

Lecture 6: Gradient Dynamics of Shallow Neural Networks:

From lazy to mean-field regime.

- Lecture Objectives:
- highlight transition from linear to non-linear learning
 - Limitations of linear models.
 - Mean-field description for shallow neural networks, beyond the linearised regime.

Lazy Dynamics [Chizat, Ouyallon, Bach, '19]

- Let $\phi(\theta)$ be a differentiable mapping
 θ_0 : initial point of the learning algorithm.
- θ parameter space $\xrightarrow{\phi}$ F function space.

- Linearised Tangent Model at initialisation:

$$\bar{\phi}(\theta) = \phi(\theta_0) + D\phi(\theta_0) \cdot (\theta - \theta_0)$$

fixed family of features

$$\hookrightarrow \bar{\phi}(\theta)(x) = \phi(\theta_0; x) + \sum_j (\theta_j - (\theta_0)_j) \cdot \nabla_{\theta_j} \phi(\theta_0; x)$$

→ observe that $\bar{\phi}$ is linear (or affine) w.r.t. θ
(but generally non-linear w.r.t. x !)

- Q: When is gradient-descent learning under these two models $(\phi, \bar{\phi})$ similar?

- Let $E(\theta) = L(\phi(\theta))$; e.g. $E(\theta) = \|\phi(\theta) - S^*\|^2$
and assume θ_0 s.t. $E(\theta_0) > 0$.

- Consider a gradient step: $\theta_1 = \theta_0 - \eta \nabla E(\theta_0)$.

Relative change in objective function:

$$\frac{E(\theta_1) - E(\theta_0)}{E(\theta_0)} \approx \frac{\langle \nabla E(\theta_0), \theta_1 - \theta_0 \rangle}{E(\theta_0)} = -\eta \frac{\|\nabla E(\theta_0)\|^2}{E(\theta_0)}$$

$$\Delta E = \frac{E(\theta_1) - E(\theta_0)}{E(\theta_0)} \approx \frac{\langle \nabla L(\theta_0), \theta_1 - \theta_0 \rangle}{E(\theta_0)} = \frac{\langle \nabla L(\theta_0), D\phi(\theta_0) \rangle}{E(\theta_0)}.$$

↳ Relative change in tangent model:

$$\Delta(D\phi) = \frac{\|D\phi(\theta_1) - D\phi(\theta_0)\|}{\|D\phi(\theta_0)\|} \stackrel{\uparrow}{\leq} \frac{\eta \|D^2\phi(\theta_0)\| \|DE(\theta_0)\|}{\|D\phi(\theta_0)\|}$$

1st order Taylor: $|f(a) - f(b)| \leq \sup |f'| |a - b|$

→ We say that a differentiable model operates in the "lazy" regime

when $(*) \quad \Delta(D\phi) \ll \Delta(E)$ i.e. tangent model evolves much slower than the loss.

→ If we consider $L(f) = \|f - f^*\|^2$, then we verify that $(*)$ is equivalent to

$$\| \phi(\theta) - f^* \| \frac{\|D^2\phi(\theta)\|}{\|D\phi(\theta)\|^2} \ll 1$$

$K_\phi(\theta)$: relative scale of ϕ
at θ

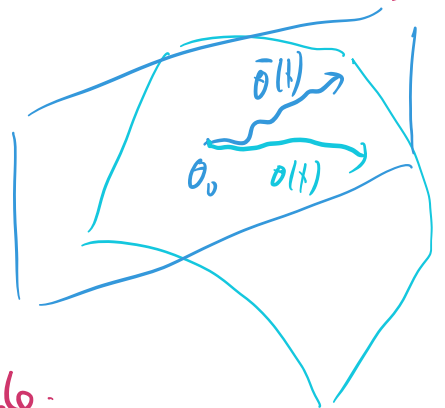
Theorem [CUB, Thm 2.3 (Simplified)]: Assume ϕ , $D\phi$ are both Lipschitz in a neighborhood of θ_0 . Let

$\theta(t)$: the (non-linear) gradient flow that solves $\dot{\theta}(t) = -\nabla L(\phi(\theta(t)))$ Non-linear
↑

