# Lecture 2: Universality Classes of Nonlinear Networks

Let's use the same strategy as last time to compute moments for a __nonlinear__ network. New things:

- allow for Gaussian biases, $\mathbb{E}[b_i] = 0$, $\mathbb{E}[b_i b_j] = C_b \delta_{ij}$
- two common examples for $\sigma(z)$ are
  - ReLU ("rectified linear unit") $\sigma(z) = \max(0, z)$
  - $\sigma(z) = \tanh(z)$

<span style="color:red">Will see that these belong to __different__ universality classes of network behavior</span>

Recall the NN forward equations:

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

We start by computing the 2-point function:

$$\mathbb{E}\left[ z_{i_1;\alpha_1}^{(1)} z_{i_2;\alpha_2}^{(1)} \right] = \mathbb{E}\left[ \left( b_{i_1}^{(1)} + W_{i_1 j_1}^{(1)} x_{j_1;\alpha_1} \right)\left( b_{i_2}^{(1)} + W_{i_2 j_2}^{(1)} x_{j_2;\alpha_2} \right) \right]$$

<span style="color:red">[cross-terms vanish since $W$ and $b$ are independent]</span>

$$= \left( C_b + \frac{C_w}{n} \, \vec{x}_{\alpha_1} \cdot \vec{x}_{\alpha_2} \right) \delta_{i_1 i_2}$$

<span style="color:red">$G_{\alpha_1 \alpha_2}^{(1)}$ : call this the __first-layer metric__</span>

A similar calculation yields $\mathbb{E}\left[ z^{(1)} z^{(1)} z^{(1)} z^{(1)} \right]_{con.} = 0$ : first-layer distribution $p(z^{(1)} | \mathcal{D})$ is Gaussian to $\mathcal{O}(\frac{1}{n})$. If we wanted, we could write this distribution as an __action__,

$$p(z^{(1)} | \mathcal{D}) = \frac{1}{[\det(2\pi G^{(1)})]^{n/2}} \exp\left( -\frac{1}{2} G_{(1)}^{\alpha_1 \alpha_2} \, \vec{z}_{\alpha_1}^{(1)} \cdot \vec{z}_{\alpha_2}^{(1)} \right)$$

<span style="color:red">inverse metric</span>

which is correct to $\mathcal{O}(\frac{1}{n})$.

An identical computation gives the layer-to-layer marginal distribution with which we build the recursion,

$$p(z^{(\ell+1)} \mid z^{(\ell)}) = \frac{1}{\sqrt{\det(2\pi \hat{G}^{(\ell+1)})^n}} \exp\left(-\frac{1}{2} \hat{G}^{\alpha_1 \alpha_2}_{(\ell+1)} \vec{z}^{(\ell+1)}_{\alpha_1} \cdot \vec{z}^{(\ell+1)}_{\alpha_2}\right)$$

where $\hat{G}^{(\ell+1)}_{\alpha_1 \alpha_2} = C_b + \frac{C_w}{n} \vec{\sigma}^{(\ell)}_{\alpha_1} \cdot \vec{\sigma}^{(\ell)}_{\alpha_2}$ is a <u>stochastic variable</u>

$$\color{red}{\sigma(\vec{z}_{\alpha_1})}$$

Since it depends on the random variables $z^{(\ell)}, z^{(\ell-1)}, \ldots z^{(1)}$

<span style="color:red">In other words, marginal distributions are Gaussian with a <u>stochastic covariance</u> matrix, which results in accumulated non-Gaussianities in $p(z^{(\ell+1)} \mid \vartheta)$</span>

To determine the recursion we integrate out the $\ell^{th}$ layer preactivations; this is analogous to one step of RG flow.

$$p(z^{(\ell+1)} \mid \vartheta) = \int \prod_i dz^{(\ell)}_i \; p(z^{(\ell+1)} \mid z^{(\ell)}) p(z^{(\ell)} \mid \vartheta)$$

We can now feel free to build up $p(z^{(\ell+1)})$ from its moments.

Starting with $\ell = 2$: $\mathbb{E}[z^{(2)}_{i_1, \alpha_1} z^{(2)}_{i_2, \alpha_2}] = \delta_{i_1 i_2} \mathbb{E}[\hat{G}^{(2)}_{\alpha_1 \alpha_2}] = \delta_{i_1 i_2}(C_b + C_w \mathbb{E}[\sigma^{(1)}_{\alpha_1} \sigma^{(1)}_{\alpha_2}])$

<span style="color:red">↑ expectation taken over $z^{(1)}$ and $z^{(2)}$</span>

<span style="color:red">↑ only over $z^{(1)}$</span>

<span style="color:red">↑ $n$ copies of same single-neuron expectation</span>

$$= \delta_{i_1 i_2}\left(C_b + C_w \langle \sigma^{(1)}_{\alpha_1} \sigma^{(1)}_{\alpha_2} \rangle_{G^{(1)}}\right)$$

<span style="color:red">↑ Gaussian expectation w.r.t. covariance $G^{(1)}$</span>

