# Field Theories and Machine Learning
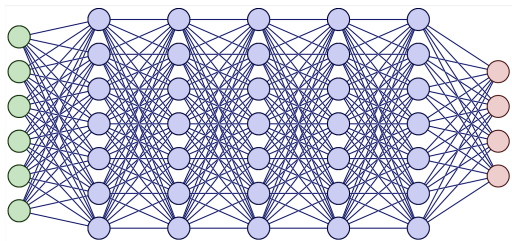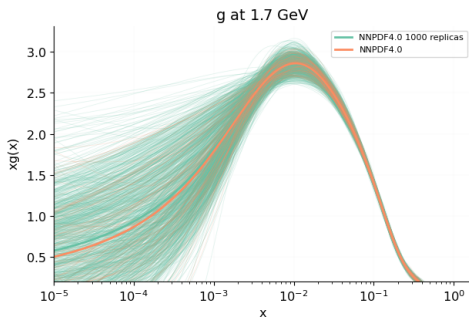


L Del Debbio
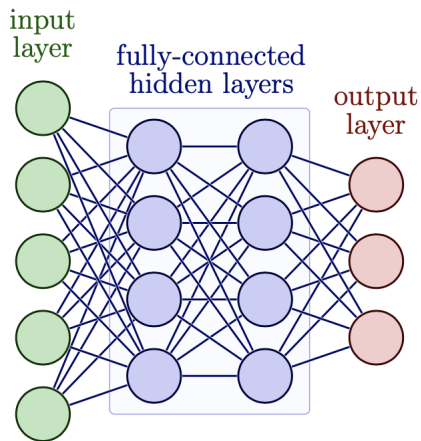
Higgs Centre for Theoretical Physics
University of Edinburgh

# NNPDF-like analysis

- generate $N_{\mathrm{rep}}$ replicas – NN initialized from random distributions

- NNs at initialization provide the prior, $f_i = N(x_i; \theta)$

- train the NNs using data

- NNs after training provide the posterior



g at 1.7 GeV

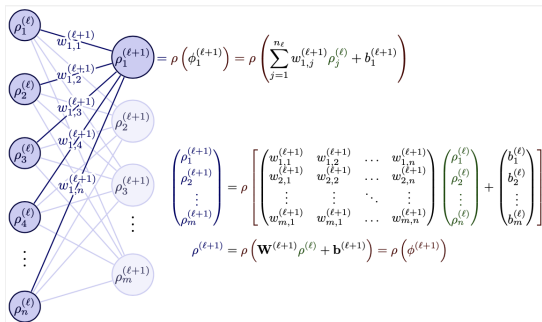# Neural Networks



MLP architecture

layers: $\ell = 1, \ldots, L$
neurons: $i = 1, \ldots, n_\ell$
weights $w_{ij}^{(\ell)}$, biases $b_i^{(\ell)}$
data: $(x_\alpha, y_\alpha), \alpha \in \mathcal{D}$

# pre-activation functions



$$\phi_{i\alpha}^{(\ell+1)} = \sum_{j=1}^{n_\ell} w_{ij}^{(\ell+1)} \rho_{j\alpha}^{(\ell)} + b_i^{(\ell+1)}$$

# statistical ensembles of NNs

initialize weights and biases using Gaussians

$$\langle b_i^{(\ell)} \rangle = 0 \,, \quad \langle b_{i_1}^{(\ell)} b_{i_2}^{(\ell)} \rangle = \delta_{i_1 i_2} C_b^{(\ell)}$$

$$\langle w_{ij}^{(\ell)} \rangle = 0 \,, \quad \langle w_{i_1 j_1}^{(\ell)} w_{i_2 j_2}^{(\ell)} \rangle = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_w^{(\ell)}}{n_{\ell-1}}$$

parameters/functions duality

$$p(\phi^{(\ell)}|\mathcal{D}) = \int [dw\, p(w)] \, [db\, p(b)] \prod_{i,\alpha} \delta \left( \phi_{i\alpha}^{(\ell)} - \sum_j w_{ij}^{(\ell)} \rho \left( \phi_{j\alpha}^{(\ell-1)} \right) - b_i^{(\ell)} \right)$$

