

Gradient Flows for Statistical Computation

Trends and Trajectories

Sam Power, University of Bristol

London Symposium on Information Theory,
University of Cambridge, 15 May 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

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')
- Variational Approximation : $\min_{\theta} \text{KL} (p_{\theta}, \pi)$
 - e.g. $\theta = (m, C)$, $p_{\theta}(\mathrm{d}x) = \mathcal{N}(\mathrm{d}x; m, C)$; "best Gaussian fit"

Optimisation over Measure Spaces

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

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

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

- (Nonparametric) Mean-Field Approximation

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

- (other objectives involving integral probability metrics, density ratio divergences, etc.)

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”

More on 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 θ 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 (sample from π , optimise parameters of dynamics)
 - Distributed Inference (sample ‘locally’, ‘tilt parameters’ for consensus)
 - See also “MCMC-Driven Learning” chapter by Bouchard-Côté+++
 - “Markovian Optimisation-Integration” framework
- IMO: Worthy of serious consideration; not just hypothetical / edge case.

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, $\Theta \subseteq \mathbf{R}^d$ 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

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

Optimisation on Metric Spaces

Conceptual Optimisation Framework

1. Abstract Optimisation Problem

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

2. Proximal Point Update

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

3. Gradient Flow

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

4. (Discretisation)

From Gradient Flows to Algorithms

Gradient Flows on Parameter Spaces

- Task: $\min_{\theta \in \Theta} f(\theta)$
- Continuous Dynamics: $\dot{\theta}_t = -\nabla_{\theta} f(\theta_t)$
- Time-Discretised Method: “Gradient Method”
- Extremely well-understood for uniformly-convex f
- Further theory when $f - \inf f \lesssim \|\nabla f\|^2$; ‘Polyak-Łojasiewicz Inequality’

Gradients on Measure Spaces?

- Suppose that we are interested in a functional $\mathcal{F} : \mathcal{P}(\mathcal{X}) \rightarrow \mathbf{R}$
- Assume that it carries the Taylor expansion

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

- Natural to decrease \mathcal{F} by pushing mass towards minima of $\left(\delta_{\mu} \mathcal{F} \right) (\mu, \cdot)$, i.e.

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

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”

Gradient Flows on Measure Spaces

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

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

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

- **Space-Time-Discretised Method:** (Mean-Field) “Langevin Monte Carlo”
- Extremely well-understood for uniformly-geodesically-convex \mathcal{F}
- Further theory for ‘well-connected’ \mathcal{F} ; e.g. Logarithmic Sobolev Inequalities

Gradient Flows on Hybrid Spaces

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

- Continuous Dynamics:

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

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

- Space-Time-Discretised Method: “Particle Gradient Descent”
- Very well-understood for uniformly-geodesically-convex \mathcal{F}
- Further theory currently missing; Open Questions

questions on the 'basic' methods?

Some New Directions

Beyond ‘Standard’ Gradient Flows

- 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’.
 - “lifting the problem to the cotangent bundle”

Enriched Objective Functions

- For parameter optimisation, consider

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

- For measure optimisation, consider

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

More Enriched Objective Functions

- For hybrid optimisation, consider

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

- N.B. These choices are not known to be “automatic” / “canonical”, but appear to make sense in many examples.

Warm-Up: Hamiltonian Flows

- A sketchy idea in isolation: conserve the ‘Hamiltonian’
- Introduce skew-symmetric matrix

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

- In abstract terms: instead of

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

take $z = (x, p)$ and do

$$\dot{z} = \mathbf{J} \nabla H(z)$$

Hamiltonian Flows in Action

- For parameter optimisation, obtain

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

- For measure optimisation, obtain (omitting entropy term)

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

- For hybrid optimisation, obtain

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

- But ... *why bother?*

Conformal Hamiltonian Flows

- A recurrent phenomenon: it can be interesting to blend Hamiltonian circulation with gradient-type damping *only on the momentum term*
- With some consistency, this appears to yield improved (and even optimal) methods
- The key matrix is then (for some $\gamma > 0$)

$$\mathbf{D}_\gamma = \begin{pmatrix} 0 & -\mathbf{I} \\ \mathbf{I} & \gamma \cdot \mathbf{I} \end{pmatrix}$$

and we will (formally) construct dynamics according to

$$\dot{z} = -\mathbf{D}_\gamma \nabla H(z)$$

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", rate-optimal for convex minimisation
- For measure optimisation, obtain

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

- \approx (Kinetic, Underdamped, ...) Langevin Monte Carlo, improving upon LMC in many cases, "plausibly" optimal
- For hybrid optimisation, obtain

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

- \approx our "Momentum Particle Gradient Descent", which empirically outperforms the original PGD

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 practical tasks can be fruitfully interpreted as optimisation?
 - Should we ever look *beyond* gradient and conformal Hamiltonian flows?

The Papers

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.

1. Introduction

In this work, we study parameter estimation for (probabilistic) latent variable models $p_\theta(y, x)$ with parameters $\theta \in \mathbb{R}^{d_\theta}$, unobserved (or latent) variables $x \in \mathbb{R}^{d_x}$, and observed variables $y \in \mathbb{R}^{d_y}$ (which we treat as fixed throughout). The type II maximum likelihood approach (Good, 1983) estimates θ by maximizing the *marginal likelihood* $p_\theta(y) := \int p_\theta(y, x) dx$. However, for most models of practical interest, the integral has no known closed-form expressions and we are unable to optimize the marginal likelihood directly.

This issue is often overcome (e.g., Neal and Hinton (1998))

¹University of Warwick ²Polygeist ³University of Bristol. Correspondence to: Jen Ning Lim <Jen-Ning.Lim@warwick.ac.uk>.

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 is the classical Expectation Maximization (EM) algorithm applicable to models for which the posterior distributions $p_\theta(\cdot | y)$ are available in closed form. As shown in Neal and Hinton (1998), its iterates coincide with those of coordinate descent applied to the free energy \mathcal{E} . Recently, Kuntz et al. (2023) sought analogues of gradient descent (GD) applicable to the free energy \mathcal{E} . In particular, building on ideas popular in optimal transport (e.g., see Ambrosio et al. (2005)), they identified an appropriate notion for the free energy’s gradient $\nabla \mathcal{E}$, discretized the corresponding gradient flow $(\dot{\theta}_t, \dot{q}_t) = -\nabla \mathcal{E}(\theta_t, q_t)$, and obtained a practical algorithm they called particle gradient descent (PGD).

In optimization, GD is well-known to be suboptimal: other practical first-order methods achieve better worst-case convergence guarantees and practical performance (Nemirovskij and Yudin, 1983; Nesterov, 1983). A common feature among algorithms that achieve the optimal ‘accelerated’ convergence rates is the presence of ‘momentum’ effects in the dynamics of the algorithm. Roughly speaking,

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.