

Gradient Flows for Statistical Computation

Trends and Trajectories

Sam Power, University of Bristol

“Advanced Langevin Methods for Bayesian Sampling”, BayesComp,
National University Of Singapore, 19 June 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

Collaborators



feel free to stop me at any point

Statistical Inference

Statistical Computing

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 incorporate a penalty term ('penalised MLE')
 - maybe use a more general loss ('M-Estimation')
- Variational Approximation : $\min_{\theta} \text{KL} (p_{\theta}, \pi)$
 - e.g. $\theta = (\mathbf{m}, \mathbf{C})$, $p_{\theta}(\mathrm{d}x) = \mathcal{N} (\mathrm{d}x; \mathbf{m}, \mathbf{C})$; "best Gaussian fit"

Optimisation over Measure Spaces

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

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

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

- (Nonparametric) Mean-Field Approximation

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

- (integral probability metrics, information-theoretic divergences, etc.)

Optimisation over Hybrid Spaces

- Latent Variable Models: sample $[x \mid \theta, y]$, optimise $[\theta \mid x; y]$ (EM)
- Unnormalised Models: sample $[x \mid \theta]$, optimise θ (CD / MC-MLE)
- Distributed Inference: sample local posterior, tilt for consensus (\sim EP)
- Opinion: more prevalent than you might expect; worth taking seriously

press me on the examples

Optimisation by Local Search

Optimisation in Metric Spaces

Metrics

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

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

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 that it carries a Taylor expansion along the path $\mu \rightsquigarrow \mu'$

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

- Natural to decrease \mathcal{F} by pushing mass towards minima of $\left(\delta_{\mu} \mathcal{F} \right) (\mu, \cdot)$, 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 $\mathcal{W}(\mu) = \frac{1}{2} \int \mu(\mathrm{d}x) \mu(\mathrm{d}y) W(x, y)$, obtain

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

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

A Special Case: Entropy

- A particularly special functional is the (shifted, negative) entropy

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

- This satisfies $\delta_\mu \mathcal{H}(\mu, x) = \log \mu(x)$, so we could decrease it by evolving

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

- Remarkably, the same path of measures is induced by instead evolving stochastically by

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

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

How do the Algorithms Look?

Gradient Flows on Parameter Spaces

- Task: $\min_{\theta \in \Theta} f(\theta)$
- ODE:

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

- Time-Discretised Method: “Gradient Method”

Gradient Flows on Measure Spaces (1)

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

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

Gradient Flows on Measure Spaces (2)

- Task: $\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\mu) + \mathcal{H}(\mu) \right\}$
- Nonlinear ('McKean-Vlasov') SDE:

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

Gradient Flows on Measure Spaces (3)

- Task: $\min_{\mu \in \mathcal{P}(\mathcal{X})} \left\{ \mathcal{F}(\mu) + \mathcal{H}(\mu) \right\}$
- Particle SDE:

$$dX_t^i = - \nabla_x \delta_\mu \mathcal{F}(\hat{\mu}_t^N, X_t^i) dt + \sqrt{2} dW_t^i, \quad 1 \leq i \leq N$$

Gradient Flows on Measure Spaces (3)

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

- Particle SDE:

$$dX_t^i = - \nabla_x \delta_\mu \mathcal{F}(\hat{\mu}_t^N, X_t^i) dt + \sqrt{2} dW_t^i, \quad 1 \leq i \leq N$$

- **Space-Time-Discretised Method:** (Mean-Field) “Langevin Monte Carlo”

Gradient Flows on Hybrid Spaces (1)

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

- ODE-PDE:

$$\begin{aligned}\dot{\theta}_t &= -\nabla_{\theta} \mathcal{F}(\theta_t, \mu_t) \\ \partial_t \mu_t &= -\nabla_{\mathcal{T}, \mu} \left\{ \mathcal{F} + \mathcal{H} \right\}(\theta_t, \mu_t)\end{aligned}$$

Gradient Flows on Hybrid Spaces (2)

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

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

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

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

Gradient Flows on Hybrid Spaces (3)

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

- ODE-pSDE:

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

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

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

Gradient Flows on Hybrid Spaces (4)

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

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

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

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

- Space-Time-Discretised Method: “Particle Gradient Descent”

Some Words on Theory

- In each case, the theoretical picture is very clear for “convex” problems
- In each case, there exists a ‘robust’ notion of convexity / connectedness which yields guarantees for a larger class of problems
- These notions are ...
 - quite well-developed on Θ ,
 - very well-developed on $\mathcal{P}(\mathcal{X})$, and
 - still under development for hybrid spaces

Some Words on Extensions

- Today, I have really spoken about ‘standard’ methods in this area
- In many applications, 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’.

Momentum-Enriched Dynamics

- A common theme involves ‘momentum’.
- Some pattern-matching indicates how to introduce this momentum.
 - Introduce a momentum and a suitable extended objective functional
 - Instead of the gradient flow, work with the *Conformal Hamiltonian Flow*
- My discussion here is brief and informal, but can be made formal.

Momentum-Enriched Algorithms

- Some pattern-matching indicates how to introduce this momentum.
 - For parameter optimisation, recover *Nesterov's Fast Gradient Method*
 - For measure optimisation, recover *Kinetic Langevin Monte Carlo*
 - For hybrid case, we propose *Momentum Particle Gradient Descent*
- Promising mixture of theoretical and practical benefits

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.

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 θ_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), θ_k converges to a stationary point θ_* 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 Lim¹ Juan Kuntz² Samuel Power³ Adam M. Johansen¹

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 underdamped 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 KL denotes the Kullback–Leibler divergence and $p_\theta(\cdot | y) := p_\theta(y, \cdot) / p_\theta(y)$ the posterior distribution. So, if (θ^*, q^*) minimizes the free energy, then θ^* 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[†] Juan Kuntz[‡] Samuel Power[§] Adam M. Johansen[†]

April 12, 2024

Abstract

We prove non-asymptotic error bounds for particle gradient descent (PGD) [26], 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, for models satisfying a condition generalizing both the log-Sobolev and the Polyak–Łojasiewicz inequalities (LSI and PLI, respectively), the flow converges exponentially fast to the set of minimizers of the free energy. 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 it to our new setting. We also generalize the Bakry–Émery Theorem and show that the LSI/PLI generalization holds for models with strongly concave log-likelihoods. For such models, we further control PGD’s discretization error, obtaining 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–Łojasiewicz inequality, Talagrand inequality, quadratic growth condition.

Error bounds for particle gradient descent, and extensions of the log-Sobolev and Talagrand inequalities

Rocco Caprio[†] Juan Kuntz[‡] Samuel Power[§] Adam M. Johansen[†]

April 12, 2024

Abstract

We prove non-asymptotic error bounds for particle gradient descent (PGD) [26], 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, for models satisfying a condition generalizing both the log-Sobolev and the Polyak–Łojasiewicz inequalities (LSI and PLI, respectively), the flow converges exponentially fast to the set of minimizers of the free energy. 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 it to our new setting. We also generalize the Bakry–Émery Theorem and show that the LSI/PLI generalization holds for models with strongly concave log-likelihoods. For such models, we further control PGD’s discretization error, obtaining 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–Łojasiewicz 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.