

# **Gradient Flows for Statistical Computation**

## **Trends and Trajectories**

Sam Power, University of Bristol

Statistics Seminar,  
University of Edinburgh, 16 November 2025

# Main ideas today

- Many statistical tasks reduce to solution of an optimisation problem
- Many common methods for these problems have ‘gradient’ structure
- Identifying these commonalities is useful for analysis, synthesis, progress

# Game Plan

- Describe a diverse variety of relevant statistical optimisation tasks
- Describe a consistent framework for solving them computationally
- Identify some ‘standard’ methods which come from this framework
  - ... and explain how some extensions can be derived
- Identify some *open questions* arising from these new methods

# Collaborators



**feel free to stop me at any point**

**ask me about references**

# **Examples of Statistical Optimisation Problems**

# Three Main Characters

- Optimisation over Parameter Spaces (“ $\Theta \subseteq \mathbf{R}^d$ ”)
- Optimisation over Measure Spaces (“ $\mathcal{P}(\mathcal{X})$ ”;  $\mathcal{X} \subseteq \mathbf{R}^d$ )
- Optimisation over ‘Hybrid’ Spaces (“ $\Theta \times \mathcal{P}(\mathcal{X})$ ”)

# Optimisation over Parameter Spaces

- Maximum Likelihood Estimation ('MLE')

$$\max_{\theta} \sum_{i \in [N]} \log p_{\theta}(y_i)$$

- maybe add a penalty term ('penalised MLE')
- maybe use a more general loss ('M-Estimation')

# Optimisation over Parameter Spaces

- Variational Approximation:

with  $\pi$  known up to a normalising constant,

$$\min_{\theta} \text{KL} (p_{\theta}, \pi)$$

- e.g.  $\theta = (m, C)$ ,  $p_{\theta}(dx) = \mathcal{N}(dx; m, C)$ ; “best Gaussian fit”

# Optimisation over Measure Spaces

- Sampling from an unnormalised distribution  $\pi \propto \exp(-V)$

$$\min_{\mu} \text{KL}(\mu, \pi)$$

$$\sim \min_{\mu} \left\{ \mathbf{E}_{\mu}[V] + \mathcal{H}(\mu) \right\}$$

$$\text{with } \mathcal{H}(\mu) = \int \mu \log \mu - \mu.$$

# Optimisation over Measure Spaces

- (Nonparametric) Mean-Field Approximation

$$\min_{\mu_1, \dots, \mu_d} \text{KL}(\mu_1 \otimes \dots \otimes \mu_d, \pi) \sim \min_{\mu} \left\{ \mathbf{E}_{\mu_1 \otimes \dots \otimes \mu_d}[V] + \sum_{i \in [d]} \mathcal{H}(\mu_i) \right\}$$

$$\text{with } \mathcal{H}(\mu) = \int \mu \log \mu - \mu.$$

# Optimisation over Measure Spaces

- ‘Quadratic Free Energy Minimisation’

$$\min_{\mu} \left\{ \mathbf{E}_{\mu}[V] + \frac{1}{2} \mathbf{E}_{\mu \otimes \mu}[W] + \mathcal{H}(\mu) \right\}$$

with  $\mathcal{H}(\mu) = \int \mu \log \mu - \mu$ .

- think of e.g. { covariance regularisation, kernel methods, ... }

# Optimisation over Hybrid Spaces

- Basic Example: Deconvolution
  - Model: draw  $X \sim p_\theta$ , but only observe  $Y \sim \mathcal{N}(X, \sigma^2)$
  - In principle, can ‘just’ do MLE ...
    - ... but here,  $p_\theta(y)$ ,  $\nabla_\theta \log p_\theta(y)$  are *likely unavailable*
  - Coupled problem: impute  $[x \mid \theta, y]$ , optimise  $[\theta \mid x; y]$
  - More generally: “*EM Algorithm*”, “*Latent Variable Models*”

# Optimisation over Hybrid Spaces

- { ‘*Energy-Based*’ / ‘*Unnormalised*’ / ‘*Pre-Normalised*’ } Models
  - Specify  $p_\theta(y) \propto \exp(-V(y; \theta))$ ; leave  $Z(\theta)$  defined implicitly
  - In principle, can ‘just’ do MLE ...
    - ... but here,  $p_\theta(y)$ ,  $\nabla_\theta \log p_\theta(y)$  are *likely unavailable*
  - Coupled problem: sample  $x \sim p_\theta$ , then optimise  $\theta$  based on  $x, y$ 
    - “*Contrastive Divergence*”, “*MC-MLE*”

