

①

Mean Field Methods for Two Layer NN - part I.

In this set of two lectures we review the work of Montanari & Mei on two layer NN treated in a "mean field asymptotic limit".

program this week:

- 1) Two layer NN and classical theory of Cybenko and Barron.
- 2) Formulation of SGD & test error.
- 3) Particle system interpretation and idea of Mean Field Analysis.

Next time:

- 4) Mean field analysis: statics & dynamics.

(2)

1) Two layer NN's.

We consider an "hypothesis class" of functions of the type

$$\hat{f}(\underline{x}; \theta) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\underline{w}_i^T \underline{x} + b_i)$$

• where $\underline{x} \in \mathbb{R}^D$, so $\hat{f}: \mathbb{R}^D \rightarrow \mathbb{R}$

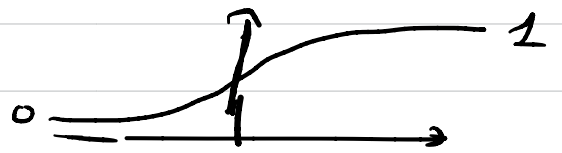
• $\theta = \{ \underline{w}_i, a_i, b_i \} = \{ \theta_i \}_{i=1}^N$, a

set of parameters (to be learned)

$$\underline{w}_i \in \mathbb{R}^D, a_i \in \mathbb{R}, b_i \in \mathbb{R}.$$

• σ = activation function of sigmoidal type

example $\sigma(x) = \frac{1}{1 + e^{-2x}}$

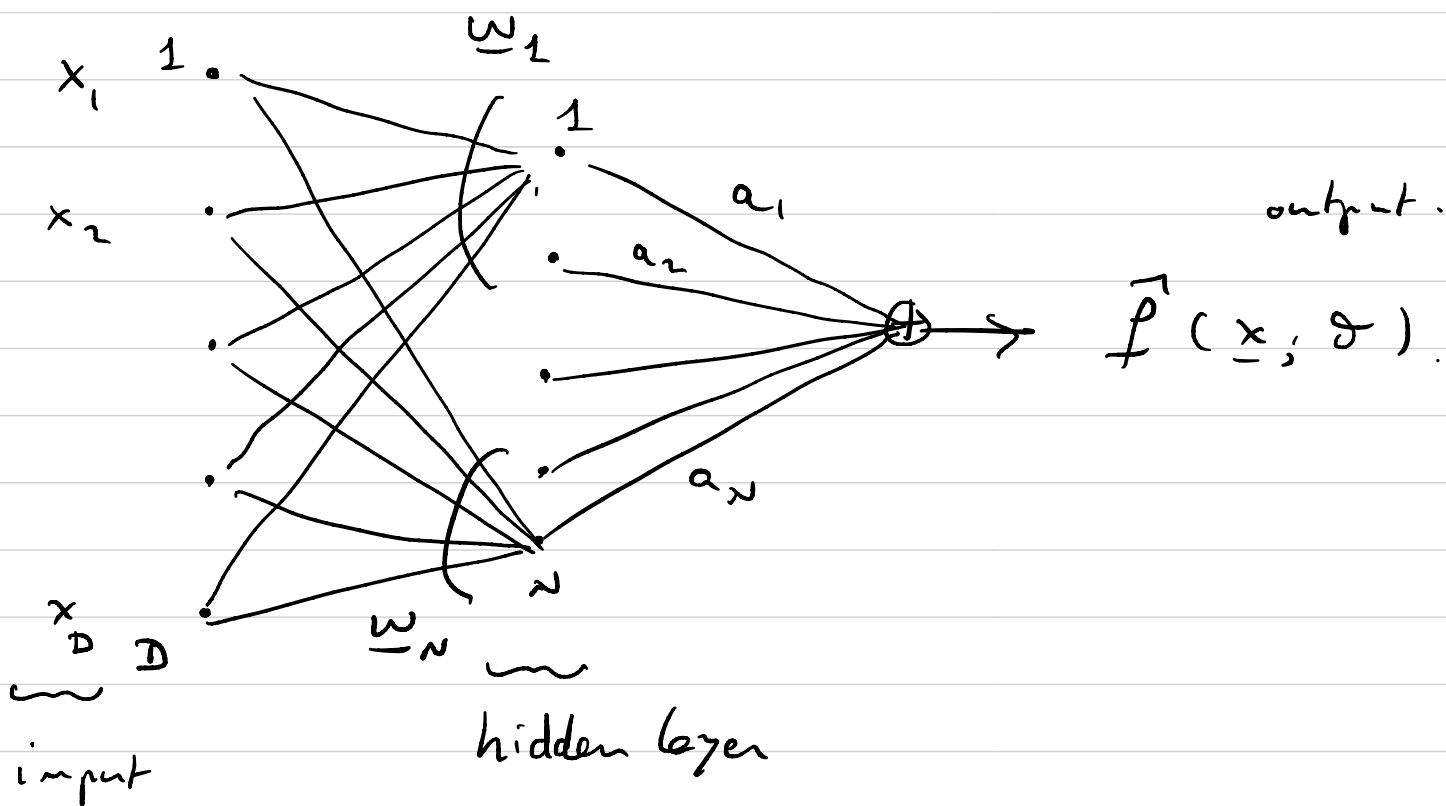


= activation of ReLU type

$$\sigma(x) = \max(0, x).$$

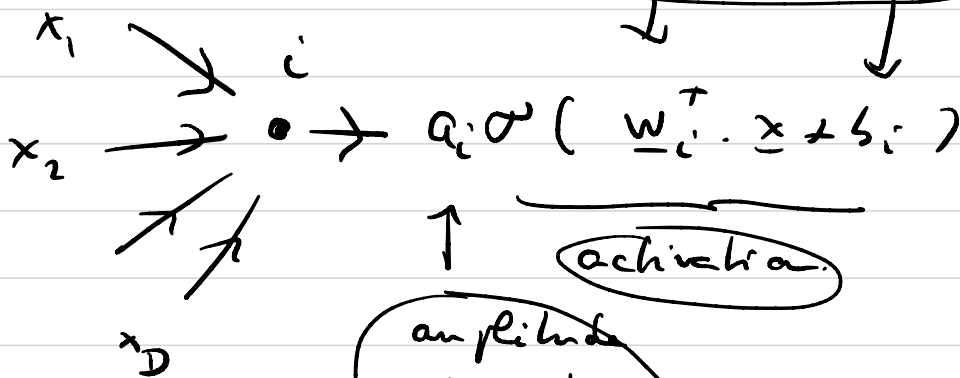


The new picture is



For each Neuron i of hidden layer:

weight associated to neuron i and bias



