

DO DEEP GENERATIVE MODELS KNOW WHAT THEY DON'T KNOW?

Авторы статьи Eric Nalisnick, Akihiro Matsukawa,
Yee Whye Teh, Dilan Gorur, Balaji Lakshminarayanan
Презентация Dayana Savostianova

НИУ ВШЭ

27 02 2020

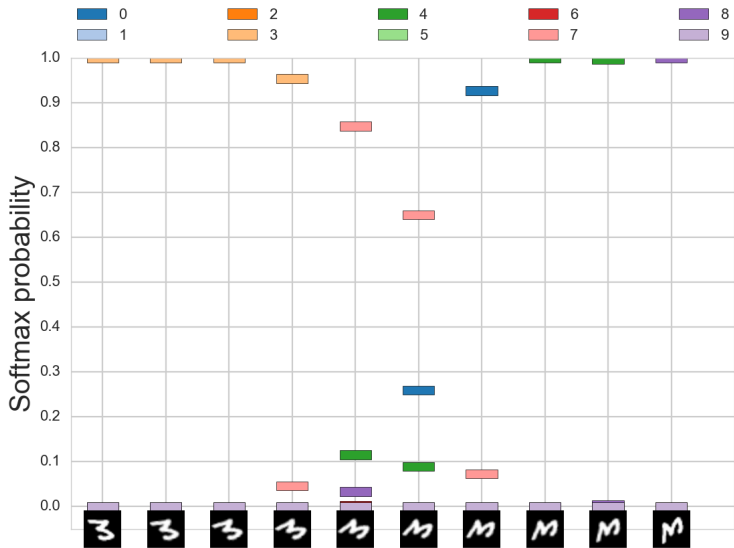


Рис.: Rotated MNIST

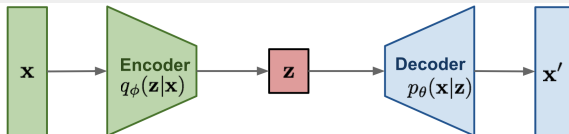
В процессе обучения мы перестаем смотреть на биты картинок как на дискретное распределение, из-за чего могут получиться неоправданно большие значения $\log q(\mathbf{x})$.

Чтобы это решить добавим шум $\mathbf{u} \sim U[0, 1]^D$, $\mathbf{y} = \mathbf{x} + \mathbf{u}$, $p(\mathbf{x})$ – реальное распределение, $q(\mathbf{x})$ – модельное распределение.

$$\begin{aligned} \int p(\mathbf{y}) \log q(\mathbf{y}) d\mathbf{y} &= \sum_{\mathbf{x}} P(\mathbf{x}) \int_{[0,1]^D} \log q(\mathbf{x} + \mathbf{u}) d\mathbf{u} \\ &\leq \sum_{\mathbf{x}} P(\mathbf{x}) \log \int_{[0,1]^D} q(\mathbf{x} + \mathbf{u}) d\mathbf{u} \\ &= \sum_{\mathbf{x}} P(\mathbf{x}) \log Q(\mathbf{x}), \end{aligned}$$

Тогда $BPD = - \sum_{\mathbf{x}} P(\mathbf{x}) \log Q(\mathbf{x})$

VAE: maximize ELBO.



**Flow-based
generative models:**
minimize the negative
log-likelihood

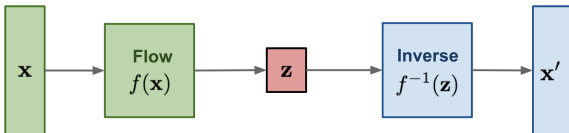


Рис.: Модели

Потоковые (Flow-based) модели в общем случае.

- $x \sim p^*(x)$, где $p^*(\cdot)$ – неизвестное реальное распределение.
- $p_\theta(x)$ – модель с параметрами θ .

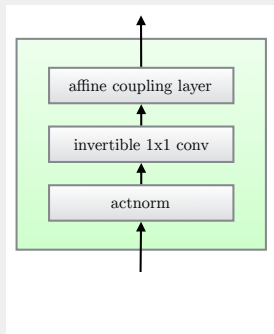
Генеративный процесс описывают следующим образом:

$z \sim p_\theta(z), x = g_\theta(z)$, где

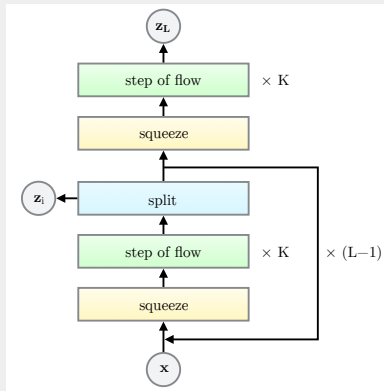
- z скрытая переменная
- $p_\theta(\cdot)$ – какое-нибудь не очень сложное распределение, например сферическая Гауссиана.
- $z = f_\theta(x) = g_\theta^{-1}(x) - f, g$ обратимые

Более того $f = f_1 \circ f_2 \circ \dots \circ f_K$ композиция (и аналогично g), такая что:

$$x \xleftrightarrow{f_1} h_1 \xleftrightarrow{f_2} h_2 \dots \xleftrightarrow{f_K} z$$



(a) One step of our flow.



