

Bayesian inference: variational forms + gradient flows

Standard Bayes:

prior - $\pi(\theta)$

likelihood - $L(x|\theta)$

posterior - $\pi_n(\theta|x) = \frac{L(x|\theta)\pi(\theta)}{Z}$

$$Z = \int_{\theta \in \Theta} L(x|\theta)\pi(\theta) d\theta$$

A classic variation formulation

$\mu(\theta) \in \mathcal{P}(\mathcal{M})$ - measures on a manifold

$\pi(\theta) \in \mathcal{P}(\mathcal{M})$

$\pi_n(\theta|x) \in \mathcal{P}(\mathcal{M})$

$l(x|\theta) = -\text{Log}[L(x|\theta)]$

$$\pi_n(\theta|x) = \arg \min_{\mathcal{F}} \left[D_{KL}(q(\theta) \parallel \pi(\theta)) + \int l(x|\theta) q(\theta) d\theta \right]$$

$$\begin{aligned} D_{KL}(q \parallel \pi) &= \sum_{\theta} q(\theta) \log \frac{q(\theta)}{\pi(\theta)} \\ &= \int_{\Theta} q(\theta) \log \frac{q(\theta)}{\pi(\theta)} d\theta \end{aligned}$$

Maxent (Jaynes)

$$\pi_n(\theta|x) = \arg \min_{\mathcal{F}} \left[-\lambda \text{Ent}(q) + \int l(x|\theta) q(\theta) d\theta \right]$$

What we will talk about today:

$$\arg \min_{\mathcal{F}} \left[\int l(x|\theta) q(\theta) d\theta + D(q \parallel \pi) \right]$$

here D is a general variational functional

Questions we will ask:

- 1) What do different variational functionals imply in terms of the solution
- 2) What do gradient flows look like for these functionals, what information do they provide?
- 3) An application Riemannian Markov chain Monte Carlo

Notation and defs.

(M, g) - m -dim Riemannian manifold with metric g on \mathcal{D}

Vol_g - volume form

$\nabla_g, \text{div}_g, \text{Hess}_g, \Delta_g$ - gradient, divergence, Hessian, Laplace-Beltrami

$\mathcal{P}(\mathcal{M})$ - space of measures on \mathcal{M}

$$\pi(\theta) = e^{-\psi} \text{vol}_g$$

$$\pi_n(\theta|x) \propto e^{-\theta} \pi, \quad \theta(\cdot) = \theta(\cdot; x) \\ = -L(x|\theta)$$

$$\mathcal{P}_2(\mathcal{M}) := \left\{ \nu \in \mathcal{P}(\mathcal{M}) : \int_{\mathcal{M}} d^2(x, x_0) d\nu(x) < \infty, \right. \\ \left. \text{for some } x_0 \in \mathcal{M} \right\}$$

Functionals:

$$\text{a) KL: } D_{\text{KL}}(\nu_1 \parallel \nu_2) = \begin{cases} \int_{\mathcal{M}} \frac{d\nu_1}{d\nu_2} \log\left(\frac{d\nu_1}{d\nu_2}(\theta)\right) d\nu_2(\theta) \\ \infty \quad \text{o.w.} \end{cases}$$

$$\text{b) } \chi^2: D_{\chi^2}(\nu_1 \parallel \nu_2) = \begin{cases} \int_{\mathcal{M}} \left[\frac{d\nu_1}{d\nu_2}(\theta) - 1 \right]^2 d\nu_2(\theta) \\ \infty \quad \text{o.w.} \end{cases}$$

c) Dirichlet energy :

$$D^u(f) = \begin{cases} \int_{\mathcal{M}} \|\nabla_g f(\theta)\|^2 d\mu(\theta), \\ +\infty \quad \text{o.w.} \end{cases}$$
$$f \in L^2(\mathcal{M}, \mu) \cap H^1(\mathcal{M})$$

d) \mathcal{W}_2^2 :

$$\mathcal{W}_2^2(\nu_1, \nu_2) = \inf_{\alpha} \int_{\mathcal{M} \times \mathcal{M}} d(x, y)^2 d\alpha(x, y)$$

$$\nu_1, \nu_2 \in \mathcal{P}_2(\mathcal{M})$$

α = transportation plan,
 $\alpha \in \mathcal{P}(\mathcal{M} \times \mathcal{M})$

