

Sampling with constraints

Xin T Tong

BIRS workshop

*Joint work with Qiang Liu, Xingchao Liu, Ruqi Zhang (UT
Austin)*

Friday 9th September, 2022

- Constrained sampling
- Review: KL gradient flow without constraint
- Moment Constraints
- Level set constraints
- *Sampling with Trustworthy Constraints: A Variational Gradient Framework* NeurIPS 2021.
- *Sampling in Constrained Domains with Orthogonal-Space Variational Gradient Descent* Under review

Standard Bayesian problem:

$$\text{Sample } (\theta) \propto p_0(\theta) \exp(-l(\theta))$$

Moment constrained Bayesian problem:

$$\text{Sample } q \quad \text{s.t. } \mathbb{E}_q[g(\theta)]$$

Equality constraint

$$\text{Sample } q \quad \text{s.t. } g(x) = 0 \text{ for } q\text{-a.s. } x$$

Type of constraint functions

- Agnostic learning: $g(\cdot) = l(\cdot)$
- Fairness: $g(\cdot) = \text{cov}(\hat{y}(x; \cdot); Z)$
- Monotonicity: $g(\cdot) = [\partial_x \hat{y}(x; \cdot)]_+$
- Safety: $\text{dist}(\hat{y}(x; \cdot); S)$

Type of questions:

- What would the solution be?
- How to obtain the distribution?
- Pareto front of l vs g

Existing fairness works: Chakraborty, Ji, Dimitrakakis

Review: unconstrained case

Markov Chain Monte Carlo (MCMC)

- Simulate a Markov Chain with π being the invariant
- Fairly well understood
- Require well specified
- Iterates tend to be dependents
- MC convergence: $O(1/\epsilon^2)$

Variational method

- Try to push a density towards π .
- Interacting particle system
- Promising on some problems.
- Understanding is much less.
- Potentially can be faster(?)

Basic formulation:

- Try to minimize $\text{KL}(q_t; \cdot)$
- Suppose we have samples from a density q_t .
- We can estimate $E_{q_t}[f]$ for any f .
- Try to push each point x in q_t with $\psi_t(x)$
- Continuity equation: $\frac{d}{dt}q_t = r(\cdot, q_t)$
- What is the optimal ψ for reducing KL?
- Solve sampling by optimization methods.

Rate of decay

$$\frac{d}{dt} \text{KL}(q_t; \cdot) = \mathbb{E}_{q_t}[h r \log \cdot - r \log q_t; \cdot]$$

Try to maximize, write $r \log \cdot = s$

$$\max_{2H} \mathbb{E}_{q_t}[h s - s q_t; \cdot] \quad \frac{1}{2} k \quad k_H^2$$

If we use $H = L_{q_t}^2$

- We obtain $\dot{q}_t = S(q_t)$.
- But how to get $S(q_t)$?
- Stein operator $A = (S + r)$

$$\frac{d}{dt}q_t = r(q_t) = r(S(q_t)) + \Delta q_t = r(A(q_t))$$

- Fokker–Plank equation (FPE) of Langevin dynamics (LD)[Jordan, Kinderlehrer, and Otto 1998]
- Algorithmic implementation (ULA):

$$q_{t+1} = q_t + S(q_t) + \sqrt{\frac{\rho}{2}} \epsilon_{t+1}$$

- Can be seen as an MCMC as well.

Use

$$\frac{d}{dt} \text{KL}(q_t; \cdot) = \mathbb{E}_{q_t} k_S - s_{q_t} k^2$$

- $\int_0^T \mathbb{E}_{q_t} k_S - s_{q_t} k^2 \text{KL}(q_0; \cdot)$
- Fisher divergence $\min_t \int_0^T \mathbb{E}_{q_t} k_S - s_{q_t} k^2 = O(1/T)$
- If the log-Sobolev inequality (LSI) holds,
 $k_S - s_{q_t} k^2 \geq c \text{KL}(q_t; \cdot)$, $\text{KL}(q_t; \cdot) = O(\exp(-ct))$.
- Can be inherited by ULA (Vempala and Wibisono 2019)

