# Lecture 1: Introduction, criticality in linear networks

Based on Roberts and Yaida, *Principles of Deep Learning Theory*, 2106.10165

Why study neural networks? Flexible, differentiable class of functions with which to perform tasks like regression.

$$\text{Data} \longrightarrow \text{affine transformation} \longrightarrow \text{nonlinearity} \longrightarrow \text{output}$$
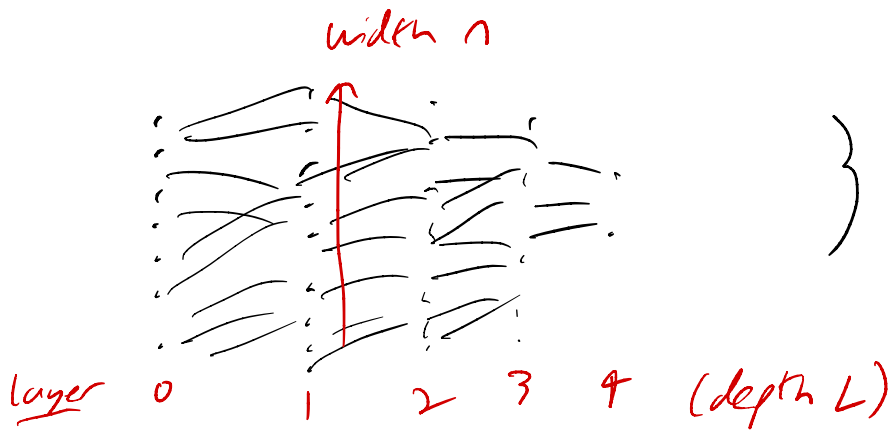
repeat many times

In equations: $z_i^{(1)} = b_i^{(1)} + W_{ij}^{(1)} x_j$

     ↑ biases    ↑ weights    ↑ data

$$z_i^{(\ell+1)} = b_i^{(\ell+1)} + W_{ij}^{(\ell+1)} \sigma(z_j^{(\ell)})$$

preactivation

(nonlinear) activation function, applied element-wise

width $n$



layer   0    1    2   3   4   (depth $L$)

} output is a function $f(x; \theta)$ parameterized by $\theta = \{ b^{(\ell)}, W^{(\ell)} \}$

In supervised learning, neural networks are trained by comparing the NN output to a desired output and updating the parameters by some form of gradient descent.

Key fact: these networks are massively overparameterized.
A standard benchmark dataset has 50k elements.
A standard "architecture" has width 256, depth 5 ⟹ 300k params.

How can you possibly avoid "overfitting" your data with so many parameters? This is one of the alluring mysteries of neural networks: somehow, they work!

3 physics analogies to keep in mind:

1. the initial choice of parameters is random, so NN's should be seen as elements of an ensemble. The statistics of this ensemble simplify as # of params $\longrightarrow \infty$, just like in stat mech.

2. the flow of information from input to output is like RG flow from UV to IR. The NN equations are recursions, so look for fixed points => criticality (lack of exponential behavior of correlation functions) (lectures 1-2)

3. there is an object called the NTK which acts like a Hamiltonian, governing updates of observables after one step of training. We can connect statistics at initialization to statistics at end of training if we can find a derivative of the NTK which is frozen during training (lecture 3)

- - - - - - - - - -

NN ensemble

In principle, 3 sources of randomness: initialization, data (drawn from some data distribution), training (e.g. stochastic gradient descent).
To simplify the analysis, we consider only initialization as random. Goal is to compute $p(f(x; \theta^*)|\mathcal{D})$, the distribution over trained networks where $\theta^*$ are the optimal parameters given training data $\mathcal{D}$.
We'll warm up by first computing $p(f(x; \theta)|\mathcal{D})$ before any training.

