

4. Training & Posterior Distribution

4.1 NNGPs

At infinite width, NNs define a class of GPs

$$\text{prior covariance: } k_{\delta_1, \delta_2} = \lim_{n \rightarrow \infty} E \left[\phi_{\delta_1}^{(L)} \phi_{\delta_2}^{(L)} \right].$$

Bayesian inference then yields the posterior distrib., as seen in L2.

↳ no fitting is required.

4.2 Training

Introduce training time t , $t=0$ initialization

Loss fun. $\alpha \in \mathcal{A}$, indices of the training set, $\mathcal{A} \subset \mathcal{D}$

$$\text{MSE: } \mathcal{L}_{\mathcal{A}} = \sum_{\alpha} \left(\phi_{i_{\alpha}}^{(L)} - y_{i_{\alpha}} \right)^2, \quad \alpha \in \mathcal{A}.$$

$$\text{Gradient Descent: } \frac{d}{dt} \theta_{\mu}(t) = - \partial_{\mu} \mathcal{L}_{\mathcal{A}}$$

sometimes, it can be generalized $\leftarrow - d_{\mu} \partial_{\nu} \mathcal{L}_{\mathcal{A}}$.

Using the chain rule:

$$\partial_{\mu} \mathcal{L}_{\mathcal{A}} = \sum_{i_{\alpha}} \frac{\partial \mathcal{L}_{\mathcal{A}}}{\partial \phi_{i_{\alpha}}^{(L)}(t)} \partial_{\mu} \phi_{i_{\alpha}}^{(L)}(t)$$

$$\Rightarrow \frac{d}{dt} \theta_{\mu} = - \sum_{i_{\alpha}} \lambda_{\mu \nu} \partial_{\nu} \phi_{i_{\alpha}}^{(L)}(t) \varepsilon_{i_{\alpha}}(t)$$

$$\varepsilon_{i_{\alpha}}(t) = \frac{\partial \mathcal{L}_{\mathcal{A}}}{\partial \phi_{i_{\alpha}}^{(L)}(t)} = 2 \left(\phi_{i_{\alpha}}^{(L)} - y_{i_{\alpha}} \right) \text{ for MSE.}$$

Flow eq. for the fields $\phi^{(L)}$.

$$\begin{aligned} \frac{d}{dt} \phi_{i\alpha}^{(L)}(t) &= \sum_{\mu} \partial_{\mu} \phi_{i\alpha}^{(L)} \frac{d}{dt} \theta_{\mu} \\ &= - \sum_{j\alpha} \left\{ \sum_{\mu, \nu} \lambda_{\mu\nu} \partial_{\mu} \phi_{i\alpha}^{(L)} \partial_{\nu} \phi_{j\alpha}^{(L)} \right\} \epsilon_{j\alpha} \\ &= - H_{i\alpha, j\alpha}(t) \epsilon_{j\alpha}(t) \quad \text{omitted the sum} \\ &\quad \uparrow \text{ } \alpha \in d! \\ \Rightarrow \frac{d \phi^{(L)}(t)}{dt} &= - H(t) E(t) \quad \text{omitted all indices.} \\ &\quad \uparrow \\ &\quad \text{Neural Tangent Kernel.} \end{aligned}$$

Exercise: Show that

$$\frac{d}{dt} O(\phi^{(L)}) = - \frac{\partial O}{\partial \phi^{(L)}} H E$$

(put back all indices & summations!).

4.3 $1/n$ Expansion

$$H_{i\alpha, j\alpha} = \delta_{ij} \theta_{\delta\alpha} + O(1/n)$$

\uparrow time-independent.

Correlations btw neurons only appear at $O(1/n)$!