computing the integral

$$p(\phi^{(\ell+1)}|\mathcal{D}) = \int d\phi^{(\ell)}\, p(\phi^{(\ell+1)}|\phi^{(\ell)})\, p(\phi^{(\ell)}|\mathcal{D})$$

$$
\begin{aligned}
p(\phi^{(\ell+1)}|\phi^{(\ell)}) &= \int \left[ dw^{(\ell+1)}\, p(w^{(\ell+1)}) \right] \left[ db^{(\ell+1)}\, p(b^{(\ell+1)}) \right] \\
&\quad \times \prod_{i,\alpha} \delta\big(\phi_{i\alpha}^{(\ell+1)} - \sum_j w_{ij}^{(\ell+1)} \rho\left(\phi_{j\alpha}^{(\ell)}\right) - b_i^{(\ell+1)}\big) \\
&= \frac{1}{|2\pi \widehat{G}^{(\ell+1)}|^{n_\ell/2}} \exp\left[ -\frac{1}{2} \left(\widehat{G}^{(\ell+1)}\right)^{-1}_{\alpha_1\alpha_2} \vec{\phi}_{\alpha_1}^{(\ell+1)} \cdot \vec{\phi}_{\alpha_2}^{(\ell+1)} \right]
\end{aligned}
$$

$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \ldots \phi_{i_{2k}\alpha_{2k}}^{(\ell+1)} \rangle = \sum_{\text{pairs}} \delta_{i_{P_1} i_{P_2}} \ldots \left\langle \left(\widehat{G}^{(\ell+1)}\right)_{\alpha_{P_1}\alpha_{P_2}} \ldots \right\rangle$$

# covariance

$$\widehat{G}^{(\ell+1)}_{\alpha_1\alpha_2} = C_b^{(\ell+1)} + \frac{C_w^{(\ell+1)}}{n_\ell} \vec{\rho}^{(\ell)}_{\alpha_1} \cdot \vec{\rho}^{(\ell)}_{\alpha_2}$$

fluctuations of $\widehat{G}$

$$\widehat{\Delta G}^{(\ell+1)}_{\alpha_1\alpha_2} = \widehat{G}^{(\ell+1)}_{\alpha_1\alpha_2} - \langle \widehat{G}^{(\ell+1)}_{\alpha_1\alpha_2} \rangle$$

$$\langle \widehat{\Delta G}^{(\ell+1)}_{\alpha_1\alpha_2} \widehat{\Delta G}^{(\ell+1)}_{\alpha_3\alpha_4} \rangle = \langle \widehat{G}^{(\ell+1)}_{\alpha_1\alpha_2} \widehat{G}^{(\ell+1)}_{\alpha_3\alpha_4} \rangle - \langle \widehat{G}^{(\ell+1)}_{\alpha_1\alpha_2} \rangle \langle \widehat{G}^{(\ell+1)}_{\alpha_3\alpha_4} \rangle$$

$$= \frac{C_w^2}{n_\ell^2} \left[ \langle \vec{\rho}^{(\ell+1)}_{\alpha_1} \cdot \vec{\rho}^{(\ell+1)}_{\alpha_2} \vec{\rho}^{(\ell+1)}_{\alpha_3} \cdot \vec{\rho}^{(\ell+1)}_{\alpha_4} \rangle - \right.$$

$$\left. - \langle \vec{\rho}^{(\ell+1)}_{\alpha_1} \cdot \vec{\rho}^{(\ell+1)}_{\alpha_2} \rangle \langle \vec{\rho}^{(\ell+1)}_{\alpha_3} \cdot \vec{\rho}^{(\ell+1)}_{\alpha_4} \rangle \right]$$

$$= \frac{1}{n_\ell} V^{(\ell+1)}_{\alpha_1\alpha_2,\alpha_3\alpha_4}$$

## propagators and all that

$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \phi_{i_2\alpha_2}^{(\ell+1)} \rangle = \delta_{i_1i_2} \langle \widehat{G}_{\alpha_1\alpha_2}^{(\ell+1)} \rangle = \delta_{i_1i_2} G_{\alpha_1\alpha_2}^{(\ell+1)}$$