Formally, since the output is given by $z^{(L)}$ for a network of depth $L$, we have $p(f) = p(z^{(L)}|\mathscr{D}) = \int \prod_{\mu=1}^{P} d\theta_\mu \, p(\theta) \, p(z^{(L)}|\theta, \mathscr{D})$

<span style="color:red">↑ deterministic, given by iteration eqn.</span>

But because NN's have an iterative layer-to-layer structure, will be easiest to marginalize over layers one at a time:

$$p(z^{(\ell+1)}|\mathscr{D}) = \int \prod_i dz_i^{(\ell)} \, p(z^{(\ell+1)}|z^{(\ell)}) \, p(z^{(\ell)}|\mathscr{D}),$$

which is a recursion with initial condition

$$p(z^{(1)}|\mathscr{D}) = \int \prod d b_i^{(1)} \, dW_{ij}^{(1)} \, p(b^{(1)}) \, p(W^{(1)}) \, \delta(z_i^{(1)} - b_i^{(1)} - W_{ij}^{(1)} x_j)$$

For simple $p(b)$, $p(w)$ (i.e. Gaussian distributions) it turns out the marginal distributions, and hence $p(z)$, are <u>perturbatively</u> <u>Gaussian</u>, with non-Gaussianities scaling as $\frac{1}{n}$. Can therefore borrow all of the tools from stat. field theory to compute expectations!

— — — — — — — —

Toy model: linear network

Take $b_i^{(\ell)} = 0$, $\sigma(z) = z$. This is literally successive multiplication by a random matrix. Output is always a linear function of input, but a highly nonlinear function of parameters (matrix entries), so there is still rich structure at initialization that will carry over to the nonlinear case.

Anticipating the nearly-Gaussian statistics, instead of computing $p(z|\mathscr{D})$ directly, we will compute its first few connected moments.

Take our weight matrix entries to be i.i.d. Gaussian;

$$\mathbb{E}[W_{ij}] = 0, \quad \mathbb{E}[W_{i_1 j_1} W_{i_2 j_2}] = \frac{C_w}{n} \delta_{i_1 i_2} \delta_{j_1 j_2}$$

<span style="color:red">(for simplicity, same width $n$ and same weight distribution in each layer)</span>

Let $\mathcal{D} = \{x_{i;\alpha}\}$ where $\alpha$ is a sample index. Then

$$p(z^{(\ell)} | \mathcal{D}) = p\left( z^{(\ell)}_{i;\alpha_1}, z^{(\ell)}_{i;\alpha_2}, \dots z^{(\ell)}_{i;\alpha_{N_{\mathcal{D}}}} \right) \quad w \quad z_\alpha \equiv z(x_\alpha)$$

Easy things first: all odd moments <u>vanish</u>

e.g. $\mathbb{E}[z_\alpha^{(\ell)}] = \mathbb{E}[W^{(\ell)} W^{(\ell-1)} \cdots W^{(1)} x_\alpha]$ <span style="color:red">(dropping indices for clarity)</span>

$$= \mathbb{E}[W^{(\ell)}] \mathbb{E}[W^{(\ell-1)}] \cdots \mathbb{E}[W^{(1)}] x_\alpha \quad$$ <span style="color:red">(each layer's weights are independent)</span>

$$= 0 \quad$$ <span style="color:red">(weights are mean-zero)</span>

First nontrivial moment is

$$\mathbb{E}[z^{(1)}_{i_1;\alpha_1} z^{(1)}_{i_2;\alpha_2}] = \mathbb{E}[W^{(1)}_{i_1 j_1} x_{j_1;\alpha_1} W^{(1)}_{i_2 j_2} x_{j_2;\alpha_2}]$$

$$= \mathbb{E}[W^{(1)}_{i_1 j_1} W^{(1)}_{i_2 j_2}] x_{j_1;\alpha_1} x_{j_2;\alpha_2}$$

$$= \frac{C_w}{n} \delta_{i_1 i_2} \delta_{j_1 j_2} x_{j_1;\alpha_1} x_{j_2;\alpha_2}$$

$$= C_w \left( \underbrace{\frac{1}{n} \vec{x}_{\alpha_1} \cdot \vec{x}_{\alpha_2}}_{\color{red}{G^{(0)}_{\alpha_1 \alpha_2}}} \right) \delta_{i_1 i_2}$$