Geodesic space: $f \in L^2(\mathcal{M}, \mu)$ with
a constant speed geodesic connecting f_0, f_1 ,
is $t \in [0, 1] \mapsto (1-t)f_0 + tf_1$

Convexity properties:

1) A functional $E: X \rightarrow \mathbb{R} \cup \{\infty\}$ is λ -geodesically
convex if $\forall x_1, x_2 \in X \exists$ a constant speed
geodesic $t \in [0, 1] \mapsto \gamma(t) \in X$, $\gamma(0) = x_1$, $\gamma(1) = x_2$

$$E(\gamma(t)) \leq (1-t) E(x_0) + t E(x_1) - \lambda \frac{t(1-t)}{2} d_x^2(x_0, x_1),$$

$$\forall t \in [0, 1],$$

2) Let $\Psi \in C^2(M)$. The following are equivalent

- (i) Ψ is λ -geodesically convex
- (ii) $\text{Hess}_g \Psi_x(v, v) \geq \lambda \quad \forall x \in M, v \in T_x M$

If (M, d) is Euclidean (i) \Leftrightarrow (ii) \Leftrightarrow (iii)

(iii) $\Psi - \frac{\lambda}{2} |\cdot|^2$ is convex.

Strong convexity in optimization

Gradient flows:

Metric space (X, d_X)

1) $X = \mathbb{R}^d$, $d_X(u, v) = \|u - v\|$

$E: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable

$$(gf) \quad \begin{aligned} \dot{x}(t) &= -\nabla E(x(t)), \quad t \geq 0 \\ x(0) &= x_0 \end{aligned}$$

The solution of (gf) is the gradient flow.

More general in integral form

$$E(x_0) = E(x(t)) + \frac{1}{2} \int_0^t |\dot{x}(\tau)|^2 d\tau + \frac{1}{2} \int_0^t |\nabla E(x(\tau))|^2 d\tau, \quad t > 0$$

Energy dissipation equality.

$$|\dot{x}(t)| := \lim_{s \rightarrow t} \frac{d_x(x(t), x(s))}{|s-t|}$$

$$|\nabla E(x)| := \limsup_{y \rightarrow x} \frac{(E(x) - E(y))_+}{d_x(x, y)}$$

Variational forms for Bayesian updating:
from stat phys/prob.

$\mu(d\theta) = \frac{1}{Z} \exp(-\phi(\theta)) \pi(d\theta)$ is
the minimizer of

$$(1) \quad J_{KL}(v) := D_{KL}(\mu || \pi) + F_{KL}(v; \phi)$$

$$F_{KL}(v; \phi) = \int_{\mathcal{M}} \phi(\theta) d\nu(\theta)$$

$$(2) \quad J_{\chi^2}(v) := D_{\chi^2}(v || \pi) + F_{\chi^2}(v; \phi)$$

$$F_{\chi^2}(v; \phi) := \int_{\mathcal{M}} \tilde{\phi}(\theta) d\nu(\theta),$$

$$\tilde{\phi} = g(\exp(\phi(\theta))), \quad g'(t) = t^{-1}, \quad t > 0$$

$$(3) \quad \text{Given posterior } \mu$$

$$(a) \quad D^*(f) = \int_{\mathcal{M}} \|\nabla_g f(\theta)\|^2 d\mu(\theta)$$

$$\min \int_{\mathcal{M}} \|\nabla_g f(\theta)\|^2 d\mu(\theta)$$

$$f \in L^2(\mathcal{M}, \mu) \text{ s.t. } \int_{\mathcal{M}} f d\mu = 1$$

Geodesic convexity & functional inequalities

$$\pi = e^{-\Psi} \text{vol}_g, \quad \mu \propto e^{-\Phi} \pi$$

Prop. 3.1: $\Psi, \Phi \in C^2(M)$, d_{KL} is

λ -geodesically convex iff

$$\text{Ric}_g(v, v) + \text{Hess}_g \Psi(v, v) + \text{Hess}_g \Phi(v, v) \geq \lambda$$