$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \phi_{i_2\alpha_2}^{(\ell+1)} \phi_{i_3\alpha_3}^{(\ell+1)} \phi_{i_4\alpha_4}^{(\ell+1)} \rangle_c = \frac{1}{n_\ell} \left[ \delta_{i_1i_2} \delta_{i_3i_4} V_{\alpha_1\alpha_2,\alpha_3\alpha_4}^{(\ell+1)} + \dots \right]$$

$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \phi_{i_2\alpha_2}^{(\ell+1)} \phi_{i_3\alpha_3}^{(\ell+1)} \phi_{i_4\alpha_4}^{(\ell+1)} \phi_{i_5\alpha_5}^{(\ell+1)} \phi_{i_6\alpha_6}^{(\ell+1)} \rangle_c = \frac{1}{n_\ell^2} \times \dots$$

# $1/n$ expansion

correlators can be expanded in $1/n$

$$G_{\alpha_1\alpha_2}^{(\ell)} = K_{\alpha_1\alpha_2}^{(\ell)} + \frac{1}{n_{\ell-1}} G_{\alpha_1\alpha_2}^{\{1\}(\ell)} + \frac{1}{n_{\ell-1}^2} G_{\alpha_1\alpha_2}^{\{2\}(\ell)} + O(\frac{1}{n_{\ell-1}^3})$$

$$V_{\alpha_1\alpha_2,\alpha_3\alpha_4}^{(\ell)} = V_{\alpha_1\alpha_2,\alpha_3\alpha_4}^{\{0\}(\ell)} + \frac{1}{n_{\ell-1}} V_{\alpha_1\alpha_2,\alpha_3\alpha_4}^{\{1\}(\ell)} + O(\frac{1}{n_{\ell-1}^2})$$

therefore

$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \phi_{i_2\alpha_2}^{(\ell+1)} \rangle = \delta_{i_1 i_2} K_{\alpha_1\alpha_2}^{(\ell+1)} + O(1/n)$$
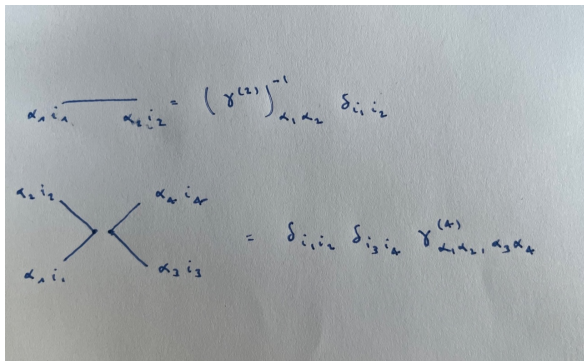
$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \phi_{i_2\alpha_2}^{(\ell+1)} \phi_{i_3\alpha_3}^{(\ell+1)} \phi_{i_4\alpha_4}^{(\ell+1)} \rangle_c = \frac{1}{n_\ell} \left[ \delta_{i_1 i_2} \delta_{i_3 i_4} V_{\alpha_1\alpha_2,\alpha_3\alpha_4}^{\{0\}(\ell+1)} + \ldots \right] + O(1/n^2)$$

# EFT

probability described by an effective action
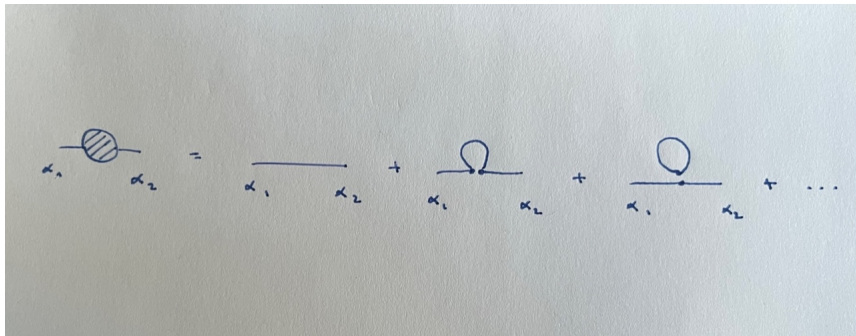
$$p(\phi|\mathcal{D}) = \frac{e^{-S(\phi)}}{Z}$$

