# Mean Field Analysis of two layer Neural Networks  II.

## 1) Recap setting from last time
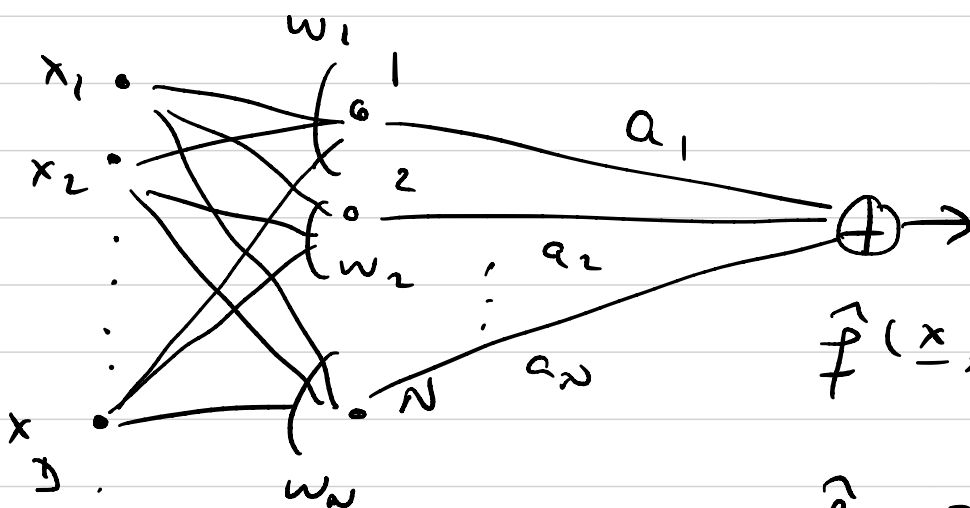
- Training data $(x_i, y_i)$, $i = 1 \cdots N$.

  e.g $\quad y_i = f(\underline{x}_i) + \xi_i$ , $\quad \xi_i \sim \mathcal{N}(0, 1)$ iid

  $\qquad x_i \sim \mathcal{D}_x$ . iid

  $\qquad f$ some "reasonable" fct.

- We want to train a two layer NN:



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

$$\hat{f} : \mathbb{R}^D \longrightarrow \mathbb{R}.$$

weights $\vartheta_i = (a_i, b_i, \underline{w}_i) \in \mathbb{R}^{D+2}$

$\qquad \vartheta = (\vartheta_1 \cdots \vartheta_N)$.

- We want to analyze SGD for true risk:

$$R_N(\theta) = \mathbb{E}_{\mathcal{D}}\left( \ell(\underline{x}, y; \theta) \right)$$

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

for $N \to +\infty$, $\mathcal{D}$ fixed.

- We argued that in this limit the problem is amenable to analysis of a PDE. The intuition comes interpretation as a "particle system".

Recap SGD iterations: $k = 1, 2, \ldots, T$

$$\vartheta^{k+1} = \vartheta^k - \delta_k \, \underline{\nabla}_\vartheta \, \ell(x_k, y_k; \vartheta^k)$$

$\underline{\nabla}_\vartheta \, \ell(x_k, y_k; \vartheta^k) =$ Stochastic gradient   since

$$\mathbb{E}_\vartheta \left( \nabla_\vartheta \, \ell(x, y; \vartheta^k) \right) = \nabla_\vartheta \, R_N(\vartheta^k) \quad \text{unbiased.}$$

We computed last time:

$$\nabla_{\vartheta_i} \ell(x_k, y_k; \vartheta^k) = \left( y_k - \frac{1}{N} \sum_{i=1}^{N} a_i^k \sigma(\underline{w}_i^k \cdot \underline{x} + s_i^k) \right)$$

$$\cdot \nabla_{\vartheta_i} \left( a_i^k \sigma(\underline{w}_i^k \cdot \underline{x}_k + s_i^k) \right)$$

Recap true risk (a generalization error):

$$R_N(\vartheta) = \mathbb{E}_{\mathscr{D}}\left[\left(y - \frac{1}{N}\sum_{i=1}^{N} a_i\, \sigma(\underline{w}_i \cdot \underline{x} + s_i)\right)^2\right]$$

$$= \mathbb{E}(y^2) - \frac{2}{N}\sum_{i=1}^{N} V(\vartheta_i) + \frac{1}{N^2}\sum_{i,j=1}^{N} U(\vartheta_i; \vartheta_j)$$

