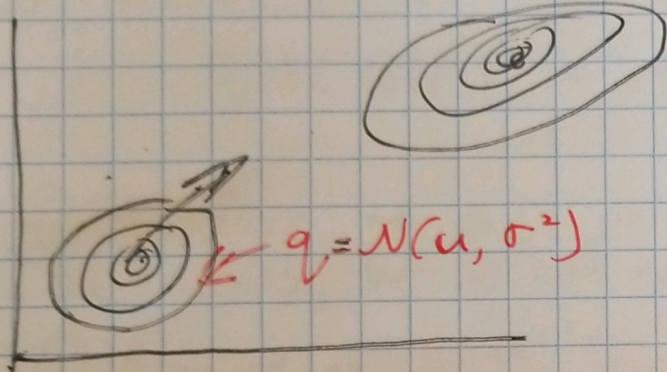


Stein Variational Gradient Descent

Qiang Liu & Dilan Wang, 2016

Intuition picture:

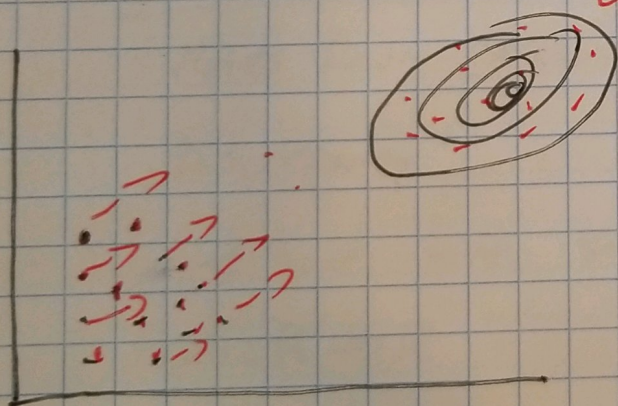
Traditional VI



Move $q_0 \rightarrow P$ by optimizing $\theta = \{u, \sigma\}$

- Also used:
- "A short introduction to Kernelized Stein Discrepancy", - Qiang Liu
- "Functional Gradient Descent", - Drew Bagnell

Stein VI:



Move $\{x_i\}_{i=1}^N \rightarrow P, x_i \sim q_0$ by:

$$x_i^{t+1} \leftarrow x_i^t + \epsilon_t \hat{\Phi}^\theta(x_i^t),$$

where $\hat{\Phi}^\theta(x_i^t)$ comes from the "kernelized Stein Discrepancy" (Defined momentarily).

"Gradient Descent on Particles"

Road Map:

- ① Introduce "Kernelized Stein Discrepancy" (KSD)
 - ② Introduce V.I. using "smooth transforms" from Family \mathcal{Q}
 - ③ Glue 1 & 2 together by showing that $\nabla_{\mathcal{Q}} \mathbb{K}[Q|P] \Rightarrow \text{KSD}$
 - ④ Realize we can't solve the Stein equations exactly, but we can operate on samples
 - ⑤ Arrive at "Stein Variational Gradient Descent"
-

① Kernelized Stein Discrepancy

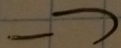
- Notes from "A short introduction to Kernelized Stein Discrepancy"

- Say we have some smooth density $p(x)$ and a function

$$f(x): \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\lim_{\|x\| \rightarrow \infty} p(x) f(x) = 0, \quad \forall x \in \mathbb{R}^d$$

- then we can derive "Stein's Identity" as follows:



$$\begin{aligned}
 p(x)f(x) &= 0 & x \in \mathbb{R}^d \\
 & & f: \mathbb{R}^d \rightarrow \mathbb{R} \\
 &= \int \nabla_x [p(x)f(x)] dx = 0 \\
 &= \int [f(x)\nabla_x p(x) + p(x)\nabla_x f(x)] dx = 0 \\
 &= \int [p(x)\nabla_x \log p(x) f(x) + p(x)\nabla_x f(x)] dx = 0 \\
 &= \mathbb{E}_{x \sim p(x)} [f(x)\nabla_x \log p(x) + \nabla_x f(x)] = 0
 \end{aligned}$$

$$\boxed{= \mathbb{E}_{x \sim p(x)} [A_p[f(x)]] = 0}$$

$$A_p[f] = f(x)\nabla_x \log p(x) + \nabla_x f(x)$$

$A_p[\cdot]$ is a functional operator, called the "Stein operator"

