\alpha\| = 1$ and we have NTK init
$$\left. \begin{array}{l} W_j^{(2)} \sim N(0,1) \\[4pt] W_j^{(1)} \sim N(0, I) \end{array} \right\} \dot{\theta}(0)$$

(i) We have
$$\left( K_{\theta(0)} \right)_{\alpha\beta} = \sum_{k=1}^{n} \partial_{W_k^{(2)}} z_\alpha^{(2)} \, \partial_{W_k^{(2)}} z_\beta^{(2)} + \langle \partial_{V_k^{(1)}} z_\alpha^{(2)}, \partial_{V_k^{(1)}} z_\beta^{(2)} \rangle$$

$$= \frac{1}{n} \sum_{k=1}^{n} \underbrace{\sigma(W_k^{(1)} \vec{x}_\alpha) \sigma(W_k^{(1)} \vec{x}_\beta)}_{k_{\alpha\beta}^{(2)}} + \underbrace{(W_k^{(2)})^2 \sigma'(W_k^{(1)} \vec{x}_\alpha) \sigma'(V_k^{(1)} \vec{x}_\beta) \vec{x}_\alpha \cdot \vec{x}_\beta}_{k_{\alpha\beta}^{(1)}}$$

$\Rightarrow K_{\theta(0)} \succeq K_{\theta(0)}^{(2)}$    We only need $K_{\theta(0)}^{(2)} \succeq \lambda_0 I$ !

Now,   $\left( K_{\theta(0)}^{(2)} \right)_{\alpha\beta} = \frac{1}{n} \sum\limits_{k=1}^{n} \sigma(z_{k;\alpha}^{(1)}) \sigma(z_{k;\beta}^{(1)}) = \frac{1}{n} \sum\limits_{k=1}^{n} K_{k;\theta(0)}^{(2)}$ ← avg. of iid matrices

<u>Idea:</u> we will write $K_{\theta(0)}^{(2)} = \mathbb{E}\{K_{1;\theta(0)}^{(2)}\} + \frac{1}{n} \sum\limits_{k=1}^{n} K_{k;\theta(0)}^{(2)} - \mathbb{E}\{K_{k;\theta(0)}^{(2)}\}$

where $\left( K_{k;\theta(0)}^{(2)} \right)_{\alpha\beta} = \sigma(z_{k\alpha}^{(1)}) \sigma(z_{k\beta}^{(1)})$   to get concentration bound on
$$K_{\theta(0)}^{(2)} - \mathbb{E}\{K_{\theta(0)}^{(2)}\}$$

**Theorem:** Matrix Bernstein Inequality

Let $Z = \sum_{j=1}^{n} S_j$, where $S_j \sim$ iid with $\mathbb{E}\{S_j\} = 0 \;\forall j$

and $\|S_j\|_{op} \leq L$. ← largest eigenvalue of $S_j$

Let $v = \max\left\{ \left\|\sum_{j=1}^{n} \mathbb{E}\{S_j S_j^T\}\right\|_{op}, \left\|\sum_{j=1}^{n} \mathbb{E}\{S_j^T S_j\}\right\|_{op} \right\}$.

Then, $\mathbb{P}\{\|Z\|_{op} > t\} \leq e^{\frac{-t^2/2}{v + Lt/3}}$

For us, $S_j = \frac{1}{n} K^{(2)}_{j; \theta(0)} - \mathbb{E}\{K^{(2)}_{j; \theta(0)}\} \Rightarrow \|S_j\|_{op} \leq \left\|\frac{1}{n} K^{(2)}_{j, \theta(0)}\right\|_{\infty} \cdot m \leq \frac{2m}{n}$

Similarly, $v \leq c\frac{m}{n}$ ← N.B. $\Rightarrow \mathbb{P}\{\|K^{(2)}_{\theta(0)} - \mathbb{E}\{K^{(2)}_{\theta(0)}\}\|_{op} > t\} \leq e^{\frac{-c\,t^2}{(1+t)\frac{m}{n}}}$  set $t = \sqrt{\frac{m}{n}}$

$\Rightarrow \|K^{(2)}_{\theta(0)} - \mathbb{E}\{K^{(2)}_{\theta(0)}\}\| \leq c\sqrt{\frac{m}{n}}$ with high probability.

So, $K^{(2)}_{\theta(0)}$ concentrates well about the mean. We now want to show the result for the expectation.

We WTS that if

- $\theta$ is not poly
- $\vec{x}_\alpha \neq \vec{x}_\beta$ if $\alpha \neq \beta$ $\forall \alpha, \beta$
- $\|x_\alpha\| = 1$ $\forall \alpha$

$\Rightarrow$ $\mathbb{E}\{K^{(2)}_{\theta(0)}\} \geq \lambda_0 I \iff \left(\mathbb{E}\{\sigma(w^{(1)} \cdot \vec{x}_\alpha) \sigma(w^{(1)} \cdot \vec{x}_\beta)\}\right)_{\alpha\beta} \geq \lambda_0 I$

Note that we can move from expectations in $\{\vec{x}_\alpha\}$ space to an infinite dimensional Hilbert space $\mathcal{H} = \{\mathcal{I}_\alpha\}$ s.t. $\mathcal{H} = L^2(\mathbb{R}^{n_0}, e^{-\frac{1}{2}\|w^{(1)}\|^2})$ ← inner products in $\mathcal{H}$ are expectations over $w^{(1)}$

$\left(\text{Also, } \mathcal{H} = \{f: \mathbb{R}^{n_0} \to \mathbb{R} \mid \mathbb{E}\{f(w)^2\} < \infty\}\right)$

So, $\mathcal{H}$ gives $\mathbb{E}\{K^{(1)}_{\theta(0)}\}_{\alpha\beta} = \langle \mathcal{I}_\alpha, \mathcal{I}_\beta \rangle_{\mathcal{H}}$ where $\mathcal{I}_\alpha^{(w)} = \sigma(w \cdot \vec{x}_\alpha)$

$\Rightarrow \mathbb{E}\{K^{(2)}_{\theta(0)}\} = \begin{bmatrix} \langle \mathcal{I}_1, \mathcal{I}_1 \rangle & \langle \mathcal{I}_1, \mathcal{I}_2 \rangle & \cdots \\ \vdots & \ddots & \end{bmatrix}$ is a <span style="color:red">Gram Matrix</span>

**Theorem:** (Gram)
For a Gram matrix $A = B^T B$, the following are equivalent:

(1) $A > 0$    (2) $\det A > 0$    (3) vol(Parallelepiped $(\{B_i\}))^2 > 0$    (4) All rows $\{B_i\}$
$A$ is PD                                     rows that generate $A$       linearly independent

We want to show that $\{\Xi_\alpha\}_{\alpha=1}^m$ are linearly independent in $\mathcal{H}$, as this will give us that $\mathbb{E}\{k^{(2)}_{\theta(0)}\} \succ 0$ by the above theorem.

As usual, suppose that
$$\sum_{\alpha=1}^m c_\alpha \Xi_\alpha = 0 \text{ in } \mathcal{H} \text{ for some } c_\alpha\text{'s.}$$

We want to show that this implies $c_\alpha = 0 \;\forall \alpha$. Now,
$$\sum_{\alpha=1}^m c_\alpha \Xi_\alpha = 0 \text{ in } \mathcal{H} \iff \forall f \in \mathcal{H}, \sum_{\alpha=1}^m c_\alpha \langle \Xi_\alpha, f \rangle_{\mathcal{H}} = 0$$
$$\iff \forall f \in \mathcal{H} \sum_{\alpha=1}^m c_\alpha \mathbb{E}\{\sigma(w_{x_\alpha}) f(w)\} = 0$$

Since the <span style="color:red">Hermite polynomials</span> are orthogonal w.r.t. weight measure $e^{-x^2}$, we can use them as an orthonormal basis for $\mathcal{H}$ to decompose $\sigma$:
$$\sigma(t) = \sum_{j=0}^\infty \frac{\sigma_j}{\sqrt{j!}} H_j(t) \qquad \text{\color{blue}($\sigma$ non-poly $\Rightarrow \sigma_k \neq 0 \;\forall k$)}$$

Let $\beta$ be arbitrary. Since our assumption holds $\forall f \in \mathcal{H}$, clearly it holds for $\{f_k(w)\}_{k=1}^\infty$, where $f_k(w) = \frac{\sigma_n}{\sqrt{k!}} H_k(w \cdot \vec{x}_\beta)$

The assumption gives
$$\forall k \in \mathbb{N}, \quad 0 = \sum_{\alpha=1}^m c_\alpha \mathbb{E}\{\sum_{j=0}^\infty \frac{\sigma_j}{\sqrt{j!}} H_j(w \cdot \vec{x}_\alpha) \frac{\sigma_k}{\sqrt{k!}} H_k(w \cdot \vec{x}_\beta)\}$$
$$\text{\color{blue}Hermite polynomials orthonormal} = \sum_{\alpha=1}^m c_\alpha (\vec{x}_\alpha \cdot \vec{x}_\beta)^k$$

As $k \to \infty$, we find that $(\vec{x}_\alpha \cdot \vec{x}_\beta) \to \delta_{\alpha\beta} \Rightarrow c_\beta = 0$.

This line of reasoning holds for all $\beta$, and so all the $c_\alpha$'s are 0. This means that the $\{\Xi_\alpha\}_{\alpha=1}^m$ are linearly independent in $\mathcal{H}$.

So, by the Gram Theorem,
$$\mathbb{E}\{k^{(2)}_{\theta(0)}\} = \text{Gram}\left(\{\Xi_\alpha\}_{\alpha=1}^m\right) \succ 0.$$

Since $k^{(2)}_{\theta(0)}$ concentrates well about its expectation and $K_{\theta(0)} \succ k^{(2)}_{\theta(0)}$, we achieve the result that the NTK $K_{\theta(0)}$ is PD at $t=0$.

$$\square$$

# Lecture 10/12 - NTK sends $\mathcal{L} \to 0$

Recall that we consider the small example

$$z_\alpha^{(2)}(t) = z_\alpha^{(2)}(\Theta(t)) = \frac{1}{\sqrt{n}} \sum_{i=1}^{\hat{n}} w_i^{(2)}(t) \sigma(w_i^{(1)}(t) \vec{x}_\alpha)$$

with $w_i^{(1)} \sim \text{Unif}([-1,1])$, $w_i^{(1)} \sim N(0, I)$, $\|\sigma\|_\infty, \|\sigma'\|_\infty, \|\sigma''\|_\infty \leq 1$.
We inspect gradient descent on MSE

$$\mathcal{L}(\Theta) = \frac{1}{2m} \sum_{j=1}^{\hat{m}} \left(z_{\alpha_j}^{(2)}(\Theta) - y_{\alpha_j}\right)^2 \qquad \Theta(t+1) = \Theta(t) - \eta \, \vec{\nabla}_\Theta \mathcal{L}(\Theta(t))$$

Assume the following:

$$\frac{d}{dt}\Theta(t) = -\eta \, \vec{\nabla}_\Theta \mathcal{L}(\Theta(t)), \qquad\qquad\qquad \underset{\substack{\text{← freeze 2nd} \\ \text{layer}}}{w_i^{(2)}(t) = w_i^{(2)}(0)}$$

We still have the NTK

$$\boxed{\left(K_{\Theta(t)}\right)_{ij} = \left(\vec{\nabla}_\Theta \, z_{\alpha_i}^{(2)}(\Theta(t))\right)^T \left(\vec{\nabla}_\Theta z_{\alpha_j}^{(2)}(\Theta(t))\right)}$$
$$\underset{m \times m \text{ Gram matrix}}{}$$

## The overall goal: Show that w.h.p. $\mathcal{L}(\Theta(t)) \xrightarrow{t \to \infty} 0$

Last time we split this into two subproblems:
- (i) $\exists \lambda_0 > 0$ s.t. $K_{\Theta(0)} \succeq \lambda_0 I$ w.h.p. $\quad$ (K_{\Theta(0)} is PD, showed this last time)
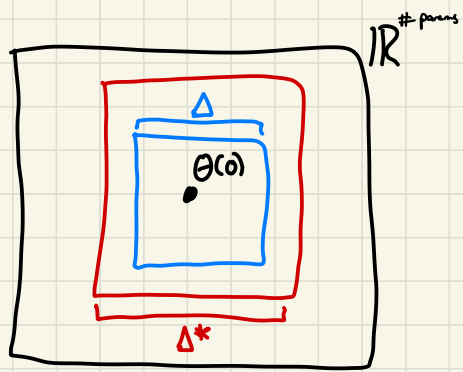
- (ii) $K_{\Theta(t)} \succeq \frac{\lambda_0}{2} I \quad \forall t \geq 0 \quad$ (K_{\Theta(t)} stays PD, show this this time)

In other words, today we want to show

$$\boxed{\forall t \geq 0, \quad \|K_{\Theta(t)} - K_{\Theta(0)}\|_{op} \leq \frac{\lambda_0}{2}}$$

The idea is as follows: (Du et. al.)



$\mathbb{R}^{\# \text{params}}$

(1) $\|\theta - \theta(0)\| \leq \Delta(n, m, \lambda_0, \dots)$ ← stay within blue box

$\Rightarrow \|K_\theta - k_{\theta(0)}\|_{op} \leq \frac{\lambda_0}{4}$ ← $K_\theta$ stays positive

We WTS the implication and that the premises hold.

(2) While $K_{\theta(t)} \geq \frac{\lambda_0}{2} I$, $\mathcal{L}(\theta(t))$ decays exponentially, $\|\frac{d}{dt}\theta(t)\|^2 \approx \mathcal{L}(t)$

We want $\Delta^* \leq \Delta$ so that we never leave the box of size $\Delta$ so that (1) will give us $K_{\theta(t)} \geq \frac{\lambda_0}{2} I$.

We first show the implication in (1).

**Lemma 1:** Let $\Delta \in (0, 1]$. If $\forall i, \|w_i^{(i)} - w_i^{(i)}(0)\| \leq \Delta$, then $\|K_\theta - k_{\theta(0)}\|_{op} \leq 2m\Delta$

If params don't change too much, $K_{\theta(t)}$ stays PD.

**Proof:** We have

$(K_\theta)_{ij} = \frac{1}{n}\sum_{k=1}^{n} w_k^{(2)}(0)^2 \underbrace{[\vec{x}_{\alpha_i} \cdot \vec{x}_{\alpha_j}]}_{|\cdot| \leq 1} \sigma'(w_k^{(i)} \cdot \vec{x}_{\alpha_i}) \sigma'(w_k^{(i)} \vec{x}_{\alpha_j})$

Note that because $\sigma'$ is bounded, we see that

$w \in \mathbb{R}^{n_0} \mapsto \sigma'(w\vec{x}_{\alpha_i}) \sigma'(w\vec{x}_{\alpha_j})$ is 2-Lipschitz.

To see this,

bounded differences $\begin{cases} |\sigma'(w\vec{x}_{\alpha_i}) \sigma'(w\vec{x}_{\alpha_j}) - \sigma'(\bar{w}\vec{x}_{\alpha_i}) \sigma'(\bar{w}\vec{x}_{\alpha_j})| \\ = |(\sigma'(w\vec{x}_{\alpha_i}) - \sigma'(\bar{w}\vec{x}_{\alpha_i}))\sigma'(w\vec{x}_{\alpha_j}) + \sigma'(\bar{w}\vec{x}_{\alpha_i})(\sigma'(w\vec{x}_{\alpha_j}) - \sigma'(\bar{w}\vec{x}_{\alpha_j}))| \\ \leq 2\|w - \bar{w}\| \end{cases}$

← largest diff.

Thus, $\|K_\theta - k_{\theta(0)}\|_{\infty} \leq 2\Delta$. Lastly, since $A \in \mathbb{R}^{m \times m}$, $\|A\|_{op} \leq m\|A\|_{\infty}$.

$\Rightarrow \|K_\theta - k_{\theta(0)}\|_{op} \leq 2m\Delta.$ □

**Corollary:** If $\|w_i^{(i)}(t) - w_i^{(i)}(0)\| \leq \frac{\lambda_0}{8m}$ $\forall t \geq 0$, $\Rightarrow K_{\theta(t)} \geq \frac{\lambda_0}{4} I$ $\forall t \geq 0$.

This tells us that we wish to set $\boxed{\Delta \equiv \frac{\lambda_0}{8m}}$

This proves the implication of (1). if $\|W_i^{(1)}(t) - W_i^{(1)}(0)\| \leq \frac{\lambda_0}{8m}$ $\forall t \geq 0$

$$\Rightarrow K_{\theta(t)} \geq \frac{\lambda_0}{4} I \quad \forall t \geq 0.$$

Now, all that is left to show is that

$$\|\theta(t) - \theta(0)\| \leq \int_0^\infty \|\frac{d}{ds}\theta(s)\| ds \leq \Delta^* \quad \text{for some } \Delta^* \leq \Delta = \frac{\lambda_0}{8m}$$

With this, we can use the Corollary to show that $K_{\theta(t)}$ stays P.D. We now show the premises.

<u>Lemma 2:</u> Fix $t \geq 0$ and suppose that $\forall s < t$, $\quad K_{\theta(s)} \geq \frac{\lambda_0}{2} I \quad (*)$

Then $\forall s < t$, $\|W_i^{(1)}(s) - W_i^{(1)}(0)\| \leq \Delta^* = \frac{2m \, \mathcal{L}(0)^{\frac{1}{2}}}{\lambda_0 \, n^{\frac{1}{2}}}$ <span style="color:blue">If $K_\theta$ stays PD, the params don't change too much</span>

<u>Proof:</u> We have $\|W_i^{(1)}(s) - W_i^{(1)}(0)\| = \|\int_0^s \frac{d}{d\tau} W_i^{(1)}(\tau) d\tau\| \leq \int_0^\infty \|\frac{d}{d\tau} W_i^{(1)}(\tau) d\tau\|$

for fixed $i, \tau$, we compute
$$\frac{d}{d\tau} W_i^{(1)}(\tau) = -\eta \, \partial_{W_i^{(1)}} \mathcal{L}(\theta(\tau)) = -\eta \, \partial_{W_i^{(1)}}\left[\frac{1}{2m}\sum_{j=1}^n \left(z_{\alpha_j}^{(2)}(\tau) - y_{\alpha_j}\right)^2\right]$$

$$= \frac{-\eta}{m} \sum_{j=1}^n \left(z_{\alpha_j}^{(2)}(\tau) - y_{\alpha_j}\right)\left(\partial_{W_i^{(1)}} z_{\alpha_j}^{(2)}(\tau) \cdot \vec{x}_{\alpha_j}\right)$$

$$= \frac{-\eta}{\sqrt{n} \, m} \sum_{j=1}^n \left(z_{\alpha_j}^{(2)}(\tau) - y_{\alpha_j}\right) W_i^{(2)}(0) \, \sigma'\left(W_i^{(1)}(\tau)\vec{x}_{\alpha_j}\right)$$

<span style="color:blue">|.| bounded by 1</span>

$$\Rightarrow \|\frac{d}{d\tau} W_i^{(1)}(\tau)\| \leq \frac{\eta}{m} \sum_{j=1}^n |z_{\alpha_j}^{(2)}(\tau) - y_{\alpha_j}|$$

We can use the Power-Mean inequality:
$$\boxed{\forall p < p', \quad \left(\frac{1}{m}\sum_{j=1}^n a_j^p\right)^{\frac{1}{p}} \leq \left(\frac{1}{m}\sum_{j=1}^n a_j^{p'}\right)^{\frac{1}{p'}}}$$

$$\Rightarrow \|\partial_{W_i^{(1)}} \mathcal{L}(\theta(\tau))\| \leq \frac{\eta}{\sqrt{n}} \mathcal{L}(\theta(\tau))^{\frac{1}{2}}. \text{ Therefore,}$$

<span style="color:blue">← apply part (2) because (*)</span>

$$\|W_i^{(1)}(s) - W_i^{(1)}(0)\| \leq \frac{\eta}{m}\int_0^\infty e^{-\frac{\eta \lambda_0 \tau}{2m}} d\tau \, \mathcal{L}(0)^{\frac{1}{2}} = \frac{\eta}{\sqrt{n}} \cdot \frac{2m}{\lambda_0 \eta} \mathcal{L}(0)^{\frac{1}{2}}$$

$$= \frac{2m \mathcal{L}(0)^{\frac{1}{2}}}{\lambda_0 \, n^{\frac{1}{2}}}.$$

This tells us to set $\boxed{\Delta^* = \frac{2m \, \mathcal{L}(0)^{\frac{1}{2}}}{\lambda_0 \, n^{\frac{1}{2}}}}$ $\square$

Here is a quick proof of (2), which we used above.

$$\boxed{\underline{\text{Recall:}}\ \frac{d}{dt}\left(z^{(2)}(t) - y\right) = -\frac{\eta}{m} K_{\theta(t)}\left(z^{(2)}(t) - y\right)}$$

To see (2) $\left(K_{\theta} \text{ PD} \Rightarrow \mathcal{L}(t) \text{ exponential decay}\right)$, note that

$$\underset{\text{MSE}}{\mathcal{L}(t)} = \frac{1}{2m}\left(z^{(2)}(t) - y\right)^{T}\left(z^{(2)}(t) - y\right) \Rightarrow \frac{d}{dt}\mathcal{L}(t) = -\frac{\eta}{m}\left(z(t) - y\right)^{T} K_{\theta(t)}\left(z(t) - y\right)$$

$$\left(K_{\theta} \geq \tfrac{\lambda_0}{2} I\right) \leq -\frac{\eta}{m} \lambda_0 \mathcal{L}(t)$$

$$\Rightarrow \mathcal{L}(t) \leq e^{-\frac{\eta \lambda_0}{m} t} \mathcal{L}(0)$$

At this point, we proved that

Lemma 1   $\|w_i^{(1)}(t) - w_i^{(1)}(0)\| \leq \Delta \Rightarrow K_{\theta(t)} \geq \tfrac{\lambda_0}{2} I$   &

Lemma 2   $\forall s < t,\ K_{\theta(s)} \geq \tfrac{\lambda_0}{2} I \Rightarrow \|w_i^{(1)}(s) - w_i^{(1)}(0)\| \leq \Delta^{*}$

Suppose that $\Delta^{*} < \Delta \iff \frac{2m\,\mathcal{L}(0)^{\frac{1}{2}}}{\lambda_0 n^{\frac{1}{2}}} < \frac{\lambda_0}{2m} \iff n \geq \frac{16 m^4 \mathcal{L}(0)}{\lambda_0^4}$

Define   $t_K = \inf\{t > 0 \text{ s.t. } K_{\theta(t)} \leq \tfrac{\lambda_0}{2} I\}$ ← first $t$ that NTK isn't PD enough

$t_{\Delta} = \inf\{t > 0 \text{ s.t. } \exists j \in \{1,\ldots,n\} \text{ s.t. } \|w_j^{(1)}(t) - w_j^{(1)}(0)\| > \Delta^{*}\}$
first $t$ s.t. weights grow a lot

$t^{*} = \min\{t_K, t_{\Delta}\}$

We claim that $\underline{t^{*} \text{ must be } \infty.}$   ?