(b) Multi-scale architecture.

Description	Function	Reverse Function	Log-determinant
Actnorm.	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{s} \odot \mathbf{x}_{i,j} + \mathbf{b}$	$\forall i, j : \mathbf{x}_{i,j} = (\mathbf{y}_{i,j} - \mathbf{b})/\mathbf{s}$	$h \cdot w \cdot \sum(\log \mathbf{s})$
Invertible 1×1 convolution. $\mathbf{W} : [c \times c]$.	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{W}\mathbf{x}_{i,j}$	$\forall i, j : \mathbf{x}_{i,j} = \mathbf{W}^{-1}\mathbf{y}_{i,j}$	$h \cdot w \cdot \log \det(\mathbf{W}) $ or $h \cdot w \cdot \sum(\log \mathbf{s})$
Affine coupling layer.	$\mathbf{x}_a, \mathbf{x}_b =$ $= \text{split}(\mathbf{x})$ $(\log \mathbf{s}, \mathbf{t}) =$ $= \text{NN}(\mathbf{x}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{y}_a =$ $\mathbf{s} \odot \mathbf{x}_a + \mathbf{t}$ $\mathbf{y}_b = \mathbf{x}_b$ $\mathbf{y} = \text{concat}(\mathbf{y}_a, \mathbf{y}_b)$	$\mathbf{y}_a, \mathbf{y}_b =$ $\text{split}(\mathbf{y})$ $(\log \mathbf{s}, \mathbf{t}) =$ $\text{NN}(\mathbf{y}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t})/\mathbf{s}$ $\mathbf{x}_b = \mathbf{y}_b$ $\mathbf{x} = \text{concat}(\mathbf{x}_a, \mathbf{x}_b)$	$\sum(\log(\mathbf{s}))$

Чуть-чуть математики

Хотим уменьшать расстояние между распределениями:

$$\text{KLD}[p^*(\mathbf{x})||p(\mathbf{x}; \theta)] = \int p^*(\mathbf{x}) \log \frac{p^*(\mathbf{x})}{p(\mathbf{x}; \theta)} d\mathbf{x} \approx -\frac{1}{N} \log p(\mathbf{X}; \theta) - \mathbb{H}[p^*]$$

Поскольку энтропия реального распределения константа, оптимизируем только левую часть:

$$\theta^* = \arg \max_{\theta} \log p(\mathbf{X}; \theta) = \arg \max_{\theta} \sum_{n=1}^N \log p(\mathbf{x}_n; \theta).$$

Хотим извлекать вероятности из плотностей, для этого надо учиться интегрировать $P(\Omega) = \int_{\Omega} p(\mathbf{x}; \theta) d\mathbf{x}$

Можно имитировать интегрирование возмущением:

$$\log \int p(\mathbf{x}_n + \delta; \theta) p(\delta) d\delta \geq \mathbb{E}_{\delta} [\log p(\mathbf{x}_n + \delta; \theta)] \approx \log p(\mathbf{x}_n + \tilde{\delta}; \theta)$$

где $\tilde{\delta}$ сэмпл из $p(\delta)$.

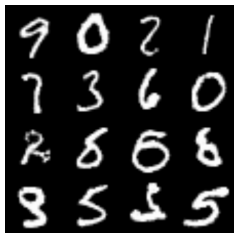
ВАЖНЫЙ ТРЮК про замену переменных:

$$p_z(f_\theta(X)) \left| \frac{df_\theta(X)}{dX} \right| = p(X)$$

$\left| \frac{df_\theta(X)}{dX} \right|$ – компенсация объема диффеоморфизма.

Тогда для Glow

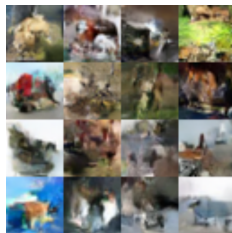
$$\theta^* = \arg \max_{\theta} \log p_X(\mathbf{X}; \theta) = \arg \max_{\phi, \psi} \sum_{n=1}^N \log p_z(f(\mathbf{x}_n; \phi); \psi) + \log \left| \frac{\partial \mathbf{f}_\phi}{\partial \mathbf{x}_n} \right|.$$