Use $H = \text{RKHS}$ with kernel k ,

- $\mathbb{E}_q \left(\int_{\mathcal{R}} (s(y) - r \log q_t(y)) k(x; y) q_t(y) dy \right)$
- A kernel embedding of \mathcal{A} into H
- Limit point meets Stein equation $\mathbb{E}_q \mathcal{A} f = 0$ for $f \in H$.
- $\mathbb{E}_q \left(\int_{\mathcal{R}} s(y) k(x; y) q_t(y) dy + \int_{\mathcal{R}} r_y k(x; y) q_t(y) dy \right)$
- Replace q_t with samples from q_t .

$$x_{i;t+1} = x_{i;t} + \frac{\gamma}{n} \sum_{j=1}^n k(x_{i;t}, x_{j;t}) r_{j;t} \log q_t(x_{j;t}) + \gamma \sum_{j=1}^n k(x_{i;t}, x_{j;t})$$

- Deterministic after initialization.
- Stein Variational Gradient Descent (SVGD) [Liu and Wang 2016]

Use

$$\frac{d}{dt} \text{KL}(q_t; \mathcal{Z}) = \int k_S - s_{q_t} k_k^2$$

$$:= \int q_t(x) q_t(y) k(x; y) (s - s_{q_t})(x)^T (s - s_{q_t})(y)$$

- Kernel Stein divergence $\min_t \int q_t k_S - s_{q_t} k_k^2 = O(1/T)$
- Is there LIS $k_S - s_{q_t} k_k^2 \leq c \text{KL}(q_t; \mathcal{Z})$?
- Actually not correct in general (Gorham and Mackey 2017)

Moment constrained

Solve

$$\min_q \text{KL}(q; \cdot); \quad s:t: \quad \mathbb{E}_q[g] = 0:$$

- Ignore the possibility $\mathbb{E}[g] = 0$; where \cdot is the solution.
- Solution: $q = \frac{1}{Z} \exp(-\beta g)$ and $\mathbb{E}[g] = 0$
- Chicken: Checking $\mathbb{E}[g] = 0$ requires samples from q
- Egg: sampling from q requires β
- Double loop: MCMC or variational, feasible but expensive

Primal dual gradient flow (PDGF)

Reformulate as

$$\min_q \max_0 fL(q; \cdot) = \text{KL}(q \parallel j) + \mathbb{E}_q[g]g:$$

Gradient ascent on \cdot :

$$\frac{d}{dt} \cdot = [\mathbb{E}_{q_t}[g]] \cdot_{t+}$$

When $H = L^2$, gradient descent on q via \cdot :

$$\dot{q}_t = r (\log \cdot_t - \log q_t) = S \cdot_t r g \cdot S q$$

When $H = \text{RKHS}$, gradient descent on q via \cdot :

$$\dot{q}_t(x) = \int (S(y) \cdot_t r g(y) + r_y) k(x; y) q_t(y) dy$$

Assume

Theorem

Suppose

$$k_{S_{q_t}} \leq k_{q_t}^2 \leq c_1 (\mathbb{E}_{q_t}[g] - \mathbb{E}[g])^2$$

LD-PDGF finds solutions $k_{S_{q_t}} \leq k_{q_t}^2 = O(1/T)$. If g is convex, satisfies log Sobolev, then linear convergence for $KL(q_t; \cdot)$

For SVGD, $k_{q_t}^2$ is replaced by kernel Stein discrepancy.

Theorem

Suppose

$$k_{S_{q_t}} \leq k_k^2 \leq c_1 (\mathbb{E}_{q_t}[g] - \mathbb{E}[g])^2$$

LD-PDGF finds solutions $k_{S_{q_t}} \leq k_k^2 = O(1/T)$.