$\underline{\text{Proof:}}$ Suppose BWOC that $t^{*} < \infty$.
  $\underline{\text{Case 1:}}\ t^{*} = t_{\Delta} \leq t_K$
  Then, $\forall t < t^{*}$, we have $\|w_i^{(1)}(t) - w_i^{(1)}(0)\| \leq \Delta^{*} < \Delta \overset{\text{Lemma 1}}{\Rightarrow} K_{\theta(t)} \geq \tfrac{\lambda_0}{2} I$
  $\to\leftarrow$ by definition of $t_K$.
  $\underline{\text{Case 2:}}\ t^{*} = t_K \leq t_{\Delta}$
  Then, $\forall t < t^{*}$, we have $K_{\theta(t)} \geq \tfrac{\lambda_0}{2} I \overset{\text{Lemma 2}}{\Rightarrow} \|w_i^{(1)}(t) - w_i^{(1)}(0)\| \leq \Delta^{*} < \Delta$
  $\to\leftarrow$ by our definitions.
  Thus, $t^{*} = \infty$.   □

So, we showed that the weights always stay within $\Delta$ and therefore that the NTK is always PD. Applying part (i) as $t \to \infty$, we have shown that $\mathcal{L}(t) \to 0$!

# Lecture 10/31- Kernels

**Def:** Let $\Omega \subset \mathbb{R}^d$. A **Kernel** on $\Omega$ is $K: \Omega \times \Omega \to \mathbb{R}$
s.t. $\forall \vec{x}_1, \ldots, \vec{x}_N, \vec{y} \in \Omega, \forall a_1, \ldots, a_N \in \mathbb{R}$
- $K(\vec{x}, \vec{y}) = K(\vec{y}, \vec{x})$
- $K$ is "positive" $\Leftrightarrow \sum_{i,j=1}^{N} a_i a_j K(\vec{x}_i, \vec{x}_j) > 0$ if $\|\vec{a}\| \neq 0$

We can think of $K$ as a infinite analog of positive definite matrices.

Ex/ 1) If $\Omega$ finite $(\vec{x}_1, \ldots, \vec{x}_n \in \mathbb{R}^d)$, $K \in \mathbb{R}^{k \times k}$, $K(\vec{x}_i, \vec{x}_j) = K_{ij}$
$\Rightarrow$
- $K$ is symmetric
- $\sum_{i,j=1}^{N} a_i a_j K(\vec{x}_i, \vec{x}_j) = \vec{a}^T K \vec{a} > 0$ if $\vec{a} \neq 0$

2) $\Omega = \mathbb{R}^d$, $K(\vec{x}, \vec{y}) = \langle \vec{x}, \vec{y} \rangle$
- dot product is commutative
- $\sum_{i,j=1}^{N} a_i a_j K(\vec{x}_i, \vec{x}_j) = \| \sum_i a_i \vec{x}_i \|^2$

3) $\Omega \in \mathbb{R}^d$, $K(\vec{x}, \vec{y}) = e^{-\|x-y\|^2/2\sigma^2}$

The general case is defined via feature maps!

**Def:** The **feature map** $\Xi: \Omega \to \mathcal{H}$ ($\leftarrow$ arbitrary Hilbert space) is given by

$K(\vec{x}, \vec{y}) = \langle \Xi(\vec{x}), \Xi(\vec{y}) \rangle_{\mathcal{H}}$

where
$$\vec{x} \underset{\underset{\Omega}{\uparrow}}{\longmapsto} \underset{}{\Xi} \underbrace{\langle \Psi_1(\vec{x}), \Psi_2(\vec{x}), \ldots \rangle}_{\substack{\text{vector in } \mathcal{H} \text{ of coefficients} \\ \text{in the ONB}}}$$

where $\{\Psi_j\}$ is an ONB of $\mathcal{H}$.

<u>Theorem</u>: Every Kernel comes from a feature map.

<u>Proof:</u> ($\Omega$ compact, $K \approx c^o$)

Fix $\mu \in P(\Omega)$ as a subset, and let $T_K: L^2(\Omega, \mu) \circlearrowleft$ be defined s.t.
<span style="color:blue">any measure on $\Omega$</span>

$$(T_K f)(\dot{x}) = \int_\Omega K(\dot{x}, \dot{y}) f(\dot{y}) d\mu(\dot{y})$$

Note that $T_K$ is compact. We can apply the spectral theorem:

$$T_K = \sum_{j=0}^\infty \lambda_j \, \varphi_j \, \varphi_j^T \quad \text{for an orthonormal basis } \{\psi_j\}$$

Moreover, $K(\dot{x}, \cdot) \in L^2(\Omega, \mu)$

$$\Rightarrow \forall \dot{x} \in \Omega, \quad K(\dot{x}, \dot{y}) = \sum_{j=0}^\infty a_j(\dot{x}) \, \varphi_j(\dot{y}) \qquad \text{<span style="color:blue">($\{\varphi_j\}$ ONB)</span>}$$

Further,
$$\lambda_j \varphi_j(\dot{x}) = (T_K \varphi_j)(\dot{x}) = \int_\Omega K(\dot{x}, \dot{y}) \varphi_j(\dot{y}) d\mu(\dot{y})$$
$$= \sum_{k=0}^\infty a_k(\dot{x}) \int_\Omega \underbrace{\varphi_k(\dot{y}) \varphi_j(\dot{y}) d\mu(\dot{y})}_{\delta_{kj}}$$
$$= a_j(\dot{x})$$
$$\Rightarrow K(\dot{x}, \dot{y}) = \sum_{j=0}^\infty \lambda_j \varphi_j(\dot{x}) \varphi_j(\dot{y}) = \langle \Xi(\dot{x}), \Xi(\dot{y}) \rangle_{\mathcal{H}_i}$$

where $\quad \Xi(\dot{x}) = \langle \sqrt{\lambda_j} \, \varphi_j(\dot{x}), \quad j = 0, 1, 2, \dots \rangle$
<span style="color:blue">$\varphi_j(\dot{x})$</span>

So, $\mathcal{H} = \ell_2$ with the ONB $\{\psi_j\}$. $\qquad\qquad \Box$

<u>Def</u> Given kernel $K$, the <span style="color:red">reproducing kernel Hilbert space (RKHS)</span> is
$$\mathcal{H}_K = T_K^{\frac{1}{2}} L^2(\Omega, \mu) = \{ \sum_{j=0}^\infty a_j \sqrt{\lambda_j} \, \varphi_j \mid a \in \ell_2 \}$$
$$\Rightarrow \langle f, g \rangle_{\mathcal{H}_K} = \langle T^{-1} f, g \rangle_{\ell^2} = \langle T^{-\frac{1}{2}} f, T^{-\frac{1}{2}} g \rangle_{\ell^2}$$

<u>Properties of RKHS:</u>

<span style="color:blue">$K \approx$ RKHS</span>

① $K(\dot{x}, \cdot) \in \mathcal{H}_K$ s.t.
$$\| K(\dot{x}, \cdot) \|_{\mathcal{H}_K}^2 = \langle \sum_{j=0}^\infty \lambda_j \varphi_j(\dot{x}) \varphi_j(\cdot), \sum_{k=0}^\infty \lambda_k \varphi_k(\dot{x}) \varphi_k(\cdot) \rangle_{\mathcal{H}_K}$$
$$= \sum_{j,k=0}^\infty \varphi_j(\dot{x}) \varphi_k(\dot{x}) \lambda_j \lambda_k \langle \varphi_j, \varphi_k \rangle_{\mathcal{H}_K}$$

$$= \sum_{j,k=0}^{\infty} \psi_j(\tilde{x}) \psi_k(\tilde{x}) \lambda_j \lambda_k \lambda_j^{-\frac{1}{2}} \lambda_k^{-\frac{1}{2}} \langle \psi_j, \psi_k \rangle_C$$

$$= \sum_{j=0}^{\infty} \lambda_j \psi_j(\tilde{x})^2 = K(\tilde{x}, \tilde{x})$$

② $\quad \forall f \in \mathcal{H}_K, \quad \langle f(\cdot), K(\tilde{x}, \cdot) \rangle_{\mathcal{H}_K} = \langle f(\cdot), \mathbb{I}(\tilde{x}) \rangle_{\mathcal{H}_K} \equiv f(\tilde{x})$

"reproducing" property

the equates point evaluation to $\langle \cdot, \cdot \rangle_{\mathcal{H}_K}$

So, $\quad f \mapsto f(\tilde{x}) \quad$ is $\quad$ **bounded** $\quad \left( \begin{array}{l} \text{linear functionals in Hilbert} \\ \text{spaces are bounded} \end{array} \right)$
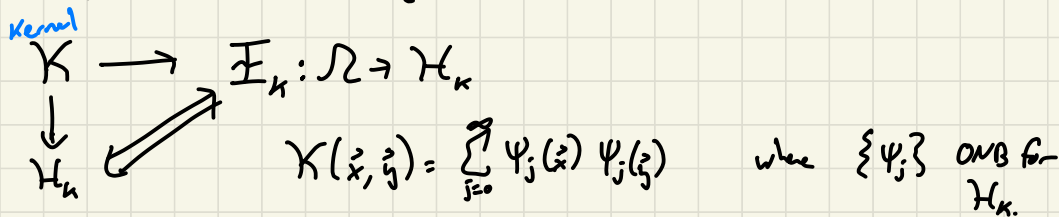
Note: you are an RKHS if and only if point evaluation is bounded.

③ $\quad \langle K(\tilde{x}, \cdot), K(\tilde{y}, \cdot) \rangle_{\mathcal{H}_K} = K(\tilde{x}, \tilde{y})$

④ $\mathcal{H}_K$ is the closure of $\left\{ \sum_{j=1}^{N} a_j K(\tilde{x}_j, \cdot) \right\}$
with respect to
$$\langle K(\tilde{x}, \cdot), K(\tilde{x}', \cdot) \rangle_{\mathcal{H}_K} = K(\tilde{x}, \tilde{x}')$$

This means that we can describe $\mathcal{H}_K$ via a dataset and function evaluations $\{ f(\tilde{x}_j) \}$.

To recap, we saw an equivalence

Kernel
$$K \longrightarrow \mathbb{I}_K : \Omega \to \mathcal{H}_K$$
$$\downarrow \qquad \nearrow$$
$$\mathcal{H}_K \qquad \qquad K(\tilde{x}, \tilde{y}) = \sum_{j=0}^{\infty} \psi_j(\tilde{x}) \psi_j(\tilde{y}) \qquad \text{where } \{\psi_j\} \text{ ONB for } \mathcal{H}_K.$$

## ML Applications

Given $\mathbb{I} = (\psi_0, \psi_1, \dots)$ with $\psi_j : \Omega \to \mathbb{R}$, we wish to find the function
$$f(\tilde{x}; \theta) = \sum_{j=0}^{\infty} \theta_j \psi_j(\tilde{x}) = \langle \theta, \mathbb{I}(\tilde{x}) \rangle \quad \text{that minimizes}$$
$$\sum_{i=1}^{n} \ell(f(\tilde{x}_i, \theta), y_i) + \frac{\lambda}{2} \| \theta \|_2^2$$


$\theta$ 
$\mathbb{I}(\tilde{x})$ 
$\mathbb{I}(\Omega)$ 
orthogonal

**Option 1:** $\ell(a,b) = \frac{1}{2}(a-b)^2 \Rightarrow \mathcal{L}_\lambda(\theta) = \frac{1}{2}\|Y - \mathcal{I}^T\theta\|^2 + \frac{\lambda}{2}\|\theta\|^2$

Yongwei method

We have
$$\vec{\nabla}_\theta \mathcal{L}_\lambda = -\mathcal{I}(Y - \mathcal{I}^T\theta) + \lambda\theta$$
and so
$$\vec{\nabla}_\theta \mathcal{L}_\lambda = 0 \iff \theta = (\mathcal{I}\mathcal{I}^T + \lambda I)^{-1}\mathcal{I}Y$$

$\underbrace{\qquad\qquad}_{\text{shitty to invert,}} \in \mathbb{R}^{\text{nfeatures} \times \text{nfeatures}}$

**Option 2:** Let's write $X(\vec{x}, \vec{y}) = \langle \mathcal{I}(\vec{x}), \mathcal{I}(\vec{y})\rangle_{\ell_2}$ and deal with

Kernel method things in $\mathcal{H}_K = \text{span}\{\mathcal{I}\}$. (So, $\mathcal{I}$ is ONB for $\mathcal{H}_K$).

We have
$$f(\vec{x};\theta) = \sum_{j=0}^\infty \theta_j \psi_j(\vec{x}) \in \mathcal{H}_K, \qquad \|\theta\|_\lambda^2 = \|f\|_{\mathcal{H}_K}^2$$

$$\Rightarrow f_* = \underset{f \in \mathcal{H}_K}{\text{argmin}} \sum_{i=1}^m \ell(f(\vec{x}_i), y_i) + \frac{\lambda}{2}\|f\|_{\mathcal{H}_K}^2$$

$\underbrace{\qquad\qquad\qquad}_{\substack{\text{only depends on} \\ \{f(\vec{x}_i)\} = \{\langle X(\vec{x}_i, \cdot), f\rangle_{\mathcal{H}_K}\}_{i=1}^m}}$

Let us consider the (finite-dim) subspace of $\mathcal{H}_K$ along the training datasets given by $\quad \Pi_X : \mathcal{H}_K \to \text{span}\{K(\vec{x}_i, \cdot)\}_{i=1}^m$

and so $\sum_{i=1}^m \ell(f(\vec{x}_i), y_i)$ depends only on $\Pi_X f$.

The minimization problem is:
$$f_* = \underset{f \in \mathcal{H}_K}{\text{argmin}} \ \mathcal{L}(\Pi_X f) + \frac{\lambda}{2}\|f\|_{\mathcal{H}_K}^2$$

However, $\|f\|_{\mathcal{H}_K}^2 = \|\Pi_X f\|_{\mathcal{H}_K}^2 + \|\Pi_X^\perp f\|_{\mathcal{H}_K}^2$.

Since $\mathcal{L}$ doesn't see $\Pi_X^\perp f$ (it only sees function eval. at data points),
$$f_* = \underset{f \in \text{span}\{K(\vec{x}_i, \cdot)\}_{i=1}^m}{\text{argmin}} \sum_{i=1}^m \ell(f(\vec{x}_i), y_i) + \frac{\lambda}{2}\|f\|_{\mathcal{H}_K}^2$$

$\underbrace{\qquad\qquad\qquad}_{\Pi_X(\mathcal{H}_K)}$

We parametrize $f(\vec{x}) = \sum_{j=1}^m a_j K(\vec{x}_j, \cdot)$ and solve for $\ell(a,b) = \frac{1}{2}(a-b)^2$

$$\Rightarrow f(\vec{x}_i) = \langle k(\vec{x}_i, \cdot), \sum_{j=1}^{3} a_j K(\vec{x}_j, \cdot) \rangle_{\mathcal{H}_k} = \sum_{j=1}^{3} a_j K(\vec{x}_i, \vec{x}_j) = K\vec{a}$$

$$K: \begin{bmatrix} K(\vec{x}_j, \vec{x}_i) & K(\vec{x}_j, \vec{x}_i) \\ \vdots & \ddots \end{bmatrix}$$

$$\Rightarrow \|f\|^2_{\mathcal{H}_k} = \langle \sum_{j=1}^{3} a_j K(\vec{x}_j, \cdot), \sum_{i=1}^{3} a_i K(\vec{x}_i, \cdot) \rangle_{\mathcal{H}_k}$$

$$= \sum_{i,j=1}^{3} a_i a_j K(\vec{x}_i, \vec{x}_j) = \vec{a}^T K \vec{a}$$

$$\Rightarrow \vec{a}_* = \underset{\vec{a}}{\text{argmin}} \; \frac{1}{2} \|Y - K\vec{a}\|_2^2 + \frac{\lambda}{2} \vec{a}^T K \vec{a}$$

$$\underbrace{\phantom{\frac{\lambda}{2} \vec{a}^T K \vec{a}}}_{\|\vec{a}\|^2_{K^{-1}}}$$

So, $\quad \nabla_a = -K(Y - K\vec{a}) + \lambda K\vec{a} = 0 \quad \Leftrightarrow \quad \vec{a}_* = (k + \lambda I)^{-1} Y$

$$\Leftrightarrow f_* = K\vec{a}_* = k(K + \lambda I)^{-1} Y$$

$$\underbrace{\phantom{k(K + \lambda I)^{-1}}}_{\in \mathbb{R}^{\# data \, \times \, \# data}}$$

<u>To sum, Kernel methods for a given Kernel $K$ yield:</u>

* $\mathbb{F}_K$ - feature map

* $\mathcal{H}_K$ - RKHS (Representer Theorem)

* $f_K$ - Gaussian Process on $\Omega$ with
$$\mathbb{E}\{f_K(\vec{x})\} = 0,$$
$$Cov(f_K(\vec{x}), f_K(\vec{\xi})) = K(\vec{x}, \vec{\xi})$$

* DPP $X_K$ on $\Omega$

# Lecture 11/2 - Quadratic Models

Last time- We considered linear models

$$z(x;\theta) = \Phi^T(x)\theta = \sum_{j=0}^{\hat{n}} \theta_j \psi_j(x)$$



$col \Phi$   $\theta_{min}$   $ker \Phi$

$\theta \in \mathbb{R}^n$
$\{\Phi(x_i)^T\theta = y_i\}$

All solutions to $\quad \mathcal{L}(\theta) = \frac{1}{2}\sum_{i=1}^{\hat{n}}(z(x_i;\theta) - y_i)^2 = 0$

are solutions of $\quad \Phi\Phi^T\theta = \Phi Y$

Today we study quadratic models $\qquad \Phi(x)$ symmetric

$$\boxed{z(x;\theta) = \Phi(x)^T\theta + \frac{\varepsilon}{2}\theta^T \Phi(x)\theta} = \sum_{j=0}^{n}\theta_j\psi_j(x) + \frac{\varepsilon}{2}\sum_{j_1,j_2=0}^{\hat{n}}\theta_{j_1}\theta_{j_2}\psi_{j_1,j_2}(x)$$
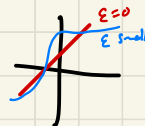
We motivate this via Taylor expansion

$$f(x;\theta) = f(x;0) + \nabla_\theta f(x;0)^T\theta + \frac{1}{2}\theta^T H_\theta f(x;0)\theta + \dots$$

With the same loss $\qquad \mathcal{L}_A(\theta) = \sum_{\alpha\in A}\frac{1}{2}(y_\alpha - z(x_\alpha;\theta))^2$

and the goal to find minima of $\mathcal{L}_A(\theta)$ to $1^{st}$ order in $\varepsilon$.

