

Улучшенный вариационный вывод с обратным авторегрессионным потоком

Плетнев Никита Вячеславович

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем

15 декабря 2020 г.

Необходимо предложить новый вид Normalizing Flows, который будет хорошо масштабироваться на латентные пространства большой размерности.

Для чего это нужно

Normalizing Flows обеспечивает общую стратегию гибкого вариационного вывода на латентных переменных, что используется при построении вариационных автокодировщиков. Однако ранее опубликованные примеры потоков плохо масштабируются на латентные пространства большой размерности: требование параллелизации накладывает ограничения на апостериорное распределение. С другой стороны, оно должно быть достаточно гибким, чтобы приближать истинное распределение.

«Игрушечный» эксперимент на синтетических данных

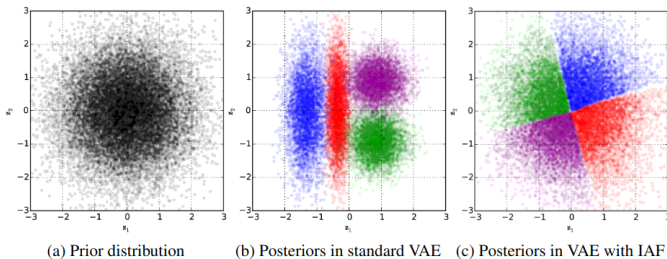


Figure 1: Best viewed in color. We fitted a variational auto-encoder (VAE) with a spherical Gaussian prior, and with factorized Gaussian posteriors **(b)** or inverse autoregressive flow (IAF) posteriors **(c)** to a toy dataset with four datapoints. Each colored cluster corresponds to the posterior distribution of one datapoint. IAF greatly improves the flexibility of the posterior distributions, and allows for a much better fit between the posteriors and the prior.

Обозначения

x — вектор наблюдаемых переменных,

z — вектор латентных переменных,

$p(x, z)$ — генеративная модель для совместного распределения,

$X = \{x_1, \dots, x_N\}$ — датасет.

Цель оптимизации

Как обычно, хотим максимизировать правдоподобие:

$$\log p(X) = \sum_{i=1}^N \log p(x_i).$$

Но в общем случае эта величина (или её производная) невычислима, поэтому применяется вариационная нижняя оценка:

$$\begin{aligned} \log p(x) &\geq \mathcal{L}(x, \theta) = \mathbb{E}_{q(z|x)}[\log p(x, z) - \log q(z|x)] = \\ &= \log p(x) - D_{KL}(q(z|x) || p(z|x)). \end{aligned}$$

- Апостериорное распределение $q(z|x)$ должно быть эффективно вычислимо и дифференцируемо;
- Если размерность латентного пространства велика, то вычисление должно быть параллелизуемо.
- Это создаёт ограничения на семейство апостериорных распределений;
- С другой стороны, оно должно быть достаточно гибким, чтобы соответствовать истинному апостериорному распределению $p(z|x)$.

Normalizing flows (Rezende and Mohamed, 2015)

Цепочка обратимых параметризованных преобразований:

$$z_0 \sim q(z_0|x), \quad z_t = f_t(z_{t-1}|x) \quad \forall t = 1..T$$

Если все якобианы вычислимы, то на последней итерации:

$$\log q(z_T|x) = \log q(z_0|x) - \sum_{t=1}^T \log \det \left| \frac{dz_t}{dz_{t-1}} \right|.$$

Сами авторы экспериментировали только с ограниченным семейством

$$f_t(z_{t-1}) = z_{t-1} + uh(w^T z_{t-1} + b).$$

(MLP с одним нейроном). Для обнаружения закономерностей высокой размерности нужна очень длинная цепочка.

Inverse Autoregressive Transformations

Основная идея

Для масштабируемости рассматриваем гауссовы версии VAE, такие как MADE, PixelCNN.

Переменная y моделируется: $y = \{y_i\}_{i=1}^D$.

$[\mu(y), \sigma(y)]$ — функция из вектора y в пару векторов (μ, σ) .

Из-за авторегрессионной структуры матрица Якоби — нижнетреугольная с нулями на диагонали.

Сэмплирование

Последовательное преобразование «шума» $\epsilon \sim \mathcal{N}(0, I)$,

$y_0 = \mu_0 + \sigma_0 \cdot \epsilon_0$, $y_i = \mu_i(y_{1:i-1}) + \sigma_i(y_{1:i-1}) \cdot \epsilon_i$.