Evolution is better understood by considering a series of fields

$$O_s(t) = O_{\delta_1, \dots, \delta_r}(t)$$

$$\text{such that: } \begin{cases} O_s(t) = \partial_{\mu} O_{s-1}(t) \partial_{\mu} \phi_{\delta_s}^{(L)}(t) \\ O_1(t) = \phi_{\delta_1}^{(L)}(t) \end{cases}$$

Then:

$$\begin{aligned}\frac{d}{dt} O_S(t) &= \partial_\mu O_S(t) \frac{d}{dt} \theta_\mu(t) \\ &= - \partial_\mu O_S(t) \partial_\mu \phi_S^{(1)}(t) \\ &= - O_{S+1}(t) \epsilon(t)\end{aligned}$$

Check:

$$\begin{aligned}O_2(t) &= O_{S_1 S_2}(t) = \partial_\mu O_{S_1}(t) \partial_\mu \phi_{S_2}^{(1)}(t) \\ &= H_{S_1 S_2}(t)\end{aligned}$$

$$\frac{d\phi^{(1)}}{dt} = -H \epsilon \quad \checkmark$$

And we have:

$$\frac{dH}{dt} = -O_3 \epsilon \quad \leftarrow \text{RHS is } O(1/n).$$

$$\frac{dO_3}{dt} = -O_4 \epsilon \quad \leftarrow \text{RHS is } O(1/n)$$

$$\frac{dO_4}{dt} = O\left(\frac{1}{n^2}\right) \rightarrow 0 \quad \text{if we are only interested in } O(1/n).$$

Expand all quantities in powers of $1/n$.

$$\phi = \phi^{(0)} + \phi^{(1)} + \dots$$

$$H = H^{(0)} + H^{(1)} + \dots$$

$$O_3 = O_3^{(1)} + \dots$$

$$O_4 = O_4^{(1)} + \dots$$

$$O(1) \quad H^{(0)}(t) = H^{(0)} \quad \text{const.}$$

$$\Rightarrow \phi^{(0)}(t) - y = e^{-H^{(0)}t} (\phi^{(0)}(0) - y).$$

$$O(1/m): \quad \frac{dO_3^{(1)}}{dt} = \frac{dO_3}{dt} = -O_4^{(0)} (\phi^{(0)}(t) - y),$$

$$\Rightarrow O_3(t) = O_3(0) - \int_0^t dt' O_4^{(0)} (\phi^{(0)}(t') - y)$$

$$= O_3(0) - O_4^{(0)} \int_0^t dt' e^{-H^{(0)}t'} (\phi^{(0)}(0) - y)$$

$$= O_3(0) - O_4^{(0)} (H^{(0)})^{-1} (1 - e^{-tH^{(0)}}) (\phi^{(0)}(0) - y).$$

$$H^{(1)}(t) = H^{(1)}(0) - \int_0^t dt' O_3^{(1)}(t') (\phi^{(0)}(t') - y).$$

can also be computed explicitly.

$$\frac{d\phi^{(1)}}{dt} = - [H^{(0)} \phi^{(1)}(t) + H^{(1)}(t) (\phi^{(0)}(t) - y)].$$

$$\hat{\phi}^{(n)} = e^{-H^{(n)}t} \hat{\phi}^{(n)}(t)$$

$$\frac{d\hat{\phi}^{(n)}}{dt} = -H^{(n)} e^{-H^{(n)}t} \hat{\phi}^{(n)}(t) + e^{-H^{(n)}t} \frac{d}{dt} \hat{\phi}^{(n)}(t)$$

$$= -H^{(n)} \hat{\phi}^{(n)}(t) - H^{(n)}(t) (\hat{\phi}^{(n)}(t) - y)$$

$$\frac{d}{dt} \hat{\phi}^{(n)} = - e^{H^{(n)}t} H^{(n)}(t) e^{-H^{(n)}t} (\hat{\phi}^{(n)}(t) - y)$$

$$\hat{\phi}^{(n)}(t) = - \int_0^t e^{+H^{(n)}t'} H^{(n)}(t') e^{-H^{(n)}t'} (\hat{\phi}^{(n)}(0) - y) dt'$$

Finally

$$\hat{\phi}(t) = y + e^{-H^{(n)}t} \left\{ 1 - \int_0^t dt' e^{+H^{(n)}t'} H^{(n)}(t') e^{-H^{(n)}t'} \right\} (\hat{\phi}(0) - y).$$

↑ analytic description of the network during the training

Quantity over/underfitting?