$\bar{\theta}(t)$: the linearised gradient flow that solves $\dot{\bar{\theta}}(t) = -\nabla L(\bar{\phi}(\bar{\theta}(t)))$

Then $\exists C_\phi$ such that for $t \in C_\phi$ it holds

$$\frac{\| \phi(\theta(t)) - \bar{\phi}(\bar{\theta}(t)) \|}{\| \phi(\theta_0) - f^* \|} \lesssim \underbrace{t^2}_{\substack{\phi \\ \text{relative scale}}} K_\phi(\theta_0)$$



Remark: Here we present a finite-time result, but under further assumptions, it can be extended to infinite-time.

Q: When does lazy regime happen?

↳ Scaled model: $\phi_\alpha(\theta) = \alpha \cdot \phi(\theta)$

$$K_{\phi_\alpha}(\theta_0) = \frac{1}{\alpha} \|\alpha \phi(\theta_0) - f^*\| \cdot \frac{\|D^2 \phi(\theta_0)\|}{\|D \phi(\theta_0)\|^2}$$

when $\phi(\theta_0) = 0$ (centering condition).

then we have $K_{\phi_\alpha}(\theta_0) = \frac{K_\phi(\theta_0)}{\alpha} \rightarrow 0$ as $\alpha \rightarrow \infty$.

$$\|\phi(\theta) - f^*\| \frac{\|D^2 \phi(\theta)\|}{\|D \phi(\theta)\|^2}$$

↳ Homogeneous Model: $\phi(\lambda\theta) = \lambda^r \phi(\theta) \quad \forall \theta, \lambda \geq 0$.

Same as before

$$K_\phi(\lambda\theta_0) = \frac{1}{\lambda^r} K_\phi(\theta_0) \quad (\text{for centered init}).$$

↳ Single hidden-layer NN: $\phi_m(\theta) = \alpha(m) \cdot \sum_{j=1}^m g(\theta_j)$, $\theta_j \sim \mu$ iid.

$$\text{ex, } g(\theta; x) = c \cdot \text{ReLU}(\langle x, a \rangle + b)$$

with $\mathbb{E}_{\mu} g(\theta) = 0$, and Dg is Lipschitz.

$$\text{we want: } K_m = \mathbb{E}_{\mu} [K_{\phi_m}(\theta)].$$

$$\hookrightarrow \mathbb{E} \|\phi_m(\theta)\|^2 = m \alpha(m)^2 \mathbb{E} \|g(\theta)\|^2 \quad (\text{since } \theta_i \text{ are iid})$$

$$\hookrightarrow D\phi_m(\theta) = \alpha(m) [Dg(\theta_1), \dots, Dg(\theta_m)]$$

$$\frac{D\phi_m(\theta) \cdot D\phi_m(\theta)^T}{m \cdot \alpha(m)^2} = \left[\frac{1}{m} \sum_{j=1}^m Dg(\theta_j) \cdot Dg(\theta_j)^T \right]$$

$$\xrightarrow{m \rightarrow \infty} \mathbb{E} [Dg(\theta) \cdot Dg(\theta)^T] \quad (\text{Law of Large Numbers})$$

$$\Rightarrow \mathbb{E} \|D\phi_m(\theta)\|^2 = \mathbb{E} \|D\phi_m(\theta) \cdot D\phi_m(\theta)^T\|$$

$$\approx m \cdot \alpha(m)^2 \underbrace{\mathbb{E} \|Dg(\theta)\|^2}$$

$$\hookrightarrow \|D^2 \phi_m(\theta)\| = \sup_{\substack{\|u\| \leq 1 \\ u \in \mathbb{R}^{d \times m}}} \alpha(m) \sum_{j=1}^m u_j^T D^2 g(\theta_j) u_j$$

$$\sup \Sigma \leq \Sigma^{\sup} \leq \alpha(m) \sup_{\theta_i} \|D^2 g(\theta_i)\| \leq \alpha(m) \cdot \text{Lip}(Dg)$$

$$\hookrightarrow \sup_x \|\nabla^2 f(x)\| \leq \text{Lip}(f).$$

$$\hookrightarrow \text{From triangle inequality: } \|\phi_m(\theta) - f^*\| \leq \|f^*\| + \|\phi_m(\theta)\|$$