$$S(\phi) = \frac{1}{2}\gamma^{(2)}_{\alpha_1\alpha_2} \vec{\phi}_{\alpha_1} \cdot \vec{\phi}_{\alpha_2} + \frac{1}{8}\gamma^{(4)}_{\alpha_1\alpha_2,\alpha_3\alpha_4} \vec{\phi}_{\alpha_1} \cdot \vec{\phi}_{\alpha_2} \, \vec{\phi}_{\alpha_3} \cdot \vec{\phi}_{\alpha_4} + \dots$$

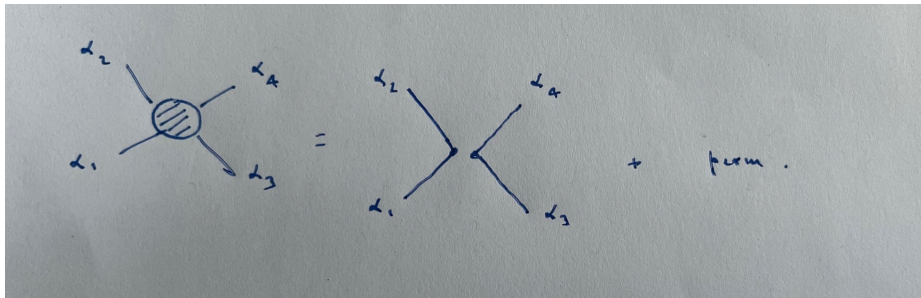# Two-point Function

couplings fixed by matching correlators in $1/n$ expansion

$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \phi_{i_2\alpha_2}^{(\ell+1)} \rangle = \delta_{i_1 i_2} K_{\alpha_1\alpha_2}^{(\ell+1)} + O(1/n)$$
$$= \delta_{i_1 i_2} \left( \gamma^{(2,(\ell+1))} \right)_{\alpha_1\alpha_2}^{-1} + O\left( \gamma^{(4,(\ell+1))} \right)$$

# Four-point Function

couplings fixed by matching correlators in $1/n$ expansion

$$\langle \phi_{i_1\alpha_1}^{(\ell+1)} \phi_{i_2\alpha_2}^{(\ell+1)} \phi_{i_3\alpha_3}^{(\ell+1)} \phi_{i_4\alpha_4}^{(\ell+1)} \rangle_c = \delta_{i_1 i_2} \delta_{i_3 i_4} \frac{1}{n_\ell} V_{\alpha_1\alpha_2,\alpha_3\alpha_4}^{\{0\}(\ell+1)} + \dots$$

$$= \delta_{i_1 i_2} \delta_{i_3 i_4} \left[ G_{\alpha_1\beta_1}^{(\ell+1)} G_{\alpha_2\beta_2}^{(\ell+1)} G_{\alpha_3\beta_3}^{(\ell+1)} G_{\alpha_4\beta_4}^{(\ell+1)} \right] \gamma_{\beta_1\beta_2,\beta_3\beta_4}^{(4,(\ell+1))} + \dots$$

## going deep – recursion relations

two-pt function at leading order

$$G_{\alpha_1\alpha_2}^{(\ell+1)} = C_b^{(\ell+1)} + C_w^{(\ell+1)} \frac{1}{n_\ell} \langle \vec{\rho}_{\alpha_1}^{(\ell)} \cdot \vec{\rho}_{\alpha_2}^{(\ell)} \rangle$$

$$K_{\alpha_1\alpha_2}^{(\ell+1)} = C_b^{(\ell+1)} + C_w^{(\ell+1)} \frac{1}{n_\ell} \langle \vec{\rho}_{\alpha_1}^{(\ell)} \cdot \vec{\rho}_{\alpha_2}^{(\ell)} \rangle \Big|_{O(1)}$$