$$\begin{cases} V(\vartheta) = -\mathbb{E}_{\mathscr{D}}\left[y\, a\, \sigma(\underline{w} \cdot \underline{x} + s)\right] \\[2em] U(\vartheta, \vartheta') = \mathbb{E}_{\mathscr{D}}\left[a\sigma(\underline{w} \cdot \underline{x} + s)\, a'\sigma(\underline{w}' \cdot \underline{x} + b')\right] \end{cases}$$
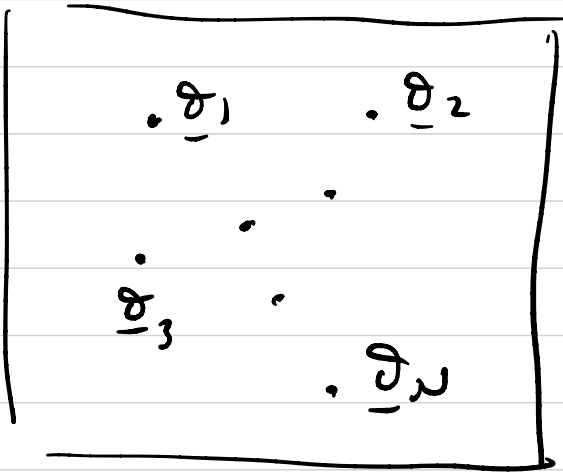
SGD can be re-written as:

$$\frac{\vartheta_i^{k+1} - \vartheta_i^{k}}{(2\delta_F/N)} = \mathcal{N}_i^{k}$$

with $\quad \mathbb{E}(\mathcal{N}_i^{k} \mid \underbrace{\underline{\vartheta}^{k}}_{\text{the past}}) = \mathbb{E}_{x_k, y_k}(\mathcal{N}_i^{k})$

$$= -\nabla_\vartheta\left(V(\vartheta_i^{k}) + \sum_{j=1}^{N} U(\vartheta_i^{k}, \vartheta_j^{k})\right)$$

## Particle system interpretation:

N particles in $\mathbb{R}^{D+2}$,

positions $\underline{\vartheta}_1^k \ldots \underline{\vartheta}_N^k$ at

time $k$.

velocities $\underline{v}_1^k \ldots \underline{v}_N^k$.

They feel potentials $V(\vartheta)$ external and

interaction potential $U(\vartheta, \vartheta')$ between all pairs.

Potential felt by a particle at $\vartheta_i$ due to

external potential and all other interaction is:

$$V(\vartheta_i) + \sum_{j=1}^{N} U(\vartheta_i, \vartheta_j).$$

We introduce the underlined empirical density

$$\rho_N^k(\vartheta) = \frac{1}{N} \sum_{j=1}^{N} \delta(\vartheta - \vartheta_j^k).$$

The potential felt by particle $i$ becomes

$$\psi(\vartheta_i^k; \rho_N^k) = V(\vartheta_i^k) + \int d\vartheta \, U(\vartheta_i^k, \vartheta) \rho_N^k(\vartheta')$$

and the true risk becomes:

$$R_N(\vartheta^k) = \mathbb{E}(y^2) + 2\int d\vartheta \, V(\vartheta) \rho_N^k(\vartheta)$$

$$+ \int d\vartheta \, d\vartheta' \, U(\vartheta, \vartheta') \rho_N^k(\vartheta) \rho_N^k(\vartheta').$$

## 2) Mean field theory : statics.

We assume that for $N \to +\infty$ we have

$$\rho_N^k (\vartheta) = \frac{1}{N} \sum_{i=1}^{N} \delta (\vartheta - \vartheta_i^k) \longrightarrow \rho (\vartheta, t).$$

Note that time steps $\frac{2\delta k}{N}$ in SGD tend

to zero so it is natural to replace $k$ by a

continuous time index $t$.

Now consider the risk and assume that

$\rho (\vartheta, t)$ has reached (for $t \to +\infty$ say) some

stationary or equilibrium density $\rho (\vartheta)$ :

$$R_N (\vartheta) \underset{N \to +\infty}{\longrightarrow} R(\rho) = \mathbb{E}[y^2] - 2 \int d\vartheta V(\vartheta) \rho(\vartheta)$$

$$+ \int d\vartheta \, d\vartheta' \, U(\vartheta, \vartheta') \rho(\vartheta) \rho(\vartheta').$$

## What is the minimizer ?

We have to minimize under the constraint that

$\int d\vartheta \, \rho(\vartheta) = 1$. Introducing a Lagrange parameter:

$$\frac{\delta}{\delta \rho(\vartheta)} \left\{ R(\rho) - \lambda \left( \int d\rho(\vartheta) - 1 \right) \right\} = 0$$

$$\Rightarrow \quad \frac{\delta R}{\delta \rho(\vartheta)} = \lambda$$

i.e. $\quad \nabla_\vartheta \left( \frac{\delta R}{\delta \rho(\vartheta)} \right) = 0$