Notation: We define $\quad \nabla_\theta z(x;\theta) \equiv \Phi^\varepsilon(x;\theta) = \Phi(x) + \varepsilon\Phi(\theta)$



$\varepsilon = 0$
$\varepsilon$ small

To solve $\quad \nabla_\theta \mathcal{L}_A(\theta) = 0$ to first order in $\varepsilon$, we have

$$\nabla_\theta \mathcal{L}_A(\theta) = \sum_{\alpha\in A}\Phi^\varepsilon(x_\alpha;\theta)(z(x_\alpha;\theta) - y_\alpha)$$

$$(1) = \sum_{\alpha\in A}(\Phi(x_\alpha) + \varepsilon\Phi(x_\alpha)\theta)\times(\Phi(x_\alpha)^T\theta + \frac{\varepsilon}{2}\theta^T\Phi(x_\alpha)\theta - y_\alpha)$$

Let's write
$$\theta_* = \overset{\text{free}}{\theta^F} + \varepsilon \overset{\text{interesting}}{\theta^I} + O(\varepsilon^2), \quad \text{where} \quad \mathbb{I}\mathbb{I}^T \theta^F = \mathbb{I} Y$$

So,
$$0 = \nabla_\theta \ell_A(\theta) \quad \text{gives}$$

$$0 = \sum_{\alpha \in A} \left( \mathbb{I}(x_\alpha) + \varepsilon \mathbb{I}(x_\alpha)(\theta^F + \varepsilon \theta^I) \right) \cdot \left( \mathbb{I}(x_\alpha)^T (\theta^F + \varepsilon \theta^I) + \frac{\varepsilon}{2}(\theta^F)^T \mathbb{I}(x_\alpha) \theta^F - y_\alpha \right)$$

$\underset{\varepsilon^0}{\overset{\text{zeroth order}}{\text{terms}}} \Rightarrow 0 = \varepsilon^0 \left[ \sum_{\alpha \in A} \overset{\text{rank-1 matrices}}{\mathbb{I}(x_\alpha)\mathbb{I}(x_\alpha)^T} \theta^F - \mathbb{I}(x_\alpha) y_\alpha \right]$ ← this is 0 because $\mathbb{I}\mathbb{I}^T \theta^F = \mathbb{I} Y$

$\underset{\varepsilon^1}{\overset{\text{first order}}{\text{terms}}} \Rightarrow 0 = \varepsilon^1 \left[ \sum_{\alpha \in A} \mathbb{I}(x_\alpha)\theta^F (\mathbb{I}(x_\alpha)^T \theta^F - y_\alpha) + \sum_{\alpha \in A} \frac{1}{2} \mathbb{I}(x_\alpha)(\theta^F)^T \mathbb{I}(x_\alpha)\theta^F \right.$
$\left. + \sum_{\alpha \in A} \mathbb{I}(x_\alpha)\mathbb{I}(x_\alpha)^T \theta^I \right]$

$$\Rightarrow \underset{y_\alpha}{\underbrace{\sum_{\alpha \in A} \mathbb{I}(x_\alpha)\theta^F}} = \overset{\text{equal}}{\underbrace{\sum_{\alpha \in A} \mathbb{I}(x_\alpha)\theta^F}} \overset{y_\alpha \text{ prediction linear model made}}{\mathbb{I}(x_\alpha)^T \theta^F} + \sum_{\alpha \in A} \frac{1}{2}\mathbb{I}(x_\alpha)(\theta^F)^T \mathbb{I}(x_\alpha)\theta^F$$
$$+ \underset{\mathbb{I}\mathbb{I}^T}{\underline{\sum_{\alpha \in A} \mathbb{I}(x_\alpha)\mathbb{I}(x_\alpha)^T \theta^I}}$$

$$\Rightarrow \boxed{\theta^I = -\underset{\text{pseudo-inverse}}{(\mathbb{I}\mathbb{I}^T)^+} \sum_{\alpha \in A} \frac{1}{2}\mathbb{I}(x_\alpha) \cdot \underset{\text{elementwise mult.}}{(\theta^F)^T} \mathbb{I}(x_\alpha)\theta^F}$$

**Interpretation:**

(1) We can write $(\theta^F)^T \mathbb{I}(x_\alpha)\theta^F = \langle \theta^F, \mathbb{I}(x_\alpha)\theta^F \rangle$

which is $\underset{\mathbb{I} \text{ not pos. def.}}{\underline{\text{almost}}} \quad \|\theta^F\|^2_{\mathbb{I}(x_\alpha)}$

(2) Also, $\theta^I = (\mathbb{I}\mathbb{I}^T)^+ (\mathbb{I} Y^I)$ ← transformed Y via these coefficients

So, if we had changed $Y \mapsto Y + \mathbb{I} Y^I$

deform Y with useful features — feature learning!

and solved least-squares with a linear model, we would get the same predictions

$\iff$

GD on nonlinear models learns label features to run linear model on

(3) Note that $(\Phi \Phi^T)\theta^I = -\sum_{\alpha \in A} \frac{1}{2} \Phi(x_\alpha)(\theta^F)^T \Phi(x_\alpha)\theta^F$

only determines $\theta^I$ on span $\{\Phi(x_\alpha)\}$: so, it is unclear what happens to $\theta^I_\perp$ ← this may depend on optimization method and allow weird things to happen

When we do **gradient flow** (continuous GD),

$$\frac{d}{dt}\theta_t = -3\nabla_\theta \mathcal{L}_A(\theta_t)$$

Recall the effective features $\qquad \overset{\nabla_\theta z(x;\theta)}{\Phi^E(x;\theta) = \Phi(x) + \varepsilon \Phi(x)\theta}$

We can write

$$\frac{d}{dt}\Phi^E(x;\theta_t) = \varepsilon \overset{\in col(\Phi)}{\Phi(x)\frac{d}{dt}\theta_t^F} + O(\varepsilon^2)$$

<u>Interpretations:</u>
- $\Phi^E$ changes!
- $\frac{d}{dt}\theta_t^F \in col(\Phi) \Rightarrow \frac{d}{dt}\Phi(x;\theta) \in span\{\Phi(x)\Phi(x_\alpha), \alpha \in A\}$

Moreover, since $\theta_t^F$ solves the linear model,

$$\frac{d}{dt}\theta_t^F = \frac{d}{dt}(\theta_t^F - \theta_*) = -3\Phi\Phi^T\underbrace{(\theta_t^F - \theta_*)}_{vector\ of\ residuals}$$

$$\to \theta_t^F - \theta_* = e^{-3t\Phi\Phi^T}(\theta_0^F - \theta_*)$$

Therefore,

$$\frac{d}{dt}\Phi^E(x;\theta_t) = \varepsilon \Phi(x)(-3\Phi\Phi^T)(\theta_t^F - \theta_*)$$
$$= \varepsilon \Phi(x)(-3\Phi\Phi^T)e^{-3t\Phi\Phi^T}(\theta_0^F - \theta_*)$$

$$\Rightarrow \Phi^E(x;\theta_t) - \Phi^E(x;\theta_0) = \varepsilon \Phi(x)(I - e^{-3t\Phi\Phi^T})*(\theta_0^F - \theta_*)$$

At $t \to \infty$,

$$\boxed{\Phi^E(x;\theta_\infty) = \Phi^E(x;\theta_0) + \varepsilon \Phi(x)(\theta_0^F - \theta_*)}$$

<u>To recap:</u>

① We got the NTK $\Phi^\varepsilon (\Phi^\varepsilon)^T$ @ all times

② Formula for what happens to $\theta_*$ to leading order in $\varepsilon$
on span$\{\Phi\}$
$\quad\curvearrowright$ what happens to $\theta_\perp$?

Next,
$$\underline{\text{grad flow}}$$
$$\frac{d}{dt}\theta_t = -3 \nabla_\theta \mathcal{L}_A(\theta_t) = -3 \sum_{\alpha \in A} \Phi^\varepsilon(x_\alpha; \theta_t) \times (z(x_\alpha; \theta_t) - y_\alpha)$$

Therefore,
$$\frac{d}{dt}\theta_t = -3 \sum_{\alpha \in A} \left( \Phi(x_\alpha) + \varepsilon E(x_\alpha)\left( I - e^{-3t\,\Phi\Phi^T} \right)(\theta_0^F - \theta_*) \right)$$
$$\times \left( \Phi(x_\alpha)^T \theta_t + \frac{\varepsilon}{2}\theta_t^T E(x_\alpha)\theta_t - y_\alpha \right)$$

$\frac{d}{dt}\theta_\infty^F$ cancels
$\downarrow$ w/ $\Phi\Phi^T \theta_0 \cdot \Phi_y$
$$\Rightarrow \frac{d}{dt}\theta_t^I = -3 \sum_{\alpha \in A} \Phi(x_\alpha)\Phi(x_\alpha)^T \theta_t^I + E(x_\alpha)\left( I - e^{-3t\,\Phi\Phi^T}\right)(\theta_0^F - \theta_*)\left(\Phi(x_\alpha)^T\theta_t^F - y_\alpha \right)$$
$$+ \frac{1}{2}\Phi(x_\alpha)\left(\theta_t^F\right)^T E(x_\alpha)\theta_t^F$$

Projecting onto the orthogonal complement of span$\{\Phi(x_\alpha)\}$
$$\Rightarrow \frac{d}{dt}\theta_{t\perp}^I = -3 \sum_{\alpha \in A} E^\perp(x_\alpha)\left( I - e^{-3t\,\Phi\Phi^T}\right)(\theta_0^F - \theta_*)\left(\Phi(x_\alpha)^T\theta_t^F - y_\alpha\right)$$

# Lecture 11/7

Recall last time: <span style="color:red">gradient flow</span> on <span style="color:red">quadratic models</span>

(#) $\frac{d}{dt}\theta(t) = -\nabla_\theta \mathcal{L}(\theta(t))$ $\qquad z(x;\theta) = \Phi^T(x)\theta + \frac{\varepsilon}{2}\theta^T \Xi(x)\theta$

where $\mathcal{L}_A(\theta) = \sum_{\alpha \in A} \frac{1}{2}(z(x_\alpha;\theta) - y_\alpha)^2$.

Note that gradient flow (#) is the limit $\zeta \to 0$ of
gradient descent $\qquad$ (##) $\theta(t+1) = \theta(t) - \zeta \nabla_\theta \mathcal{L}(\theta(t))$

We have seen that $\zeta$ small vs. $\zeta$ large can make qualitative
differences.

<u>Today</u>: We consider "large" $\zeta$ in quadratic approximations to 1-layer
ReLU nets:
$$z(x;\theta) = \sum_{i=1}^m \frac{v_i}{\sqrt{m}}\,\Theta\!\left(\frac{u_i^T}{\sqrt{d}}x\right),$$
<span style="color:blue">$\Theta$: ReLU
$\hat{u} \in \mathbb{R}^d$ input
$\hat{v} \in \mathbb{R}^{d \cdot m}$ weights</span>

Writing the quadratic approx. <span style="color:red">←evaluated at t=0</span>
$$z(x;\theta) \approx z^0 + \sum_{i=1}^m \nabla_{u_i}^0(u_i - u_i(0)) + \partial_{v_i}^0(v_i - v_i(0)) + (u_i - u_i(0))^T \widetilde{H}_i^0 (v_i - v_i(0))$$
<span style="color:blue">$\in \mathbb{R}^d$ $\qquad$ $\in \mathbb{R}$ $\qquad$ $\in \mathbb{R}^d$ $\qquad$ $\in \mathbb{R}$</span> $\qquad$ <span style="color:blue">$\in \mathbb{R}^{d \times 1}$</span>

where $\quad z^0 = z(x;\theta(0)), \quad \nabla_{u_i}^0 = \nabla_{u_i}z(x;\theta(0)), \quad \partial_{v_i}^0 = \partial_{v_i}z(x;\theta(0))$

$\qquad H_i^0 = \partial_{v_i}\nabla_{u_i}z(x;\theta(0))$ ] <span style="color:blue">no second derivatives $\partial_{u_i}^2$ or $\partial_{v_i}^2$ because ReLU assumption!</span>

<u>Goal</u>: Following Zhu et. al, we consider <u>one</u> training datapoint $(x,y)$
and show that the "catapult phase" occurs.
Explicitly,
$$\lambda(u,v) = \|\nabla_\theta z(x;\theta)\|^2 = \sum_{i=1}^m \|\nabla_{u_i}z(x;\theta)\|^2 + (\partial_{v_i}z(x;\theta))^2$$
<span style="color:blue">"NTK"
is $1 \times 1$</span> and $\lambda(t) = \lambda(u(t),v(t)), \quad \mathcal{L}(t) = \mathcal{L}(u(t),v(t))$

We have the following "phase diagram" for optimization:

"Thm" [Zhu] When $n \gg 1$

$0 < \eta < \frac{2}{\lambda(0)}$ : optimization "looks linear" in the sense that
$$\mathcal{L}(t) \asymp c(1-\varepsilon)^t, \quad \lambda(t) \approx \lambda(0)$$

$\frac{2}{\lambda(0)} < \eta < \frac{4}{\lambda(0)}$ : "catapult phase"

    loss grows exponentially   if $t \in [0, T_1)$: $\overset{\sim \log n}{}$ $\mathcal{L}(t) \asymp c(1+\varepsilon)^t, \quad \lambda(t) \approx \lambda(0)$   this looks like flatten NTK

    loss settles   if $t \in [T_1, T_2]$: $\overset{\sim \log n}{}$ $\mathcal{L}(t) = \Theta(n)$ plateaus, $\quad \lambda(t+1) < \lambda(t)$

    loss shrinks exponentially   if $t \in [T_2, \infty)$: $\mathcal{L}(t) \asymp (1-\varepsilon)^t, \quad \lambda(t) \to \lambda(\infty)$ small

$\frac{4}{\lambda(0)} < \eta$ : optimization diverges $\mathcal{L}(t) \asymp (1+\varepsilon)^t \; \forall t$

Interpretation of catapult phase:

* $\mathcal{L}(t) = (1+\varepsilon)^t \Rightarrow \theta(t)$ leaves the region around $\theta(0)$
* $\lambda(t+1) < \lambda(t) \Rightarrow$ find a "flat part" of parameter space. Since $\mathcal{H}_\theta \mathcal{L}$ and the NTK $(\lambda)$ are isospectral, $\mathcal{H}_\theta \mathcal{L} = \nabla_\theta z (\nabla_\theta z)^T$ has the same nonzero eigenvalues as $\lambda = (\nabla_\theta z)^T \nabla_\theta z$
  $\Rightarrow$ max eigenvalue keeps decreasing

The **Key step** is to derive a closed set of equations for two "order parameters", which are

residual → $\left\{ z(t+1) - y = f\left(z(t) - y, \lambda(t)\right) \right.$    coupled recursion of two parameters

NTK → $\left. \lambda(t+1) = f\left(z(t) - y, \lambda(t)\right) \right\}$

Let me read the boxed proposition.

**Prop. ☆**

↑ we will prove this at the end

$$z(t+1) - y = (z(t) - y) \overset{\text{linear part}}{\left[1 - \eta\lambda(t) + \frac{\|x\|^2}{nd}\eta^2 z(t)(z(t)-y)\right]} \quad (1)$$

$$\lambda(t+1) = \underset{\text{linear part}}{\lambda(t)} + \eta\frac{\|x\|^2}{nd}(z(t)-y)^2\left[\eta\lambda(t) - \frac{4 z(t)}{z(t)-y}\right] \quad (2)$$

First, though, we will prove the theorem from this proposition.

(1) Since $m \gg 1$, the two quadratic terms above scale like $\sim \frac{1}{m}$
  <u>unless</u> the residuals scale with $z(t) \sim \sqrt{m}$

  (note that we can think of $\epsilon$ from the previous lecture
   to be like $\frac{1}{m}$ (the thing that scales the Hessian)

  $\Rightarrow$ early dynamics (before $z(t)$ gets too big) are always $\approx$ linear
                                                    or converge

(2) So, if $\quad 3 < \frac{2}{\lambda(0)}, \quad |z(t) - y| \approx C e^{-t} \ll \sqrt{m} \ \forall t$, yielding the first
  "if you look linear and are driven linearly, you behave linearly"                    phase

(3) If $3 > \frac{2}{\lambda(0)}$, we diverge linearly with $|z(T_1)| \approx C e^{T_1} \sim \sqrt{m} \Rightarrow T_1 = O(\log m)$
                                                      $\lambda(t) \propto \lambda(0)$
  Around time $t = T_1$, the recursion in Prop. ☆ yields
  $$\lambda(t+1) \approx \lambda(t) + 3 \frac{\|x\|^2}{d}\left[3 \lambda(0) - 4\right]$$
  So, if $3 < \frac{4}{\lambda(0)}, \quad \lambda(t+1) < \lambda(t)$ decreases and $|1 - 3\lambda(t)|$ gets smaller.

  $\Rightarrow z(t) - y$ stops growing until $|1 - 3\lambda(t)| < 1$, and we re-enter the
  linear regime with $\lambda(t) \to 0$ exponentially.

This yields the result! The residuals and NTK fight each other
                    in the quadratic case.

Now, we prove the recursion.