$$= C_b^{(\ell+1)} + C_w^{(\ell+1)} \frac{1}{n_\ell} \langle \vec{\rho}_{\alpha_1}^{(\ell)} \cdot \vec{\rho}_{\alpha_2}^{(\ell)} \rangle_{K^{(\ell)}}$$

$$\frac{1}{n_\ell} \langle \vec{\rho}_{\alpha_1}^{(\ell)} \cdot \vec{\rho}_{\alpha_2}^{(\ell)} \rangle_{K^{(\ell)}} = \int \prod_\alpha d\phi_\alpha \frac{e^{-\frac{1}{2}\left(K^{(\ell)}\right)^{-1}_{\beta_1\beta_2} \phi_{\beta_1}\phi_{\beta_2}}}{\left|2\pi K^{(\ell)}\right|^{1/2}} \rho(\phi_{\alpha_1})\rho(\phi_{\alpha_2})$$

# evaluating the integral

if $\alpha_1 = \alpha_2$, only one input

$$\int \prod_\alpha d\phi_\alpha \frac{e^{-\frac{1}{2}\left(K^{(\ell)}\right)^{-1}_{\beta_1\beta_2}\phi_{\beta_1}\phi_{\beta_2}}}{\left|2\pi K^{(\ell)}\right|^{1/2}} \rho(\phi_{\alpha_1})\rho(\phi_{\alpha_1}) =$$

$$= \int d\phi_{\alpha_1} \frac{e^{-\frac{1}{2}\left(K^{(\ell)}_{\alpha_1\alpha_1}\right)^{-1}\phi^2_{\alpha_1}}}{\left(2\pi K^{(\ell)}_{\alpha_1\alpha_1}\right)^{1/2}} \rho(\phi_{\alpha_1})^2 =$$

$$= g(K^{(\ell)}_{\alpha_1\alpha_1})$$

solving the recursion

$$K_{\alpha\alpha}^{(\ell)} = \langle \frac{1}{n_\ell} \sum_{i=0}^{n_\ell} \left( \phi_{i,\alpha}^{(\ell)} \right)^2 \rangle$$

$K_{\alpha\alpha}$, introduce the auxiliary function

$$g(K) = \int d\phi \, \frac{1}{\sqrt{2\pi K}} \exp\left(-\frac{1}{2K}\phi^2\right)\rho(\phi)^2$$

fixed point: $K^* = C_b + C_w g(K^*)$
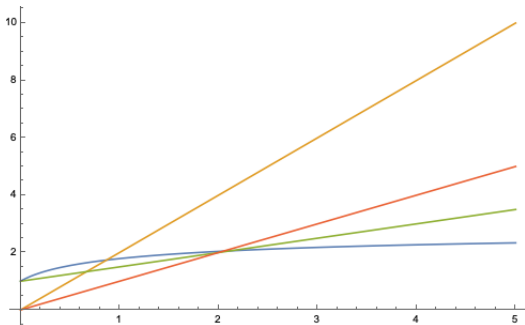
$$K_{\alpha\alpha} = K^* + \Delta K_{\alpha\alpha} \implies \Delta K_{\alpha\alpha}^{(\ell+1)} = \chi_\parallel(K^*)\Delta K_{\alpha\alpha}^{(\ell)}$$

$$\chi_\parallel(K) = C_w g'(K)$$

criticality: $\chi_\parallel(K^*) = 1$

# examples

graphical solution for $C_b = 1, C_w = 2$

## more recursions

for the off-diagonal elements of $K$: $\alpha_1 \neq \alpha_2$

$$K_{\alpha_1\alpha_2}^{(\ell+1)} = C_b^{(\ell+1)} + C_w^{(\ell+1)} \times$$
$$\times \int d\phi_{\alpha_1} d\phi_{\alpha_2} e^{-\frac{1}{2}(\phi_{\alpha_1}\phi_{\alpha_2})\bar{K}^{-1}\begin{pmatrix}\phi_{\alpha_1}\\\phi_{\alpha_2}\end{pmatrix}} \rho(\phi_{\alpha_1})\rho(\phi_{\alpha_2})$$

$$= H(K_{\alpha_1\alpha_1}^{(\ell)}, K_{\alpha_2\alpha_2}^{(\ell)}, K_{\alpha_1\alpha_2}^{(\ell)})$$

$$\bar{K} = \begin{pmatrix} K_{\alpha_1\alpha_1} & K_{\alpha_1\alpha_2} \\ K_{\alpha_2\alpha_1} & K_{\alpha_2\alpha_2} \end{pmatrix}$$