Typically we have training data (\underline{x}_k, y_k) and we want to adjust or learn weights θ .

(4)

Representation Power of NN's.

We state two basic thems about the representability of functions by these two layer or single-hidden layer NN's.

Cybenko 1989 thm:

Let $\mathbb{E}_{x \sim \mathcal{D}} ((f(x))^2) < \infty$ i.e. $\int dx \mathcal{D}(x) |f(x)|^2 < \infty$

Assume $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ continuous with

$$\lim_{x \rightarrow +\infty} \sigma(x) = +1, \quad \lim_{x \rightarrow -\infty} \sigma(x) = 0.$$

Then $\forall \epsilon > 0$, $\exists N(\epsilon)$ large enough s.t

$$\mathbb{E}_{x \sim \mathcal{D}} \left[\left(f(x) - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\underline{w}_i^T \underline{x} + b_i) \right)^2 \right] < \epsilon$$

(5)

The next question is how big should $N(\varepsilon)$ be?

Berman 1983 Thm.

Let $\mathcal{D}(x)$ have support in $\text{Ball}(\underline{0}, r)$.

$$\text{let } f(x) = \int_{\mathbb{R}^d} e^{i \omega \cdot x} \underbrace{F(\omega)}_{\text{Fourier transf.}} d\omega$$

and σ continuous with $\lim_{x \rightarrow +\infty} \sigma(x) = 1$, $\lim_{x \rightarrow -\infty} \sigma(x) = 0$.

then we can take

$$N(\varepsilon) = \frac{1}{\varepsilon} \left(2r \int \|\omega\|_2 F(\underline{\omega}) d\underline{\omega} \right)^2.$$

2) Formulation of SGD for two layer NN's.