- Gives an infinite number of identities indexed by choice of function f

"Reverse" of Stein's identity also holds; ie

$$\mathbb{E}_{x \sim p(x)} [A_p[f(x)]] = 0 \quad \text{Stein Identity}$$

$$\mathbb{E}_{x \sim q(x)} [A_p[f(x)]] \neq 0 \quad \text{"Reverse" Stein Identity}$$

(See "A short introduction to KSA for proof")

→ The reverse Stein identity therefore gives a measure of how "different" p & q are.

We can then ask:

Which f maximizes the violation of Stein's identity?

This leads to the notion of the "Stein Discrepancy":

$$\sqrt{S(q, p)} = \max_{f \in \mathcal{F}} E_{x \sim q} [A_p(f(x))] \rightarrow S(q, p) = \max_{f \in \mathcal{F}} [E_{x \sim q} [A_p(f(x))]]^2$$

Choice of set \mathcal{F} is critical;

→ must be computationally tractable & have discriminative power

Previous work (Liu et al) has shown that

if you choose \mathcal{F} to be the unit ball in the "reproducing kernel Hilbert space"

\mathcal{H}_k , the optimization has a closed form

solution:

$$S(q, p) = \max_{f \in \mathcal{H}_k} [E_{x \sim q} [A_p(f(x))]] \text{ s.t. } \|f\|_{\mathcal{H}_k} \leq 1$$

$$f_{w,p}^*(x) = \frac{\hat{f}(x)}{\|\hat{f}\|_{\mathcal{H}_k}}$$

$$\hat{f}(\cdot) = E_{x \sim q} [A_p(k(x, \cdot))]$$

$$S(q, p) = \|\hat{f}\|_{\mathcal{H}_k}^2$$

"Reproducing" property

A kernel:
 x be a non-empty set.
 $k(\cdot, \cdot) : x \times x \rightarrow \mathbb{R}$
 $\phi : x \rightarrow \mathcal{H}$
 $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$
 Reproducing property (1 Def'n)
 $k(x, x') = \langle k(\cdot, x), k(\cdot, x') \rangle = \langle \phi(x), \phi(x') \rangle$

"A Hilbert space allows us to apply concepts from finite-dimensional linear algebra to infinite-dimensional function spaces"
 - Internet

PART 2 Variational Inference Using Smooth Transforms:

Standard VI approximates $p(x)$ via a simpler distribution from some set $Q = \{q(x)\}$ by minimizing KL:

$$q^* = \underset{q \in Q}{\operatorname{argmin}} \{KL(q||p)\} = E_q \log q(x) - E_q \log p(x) + C$$

Consider set Q formed by smooth transforms of reference dist:

Indexed by T \rightarrow

$$Q = \{z = T(x) \mid x \sim q_0(x), T: x \rightarrow z\}$$

Then; by change of variables:

$$q_{[T]}(z) = q(T^{-1}(z)) \cdot \underbrace{|\det(\nabla_z T^{-1}(z))|}_{\text{Jacobian}}$$

How to choose T ?

\rightarrow Give T a parametric form and opt. parameters
 \rightarrow "Normalizing Flow"

\rightarrow Or, we can iteratively select T 's that minimize the KL.

\rightarrow Can be thought of as "functional gradient descent"

PART 3 / Functional Gradient Descent

Regular gradient descent:

$$L(w) = \sum_{i=1}^N (y_i - w^T x_i)^2$$

Goal: $\min_w L(w)$

- Procedure:
- compute $\nabla_w L$
 - select step size α_t
 - $w_{t+1} = w_t - \alpha_t \nabla_w L(w)$

Functional Gradient Descent

$$L[f] = \sum_{i=1}^N (y_i - f(x_i))^2$$

Goal: $\min_f L[f]$

- Procedure
- compute $\nabla_f L[f]$ "Functional Gradient"
 - select step size α_t
 - $f_{t+1} = f_t - \alpha_t \nabla_f L$

