


Introduction to Neural Networks

Def: (Fully Connected)

Fix $L \geq 1$, $n_0, \dots, n_{L+1} \geq 1$, $\sigma: \mathbb{R} \rightarrow \mathbb{R}$

$$\begin{aligned} x \in \mathbb{R}^{n_0} &\mapsto z^{(1)}(x) = W^{(1)}x + b^{(1)} \in \mathbb{R}^{n_1} \\ &\mapsto z^{(2)}(x) = W^{(2)}\sigma(z^{(1)}(x)) + b^{(2)} \in \mathbb{R}^{n_2} \\ &\vdots \\ &\mapsto z^{(L+1)}(x) = W^{(L+1)}\sigma(z^{(L)}(x)) + b^{(L+1)} \in \mathbb{R}^{n_{L+1}} \end{aligned}$$

$W_{ij}^{(l)} \sim$ "weights" $b_i^{(l)} \sim$ "biases"

Typical Use:

① Data Acquisition: $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^{n_{\text{data}}}$
 $y^{(i)} = f(x^{(i)})$, f "unknown"

② Model Selection: Fix L, n_l, σ .

③ Initialization: Choose

$$\Theta = \{W^{(l)}, b^{(l)}\} \text{ @ random}$$

④ Optimization: Do GD:

$$\Theta(t+1) = \Theta(t) - \eta_t \nabla_{\Theta} \mathcal{L}(\Theta(t))$$

$$\mathcal{L}(\Theta) = \frac{1}{n_{\text{data}}} \sum_{k=1}^{n_{\text{data}}} (y^{(k)} - z(x^{(k)}; \Theta))^2$$



Big Questions

1) Success of non-convex optimization: the loss $L(\theta)$ is "very" non-convex in θ but GD-type optimization typically finds global opt.

Key: NNs are overparameterized

$$\frac{\# \text{params}}{\# \text{data}} \gg 1$$

$$\Rightarrow \left\{ \nabla_{\theta} L(x_i^{\text{data}}) \right\}_{i=1}^{n_{\text{data}}} \subseteq \mathbb{R}^{\# \text{params}}$$

linearly indep

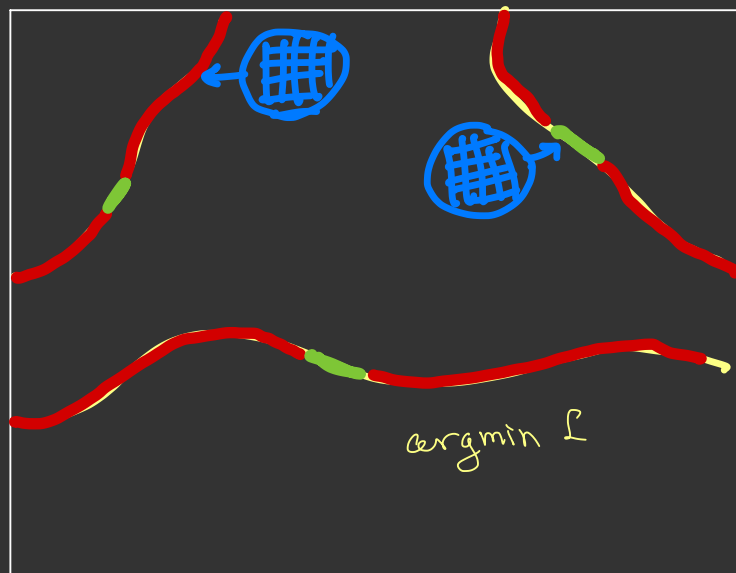
$$\Rightarrow \nabla_{\theta} \frac{1}{n_{\text{data}}} \sum_{i=1}^{n_{\text{data}}} L(x_i; \theta) = 0$$

$$\Leftrightarrow \nabla_{\theta} L(x_j; \theta) = 0$$

Neural Tangent Kernel

2) Implicit/Algorithmic Regularization

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



Q: How do diff opt methods influence which global min is found?