Now $\quad \frac{\delta R}{\delta \rho(\vartheta)} = 2 V(\vartheta) + 2 \int d\vartheta' \rho(\vartheta') U(\vartheta, \vartheta')$

so the equilibrium condition is

$$\nabla_\vartheta \left\{ V(\vartheta) + \int d\vartheta' \rho(\vartheta') U(\vartheta, \vartheta') \right\} = 0$$

This can be interpreted as saying that the
mean field potential defined as

$$\psi(\vartheta; \varsigma) \equiv V(\vartheta) + \int d\vartheta' \varsigma(\vartheta') U(\vartheta, \vartheta')$$

$\underbrace{\qquad}$ external
pot

$\underbrace{\qquad\qquad\qquad}$ interaction energy due
to all other particles.

should be constant; i.e the force acting at $\vartheta$

on the "fluid" vanishes: $-\nabla_{\vartheta} \psi(\vartheta; \varsigma) = 0$.

3) <u>Mean field theory : dynamics.</u>

As said above SGD is (with $\varepsilon = \frac{2\delta_k}{N}$)

$$
\begin{cases}
\dfrac{\theta_i^{k+1} - \theta_i^{k}}{\varepsilon} = \mathcal{N}_i^{k} = \text{stochastic velocity} \\[4pt]
\qquad\qquad \underline{\text{"overdamped motion"}} \\[10pt]
\mathbb{E}(\mathcal{N}_i^{k} \mid \text{past}) = -\nabla_{\underset{v}{\theta_i^{k}}} \psi(\theta_i^{k} ; \rho_N^{t})
\end{cases}
$$

We deduce from here somewhat heuristically the

continuity equation for $N \to +\infty$ :

$$
\frac{\partial}{\partial t}\rho(\theta ; t) - \underbrace{\nabla_{\theta}}_{\text{Divergence}} \cdot \Big( \rho(\theta, t) \overbrace{\underbrace{\nabla_{\theta}}_{\text{gradient}} \psi(\theta, \rho)}^{\text{velocity}} \Big) = 0
$$

current

Non Linear PDE because $\psi(\theta, \rho) = V(\theta) + \int d\theta' \rho(\theta', t) U(\theta, \theta')$

It is often attributed to Vlasov and McKean.

$\uparrow$ Plasma physics $\qquad\qquad \uparrow$ Stochastic Processes.

Heuristic derivation:

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

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

$$\approx \frac{1}{N} \sum_{i=1}^{N} \nabla_\vartheta \delta(\vartheta - \vartheta_i^k) \cdot (\vartheta_i^k - \vartheta_i^{k+1})$$

$$\Rightarrow \frac{\rho_N^{k+1}(\vartheta) - \rho_N^k(\vartheta)}{\varepsilon} \approx -\frac{1}{N} \sum_{i=1}^{N} \nabla_\vartheta \delta(\vartheta - \vartheta_i^k) \cdot \underbrace{\left( \frac{\vartheta_i^{k+1} - \vartheta_i^k}{\varepsilon} \right)}_{\displaystyle v_i^k}$$

$$- \nabla_{\vartheta_i^k} \psi(\vartheta_i^k, \rho_N^k)$$

For $N \to +\infty$, time becomes continuous as $\varepsilon \to 0$ :

$$\boxed{\frac{\partial}{\partial t} \rho(\vartheta, t) = \nabla_\vartheta \cdot \left( \rho(\vartheta, t) \, \nabla_\vartheta \psi(\vartheta, \rho) \right)} \, .$$

# Theorem (Montanari - Mei).

Here stated informally; under suitable hypothesis it is possible to prove that:

$$
\sup_{k \in [0, \frac{T}{\varepsilon}] \cap \mathbb{N}} \left| R_N(\theta^k) - R(\rho_{k\varepsilon}) \right|
$$

$$
\leq c \, e^{cT} \sqrt{\max(\tfrac{1}{N}, \varepsilon)} \left( \sqrt{D + \log \tfrac{N}{\varepsilon}} + z \right)
$$

holds with probability $1 - e^{-z^2}$, where $\rho_{k\varepsilon}$ is the solution $\rho(\theta, t)$ of PDE for $t = k\varepsilon$.

Remark: This theorem says how far the solution of the PDE (or mean field solution) is from the solution of underlying SGD.

## 4) Variational formulation of the Nonlinear PDE.

A powerful tool to analyze the solution of the PDE is a variational formulation.

Consider a simpler toy system first: Take gradient flow in $\mathbb{R}^d$:

$$\dot{\underline{x}}(t) = -\underline{\nabla} F(\underline{x}). \qquad \underline{x} \in \mathbb{R}^d.$$

The Euler discretisation would be the "forward scheme"

$$\underline{x}_{k+1} = \underline{x}_k - \varepsilon \underline{\nabla} F(\underline{x}_k)$$

Although iterations look easy to obtain this is pretty unstable and is not the best way to approximate the continuous time evolution.