# Additional Comments on Hybrid Spaces

- Increasingly, clear that many problems have this two-scale structure
  - Adaptive MCMC, Distributed Inference, ... (ask me!)
  - see also “MCMC-Driven Learning” chapter by Bouchard-Côté++
- IMO: worthy of serious attention

**last chance to ask about examples**

# Metric Structures in Statistical Optimisation

# Metrics

- Nothing too fancy - just want enough structure to ‘do good calculus’
- For parameter optimisation,  $\Theta \subseteq \mathbf{R}^d$  with Euclidean metric.
- For measure optimisation,  $\mathcal{P}(\mathcal{X})$  with 2-Wasserstein metric.
- For hybrid optimisation,  $\Theta \times \mathcal{P}(\mathcal{X})$  with ‘hybrid’ metric

$$d_{\text{hyb}} \left( (\theta, \mu), (\theta', \mu') \right) = \sqrt{\|\theta - \theta'\|_2^2 + \mathcal{T}_2^2(\mu, \mu')}$$

# **Optimisation on Metric Spaces**

# Conceptual Optimisation Framework

OPT

$$\min_{x \in \mathcal{X}} f(x)$$

PPM

$$x_0 \mapsto \arg \min_{x \in \mathcal{X}} \left\{ f(x) + \frac{1}{2h} \cdot d(x, x_0)^2 \right\}$$

FLOW

$$\dot{x}_t = -\nabla f(x_t)$$

**Specify a Metric Structure**

**Receive an Optimisation Algorithm**

# A Word on Gradients

- For parametric optimisation, the gradient is the ‘usual’ gradient
- For measure-valued optimisation, the gradient is a ‘Wasserstein’ gradient
- For hybrid optimisation, you take part of each, ‘as expected’

# Wasserstein Gradients: Some Intuition

- Suppose that we are interested in a functional  $\mathcal{F} : \mathcal{P}(\mathcal{X}) \rightarrow \mathbf{R}$
- Assume a Taylor expansion along  $\mu \rightsquigarrow \mu'$

$$\frac{\mathcal{F}((1-t) \cdot \mu + t \cdot \mu') - \mathcal{F}(\mu)}{t} \approx \int (\delta_\mu \mathcal{F})(\mu, x) \cdot \left\{ \mu'(\mathrm{d}x) - \mu(\mathrm{d}x) \right\}$$

- Decrease  $\mathcal{F}$  by pushing mass towards minima of  $(\delta_\mu \mathcal{F})(\mu, \cdot)$ , i.e.

$$\dot{X}_t = -\nabla_x \delta_\mu \mathcal{F}(\mu_t, X_t)$$

# Wasserstein Gradients: An Example

- For  $\mathcal{V}(\mu) = \mathbf{E}_\mu[V]$ , obtain

$$\delta_\mu \mathcal{V}(\mu, x) = V(x)$$

$$\nabla_{\mathcal{T}} \mathcal{V}(\mu, x) = \nabla V(x)$$

# Wasserstein Gradients: And Another

- For  $\mathcal{W}(\mu) = \frac{1}{2} \int \mu(dx) \mu(dy) W(x, y)$ , obtain

$$\delta_\mu \mathcal{W}(\mu, x) = \int \mu(dy) \cdot W(x, y)$$

$$\nabla_{\mathcal{T}} \mathcal{W}(\mu, x) = \int \mu(dy) \cdot \nabla_x W(x, y)$$

# A Special Case: Entropy

- Recall

$$\mathcal{H}(\mu) = \int (\mu \log \mu - \mu) \quad \rightsquigarrow \quad \delta_\mu \mathcal{H}(\mu, x) = \log \mu(x)$$

- Hence, descend  $\mathcal{H}$  by microscopically evolving

$$\dot{X}_t = -\nabla_x \log \mu_t(X_t)$$

- At the *population level*, equivalent to do

$$dX_t = \sqrt{2} dW_t$$

- “*The gradient flow of the entropy can be realised by Brownian motion*”