У нас есть параллелизуемое обратное преобразование:

$$\epsilon = \frac{y - \mu(y)}{\sigma(y)}.$$

И вычисляемый якобиан: $\log \det \left| \frac{d\epsilon}{dy} \right| = - \sum_{i=1}^D \log \sigma_i(y)$

Схема

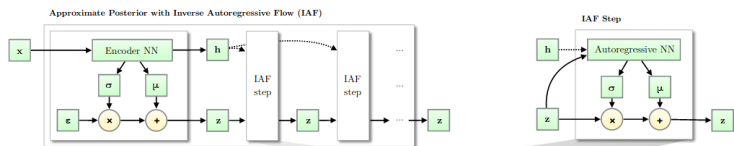


Figure 2: Like other normalizing flows, drawing samples from an approximate posterior with Inverse Autoregressive Flow (IAF) consists of an initial sample z drawn from a simple distribution, such as a Gaussian with diagonal covariance, followed by a chain of nonlinear invertible transformations of z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

Формулы

$$z_0 = \mu_0 + \sigma_0 \odot \epsilon; \quad z_t = \mu_t + \sigma_t \odot z_{t-1}.$$

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

Алгоритм

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

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

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

$\text{sum}(\cdot)$: sum over vector elements

$\text{sigmoid}(\cdot)$: 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} '

$[\mu, \sigma, \mathbf{h}] \leftarrow \text{EncoderNN}(\mathbf{x}; \theta)$

$\epsilon \sim \mathcal{N}(0, I)$

$\mathbf{z} \leftarrow \sigma \odot \epsilon + \mu$

$l \leftarrow -\text{sum}(\log \sigma + \frac{1}{2}\epsilon^2 + \frac{1}{2}\log(2\pi))$

for $t \leftarrow 1$ **to** T **do**

$[\mathbf{m}, \mathbf{s}] \leftarrow \text{AutoregressiveNN}[t](\mathbf{z}, \mathbf{h}; \theta)$

$\sigma \leftarrow \text{sigmoid}(\mathbf{s})$

$\mathbf{z} \leftarrow \sigma \odot \mathbf{z} + (1 - \sigma) \odot \mathbf{m}$

$l \leftarrow l - \text{sum}(\log \sigma)$

end

Преобразование как в LSTM для стабильности вычислений.

Эксперимент на MNIST

Model	VLB	$\log p(\mathbf{x}) \approx$
Convolutional VAE + HVI [1]	-83.49	-81.94
DLGM 2hl + IWAE [2]		-82.90
LVAE [3]		-81.74
DRAW + VGP [4]	-79.88	
Diagonal covariance	-84.08 (± 0.10)	-81.08 (± 0.08)
IAF (Depth = 2, Width = 320)	-82.02 (± 0.08)	-79.77 (± 0.06)
IAF (Depth = 2, Width = 1920)	-81.17 (± 0.08)	-79.30 (± 0.08)
IAF (Depth = 4, Width = 1920)	-80.93 (± 0.09)	-79.17 (± 0.08)
IAF (Depth = 8, Width = 1920)	-80.80 (± 0.07)	-79.10 (± 0.07)

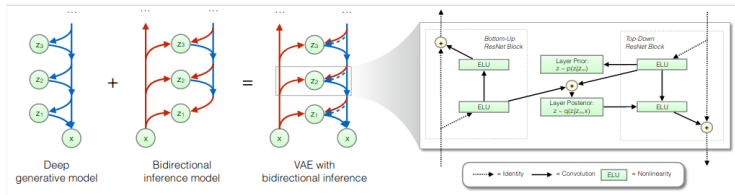


Figure 3: Overview of our ResNet VAE with bidirectional inference. The posterior of each layer is parameterized by its own IAF.

Эксперимент на CIFAR-10: латентные вектора в памяти

Table 2: Our results with ResNet VAEs on CIFAR-10 images, compared to earlier results, in *average number of bits per data dimension* on the test set. The number for convolutional DRAW is an upper bound, while the ResNet VAE log-likelihood was estimated using importance sampling.

Method	bits/dim \leq
<i>Results with tractable likelihood models:</i>	
Uniform distribution (van den Oord et al., 2016b)	8.00
Multivariate Gaussian (van den Oord et al., 2016b)	4.70
NICE (Dinh et al., 2014)	4.48
Deep GMMs (van den Oord and Schrauwen, 2014)	4.00
Real NVP (Dinh et al., 2016)	3.49
PixelRNN (van den Oord et al., 2016b)	3.00
Gated PixelCNN (van den Oord et al., 2016c)	3.03
<i>Results with variationally trained latent-variable models:</i>	
Deep Diffusion (Sohl-Dickstein et al., 2015)	5.40
Convolutional DRAW (Gregor et al., 2016)	3.58
ResNet VAE with IAF (Ours)	3.11

При этом наш метод сэмплирует изображение за 0.05 секунд, а PixelNet — за 52 секунды.

Предложены новый тип потоков — Inverse Autoregressive Flow, который хорошо масштабируется в латентные пространства большой размерности.

Проведены эксперименты, показывающие эффективность подхода: получены близкие к наилучшим результаты при намного более быстром сэмплинговании.

- [1] Diederik P. Kingma, Tim Salimans: Improved Variational Inference with Inverse Autoregressive Flow, <https://arxiv.org/pdf/1606.04934.pdf>