# Lecture 16 : Mean-field limit

## Derivation and meaning of continuity equation:

$$\partial_t \rho_t = \text{div}\left( \rho_t \cdot \nabla \Psi_{\rho_t} \right) \qquad (\spadesuit)$$

Holds in "weak sense", i.e. for any "nice" (e.g. bounded, differentiable, with bounded gradient) test function $\varphi : \mathbb{R}^d \to \mathbb{R}$,

$$\int \varphi(\theta) \; \partial_t \rho_t(\theta) \, d\theta = \int \varphi(\theta) \cdot \text{div}(\rho_t \cdot \nabla \Psi_{\rho_t})(\theta) \, d\theta \quad (\spadesuit\spadesuit)$$

(b/c differentiable solution to $(\spadesuit)$ may not exist).

Note that for $\bar{\Theta}_t \sim \rho_t$,

$$\text{LHS } (\spadesuit) = \frac{\partial}{\partial t} \mathbb{E}\left[ \varphi(\bar{\Theta}_t) \right]$$

(diff. under integral)

$$= \mathbb{E}\left[ \langle \nabla \varphi(\bar{\Theta}_t), \tfrac{d}{dt} \bar{\Theta}_t \rangle \right]$$

(gradient flow for $\bar{\Theta}_t$)

$$= \int \langle \nabla \varphi(\theta), \; -\nabla \Psi_{\rho_t}(\theta) \rangle \, d\rho_t(\theta)$$

( integration by parts)

$$\text{RHS } (\spadesuit) = -\int \langle \nabla \varphi(\theta), \; \nabla \Psi_{\rho_t}(\theta) \rangle \, d\rho_t(\theta) \Bigg\} =$$

# Non-asymptotic convergence to the mean-field limit :

"Propagation of chaos" [Kac '56], [McKean '69],
                                          [Sznitman '91]

want to compare:

- $\left(\Theta_i^{(k)}\right)_{k=0,1,2,\dots}$ : GD iterates given by
$$\Theta_i^{(k+1)} \leftarrow \Theta_i^{(k)} - h\nabla L(\Theta^{(k)})$$

- $\left(\bar{\Theta}_i^t\right)_{t\geq 0}$ : mean-field iterates given by
$$d\bar{\Theta}_i^t = -\nabla L_{\rho_t}(\bar{\Theta}_i^t)\,dt$$
where $\rho_t = \text{law}(\bar{\Theta}_i^t)$

## Note :

$$\Theta_i^{(k)} = \Theta_i^{(0)} + 2h\sum_{\ell=0}^{k-1} F_i\left(\Theta^{(\ell)}; (x_{\ell+1}, y_{\ell+1})\right)$$

$$\bar{\Theta}_i^t = \Theta_i^{(0)} + 2\int_0^t G\left(\bar{\Theta}_i^s; \rho_s\right)ds$$

for $F_i\left(\Theta; (x,y)\right) \triangleq \left(y - f_\theta(x)\right)\cdot\nabla_{\theta_i}\sigma(x;\theta_i)$

$$G(\theta,\rho) \triangleq -\nabla\Psi_\rho(\theta)$$

Our goal: upper bound $\|\bar{\Theta}_i^{kh} - \Theta_i^{(k)}\|$

To do so, will bound by a <u>self-similar</u> expression of the form

(small terms) $+ \int_0^{kh} \|\bar{\Theta}_i^s - \Theta_i^{(\lfloor s/h\rfloor)}\|\,ds$

This will imply (by Grönwall's inequality), the desired bound

Let $[S] = h \cdot \lfloor s/h \rfloor$

$$\left\| \bar{\Theta}_i^{kh} - \Theta_i^k \right\|$$

$$= 2 \left\| \int_0^{kh} G(\bar{\Theta}_i^s ; \rho_s) \, ds - h \sum_{\ell=0}^{k-1} F_i(\Theta^{(\ell)}, (x_{\ell+1}, y_{\ell+1})) \right\|$$

$$\leq 2 \left\| \int_0^{kh} \left[ G(\breve{\Theta}_i^s ; \rho_s) - G(\bar{\Theta}_i^{[s]} ; \rho_{[s]}) \right] ds \right\| \quad ①$$

$$+$$

$$2 \left\| \int_0^{kh} \left[ G(\bar{\Theta}_i^{[s]} ; \rho_{[s]}) - G(\Theta_i^{(\lfloor s/h \rfloor)} ; \rho_{[s]}) \right] ds \right\| \quad ②$$

$$+$$

$$2 \left\| h \sum_{\ell=0}^{k-1} \left[ G(\Theta_i^{(\ell)} ; \rho_{\ell h}) - F_i(\Theta^{(\ell)} ; (x_{\ell+1}, y_{\ell+1})) \right] \right\| \quad ③$$

---

① (easy):

Small because $G$ is Lipschitz by assumption, and can show $\rho$ varies smoothly over time so that $\rho_s$ and $\rho_{[s]}$ are close