<span style="color:red">$G^{(0)}_{\alpha_1 \alpha_2} \equiv$ covariance btw. two inputs</span>

Note this is diagonal in neural indices: no covariance btw. neurons 1 and 2

Now write a recursion for 2-point correlator in layer $\ell$, with the ansatz $\mathbb{E}[z^{(\ell)}_{i_1;\alpha_1} z^{(\ell)}_{i_2;\alpha_2}] = G^{(\ell)}_{\alpha_1 \alpha_2} \delta_{i_1 i_2}$. Our initial condition is

$$G^{(1)}_{\alpha_1 \alpha_2} = C_w G^{(0)}_{\alpha_1 \alpha_2}.$$

Continuing, $\mathbb{E}[z_{i_1;\alpha_1}^{(\ell+1)} z_{i_2;\alpha_2}^{(\ell+1)}] = \mathbb{E}[W_{i_1 j_1}^{(\ell+1)} z_{j_1;\alpha_1}^{(\ell)}, W_{i_2 j_2}^{(\ell+1)} z_{j_2;\alpha_2}^{(\ell)}]$

**Key point:** $z^{(\ell)}$ only depends on $W^{(\ell)}, W^{(\ell-1)}, \ldots$, so is statistically independent from deeper layers, including $W^{(\ell+1)}$

$$= \mathbb{E}[W_{i_1 j_1}^{(\ell+1)} W_{i_2 j_2}^{(\ell+1)}] \mathbb{E}[z_{j_1;\alpha_1}^{(\ell)}, z_{j_2;\alpha_2}^{(\ell)}]$$

$$= \frac{C_W}{n} \delta_{i_1 i_2} \delta_{j_1 j_2} \underbrace{G_{\alpha_1 \alpha_2}^{(\ell)} \delta_{j_1 j_2}}_{} \qquad \text{(ansatz)}$$

$Tr(\mathbb{1}_{n\times n}) = n$

$$= C_W G_{\alpha_1 \alpha_2}^{(\ell)} \delta_{i_1 i_2}$$

So our ansatz is consistent, and the 2-point recursion is

$$G_{\alpha_1 \alpha_2}^{(\ell+1)} = C_W G_{\alpha_1 \alpha_2}^{(\ell)} \qquad (\ell = 0, 1, \ldots)$$

Solution: $G_{\alpha_1 \alpha_2}^{(\ell)} = C_W^{\ell} G_{\alpha_1 \alpha_2}^{(0)}$

We can also sum over $i_1, i_2$ in our ansatz to find $G_{\alpha_1 \alpha_2}^{(\ell)} = \frac{1}{n} \mathbb{E}[\vec{z}_{\alpha_1} \cdot \vec{z}_{\alpha_2}]$

$\Rightarrow G_{\alpha_1 \alpha_2}^{(\ell)}$ is a covariance at layer $\ell$, and blows up exponentially with depth if $C_W > 1$, or shrinks exponentially if $C_W < 1$.

Can tune the network to <u>criticality</u> by choosing $\boxed{C_W = 1}$, in which case $G_{\alpha_1 \alpha_2}^{\circ} \equiv G_{\alpha_1 \alpha_2}^{(0)}$ is a <u>fixed point</u> of the $G$ recursion. Input covariance is preserved during propagation through the network.

(Very) surprising fact about NN's: this tuning is sufficient to prevent exponential behavior of <u>all</u> higher-point correlators, and hence the full $p(z^{(\ell)} | \mathcal{D})$

4-point recursion (briefly!): start at layer 1,
set $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 \equiv \alpha$

$$\mathbb{E}[z_{i_1}^{(1)} z_{i_2}^{(1)} z_{i_3}^{(1)} z_{i_4}^{(1)}] = \mathbb{E}[W_{i_1 j_1} W_{i_2 j_2} W_{i_3 j_3} W_{i_4 j_4}] x_{j_1} x_{j_2} x_{j_3} x_{j_4}$$