$\boxed{\text{Proof of Prop. ☆}}$ —

Recall that

$$z(t) = z^0 + \sum_{i=1}^{\hat{n}} \nabla_{u_i}^0 \left(u_i(t) - u_i(0)\right) + \partial_{v_i}^0 \left(v_i(t) - v_i(0)\right) + \left(u_i(t) - u_i(0)\right)^T H_i^0 \left(v_i(t) - v_i(0)\right)$$

Taking a gradient,

$$\nabla_{u_i} z(t) = \nabla_{u_i}^0 + H_i^0\left(v_i(t) - v_i(0)\right)$$

We can also write at

$$\nabla_{u_i}^0 = \nabla_{u_i}\left[\sum_{j=1}^{3} \frac{1}{\sqrt{m}} v_j(0) \sigma\left(u_j(0)^T \frac{x}{\sqrt{d}}\right)\right] = \frac{x}{\sqrt{md}} v_i(0) \, \mathbb{1}_{\{u_i(0)^T x \geq 0\}}$$

Also,

$$\partial_{v_i} z(t) = \partial_{v_i}^0 + (u_i(t) - u_i(0))^T H_i^0$$

and
$$\partial_{v_i}^0 = \partial_{v_i}\left[\sum_{j=1}^{\tilde{z}} \frac{1}{\sqrt{m}} v_j(0)\, \sigma\left(\frac{u_j(0)^T x}{\sqrt{d}}\right)\right] = \frac{u_i(0)^T x}{\sqrt{md}}\, \mathbb{1}_{\{u_i(0)^T x \geq 0\}}$$

The mixed derivative is
$$H_i^0 = \frac{x}{\sqrt{md}}\, \mathbb{1}_{\{u_i(0)^T x \geq 0\}}$$

So, we can compute the residual

$$z(t+1) - y = -y + z^0 + \sum_{i=1}^{\tilde{z}} \nabla_{u_i}^0 (u_i(t+1) - u_i(0)) + \partial_{v_i}^0 (v_i(t+1) - v_i(0))$$

$$+ (u(t+1) - u(0))^T H_i^0 (v_i(t+1) - v_i(0))$$

$$= -y + z^0 + \sum_{i=1}^{\tilde{z}} \nabla_{u_i}^0 \left(u_i(t) - u_i(0) - 3\nabla_{u_i}\mathcal{L}(t)\right) + \partial_{v_i}^0\left(v_i(t) - v_i(0) - 3\,\partial_{v_i}\mathcal{L}(t)\right)$$

$$+ \left((u_i(t) - u_i(0) - 3\nabla_{u_i}\mathcal{L}(t)\right) H_i^0 \left(v_i(t) - v_i(0) - 3\partial_{v_i}\mathcal{L}(t)\right)$$

$$= z(t) - y - 3\left[\sum_{i=1}^{\tilde{z}} \nabla_{u_i}^0\left(\nabla_{u_i}[z(t)](z(t) - y)\right) + \partial_{v_i}^0\left(\partial_{v_i}[z(t)](z(t) - y)\right) \right.$$
$$\left. + (\nabla_{u_i}[z(t)])^T H_i^0 (v_i(t) - v_i(0))(z(t) - y) + (u_i(t) - u_i(0))^T H_i^0 (\partial_{v_i}[z(t)])\right]$$

$$+ 3^2 (z(t) - y)^2 \times (\nabla_{u_i}[z(t)])^T H_i^0 \partial_{v_i}[z(t)]$$

$$= z(t) - y - (z(t) - y) 3\, \lambda(t) + 3^2 (z(t) - y)^2 \times \underbrace{(\nabla_{u_i}[z(t)])^T H_i^0 \partial_{v_i}[z(t)]}_{\color{blue}{\frac{\|x\|^2}{md} z(t)}}$$

$$= (z(t) - y)\left[1 - 3\lambda(t) + \frac{\|x\|^2}{md} 3^2 z(t)(z(t) - y)\right]$$

$\square$

## Open problems for quadratic models!

* $\theta \neq \text{ReLU}$, one datapoint (perhaps $\theta$ s.t. $\theta'$ monotone)

* # data $\geq 2$, $d \geq 2$ (Zhu et. al do $d=1$, #data $=2$)

* $z(T_1) \sim \sqrt{m}$, same as mean field scaling?? <span style="color:red">★ relationship between growing fischer orbits & feature learning/catapult??</span>

* What happens for $\mathcal{I}$, $\mathcal{I}$ random?

* Do cubic models have another "catapult phase"?

# Lecture 11/9 - Implicit Bias

**Recall:** Last time we considered $z(x;\theta) = \mathbb{E}(x)^T \theta + \frac{\varepsilon}{2} \theta^T \mathbb{E}(x)\theta$ for small $\varepsilon$ expansions for GD and GF.

**Today:** [Woodworth]

Consider quadratic models of the form

$$\dot{\theta} = \begin{pmatrix} \theta^+ \\ \theta^- \end{pmatrix} \in \mathbb{R}^{2d} \quad \text{where} \quad z(x;\theta) = \langle \beta_\theta ; x \rangle = \langle \underset{\uparrow}{\theta_+^{\circ 2}} - \theta_-^{\circ 2}, x \rangle$$

*elementwise square*

$$= 0^T \cdot \theta + \frac{2}{2} \theta^T \begin{pmatrix} \text{Diag}(x) & 0 \\ 0 & \text{Diag}(x) \end{pmatrix} \theta$$

$$\mathbb{E} = 0, \quad \mathbb{E} = \begin{pmatrix} \text{Diag } x & 0 \\ 0 & \text{Diag } x \end{pmatrix}, \quad \varepsilon = 2$$

The **motivation** for this is that we are able to express all linear functions with a nonlinear parameterization.

We train by gradient flow (GF)

$$\frac{d}{dt}\theta(t) = -\nabla_\theta \mathcal{L}(\theta(t)), \quad \mathcal{L}(\theta(t)) = \sum_{n=1}^N \frac{1}{2}(z(x_n;\theta) - y_n)^2$$

$$\theta(0) = \alpha \begin{pmatrix} \theta_o \\ \theta_o \end{pmatrix} \Rightarrow \beta_\theta(0) = \theta_+^{\circ 2}(0) - \theta_-^{\circ 2}(0) = 0$$



$\mathbb{R}^d$

span $\{x_n\}_{n=1}^3$

GF on linear model

Which does model avoid go?

$\beta_\theta(0)$

$I = \{\beta \mid \langle x_n, \beta \rangle = y_n \; \forall n\}$
- set of all interpolants.
- points that perfectly fit the training data

**The question:** As a function of "scale" $\alpha$, "shape" $\theta_o$, which minimum on $I$ does GF find?

$$\begin{pmatrix} \text{Mean-field is} & \alpha \to 0 \\ \text{NTK parameterization is} & \alpha \to \infty \end{pmatrix} \text{"Implicit bias"}$$

**Theorem:** If GF with some initialization converges to a minimum loss of $L$, the minimum is given by

$$\beta^{*}_{\alpha,\theta_{0}} = \underset{\beta \in \mathbb{R}^{d}}{\arg\min}\ Q_{\alpha,\theta_{0}}(\beta) \quad \text{subject to} \quad X^{T}\beta = y.$$

"implicit" bias

$Q_{\alpha\theta_{0}}$ is strictly convex:

$$Q_{\alpha\theta_{0}}(\beta) = \sum_{i=1}^{d} \alpha^{2}\theta_{0,i}^{2}\ q\left(\frac{\beta_{i}}{\alpha^{2}\theta_{0,i}^{2}}\right)$$

where

$$q(z) = 2 - \sqrt{4+z^{2}} + z\ \text{arcsinh}\left(\frac{z}{2}\right)$$

Reduce nonlinear optimization to linear optimization with explicit penalty!

**Interpretation:**

(1) This says that we have implicit regularization $Q_{\alpha,\theta_{0}}$ st. GF returns

$$\underset{\beta}{\arg\min}\left\{ \sum_{n=1}^{N} (\langle\beta,x_{n}\rangle - y_{n})^{2} + \lambda\, Q_{\alpha,\theta_{0}}(\beta) \right\} \quad \text{with } \lambda \to 0$$

minimize loss firstly then regularizer

(2) $\alpha \to 0$  causes  $Q_{\alpha,\theta_{0}}(\beta) \to \|\beta\|_{1}$  implicit $L_{1}$ regularization
"feature selection"

(3) $\alpha \to \infty$  causes  $\alpha^{2}Q_{\alpha,\theta_{0}}(\beta) \to \frac{1}{4}\sum_{i=1}^{d}\frac{\beta_{i}^{2}}{\theta_{i,0}^{2}}$  implicit weighted $L_{2}$ regularization

(4) For $\alpha \in (0,\infty)$, $Q$ somehow interpolates between the two.

**Proof:**

write diff eq

**Lemma 1:** We have $\frac{d}{dt}\theta(t) = -2\left((X,-X)^{T}\vec{r}(t)\right)\odot\theta(t)$

$\in \mathbb{R}^{2d\times N}$   $\in \mathbb{R}^{N}$  element-wise product  $\in \mathbb{R}^{2d}$

where $\vec{r}(t) = \begin{pmatrix} r_{1}(t) \\ \vdots \\ r_{N}(t) \end{pmatrix}$,   $r_{n}(t) = \langle\beta(t),x_{n}\rangle - y_{n}$   residuals

solve diff eq

**Lemma 2:** We solve

$$\beta_{\alpha,\theta_{0}}(\infty) = 2\alpha^{2}\theta_{0}^{\odot 2}\odot\sinh\left(-4X^{T}\int_{0}^{\infty}\vec{r}(s)\,ds\right)$$

$\in\mathbb{R}^{d}$    $\in\mathbb{R}^{d\times N}$  $\in\mathbb{R}^{N}$

conclude **Lemma 3:** Show Lemma 2 $\Rightarrow Q_{\alpha, \theta_o}(\beta)$ as stated.

## Proof of Lemma 1:

Fix $w = (w_1, \ldots, w_d)$. We can see $\partial_{w_k} \langle w^{\odot 2}, x \rangle = 2 x_k w_k$

$$\Rightarrow \vec{\nabla}_w \langle \vec{w}^{\odot 2}, \vec{x} \rangle = 2 \vec{x} \odot \vec{w}$$

Thus,

$$\frac{d}{dt} \theta_\pm (t) = -\vec{\nabla}_{\theta_\pm} \mathcal{L}(\theta(t)) = -\sum_{n=1}^{N} r_n(t) \vec{\nabla}_{\theta_\pm} \langle \beta, \vec{x}_n \rangle = -\sum_{n=1}^{N} r_n(t) \vec{\nabla}_{\theta_\pm} \langle \theta_+^{\odot 2} - \theta_-^{\odot 2}, \vec{x}_n \rangle$$

$$= -\sum_{n=1}^{N} r_n(t) \left( \pm 2 \vec{x}_n \odot \theta_\pm(t) \right) = \left( \mp \sum_{n=1}^{N} 2 \vec{x}_n r_n(t) \right) \odot \theta_\pm(t)$$

$$= \left( -2 (\pm x)^T r(t) \right) \odot \theta_\pm(t)$$

So, $$\frac{d}{dt} \theta(t) = \begin{pmatrix} \frac{d}{dt} \theta_+(t) \\ \frac{d}{dt} \theta_-(t) \end{pmatrix} = \begin{pmatrix} -2 x^T \vec{r}(t) \odot \theta_+(t) \\ 2 x^T \vec{r}(t) \odot \theta_-(t) \end{pmatrix} = -2 \left( (x, -x)^T \vec{r}(t) \right) \odot \theta(t)$$

function of $\theta$

concatenation $\square$

## Proof of Lemma 2:

We have    $\beta(t) = \theta_+^{\odot 2}(t) - \theta_-^{\odot 2}(t)$

By Lemma 1,

$$\theta(t) = \theta(0) \odot e^{-2(x, -x)^T \int_0^t \vec{r}(s) ds}$$

$$\Rightarrow \theta_\pm(t) = \alpha \theta_0 \odot e^{\mp 2 x^T \int_0^t \vec{r}(s) ds}$$

$$\Rightarrow \beta(t) = \alpha^2 \theta_0^{\odot 2} \odot \left( e^{-4 x^T \int_0^t \vec{r}(s) ds} - e^{4 x^T \int_0^t \vec{r}(s) ds} \right)$$

$$= 2 \alpha^2 \theta_0^{\odot 2} \odot \sinh\left( -4 x^T \int_0^t \vec{r}(s) ds \right) \qquad \square$$

**Note:** $-4 x^T \int_0^\infty \vec{r}(s) ds \in \text{col}(x^T)$   is some vector in data span

This moves us orthogonally to the interpolant hyperplane $\mathcal{I}$, acting as a Lagrange Multiplier.

## Proof of lemma 3:

Suppose $\beta^*_{\alpha,\theta_0} \equiv \beta_{\alpha;\theta_0}(\infty)$ is a global minimum of $\mathcal{L}$.

$$\Rightarrow \langle \beta^*_{\alpha,\theta_0}, \vec{x}_n \rangle = y_n \quad \forall n$$

Let's write $f_{\alpha,\theta_0}(\beta) \equiv 2\alpha^2 \theta_0^{\theta_2} \odot \sinh(\beta)$ for notation

The KKT conditions (optimality for Lagrange multipliers) for

$$\beta^* = \arg\min_\beta Q_{\alpha,\theta_0}(\beta) \text{ s.t. } X\beta = Y$$

are $\quad X\beta^* = Y \quad$ and $\quad \exists v$ s.t. $\vec{\nabla}_\beta Q_{\alpha,\theta_0}(\beta^*) = X^T v$

<span style="color:blue">grad of constraints lies in column space of $X$ (i.e. is orthogonal to independent lies)</span>

But if we have

$$\vec{\nabla}_\beta Q_{\alpha,\theta_0}\left(f_{\alpha,\theta_0}(x^T v)\right) = X^T v$$

then KKT constraints are satisfied. So, we want

$$\left(\vec{\nabla}_\beta Q_{\alpha,\theta_0}\right) \circ f_{\alpha,\theta_0} = \text{Identity} \Leftrightarrow \vec{\nabla}_\beta Q_{\alpha,\theta_0}(\beta) = f_{\alpha,\theta_0}^{-1}(\beta) \Leftrightarrow Q_{\alpha,\theta_0}(\beta) = \vec{\nabla}^{-1}(f_{\alpha,\theta_0}^{-1}(\beta))$$

<span style="color:blue">$\nabla Q$ evaluated at $\beta$</span>

Since $f_{\alpha,\theta_0} = \alpha^2 \theta_0^{\theta_2} \odot \sinh(\beta)$, we find

$$Q_{\alpha,\theta_0}(\beta) = \sum_{i=1}^d \alpha^2 \theta_{0i}^2 \, q\left(\frac{\beta_i}{\alpha^2 \theta_{0i}^2}\right) \quad \text{where} \quad q(z) = 2 - \sqrt{4+z^2} + z \,\text{arcsinh}\left(\frac{z}{2}\right)$$

□

---

## Open problems for quadratic models!

---

* Implicit bias of general quadratic model?
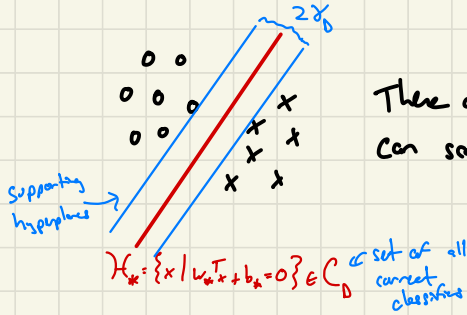  <span style="color:blue">(how does optimizer & $\Sigma, \Xi$ interact)</span>
  <span style="color:blue">- $\Sigma = 0$, $\Xi$ general $\quad$ or $\quad \Sigma = 0$, $\Xi(x)$ simultaneously diagonalizable $\quad$ or $\quad \Xi$ expansion in $\varepsilon$</span>
* Catapult phase for general quadratic models
* Convergence of gradient flow

# Lecture 11/14- Implicit Bias II

Consider a dataset $D = \{(x_i, y_i)\}_{i=1}^{N}$, $x_i \in \mathbb{R}^d$, $y_i \in \{\pm 1\}$  ← same sign
that is "linearly separable". i.e. $\exists b_* \in \mathbb{R}$, $w_* \in \mathbb{R}^d$ s.t. $y_i(\hat{w}_*^T \hat{x}_i + b_*) > 0$



There are $\infty$ many classifiers, since we can scale $w_*$ and $b_*$ to get the same classifier.

$\mathcal{H}_* : \{x \mid w_*^T x + b_* = 0\} \in C_D$  ← set of all correct classifiers

## Goal: Find implicit bias of GF

$$\begin{cases} \frac{d}{dt} w(t) = -\vec{\nabla}\mathcal{L}(w(t)) \\ \mathcal{L}(w(t)) = \sum_{i=1}^{N} \ell(y_i w^T x_i) \end{cases} \quad \text{where } \ell(u) = e^{-u}, \log(1 + e^{-u}), \ldots$$

← turn off the bias

## Margins, Support Vectors

Given a classifier $x \rightarrow y(x) = \text{sgn}(w^T x + b)$ with $(w, b) \in C_D$,
the **margin** is
$$\gamma_{(x_i, y_i)}(w, b) = \text{"margin on } (x_i, y_i)\text{"}$$



$$= \frac{y_i(w^T x_i + b)}{\|w\|} \overset{\text{because } (w,b) \in C_D}{=} \frac{|w^T x_i + b|}{\|w\|}$$

$$= \text{dist}(x_i, \text{decision boundary})$$

We define the **margin on the dataset**
by $\gamma_D(w, b) = \min_{(x_i, y_i) \in D} \gamma_{(x_i, y_i)}(w, b)$

We define the **max-margin classifier** $\hat{w}, \hat{b}$ as a
classifier that maximizes $\max_{(w,b)} \gamma_D(w, b)$ (#)

Note that $\forall (w, b) \in C_D$, $\dfrac{y_i (w^T x_i + b)}{\|w\|}$ is invariant to

the transformation $(w, b) \to C(w, b)$ for some $c > 0$. So, $\forall (w, b) \in C_D$
we can find $(\tilde{w}, \tilde{b}) \in C_D$ s.t.

$$\gamma_D (w, b) = \gamma_D (\tilde{w}, \tilde{b}) \quad \text{and} \quad \min_i \; y_i (\tilde{w}^T x_i + \tilde{b}) = 1 \iff y_i (\tilde{w}^T x_i + \tilde{b}) \geq 1 \; \forall i$$

So, (#) can be viewed with the numerator of $\gamma_D$ as a constraint
in the form
$$\max_w \; \frac{1}{\|w\|} \quad \text{s.t.} \quad y_i (w^T x_i + b) \geq 1 \; \forall i$$

$$\iff \boxed{\min_w \; \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad y_i (w^T x_i + b) \geq 1 \; \forall i}$$

max-margin
classifier objective (##)

Since this is a convex objective over convex region, we can find
a dual problem

dual variable
in $\mathbb{R}^{\# \text{ constraints}} = \mathbb{R}^d$

$$\mathcal{L}(\tilde{w}, b, \tilde{a}) = \frac{1}{2} \|\tilde{w}\|^2 - \sum_{i=1}^n a_i (y_i (w^T x_i + b) - 1)$$

So, solutions to (##) must have

$\vec{\nabla}\mathcal{L} = 0$,      $y_i (w^T x_i + b) - 1 \geq 0$      $a_i \geq 0$      $a_i (y_i (w^T x_i + b) - 1) = 0$

(stationary    primal feasibility    dual feasibility    boundary constraints
point)                                                     (we are tight and on the
                                                            boundary at either primal or dual)

$\Rightarrow 0 = \vec{\nabla}_w \mathcal{L} = w - \sum_{i=1}^n a_i y_i x_i$,    $0 = \vec{\nabla}_b \mathcal{L} = \sum_{i=1}^n a_i y_i$

The boundary constraint gives $\forall i$, $a_i = 0$ or $\gamma_{(x_i, y_i)}(w, b) = 1$
So, define $S = \{ i \mid a_i \neq 0 \} \Rightarrow \underline{\{ x_i \text{ lie} \in S \}}$

support vectors

we are on one of
the supporting, blue
hyperplanes

The gradient constraint gives $w = \sum_{i=1}^n a_i y_i x_i \in \text{span} \{ x_i, i \in S \}$.
So, the max-margin classifier $\hat{w}$ is defined by the support vectors! If we get
new, easier data, we don't change anything. Particularly, for any new point $\tilde{x}$,

$$y(\tilde{x}; \hat{w}, \hat{b}) = \text{sgn}(\hat{w}^T \tilde{x} + \hat{b}) = \text{sgn}\left( \sum_{i \in S} a_i y_i \, x_i^T \tilde{x} + \hat{b} \right)$$

hopefully small
# of support
vectors

dot product
kernel!
If we replace with
feature representations,
we get Kernel SVM.

**Theorem:** Given any init, as $t \to \infty$

- $\|w(t)\| \to \infty$
- $\mathcal{L}(w(t)) \to 0$
- $\frac{w(t)}{\|w(t)\|} = \frac{\hat{w}}{\|\hat{w}\|} + O\left(\frac{1}{\log t}\right)$ where $\hat{w}$ is the max-margin classifier

**Proof:**

First, suppose WOLOG that all points $(x_i, y_i) \to (y_i x_i, 1)$. This works since we set the bias to 0.