$$\forall x \in M, \forall v \in T_x M$$

Prop. 3.2: d_{χ^2} is λ -geodesically iff the following hold

1. $\text{Ric}_g(v, v) + \text{Hess}_g \Psi(v, v) + \frac{1}{m+1} \langle \nabla_g \Psi, v \rangle^2 \geq 0$

$$\forall x \in M, \forall v \in T_x M$$

2. Φ is a λ -geodesically convex fcn.

Convexity notions:

Convexity implies convergence / flows of

$\pi \rightarrow \mu$ (posterior)

$$\pi \rightarrow \mu_t \rightarrow \mu.$$

$$D_{KL}(\mu_t \parallel \mu) \leq e^{-\lambda t} D_{KL}(\pi \parallel \mu)$$

$$W_2(\mu_t \parallel \mu) \leq \sqrt{2} e^{-\lambda t} D_{KL}(\pi \parallel \mu)$$

A measure μ has Poincaré inequality λ if $\forall f \in L^2(\mathcal{X}, \mu)$ & $\int_{\mathcal{X}} f d\mu = 0$

$$\|f\|_{\mu}^2 \leq \frac{1}{\lambda} D^{\mu}(f)$$

and D^{μ} is 2λ -geodesically convex \Leftrightarrow Poincaré constant λ .

Spectral gap:

Poincaré constant is the smallest non-trivial eigenvalue of

$$-\Delta_g^\mu f := -\frac{1}{2} \operatorname{div}_g (e^{-\rho-\gamma} \nabla_g f)$$

$$\lambda_2 := \min_{f \in L^2(m, \mu)} \frac{D^\mu(f)}{\|f - f_\mu\|_\mu^2}$$

$$f_\mu := \int m f d\mu$$

PDEs & Diffusions

1) J_{KL} + Wasserstein flow
gradient flow $t \in (0, \infty) \mapsto \mu_t$ of
 $D_{KL}(\cdot \| \mu)$ starting at $\mu_0 = \pi \in \mathcal{P}(M)$

$$\rho_t := \frac{d\mu_t}{d\mu}, \quad \theta_t := \frac{d\mu_t}{d\text{vol}_g}$$

satisfy Fokker-Planck

$$\frac{\partial \rho}{\partial t} = \Delta_g^\mu \rho$$

$$\frac{\partial \theta}{\partial t} = \Delta_g \theta + \text{div}_g(\theta(\nabla_g \varphi + \nabla_g \gamma))$$

$$dX_t = -\nabla_g(\gamma(X_t) + \varphi(X_t)) dt + \sqrt{2} dB_t'$$

2) J_{X^2} & Wasserstein flow

$$\tilde{\rho}_t := \frac{d\mu_t}{d\pi}$$

J_{X^2} - Wasserstein flow $t \in (0, \infty) \rightarrow \mathcal{M}_+$

$$\frac{\partial \tilde{\rho}}{\partial t} = \Delta_g^\pi \tilde{\rho}^2 + \operatorname{div}(\tilde{\rho} \nabla_g \theta)$$

$$\Delta_g^\pi := \Delta_g f - \langle \nabla_g f, \nabla_g \Psi \rangle$$

$$\operatorname{div}_g^\pi F := \operatorname{div}_g F - \langle F, \nabla_g \Psi \rangle$$

