

[Paper review 33]

Improved Variational Inference with Inverse Autoregressive Flow

(Kingma, et al. 2016)

[Contents]

1. Abstract
2. Introduction
3. Variational inference and Learning
 1. Requirements for computational tractability
 2. Normalizing Flow
4. Inverse Autoregressive Transformations
5. Inverse Autoregressive Flow (IAF)
6. Summary

1. Abstract

NF (Normalizing Flow)

- strategy for flexible VI

"Inverse Autoregressive Flow (IAF)"

- new type of NF, which scales well to high-dimensional latent spaces
- consists of chain of invertible transformations
 - (each transformation is based on a "Autoregressive NN")

2. Introduction

SVI : scalable posterior inference, using stochastic gradient update

VAE : NN with inference network & generative network

IAF : scales well to high-dimensional latent space

IAF (Inverse Autoregressive Flow)

- lie in Gaussian autoregressive functions

(normally used for density estimation)

- input : variable with some specific ordering
- output : mean and std for each element

ex) RNN, MADE, PixelCNN, WaveNet

- such functions can....
 - be turned into "invertible" nonlinear transformations
 - with a "simple Jacobian determinant"
- flexibility + known determinant = Normalizing Flow

demonstrate this method by "improving the inference network" of deep VAE

3. Variational inference and Learning

notation

- x : observed variable
- z : latent variable
- $p(x, z)$: joint pdf → "generative model"

Maximize ELBO

- $\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z} | \mathbf{x})] = \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}) \log p(\mathbf{x}) - D_{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x}))$
- re-parameterization trick in $q(\mathbf{z} | \mathbf{x})$

Models with multiple latent variable

- factorize $q(\mathbf{z} | \mathbf{x})$
(factorize into partial inference models with some ordering)
- $q(\mathbf{z}_a, \mathbf{z}_b | \mathbf{x}) = q(\mathbf{z}_a | \mathbf{x}) q(\mathbf{z}_b | \mathbf{z}_a, \mathbf{x})$

3-1. Requirements for computational tractability

have to efficiently optimize ELBO!

Computationally efficient to...

- 1) compute and differentiate $q(z | x)$
- 2) sample from it
(since both these operations need to be performed for each datapoint in a minibatch at every iteration of optimization)

example) diagonal posterior

- $q(\mathbf{z} \mid \mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\sigma}^2(\mathbf{x}))$,
- but not much flexible....

3-2. Normalizing Flow

NF

- in the context of SGVI (Stochastic Gradient Variational Inference)
- build flexible posterior distribution using "iterative procedure"

$$\mathbf{z}_0 \sim q(\mathbf{z}_0 \mid \mathbf{x}), \quad \mathbf{z}_t = \mathbf{f}_t(\mathbf{z}_{t-1}, \mathbf{x}) \quad \forall t = 1 \dots T$$

$$\log q(\mathbf{z}_T \mid \mathbf{x}) = \log q(\mathbf{z}_0 \mid \mathbf{x}) - \sum_{t=1}^T \log \det \left| \frac{d\mathbf{z}_t}{d\mathbf{z}_{t-1}} \right|$$

$$\mathbf{f}_t(\mathbf{z}_{t-1}) = \mathbf{z}_{t-1} + \mathbf{u}h(\mathbf{w}^T \mathbf{z}_{t-1} + b)$$

- $\mathbf{u}h(\mathbf{w}^T \mathbf{z}_{t-1} + b)$ can be interpreted as a MLP with "bottleneck hidden layer" with a single unit

→ problem : long chain of transform is needed in high-dimension

4. Inverse Autoregressive Transformations

for NF that scales well to high-dimensional space...

- consider Gaussian version of autoregressive AE
(ex. MADE, PixeCNN)

Notation

- $[\boldsymbol{\mu}(\mathbf{y}), \boldsymbol{\sigma}(\mathbf{y})]$: function of the vector \mathbf{y} , to the vectors $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$
- $[\mu_i(\mathbf{y}_{1:i-1}), \sigma_i(\mathbf{y}_{1:i-1})]$: predicted mean and standard deviation of the i -th element of \mathbf{y}
(Due to the autoregressive structure, Jacobian is lower triangular with zeros on the diagonal:
 $\partial[\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i] / \partial \mathbf{y}_j = [0, 0]$ for $j \geq i$)
- $\epsilon \sim \mathcal{N}(0, \mathbf{I})$: sample from noise vector
- \mathbf{y}
 - $y_0 = \mu_0 + \sigma_0 \odot \epsilon_0$
 - $y_i = \mu_i(\mathbf{y}_{1:i-1}) + \sigma_i(\mathbf{y}_{1:i-1}) \cdot \epsilon_i$ where $i > 0$

Computation involved in this transformation : proportional to dimension D

- but inverse transformation is interesting in NF!

$$\epsilon_i = \frac{y_i - \mu_i(\mathbf{y}_{1:i-1})}{\sigma_i(\mathbf{y}_{1:i-1})}$$

2 key observations

- 1) inverse transformation can be parallelized : $\epsilon = (\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})) / \boldsymbol{\sigma}(\mathbf{y})$
(individual element ϵ_i do not depend on each other!)
- 2) inverse autoregressive operation has a simple Jacobian determinant
due to the autoregressive structure ($\partial [\mu_i, \sigma_i] / \partial y_j = [0, 0]$ for $j \geq i$)
 - $\partial \epsilon_i / \partial y_j = 0$ for $j > i$
 - $\partial \epsilon_i / \partial y_i = \sigma_i$ (simple diagonal)

Thus, $\log \det \left| \frac{d\epsilon}{dy} \right| = \sum_{i=1}^D -\log \sigma_i(\mathbf{y})$