$\hookrightarrow$ coupled equations

RG-style evolution equation to go deep into the network

# RG interpretation

- distribution in each layer at initialization is described by an EFT

- for $n \to \infty$, NN defines a GP

- EFT couplings obey RG relations as we go *deep*

- taken into account in the choice of the architecture/initialization
  parameters $C_w$ and $C_b$

# Bayesian Learning with NN

Divide the data into a training set and the rest

$$\mathcal{D} = \mathcal{A} \cup \mathcal{B}$$

Bayes:

$$p(\theta|\mathcal{A}, \mathcal{H}) = \frac{p(\mathcal{A}|\theta, \mathcal{H})p(\theta|\mathcal{H})}{p(\mathcal{A}|\mathcal{H})}$$

likelihood:

$$p(A|\theta, \mathcal{H}) = \exp\left(-\mathcal{L}(\theta, \mathcal{A})\right)$$

estimators:

$$\theta^*_{\mathsf{MAP}} = \arg\max_\theta p(\theta|\mathcal{A}, \mathcal{H})$$
$$\theta^*_{\mathsf{MLE}} = \arg\max_\theta p(y_A|\theta, \mathcal{H})$$

# Bayesian Learning at $\infty$ Width

$$p(\phi_B^{(L)}|y_a) = \frac{p(y_A, \phi_B^{(L)})}{p(y_A)}$$

leading order in $1/n$

$$p(y_A, \phi_B^{(L)}) \propto \exp\left(-\frac{1}{2}(y_A, \phi^{(L)})^T K^{-1} \begin{pmatrix} y_A \\ \phi_B^{(L)} \end{pmatrix}\right)$$

where

$$K^{-1} = \begin{pmatrix} K_{\alpha\alpha}^{-1} & K_{\alpha\beta}^{-1} \\ K_{\beta\alpha}^{-1} & K_{\beta\beta}^{-1} \end{pmatrix}, \quad K_{\beta\beta}^{-1} = K_{\beta\beta} - K_{\beta\alpha}(K_{\alpha\alpha})^{-1}K_{\alpha\beta}$$

$$p(\phi_B^{(L)}|y_a) \propto \exp\left(-\frac{1}{2}(\phi^{(L)} - m_\beta^\infty)^T K_{\beta\beta}^{-1}(\phi^{(L)} - m_\beta^\infty)\right)$$

$$m_\beta^\infty = K_{\beta\alpha}(K_{\alpha\alpha})^{-1}y_\alpha$$

## Lack of Representation Learning

at infinite width

$$p(\phi_{\beta,1}^{(L)}, \ldots, \phi_{\beta,n_L}^{(L)} | y_A) = p(\phi_{\beta,1}^{(L)} | y_A) \ldots p(\phi_{\beta,n_L}^{(L)} | y_A)$$

different neurons on the output layer are statistically independent

$$p(\phi_D^{(L-1)} | y_A) = \frac{p(y_A | \phi_D^{(L-1)}) p(\phi_D^{(L-1)})}{p(y_A)}$$

$$p(y_A | \phi_D^{(L-1)}) \propto \exp\left(-\frac{1}{2} y^T (K_{\alpha\alpha})^{-1} y\right) = p(y_A)$$

$$p(\phi_D^{(L-1)} | y_A) = p(\phi_D^{(L-1)})$$

# Bayesian Learning at finite width

consider a single datapoint, quartic interactions introduce correlations

$$p(\phi_1, \ldots, \phi_m) \propto \exp\left(-\frac{\gamma^{(2)}}{2}\phi^T\phi + \frac{\gamma^{(4)}}{8}(\phi^T\phi)^2\right)$$