Note that GF and the definition of $\mathcal{L}$ gives

$$\frac{d}{dt} w_*^T w(t) = w_*^T \left(-\vec{\nabla}_w \mathcal{L}(w(t))\right) = - w_*^T \sum_{i=1}^{n} y_i x_i \, \ell'(w(t)^T x_i y_i)$$

$$= - \sum_{i=1}^{n} \underbrace{(y_i w_*^T x_i)}_{\substack{\text{independent of } t \\ >0 \text{ because} \\ w_* \in C_0}} \underbrace{\ell'(w(t)^T x_i y_i)}_{\substack{\text{uniform bound} \\ <0 \text{ because} \\ \ell(\cdot) \sim e^{-u}}} > 0$$

So, if the data is linearly separable, $\frac{d}{dt} w_*^T w(t) > 0$

Suppose BWOC that $\|w(t)\| \le R \;\; \forall t \ge 0$ (bounded). Then, $\exists \delta > 0$ s.t.

$$\ell'(y_i w(t)^T x_i) \le -\delta < 0 \implies \frac{d}{dt} w_*^T w(t) \ge \delta n \min_i (y_i w_*^T x_i)$$

This is a contradiction, since the derivative is uniformly bounded from below, and so must diverge. Therefore, $\boxed{\|w(t)\| \to \infty}$

Now, GF grants that $\frac{d}{dt} \mathcal{L}(w(t)) = -\|\vec{\nabla}_w \mathcal{L}(w(t))\|^2$

$$\le -\left(w_*^T \vec{\nabla}_w \mathcal{L}(w(t))\right)^2$$

Similar logic gives $\boxed{\mathcal{L}(w(t)) \to 0}$.

Lastly, let $r(t) = w(t) - \hat{w} \log(t) - \tilde{w}$, where

$$\forall i \in S, \; e^{-x_i^T \tilde{w}} \equiv a_i \implies \hat{w} = \sum_{i=1}^{n} e^{-x_i^T \tilde{w}} x_i$$

set $y_i = 1$ cutely

Let $\theta \equiv \min_{i \notin S} x_i^T \hat{w} > 0$ } min non-supported margin

We want to show that $\|r(t)\|$ is bounded.
We can do this by showing the derivative is integrable.

Now, $\frac{d}{dt} \frac{1}{2} \|r(t)\|^2 = \left(\frac{d}{dt} r(t)\right)^T r(t) = \left(-\tilde{\nabla} \mathcal{L}(w(t)) - \frac{1}{t} \hat{w}\right)^T r(t)$

$$= \left(\sum_{i=1}^{\hat{n}} x_i \, e^{-w(t)^T x_i} - \frac{1}{t} \hat{w}\right) r(t)$$

$$= \left(\sum_{i=1}^{\hat{n}} x_i \, e^{(-r(t) - \hat{w} \log(t) - \tilde{w})^T x_i} - \frac{1}{t} \sum_{i \in S} e^{-x_i^T \tilde{w}} x_i\right)^T r(t)$$

$$= \left(\sum_{i=1}^{\hat{n}} x_i \left(\frac{1}{t}\right)^{\hat{w}^T x_i} \overbrace{e^{-r(t)^T x_i}}^{\rightarrow \, \text{nega!}} e^{-\tilde{w}^T x_i} - \frac{1}{t} \sum_{i \in S} e^{-x_i^T \tilde{w}} x_i\right)^T r(t)$$

Collecting terms with $i \in S$,

$$\frac{1}{t} \sum_{i \in S} x_i \, e^{-x_i^T \tilde{w}} \left(e^{-r(t)^T x_i} - 1\right)^T r(t)$$

$$= \frac{1}{t} \sum_{i \in S} x_i^T r(t) \, e^{-x_i^T \tilde{w}} \underbrace{\left(e^{-x_i^T r(t)} - 1\right)}_{\color{blue}{c \cdot z \, (e^{-z} - 1) \leq 0 \, !}}$$

For $i \notin S$,

$$\left\| \sum_{i \notin S} x_i \left(\frac{1}{t}\right)^{\gamma_{(x_i, y_i)}(\hat{w})} e^{-\text{const.}} \right\| \leq n \cdot C \cdot \frac{1}{t^{\theta}} \implies \int_0^{\infty} \frac{d}{dt} \|r(t)\|^2 < \infty.$$

$\color{blue}{\theta > 1 \Rightarrow \text{bounded integral}}$

So, $\|r(t)\|$ is bounded $\implies$ $\boxed{\dfrac{w(t)}{\|w(t)\|} = \dfrac{\hat{w}}{\|\hat{w}\|} + O\left(\dfrac{1}{\log t}\right)}$ $\quad\quad$ $\square$

<u>Things we see:</u>
- longer gradient signal for small margins
- optimization moves in support vector directions
  - for these directions, the optimization has unique solution

<u>Remarks</u> $\quad\quad$ $\color{blue}{(\text{Conv, ReLU, avg}) \text{ are pos. homogenous}}$
1) This thing works for any homogenous classifier $\left(\begin{array}{c}\text{scaling changes score by a} \\ \text{power of the scalar}\end{array}\right)$
2) Convergence is slow $\quad \frac{1}{\log(t)}$

<u>Open problems</u>
* Include a bias? $\quad\quad$ * new losses? $\quad\quad$ * quadratic models

$\quad\quad$ Paper: Implicit bias of gD a separable data

# Lecture 11/16- SGD Implicit Bias

Suppose that we are given a model $z = (x; \theta)$, $\theta \in \mathbb{R}^n$ which we train by SGD:

$$\mathcal{L}(\theta) = \frac{1}{m} \sum_{k=1}^{\widehat{n}} \ell(x_k; \theta) \qquad \theta(t+1) = \theta(t) - \zeta \vec{\nabla} \mathcal{L}^B(\theta)$$

$$\mathcal{L}^B(\theta) = \frac{1}{|B|} \sum_{k=1}^{\widehat{n}} \mathbb{1}_{\{x_k \in B\}} \ell(x_k; \theta) \qquad x_k \in B \text{ w.p. } \frac{|B|}{m} \text{ independently}$$

The goal: We wish to understand the implicit bias of SGD.

tl;dr: SGD prefers "wider minima" or "flatter parts" of $\mathbb{R}^n$

## Yaida

Yaida uses the dynamical systems perspective that $\theta \sim \mathbb{P}_{ss}$ "steady state" and for any observable $O : \mathbb{R}^n \to \mathbb{R}$,

$$\boxed{\langle O(\theta) \rangle = \langle [[O(\theta - \zeta \vec{\nabla}\mathcal{L}^B(\theta))]] \rangle}$$

"fluctuation/dissipation" relationship

where $\langle \cdot \rangle$ is an average w.r.t. the $\mathbb{P}_{ss}$ distribution, and $[[\cdot]]$ is an average w.r.t. batching B.

The philosophy about this is to use FDR to
  (i) Taylor expand the RHS
  (ii) Collect powers of $\zeta$
  (iii) Compute properties of the steady-state $\mathbb{P}_{ss}$.

Def: Let $\widetilde{C}_{ij}(\theta) = (2^{nd}$ moments of $\vec{\nabla}\mathcal{L}^B$ w.r.t. B) $\leftarrow$ Second moment
$$= [[\partial_{\theta_i} \mathcal{L}^B(\theta) \, \partial_{\theta_j} \mathcal{L}^B(\theta)]]$$

<u>Lemma:</u> In the steady state distribution,

    (i) $\langle \vec{\nabla} \mathcal{L}(\theta) \rangle = 0$     (ii) $\langle \theta \cdot \vec{\partial} \mathcal{L}(\theta) \rangle = \langle \frac{1}{2} \zeta \, \text{tr}(\tilde{C}) \rangle \geq 0$    *← diagonal terms are squares*

            *no net gradient in steady state*                   *alignment is related to*    *batch-averaged loss covariances*

<u>Proof:</u> Consider the identity observable $O(\theta) = \theta$. Then, FDR gives

$$\langle \theta \rangle = \langle [[\theta - \zeta \vec{\nabla} \mathcal{L}^B(\theta)]] \rangle = \langle \theta \rangle - \zeta \langle [[\vec{\nabla} \mathcal{L}^B(\theta)]] \rangle$$

However, $\forall \theta$ we have        *|B| in expectation*

$$[[\vec{\nabla} \mathcal{L}^B(\theta)]] = [[ \frac{1}{|B|} \sum_{i=1}^{m} \mathbb{1}_{x_k \in B} \, \vec{\nabla} \ell(x_k; \theta) ]]$$

$$= \frac{1}{m} \sum_{i=1}^{m} \vec{\nabla} \ell(x_k; \theta) = \vec{\nabla} \mathcal{L}(\theta) \Rightarrow (\because).$$

So,    $\langle \theta \rangle = \langle \theta \rangle - \zeta \langle \vec{\nabla} \mathcal{L}(\theta) \rangle \Rightarrow \langle \vec{\nabla} \mathcal{L}(\theta) \rangle = 0.$

Next, if   $O(\theta) = \frac{1}{2} \theta_i^2$, FDR gives

$$\langle \frac{1}{2} \theta_i^2 \rangle = \langle [[ \frac{1}{2} (\theta_i - \zeta \partial_{\theta_i} \mathcal{L}^B(\theta))^2 ]] \rangle$$

$$= \langle [[ \frac{1}{2} \theta_i^2 - \zeta \theta_i \partial_{\theta_i} \mathcal{L}^B(\theta) + \frac{1}{2} \zeta^2 (\partial_{\theta_i} \mathcal{L}^B(\theta))^2 ]] \rangle$$

$$= \langle \frac{1}{2} \theta_i^2 \rangle - \zeta \langle \theta_i \partial_{\theta_i} \mathcal{L}^B(\theta) \rangle + \frac{\zeta^2}{2} \langle [[ (\partial_{\theta_i} \mathcal{L}^B(\theta))^2 ]] \rangle$$

$$\Rightarrow \zeta \langle \theta_i \partial_{\theta_i} \mathcal{L}(\theta) \rangle = \frac{\zeta^2}{2} \langle \tilde{C}_{ii} \rangle.$$

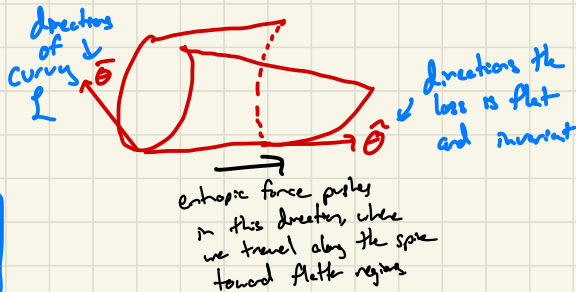Summing over $i$'s, we get (ii).            $\square$

# <u>Wei-Schwab-</u>

Typically mean, and so there $\exists$ mm directions in which $\mathcal{L}$ is flat.

Consider $\mathcal{L}(\theta) = \sum_{i=1}^{n} \bar{\theta}_i^2 \lambda_i (\hat{\theta}_{n+1}, \dots, \hat{\theta}_N).$    The landscape looks like

          *≥0* above

     $\theta = (\bar{\theta}, \hat{\theta}) \in \mathbb{R}^N$    *↑ # params*



*directions of $\downarrow$ curvic $\mathcal{L}$*   $\bar{\theta}$

*directions the loss is flat and invariant*

$\hat{\theta}$

At fixed $\hat{\theta}$ directions, the loss function as a function of $\bar{\theta}$ looks like a curved landscape with

$$\text{Hess}_{\bar{\theta}}(\mathcal{L}) = 2 \begin{bmatrix} \lambda_{n+1}(\hat{\theta}) & & 0 \\ & \ddots & \\ 0 & & \lambda_N(\hat{\theta}) \end{bmatrix}$$

*entropic force pushes in this direction where we travel along the spine toward flatter regions*

We have the intuition: $\hat{\theta}$ directions evolve "slowly" w.r.t. $\bar{\theta}$ directions since $\bar{\theta}$ directions drive the loss down.

So, we assume $\bar{\theta} \sim \mathbb{P}_{SS}$ and ask then about one step in $\hat{\theta}$ directions.

This assumes that $\bar{\theta}$ already equilibrates before substantial movement in $\hat{\theta}$ directions.

We sketch below some helpful claims and lemmas:

"Claim": Write $C_{ij} = \text{Cov}_{mb}\left(\partial_{\theta_i} \mathcal{L}^B(\theta), \partial_{\theta_j}\mathcal{L}^B(\theta)\right)$. ← Variance
At late times, $\qquad C = \alpha \, H_{\bar{\theta}}(\mathcal{L}) \qquad \alpha > 0$

Proof: Lol we don't prove this. Use it as a tool for later tho.  ▢

Note: For each $i = 1, \ldots, n,$ $\qquad \widetilde{C}_{ii} = C_{ii} + \left(\partial_{\theta_i}\mathcal{L}^B(\theta)\right)^2 \geq C_{ii}$.
This is the relation between second moments $\widetilde{C}$ and covariances $C$.

Corollary to note: $\langle \bar{\theta}_i^2 \rangle \geq \frac{\alpha}{4}\frac{3}{} + O(\zeta^2)$
↑ how high up the walls we walk in $\bar{\theta}$ $\mathbb{P}_{SS}$ directions

Proof: We apply FDR with the observable $O(\theta) = \langle \theta_i^2 \rangle$. Up to $O(\zeta^2)$,
we get $\qquad \langle \bar{\theta}_i \, \partial_{\bar{\theta}_i}\mathcal{L} \rangle = \frac{3}{4}\langle \widetilde{C}_{ii} \rangle$ $\qquad$ by lemma (ii) evaluated component-wise.
take a derivative $= \langle 2\bar{\theta}_i^2 \lambda_i(\hat{\theta}) \rangle = \frac{3}{4}\langle \widetilde{C}_{ii} \rangle$
$\qquad\qquad \parallel$
$\qquad 2\lambda_i(\hat{\theta})\langle\bar{\theta}_i^2\rangle \qquad \Rightarrow \langle\bar{\theta}_i^2\rangle = \frac{3}{4}\frac{\langle\widetilde{C}_{ii}\rangle}{\lambda_i(\hat{\theta})} \geq \frac{3}{4}\frac{\alpha \cdot H_{ii}(\hat{\theta})}{\lambda_i(\hat{\theta})}$  ← "claim"
$\Rightarrow \langle\bar{\theta}_i^2\rangle \geq \frac{3}{4}\alpha$.  $\qquad\qquad\qquad\qquad$ ▢

Note that in what we do above, we think of $\hat{\theta}$ as fixed as we determine the overall effect of $\mathbb{P}_{SS}$. Only after this do we consider a step in $\hat{\theta}$'s.

**Prop:** $\mathbb{E}$ w.r.t. both $P_{ss}$ and initializations

$$\mathbb{E}\left\{ tr\left( H_{\bar{\theta}}(\mathcal{L}(\theta(t+1)))\right) - tr\left( H_{\bar{\theta}}(\mathcal{L}(\theta(t)))\right)\right\} \leq 0$$

So, we go to places with smaller $\lambda$'s over time (flatter region in loss). We call this the **entropic force**.

**Proof:** Fix $i \in 1, ..., n$, $\hat{\theta}(t)$. We have

after $\bar{\theta}$'s equilibrate ($P_{ss}$ avg.), take a step in $\hat{\theta}$

$$\left\langle \left[\left[ \lambda_i(\theta(t+1)) - \lambda_i(\theta(t))\right]\right]\right\rangle$$

$$= \left\langle \left[\left[ \lambda_i(\hat{\theta}(t) - \gamma \vec{\nabla}_{\hat{\theta}} \mathcal{L}^{\bar{\theta}}(\theta(t))) - \lambda_i(\theta(t))\right]\right]\right\rangle$$

Taylor expand

$$= \left\langle \left[\left[ -\gamma \vec{\nabla}_{\hat{\theta}} \lambda_i(\hat{\theta}(t)) \cdot \vec{\nabla}_{\hat{\theta}} \mathcal{L}^{\bar{\theta}}(\theta(t)) + O(\gamma^2)\right]\right]\right\rangle$$

$$= -\gamma \left\langle \vec{\nabla}_{\hat{\theta}} \lambda_i(\hat{\theta}(t)) \cdot \vec{\nabla}_{\hat{\theta}} \mathcal{L}(\theta(t))\right\rangle + O(\gamma^2)$$

take derivative

$$= -\gamma \vec{\nabla}_{\hat{\theta}} \lambda_i(\hat{\theta}(t)) \sum_{j=1}^{\hat{n}} \langle \bar{\theta}_j^2 \rangle \vec{\nabla}_{\hat{\theta}} \lambda_j(\theta(t)) + O(\gamma^2)$$

Corollary

$$\leq -\frac{\alpha}{4} \gamma \vec{\nabla}_{\hat{\theta}} \lambda_i(\theta(t)) \sum_{j=1}^{\hat{n}} \vec{\nabla}_{\theta} \lambda_j(\hat{\theta}(t)) + O(\gamma^2)$$

Summing this over all $i$, we get our result.

Since we are riding up the walls (see corollary), we are moving in a direction to give us more room to ride up the walls.
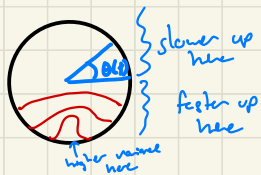
$\square$

# Challenge Problem

Suppose $\theta(t) \in [0, 2\pi)$.
We have the dynamics $\theta(t+dt) = \theta(t) \pm dt$ with probability $\frac{1}{2}$.

The steady state is given by $dP_{ss}(\theta) = \frac{1}{2\pi} d\theta$     uniform distribution

Suppose now that we have the dynamics



$$\theta(t+dt) - \theta(t) = \begin{cases} \pm\, dt & \theta < \pi \\[2mm] \pm 2\, dt & \theta \geq \pi \end{cases}$$

We expect more probability mass up top, since we bounce around the bottom $2\times$ faster.

The answer is $d\mathbb{P}_{ss}(\theta) = \frac{\sqrt{2}}{\sqrt{2}+1} \cdot \frac{1}{\pi} d\theta \, \mathbb{1}_{\theta < \pi} + \frac{1}{\sqrt{2}+1} \cdot \frac{1}{\pi} d\theta \, \mathbb{1}_{\theta > \pi}$

far away from the boundary.

The result we see is that we spend less time where the variance is larger.
Coming back to SGD, we have

$$\theta(t+1) - \theta(t) = -\eta\, \vec{\nabla}\mathcal{L}^B(\theta) = \underbrace{-\eta\, \vec{\nabla}\mathcal{L}(\theta)}_{\substack{\text{mean,} \\ \text{GF drift}}} + \underbrace{\eta\left(\vec{\nabla}\mathcal{L}(\theta) - \vec{\nabla}\mathcal{L}^B(\theta)\right)}_{\substack{\sim \mathcal{N}(0, C(\theta)) \\ \text{state-dependent diffusion term}}}$$

$\implies$ an implicit bias of SGD is that, in addition to minimizing loss (which the mean takes care of), to also $\boxed{\text{minimize } \operatorname{tr}(C(\theta))}$

We look for areas where the between-batch loss variance is low. This can be thought of as finding flat/isotropic/nice regions of the spine of the loss.

Fill in Comment notes here

# Lecture 11/28- Entropy + Widths

First, observe that **generalization** only makes sense given a priori <u>complexity information</u> about the function $f$ we want to learn.

To see this precisely, note that $\forall \Omega \subseteq \mathbb{R}^n$, $m \geq 1$, there exists $f: \Omega \to \mathbb{R}$ s.t. we cannot learn $f$ from any dataset of size $m$.

<u>Proof:</u> Discretize $\Omega = \coprod_{j=1}^{\infty} \Omega_j$ with $n > m$ and $f|_{\Omega_j}$ i.i.d. random.  □

We need better notions to talk about how complex a function is to encode/learn.

class of functions

<u>Def:</u> A **model class** $K$ is a <u>compact</u> subset of a Banach space $(X, \|\cdot\|_X)$.

<u>Some examples:</u>

① $K = \{ f: \Omega \to \mathbb{R} \mid \int_\Omega |f(x)|^2 + \|\nabla f\|^2 dx \leq 1 \} \subset L^2(\Omega)$

② $K = \{ f: \Omega \to \mathbb{R} \mid \|f\|_{Lip} \leq 1 \} \subset C^0(\Omega)$

⟵ Lipschitz constant

The question is: Given any method for "learning" $f \in K$, how do you measure "how well you did"?