- Recall the true risk or population risk:

$$R_N(\underline{\theta}) = \mathbb{E}_{\mathcal{D}} [l(y, \underline{x}; \underline{\theta})]$$

here we take

$$l(y, \underline{x}; \underline{\theta}) = (y - \hat{f}(\underline{x}; \underline{\theta}))^2$$

and $\mathcal{D}(\underline{x}, y)$ supported on $\mathbb{R}^D \times \mathbb{R}$.

- As we do not know \mathcal{D} we now try as a "proxy" for minimization to run SGD as explained in previous lectures. There is no guarantee that we truly minimize of course here since l is not convex w.r.t θ . But this is an algorithm.

Stochastic Gradient Descent algo:

$k = 1, 2, \dots, T$ time steps.

at time k take a fresh sample from training set (\underline{x}_k, y_k) and update

$$\underline{\theta}^{k+1} = \underline{\theta}^k + \delta_k \nabla_{\underline{\theta}} l(y_k, x_k; \underline{\theta}^k)$$

where $\nabla_{\underline{\theta}} l(y_k, x_k; \underline{\theta}^k)$ is an unbiased "stochastic gradient" in the sense that

$$\mathbb{E}_{(x_k, y_k)} \left(\nabla_{\underline{\theta}} l(y_k, x_k; \underline{\theta}^k) \right)$$

$$= \nabla_{\underline{\theta}} \mathbb{E}_{\underline{\theta}} l(y, x; \underline{\theta}^k)$$

$$= \nabla_{\underline{\theta}} R_{\mathcal{D}}(\underline{\theta}) \quad (\text{gradient of true risk})$$

(8)

This can be interpreted as run SGD on empirical risk $L_S(\theta) = \frac{1}{T} \sum_{k=1}^T \ell(y_k, x_k; \theta)$

with a mini-batch of size 1 at each time step k .

†.

Computation of gradient:

$\underline{\theta} = (\underline{\theta}_i)_{i=1}^N$, For each $\underline{\theta}_i = (\underline{w}_i, a_i, b_i)$

we have:

$$\nabla_{\underline{\theta}_i} \ell(y_k, x_k; \underline{\theta}) = \nabla_{\underline{\theta}_i} (y_k - \hat{f}(x_k, \underline{\theta}))^2$$

$$= \nabla_{\underline{\theta}_i} \left(y_k - \frac{1}{N} \sum_{j=1}^N a_j \sigma(\underline{w}_j^T \cdot x + b_j) \right)^2$$

$$= 2 \left(y_k - \frac{1}{N} \sum_{j=1}^N a_j \sigma(\underline{w}_j^T \cdot x + b_j) \right) \cdot \nabla_{\underline{\theta}_i} \left(\frac{1}{N} a_i \sigma(\underline{w}_i^T \cdot x + b_i) \right)$$

(9)

The SGD update are thus for each

$k=1 \dots T$ and $i=1 \dots n$;

$$\underline{\theta}_i^{k+1} = \underline{\theta}_i^k - \frac{2\delta_k}{N} \left(\gamma_k - \frac{1}{N} \sum_{j=1}^N a_j^k \sigma \left(\underline{w}_j^{TK} \underline{x}_k + b_j^k \right) \right) \nabla_{\underline{\theta}_i} \left(a_i^k \sigma \left(\underline{w}_i^{TK} \underline{x}_k + b_i^k \right) \right)$$

More explicitly we update for each k the $\underline{w}_i, a_i, b_i$:

$$\underline{w}_i^{k+1} = \underline{w}_i^k - \frac{2\delta_k}{N} \left(\gamma_k - \frac{1}{N} \sum_{j=1}^N a_j^k \sigma \left(\underline{w}_j^{TK} \underline{x}_k + b_j^k \right) \right) a_i^k \underline{x}_k \sigma' \left(\underline{w}_i^{TK} \underline{x}_k + b_i^k \right)$$

$$a_i^{k+1} = a_i^k - \frac{2\delta_k}{N} \left(\gamma_k - \frac{1}{N} \sum_{j=1}^N a_j^k \sigma \left(\underline{w}_j^{TK} \underline{x}_k + b_j^k \right) \right) \sigma \left(\underline{w}_i^{TK} \underline{x}_k + b_i^k \right)$$

$$b_i^{k+1} = b_i^k - \frac{2\delta_k}{N} \left(\gamma_k - \frac{1}{N} \sum_{j=1}^N a_j^k \sigma \left(\underline{w}_j^{TK} \underline{x}_k + b_j^k \right) \right) a_i^k \sigma' \left(\underline{w}_i^{TK} \underline{x}_k + b_i^k \right)$$