# **From Gradient Flows to Algorithms**

# Gradient Flows on Parameter Spaces

$$\min_{\theta \in \Theta} f(\theta)$$

$$\dot{\theta}_t = - \nabla_{\theta} f(\theta_t)$$

“*Gradient Method*”

# Gradient Flows on Parameter Spaces

$$\max_{\theta \in \Theta} \log p_\theta(y)$$

$$\dot{\theta}_t = \nabla_\theta \log p_{\theta_t}(y)$$

*“MLE by Gradient Flow”*

# Gradient Flows on Parameter Spaces

- Time-Discretised Method: “*Gradient Method*”
- Extremely well-understood when  $f$  is uniformly-convex
- Further theory for ‘gradient-dominated’  $f$ ; ‘*Polyak-Łojasiewicz Inequality*’

# Gradient Flows on Measure Spaces

$$\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\mu) + \mathcal{H}(\mu) \right\}$$

$$\partial_t \mu_t = - \nabla_{\mathcal{T}, \mu} \left\{ \mathcal{F} + \mathcal{H} \right\} (\mu_t)$$

# Gradient Flows on Measure Spaces

$$\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\mu) + \mathcal{H}(\mu) \right\}$$

$$\partial_t \mu_t = \operatorname{div}_x \left( \mu_t \nabla_x \delta_\mu \mathcal{F}(\mu_t) \right) + \Delta_x \mu_t$$

# Gradient Flows on Measure Spaces

$$\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\mu) + \mathcal{H}(\mu) \right\}$$

$$dX_t = -\nabla_x \delta_\mu \mathcal{F}(\mu_t, X_t) dt + \sqrt{2} dW_t$$

$$\mu_t = \text{Law}\left(X_t\right)$$

# Gradient Flows on Measure Spaces

$$\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\mu) + \mathcal{H}(\mu) \right\}$$

$$dX_t^i = -\nabla_x \delta_\mu \mathcal{F}\left(\hat{\mu}_t^N, X_t^i\right) dt + \sqrt{2} dW_t^i$$

$$\hat{\mu}_t^N = \text{Emp\_Meas}\left(\{X_t^1, \dots, X_t^N\}\right)$$

# Gradient Flows on Measure Spaces

$$\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathbf{E}_\mu [V] + \frac{1}{2} \mathbf{E}_{\mu \otimes \mu} [W] + \mathcal{H}(\mu) \right\}$$

$$dX_t = -\nabla V(X_t) dt - \left( \int \nabla_1 W(X_t, y) \mu_t(dy) \right) dt + \sqrt{2} dW_t$$

$$\mu_t = \text{Law}(X_t)$$

# Gradient Flows on Measure Spaces

$$\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathbf{E}_{\mu}[V] + \frac{1}{2} \mathbf{E}_{\mu \otimes \mu}[W] + \mathcal{H}(\mu) \right\}$$

$$\mathrm{d} X_t^i = -\nabla V\!\left(X_t^i\right)\,\mathrm{d} t-\frac{1}{N-1}\sum_{j\neq i}\nabla_1 W\!\left(X_t^i,X_t^j\right)\,\mathrm{d} t+\sqrt{2}\mathrm{d} W_t^i$$

# Gradient Flows on Measure Spaces

- **Space-Time-Discretised Method:** (Mean-Field) “*Langevin Monte Carlo*”
- Extremely well-understood when  $\mathcal{F}$  is uniformly-geodesically-convex
- Further theory for ‘well-connected’  $\mathcal{F}$ ; “*Sobolev-type*” inequalities

# Gradient Flows on Hybrid Spaces

$$\min_{\theta \in \Theta, \mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\theta, \mu) + \mathcal{H}(\mu) \right\}$$

$$\dot{\theta}_t = -\nabla_{\theta}\mathcal{F}(\theta_t, \mu_t)$$

$$\partial_t \mu_t = -\nabla_{\mathcal{T}, \mu} \left\{ \mathcal{F} + \mathcal{H} \right\} (\theta_t, \mu_t)$$