# Entropy (Kolmogorov '30s)

model class
c to learn

<u>Def:</u> Let $\varepsilon_n(K) = n^{th}$ entropy # of $K = \inf \left\{ \varepsilon > 0 \mid \exists \text{ covering of } K \text{ by } 2^n \text{ balls of radius } \varepsilon \right\}$

# Intuitions

① $K$ compact $\Rightarrow$ finite cover $\Rightarrow$ $\varepsilon_n(K) < \infty$ $\forall n$

② $\varepsilon_n(K) =$ error in $\|\cdot\|_X$ of best $n$-bit compression of $k$

$\underbrace{K \subset \overset{2^n}{\underset{i=1}{\cup}} N_\varepsilon(f_i)}_{\text{def. of } \varepsilon_n(k)}$ yields a bijection $\{f_i\} \longleftrightarrow \{0,1\}^n$ where $f \in K \mapsto$ $\begin{smallmatrix}\text{nearest}\\ \text{ball center}\\ \text{in } K\end{smallmatrix}$

③ $\varepsilon_n(K)$ typically can be computed as $n \to \infty$, but this only tells us how hard a function is to learn, **not** how well a learning procedure does (not yet).

# Stable Width

**Def:** An <span style="color:red">$n$-param approximation scheme</span> for learning $K$ is a pair of functions

"param extraction" $\qquad a_n : K \to \mathbb{R}^n$

"reconstruction" $\qquad M_n : \mathbb{R}^n \to K$



**Def:** The <span style="color:red">error</span> of $(a_n, M_n)$ is
$$E_{a_n, M_n}(K) = \sup_{f \in K} \|f - M_n(a_n(f))\|_X \qquad \text{worst reconstruction error over } f \in K$$

**Def:** The <span style="color:red">stable $n$-width of $K$</span> is
$$S_n(K) = \inf_{a_n, M_n \in \text{Lip-2}} E_{a_n, M_n}(K) \qquad \text{best error we can do}$$

Note that intuitively, Lipschitz $\Leftrightarrow$ "numerically stable" in the sense that it excludes space-filling curves.

The amazing result is that $\mathcal{E}_n(k)$ and $\mathcal{S}_n(k)$ are equivalent!
We prove this below. First, recall the following results:

**Theorem:** (Johnson-Lindenstrauss Lemma)
Let $\varepsilon \in (0,1)$. For any $x_1, \ldots, x_k \in X$, $\exists$ a 1-Lipschitz (and linear!) function
$A: X \to \mathbb{R}^m$ s.t. $\forall i,j$ $(1-\varepsilon)\|x_i - x_j\|_X \le \|Ax_i - Ax_j\|_{\mathbb{R}^m} \le \|x_i - x_j\|_X$
as long as $m > \frac{8}{\varepsilon^2} \log(k)$.


**Theorem:** (Kirzbraun Extension Theorem)
If $f: U_1 \to \mathcal{H}_2$, $U_1 \subseteq \mathcal{H}_1$ is Lipschitz, then $\exists F: \mathcal{H}_1 \to \mathcal{H}_2$ s.t.
$$F|_{U_1} = f \quad \text{and} \quad \|F\|_{Lip} = \|f\|_{Lip} \quad \text{\small same Lipschitz constant}$$

With this machinery, we can prove both directions.


$(\Longrightarrow)$

**Theorem:** $\forall n \quad \mathcal{S}_{32n}(k) \le 3\mathcal{E}_n(k)$

**Proof:** Fix $n$. Choose $\{f_i, i \in [2^n]\} \subseteq K$ s.t. $K \subseteq \overset{2^n}{\underset{i=1}{\bigcup}} N_{\mathcal{E}_n(k)}(f_i)$  ← definition of $\mathcal{E}_n(k)$
Applying JL on these ball centers with $\varepsilon = \frac{1}{2}$, $k = 2^n$, $x_i = f_i$, we get
$a: K \to \mathbb{R}^{32n}$ s.t. $\forall i,j$, $\frac{1}{2}\|f_i - f_j\|_X \le \|a(f_i) - a(f_j)\|_{\mathbb{R}^{32n}} \le \|f_i - f_j\|_X$

Note that over $U_1 = \{a(f_i)\} \subseteq \mathbb{R}^{32n}$, a function $M_1: U_1 \to X$
that inverts $a$ on the ball centers (i.e. $M_1(a(f_i)) = f_i$) is 2-Lipschitz
by the JL inequality. So, by the extension theorem, there exists $M: \mathbb{R}^{32n} \to X$
that is 2-Lipschitz with $M(a(f_i)) = f_i$ $\forall i$. So, $\forall f \in K$,

$$\underbrace{\|f - M(a(f))\|_X}_{} \le \underbrace{\|f - f_i\|_X}_{\le \mathcal{E}_n(k)} + \underbrace{\|f_i - M(a(f_i))\|_X}_{= 0} + \underbrace{\|M(a(f_i)) - M(a(f))\|_X}_{\le 2\mathcal{E}_n(k) \text{ because } \|M\|_{Lip} = 2, \|a\|_{Lip} = 1 \Rightarrow \|M \circ a\|_{Lip} = 2} \quad \text{Triangle ineq.}$$

$$\le \mathcal{E}_n(k) + 0 + 2\mathcal{E}_n(k) = 3\mathcal{E}_n(k)$$
Since this holds for all $f \in K$,
$$\mathcal{S}_{32n}(k) \le \mathcal{E}_{a,m}(k) \le 3\mathcal{E}_n(k)$$

$\square$

$(\Leftarrow)$

__Theorem:__ Fix $r > 0$. Then, $\delta_n(K) \lesssim n^{-r} \Rightarrow \varepsilon_n(K) \lesssim \left(n/\log n\right)^{-r}$

($\varepsilon, \delta$ go to 0 together)

__Proof:__ Fix $n$ and consider a near-optimal $(a_n, M_n)$ s.t. $\delta = \mathcal{E}_{a_n, M_n}(K)$ and $\delta_n(K) \leq \delta \leq 2\delta_n(K)$. Suppose $a_n(K) \subset N_R(?) \subset \mathbb{R}^n$.

         image of $K$

Let $\left\{N_{2\delta}(f_i)\right\}_{i=1}^{P_\delta(K)}$ be a maximal $2\delta$-packing of $K$.

         ($P_\delta(K)$ is max # of disjoint balls of radius $2\delta$ fitting in $K$)

Note that $\left\{N_{4\delta}(f_i)\right\}_{i=1}^{P_\delta(K)}$ is a covering of $K$ (if not, we could have fit another $2\delta$ ball in the packing). We analyze the functions of $a_n, M_n$ at each ball center $f_i$. Note that $\forall i, j \in [P_\delta(K)]$,

$$\|M_n(a_n(f_i)) - M_n(a_n(f_j))\|_X \geq 2\delta \Rightarrow \|a_n(f_i) - a_n(f_j)\|_{\mathbb{R}^n} \geq \delta \quad \text{by } M_n \text{ 2-Lipschitz.}$$

Thus, $\left\{N_\delta(a_n(f_i)), i \in [P_\delta(K)]\right\}$ is a $\delta$-packing of $N_R(?)$ in $\mathbb{R}^n$.

Hence, $P_\delta(K) \leq \left(\frac{6R}{\delta}\right)^n = 2^{n \log\left(\frac{c}{\delta}\right)}$ for some $c$    we know the volume!

So, $\mathcal{E}_{n \log\left(\frac{c}{\delta}\right)}(K) \leq 4\delta \leq 8\delta_n(K)$    still shaky about these 3 lines ...

Then, if $\delta_n(K) \lesssim n^{-r}$, $\varepsilon_{n \log n}(K) \lesssim n^{-r}$.

         $\underbrace{\phantom{\delta_n(K)}}_{N}$    $\underbrace{\phantom{\varepsilon_{n\log n}}}_{\frac{N}{\log N}}$     □

We combine these as follows:

⭐ __Theorem:__ (Carl, Cohen, DeVore, ...)

   differ by a universal constant, ∠ grow the same

     When $K \subset\subset X$ and $X$ is a Hilbert space, $\varepsilon_n(K) \asymp \delta_n(K)$ as $n \to \infty$.

     __Proof:__ Results of the two above theorems as $n \to \infty$.     □

__Open problems!__

✳ Add a dataset of size $m$ (restrict $a$ to something factorable over an evaluation map at $m$ points)

✳ How regular (Lipschitz?) are NN functions?    ✳ Solidify relationship between above and statistical learning ($\varepsilon_n(K)$ is basically VC dim.)

# Lecture 11/30 - Path Counting for ReLU

Consider a ReLU FCNN
$$z_{i;\alpha}^{(\ell+1)} = \begin{cases} \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma(z_{j;\alpha}^{(\ell)}) & \ell \geq 1 \\ \sum_{j=1}^{n_\ell} W_{ij}^{(\ell)} x_{j;\alpha} & \ell = 0 \end{cases}$$

with widths $n_1, \ldots, n_L$ and
$$\sigma(t) = t \mathbb{1}_{t \geq 0}, \qquad W_{ij}^{(\ell+1)} = \sqrt{\frac{2}{n_\ell}} \hat{W}_{ij}^{(\ell+1)}, \qquad \hat{W}_{ij}^{(\ell+1)} \sim \mu \text{ i.i.d}$$

where the distribution $\mu$ is symmetric about the mean with variance 1 and finite moments.

$$\underline{\text{path counting}}$$

The goal is to explain a $\underline{\text{combinatorial approach}}$ to study any statistic of random ReLU at a single input $x_\alpha \neq 0$.

**Def:** For each $n \geq 1$ write $[n] = \{1, \ldots, n\}$.
The **space of paths** in a FCNN with widths $n_0, \ldots, n_{L+1}$ is
$$\Gamma = [n_0] \times \cdots \times [n_{L+1}]$$
i.e. $\gamma \in \Gamma$ is $\gamma = (\gamma(0), \gamma(1), \ldots, \gamma(L+1))$, $\gamma(\ell) \in [n_\ell]$ $\forall \ell$.



**Notation:** For each $\ell = 1, \ldots, L+1$, let
$$W_\gamma^{(\ell)} = W_{\gamma(\ell), \gamma(\ell-1)}^{(\ell)}, \qquad z_\gamma^{(\ell)} = z_{\gamma(\ell); \alpha}^{(\ell)} \qquad \#\Gamma = \prod_{\ell=0}^{L+1} n_\ell$$

$$\Gamma_{p,q} = \{\gamma \in \Gamma : \gamma(0) = p, \ \gamma(L+1) = q\}$$

**Prop.** We have
$$z_{a;i}^{(L+1)} = \sum_{p=1}^{n_0} x_{p;a} \sum_{\gamma \in \Gamma_{p,a}} W_\gamma^{(L+1)} \prod_{l=1}^{L} W_\gamma^{(l)} \xi_\gamma^{(l)} \quad \text{(A)}$$

where $\xi_\gamma^{(l)} = \mathbb{1}_{\{z_\gamma^{(l)} \geq 0\}}$ indicates if the neuron $\gamma(l)$ is on.

**Proof:** Given $\vec{v} = (v_1, ..., v_n) \in \mathbb{R}^n$, $\sigma(\vec{v}) = D_v \vec{v}$, where
$D_v = \text{Diag}(\mathbb{1}_{\{v_1 \geq 0\}}, ..., \mathbb{1}_{\{v_n \geq 0\}})$. Thus,
$$\vec{z}_a^{(L+1)} = W^{(L+1)} \sigma(W^{(L)} \sigma(... \sigma(W^{(1)} \vec{z}_a)...)) = W^{(L+1)} D^{(L)} W^{(L)} ... D^{(1)} W^{(1)} \vec{x}_a$$

where $D^{(l)} = \text{Diag}(\mathbb{1}_{\{z_{i;a}^{(l)} \geq 0\}}, i=1, ..., n_l)$. Hence,
$$z_{a;i}^{(L+1)} = \sum_{p=1}^{n_0} x_{p;a} \left( W^{(L+1)} D^{(L)} W^{(L)} ... D^{(1)} W^{(1)} \right)_{pa}$$
$$= \sum_{p=1}^{n_0} x_{p;a} \sum_{i_1=1}^{n_1} \sum_{i_2=1}^{n_2} ... \sum_{i_L=1}^{n_L} D_{i_1 i_1}^{(1)} W_{i_1 p}^{(1)} \cdot ... \cdot D_{i_L i_L}^{(L)} W_{i_L i_{L-1}}^{(L)} \cdot W_{a i_L}^{(L-1)}$$
$$= \sum_{p=1}^{n_0} x_{p;a} \sum_{\gamma \in \Gamma_{p,a}} W_\gamma^{(L+1)} \prod_{l=1}^{L} W_\gamma^{(l)} \xi_\gamma^{(l)} \qquad \square$$

**Prop:** At init,
$$\mathbb{1}_{\{z_{i;a}^{(l)} \geq 0\}} \overset{d}{=} \text{Bernoulli}\left(\tfrac{1}{2}\right) \text{ i.i.d}$$

and independent of any __even__ function of the weights $W^{(l)}$'s.

<span style="color:blue">i.e. symmetric under $W^{(l)} \mapsto -W^{(l)}$</span>

**Proof:**

   Idea: given $\vec{z}_a^{(l)}$, i.i.d weights symmetric about 0 means
$$\mathbb{1}_{\{z_{i;a}^{(l+1)} \geq 0\}} \overset{d}{=} \text{Bernoulli}(\tfrac{1}{2}) \text{ given } \vec{z}_a^{(l)}.$$
   Since this distribution is the same regardless of $\vec{z}_a^{(l)}$, we are done.

   For the weight-wise independence, check the paper :-) 

$\square$

**Corollary:** We have
$$z_\alpha^{(L+1)} \stackrel{d}{=} W^{(L+1)} \hat{D}^{(L)} W^{(L)} \cdots \hat{D}^{(1)} W^{(1)} x_\alpha$$
where $\hat{D}_{ii}^{(\ell)} \sim \text{Bernoulli}(\tfrac{1}{2})$ i.i.d.

**Proof:** Duh.                                                    ☐

**Corollary:** We have
$$\frac{\partial z_{a;i\alpha}^{(L+1)}}{\partial x_{p;\alpha}} = \sum_{\gamma \in \Gamma_{p,a}} W_\gamma^{(L+1)} \prod_{\ell=1}^{L} W_\gamma^{(\ell)} \xi_\gamma^{(\ell)}$$

**Proof:** Duh.                                                    ☐

**Lemma:** For any $n_0, \ldots, n_{L+1}$,  $\mathbb{E}\left\{\left(\frac{\partial z_{a;i\alpha}^{(L+1)}}{\partial x_{p;\alpha}}\right)^2\right\} = \frac{2}{n_0}.$

**Proof:** Let $A = \mathbb{E}\left\{\left(\frac{\partial z_{a;i\alpha}^{(L+1)}}{\partial x_{p;\alpha}}\right)^2\right\}$. Then

$$A = \mathbb{E}\left\{ \sum_{\gamma_1, \gamma_2 \in \Gamma_{p,a}} \prod_{\kappa=1}^{2}\left( W_{\gamma_\kappa}^{(L+1)} \prod_{\ell=1}^{L} W_{\gamma_\kappa}^{(\ell)} \xi_{\gamma_\kappa}^{(\ell)} \right)\right\}$$

$\underline{\text{even function in } W^{(\ell)}, \text{ so } \xi\text{'s independent}}$

$$\Rightarrow A = \sum_{\gamma_1, \gamma_2 \in \Gamma_{p,a}} \mathbb{E}\left\{ \prod_{\kappa=1}^{2} W_{\gamma_\kappa}^{(L+1)}\right\} \cdot \prod_{\ell=1}^{L} \mathbb{E}\left\{\prod_{\kappa=1}^{2} W_{\gamma_\kappa}^{(\ell)}\right\} \mathbb{E}\left\{\prod_{\kappa=1}^{2} \xi_{\gamma_\kappa}^{(\ell)}\right\}$$

But note that $\mathbb{E}\left\{\prod_{\kappa=1}^{2} W_{\gamma_\kappa}^{(\ell)}\right\} = \frac{2}{n_{\ell-1}} \delta_{\gamma_1(\ell)\,\gamma_2(\ell)} \delta_{\gamma_1(\ell-1)\,\gamma_2(\ell-1)}$

↳ variance       ;if not the same, this is a product of i.i.d. mean $0$ weights

So, we sum over $\gamma = \gamma_1 = \gamma_2$. For this, $\mathbb{E}\left\{\prod_{\kappa=1}^{2} \xi_{\gamma_\kappa}^{(\ell)}\right\} = \mathbb{E}\left\{\xi_\gamma^{(\ell)2}\right\} = \frac{1}{2}.$

All together,
$$A = \sum_{\gamma \in \Gamma_{p,a}} \frac{2}{n_L} \prod_{\ell=1}^{L} \frac{2}{n_{\ell-1}} \cdot \frac{1}{2} = \frac{2}{n_0} \frac{1}{\prod_{\ell=1}^{L} n_\ell} \sum_{\gamma \in \Gamma_{p,a}} 1$$

$$= \frac{2}{n_0} \mathcal{E}\{1\} \quad \text{← expectation over uniform measure in path space}$$

where $\mathcal{E}$ is an average over choices of random $\gamma \in \Gamma_{p,q}$ with
$\gamma(0) = p$, $\gamma(L+1) = q$, $\gamma(\ell) \sim \text{Unif}([n_\ell])$ independently.

Clearly, $A = \frac{2}{n_0}$

$\square$

---

<u>Theorem:</u> (Boris Spitkin)
When $n_1, \ldots, n_L$ are large, let $\beta = S \sum_{\ell=1}^{L} \frac{1}{n_\ell}$ $\quad$ (S · aspect ratio $(r = \frac{L}{n})$)
Then, $\left( \frac{\partial z_{q;\alpha}^{(L+1)}}{\partial x_{p;\alpha}} \right)^2 = \exp\left( \mathcal{N}\left(-\frac{\beta}{2}, \beta\right) + O\left(\frac{\beta}{n}\right) \right)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \sim O\left(\frac{L}{n^2}\right)$
<u>Exponentially</u> sensitive in the aspect ratio!