If we let $G^{(\ell)}_{\alpha_1 \alpha_2} \equiv \mathbb{E}[\hat{G}^{(\ell)}_{\alpha_1 \alpha_2}]$ and take the ansatz $\mathbb{E}[z^{(\ell)}_{i_1, \alpha_1} z^{(\ell)}_{i_2, \alpha_2}] = \delta_{i_1 i_2} G^{(\ell)}_{\alpha_1 \alpha_2}$, we have the recursion

$$\boxed{G^{(\ell+1)}_{\alpha_1 \alpha_2} = C_b + C_w \langle \sigma^{(\ell)}_{\alpha_1} \sigma^{(\ell)}_{\alpha_2} \rangle_{G^{(\ell)}} + \mathcal{O}\left(\frac{1}{n}\right)}$$

<span style="color:red">for $\ell > 2$, $\mathbb{E}[\sigma^{(\ell)} \sigma^{(\ell)}]$ is Gaussian to $\mathcal{O}\left(\frac{1}{n}\right)$</span>

We can derive criticality conditions in the $n \to \infty$ limit.  3

Let $\lim\limits_{n \to \infty} G = K$, so our recursion is $K_{\alpha\beta}^{(\ell+1)} = C_b + C_w \langle \sigma_\alpha \sigma_\beta \rangle_{K^{(\ell)}}$

If $\sigma$ is nonlinear, e.g. $\sigma(z) = z^2$, we have

$$K_{\alpha\beta}^{(\ell+1)} = C_b + C_w \langle z_\alpha^2 z_\beta^2 \rangle_{K^{(\ell)}} = C_b + C_w \left( K_{\alpha\alpha}^{(\ell)} K_{\beta\beta}^{(\ell)} + 2(K_{\alpha\beta}^{(\ell)})^2 \right)$$

<span style="color:red">Operator mixing under RG!</span>

So unlike in linear networks, we have to consider the whole $2 \times 2$ $K$ matrix, rather than just a single input.

The diagonal component is easy because it decouples:

$$K_{00}^{(\ell+1)} = C_b + C_w \langle \sigma^2 \rangle_{K^{(\ell)}} = C_b + C_w \left[ \frac{1}{\sqrt{2\pi K^{(\ell)}}} \int_{-\infty}^{\infty} dz \, e^{-\frac{z^2}{2K^{(\ell)}}} \sigma(z)^2 \right]$$

<span style="color:red">$g(K)$: single-variable Gaussian expectation</span>

Find a fixed point $K_{00}^*$ by linearizing: $K_{00}^{(\ell)} = K_{00}^* + \Delta K_{00}^{(\ell)}$

To first order in $\Delta$: $\Delta K_{00}^{(\ell+1)} = \chi_{||}(K_{00}^*) \times \Delta K_{00}^{(\ell)}$

<span style="color:red">$C_w g'(K) = \frac{C_w}{2K^2} \langle \sigma(z)\sigma(z)(z^2 - K) \rangle_K$</span>

To ensure we don't move away from the fixed point, our first criticality condition is $\boxed{\chi_{||}(K_{00}^*) = 1}$

Note: when $\sigma(z) = z$, $g'(K) = 1$, so we recover $C_w = 1$ from last time.

If $C_b \neq 0$, we have $K_{00}^{(\ell+1)} = C_b + K_{00}^{(\ell)} \Rightarrow K_{00}$ grows <u>linearly</u>

<span style="color:red">$\hookrightarrow$ Semi-critical: fixed point at infinity, so should rather choose $C_b = 0$.</span>

The constancy of $K_{00}$ ($\Leftrightarrow \chi_{||}=1$) ensures that norms of preactivations don't change exponentially. The off-diagonal components of $K$ measure changes to nearby inputs. Requiring that these also not change exponentially gives a second criticality condition.

Derivation is long but straightforward: find a decomposition of $K_{\alpha\beta}$ that diagonalizes the RG evolution, and linearize about a fixed point.

$\Rightarrow$ require $\chi_{\perp}(K^*)=1$, where $\chi_{\perp}(K) = C_w \langle \sigma'(z)^2 \rangle_K$

To summarize:

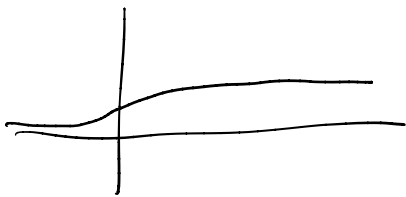criticality $\Rightarrow$ $\begin{array}{c} \chi_{||}=1 \\ \chi_{\perp}=1 \end{array}$ $\Rightarrow \dfrac{\chi_{\perp}}{\chi_{||}}=1$, solve for $C_b$ and $C_w$

$$\boxed{\left[\frac{2K^2 \langle \sigma'(z)^2 \rangle_K}{\langle \sigma(z)^2 (z^2-K) \rangle_K}\right]_{K=K^*} = 1, \qquad \begin{array}{l} C_w = \dfrac{1}{\langle \sigma'(z)^2 \rangle_{K^*}} \\[4mm] C_b = K^* - \dfrac{\langle \sigma(z)^2 \rangle_{K^*}}{\langle \sigma'(z)^2 \rangle_{K^*}} \end{array}}$$