# Gradient Flows on Hybrid Spaces

$$\min_{\theta \in \Theta, \mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\theta, \mu) + \mathcal{H}(\mu) \right\}$$

$$\dot{\theta}_t = -\nabla_\theta \mathcal{F}(\theta_t, \mu_t)$$

$$\partial_t \mu_t = \operatorname{div}_x \left( \mu_t \nabla_x \delta_\mu \mathcal{F}(\theta_t, \mu_t) \right) + \Delta_x \mu_t$$

# Gradient Flows on Hybrid Spaces

$$\min_{\theta \in \Theta, \mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\theta, \mu) + \mathcal{H}(\mu) \right\}$$

$$\dot{\theta}_t = -\nabla_\theta \mathcal{F}(\theta_t, \mu_t)$$

$$dX_t = -\nabla_x \delta_\mu \mathcal{F}(\theta_t, \mu_t, X_t) dt + \sqrt{2} dW_t$$

$$\mu_t = \text{Law}(X_t)$$

# Gradient Flows on Hybrid Spaces

$$\min_{\theta \in \Theta, \mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\theta, \mu) + \mathcal{H}(\mu) \right\}$$

$$\dot{\theta}_t = -\nabla_{\theta}\mathcal{F}(\theta_t, \hat{\mu}_t^N)$$

$$dX_t^i = -\nabla_x \delta_\mu \mathcal{F}(\theta_t, \hat{\mu}_t^N, X_t^i) \, dt + \sqrt{2} dW_t^i$$

$$\hat{\mu}_t^N = \texttt{Emp\_Meas}\left(\{X_t^1, \dots, X_t^N\}\right)$$

# Gradient Flows on Hybrid Spaces

$$\max_{\theta \in \Theta, \mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathbf{E}_{\mu} \left[ \log p_{\theta}(X, y) \right] - \mathcal{H}(\mu) \right\}$$

$$\dot{\theta}_t = \int \nabla_{\theta} \log p_{\theta_t}(x, y) \, \mu_t(dx)$$

$$dX_t = \nabla_x \log p_{\theta_t}(X_t, y) \, dt + \sqrt{2} dW_t$$

$$\mu_t = \text{Law}(X_t)$$

# Gradient Flows on Hybrid Spaces

$$\max_{\theta \in \Theta, \mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathbf{E}_{\mu} \left[ \log p_{\theta}(X, y) \right] - \mathcal{H}(\mu) \right\}$$

$$\dot{\theta}_t = \frac{1}{N} \sum_i \nabla_{\theta} \log p_{\theta_t}(X_t^i, y)$$

$$dX_t^i = \nabla_x \log p_{\theta_t}(X_t^i, y) dt + \sqrt{2} dW_t^i$$

# Gradient Flows on Hybrid Spaces

- Space-Time-Discretised Method: “*Particle Gradient Descent*”
- Very well-understood when  $\mathcal{F}$  is uniformly-geodesically-convex
- Further theory currently missing / in development; Open Questions

**questions about ‘core’ methods?**

# **Some New Directions**

# Beyond ‘Standard’ Gradient Flows

- Typically, the ‘standard’ gradient flow is sub-optimal.
- This is true in both continuous and in discrete time.
- Some intuition has developed for how ‘optimal’ improvements look.
- A common (though not universal) theme seems to involve ‘momentum’.

# Beyond ‘Standard’ Gradient Flows

- We start off with some notion of an objective function.
- We then identify dynamics which can minimise that function for us.
- How shall this pipeline change when ‘momentum’ enters the picture?
  - *New objective function, new dynamics.*
  - Minimisers of new objective map onto minimisers of old objective.

# Enriched Objective Functions

- For parameter optimisation, consider

$$\min_{(\theta, \varphi) \in \mathcal{T}^{\star\Theta}} \left\{ h(\theta, \varphi) := f(\theta) + \frac{1}{2} \cdot \|\varphi\|_2^2 \right\}$$

# Enriched Objective Functions

- For measure optimisation, consider

$$\min_{\nu \in \mathcal{P}(\mathcal{T}^\star \mathcal{X})} \left\{ H(\nu) := \mathcal{F}(\nu) + \mathcal{H}(\nu) + \mathbf{E}_\nu \left[ \frac{1}{2} \cdot \|P\|^2 \right] \right\}$$