The backward Euler scheme turns out to be more stable (although harder to implement):

$$X_{k+1} = X_k - \varepsilon \underline{\nabla} f(x_{k+1})$$

The point now is that this scheme can be formulated in a variational way. Remark that (at least for $F$ convex):

$$\underline{X}_{k+1} = \underset{\underline{x}}{argmin} \left\{ F(\underline{x}) + \frac{\| \underline{x} - \underline{x}_k \|^2}{2\varepsilon} \right\}.$$

So here we minimize $F(\underline{x})$ but under the constraint of not moving too far away from current position $x_k$. ( Euclidean norm penalty ).

$\underline{Proof}$  $\underline{\nabla}_x \left\{ f(\underline{x}) + \frac{\| x - x_k \|^2}{\varepsilon} \right\} = \nabla f(\underline{x}) + \frac{x - x_k}{\varepsilon}$

so the extremum is at:

$$\underline{x} = x_k - \varepsilon \nabla f(\underline{x}) \implies \boxed{x_{k+1} = x}$$

It is a Minimum if $f$ is convex.

This scheme also guarantees that we reach a minimum of $f$:

$$f(x_{k+1}) + \frac{dist(x_{k+1}, x_k)}{2\epsilon} \leq f(x_k) + \frac{dist(x_k, x_k)}{2\epsilon}$$

$\uparrow$

Since $x_{k+1}$ is the argmin!

$\Rightarrow$

$$f(x_{k+1}) \leq f(x_k) - \frac{dist(x_{k+1}, x_k)}{2\epsilon}$$

If $x_{k+1} \neq x_k$ then $f(x_{k+1}) < f(x_k)$ we strictly "improve". The sequence $(x_k)_{k \in \mathbb{N}}$ converges as long as $f$ is bounded below.

Going back to the PDE, it is possible

to show that :

$$\frac{\partial}{\partial t} \rho(\underline{\vartheta}, t) = \underline{\nabla}\left(\rho(\underline{\vartheta}, t)\, \psi(\underline{\vartheta}, \rho_t)\right)$$

with $\qquad \psi(\underline{\vartheta}, \rho_t) = \dfrac{\delta R(\rho)}{\delta \rho_t(\underline{\vartheta})}$

is equivalent to $\qquad \rho(\underline{\vartheta}, t) = \lim\limits_{\varepsilon \to 0} \rho_{(k+1)\varepsilon}$

$$\rho_{(k+1)\varepsilon} = \text{argmin}\left[ R(\rho) + \frac{W_2^2(\rho, \rho_{k\varepsilon})}{2\varepsilon} \right]$$

where $W_2^2(\rho, \rho')$ is the squared Wasserstein

distance between two distributions :

$$W_2^2(\rho, \rho') \equiv \inf \quad \mathbb{E}\left(\| \underline{x} - \underline{y} \|^2\right)$$

$$\text{couplings } \gamma$$
$$\text{of } \rho, \rho'$$

where the set of couplings $d\gamma(x, y)$ satisfy

$$\int_{Y} d\gamma(x, y) = \rho(x)$$

$$\int_{X} d\gamma(x, y) = \rho'(y)$$

(i.e Marginals equal $\rho$ & $\rho'$).

Idea of proof:

Go back to discretized problem and use that

the Wasserstein distance between

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

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

is $\inf_{\pi} \sum_{i=1}^{N} \| \vartheta_i^{k+1} - \vartheta_i^{k} \|^2 = \sum_{i=1}^{N} \| \vartheta_i^{k+1} - \vartheta_i^{k} \|^2$

↑
for small moves at each step.

Then taking

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

we have

$$R(\rho) + \frac{W_2^2(\rho, \rho_{k\epsilon})}{2\epsilon}$$

$$\approx R_N(\vartheta) + \sum_{i=1}^{N} \frac{\|\vartheta_i - \vartheta_i^k\|^2}{2\epsilon}$$

Minimization over $\vartheta$ gives back SGD in

Backward Euler Scheme:

$$\frac{\vartheta_i - \vartheta_i^k}{\epsilon} = -\nabla_{\vartheta_i} \left( V(\vartheta_i) + \sum_{j=1}^{N} U(\vartheta_i, \vartheta_j^k) \right)$$

$$\Rightarrow \vartheta_i^{k+1} = \vartheta_i^k - \epsilon \nabla_{\vartheta^{k+1}} V(\vartheta^{k+1}) + \sum_{j=1}^{N} U(\vartheta_i^{k+1}, \vartheta_j^k)$$

For $\epsilon \to 0$ this gives back continuity equation.