We put everything together:

$$K_m \leq \left(C_1 + \sqrt{m} \cdot \alpha(m) \cdot C_2 \right) \cdot \frac{\alpha(m) \cdot C_3}{m \cdot \alpha(m)^2 \cdot C_4}$$

$$K_m \leq \frac{\tilde{C}_1}{\sqrt{m}} + \frac{\tilde{C}_2}{m \cdot \alpha(m)}$$

Conclusion: If $\frac{m \cdot \alpha(m)}{\sqrt{m}} \rightarrow \infty$ as m grows, then this model will become "lazy" in the overparametrised (m large) regime!

In particular, NTK considers $\alpha(m) = \frac{1}{\sqrt{m}}$, so $m \alpha(m) = \sqrt{m}$!

- ⊕ It will guarantee global convergence; optimisation is easy!
- ⊖ We are in essence giving up on the non-linear nature of the model: we are learning only a linear combination

of fixed feature maps. Its

↳ no "representation" learning

↳ Associated functional space is the RKHS with kernel generated by the NTIC (features given by $\nabla \phi(\theta_0)$)

↳ even a single neuron $f^*(x) = g(\theta^*; x)$ is not in the RKHS [Bach '15].

Q: What happens when $m \cdot \alpha(m) = \Theta(1)$, i.e. $\alpha(m) = 1/m$?

Shallow Neural Networks and Particle Interaction Systems

→ Recall our shallow model $\phi(\theta_1 \dots \theta_m; x) = \frac{1}{m} \sum_{j=1}^m g(\theta_j; x)$

$$g(\theta; x) = c \cdot \sigma(\langle x, a \rangle + b) \quad \theta = (a; b; c) \in \underbrace{\mathbb{R}^d \times \mathbb{R} \times \mathbb{R}}_{\mathbb{R}^{d+2}} := \mathcal{D}$$

\downarrow sigmoid, ReLU, etc. \downarrow input weight \downarrow bias
 \downarrow output weight

→ Consider a least-squares regression:

$$\min_{\substack{\theta_1 \dots \theta_m \\ \vec{\theta}}} \mathcal{L}(\vec{\theta}) = \|\phi(\vec{\theta}) - f^*\|_J^2 + \lambda \cdot V(\vec{\theta}), \quad V(\vec{\theta}) = \frac{1}{m} \sum_j V(\theta_j)$$

$$\|\underline{f}\|_J^2 = \mathbb{E}_{x \sim \mu} [f(x)]^2$$

→ By developing the square, we get

$$\mathcal{L}(\vec{\theta}) = \|\underline{f^*}\|^2 - 2 \langle \underline{\phi(\vec{\theta})}, \underline{f^*} \rangle + \|\underline{\phi(\vec{\theta})}\|^2 + \lambda V(\vec{\theta})$$

→ Introduce the functions

$$F: \mathcal{D} \rightarrow \mathbb{R}$$

$$\langle g, f \rangle := \mathbb{E}_x [f(x) \cdot g(x)]$$

$$\left(= \frac{1}{n} \sum_{i=1}^n f(x_i) g(x_i) \right) \text{ empirical loss}$$

$$\theta \mapsto F(\theta) = \langle g(\theta), \underline{f^*} \rangle - \frac{1}{2} \lambda \cdot V(\theta)$$

$$K: \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$$

$$(\theta, \theta') \mapsto K(\theta, \theta') = \langle g(\theta), g(\theta') \rangle$$

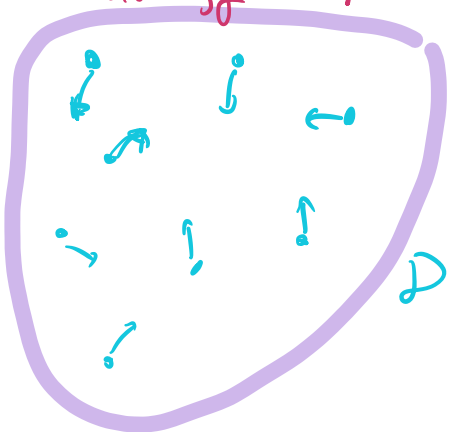
↳ observe that K is symmetric and psd operator.