# Enriched Objective Functions

- For hybrid optimisation, consider

$$\min_{(\theta, \varphi), \nu} \left\{ H(\theta, \varphi, \nu) := \mathcal{F}(\theta, \nu) + \mathcal{H}(\nu) + \frac{1}{2} \cdot \|\varphi\|_2^2 + E_\nu \left[ \frac{1}{2} \cdot \|P\|^2 \right] \right\}$$

# Enriched Objective Functions

$$\min_{(\theta,\varphi) \in \mathcal{T}^{\star}\Theta} \left\{ h\left(\theta,\varphi\right) := f(\theta) + \frac{1}{2} \cdot \|\varphi\|_2^2 \right\}$$

$$\min_{\nu \in \mathcal{P}(\mathcal{T}^{\star}\mathcal{X})} \left\{ H(\nu) := \mathcal{F}(\nu) + \mathcal{H}(\nu) + \mathbf{E}_{\nu}\left[\frac{1}{2} \cdot \|P\|^2\right] \right\}$$

$$\min_{(\theta,\varphi),\nu} \left\{ \mathsf{H}\left(\theta,\varphi,\nu\right) := \mathcal{F}(\theta,\nu) + \mathcal{H}(\nu) + \frac{1}{2} \cdot \|\varphi\|_2^2 + \mathbf{E}_{\nu}\left[\frac{1}{2} \cdot \|P\|^2\right] \right\}$$

# Warm-Up: Hamiltonian Flows

- An odd idea in isolation:
  - don't descend the objective, but instead conserve the '*Hamiltonian*'
  - Introduce skew-symmetric matrix

$$\mathbf{J} = \begin{pmatrix} 0 & \mathbf{I} \\ -\mathbf{I} & 0 \end{pmatrix}$$

- In abstract terms: with  $z = (x, p)$

go from  $\dot{x} = -\nabla f(x)$  to  $\dot{z} = \mathbf{J} \nabla H(z),$

# Hamiltonian Flows in Action

- For parameter optimisation, obtain

$$\dot{\theta}_t = \varphi_t, \quad \dot{\varphi}_t = -\nabla f(\theta_t)$$

# Hamiltonian Flows in Action

- For measure optimisation, obtain (omitting entropy term)

$$dX_t = P_t dt, \quad dP_t = -\nabla_x \delta_\mu \mathcal{F}(\mu_t, X_t) dt$$

with ‘expected’ modifications for particle version

# Hamiltonian Flows in Action

- For hybrid optimisation, obtain

$$\dot{\theta}_t = \varphi_t,$$

$$\dot{\varphi}_t = -\nabla_{\theta}\mathcal{F}(\theta_t, \mu_t)$$

$$dX_t = P_t dt,$$

$$dP_t = -\nabla_x \delta_\mu \mathcal{F}(\theta_t, \mu_t, X_t) dt$$

with ‘expected’ modifications for particle version

**okay, but why bother with this?**

# *Conformal Hamiltonian Flows*

- Experience suggests to (linearly) blend
  - Hamiltonian circulation, and
  - gradient-type damping *only on the momentum term*
- This routinely yields improved methods

# *Conformal Hamiltonian Flows*

- Blend Hamiltonian circulation with gradient-type *momentum damping*
- The key matrix is then (for some  $\gamma > 0$ )

$$D_\gamma = \begin{pmatrix} 0 & -I \\ I & \gamma \cdot I \end{pmatrix}$$

and we will (formally) construct dynamics according to

$$\dot{z} = -D_\gamma \nabla H(z)$$

**so what does this look like?**

# Conformal Hamiltonian Flows in Action

- For parameter optimisation, obtain

$$\dot{\theta}_t = \varphi_t, \quad \dot{\varphi}_t = -\nabla f(\theta_t) - \gamma \cdot \varphi_t$$

- $\approx$  Nesterov's “*Fast Gradient Method*”