$$\frac{1}{\gamma^{(2)}} = G^{(\ell)} - \frac{m+2}{n_{\ell-1}}\frac{V^{(\ell)}}{G^{(\ell)}} + O\left(1/n^2\right)$$

$$\gamma^{(4)} = \frac{1}{n_{\ell-1}}\frac{V^{(\ell)}}{(G^{(\ell)})^4} + O\left(1/n^2\right)$$

conditional distribution for $\phi_2$

$$p(\phi_2|\check{\phi}_1) \propto \exp\left(-\frac{\gamma^{(2)}}{2}\phi_2^2 + \frac{\gamma^{(4)}}{8}(\phi_2^4 + 2\phi_2^2\check{\phi}_1^2)\right)$$

$$\int d\phi_2 \, p(\phi_2|\check{\phi}_1)\phi_2^2 = G^{(\ell)} + \frac{1}{2}\left(\check{\phi}_1^2 - G^{(\ell)}\right)\frac{V^{(\ell)}}{n_{\ell-1}(G^{(\ell)})^2} + O(1/n^2)$$

## training and NTK

gradient descent

$$\frac{d}{dt}\theta_\mu(t) = -\lambda_{\mu\nu}\frac{\partial}{\partial\theta_\mu}\mathcal{L}_A$$

evolution of $O(\phi)$

$$\frac{d}{dt}O(t) = -\frac{\partial O}{\partial\phi_{i\delta}}\frac{\partial\phi_{i\delta}}{\partial\theta_\mu}\lambda_{\mu\nu}\frac{\partial\mathcal{L}_A}{\partial\phi_{j\alpha}}\frac{\partial\phi_{j\alpha}}{\partial\theta_\nu}$$
$$= -\left.\frac{\partial O}{\partial\phi_{i\delta}}\right|_{\phi(t)} H_{i\delta,j\alpha}(t)\,\varepsilon_{j\alpha}(t)$$

in particular

$$\boxed{\frac{d}{dt}\phi_{i\delta}(t) = -H_{i\delta,j\alpha}(t)\varepsilon_{j\alpha}(t)}$$

## Forward Equation for the NTK

$$H^{(\ell)} = \lambda_{\mu\nu} \frac{\partial \phi^{(\ell)}}{\partial \theta_\mu} \frac{\partial \phi^{(\ell)}}{\partial \theta_\nu}$$

learning rate tensor

$$\lambda_{b_{i_1}^{(\ell)} b_{i_2}^{(\ell)}} = \delta_{i_1 i_2} \lambda_b^{(\ell)}, \quad \lambda_{W_{i_1 j_1}^{(\ell)} W_{i_2 j_2}^{(\ell)}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda_W^{(\ell)}}{n_{\ell-1}}$$

then

$$H_{i_1 \alpha_1, i_2 \alpha_2}^{(\ell)} = \sum_{\ell'=1}^{\ell} \left( \lambda_b^{(\ell')} \frac{\partial \phi_{i_1 \alpha_1}^{(\ell')}}{\partial b_j^{(\ell')}} \frac{\partial \phi_{i_2 \alpha_2}^{(\ell')}}{\partial b_j^{(\ell')}} + \frac{\lambda_w^{(\ell')}}{n_{\ell-1}} \frac{\partial \phi_{i_1 \alpha_1}^{(\ell')}}{\partial W_{jk}^{(\ell')}} \frac{\partial \phi_{i_2 \alpha_2}^{(\ell')}}{\partial W_{jk}^{(\ell')}} \right)$$

$$H_{i_1i_2\alpha_1\alpha_2}^{(\ell+1)} = \sum_{j=1}^{n_{\ell+1}} \left( \lambda_b^{(\ell+1)} \frac{\partial \phi_{i_1\alpha_1}^{(\ell+1)}}{\partial b_j^{(\ell+1)}} \frac{\partial \phi_{i_2\alpha_2}^{(\ell+1)}}{\partial b_j^{(\ell+1)}} + \frac{\lambda_w^{(\ell+1)}}{n_{\ell-1}} \frac{\partial \phi_{i_1\alpha_1}^{(\ell+1)}}{\partial W_{jk}^{(\ell+1)}} \frac{\partial \phi_{i_2\alpha_2}^{(\ell+1)}}{\partial W_{jk}^{(\ell+1)}} \right) +$$