The generalization error is the average error done when a new sample is presented, here this is the true error/risk:

$$R_N(\theta) = \mathbb{E}_{\theta} \left(y - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\underline{w}_i^T \underline{x} + b_i) \right)^2$$

where $(\underline{x}, y) \sim \mathcal{D}$ is the new sample presented (not in training set).

$$\begin{aligned} R_N(\theta) &= \mathbb{E}(y^2) - \frac{2}{N} \mathbb{E} \left[y \sum_{i=1}^N a_i \sigma(\underline{w}_i^T \underline{x} + b_i) \right] \\ &\quad + \frac{1}{N^2} \mathbb{E} \left[\sum_{i,j=1}^N a_i a_j \sigma(\underline{w}_i^T \underline{x} + b_i) \sigma(\underline{w}_j^T \underline{x} + b_j) \right] \\ &= \mathbb{E}(y^2) - \frac{2}{N} \underbrace{\sum_{i=1}^N \mathbb{E} \left[y a_i \sigma(\underline{w}_i^T \underline{x} + b_i) \right]}_{\equiv V(\delta_i) \equiv V(\underline{w}_i, a_i, b_i)} \\ &\quad + \frac{1}{N^2} \underbrace{\sum_{i,j=1}^N \mathbb{E} \left[a_i a_j \sigma(\underline{w}_i^T \underline{x} + b_i) \sigma(\underline{w}_j^T \underline{x} + b_j) \right]}_{\equiv U(\delta_i, \delta_j) \equiv U(\underline{w}_i, a_i, b_i, \underline{w}_j, a_j, b_j)} \end{aligned}$$

So

$$R_N(\theta) = \bar{\kappa}(\gamma^2) - \frac{2}{N} \sum_{i=1}^N V(\theta_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\theta_i, \theta_j)$$

$$V(\theta) = -\bar{\kappa}_\theta \left[\gamma a \sigma(\underline{\omega} \cdot \underline{x} + b) \right] \quad \theta = (\underline{\omega}, a, b)$$

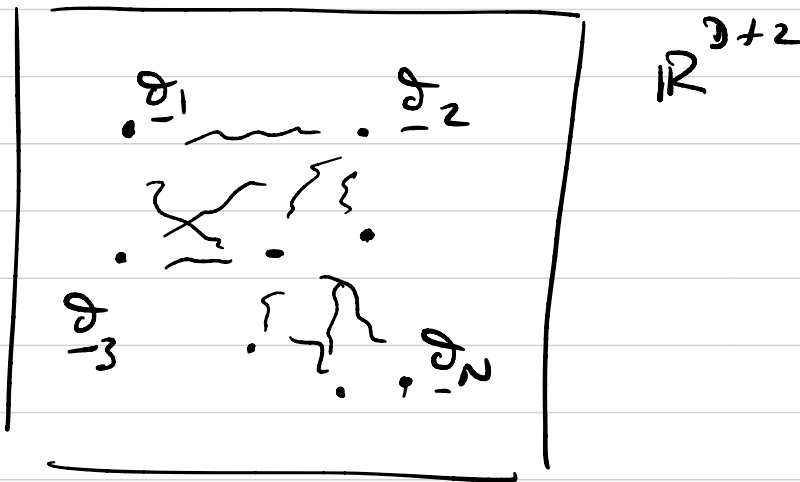
$$U(\theta, \theta') = \bar{\kappa}_\theta \left[a \sigma(\underline{\omega} \cdot \underline{x} + b) a' \sigma(\underline{\omega}' \cdot \underline{x} + b') \right]$$

$$\theta = (\underline{\omega}, a, b) \quad \theta' = (\underline{\omega}', a', b')$$

We will now discuss a "particle system" interpretation of $R_N(\theta)$ and then of SGD. This interpretation allows to take inspiration from the Mean Field Methods in physics to analyse the $N \rightarrow \infty$ limit of the two layer net's (D fixed).