# Conformal Hamiltonian Flows in Action

- For measure optimisation, obtain

$$dX_t = P_t dt, \quad dP_t = -\nabla_x \delta_\mu \mathcal{F}(\mu_t, X_t) dt - \gamma \cdot P_t dt + \sqrt{2 \cdot \gamma} dW_t$$

with ‘expected’ modifications for particle version

- $\approx$  (*Kinetic, Underdamped, ...*) Langevin Monte Carlo

# Conformal Hamiltonian Flows in Action

- For hybrid optimisation, obtain

$$\begin{aligned}\dot{\theta}_t &= \varphi_t, & \dot{\varphi}_t &= -\nabla_{\theta}\mathcal{F}(\theta_t, \mu_t) - \gamma \cdot \varphi_t \\ dX_t &= P_t dt, & dP_t &= -\nabla_x \delta_{\mu} \mathcal{F}(\theta_t, \mu_t, X_t) dt - \gamma \cdot P_t dt + \sqrt{2 \cdot \gamma} dW_t\end{aligned}$$

with ‘expected’ modifications for particle version

- $\approx$  Momentum Particle Gradient Descent

# **Recap and Open Questions**

# Main ideas today

- Optimisation problems are widespread in statistical tasks
  - ... and often involve more than ‘just’ fixed-dimensional parameters.
- It is often possible to solve such problems “with gradient descent”
  - ... and we can even systematically concoct improvements on GD.
- Identifying these commonalities is useful for analysis, synthesis, progress
  - ... and many interesting questions still remain.

# Some Open Questions

- For optimisation problems on hybrid spaces,
  - Can we strengthen the theory *outside of the uniformly-convex case*?
  - Can we develop good principles for *numerical discretisation*?
  - What more shall be learned from “*pure*” optimisation and sampling?
- For momentum-enrichment,
  - How should we *systematically* construct ‘enriched’ objective functions?

# Some Further Questions

- In general,
  - Which other tasks can be fruitfully *interpreted through optimisation*?
  - Should we ever look *beyond* gradient and conformal Hamiltonian flows?

---

# Particle algorithms for maximum likelihood training of latent variable models

---

Juan Kuntz

Jen Ning Lim

Adam M. Johansen

Department of Statistics, University of Warwick.

## Abstract

Neal and Hinton (1998) recast maximum likelihood estimation of any given latent variable model as the minimization of a free energy functional  $F$ , and the EM algorithm as coordinate descent applied to  $F$ . Here, we explore alternative ways to optimize the functional. In particular, we identify various gradient flows associated with  $F$  and show that their limits coincide with  $F$ 's stationary points. By discretizing the flows, we obtain practical particle-based algorithms for maximum likelihood estimation in broad classes of latent variable models. The novel algorithms scale to high-dimensional settings and perform well in numerical experiments.

(S2) obtain the corresponding *posterior distribution*,

$$p_{\theta_*}(x|y) := \frac{p_{\theta_*}(x, y)}{p_{\theta_*}(y)}.$$

Perhaps the most well-known method for tackling (S1,2) is the *expectation maximization* (EM) algorithm (Dempster et al., 1977): starting from an initial guess  $\theta_0$ , alternate,

(E) compute  $q_k := p_{\theta_k}(\cdot|y)$ ,

(M) solve for  $\theta_{k+1} := \arg \max_{\theta \in \Theta} \int \ell(\theta, x) q_{k+1}(x) dx$ ,

where  $\ell(\theta, x) := \log(p_\theta(x, y))$  denotes the log-likelihood. Under general conditions (McLachlan, 2007, Chap. 3),  $\theta_k$  converges to a stationary point  $\theta_*$  of the marginal likelihood and  $q_k$  to the corresponding posterior  $p_{\theta_*}(\cdot|y)$ . In cases where the above steps are not analytically tractable, it is common to approximate (E) using Monte Carlo (or Markov chain Monte Carlo if  $p_\theta(\cdot|y)$  cannot be sampled

---

# Momentum Particle Maximum Likelihood

---

Jen Ning Lim<sup>1</sup> Juan Kuntz<sup>2</sup> Samuel Power<sup>3</sup> Adam M. Johansen<sup>1</sup>

## Abstract

Maximum likelihood estimation (MLE) of latent variable models is often recast as the minimization of a free energy functional over an extended space of parameters and probability distributions. This perspective was recently combined with insights from optimal transport to obtain novel particle-based algorithms for fitting latent variable models to data. Drawing inspiration from prior works which interpret ‘momentum-enriched’ optimization algorithms as discretizations of ordinary differential equations, we propose an analogous dynamical-systems-inspired approach to minimizing the free energy functional. The result is a dynamical system that blends elements of Nesterov’s Accelerated Gradient method, the under-damped Langevin diffusion, and particle methods. Under suitable assumptions, we prove that the continuous-time system minimizes the functional. By discretizing the system, we obtain a practical algorithm for MLE in latent variable models. The algorithm outperforms existing particle methods in numerical experiments and compares favourably with other MLE algorithms.