Monitor bias & variance in a closure test?

Eigenvectors/eigenvalues of $H^{(n)}$?

Linear Networks. (Fixed weights $f, l < L$).

$$\phi^{(L)}(t) = \sum_{\mu} \phi^{(L)}_{\mu} \theta_{\mu}(t)$$

$$\dot{\theta}_{\mu}(t) = -\eta \sum_{\alpha} \phi_{\alpha}^{(L)} \left(\sum_{\nu} \phi_{\nu}^{(L)} \theta_{\nu}(t) - y_{\alpha} \right)$$

$$\Rightarrow \theta_{\mu}(t) = -\sum_{\mu} \phi_{\mu}^{(L)} H^{-1} (1 - e^{-H\eta t}) (\phi^{(L)}(0) - \eta) + \theta_{\mu}(0)$$

$$\Rightarrow \phi_{\alpha_1}^{(L)}(t) = (1 - e^{-\eta H t})_{\alpha_1 \alpha_2} y_{\alpha_2} + (e^{-\eta H t})_{\alpha_1 \alpha_2} \phi_{\alpha_2}^{(L)}(0)$$

$$\phi_{\beta}^{(L)}(t) = \phi_{\beta}^{(L)}(0) + H_{\beta \alpha_1} H^{-1}_{\alpha_1 \alpha_2} (1 - e^{-\eta H t})_{\alpha_2 \alpha_3} (\eta - \phi_{\alpha_3}^{(L)}(0))_{\alpha_3}$$

$$\lambda_{b_i^{(k)}, b_i^{(k)}} = \delta_{i,i_2} \lambda_b^{(k)}, \quad \lambda_{w_{ij}^{(k)}, w_{ij}^{(k)}} = \delta_{i,i_2} \delta_{j,j_2} \frac{\lambda_w^{(k)}}{n_{l-1}}$$

$$H_{i_1, \alpha_1; i_2, \alpha_2}^{(k+1)} = \sum_{j=1}^{n_{k+1}} \left\{ \lambda_b^{(k+1)} \frac{\partial \phi_{i_1, \alpha_1}^{(k+1)}}{\partial b_j^{(k+1)}} \frac{\partial \phi_{i_2, \alpha_2}^{(k+1)}}{\partial b_j^{(k+1)}} + \frac{\lambda_w^{(k+1)}}{n_k} \sum_{k=1}^{n_k} \frac{\partial \phi_{i_1, \alpha_1}^{(k+1)}}{\partial w_{jk}^{(k+1)}} \frac{\partial \phi_{i_2, \alpha_2}^{(k+1)}}{\partial w_{jk}^{(k+1)}} \right\} +$$

$$+ \sum_{j, k=1}^{n_k} \frac{\partial \phi_{i_1, \alpha_1}^{(k+1)}}{\partial \phi_{j, \alpha_1}^{(k)}} \frac{\partial \phi_{i_2, \alpha_2}^{(k+1)}}{\partial \phi_{j, \alpha_2}^{(k)}} H_{j, \alpha_1; j, \alpha_2}^{(k)}.$$

If we only consider the parameters in the last layer, $b_i^{(k)} = 0$

$$H_{i_1, \alpha_1; i_2, \alpha_2}^{(k)} = \frac{\lambda_w}{n_{k-1}} \sum_{j, k} \delta_{i_1, j} \delta_{i_2, j} \rho_{k, \alpha_1}^{(k-1)} \rho_{k, \alpha_2}^{(k-1)}$$

$$= \frac{\lambda_w}{n_{k-1}} \delta_{i_1, i_2} \rho_{\alpha_1}^{(k-1)} \rho_{\alpha_2}^{(k-1)}$$

if $\lambda_w = c_w$ then $H^{(k)} = G^{(k)}$

→ reproduce NN GP.