

# Deep learning theory Lecture 13

Consider single layer neural networks of the form

$$f_\theta(x) = \frac{1}{m} \sum_{j=1}^m \alpha_j \sigma(w_j x)$$

Note the different scaling w.r.t. the NTK regime.  
We assume the following statistical model

$$y_j = f^*(x_j) + \varepsilon_j \quad \varepsilon_j \sim N(0, \tau^2)$$

We aim to minimize the MSE

$$R(\theta) = \frac{1}{2} \int_{\mathbb{R}^d} |f_\theta(x) - f^*(x)|^2 P(dx) + \cancel{\tau^2}$$

excess noise

This can be written as

$$\begin{aligned} f_\theta(x) &= \frac{1}{m} \sum_{j=1}^m \alpha_j \sigma(w_j x) \\ &= \int_{\mathbb{R}^{d+1+m}} \alpha \sigma(\omega x) \left( \frac{1}{m} \sum_{j=1}^m \delta(\alpha - \alpha_j) \delta(\omega - \omega_j) \right) d\omega d\alpha \\ &= \int_{\mathbb{R}^{d+2}} \alpha \sigma(\omega x) \mu_n(d\alpha, d\omega) \end{aligned}$$

where  $\mu_n(d\alpha d\omega) = \frac{1}{m} \sum_{j=1}^m \delta(\omega - \omega_j) \delta(\alpha - \alpha_j) d\omega d\alpha$

(mean field)

Note: we can describe the network as a linear function of  $\mu_n$  (the state of the network)

$$f[\mu](x) = \int a \sigma(w \cdot x) \mu(dw)$$

furthermore,  $\mu$  is invariant under permutation of neuron indices, and

$$R[\mu] = \int \left( \int a \sigma(w \cdot x) \mu(dw) - f_*(x) \right)^2 \mathbb{P}(dx)$$

is convex in  $\mu$  (composition of convex functions).

Furthermore, note that at initialization the limit  $n \rightarrow \infty$  is well-defined: if  $a_j, w_j \stackrel{iid}{\sim} p(a, w)$

$$\mu_n \xrightarrow[n \rightarrow \infty]{w} \mathbb{P}$$

The price to pay is that the network lives in an  $\infty$ -dimensional space.

### The dynamics

The dynamics of each "particle" is

$$\frac{d}{dt} \alpha_j = - \alpha_j \int \left( \frac{1}{m} \sum \alpha_c \sigma(\omega_c x) - f_*(x) \right)^2 P(dx)$$

$$= - \frac{1}{m} \int \sigma(\omega_j x) \left( \frac{1}{m} \sum \alpha_c \sigma(\omega_c x) - f_*(x) \right) P$$

$$= - \frac{1}{m} \int \sigma(\omega_j x) \left( \int \alpha \sigma(\omega x) \mu(d\omega) - f_*(x) \right) P + \eta \frac{d\beta_j}{dt}$$

$$\frac{d}{dt} \omega_j = - \frac{1}{m} x \int \alpha_j \sigma'(\omega_j x) \left( \int \alpha \sigma(\omega x) \mu(d\omega) - f_*(x) \right) P$$

$$+ \eta \frac{d\beta_j}{dt}$$

Changing time we have

$$\frac{d}{dt} \left( \begin{matrix} \alpha_j \\ \omega_j \end{matrix} \right) = \mathcal{V}[\mu] \left( \begin{matrix} \alpha_j \\ \omega_j \end{matrix} \right) = D_{(\omega)} \left( \alpha \sigma(\omega \cdot x) \left( \dots \right) P(dx) \right) + \eta \frac{d\beta_j}{dt}$$

$$\mathcal{V}[\mu] \left( \begin{matrix} \alpha \\ \omega \end{matrix} \right) = \begin{pmatrix} \int \sigma(\omega \cdot x) \left( \int \alpha' \sigma(\omega' \cdot x) \mu(d\omega' d\omega) - f_*(x) \right) P(dx) \\ x \int \alpha' \sigma'(\omega \cdot x) \left( \int \alpha' \sigma(\omega' \cdot x) \mu(d\omega' d\omega) - f_*(x) \right) P(dx) \end{pmatrix}$$

This vector field can be evaluated at any point ( $\omega$ )

The equation that evolves the measure  $\mu$  is the continuity equation :

$$\partial_t \mu_+ = - \operatorname{div} (\mu_+ \cdot \sigma(\omega)) + \frac{1}{2} \eta \Delta \mu_+$$

in this case  $\sigma(\omega) = \sigma[\mu_+](\omega)$  (mean-field)

$$\begin{aligned} \partial_t \mu_+ &= - \operatorname{div} \left( \mu_+ \left( - D_\sigma \int \omega \sigma(\omega \cdot x) (\int \omega \sigma(\omega \cdot x) \mu_+ - f_*) P(dx) \right) \right. \\ &\quad \left. + \frac{1}{2} \eta \Delta \mu_+ \right) \\ &= - \operatorname{div} \left( \mu_+ \left( - D_\sigma \frac{\delta}{\delta \mu} R^*(\mu_+) \right) \right) \quad (\text{MFPDE}) \\ &\quad \text{, variational derivative} \end{aligned}$$

$\Rightarrow$  Gradient flow in  $\mathcal{W}_2$  space. of

$$\frac{1}{2} \int \left( \int \omega \sigma(\omega \cdot x) \mu_+ (d\omega, dw) - f_* \right)^2 P(dx) + \eta \int \mu_+ \log(\mu_+)$$

For a solution  $\mu_+$  of MFPDE we write the flow map

$$\Phi_t[\mu_0] = \mu_+$$

Assumptions:  $\sigma$  is bounded, Lipschitz and has Lipschitz derivative  $\sigma'$

Proposition: Let  $\mu_0^{(n)} \xrightarrow{n \rightarrow \infty} \mu_0 \in P_2$  then for all  $T > 0$

$$\Phi_+[\mu_0^{(n)}] \xrightarrow{n \rightarrow \infty} \Phi_+[\mu_0] \quad \forall t \in [0, T]$$

Theorem 1: Let  $\mu_0$  have full support, then if

$$\mu_t \xrightarrow{t \rightarrow \infty} \mu_* \quad \mu_* \in \underset{\mu \in \mathcal{M}_+^1(\Theta)}{\operatorname{argmin}} R(\mu)$$

Theorem 2: When  $\eta > 0$  we have  $\mu_t \xrightarrow{t \rightarrow \infty} \mu_* \in \operatorname{argmin} R_\eta(\mu)$