

[Paper review 26]

Variational Inference with Normalizing Flows

(Danilo Jimenez Rezende, Shakir Mohamed, 2016)

[Contents]

1. Abstract
2. Introduction
3. Amortized Variational Inference
 1. Stochastic Backpropagation
 2. Inference Networks
 3. Deep Latent Gaussian Models (DLGM)
4. Normalizing Flows (NF)
 1. Finite Flow
5. Inference with NFs
 1. Invertible Linear-time Transformations
 1. Planar Flows
 2. Radial Flows
 2. Flow-Based Free Energy Bound
 3. Algorithm Summary

1. Abstract

choice of approximate posterior distribution q in VI :

- had been simple families
(ex. mean-field or other simple structured approximations)
- these restrictions → not good performance

Introduce a new approach, "Normalizing Flow"

- flexible, complex, and scalable

2. Introduction

limitations of variational methods : choice of posterior approximation are often limited

→ thus, richer approximation is needed

Methods for richer approximation

- ex1) structured mean field approximations that incorporate basic form of dependency within the approximate posterior
- ex2) mixture model (limit : potential scalability... have to compute each for the mixture component)

We will

- 1) review the current est practice (based on "amortized VI ")
- 2) make following contributions
 - a) propose a method using normalizing flow (NF)
 - b) show that NF admit infinitesimal flows

3. Amortized Variational Inference

current best practice in VI uses...

- 1) mini-batches
- 2) stochastic gradient descent (SGD)

→ to deal with very large dataset

for successful variational approach, we need to...

- 1) efficient computation of the derivatives of the expected log-likelihood,
 $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(\mathbf{x} | \mathbf{z})]$
 - solution 1) MC estimation
 - solution 2) inference networks
 - (solution 1+2 = "amortized VI")
- 2) choosing the richest, computationally-feasible approximate posterior distribution, $q(\cdot)$
 - solution) Normalizing Flow!

3.1 Stochastic Backpropagation

compute $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(\mathbf{x} | \mathbf{z})]$ (expected log likelihood) ... with MC estimation!

also called "doubly-stochastic estimation".. why double?

- 1) stochasticity from the mini-batch
- 2) stochasticity from the MC approximation of the expectation

"continuous latent variables" + "MC approximation"

= Stochastic Gradient Variational Bayes (SGVB)

SGVB involves 2 steps

- 1) Reparameterization
 $z \sim \mathcal{N}(z | \mu, \sigma^2) \Leftrightarrow z = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$
- 2) Backprop with MC
 $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [f_{\theta}(z)] \Leftrightarrow \mathbb{E}_{\mathcal{N}(\epsilon|0,1)} [\nabla_{\phi} f_{\theta}(\mu + \sigma\epsilon)]$

3.2 Inference Networks

Inference Network

- def) model that learns an INVERSE MAP from observation(x) to latent variables(z)
- $q_{\phi}(\cdot)$ is represented using Inference Networks!
- why Inference Network?
 - we avoid the need to compute per data point variational parameters, but can instead compute a set of global variational parameters ϕ valid for inference at both training and test time.
- simplest Inference Network : "DIAGONAL GAUSSIAN densities"

$$q_{\phi}(\mathbf{z} | \mathbf{x}) = \mathcal{N}\left(\mathbf{z} | \boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}\left(\boldsymbol{\sigma}_{\phi}^2(\mathbf{x})\right)\right)$$

3.3 Deep Latent Gaussian Models (DLGM)

hierarchy of L layers of Gaussian latent variables z_l for layer l

$$p(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_L) = p(\mathbf{x} | f_0(\mathbf{z}_1)) \prod_{l=1}^L p(\mathbf{z}_l | f_l(\mathbf{z}_{l+1}))$$

- prior over latent variables : $p(\mathbf{z}_l) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- observation likelihood : $p_{\theta}(\mathbf{x} | \mathbf{z})$ by NN

DLGMs

- use continuous latent variable
- model class perfectly suited to fast amortized VI (using ELBO & stochastic back-prop)
- end-to-end system of DLGM \approx encoder-decoder architecture

4. Normalizing Flows (NF)

optimal variational distribution

- $\mathbb{D}_{\text{KL}}[q||p] = 0$
(= $q_{\phi}(\mathbf{z} | \mathbf{x}) = p_{\theta}(\mathbf{z} | \mathbf{x})$)
- $q_{\phi}(\mathbf{z} | \mathbf{x})$ should be highly flexible

NF describes the transformation of probability density through "A SEQUENCE OF INVERTIBLE MAPPINGS"