$$+ \sum_{j_1,j_2=1}^{n_\ell} \frac{\partial \phi_{i_1\alpha_1}^{(\ell+1)}}{\partial \phi_{j_1\alpha_1}^{(\ell)}} \frac{\partial \phi_{i_2\alpha_2}^{(\ell+1)}}{\partial \phi_{j_2\alpha_2}^{(\ell)}} H_{j_1j_2\alpha_1\alpha_2}^{(\ell)}$$

$$= \delta_{i_1i_2} \left[ \lambda_b^{(\ell+1)} + \lambda_W^{(\ell+1)} \left( \frac{1}{n_\ell} \sum_{j=1}^{n_\ell} \rho_{j\alpha_1}^{(\ell)} \rho_{j\alpha_2}^{(\ell)} \right) \right] +$$

$$+ \sum_{j_1,j_2=1}^{n_\ell} W_{i_1j_1}^{(\ell+1)} W_{i_2j_2}^{(\ell+1)} \rho_{j_1\alpha_1}^{\prime(\ell)} \rho_{j_2\alpha_2}^{\prime(\ell)} H_{j_1j_2\alpha_1\alpha_2}^{(\ell)}$$

$1/n$ expansion

$$H = \Theta + \frac{1}{n} H^{[1]} + \dots$$

## Gradient Training at Infinite Width

at the minimum of $\mathcal{L}(\theta, A)$

$$\phi^{(L)}(T) = \phi^{(L)} - \Theta_{\beta\alpha}^{(L)}(\Theta_{\alpha\alpha}^{(L)})^{-1}\left(\phi^{(L)} - y\right)$$

hence

$$m^{\infty} = \Theta_{\beta\alpha}^{(L)}(\Theta_{\alpha\alpha}^{(L)})^{-1}y$$

recall

$$K_{\delta_1\delta_2}^{(L)} = C_b^{(L)} + C_W^{(L)}\langle\rho_{\delta_1}\rho_{\delta_2}\rangle_{K^{(L-1)}}$$
$$\Theta_{\delta_1\delta_2}^{(L)} = \lambda_b^{(L)} + \lambda_W^{(L)}\langle\rho_{\delta_1}\rho_{\delta_2}\rangle_{K^{(L-1)}} + C_W^{(L)}\langle\rho'_{\delta_1}\rho'_{\delta_2}\rangle_{K^{(L-1)}}\Theta_{\delta_1\delta_2}^{(L-1)}$$

$$\implies \lambda_b^{(L)} = C_b^{(L)}, \quad \lambda_W^{(L)} = C_W^{(L)}$$

# EFT again

introduce an auxiliary field $L_{i\alpha}(t)$

$$p(\varphi, L|\mathcal{D}) = \frac{1}{Z} \exp\left[-S(\phi) - \int dt\, L_{i\alpha}(t) \left(\frac{d}{dt}\varphi_{i\delta}(t) + H_{i\delta,j\alpha}(t)\varepsilon_{j\alpha}(t)\right)\right]$$

with $\varphi_{i\alpha}(0) = \phi_{i\alpha}$

EOM:

$$\frac{d}{dt}\phi_{i\delta}(t) = -H_{i\delta,j\alpha}(t)\varepsilon_{j\alpha}(t)$$

hence

$$\frac{d}{dt}\mathcal{L}_A = -\varepsilon_{i\delta}H_{i\delta,j\alpha}(t)\varepsilon_{j\alpha}(t)$$

use this theory to compute correlators during training

# Summary

- distribution of $\phi_{i\alpha}$ can be described by EFT

- power counting scheme in $1/n$ to have a predictive framework

- NN at infinite width $\longrightarrow$ made contact with GP

- $1/n$ corrections to go beyond Gaussianity

- deep networks and RG equations

- training and NTK analysis