→ The loss becomes

$$\begin{aligned} \mathcal{L}(\theta_1 \dots \theta_m) &= C - \underbrace{\frac{2}{m} \langle \sum_{j=1}^m g(\theta_j), f \rangle}_{\text{irrelevant}} + \underbrace{\left\| \sum_{j=1}^m g(\theta_j) \right\|^2}_{\text{interaction kernel}} + \frac{1}{m} \sum_{j=1}^m F(\theta_j) \\ &= C - \frac{2}{m} \sum_{j=1}^m F(\theta_j) + \frac{1}{m^2} \sum_{j,j'=1}^m K(\theta_j, \theta_{j'}) \end{aligned}$$

We have written $\left[\begin{array}{l} \text{an energy of a system} \\ \text{external field/force} \end{array} \right]$ of m interacting particles. $\left[\begin{array}{l} \text{interaction} \\ \text{kernel} \end{array} \right]$

analogy between neurons \leftrightarrow particles.



→ Scaled gradient Flow wrt $\theta_1 \dots \theta_m$ gives the associated Lagrangian dynamics:

$$\begin{aligned} \dot{\theta}_j &= -\frac{m}{2} \nabla_{\theta_j} \mathcal{L}(\theta_1 \dots \theta_m) \quad j=1 \dots m \\ &= + \nabla F(\theta_j) - \frac{1}{m} \sum_{j'=1}^m \nabla K(\theta_j, \theta_{j'}) \end{aligned}$$

↳ This Lagrangian description is increasingly complex (defined over $\mathcal{D}^{\otimes m}$)

↳ It is exact at particle level — but the resulting optimization landscape is very complex. (it has bad local minima).

↳ Can we instead obtain a simple collective behavior?

→ let's instead consider the Eulerian perspective

$$\vec{\theta} = (\theta_1, \dots, \theta_m) \in \mathcal{D}^{(m)} \longleftrightarrow \hat{\mu}_m = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j} \in \mathcal{P}(\mathcal{D})$$

\longleftrightarrow Eulerian view

Lagrangian view
 \mathcal{D}^m Finite-dimensional
 Euclidean space

Infinite-Dimensional
 Non-Euclidean

Space of
 probability
 measures over \mathcal{D}

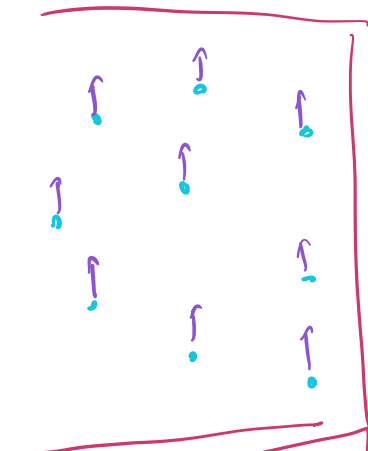
→ The model $\phi(\vec{\theta}; x) = \frac{1}{m} \sum_{j=1}^m g(\theta_j; x)$ becomes

$$\phi(\vec{\theta}; x) = \int_{\mathcal{D}} g(\theta; x) \hat{\mu}_m(d\theta)$$

non
 linear
 depends
 between
 ϕ and $\vec{\theta}$

→ ϕ is linear w.r.t. $\hat{\mu}_m$!!

→ The energy of the system expressed in terms of μ is



$$\mu = \delta_{\theta_0}$$

$$\int \chi(\theta) d\mu(\theta) = \chi(\theta_0)$$

$$\phi(x) = \int g(\theta, x) \mu(d\theta)$$

$$\mu_1 \left\{ \begin{array}{l} \mu = \frac{1}{2}(\mu_1 + \mu_2) \end{array} \right.$$

$$= \frac{2}{m} \sum_{j=1}^m F(\theta_j) + \frac{1}{m^2} \sum_{j,j'=1}^m K(\theta_j, \theta_{j'})$$

$$\mathcal{L}[\mu] = -2 \int_{\mathcal{D}} F(\theta) \mu(d\theta) + \int_{\mathcal{D} \times \mathcal{D}} K(\theta, \theta') \mu(d\theta) \mu(d\theta')$$

$$\phi(x) = \frac{1}{2}(\phi_1(x) + \phi_2(x))$$

$L(\mu)$ is now quadratic w.r.t. μ ; in fact it is convex w.r.t. μ (recall that K is psd).