4.1 Finite Flows

setting

- $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$, where $f^{-1} = g$
- $g \circ f(\mathbf{z}) = \mathbf{z}$.
- $\mathbf{z}' = f(\mathbf{z})$

variable transformation

$$-q(\mathbf{z}') = q(\mathbf{z}) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{z}'} \right| = q(\mathbf{z}) \left| \det \frac{\partial f}{\partial \mathbf{z}} \right|^{-1}$$

successive application

$$\mathbf{z}_K = f_K \circ \dots \circ f_2 \circ f_1(\mathbf{z}_0)$$

$$\ln q_K(\mathbf{z}_K) = \ln q_0(\mathbf{z}_0) - \sum_{k=1}^K \ln \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|$$

expectation

$$\mathbb{E}_{q_K}[h(\mathbf{z})] = \mathbb{E}_{q_0}[h(f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0))]$$

- does not depend on q_k

summary

- use simple factorized distribution (ex. independent Gaussian)
- apply NF of different lengths to get increasingly complex distribution

5. Inference with NFs

we must ...

- 1) specify a class of invertible transformations
- 2) efficient mechanism for computing the determinant of Jacobian

Therefore we require NF that allow for low-cost computation of the determinant, or where Jacobian is not needed!

5.1 Invertible Linear-time Transformations

linear time transformation

= we can compute the log det-Jacobian term in $O(D)$ time

5.1.1 Planar Flows

form : $f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^\top \mathbf{z} + b)$

- $\lambda = \{\mathbf{w} \in \mathbb{R}^D, \mathbf{u} \in \mathbb{R}^D, b \in \mathbb{R}\}$
- $h(\cdot)$: smooth element-wise non-line with derivative $h'(\cdot)$
- $\left| \det \frac{\partial f}{\partial \mathbf{z}} \right| = \left| \det(\mathbf{I} + \mathbf{u}\psi(\mathbf{z})^\top) \right| = |1 + \mathbf{u}^\top \psi(\mathbf{z})|$
(where $\psi(\mathbf{z}) = h'(\mathbf{w}^\top \mathbf{z} + b) \mathbf{w}$)

$\mathbf{z}_K = f_K \circ \dots \circ f_2 \circ f_1(\mathbf{z}_0)$

- before) $\ln q_K(\mathbf{z}_K) = \ln q_0(\mathbf{z}_0) - \sum_{k=1}^K \ln \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|$
- after) $\ln q_K(\mathbf{z}_K) = \ln q_0(\mathbf{z}) - \sum_{k=1}^K \ln |1 + \mathbf{u}_k^\top \psi_k(\mathbf{z}_{k-1})|$

5.1.2 Radial Flows

form : $f(\mathbf{z}) = \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0)$

- $\left| \det \frac{\partial f}{\partial \mathbf{z}} \right| = [1 + \beta h(\alpha, r)]^{d-1} [1 + \beta h(\alpha, r) + \beta h'(\alpha, r)r]$

under certain conditions...

5.1.1) Planar flows and 5.1.2) Radial Flows can be invertible!

5.2 Flow-Based Free Energy Bound

approximate our posterior distribution, with a flow of length K

$q_\phi(\mathbf{z} | \mathbf{x}) := q_K(\mathbf{z}_K)$

$$\begin{aligned} \mathcal{F}(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log q_\phi(\mathbf{z} | \mathbf{x}) - \log p(\mathbf{x}, \mathbf{z})] \\ &= \mathbb{E}_{q_0(\mathbf{z}_0)} [\ln q_K(\mathbf{z}_K) - \log p(\mathbf{x}, \mathbf{z}_K)] \\ &= \mathbb{E}_{q_0(\mathbf{z}_0)} [\ln q_0(\mathbf{z}_0)] - \mathbb{E}_{q_0(\mathbf{z}_0)} [\log p(\mathbf{x}, \mathbf{z}_K)] - \mathbb{E}_{q_0(\mathbf{z}_0)} \left[\sum_{k=1}^K \ln |1 + \mathbf{u}_k^\top \psi_k(\mathbf{z}_{k-1})| \right] \end{aligned}$$

- do not need $q_k(\cdot)$, only need $q_0(\cdot)$

5.3 Algorithm Summary

Algorithm 1 Variational Inf. with Normalizing Flows

Parameters: ϕ variational, θ generative

while not converged **do**

$\mathbf{x} \leftarrow \{\text{Get mini-batch}\}$

$\mathbf{z}_0 \sim q_0(\bullet|\mathbf{x})$

$\mathbf{z}_K \leftarrow f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0)$

$\mathcal{F}(\mathbf{x}) \approx \mathcal{F}(\mathbf{x}, \mathbf{z}_K)$

$\Delta\theta \propto -\nabla_{\theta}\mathcal{F}(\mathbf{x})$

$\Delta\phi \propto -\nabla_{\phi}\mathcal{F}(\mathbf{x})$

end while