5. Inverse Autoregressive Flow (IAF)

based on $\epsilon = (\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})) / \boldsymbol{\sigma}(\mathbf{y})$ (Inverse Autoregressive Transformations)

Algorithm 1: Pseudo-code of an approximate posterior with Inverse Autoregressive Flow (IAF)

Data:

\mathbf{x} : a datapoint, and optionally other conditioning information

θ : neural network parameters

EncoderNN($\mathbf{x}; \theta$): encoder neural network, with additional output \mathbf{h}

AutoregressiveNN[*]($\mathbf{z}, \mathbf{h}; \theta$): autoregressive neural networks, with additional input \mathbf{h}

sum(.): sum over vector elements

sigmoid(.): element-wise sigmoid function

Result:

\mathbf{z} : a random sample from $q(\mathbf{z}|\mathbf{x})$, the approximate posterior distribution

l : the scalar value of $\log q(\mathbf{z}|\mathbf{x})$, evaluated at sample ' \mathbf{z} '

```

 $[\boldsymbol{\mu}, \boldsymbol{\sigma}, \mathbf{h}] \leftarrow \text{EncoderNN}(\mathbf{x}; \theta)$ 
 $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
 $\mathbf{z} \leftarrow \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}$ 
 $l \leftarrow -\text{sum}(\log \boldsymbol{\sigma} + \frac{1}{2} \boldsymbol{\epsilon}^2 + \frac{1}{2} \log(2\pi))$ 
for  $t \leftarrow 1$  to  $T$  do
   $[\mathbf{m}, \mathbf{s}] \leftarrow \text{AutoregressiveNN}[t](\mathbf{z}, \mathbf{h}; \theta)$ 
   $\boldsymbol{\sigma} \leftarrow \text{sigmoid}(\mathbf{s})$ 
   $\mathbf{z} \leftarrow \boldsymbol{\sigma} \odot \mathbf{z} + (1 - \boldsymbol{\sigma}) \odot \mathbf{m}$ 
   $l \leftarrow l - \text{sum}(\log \boldsymbol{\sigma})$ 
end

```

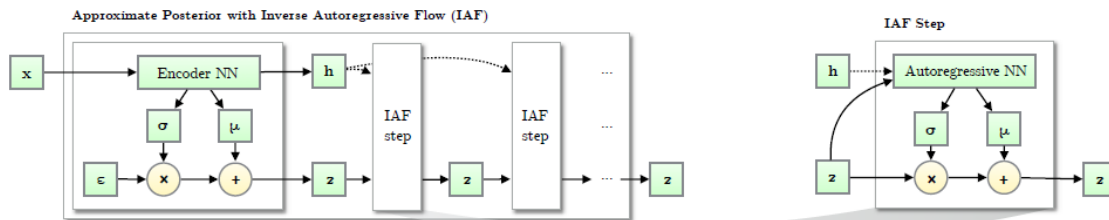


Figure 2: Like other normalizing flows, drawing samples from an approximate posterior with Inverse Autoregressive Flow (IAF) consists of an initial sample \mathbf{z} drawn from a simple distribution, such as a Gaussian with diagonal covariance, followed by a chain of nonlinear invertible transformations of \mathbf{z} , each with a simple Jacobian determinants.

chain of T : $\mathbf{z}_t = \boldsymbol{\mu}_t + \boldsymbol{\sigma}_t \odot \mathbf{z}_{t-1}$

autoregressive w.r.t. \mathbf{z}_{t-1}

- Jacobians $\frac{d\mu_t}{d\mathbf{z}_{t-1}}$ and $\frac{d\sigma_t}{d\mathbf{z}_{t-1}}$ are triangular with zeros on the diagonal.

Thus, $\frac{d\mathbf{z}_t}{d\mathbf{z}_{t-1}}$ is triangular with σ_t on the diagonal, with determinant $\prod_{i=1}^D \sigma_{t,i}$

$$\log q(\mathbf{z}_T | \mathbf{x}) = - \sum_{i=1}^D \left(\frac{1}{2} \epsilon_i^2 + \frac{1}{2} \log(2\pi) + \sum_{t=0}^T \log \sigma_{t,i} \right)$$

Autoregressive Neural Network

$$\begin{aligned} [\mathbf{m}_t, \mathbf{s}_t] &\leftarrow \text{AutoregressiveNN}[t](\mathbf{z}_t, \mathbf{h}; \boldsymbol{\theta}) \\ \boldsymbol{\sigma}_t &= \text{sigmoid}(\mathbf{s}_t) \\ \mathbf{z}_t &= \boldsymbol{\sigma}_t \odot \mathbf{z}_{t-1} + (1 - \boldsymbol{\sigma}_t) \odot \mathbf{m}_t \end{aligned}$$

Autoregressive NN form a rich family of nonlinear transformations for IAF!

6. Summary

(by Coursera)

Inverse Autoregressive Flow (IAF)

The inverse autoregressive flow reverses the dependencies to make the forward pass parallelisable, but the inverse pass sequential.

It uses the same equations:

$$x_i = \mu_i + \exp(\sigma_i) z_i \quad i = 1, \dots, D$$

but has the scale and shift functions depend on the z_i instead of the x_i :

$$\mu_i = f_{\mu_i}(z_1, \dots, z_{i-1}) \quad \sigma_i = f_{\sigma_i}(z_1, \dots, z_{i-1}).$$

Note that now the forward equation (determining \mathbf{x} from \mathbf{z}) can be parallelised, but the reverse transformations require determining z_1 , followed by z_2 , etc. and must hence be solved in sequence.