Constraint Controlled gradient flow (CCGF)

Try to solve

$$\max_{s, q_t} \mathbb{E}_{q_t}[h(s, q_t; \theta)] - \frac{1}{2} k \|k_H\|^2; \quad s: \frac{d}{dt} \mathbb{E}_{q_t} g = \mathbb{E}_{q_t} \nabla_r g - \mathbb{E}_{q_t}[g]$$

Solve quadratic opt.

$$\min_0 \max_{s, q_t} \mathbb{E}_{q_t}[h(s, q_t; \theta)] - \frac{1}{2} k \|k_H\|^2 + (\mathbb{E}_{q_t} \nabla_r g - \mathbb{E}_{q_t}[g])$$

We have $t = s$ $\nabla_r g = s_q$ (LD case)

$$t = \max \frac{\mathbb{E}_{q_t}[g] + h(s, q_t; \theta) \nabla_r g}{k \|g\|_{q_t}^2}; 0$$

Or $t(x) = \int_{\mathcal{R}} (s, q_t; \theta) \nabla_r g = s_q(y) k(x; y) q_t(y) dy$ (SVGD case).

$$t = \max \frac{\mathbb{E}_{q_t}[g] + h(s, q_t; \theta) \nabla_r g}{k \|g\|_k^2}; 0$$

Theorem

Suppose t is bounded by a constant, LD-CCGF finds solutions $k_{S_{q_t}} \leq k_{q_t}^2 = O(1/T)$. If g is convex, μ satisfies log Sobolev, then linear convergence for $KL(q_t; \mu)$.

For SVGD, $k_{S_{q_t}} \leq k_{q_t}^2$ is replaced by kernel Stein discrepancy.

Theorem

Suppose t is bounded by a constant, SVGD-CCGF finds solutions $k_{S_{q_t}} \leq k_k^2 = O(1/T)$.

Algorithm 3 Primal-Dual Method

Initialize the particles $\{\theta_{i,0}\}_{i=1}^n$ and λ_0 .

for iteration t **do**

If Langevin, update $\theta_{i,t+1} = \theta_{i,t} + h(\nabla \log p_0^*(\theta_{i,t}) - \lambda_t \nabla g(\theta_{i,t})) + \sqrt{2h} \xi_{i,t}$.

If SVGD, update

$$\theta_{i,t+1} = \theta_{i,t} + \frac{h}{n} \sum_{j=1}^n [(\nabla \log p_0^*(\theta_{j,t}) - \lambda_t \nabla g(\theta_{j,t})) k_t(\theta_{j,t}, \theta_{i,t})] + \nabla_{\theta_{j,t}} k_t(\theta_{j,t}, \theta_{i,t}).$$

 Update λ_t by $\lambda_{t+1} = \max(\lambda_t + \frac{h}{n} \sum_{i=1}^n [g(\theta_{i,t+1})], 0)$.

end for

Algorithm 4 Constraint Controlled Method

Initialize the particles $\{\theta_{i,0}\}_{i=1}^n$.

for iteration t **do**

If Langevin, update

$$\lambda_t = \max \left(\frac{\sum_{j=1}^n \alpha g(\theta_{j,t}) + [(\nabla \log p_0^*(\theta_{j,t}))^\top \nabla g(\theta_{j,t}) + \nabla^\top \nabla g(\theta_{j,t})]}{\sum_{j=1}^n [\|\nabla g(\theta_{j,t})\|^2]}, 0 \right),$$

 update $\theta_{i,t+1} = \theta_{i,t} + h(\nabla \log p_0^*(\theta_{i,t}) - \lambda_t \nabla g(\theta_{i,t})) + \sqrt{2h} \xi_{i,t}$.