Q: Too good to be true?

A: Not so fast: L is convex with respect to the geometry of linear mixtures.

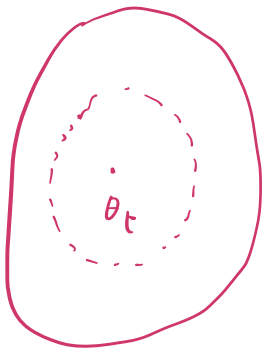
However, gradient descent dynamics correspond to a different geometry.

Q: What is the relationship between GD/GF and metric?

Proximal View point of GD.

The standard (Euclidean) GD step satisfies

$$\begin{aligned}\theta_{t+1} &= \theta_t - \eta \nabla f(\theta_t) \\ &= \operatorname{argmin}_{\theta} \left\{ \underbrace{f(\theta_t) + \langle \nabla f(\theta_t), \theta - \theta_t \rangle}_{\text{linear approx of } f \text{ at } \theta_t} + \underbrace{\frac{1}{2\eta} \|\theta - \theta_t\|_2^2}_{\text{proximity term}} \right\}\end{aligned}$$



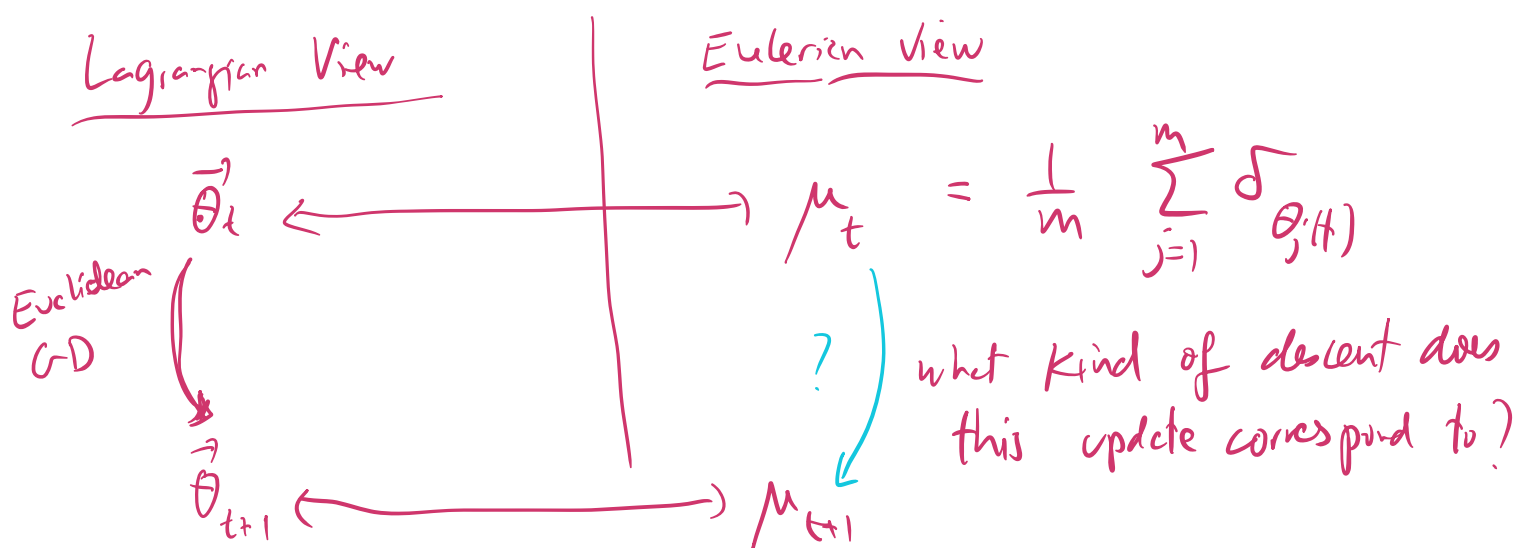
• How do we measure proximity? L_2 metric is one choice, but not always the "good" choice!