Random Neural Nets

Fix $L \geq 1$, $n_0, \dots, n_{L+1} \geq 1$, $\sigma: \mathbb{R} \rightarrow \mathbb{R}$
 $x \in \mathbb{R}^{n_0} \mapsto z^{(1)}(x) = W^{(1)} x \in \mathbb{R}^{n_1}$
 $\mapsto z^{(2)}(x) = W^{(2)}(z^{(1)}(x)) \in \mathbb{R}^{n_2}$
 \vdots
 $\mapsto z^{(L+1)}(x) = W^{(L+1)} \sigma(z^{(L)}(x)) \in \mathbb{R}^{n_{L+1}}$

with $W_{ij}^{(l)} \sim \mathcal{N}(0, \frac{C_w}{n_{l-1}})$

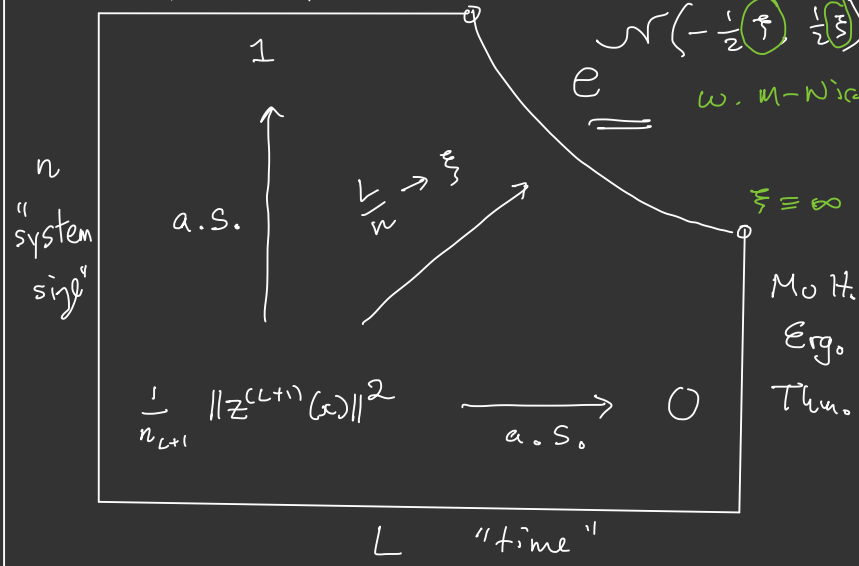
Goal: Fix n_0, n_{L+1}, σ . Describe $x \mapsto z^{(L+1)}(x)$ when $n_1, \dots, n_L \approx n \gg 1$

Metaclaim: When n, L large the distribution $z^{(L+1)}(x)$ depends on

$$\frac{L}{n} \left(= \frac{1}{n_1} + \dots + \frac{1}{n_L} \right)$$

and universality class of σ .

Intuition: $\sigma(t) = t$, $C_w = 1$
 $z^{(L+1)} = W^{(L+1)} \dots W^{(1)} x \propto \|x\|^2 = n_0$
Free Prob / NTK $\xi \equiv 0$



$$\begin{aligned} & \frac{1}{n_{L+1}} \|W^{(L+1)} \dots W^{(1)} x\|^2 \\ &= \frac{1}{n_{L+1}} \|W^{(L+1)} \dots W^{(2)} \frac{W^{(1)} x}{\|W^{(1)} x\|} \|^2 \|W^{(1)} x\|^2 \\ & \stackrel{d}{=} \prod_{l=1}^{L+1} \|W^{(l)} a_l\|^2 \approx \left(1 + o\left(\frac{1}{\sqrt{n}}\right)\right) \\ & \quad o\left(\sqrt{\frac{L}{n}}\right) \end{aligned}$$

Thm: Fix L , n_0, n_{L+1} . As $n_1, \dots, n_L \rightarrow \infty$

$$z^{(L+1)}(x) \in \mathbb{R}^{n_{L+1}} \longrightarrow GP(\underline{0}, \underline{K}^{(L+1)})$$

$$\lim_{n_1, \dots, n_L \rightarrow \infty} \text{Cov}(z_i^{(L+1)}(x), z_j^{(L+1)}(x'))$$

$$= \delta_{ij} K^{(L+1)}(x, x')$$

$$\underline{K}^{(L+1)}(x, x') = \underline{C}_W \underline{\mathbb{E}}_{K^{(L)}} [\sigma(z_1^{(L)}(x)) \sigma(z_1^{(L)}(x'))]$$

Ex: $\sigma(t) = \text{ReLU}(t) = t \mathbb{1}_{t > 0}$

$$\underline{K}^{(L+1)}(x, x') = C_W \underline{\mathbb{E}}_{K^{(L)}} [(z_1^{(L)}(x))^2 \mathbb{1}_{z_1^{(L)}(x) > 0}]$$

$$= \frac{C_W}{2} K^{(L)}(x, x')$$

$C_W = 2$: $K^{(L)}(x, x') \simeq L \textcircled{0}$

Ex: $\sigma(t) = \tanh(t)$

$C_W = 1$: $K^{(L)}(x, x') \simeq \frac{1}{L} \textcircled{1}$

All This: $\frac{L}{n} \rightarrow 0$

complexity

$$z^{(L+1)}(x) = W^{(L+1)} \sigma(W^{(L)} \sigma(\dots \sigma(W^{(2)} \sigma(W^{(1)} x) \dots))$$