<span style="color:red">Wick's Theorem gives this in terms of $C_W$:</span>

$$\frac{C_W^2}{n^2}\left(\delta_{i_1 i_2}\delta_{j_1 j_2}\delta_{i_3 i_4}\delta_{j_3 j_4} + (2 \text{ more contractions})\right)$$

$$= C_W^2 \left(\delta_{i_1 i_2}\delta_{i_3 i_4} + \delta_{i_1 i_3}\delta_{i_2 i_4} + \delta_{i_1 i_4}\delta_{i_2 i_3}\right)(G_2^{(0)})^2$$

<span style="color:red">$\equiv G_{\alpha\alpha}^{(0)}$</span>

This is precisely the structure of a Wick contraction, so can subtract off 2-point correlators to find

$$\mathbb{E}[z_{i_1}^{(1)} z_{i_2}^{(1)} z_{i_3}^{(1)} z_{i_4}^{(1)}]_{\text{conn.}} = 0 \quad \textcolor{red}{\Leftarrow \text{ Gaussian up to 4th moment, in \underline{first} layer}}$$

However, non-Gaussianities get generated in deeper layers. Taking the ansatz $\mathbb{E}[z_{i_1}^{(\ell)} z_{i_2}^{(\ell)} z_{i_3}^{(\ell)} z_{i_4}^{(\ell)}] = G_4^{(\ell)}\left(\delta_{i_1 i_2}\delta_{i_3 i_4} + (2 \text{ sim})\right)$

compute $\mathbb{E}[z^{(\ell+1)} z^{(\ell+1)} z^{(\ell+1)} z^{(\ell+1)}]$ to derive the recursion

$G_4^{(\ell+1)} = C_W^2\left(1 + \frac{2}{n}\right) G_4^{(\ell)}$. Has a closed-form solution, but instead let's expand perturbatively in $\frac{1}{n}$:

$$\mathbb{E}[z^{4\,(\ell)}]_{\text{conn.}} \propto G_4^{(\ell)} - (G_2^{(\ell)})^2 = \frac{2(\ell-1)}{n}(G_2^{(\ell)})^2$$

At criticality ($C_W = 1$), connected 4-pt. grows linearly with depth:
$$\mathbb{E}[z^{4(\ell)}]_{\text{conn.}} \propto \frac{2\ell}{n}(G_2^{(\ell)})^2 \quad \textcolor{red}{\Leftarrow \text{ marginally relevant}}$$

Some final comments:

- as $n \to \infty$, $4^{th}$ and higher cumulants vanish: $p(z^{(l)}|\partial)$ is purely Gaussian.

- as $L \to \infty$, combinatorial factors eventually grow and spoil criticality. "Infinite size" limit is $n \to \infty$, not $L \to \infty$

- for Gaussian inits, $\frac{L}{n}$ appears as the cutoff of the effective theory. (Not true for other inits in general!)

- there is a closed-form expression for single-input $p(z|x_\alpha)$:

$$p(z^{(l)}|x) \propto G_{0,1}^{l,0}\left(\frac{\vec{z}^{(l)} \cdot \vec{z}^{(l)}}{2^l n C_W^l G_2^{(0)}}\bigg|_{0,0,\ldots 0} \overline{\phantom{---}}\right) \quad [2104.11734]$$

<span style="color:red">↑ depth dependence encoded in Meijer G-function</span>

- if we take our weight matrices to be <u>orthogonal</u>, distributed under the Haar measure on $O(n)$

$$p(z^{(l)}|x) \propto \delta\left(\vec{z}^{(l)} \cdot z^{(l)} - \vec{x} \cdot \vec{x}\right) \quad [\text{YK, Hannah Day, D. Roberts}]$$

Duh. $O(n)$ preserves norm, so just a random rotation on the $(n-1)$-sphere. But expand perturbatively:

$$\mathbb{E}[z^{4(l)}]_{conn.} = -\frac{2}{n}(G_2^{(0)})^2 \quad$$ <span style="color:red">independent of $l$: no $\frac{L}{n}$ cutoff! Exactly marginal.</span>