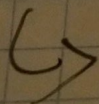
Our problem:

$$T(x) = x + f(x)$$

$$L[f] = KL[\tilde{q}_{T(x)}(x) \parallel p(x)]$$

$T(x) = x + f(x)$

Goal: $\min_f L[f]$



Functional Gradients

C.C.F. "Functional Gradient Descent"
Note by Drew Bagnell

A gradient can be thought of as:

- A vector of partial derivatives
- The direction of steepest ascent

★ - The linear approximation of the function:

$$f(x_0 + \epsilon) = f(x_0) + \epsilon \underbrace{\nabla f(x_0)}_{\text{Gradient}} + O(\epsilon^2)$$

A "Functional" gradient is similar to the third definition, and is defined implicitly as the linear term of change in a function due to a small perturbation ϵ to its input.

Formal definition:

For any functional $F[f]$ of $f \in \mathcal{H}^d$ its (functional) gradient $\nabla_f F[f]$ is a function in \mathcal{H}^d s.t.:

$$F[f + \epsilon g] = F[f] + \epsilon \langle \nabla_f F[f], g \rangle_{\mathcal{H}^d} + O(\epsilon^2)$$

For any $g \in \mathcal{H}^d$ and $\epsilon \in \mathbb{R}$

PART 4

Procedure: (Attempt 1)

- Compute $\nabla_x [L] = \nabla_x \text{KL}(q_T \| p) = -\Phi_{q,p}^*(x)$

- select step size ϵ_t

$F=0$ \rightarrow Lemma 3. in paper

\rightarrow step

$$F_{t+1} \rightarrow F_t + \epsilon_t \Phi_{q_t, p}^*(\cdot)$$

$$T_{t+1}(x) = x + F_{t+1}$$

$$T_{t+1}(x) = x + \epsilon_t \Phi_{q_t, p}^*(\cdot)$$

$$\Phi_{q_t, p}^*(\cdot) = E_{x \sim q_t} [k(\cdot, x) \nabla_x \log p(x)^T + \nabla_x k(\cdot, x)]$$

The evaluation at 0 is critical. Without it,

$$T_{t+1}(x) = x + \epsilon_t \Phi_{q_t}^*(x) + \epsilon_{t-1} \Phi_{q_{t-1}}^*(x) + \dots + \epsilon_0 \Phi_{q_0}^*(x)$$

which

is computationally infeasible.

But this means the step is only valid from $F=0$, i.e. $\tau(x) = x + F$,

IE, this isn't allowed:

$$x_0 \sim q_0$$

$$z_1 = T_0(x_0)$$

$$z_2 = T_1(z_1)$$

\vdots

$$z_t = T_t(z_t)$$

$$x_0 \sim q_0$$

$$\Rightarrow (T_t \circ T_{t-1} \circ T_{t-2} \dots T_0)(x_0)$$

PART 5

We can get around this by moving the particles at each step so they are constantly sampled from $q_{[T^t]}$

We then empirically approximate $\hat{\Phi}^p(x)$ via:

$$\hat{\Phi}^p(\cdot) \approx \frac{1}{n} \sum_{j=1}^N [K(\cdot, x_j^t) \nabla_{x_j^t} \log p(x_j^t) + \nabla_{x_j^t} K(\cdot, x_j^t)]$$
$$x_j^t \sim q_{[T^t]}(x)$$

This leads to the final algorithm, "Stein Variational (Functional) Gradient Descent":

Input:

- Target density function $p(x)$
- Initial particles $\{x_i^0\}_{i=1}^N$

Output: - Particles $\{x_i\}_{i=1}^N$ that approximate $p(x)$

Procedure:

- for iteration t :

- Make smooth transform function $T^t(\cdot)$:

$$- T^t(x) = x + \hat{\Phi}^p(x)$$

$$- \hat{\Phi}^p(\cdot) = \frac{1}{N} \sum_{j=1}^N [K(\cdot, x_j^t) \nabla_{x_j^t} \log p(x_j^t) + \nabla_{x_j^t} K(\cdot, x_j^t)]$$

- Update particles:

$$x_i^{t+1} = T^t(x_i^t)$$