If SVGD, update

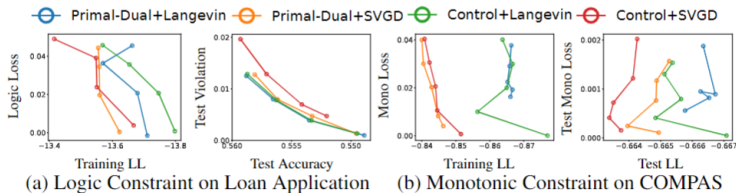
$$\lambda_t = \max \left(\frac{\sum_{i,j=1}^n \alpha g(\theta_{i,t}) + [\nabla g(\theta_{j,t})^\top (\nabla \log p_0^*(\theta_{i,t}) + \nabla_{\theta_{i,t}}) k_t(\theta_{i,t}, \theta_{j,t})]}{\sum_{i,j=1}^n [\nabla g(\theta_{i,t})^\top \nabla g(\theta_{j,t}) k_t(\theta_{i,t}, \theta_{j,t})]}, 0 \right),$$

 update

$$\theta_{i,t+1} = \theta_{i,t} + \frac{h}{n} \sum_{j=1}^n [(\nabla \log p^*(\theta_{j,t}) - \lambda_t \nabla g(\theta_{j,t})) k_t(\theta_{j,t}, \theta_{i,t}) + \nabla_{\theta_{j,t}} k_t(\theta_{j,t}, \theta_{i,t})].$$

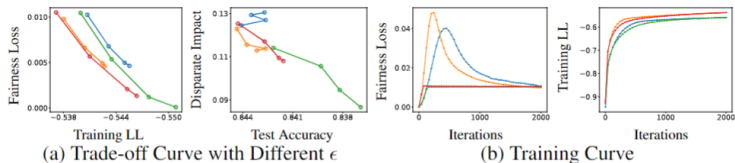
end for

Logic and Monotonicity constrained logistic regression.



Fairness constrained Neural Network

Primal-Dual+Langevin Primal-Dual+SVGD Control+Langevin Control+SVGD



Equality constrained

Formulation of problem

- Minimize $\text{KL}(q; \cdot)$ so that q is supported on $G_0 = \{x : g(x) = 0\}$
- Ill-posed: q is singular w.r.t. \cdot .
- Try to sample the conditional measure $q_0(\cdot) = [jg = 0]$.
- Hausdorff density $\int_{G_0} g(x) dx$ on G_0 .

Sampling on manifolds

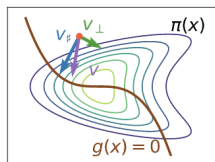
- Several existing MCMC (Girolami, Brubaker, Lelievre...)
- Assume MCMC start and stay on G_0
- Often require explicit knowledge of G_0 (parameterization, geodesic, projection)
- Not so friendly for large scale ML models.

Deriving algorithm

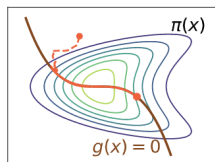
Try to solve

$$\max E_{q_t}[hS \quad S_{q_t; v}] \quad \frac{1}{2} k v k_H^2;$$

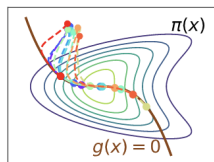
$$s.t.: \frac{d}{dt} g(x_t) = v^T(x) r g(x) = \quad (g(x))$$



(a) O-Gradient



(b) O-Langevin



(c) O-SVGD

Deriving algorithm

Along $r g$

- Use $\psi(z) = \text{sign}(z)|z|^{1+\epsilon}$

- The component along $r g$: $v_j = \frac{(g(x))r g(x)}{kr g(x)k^2}$

Along the orthogonal direction:

- Projection: $D = I - \frac{r g r g^T}{kr g k^2}$

- $v_\gamma = Du, \max_u E_{q_t}[(D(s \quad s_{q_t}))^T u] \quad \frac{1}{2}kDu k_H^2:$

- LD: $v_\gamma = D(s \quad s_{q_t})$

- SVGD:

$$\begin{aligned} v_\gamma(x) &= \int_{\mathcal{Z}} D(x)k(x;y)D(y)(s \quad s_{q_t})(y)q_t(y)dy \\ &= \int_{\mathcal{Z}} k_\gamma(x;y)(s \quad s_{q_t})(y)q_t(y)dy \end{aligned}$$

Implementation

- LD: $v_T = D(s \quad s_{q_t})$ cannot be implemented directly by $dx_t = (v_J(x_t) + D(x_t)s(x_t))dt + \sqrt{2D(x_t)}dW_t$:
- Consider adding a correction drift r

Theorem

When $r(x) = r - D(x)$,

$$dx_t = (v_J(x_t) + D(x_t)s(x_t))dt + \sqrt{2D(x_t)}dW_t \quad (1)$$

its FPE matches the orthogonal density flow. Moreover, i) the value $g(x_t)$ has deterministic decay $\frac{d}{dt}g(x_t) = -r g(x_t)$; ii) for any f with $r f - r g = 0$, the generator of x_t matches the Langevin ones $Lf(x) = r f''(x)s(x) + \Delta f(x)$.