• The generalisation of GD on more general metrics is generally called Mirror Descent (Nemirovski, Yudin '83).

↳ In essence,
$$\theta_{t+1} = \operatorname{argmin}_{\theta} \left\{ f(\theta) + \frac{1}{2\eta} \underbrace{D(\theta; \theta_t)}_{\substack{\text{Divergence} \\ \text{(Bregman) function} \\ \text{measures proximity}}} \right\}$$

between $\dot{\theta}_t$ and θ .

Back at our measure space:



- Recall that $\dot{\theta}_j = +\nabla F(\theta_j) - \frac{1}{m} \sum_{j'=1}^m \nabla K(\theta_j, \theta_{j'})$
- Using $\mu_t = \frac{1}{m} \sum_j \delta_{\theta_j(t)}$, we obtain

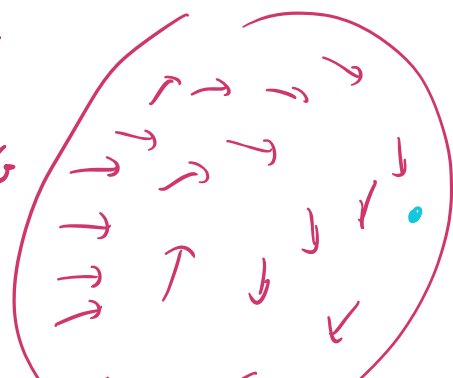
$$\dot{\theta}_j(t) = \nabla F(\theta_j(t)) - \int_D \nabla K(\theta_j(t), \theta) \mu_t(d\theta)$$

$$= -\nabla G(\theta_j(t); \mu_t), \text{ with}$$

$$G(\theta; \mu) = -F(\theta) + \int_D K(\theta, \theta') \mu(d\theta')$$

In physics, G is called the instantaneous potential of the system. Its gradient $\nabla G(\theta; \mu_t)$ defines a velocity field over D that any particle "feels".

→ Consider now a test function (smooth)
 $\chi: D \rightarrow \mathbb{R}$ and



$$\int_D \chi(\theta) \mu_t(d\theta) = \frac{1}{m} \sum_{j=1}^m \chi(\theta_j(t))$$

Time derivatives.

$$\begin{aligned} \int_D \chi(\theta) \partial_t \mu_t(d\theta) &= \frac{1}{m} \sum_{j=1}^m \nabla \chi(\theta_j(t)) \cdot \dot{\theta}_j(t) \\ &= -\frac{1}{m} \sum_{j=1}^m \langle \nabla \chi(\theta_j(t)), \nabla G(\theta_j(t); \mu_t) \rangle \\ &= - \int_D \langle \nabla \chi(\theta), \nabla G(\theta; \mu_t) \rangle \mu_t(d\theta) \end{aligned}$$

→ This is known as a continuity / transport equation.
(Mass is conserved).

→ The associated PDE is written as

$$\partial_t \mu_t = \operatorname{div} (\nabla G(\theta; \mu_t) \mu_t) \quad \boxed{\text{Liouville equation}}.$$

Q: Do these dynamics correspond also to a gradient flow?

A: Yes! The proximal interpretation is given in terms of a so-called Wasserstein Gradient Flow:

$$\mu_{t+1} = \operatorname{argmin}_{\mu \in P(D)} \left\{ \mathcal{L}[\mu] + \frac{1}{2\eta} W_2^2(\mu, \mu_t) \right\}$$

[JKO Scheme].
Jordan, Klein, Otto

Remarks :

→ This description is exact.

→ Add noise to the gradient updates (GF → Langevin dynam.)
translates into another PDE in the space of

measures (Liouville eq \rightarrow McKean-Vlasov eq).

This connection between Wasserstein gradient Flow and training of shallow NN was made by
 [Chizat, Bach] [Rotskoff & Vanden-Eijnden] [Mei, Montanari, Nguyen]
 [Sirignano & Spiliopoulos] - all in 2018.

