# Gradient Flows
# for Statistical Computation
## Trends and Trajectories

Sam Power, University of Bristol

# 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

- Post-Bayes as a rich source of motivating applications

# 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 ("$\mathscr{P}\left(\mathscr{X}\right)$"; $\mathscr{X} \subseteq \mathbf{R}^d$)

- Optimisation over 'Hybrid' Spaces ("$\Theta \times \mathscr{P}\left(\mathscr{X}\right)$")

# Optimisation over Parameter Spaces

- Maximum Likelihood Estimation ('MLE'): $\displaystyle\max_\theta \sum_{i\in[N]} \log p_\theta(y_i)$

  - maybe incorporate a penalty term ('penalised MLE')

  - maybe use a more general loss ('M-Estimation')

- Variational Approximation : $\displaystyle\min_\theta \mathsf{KL}\left(p_\theta, \pi\right)$

  - e.g. $\theta = (\mathrm{m}, \mathrm{C})$, $p_\theta(\mathrm{d}x) = \mathcal{N}\left(\mathrm{d}x; \mathrm{m}, \mathrm{C}\right)$; "best Gaussian fit"

# Optimisation over Measure Spaces

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

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

with $\mathscr{H}(\mu) = \int \left(\mu \log \mu - \mu\right)$ (special).

- (Nonparametric) Mean-Field Approximation

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

- (integral probability metrics, information-theoretic divergences, etc. - see Zheyang's talk!)

# Optimisation over Hybrid Spaces

- Latent Variable Models: impute $\left[x \mid \theta, y\right]$, optimise $\left[\theta \mid x; y\right]$ (EM)

- Unnormalised Models: sample $\left[x \mid \theta\right]$, optimise $\theta$ (CD / MC-MLE)

- Distributed Inference: sample local posterior, tilt for consensus (~EP)

- Opinion: more prevalent than you might expect; worth taking seriously

# 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, $\mathscr{P}(\mathscr{X})$ can carry transport ('Wasserstein') metric.

- For hybrid optimisation, $\Theta \times \mathscr{P}(\mathscr{X})$ can carry 'hybrid' metric

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

# Conceptual Optimisation Framework

OPT $\qquad \min\limits_{x \in \mathscr{X}} f(x)$

PPM $\qquad x_0 \mapsto \arg\min\limits_{x \in \mathscr{X}} \left\{ f(x) + \frac{1}{2h} \cdot \mathsf{d}\left(x, x_0\right)^2 \right\}$

FLOW $\qquad \dot{x}_t = -\nabla f\left(x_t\right)$

# Specify a Metric Structure

# Receive an Optimisation Algorithm

# Gradient Flows on Parameter Spaces

- Task: $\min\limits_{\theta \in \Theta} f(\theta)$

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

- Time-Discretised Method: "Gradient Method"

# Gradient Flows on Measure Spaces

- Task: $\min\limits_{\mu \in \mathscr{P}(\mathscr{X})} \left\{ \mathscr{F}(\mu) + \mathscr{H}(\mu) \right\}$

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

- SDE: $\mathrm{d}X_t = -\nabla_x \delta_\mu \mathscr{F}(\mu_t, X_t)\, \mathrm{d}t + \sqrt{2}\mathrm{d}W_t$

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

# Gradient Flows on Hybrid Spaces

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

- ODE-PDE: $\dot{\theta}_t = -\nabla_\theta\mathscr{F}\left(\theta_t,\mu_t\right),\,\partial_t\mu_t = -\nabla_{\mathscr{T},\mu}\left\{\mathscr{F}+\mathscr{H}\right\}\left(\theta_t,\mu_t\right)$

- ODE-SDE: $\dot{\theta}_t = -\nabla_\theta\mathscr{F}\left(\theta_t,\mu_t\right),\,\mathrm{d}X_t = -\nabla_x\delta_\mu\mathscr{F}\left(\theta_t,\mu_t,X_t\right)\,\mathrm{d}t + \sqrt{2}\mathrm{d}W_t$

- Space-Time-Discretised Method: "Particle Gradient Descent"

# A Word 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 $\Theta$,

  - very well-developed on $\mathscr{P}\left(\mathscr{X}\right)$, and

  - still under development for hybrid spaces

# 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.