$$dX_t = - \left(\tilde{\rho}(t, X_t) \nabla_g \Psi(X_t) + \nabla_g \theta(X_t) \right) dt + \sqrt{2 \tilde{\rho}(t, X_t)} dB_t^g, \quad u_0 \sim \rho_0$$

$$\frac{\partial \theta}{\partial t} = - \operatorname{div}_g \left(\theta (-\tilde{\rho} \nabla_g \Psi - \nabla_g \theta) \right) + \Delta_g (\tilde{\rho} \theta)$$

$$\beta = \frac{1}{2} \exp(-\Psi) \theta$$

$$\Delta_g^\pi \beta^2 = e^\Psi \left(\Delta_g (\beta^2 e^{-\Psi}) + \operatorname{div}_g \left((\beta^2 e^{-\Psi}) \nabla_g \Psi \right) \right)$$

Stop and talk about
Bayes updating vs. PDE
idea.

Semi supervised learning -

$$\{(x_i, y_i)\}_{i=1}^n, \{(\hat{x}_i)\}_{i=n+1}^{n+m}$$

Use a graph diffusion to propagate
labels from labeled points to
unlabeled points

$$G = (X, W), \quad W_{ij} = K\left(\frac{|x_i - x_j|}{\nu}\right)$$

$$L = D - W,$$

$D_{ii} = \sum_j W_{ij}$

compact
Kernel
fcn.

Probit model

$$G(u) := \frac{1}{2} \langle L^x u, u \rangle - \sum_j \log(H(y_j u_i; y))$$

$$H(w; \gamma) = \int_{-\infty}^{\infty} \exp(-t^2 / 2\gamma^2) dt$$

$$\hat{u}_{in} = \arg \min_{\substack{u \in \mathbb{R}^n \\ \sum_i u_i = 0}} G(u)$$

$$\hat{u} = \text{sign}(\hat{u}_{in})$$

Prior - $U := \{u \in \mathbb{R}^n : \sum u_i = 0\}$

$$\frac{d\pi}{du}(u) \propto \exp(-\frac{1}{2} \langle L^\alpha u, u \rangle)$$

$$:= \exp(-\psi u)$$

Likelihood:

$$y_i = \text{sign}(u_i + \eta_i)$$

$$\eta_j \stackrel{iid}{\sim} N(0, \gamma^2)$$

$$\mathcal{L}(u; y) := - \sum_j \log(H(y; u_j; \gamma))$$

$u \in U$

Posterior :

$$\frac{du^y}{d\pi}(v) = \exp(-\phi(v; y))$$

$$\frac{du^y}{dv}(v) \propto \exp(-\psi(v) - \phi(v; y))$$

$$J_{KL}(v) := D_{KL}(v \| \pi) + \int_{\mathbb{R}^n} \phi(u; y) d\nu(u)$$

$$v \in \mathcal{P}(\mathbb{R}^n)$$

1) π is Gaussian with $\Sigma = L^{-\alpha}$
 $\lambda(L)$ is the smallest n.t. eigenvalue

2) $\phi(\cdot; y)$ is convex


$$F_{KL}(v) = \int_{\mathbb{R}^n} \phi(u; y) d\nu(u) \text{ is}$$

Θ -geodesically convex.

diffusion equation

$$dX_t = (-L^\alpha X_t - \nabla \phi(X_t, y)) dt + \sqrt{2} dB_t$$

as α gets large L^α has issues


$$dX_t = -(X_t + L^{-\alpha} \nabla \phi(X_t, y)) dt + \sqrt{2L^{-\alpha}} dB_t$$

$$\nabla_g = G^{-1} \nabla, \quad B^g \stackrel{?}{=} \sqrt{G^{-1}} B$$

$$dX_t = -\nabla_g(\phi(X_t, y)) + \frac{1}{2} \langle L^\alpha X_t, X_t \rangle dt + \sqrt{2} dB_t^g$$