<u>Proof:</u> Nope :̈

$\square$

---

<u>Exercise:</u> Show that $\quad \mathbb{E}\left\{ \left( \frac{\partial z_{q;\alpha}^{(L+1)}}{\partial x_{p;\alpha}} \right)^4 \right\} \simeq \frac{\text{const}}{n_0^2} \exp\left( S \sum_{\ell=1}^{L} \frac{1}{n_\ell} \right)$

---

<u>Lemma:</u> Consider the on-diagonal NTK
$$\Theta_{\alpha\alpha}^{(L+1)} = \| \dot{\nabla}_\theta z_\alpha^{(L+1)} \|^2 \quad \text{when } n_{L+1} = 1$$
$$= \sum_{\ell=1}^{L+1} \sum_{i=1}^{n_\ell} \sum_{j=1}^{n_{\ell-1}} \left( \frac{\partial z_{1;\alpha}^{(L+1)}}{\partial \hat{W}_{ij}^{(\ell)}} \right)^2$$

$\Rightarrow \mathbb{E}\{ \Theta_{\alpha\alpha}^{(L+1)} \} = 2L \frac{\|\vec{x}_\alpha\|^2}{n_0}$ $\qquad$ <u>call this B</u>

<u>Proof:</u> Note that $\frac{\partial z_{1;\alpha}^{(L+1)}}{\partial \hat{W}_{ij}^{(\ell)}} = \sum_{p=1}^{n_0} x_{p;\alpha} \cdot \sum_{\substack{\delta \in \Gamma_{p,1} \\ \gamma(\ell) = i \\ \gamma(\ell-1) = j}} \frac{\sqrt{\frac{2}{n_L}} \hat{W}_\gamma^{(L+1)} \prod_{\ell=1}^{L} \sqrt{\frac{2}{n_{\ell-1}}} \hat{W}_\gamma^{(\ell)} \xi_\gamma^{(\ell)}}{\hat{W}_{ij}^{(\ell)}}$ $\qquad$ by $\star$

Then, $\mathbb{E}\{B\} = \mathbb{E}\left\{ \sum_{p_1,p_2=1}^{n_0} x_{p_1;\alpha} x_{p_2;\alpha} \cdot \sum_{\substack{\delta_1, \delta_2 \in \Gamma_{p_1} \\ \gamma_k \ni \cup_{ij}^{(\ell)}}} \frac{\frac{2}{n_L} \prod_{\kappa=1}^{2} \hat{W}_{\gamma_\kappa}^{(L+1)} \prod_{\ell=1}^{L} \frac{2}{n_{\ell-1}} \prod_{\kappa=1}^{2} \hat{W}_{\gamma_\kappa}^{(\ell)} \xi_{\gamma_\kappa}^{(\ell)}}{\left( \hat{W}_{ij}^{(\ell)} \right)^2} \right\}$

$$= \sum_{p_1, p_2}^{n_0} x_{p_1, \alpha} \, x_{p_2, \alpha} \sum_{\substack{\gamma_1, \gamma_2 \in \Gamma_{p,1} \\ \gamma_k \ni \hat{w}_{i,j}^{(\ell)}}} \prod_{\ell \neq L} \frac{2}{n_{\ell-1}} \; \mathbb{E}\left\{ \prod_{k=1}^{2} \hat{w}_{\gamma_k}^{(\ell)} \right\} \mathbb{E}\left\{ \gamma_{\gamma_k}^{(\ell)} \right\} \mathbb{E}\left\{ \frac{2}{n_\ell} \prod_{k=1}^{2} \hat{w}_{\gamma_k}^{(\ell, n)} \right\}$$

$$\mathbb{E}\left\{ \frac{2}{n_{L-1}} \prod_{k=1}^{2} \gamma_{\delta_k}^{(L)} \right\}$$

$$= \sum_{p=1}^{n_0} x_{p, \alpha}^{2} \cdot 2 \prod_{\ell=0}^{L} \frac{1}{n_\ell} \sum_{\substack{\gamma \in \Gamma_{p,1} \\ \hat{w}_{i,j}^{(\ell)} \in \gamma}} 1 = \frac{2 \|\vec{x}_\alpha\|^2}{n_0} \cdot \frac{1}{\prod_{\ell=1}^{L} n_\ell} \; \#\left\{ \gamma \in \Gamma_{p,1} : w_{i,j}^{(\ell)} \in \gamma \right\}$$

$$= \frac{2 \|\vec{x}_\alpha\|^2}{n_0} \left( n_\ell n_{\ell-1} \right)^{-1}$$

So,
$$\mathbb{E}\left\{ \Theta_{\alpha\alpha}^{(L+1)} \right\} = \sum_{\ell=1}^{L} \sum_{i=1}^{n_\ell} \sum_{j=1}^{n_{\ell-1}} \frac{2 \|\vec{x}_\alpha\|^2}{n_0} \left( n_\ell n_{\ell-1} \right)^{-1} = \frac{2L \|\vec{x}_\alpha\|^2}{n_0}$$

$$\square$$

**Lemma:** Consider the off-diagonal NTK

$$\Theta_{\alpha\beta}^{(L+1)} = \left(\dot{\nabla}_\theta z_\alpha^{(L+1)}\right)^T \left(\dot{\nabla}_\theta z_\beta^{(L+1)}\right) \quad \text{where} \quad n_{L+1} = 1$$

$$= \sum_{l=1}^{L+1} \sum_{i=1}^{n_\ell} \sum_{j=1}^{n_{\ell-1}} \frac{\partial z_{1;\alpha}^{(L+1)}}{\partial \hat{W}_{ij}^{(\ell)}} \frac{\partial z_{1;\beta}^{(L+1)}}{\partial \hat{W}_{ij}^{(\ell)}} \quad \Big\} = 0 \text{ when } \ell = L+1 ??$$

**Proof:** Note that

$$\frac{\partial z_{1;\alpha}^{(L+1)}}{\partial \hat{W}_{ij}^{(\ell)}} = \sum_{p=1}^{n_0} x_{p;\alpha} \cdot \sum_{\substack{\delta \in \Gamma_{p,1} \\ \gamma(\ell)=i \\ \gamma(\ell-1)=j}} \frac{\sqrt{\frac{2}{n_L}} \hat{W}_\gamma^{(L+1)} \prod_{\ell'=1}^{L} \sqrt{\frac{2}{n_{\ell'-1}}} \hat{W}_\gamma^{(\ell')} \xi_{\gamma,\alpha}^{(\ell')}}{\hat{W}_{ij\,\gamma}^{(\ell)}} \quad \text{by} \; \textcolor{blue}{\bigstar}$$

and similarly

$$\frac{\partial z_{1;\beta}^{(L+1)}}{\partial \hat{W}_{ij}^{(\ell)}} = \sum_{p=1}^{n_0} x_{p;\beta} \cdot \sum_{\substack{\delta \in \Gamma_{p,1} \\ \gamma(\ell)=i \\ \gamma(\ell-1)=j}} \frac{\sqrt{\frac{2}{n_L}} \hat{W}_\gamma^{(L+1)} \prod_{\ell'=1}^{L} \sqrt{\frac{2}{n_{\ell'-1}}} \hat{W}_\gamma^{(\ell')} \xi_{\gamma,\beta}^{(\ell')}}{\hat{W}_{ij}^{(\ell)}}$$

Then, 
$$\mathbb{E}\{\Theta_{\alpha\beta}^{\}} = \mathbb{E}\left\{ \sum_{p_\alpha,p_\beta=1}^{n_0} x_{p;\alpha} x_{p;\alpha} \cdot \sum_{\substack{\gamma_\alpha \in \Gamma_{p_\alpha,1} \\ \gamma_\beta \in \Gamma_{p_\beta,1}}} \frac{\frac{2}{n_L} \hat{W}_{\gamma_\alpha}^{(L+1)} \hat{W}_{\gamma_\beta}^{(L+1)} \prod_{\ell=1}^{L} \frac{2}{n_{\ell-1}} \hat{W}_{\gamma_\alpha}^{(\ell)} \xi_{\gamma_\alpha}^{(\ell)} \hat{W}_{\gamma_\beta}^{(\ell)}}{\left(\hat{W}_{ij}^{(\ell)}\right)^2} \xi_{\gamma_\beta\,\beta}^{(\ell)} \right\}$$

$$\gamma_\alpha(\ell) = \gamma_\beta(\ell) =: i$$
$$\gamma_\alpha(\ell-1) = \gamma_\beta(\ell-1) =: j$$

$$\textcolor{blue}{\text{= b.c. } \xi \text{ ind. from each fn. at weights}}$$

$$= \sum_{\substack{\rho_\alpha \rho_\beta = 1}}^{n_0} x_{p_\alpha\alpha} x_{p_\beta\beta} \sum_{\substack{\gamma_\alpha \in \Gamma_{p_\alpha,1} \\ \gamma_\beta \in \Gamma_{p_\beta,1}}} \frac{2}{n_L} \mathbb{E}\{\hat{W}_{\gamma_\alpha}^{(L+1)} \hat{W}_{\gamma_\beta}^{(L+1)}\} \prod_{\ell' \neq \ell} \frac{2}{n_{\ell'-1}} \mathbb{E}\{\hat{W}_{\gamma_\alpha}^{(\ell')} \hat{U}_{\gamma_\beta}^{(\ell')}\} \mathbb{E}\{\xi_{\gamma_\alpha}^{(\ell')} \xi_{\gamma_\beta}^{(\ell')}\}$$

$$\gamma_\alpha(\ell) = \gamma_\beta(\ell) =: i$$
$$\gamma_\alpha(\ell-1) = \gamma_\beta(\ell-1) =: j$$

$$\frac{2}{n_{\ell-1}} \mathbb{E}\{\xi_{\gamma_\alpha\alpha}^{(\ell)} \xi_{\gamma_\beta\beta}^{(\ell)}\} \qquad \textcolor{blue}{\underbrace{\qquad}_{\text{= between layers}}}$$

For any paths $\gamma_\alpha \in \Gamma_{p_\alpha,1}$, $\gamma_\beta \in \Gamma_{p_\beta,1}$ that are distinct, then their contributions disappear. Precisely, suppose that $\gamma_\alpha(\hat{\ell}) \neq \gamma_\beta(\hat{\ell}) \in [n_{\hat{\ell}}]$ for some $\hat{\ell} \neq \ell$. Then, $\hat{W}_{\gamma_\alpha(\hat{\ell}),\,\gamma_\alpha(\hat{\ell}-1)}^{(\hat{\ell})} \perp\!\!\!\perp \hat{W}_{\gamma_\beta(\hat{\ell}),\,\gamma_\beta(\hat{\ell}-1)}^{(\hat{\ell})}$ and so $\mathbb{E}\{\hat{W}_{\gamma_\alpha}^{(\hat{\ell})} \hat{U}_{\gamma_\beta}^{(\hat{\ell})}\} = \mathbb{E}\{\hat{W}_{\gamma_\alpha}^{(\hat{\ell})}\} \mathbb{E}\{\hat{U}_{\gamma_\beta}^{(\hat{\ell})}\} = 0$.

$\gamma_\alpha$ and $\gamma_\beta$ cannot disagree at $\ell$ either, since these do not contribute to the derivative w.r.t. $\hat{W}_{ij}^{(\ell)}$ (we required $\gamma_\alpha(\ell) = \gamma_\beta(\ell) =: i$, $\gamma_\alpha(\ell-1) = \gamma_\beta(\ell-1) =: j$). Thin layer at $\hat{\ell}=1$ requires that $p_\alpha = p_\beta$ as well. So, we sum over identical paths to get

$$= \sum_{p=1}^{n_0} x_{p;\alpha} x_{p;\beta} \sum_{\substack{\gamma \in \Gamma_{p,1} \\ \gamma(\ell) = i \\ \gamma(\ell-1) = j}} \frac{2}{n_L} \underbrace{\mathbb{E}\{\hat{W}_\gamma^{(L+1)^2}\}}_{= \text{Var } \hat{W} = 1} \prod_{\ell' \neq \ell} \underbrace{\mathbb{E}\{\hat{W}_\gamma^{(\ell')^2}\}}_{=1} \prod_{\ell=1}^{L} \frac{2}{n_{\ell-1}} \mathbb{E}\{\xi_{\gamma;\alpha}^{(\ell)} \xi_{\gamma;\beta}^{(\ell)}\}$$

$$= \sum_{p=1}^{n_0} x_{p\alpha} x_{p\beta} \cdot \frac{2^{L+1}}{\prod_{\ell=0}^{L} n_{\ell'}} \cdot \sum_{\substack{\gamma \in \Gamma_{p,1} \\ \gamma(\ell)=i \\ \gamma(\ell-1)=j}} \prod_{\ell=1}^{L} \mathbb{E}\{\xi_{\gamma;\alpha}^{(\ell)} \xi_{\gamma;\beta}^{(\ell)}\}$$

Due to symmetry over paths (since $W_{ij}^{(\ell)} \sim \mu$), $\prod_{\ell=1}^{L} \mathbb{E}\{\xi_{\gamma;\alpha}^{(\ell)} \xi_{\gamma;\beta}^{(\ell)}\}$ is the same $\forall \gamma$.

$$= \sum_{p=1}^{n_i} x_{p\alpha} x_{p\beta} \cdot \frac{2^{L+1}}{\prod_{\ell=0}^{L} n_{\ell'}} \cdot \prod_{\ell=1}^{L} \mathbb{E}\left\{ \mathcal{Z}_{\gamma;\alpha}^{(\ell)} \mathcal{Z}_{\gamma;\beta}^{(\ell')} \right\} \cdot \left| \left\{ \gamma \in \Gamma_{p,1} \mid \gamma(L)=i, \ \gamma(L-1)=j \right\} \right|$$

<span style="color:blue">$$\frac{\prod_{\ell=1}^{L} n_{\ell'}}{n_L \cdot n_{L-1}}$$</span>

$$= \sum_{p=1}^{\hat{n}_0} x_{p\alpha} x_{p\beta} \cdot \frac{2^{L+1}}{n_0 \, n_{L-1} \, n_L} \cdot \prod_{\ell=1}^{L} \mathbb{E}\left\{ \mathcal{Z}_{\gamma;\alpha}^{(\ell)} \mathcal{Z}_{\gamma;\beta}^{(\ell')} \right\} = \vec{x}_\alpha \cdot \vec{x}_\beta \cdot \frac{2^{L+1}}{n_0 \, n_{L-1} \, n_L} \prod_{\ell=1}^{L} \mathbb{E}\left\{ \mathcal{Z}_{\gamma;\alpha}^{(\ell)} \mathcal{Z}_{\gamma;\beta}^{(\ell')} \right\}$$

This gives

$$\mathbb{E}\left\{ \Theta_{\alpha\beta}^{(L+1)} \right\} = \sum_{\ell=1}^{L} \sum_{i=1}^{n_L} \sum_{j=1}^{n_{L-1}} \vec{x}_\alpha \cdot \vec{x}_\beta \cdot \frac{2^{L+1}}{n_0 \, n_{L-1} \, n_L} \cdot \prod_{\ell=1}^{L} \mathbb{E}\left\{ \mathcal{Z}_{\gamma;\alpha}^{(\ell)} \mathcal{Z}_{\gamma;\beta}^{(\ell')} \right\}$$

$$= \vec{x}_\alpha \cdot \vec{x}_\beta \cdot \frac{2^{L+1}}{n_0} \cdot \sum_{\ell=1}^{L} \prod_{\ell=1}^{L} \mathbb{E}\left\{ \mathcal{Z}_{\gamma;\alpha}^{(\ell)} \mathcal{Z}_{\gamma;\beta}^{(\ell')} \right\}$$

<span style="color:blue">this sum goes to L instead of L+1. Why??</span>

$$\boxed{\color{red}{= \vec{x}_\alpha \cdot \vec{x}_\beta \cdot \frac{L}{n_0} \cdot 2^{L+1} \prod_{\ell=1}^{L} \mathbb{E}\left\{ \mathcal{Z}_{\gamma;\alpha}^{(\ell)} \mathcal{Z}_{\gamma;\beta}^{(\ell)} \right\}}}$$

<span style="color:blue">P both inputs turn on ALL neurons in a path</span>

We have $K \in \mathbb{R}^{n \times n}$ as the expected NTK an init.
Thus define $\mu_{max} = \lambda_{max}(K)$. Since the top eigenvalue is $\leq$ any matrix norm, we get

$$\mu_{max} \leq \max_\beta \sum_{\alpha=1}^{\hat{n}} |K_{\alpha\beta}| = \max_\beta \sum_{\alpha=1}^{\hat{n}} |\vec{x}_\alpha \cdot \vec{x}_\beta| \frac{L}{n_0} 2^{L+1} \cdot \prod_{\ell=1}^{L} \mathbb{E}\left\{ \mathcal{Z}_{\gamma\alpha}^{(\ell)} \mathcal{Z}_{\gamma\beta}^{(\ell)} \right\}$$

We have $\mathbb{E}\left\{ \mathcal{Z}_{\gamma\alpha}^{(1)} \mathcal{Z}_{\gamma\beta}^{(1)} \right\} = \frac{1}{2} - \frac{1}{2\pi} \arccos\left( \frac{\vec{x}_\alpha \cdot \vec{x}_\beta}{\|\vec{x}_\alpha\| \|\vec{x}_\beta\|_2} \right)$

and $\mathbb{E}\left\{ \mathcal{Z}_{\gamma\alpha}^{(\ell)} \mathcal{Z}_{\gamma\beta}^{(\ell)} \right\} \leq \frac{1}{2}$

<span style="color:blue">row/col. norm</span> $\Rightarrow \mu_{max} \leq \max_\beta \sum_{\alpha=1}^{\hat{n}} |\vec{x}_\alpha \cdot \vec{x}_\beta| \frac{2L}{n_0} \left( 1 - \frac{1}{\pi} \arccos\left( \frac{\vec{x}_\alpha \cdot \vec{x}_\beta}{\|\vec{x}_\alpha\| \|\vec{x}_\beta\|} \right) \right)$

Also, $M_{max}^2 \leq \sum_{\alpha=1}^{\hat{n}} \sum_{\beta=1}^{\hat{n}} (\vec{x}_\alpha \cdot \vec{x}_\beta)^2 \frac{4L^2}{n_o^2} \left(1 - \frac{1}{\pi} \arccos\left(\frac{\vec{x}_\alpha \cdot \vec{x}_\beta}{\|\vec{x}_\alpha\| \|\vec{x}_\beta\|}\right)\right)^2$

Let $C_D = \min \left\{ \dfrac{\max_\beta \sum_{\alpha=1}^{\hat{n}} |\vec{x}_\alpha \cdot \vec{x}_\beta| \frac{2L}{n_o}\left(1 - \frac{1}{\pi}\arccos\left(\frac{\vec{x}_\alpha \cdot \vec{x}_\beta}{\|\vec{x}_\alpha\|\|\vec{x}_\beta\|}\right)\right)}{\sqrt{\sum_{\alpha=1}^{\hat{n}}\sum_{\beta=1}^{\hat{n}} (\vec{x}_\alpha \cdot \vec{x}_\beta)^2 \frac{4L^2}{n_o^2}\left(1 - \frac{1}{\pi}\arccos\left(\frac{\vec{x}_\alpha \cdot \vec{x}_\beta}{\|\vec{x}_\alpha\|\|\vec{x}_\beta\|}\right)\right)^2}} \right\}$

be dataset dependent. Then, $M_{max} \leq C_D$.

the following

**Proposition 3** (Pure weight moments for $K_N, \Delta K_N$). We have

$$\mathbb{E}[K_w] = \frac{d}{n_0}\|x\|_2^2.$$

Moreover,

$$\mathbb{E}[K_w^2] \simeq \frac{d^2}{n_0^2}\|x\|_2^4 \exp(5\beta)\left(1 + O\left(\sum_{i=1}^d \frac{1}{n_i^2}\right)\right), \qquad \beta := \sum_{i=1}^d \frac{1}{n_i}.$$

Finally,

Suppose $\mathbb{E}\{k_{\alpha\alpha}^2\} \leq C_1 \frac{4d^2}{n_o^2}\|\vec{x}_\alpha\|^4 e^{5\beta}$

$\Rightarrow \text{Var}(k_{\alpha\alpha}) \leq (C_1 e^{5\beta} - 1)\frac{4d^2}{n_o^2}\|\vec{x}_\alpha\|^4$

$\Rightarrow \sigma \leq \sqrt{C_1 e^{5\beta} - 1} \cdot \frac{2d}{n_o}\|\vec{x}_\alpha\|^2$

$\Rightarrow \mathbb{P}\left\{ k_{\alpha\alpha} \geq \left(1 + \sqrt{\frac{m}{\delta}}\sqrt{C_1 e^{5\beta} - 1}\right) \cdot \frac{2d}{n_o}\|\vec{x}_\alpha\|^2\right\} \leq \frac{\delta}{m}$

$\Rightarrow \mathbb{P}\left\{T_r(k) \leq \left(1 + \sqrt{\frac{m}{\delta}}\sqrt{C_1 e^{5\beta} - 1}\right) \cdot \frac{2d}{n_o}\sum\|\vec{x}_\alpha\|^2\right\} \geq 1 - \delta$

$\Rightarrow R = \left(1 + \sqrt{\frac{m}{\delta}}\sqrt{C_1 e^{5\beta} - 1}\right) \cdot \frac{2d}{n_o}\sum\|\vec{x}_\alpha\|^2$

Then, Matrix Chernoff gives

$$\mathbb{P}\left\{\lambda_{max}(K) \geq (1+\varepsilon)\mu_{max}\right\} \leq m\left(\frac{e^{\varepsilon}}{(1+\varepsilon)^{1+\varepsilon}}\right)^{\mu_{max}/R}$$

$$\Rightarrow \mathbb{P}\left\{\lambda_{max}(k) \geq (1+\varepsilon)C_{\Delta}\right\} \leq m\left(\frac{e^{\varepsilon}}{(1+\varepsilon)^{1+\varepsilon}}\right)^{C_m/R}$$

where

$$C_{\Delta} = \min\left\{\frac{\max_{\beta}\sum_{\alpha=1}^{\hat{n}}|\vec{x}_{\alpha}\cdot\vec{x}_{\beta}|\frac{2L}{n_0}\left(1-\frac{1}{\pi}\arccos\left(\frac{\vec{x}_{\alpha}\cdot\vec{x}_{\beta}}{\|\vec{x}_{\alpha}\|\|\vec{x}_{\beta}\|}\right)\right)}{\sqrt{\sum_{\alpha=1}^{\hat{n}}\sum_{\beta=1}^{\hat{n}}\left(\vec{x}_{\alpha}\cdot\vec{x}_{\beta}\right)^2\frac{nL^2}{n_0^2}\left(1-\frac{1}{\pi}\arccos\left(\frac{\vec{x}_{\alpha}\cdot\vec{x}_{\beta}}{\|\vec{x}_{\alpha}\|\|\vec{x}_{\beta}\|}\right)\right)^2}}\right\}$$

$$C_m = \frac{2L}{n_0}\max_{\alpha}\left\{\|\vec{x}_{\alpha}\|^2\right\}$$

$$R = \left(1+\sqrt{\frac{m}{\delta}}\sqrt{C_1 e^{5\beta}-1}\right)\cdot\frac{2d}{n_0}\sum\|\vec{x}_{\alpha}\|^2$$

# Lecture 12/5- Linear Regions

Consider a FC ReLU net,

$$z_i^{(\ell+1)}(\vec{x}) = b_i^{(\ell+1)} + \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \sigma\left(z_j^{(\ell)}\right)$$

with $n_{L+1} = 1$. Note that $\vec{x} \in \mathbb{R}^{n_0} \mapsto z^{(L+1)}(x) \in \mathbb{R}$ is continuous, piecewise linear.