by constructing an objective defined over an extended space, whose optima are in one-to-one correspondence with those of the MLE problem. To this end, we define the ‘free energy’ functional:

$$\mathcal{E}(\theta, q) := \int \log\left(\frac{q(x)}{p_\theta(y, x)}\right) q(x) dx. \quad (1)$$

The infimum of  $\mathcal{E}$  over the *extended space*  $\mathbb{R}^{d_\theta} \times \mathcal{P}(\mathbb{R}^{d_x})$  (where  $\mathcal{P}(\mathbb{R}^{d_x})$  denotes the space of probability distributions over  $\mathbb{R}^{d_x}$ ) coincides with negative of the log marginal likelihood’s supremum:

$$\mathcal{E}^* := \inf_{(\theta, q)} \mathcal{E}(\theta, q) = - \sup_{\theta \in \mathbb{R}^{d_\theta}} \log p_\theta(y). \quad (2)$$

To see this, note that

$$\mathcal{E}(\theta, q) = -\log p_\theta(y) + \text{KL}(q, p_\theta(\cdot | y)), \quad (3)$$

where  $\text{KL}$  denotes the Kullback–Leibler divergence and  $p_\theta(\cdot | y) := p_\theta(y, \cdot) / p_\theta(y)$  the posterior distribution. So, if  $(\theta^*, q^*)$  minimizes the free energy, then  $\theta^*$  maximizes the marginal likelihood and  $q^*(\cdot) = p_{\theta^*}(\cdot | y)$ . In other words, by minimizing the free energy, we solve our MLE problem.

This perspective motivates the search for practical procedures that minimize the free energy  $\mathcal{E}$ . One such example

# Error Bounds for Particle Gradient Descent, and Extensions of the log-Sobolev and Talagrand Inequalities

**Rocco Caprio**

*Department of Statistics  
University of Warwick  
Coventry, CV4 7AL, UK*

ROCCO.CAPRIO@WARWICK.AC.UK

**Juan Kuntz**

**Samuel Power**  
*School of Mathematics  
University of Bristol  
Bristol, BS8 1UG, UK*

SAM.POWER@BRISTOL.AC.UK

**Adam M. Johansen**

*Department of Statistics  
University of Warwick  
Coventry, CV4 7AL, UK*

A.M.JOHANSEN@WARWICK.AC.UK

**Editor:** Alexandre Bouchard

## Abstract

We derive non-asymptotic error bounds for particle gradient descent (PGD, Kuntz et al. (2023)), a recently introduced algorithm for maximum likelihood estimation of large latent variable models obtained by discretizing a gradient flow of the free energy. We begin by showing that the flow converges exponentially fast to the free energy's minimizers for models satisfying a condition that generalizes both the log-Sobolev and the Polyak–Lojasiewicz inequalities (LSI and PLI, respectively). We achieve this by extending a result well-known in the optimal transport literature (that the LSI implies the Talagrand inequality) and its counterpart in the optimization literature (that the PLI implies the so-called quadratic growth condition), and applying the extension to our new setting. We also generalize the Bakry–Émery Theorem and show that the LSI/PLI extension holds for models with strongly concave log-likelihoods. For such models, we further control PGD's discretization error and obtain the non-asymptotic error bounds. While we are motivated by the study of PGD, we believe that the inequalities and results we extend may be of independent interest.

**Keywords:** latent variable models, maximum marginal likelihood, gradient flows, log-Sobolev inequality, Polyak–Lojasiewicz inequality, Talagrand inequality, quadratic growth condition.

# Some Take-Aways

- Optimisation problems are widespread in statistical tasks
  - ... and often involve more than ‘just’ fixed-dimensional parameters.
- It is often possible to solve such problems “with gradient descent”
  - ... and we can even systematically concoct improvements on GD.
- Identifying these commonalities is useful for analysis, synthesis, progress
  - ... and many interesting questions still remain.