②  :

Again by Lipschitzness of $G$,

$$\left\| G(\bar{\Theta}_i^s ; \rho_{[s]}) - G(\Theta_i^{(\lfloor s/h \rfloor)} ; \rho_{[s]}) \right\| \lesssim \left\| \bar{\Theta}_i^s - \Theta_i^{(\lfloor s/h \rfloor)} \right\|,$$

so  ②  is bounded by

$$\int_0^{kh} \underbrace{\left\| \bar{\Theta}_i^s - \Theta_i^{(\lfloor s/h \rfloor)} \right\|}_{\substack{\text{looks analogous} \\ \text{to what we want} \\ \text{to bound on LHS...}}} ds$$

③ :

$$\sum_{l=0}^{k-1} \left[ G(\Theta_i^{(l)} ; \rho_{lh}) - F_i(\Theta^{(l)} ; (x_{l+1}, y_{l+1})) \right]$$

Key idea: this has expectation $G(\Theta_i^{(l)} ; \hat{\rho}_l)$,

where $\hat{\rho}_l$ is empirical dist $\frac{1}{N} \sum_{i=1}^{N} \delta_{\Theta_i^{(l)}}$

over many steps $l$, the total deviation between

$F_i(\Theta^{(l)} ; (x_{l+1}, y_{l+1}))$'s and $G(\Theta_i^{(l)} ; \hat{\rho}_l)$'s is

of order $h\sqrt{kp}$ by <u>martingale concentration</u>

Remains to bound

$$\sum_{l=0}^{k-1} \left[ G(\theta_{:,i}^{(l)}; \rho_{lh}) - G(\theta_{:,i}^{(l)}; \hat{\rho}_l) \right]$$

$$= \frac{1}{N} \sum_{l=0}^{k-1} \sum_{j=1}^{N} \left[ \underset{\bar{\theta}}{\mathbb{E}}\, U(\theta_{:,i}^{(l)}, \bar{\theta}_j^{lh}) - U(\theta_{:,i}^{(l)}, \theta_j^{(l)}) \right]$$

again, by martingale concentration we can essentially replace $\underset{\bar{\theta}}{\mathbb{E}}\, U(\theta_{:,i}^{(l)}, \bar{\theta}_j^{lh})$ (deterministic)

with $U(\theta_{:,i}^{(l)}, \bar{\theta}_j^{lh})$ (random)

Then we use Lipschitzness of $U$ to get

$$\frac{1}{N} \sum_{l=0}^{k-1} \sum_{j=1}^{N} \left\| U(\theta_{:,i}^{(l)}, \bar{\theta}_j^{lh}) - U(\theta_{:,i}^{(l)}, \theta_j^{(l)}) \right\|$$

$$\leq \frac{1}{N} \sum_{l=0}^{k-1} \sum_{j=1}^{N} \left\| \underbrace{\bar{\theta}_j^{lh} - \theta_j^{(l)}}_{} \right\|$$

once again, a term
that looks similar to what we want to bound

<u>When data distribution has symmetries,</u>
<u>PDE simplifies considerably :</u>

Suppose training data $\{(x_i, y_i)\}$ satisfy $x_i \sim N(0, Id)$

and $y_i = \varphi(\Pi x)$ for $\Pi$ a projection to a low-dim subspace $V^\circ$.

Then joint dist over $(x, y)$ invariant under rotations of $x$ that preserve $V^\circ$, i.e. $Rv \in V^\circ \; \forall \; v \in V^\circ$.

<u>Observation:</u> Let $R$ be such a rotation. If

$\rho_0$ and $\rho_0'$ are two different initializations of the weights related by $\rho_0' = R_\# \rho_0$ ( i.e. to sample $(a', w')$ from $\rho_0'$, sample $(a, w)$ from $\rho_0$ and take $a' = a$, $w' = Rw$), then $\rho_t' = R_\# \rho_t$.

So if $\rho_0$ rotation-invariant, $\rho_t$ is invariant to rotations preserving $V^\circ$, for any $t \geq 0$!

$\rho_t$ thus completely specified by distribution on

$$\left( a, \underbrace{\Pi w}_{= s}, \underbrace{\| \Pi^\perp w \|_2}_{= r} \right),$$

i.e. we get a $\boxed{\dim(V^\sigma)+2}$ -dimensional PDE!

Denote dist. on $(a, \vec{s}, r)$ by $\bar{\rho}_t$ .

$$\partial_+ \bar{\rho}_t = \operatorname{div}\left( \bar{\rho}_t \cdot \nabla_{\vec{s}} \, \underline{\Psi}_{\bar{\rho}_t} \right) +$$

$$\partial_a \left( \bar{\rho}_t \cdot \partial_a \, \underline{\Psi}_{\bar{\rho}_t} \right) +$$

$$\frac{1}{r} \partial_r \left( r \cdot \bar{\rho}_t \cdot \partial_r \, \overline{\Psi}_{\bar{\rho}_t} \right) .$$