

Lecture 12 : Mean-field limit

Derivation and meaning of continuity equation:

$$\partial_t \rho_+ = \operatorname{div}(\rho_+ \cdot \nabla \bar{\Psi}_{\rho_+}) \quad (\dagger)$$

Holds in "weak sense", i.e. for any "nice" (e.g. bounded, differentiable, with bounded gradient)

test function $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\int (\varphi(\theta) \partial_t \rho_+(\theta)) d\theta = \int \varphi(\theta) \cdot \operatorname{div}(\rho_+ \cdot \nabla \bar{\Psi}_{\rho_+})(\theta) d\theta \quad (\ddagger)$$

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

Note that for $\bar{\Theta}_+ \sim \rho_+$,

$$\text{LHS } (\ddagger) = \frac{\partial}{\partial t} \mathbb{E} [\varphi(\bar{\Theta}_+)]$$

(diff. under integral)

$$= \mathbb{E} \left[\nabla \varphi(\bar{\Theta}_+) \cdot \frac{d}{dt} \bar{\Theta}_+ \right]$$

(gradient flow
for $\bar{\Theta}_+$)

$$= \int \langle \nabla \varphi(\theta), -\nabla \bar{\Psi}_{\rho_+}(\theta) \rangle d\rho_+(\theta)$$

(integration by parts)

$$\text{RHS } (\ddagger) = - \int \langle \nabla \varphi(\theta), \nabla \bar{\Psi}_{\rho_+}(\theta) \rangle d\rho_+(\theta)$$

Non-asymptotic convergence to the mean-field limit :

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

want to compare:

- $(\Theta_i^{(k)})_{k=0,1,2,\dots}$: GD iterates given by

$$\Theta_i^{(k+1)} \leftarrow \Theta_i^{(k)} - h \nabla L(\Theta^{(k)})$$
- $(\bar{\Theta}_i^+)_{t \geq 0}$: mean-field iterates given by

$$d\bar{\Theta}_i^+ = - \nabla L_p(\bar{\Theta}_i^+) dt$$

 where $\rho_t = \text{law}(\bar{\Theta}_i^+)$

Note:

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

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

$$\text{for } F_i(\Theta; (x, y)) \triangleq (y - f_\Theta(x)) \cdot \nabla_{\Theta_i} \sigma(x; \Theta_i)$$

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

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

To do so, will bound by $\underbrace{\text{small terms}}_{(small \text{ terms})} + \int_0^{kh} \underbrace{\text{self-similar expression of the form}}_{\|\bar{\Theta}_i^s - \Theta_i^{(ls/h)}\|} ds$

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

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

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

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

$$\leq 2 \left\| \int_0^{kh} [G(\bar{\Theta}_i^s; \beta_s) - G(\bar{\Theta}_i^{[s]}; \beta_{[s]})] ds \right\| \textcircled{1}$$

+

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

+

$$2 \left\| h \sum_{l=0}^{k-1} [G(\Theta_i^{(l)}; \beta_{lh}) - F_i(\Theta_i^{(l)}; (x_{l+1}, y_{l+1}))] \right\| \textcircled{3}$$

\textcircled{1} (easy):

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

② :

Again by Lipschitzness of G ,

$$\|G(\bar{\theta}_i^s; \beta_{[s]}) - G(\theta_i^{(ls/h)}; \beta_{[s]})\| \leq \|\bar{\theta}_i^s - \theta_i^{(ls/h)}\|,$$

so ② is bounded by

$$\int_0^{kh} \|\bar{\theta}_i^s - \theta_i^{(ls/h)}\| ds$$

looks analogous
to what we want
to bound on LHS...

③ : $\sum_{l=0}^{k-1} \left[G(\theta_i^{(l)}; \beta_{lh}) - F_i(\theta_i^{(l)}; (x_{l+1}, y_{l+1})) \right]$

Key idea: this has expectation $G(\theta_i^{(l)}; \hat{\beta}_i)$,

where $\hat{\beta}_i$ is empirical dist $\frac{1}{N} \sum_{j=1}^N \delta_{\theta_i^{(j)}}$

Over many steps l , the total deviation between

$F_i(\theta_i^{(l)}; (x_{l+1}, y_{l+1}))$'s and $G(\theta_i^{(l)}; \hat{\beta}_i)$'s is
of order $h\sqrt{kp}$ by martingale concentration

Remains to bound

$$\sum_{l=0}^{k-1} \left[G(\theta_i^{(l)}; p_{lh}) - G(\theta_i^{(l)}; \hat{p}_l) \right]$$

$$= \frac{1}{N} \sum_{l=0}^{k-1} \sum_{j=1}^N \left[\underset{\bar{\Theta}}{\oplus} V(\theta_i^{(l)}, \bar{\theta}_j^{lh}) - V(\theta_i^{(l)}, \theta_j^{(l)}) \right]$$

again, by martingale Concentration we can

essentially replace $\underset{\bar{\Theta}}{\oplus} V(\theta_i^{(l)}, \bar{\theta}_j^{lh})$ (deterministic)

with $V(\theta_i^{(l)}, \bar{\theta}_j^{lh})$ (random)

Then we use Lipschitzness of V to get

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

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

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

When data distribution has symmetries,

PDE simplifies considerably :

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

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

Then joint dist over (x, y) invariant under rotations of x that preserve V° , i.e. $R \in V^\circ \perp \vee V^\circ$.

Observation: Let R be such a rotation. If p_0 and p'_0 are two different initializations of the weights related by $p'_0 = R \# p_0$ (i.e. to sample (α', w) from p'_0 , sample (α, w) from p_0 and take $\alpha' = \alpha, w' = R w$), then $p'_t = R \# p_t$.

So if p_0 rotation-invariant, p_t is invariant to rotations preserving V° , for any $t \geq 0$!

p_t thus completely specified by distribution on $(\underbrace{\alpha}_{\in S}, \underbrace{\Pi w}_{\in \mathbb{R}^d}, \underbrace{\|\Pi^\perp w\|_2}_{\in \mathbb{R}})$,

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

Denote dist. or (a, \vec{s}, r) by $\bar{\rho}_+$.

$$\begin{aligned} \partial_t \bar{\rho}_+ &= \operatorname{div}(\bar{\rho}_+ \cdot \nabla_{\vec{s}} \Psi_{\bar{\rho}_+}) + \\ &\quad \partial_a (\bar{\rho}_+ \cdot \partial_a \Psi_{\bar{\rho}_+}) + \\ &\quad \frac{1}{r} \partial_r (r \cdot \bar{\rho}_+ \cdot \partial_r \Psi_{\bar{\rho}_+}). \end{aligned}$$