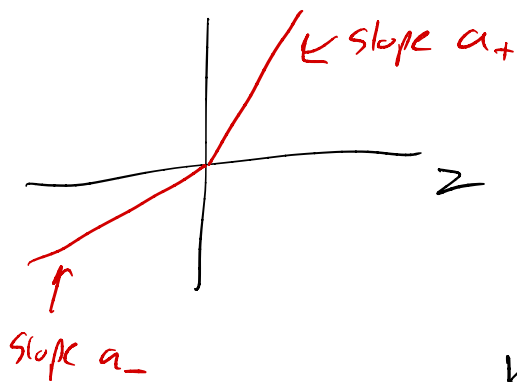<span style="color:red">Note: not every activation function has a consistent solution!</span>

e.g. $\sigma(z) = \dfrac{1}{1+e^{-z}}$



Used historically... but criticality implies $C_b = -\left(\dfrac{\sigma(0)^2}{\sigma'(0)^2}\right) < 0$ and a negative variance is unphysical. The problem is $\sigma(0) \neq 0$!

The criticality equations can be solved numerically, or by inspection. Different values of $K^*$ correspond to different universality classes.

---

- Scale-invariant (i.e. piecewise-linear w/ $\sigma(0)=0$)



$$\Rightarrow K^* = \frac{2}{a_+^2 + a_-^2}\left[\frac{1}{n}\vec{x}\cdot\vec{x}\right]$$

when $C_b = 0$ and $C_w = \frac{2}{a_+^2 + a_-^2}$,

$K^*$ is constant as a function of depth: a line of nontrivial fixed points corresponding to different input norms.

Linear activation has $a_+ = a_- = 1 \Rightarrow C_w = 1$, as we found before.

ReLU has $a_- = 0$, $a_+ = 1 \Rightarrow C_w = 2$, to compensate for the fact that $\max(0, z)$ knocks out half the activations on average.

- $K^* = 0$: consider a smooth activation function

$$\sigma(z) = \sum_{p=0}^{\infty} \frac{\sigma_p}{p!} z^p$$

$$\Rightarrow \langle \sigma^2(z) \rangle_K = \sigma_0^2 + (\sigma_1^2 + 2\sigma_0\sigma_2)K + \mathcal{O}(K^2)$$

$$\Rightarrow K^* = C_b + C_w \langle \sigma^2 \rangle_{K^*} \text{ has a solution, } K^* = 0, \text{ iff } \sigma_0 = 0$$

and $C_b = a$

Linearizing $x_{\parallel}$ and $x_{\perp}$ about $z=0$ and expanding in $K$, we find $C_w = \frac{1}{\sigma_1^2}$. For $\sigma(z) = \tanh z$, $\sigma_1 = 1$, and $C_w = 1$.

In the $K^* = 0$ universality class, $K$ decays to $0$, but only like a power law. Linearizing the recursions gives $K^{(l)} = \frac{1}{2l}$ for tanh (and $\# \times \frac{1}{l}$ for other activations), up to $\mathcal{O}(\frac{1}{l^2})$ corrections. In other words, $K$ is marginally <u>irrelevant</u>.

---

## Fluctuations

Using the same techniques as before, we can derive recursions for the 4-point correlator. Non-Gaussianity comes from both $\mathbb{E}[(\hat{G} - G)^2]$ <u>and</u> $\det(2\pi \hat{G})$, this is a long calculation.

For a single input, let's call the coefficient of the Wick tensor structure $V$ (meant to remind us of a 4-point vertex).

Recursion is $V^{(l)} = \chi_{\parallel}^2 (K^{(l)}) V^{(l)} + C_w^2 \left[ \langle \sigma^4 \rangle_{K^{(l)}} - \langle \sigma^2 \rangle_{K^{(l)}}^2 \right]$

Amazingly, tuning $C_w$ to criticality tames exponential behavior of $V^{(l)}$ too!

- Scale-invariant: $V^{(l)} = (l-1)(\#)(K^*)^2$ (marginally relevant)
- $K^* = 0$ : $V^{(l)} = (\#)(\frac{1}{l})$ (marginally irrelevant)

Note that we can assign a power-counting dimension to $z$. $[K] = 2$ and $[V] = 4$, so the dimensionless correlator is $\frac{V}{K^2} \overset{(\mathcal{O}(\frac{1}{n}))}{\sim} \frac{l}{n}$ for <u>both</u> universality classes!

This is the same linear growth of fluctuations we
saw in a linear network, but now we know it also
characterizes nonlinear networks.

Caveat: with orthogonal initializations, $\frac{V}{K^2}$ is
constant with $\ell$ for $K^a = 0$, but scales as $\frac{\ell}{n}$
for scale-invt., except linear activations which give

$\frac{V}{K^2}$ constant with depth.

$\Rightarrow K^a = 0$ is linearizing an activation function at
large depths, so nonlinear networks sort of behave linearly.