Thm: (H) When $n_1, \dots, n_L \asymp n \gg 1$

$$\kappa(z_i^{(L+1)}(x_1), \dots, z_i^{(L+1)}(x_k))$$

$$= \begin{cases} 0, & k \text{ odd} \\ O\left(\frac{1}{n^{\frac{k}{2}-1}}\right), & k \text{ even} \end{cases}$$

$$\Rightarrow \underline{\kappa}_{2k}^{(L+1)}(x) \equiv \kappa(z_1^{(L+1)}(x), \dots, z_1^{(L+1)}(x))$$

$$= O_L\left(\frac{1}{n^{k-1}}\right) \stackrel{2k \text{ times}}{=}$$

Obtain recursions for $\kappa_{2k}^{(L+1)}(x)$ in terms of $\{K_{2j}^{(L)}(x), j \leq k\}$ and find: $k = 2, 3, 4$

$$\kappa_{2k}^{(L+1)}(x) = C_\sigma \frac{L^{k-1}}{n^{k-1}} + O(n^{-k})$$

$$= C_\sigma \left(\frac{L}{n}\right)^{k-1} + O(n^{-k})$$

Coro:

$$\text{Cov} \left(\left(z_1^{(L+1)}(x) \right)^2, \left(z_2^{(L+1)}(x) \right)^2 \right) = C_\sigma \cdot \frac{L}{n}$$

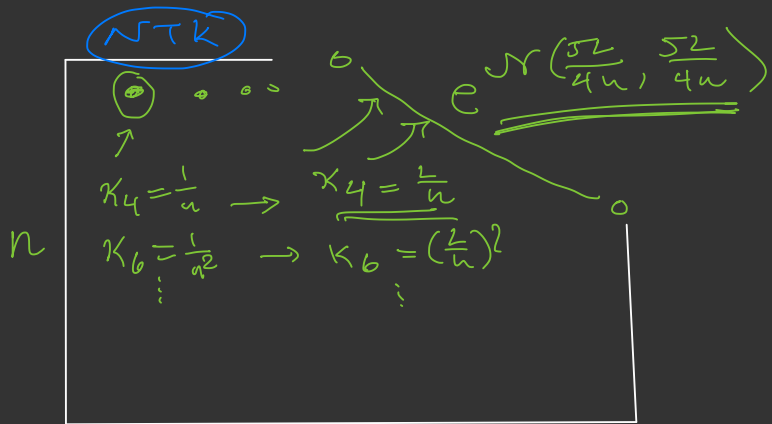
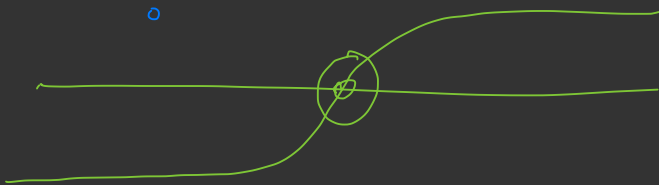
$$\text{Var} \left(\left\| \frac{\partial z_1^{(L+1)}(x)}{\partial x} \right\|^2 \right) = C_\sigma \cdot \frac{L}{n}$$

\Leftrightarrow (EVGP happens in FC iff $\frac{L}{n} \gg 1$)

$$\| \text{Hess}_\theta z_1^{(L+1)}(x; \theta) \| \simeq C_\sigma \cdot \frac{L}{n}$$

\Leftrightarrow $\left(\frac{L}{n} \gg 0 \right)$ necessary in FC for feature learning

⋮



make $\frac{L}{n}$ small

make $\frac{L}{n}$ large

$$\sigma(t) = \frac{(a_- \mathbb{1}_{t < 0} + a_+ \mathbb{1}_{t > 0})t}{C_w} \quad (a_- - a_+ = c)$$

$$C_w = \frac{2}{a_-^2 + a_+^2}$$

$$K^{(L)}(x, x) = L^{-\alpha}$$

$$\partial_x \partial_{x'} \Big|_{\underline{x}=\underline{x}'} K^{(L)}(x, x') = L^{-\beta}$$

X