Particle system interpretation,

- $R_N(\underline{\sigma})$ is the energy fct of a system of N particles at positions $\underline{\sigma}_1, \underline{\sigma}_2, \dots, \underline{\sigma}_N \in \mathbb{R}^{D+2}$



- $V(\underline{\sigma})$ external potential energy felt by a particle at position $\underline{\sigma}$ (e.g. gravity potential, electric potential, ...)
- $U(\underline{\sigma}, \underline{\sigma}')$ two body interaction potential felt by two particles at positions $\underline{\sigma}$ & $\underline{\sigma}'$.

Here all particles interact with a strength of order $\frac{1}{N}$.

pairs

This allows to use MeanField Methods (analogous to Curie-Weiss Ising spin system)

• SGD: is a kind of dynamics in discrete time steps of the system of particles.

• SGD can be rewritten as:

$$\frac{\theta_i^{k+1} - \theta_i^k}{\varepsilon} = + \nabla_{\theta_i} \left\{ \gamma_k a_i \sigma(\underline{w}_i^k \cdot \underline{x} + b_i^k) \right. \\ \left. - \frac{1}{2} \sum_{j=1}^N a_i \sigma(\underline{w}_i^k \cdot \underline{x} + b_i^k) a_j \sigma(\underline{w}_j^k \cdot \underline{x} + b_j^k) \right\}$$

where $\varepsilon = \frac{2\delta k}{N}$.

Now the left hand side is a "velocity at time k".

and

$$\nabla_{\theta_i} \left(\underbrace{V(\theta_i^k) + \frac{1}{N} \sum_{j=1}^N U(\theta_i^k, \theta_j^k)}_{\text{potential felt by particle } i \text{ due to external pot \& all interactions at time } k} \right) = - \nabla_{\theta_i}$$

potential felt by particle i due to external pot & all interactions at time k .

$$= - \nabla \text{Potential} = \underline{\underline{\text{Force}}}$$

\Rightarrow SGD = Overdamped dynamics.

- Idea of Mean Field Theory developed next time:

All equations above can be expressed by introducing an empirical particle density at time k .

$$\rho_N^k(\vartheta) = \frac{1}{N} \sum_{i=1}^N \delta(\vartheta - \vartheta_i^k)$$

In the limit $N \rightarrow +\infty$ we get a continuous particle density at continuous time t :

$$\rho_N^k(\vartheta) \rightarrow \rho(\vartheta, t)$$

and SGD equations become a PDE for

$\rho(\vartheta, t)$. This PDE is often called the

Vlasov-McKean equation. It involves a single density

$\rho(\vartheta, t)$ and potential $V(\vartheta) + \int d\vartheta' U(\vartheta, \vartheta') \rho(\vartheta, t)$

can be the "Mean Field Potential".

• In a nutshell: in the continuous limit we have a "fluid" with density $\rho(\vartheta, t)$ and velocity field $\underline{v}(\vartheta, t) = -\underline{\nabla}_{\vartheta} \psi(\vartheta; \rho(\vartheta, t))$

$$\psi(\vartheta; \rho(\vartheta, t)) = V(\vartheta) + \int d\vartheta' U(\vartheta, \vartheta') \rho(\vartheta', t)$$

• Intervention of man imposes the continuity equation:

$$\frac{\partial \rho(\vartheta, t)}{\partial t} + \underline{\nabla}_{\vartheta} \cdot \left(\rho(\vartheta, t) \underline{\nabla}_{\vartheta} \psi(\vartheta; \rho(\vartheta, t)) \right) = 0$$

(this is the continuous time limit of SFD).

• True risk becomes a functional of $\rho(\vartheta, t)$:

$$R(\rho) = \int d\vartheta V(\vartheta) \rho(\vartheta, t) + \int d\vartheta \int d\vartheta' U(\vartheta, \vartheta') \rho(\vartheta, t) \rho(\vartheta', t)$$

• The PDE above can be reinterpreted as the solution of a suitable variational problem \rightarrow allow rigorous analysis.