One question we can ask is <u>how many pieces we get</u> in best/worst/avg cases?

We can use this result as a very rough measure of the complexity of ReLU nets.

## Examples

### Example 1

$n_0 = L = 1$

$$z^{(2)}(x) = b^{(2)} + \sum_{j=1}^{n_1} W_j^{(2)} \sigma\left(W_j^{(1)} x + b_j^{(1)}\right)$$

As on the pset, we define breakpoints $\xi_j = -\dfrac{b_j^{(1)}}{W_j^{(1)}}$

Since $\dfrac{dz^{(2)}}{dx}$ is constant between breakpoints

$$\# \text{ pieces} \leq n_1 + 1$$

### Example 2

$n_0 \geq 2, \ L = 1$



$$z^{(2)}(\vec{x}) = b^{(2)} + \sum_{j=1}^{n_1} \sigma\left(\vec{W}_j^{(1)} \cdot \vec{x} + b_j^{(1)}\right)$$

For each $j = 1, \ldots, n_1$, define

$$H_{j,\pm}^{(1)} = \left\{ x \in \mathbb{R}^{n_0} \mid \text{sgn}\left(\vec{W}_j^{(1)} \cdot \vec{x} + b_j^{(1)}\right) = \pm 1 \right\}$$

In $\mathbb{R}^2$, this makes a planar subdivision:



$\uparrow$ = direction at an

Note that in each component of $\mathbb{R}^{n_0} \setminus \bigcup_{j=1}^{\hat{n}_1} \partial H_{j,i}^{(1)}$ <span style="color:blue">(the cells of the hyperplane arrangement)</span> each neuron is either on or off. Thus, $\vec{\nabla}_x z^{(2)}(x)$ is constant on each cell of the hyperplane arrangement, and so

$$\begin{array}{ll} \# \text{ pieces} \leq & \# \text{ cells in arrangement of} \\ & n_1 \text{ hyperplanes in } \mathbb{R}^{n_0} \end{array} \leq \underbrace{\sum_{i=0}^{n_0} \binom{n_1}{i}}_{} = \begin{cases} n_1^{n_0}/n_0! & n_1 \ggg n_0 \\ 2^{n_1} & n_1 \leq n_0 \end{cases}$$

<span style="color:blue">Zaslavsky's Theorem, equality when in general position</span>

<span style="color:blue"># General Setting</span>

**Def:** A <span style="color:red">linear region</span> is a <u>maximal</u> $n_0$-dimensional connected set on which $\vec{\nabla}_x z^{L+1}(\vec{x})$ is constant.

<span style="color:blue">(Worst Case)</span>

<u>Lemma</u>: For any $n_0, L$, any $n_1, \dots, n_L$, $\quad \#$ linear regions $\leq 3^{\# \text{ neurons}}$ <span style="color:blue">← each neuron partitions space in 3 parts</span>

<u>Proof:</u> For each assignment of neuron on/offs

$$\vec{\varepsilon} = (\varepsilon_i^{(\ell)} \dots) \in \{-1, 0, 1\}^{[n_1] \times \dots \times [n_L]}$$

define $\quad P(\vec{\varepsilon}) = \{ x \in \mathbb{R}^{n_0} \mid \operatorname{sgn}(z_i^{(\ell)}(\vec{x})) = \varepsilon_i^{(\ell)} \}$.

Each $P(\vec{\varepsilon})$ is a region of input space with the same signs of preactivations, and so $\vec{\nabla}_x z^{(L+1)}(\vec{x})$ is constant on each $P(\vec{\varepsilon})$. They also partition $\mathbb{R}^{n_0}$ (disjoint union), i.e. $\mathbb{R}^{n_0} = \coprod_{\vec{\varepsilon}} P(\vec{\varepsilon})$.

We want to show that each $P(\vec{\varepsilon})$ is a connected set. In fact, we will show that each $P(\vec{\varepsilon})$ is a <u>convex polytope</u>!

Write $\quad P(\vec{\varepsilon}) = \bigcap_{\ell=1}^{L} P^{(\ell)}(\vec{\varepsilon}) \quad$ and $\quad P^{\ell}(\vec{\varepsilon}) = \bigcap_{i=1}^{n_\ell} P_i^{(\ell)}(\vec{\varepsilon})$, where

$$P_i^{(\ell)}(\vec{\varepsilon}) = \{ x \in \mathbb{R}^{n_0} \mid \operatorname{sgn}(z_i^{(\ell)}(\vec{x})) = \varepsilon_i^{(\ell)} \}$$

Note that $P_i^{(1)}(\vec{\varepsilon})$ is a convex polytope because each

$$P_i^{(1)} = \{ x \in \mathbb{R}^{n_0} \mid \operatorname{sgn}(\vec{w}_i^{(1)} \cdot \vec{x} + \vec{b}^{(1)}) = \varepsilon_i^{(1)} \}$$

is either a half-space or a hyperplane. So, $P^{(i)}(\hat{\varepsilon}) = \bigcap_{i=1}^{n} P_i^{(1)}(\hat{\varepsilon})$ is a convex polytope.

Next, note that on $P^{(1)}(\hat{\varepsilon})$, if $\dim(P^{(1)}(\hat{\varepsilon})) = n_0$, $\hat{\nabla}_x z^{(2)}(\hat{z})$ is constant. So, $P^{(1)}(\hat{\varepsilon}) \wedge P_i^{(2)}(\hat{\varepsilon})$ is the intersection of $P^{(1)}(\hat{\varepsilon})$ with a hyperplane or a half-space. Thus,

$$P^{(2)}(\hat{\varepsilon}) = \bigcap_{i=1}^{n_2} P^{(1)}(\hat{\varepsilon}) \wedge P_i^{(2)}(\hat{\varepsilon})$$

is a convex polytope. Repeat inductively to see that $P(\hat{\varepsilon})$ is a convex polytope, and is therefore connected. Since there are $3^{\# \text{neurons}}$ of possible $\hat{\varepsilon}$'s, each of which makes a new (possibly empty) region, the result follows.



$P^{(1)}(\hat{z})$ for some $\hat{z}$

second layer neuron draws a hyperplane that bends in each different cell

□

<u>Upshot:</u> $L \geq 2 \implies$ # pieces grows quickly because "bent hyperplane" wiggle

<u>Open Problems:</u>

* $\mathbb{P}\{\exists$ bounded bent hyperplane$\} = ?$     * for $n_0 = 2, L=1$   $\mathbb{E}\{$# sides of polygon containing the origin$\}$?

(Worst case — exponential in # neurons — can exist)

<u>Theorem A:</u> (Telgarsky)

Suppose $n_0 = 1$. Then, $\exists$ a ReLU net with large enough $L$ s.t.
  • depth = $2L$     • # neurons = $3L-1$     • # linear regions = $2^L$

<u>Proof:</u> Define $f(x) = \sigma(2\sigma(x) - 4\sigma(x-\frac{1}{2}))$
  Let $f_\ell(x) = f \circ f \circ \dots \circ f$
                   $\underbrace{\qquad}_{\ell \text{ times}}$

So, $f_\ell$ is a ReLU net with $L$ spikes and so $2 + 2^L$ regions.   □

(Avg. Case)

Theorem B: (Hanin-Rolnick)

Suppose $w_{ij}^{(l+1)} \sim N(0, \frac{2}{n_l})$, $b_i^{(l+1)} \sim N(0, C_b)$, $n_0 = 1$,

Then, $\mathbb{E}\{\# \text{ linear regions in } [a,b]\} \leq C \cdot |a-b| \cdot \# \text{ neurons}$

Proof idea: look up <span style="color:red">co-area formula</span> !



Open Problems:
* Count # of global regions (set $[a,b]$ to $\mathbb{R}$).

# Lecture 12/7- Bayesian Interpolation w/ Linear Nets

NNs have many large parameters:
- depth $L$
- width $n_\ell$
- input dim. $n_0$
- # train datapoints $P$

We want to ask how do $L, n_\ell, n_0, P$ influence "model quality",
i.e. feature learning, robustness, generalization, etc.

There are some **challenges** with any analysis:

① model is nonlinear in its parameters

② limits as $P, L, n_0, n_\ell \to \infty$ in different orders don't commute

# Examples of non-commuting limits

## Ex 1/ (Marchenko-Pastur)

Suppose $X \sim \mathbb{R}^{P \times n_0}$ with $X_{ij} \sim N(0,1)$ and

sample covariance determines what linear regressors do $\to$ $\Sigma_{n_0, P} = \frac{1}{n_0} X X^T \in \mathbb{R}^{P \times P}$

Since $\Sigma_{n_0, P}$ is PSD, write $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_P \geq 0$ as eigenvalues of $\Sigma_{n_0, P}$
and

$$\mu_{n_0, P} \equiv \frac{1}{P} \sum_{j=1}^{P} \delta_{\lambda_j} \quad \leftarrow \text{counting measure on eigenvalues}$$

**Theorem:** If $n_0, P \to \infty$ with $P/n_0 \to \alpha \in (0,1)$, then

$$\mu_{n_0, P} \xrightarrow[\text{converges in distribution weakly almost surely}]{w} \mu_{MP; \alpha} \quad, \text{ where}$$

where
$$d\mu_{MP;\alpha}(x) = \frac{1}{2\pi} \frac{\sqrt{(\lambda_+ - x)(x - \lambda_-)}}{\alpha x} \mathbb{1}_{[\lambda_-, \lambda_+]}(x), \qquad \lambda_\pm = (1 \pm \sqrt{\alpha})^2$$



## Ex 2/ Deep Linear Network

Consider $z(\vec{x}; \theta) = W^{(L+1)} \cdots W^{(1)} \vec{x} = \vec{\theta}^T \vec{x}$ where $W_{ij}^{(\ell)} \sim \mathcal{N}\left(0, \frac{1}{n_{\ell-1}}\right)$.

$\underbrace{\phantom{W^{(L+1)} \cdots W^{(1)}}}_{\in \mathbb{R}^{n_0}}$

We have $\vec{\theta} = \frac{\vec{\theta}}{\|\vec{\theta}\|} \|\vec{\theta}\|$, but $\frac{\vec{\theta}}{\|\vec{\theta}\|} \sim \text{Unif}(S^{n_0-1}) \perp \|\vec{\theta}\|$

$\underset{\text{uniform on sphere}}{\underbrace{\phantom{xxxxxx}}}$

Recall the following fact: if $W \in \mathbb{R}^{n \times m}$ has $W_{ij} \sim \mathcal{N}(0, \sigma^2)$,
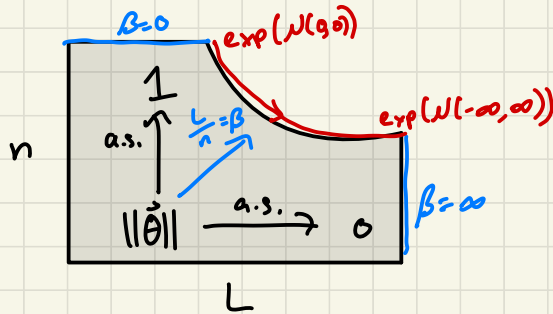then $W = UWV$ for $U \in O(n)$, $V \in O(m)$
(rotation/reflection invariant!)

So, $\|W^{(L+1)} \cdots W^{(1)}\| = \underbrace{\|W^{(L+1)}\|}_{\substack{\text{distributed as} \\ \text{chi-squared!}}} \left\| \frac{W^{(L+1)}}{\|W^{(L+1)}\|} W^{(L)} \cdots W^{(1)} \right\|$

$\overset{d}{=} \left(\frac{1}{n_L} \chi^2_{n_L}\right)^{\frac{1}{2}} \left\| \frac{W^{(L+1)}}{\|W^{(L+1)}\|} W^{(L)} \cdots W^{(1)} \right\|$

$\underset{\substack{W^{(L)} \text{ rotationally invariant,} \\ \text{can replace } \frac{W^{(L+1)}}{\|W^{(L+1)}\|} \\ \text{with } e_1}}{\phantom{xxxxxxxxxx}}$

$\overset{d}{=} \cdots \overset{d}{=} \left(\prod_{\ell=0}^{L} \frac{1}{n_\ell} \chi^2_{n_\ell}\right)^{\frac{1}{2}}$

$\underset{\substack{\text{indep. w/ mean 1} \\ \text{and var. } \frac{1}{n_\ell}}}{\underbrace{\phantom{xxxxxxxx}}}$

So, as $n \to \infty$, $\|\vec{\theta}\| \to 1$ almost surely.
However, we can also do
$$\|\vec{\theta}\| = \exp\left(\frac{1}{2} \sum_{\ell=1}^{L} \log\left(\frac{1}{n_\ell} \chi^2_{n_\ell}\right)\right)$$
$$\overset{n, L \gg 1}{\approx} \exp\left(\mathcal{N}\left(-\frac{L}{4n}, \frac{L}{4n}\right)\right)$$

So, as $L \to \infty$, $\|\dot{\theta}\| \to \exp(N(-\infty, \infty)) = 0$.
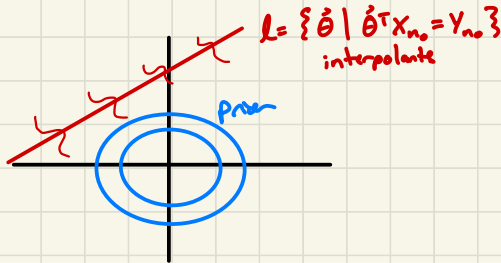This looks like the picture



# Bayesian Interpolation (Hanin + Alex Zlokapa)

__Model:__ $z(\vec{x}; \vec{\theta}) = W^{(L+1)} \dots W^{(1)} \vec{x} = \hat{\theta}^T \vec{x}$

__Data:__ $X_{n_0} = (\vec{x}_{j, n_0}) \in \mathbb{R}^{n_0 \times P}$, $\quad Y_{n_0} = (y_{j, n_0}) \in \mathbb{R}^{1 \times P}$

__Prior:__ $W_{ij}^{(\ell)} \sim N(0, \frac{\sigma^2}{n_{\ell-1}})$

__NLL:__ $\mathcal{L}_D(\hat{\theta}) = \frac{1}{2} \|\hat{\theta}^T X_{n_0} - Y_{n_0}\|_2^2$    likelihood $\propto \exp\left(-\frac{\beta}{2} \mathcal{L}_D(\hat{\theta})\right)$

$\ell = \{\hat{\theta} \mid \hat{\theta}^T X_{n_0} = Y_{n_0}\}$
interpolants



The Bayesian inference on each model is

$$d\mathbb{P}_{post}(\theta \mid X_{n_0}, Y_{n_0}, L, n_\ell, \sigma^2) = \lim_{\beta \to \infty} \frac{d\mathbb{P}_{prior}(\theta \mid L, n_\ell, \sigma^2) \times \exp\left(-\frac{\beta}{2} \mathcal{L}_D(\hat{\theta})\right)}{Z_\beta(X_{n_0}, Y_{n_0} \mid L, n_\ell, \sigma^2)}$$

Bayesian evidence

$$\propto \delta_\ell(\hat{\theta}) \mathbb{P}_{prior}(\hat{\theta})$$
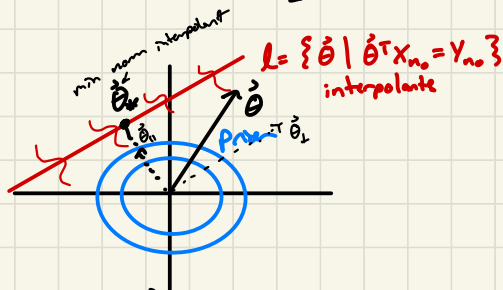
partition function, normalizes
distribution of posterior
$= \mathbb{P}\{\text{data} \mid \text{model}\}$

We can then perform **Bayesian model selection** to maximize $Z_\infty(X_{n_0}, Y_{n_0} \mid L, n_\ell, \sigma^2) \iff$ MLE on space of NNs

<u>maximize volume of interpolating models $\hat\theta$ over $L, n_\ell, \sigma^2$ architecture</u>

## <u>The <u>results</u> thidr</u>

* $P_{post}(\hat\theta)$, $Z_\infty$ are <u>exactly computable</u> (not asymptotically!)

* Effective depth $P \cdot \frac{L}{n}\left(= P \sum_{\ell=1}^{L} \frac{1}{n_\ell}\right) = \lambda_{post}$ determines posterior!

* $\lambda_{post} \to \infty \implies$ optimal feature learning from data-agnostic priors $(\theta^2 \geq 1)$

<u>Claim:</u> Any $\hat\theta$ can be decomposed into $\vec\theta_{\parallel} + \hat\theta_\perp$, where
$\hat\theta_{\parallel} \in Col(X_{n_0})$ and $\vec\theta_\perp \in Col(X_{n_0})^\perp$



We claim if $\hat\theta \sim P_{post}$, then
$$\vec\theta = \hat\theta_* + u \|\hat\theta_\perp\|, \qquad \hat\theta_{\parallel} = \hat\theta_*$$
where $u \sim Unif(S\, Col(X_{n_0})^\perp)$ independently of $\|\hat\theta_\perp\|$

<u>Interpolation</u> For a test point $\vec x$, $\vec x = \vec x_{\parallel} + \vec x_\perp$ by projection onto $X_{n_0}$. Then,
$$f(\vec x) = \hat\theta^T \vec x = \hat\theta_*^T \vec x_{\parallel} + (u \vec x_\perp)\|\hat\theta_\perp\|$$
$$\approx N\left(\hat\theta_*^T \vec x_{\parallel}, \frac{\|x_\perp\|^2}{n_0 - P}\|\theta_\perp\|^2\right)$$

this is how Bayesian inference can learn features!

So, $\|\hat\theta_\perp\|$ controls overall prediction scale in unseen directions!

**Theorem:** Suppose $n_0, P \to \infty$ with $P/n_0 \to \alpha \in (0,1)$ s.t.

$$\|\hat{\theta}_{*, n_0}\| \xrightarrow{d} \|\hat{\theta}_*\|$$

Then, $\theta_*^2 = \underset{\sigma^2}{\text{argmax}} \underset{\substack{n_0, P \to \infty \\ P/n_0 \to \alpha}}{\lim} Z_\infty(X_{n_0}, Y_{n_0} \mid L, n_\ell, \sigma^2)$  $\Big\}$ best data-dependent prior doesn't depend on architecture

gives $\underset{P, n_0 \to \infty}{\lim} \mathbb{P}_{post}(\|\hat{\theta}_*\| \mid \sigma^2 = \sigma_*^2, L, n_\ell) = \delta_{\frac{\|\hat{\theta}_*\|^2}{\alpha}}$

Furthermore,

$$\underset{\substack{P, n_0 \to \infty \\ P/n_0 \to \alpha}}{\lim} \mathbb{P}_{post}(\|\hat{\theta}_\perp\| \mid \sigma^2 = 1, L, n_\ell) = \delta_{\frac{\|\hat{\theta}_*^2\|}{\alpha} z(\lambda)}$$  $\Big\}$ in $\infty$ effective depth, we match best prior in data-agnostic way

As $\lambda \to \infty$, $z(\lambda) \to 1$.