Define orthogonal space (OS) Fisher divergence

$$F_{\mathcal{F}}(q; \cdot) = kD(s \quad s_q)k_q^2 \text{ or } kD(s \quad s_q)k_k^2$$

Theorem

Suppose $g(x)$ is bounded for the initial distribution, and it's "regular", $KL(q_0; \cdot) < 1$, then

$M_T = \max_{f \in \mathcal{F}} \int f(x)g(x); X \quad q_T g = O(T^{-1})$; also convergence in OS-Fisher $\min_{t \leq T} F_{\mathcal{F}}(q_t; \cdot) = O(\log T/T)$.

But is OS-Fisher useful?

Simpler formulation

The distribution $\Pi_Z = (\int g(x) = z)$ is too abstract.

Theorem

Suppose g has Lipschitz density. Then the weak limit of $\int g(x) \delta_{z(x)} / \int (x) \exp(-\frac{1}{2}(g(x) - z)^2)$ as $\sigma \rightarrow 0$ concentrates on $G_Z = \{x : g(x) = z\}$ and is a version of $\delta_{z(x)}$. Moreover,

$$E_{\Pi_Z}[A] = 0; \quad \forall \delta \in \mathcal{T}_g$$

- This gives a Stein equation $E_q[A] = 0$
- The tangent bundle of G_Z is a subset of \mathcal{T}_g
- $E_q[A] = \int A \overline{F_{\mathcal{T}}(q; \delta)}$ when $k = 1$.
- $E_q[A]$ or $F_{\mathcal{T}}(q; \delta)$ do not require q being on G_Z
- This only check the OS directions.
- Checking how far is q away from G_Z is easy.

Theorem

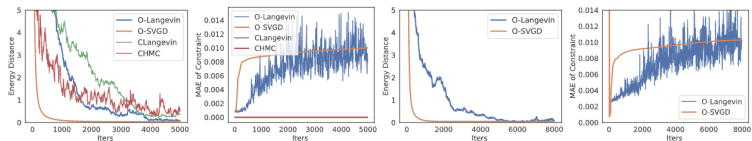
Suppose that Π_Z satisfies χ^2 -Poincare Inequality for χ^2 , and q is supported on $f_X: \chi^2(x) \leq g$. Then for any function f such that $\int f^2 \leq 1$, the following holds

$$\left| \int f q - \int f \Pi_0 \right| \leq \sqrt{\frac{1}{\chi^2(q)}} + \max_{\chi^2} \left| \int f \Pi_Z - \int f \Pi_0 \right|$$

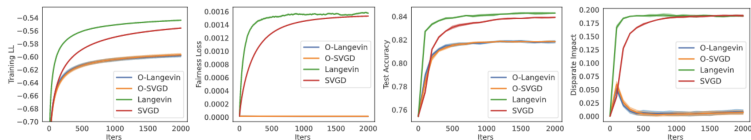
- Decomposition of mean difference/TV
- Only in L^2 case
- Poincare inequality with Euclidean-inherent distance
- Can be used for q supported on \mathbb{R}^d .

Numerical examples

Toy example (Intialized on/off manifold)



Income prediction



Agonistic Bayesian Image classification

	Test Error (\downarrow)	ECE (\downarrow)	AUROC (\uparrow)
SGLD	15.00	2.21	89.41
Tempered SGLD	4.73	0.83	97.63
O-Langevin	4.46	0.87	98.68
SVGD	6.11	0.93	93.55
O-SVGD	4.92	0.77	94.69