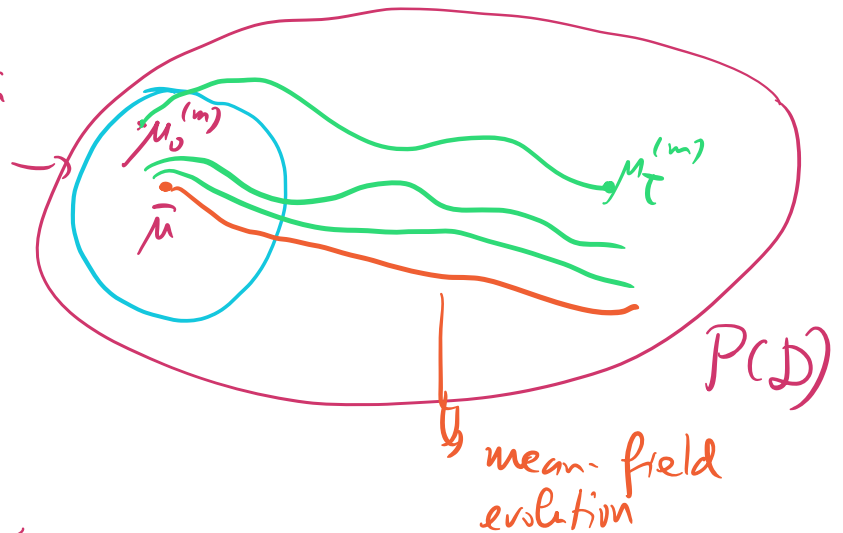
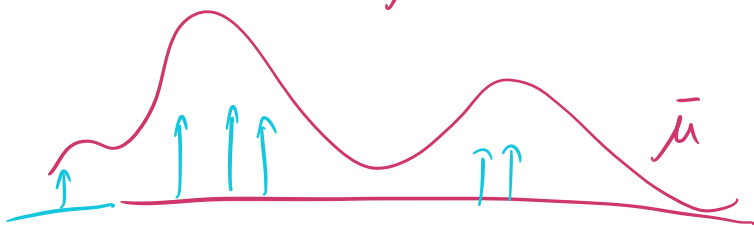
Mean-field regime

• Consider the evolution of the system as $m \rightarrow \infty$.

$\mu_t^{(m)}$: state after time t , with $\theta_j(0) \sim \bar{\mu}$

$$\mu_0^{(m)} = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j(0)}, \quad \theta_j(0) \sim \bar{\mu}$$

is the empirical measure associated with $\bar{\mu}$.



μ_t solves

$$\partial \mu_t = \text{div}(\nabla G(\theta; \mu_t) \mu_t)$$

with initial condition

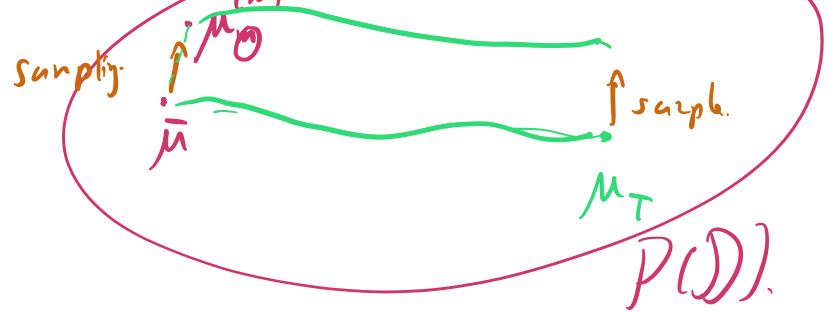
$$\mu_0 = \bar{\mu}$$

Theorem [CB, REVE, MMA, SS]

For any fixed $T > 0$, $\mu_T^{(m)}$ converges weakly to μ_T as $T \rightarrow \infty$, where μ_t solves.

In other words:

Dynamics and sampling commute in the limit of $m \rightarrow \infty$.



Two main questions

- (i) Under what conditions does this PDE converge to the global minimum of \mathcal{L} ? (convergence in time).
- (ii) How are the dynamics affected in terms of overperm. (convergence / fluctuations in m).

[Existing positive results for global convergence, but they are all qualitative in m and also in t (no rates!).