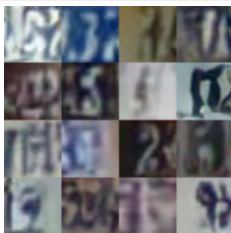
(a) MNIST



(b) FashionMNIST

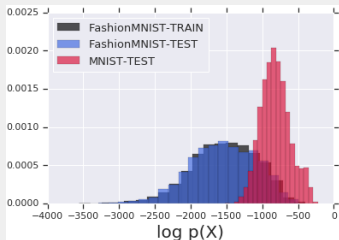


(c) CIFAR-10

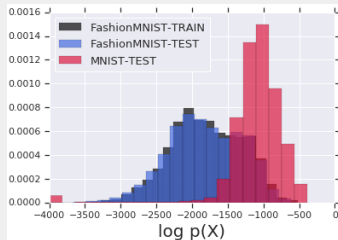


(d) SVHN

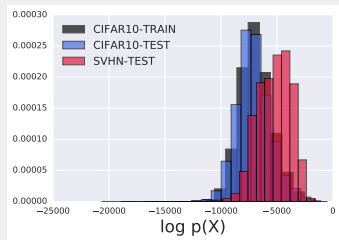
PixelCNN и VAE



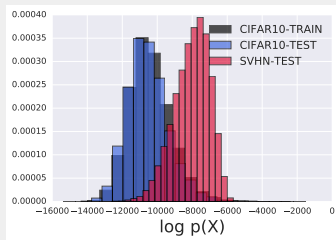
(a) PixelCNN: FashionMNIST vs MNIST



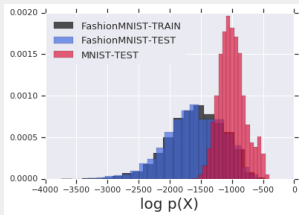
(b) VAE: FashionMNIST vs MNIST



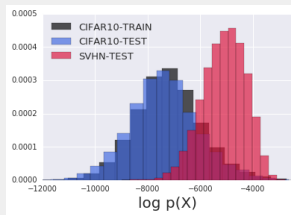
(c) PixelCNN: CIFAR-10 vs SVHN



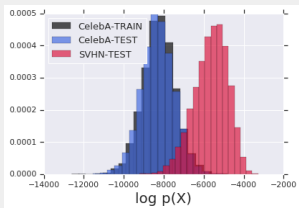
(d) VAE: CIFAR-10 vs SVHN



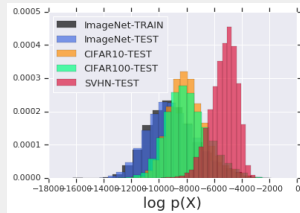
(a) Train on FashionMNIST, Test on MNIST



(b) Train on CIFAR-10, Test on SVHN

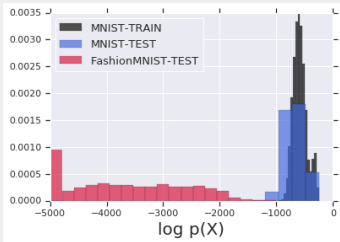


(c) Train on CelebA, Test on SVHN

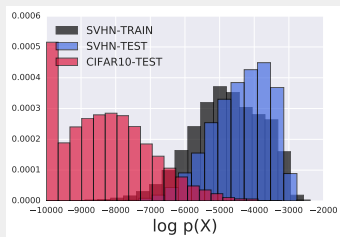


(d) Train on ImageNet,
Test on CIFAR-10 / CIFAR-100 / SVHN

Несимметричность



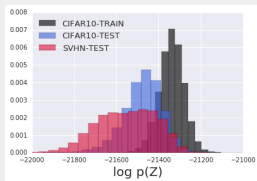
(a) Train on MNIST, Test on FashionMNIST



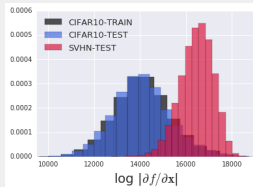
(b) Train on SVHN, Test on CIFAR-10

CV vs. NVP

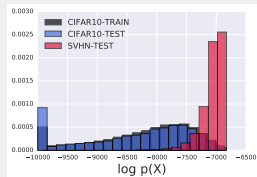
- NVP – Non Volume Preserving (классический вариант)
- CV – Constant Volume (отличается отсутствием скейлинга в ACL)



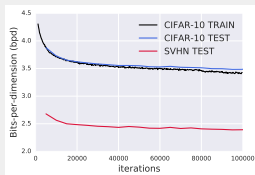
(a) CIFAR-10: $\log p(\mathbf{z})$



(b) CIFAR-10: Volume



(c) CV-Glow Likelihoods



(d) LL vs Iter.

Хотим анализировать, но сложно:

$$\mathbb{E}_q[\log p(\mathbf{x}; \theta)] - \mathbb{E}_{p^*}[\log p(\mathbf{x}; \theta)] > 0$$

Берем Тейлора

$$\begin{aligned} \log p(\mathbf{x}; \theta) &\approx \log p(\mathbf{x}_0; \theta) + \\ &+ \nabla_{\mathbf{x}_0} \log p(\mathbf{x}_0; \theta)^T (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} \text{Tr} \{ \nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \theta) (\mathbf{x} - \mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0)^T \} \end{aligned}$$

Вот сомнительная часть:

Утверждается что раз Тейлор работает только около \mathbf{x}_0 , тогда

$$\mathbb{E}_q[\mathbf{x}] = \mathbb{E}_{p^*}[\mathbf{x}] = \mathbf{x}_0$$

Тогда получаем

$$\begin{aligned} 0 < \mathbb{E}_q[\log p(\mathbf{x}; \theta)] - \mathbb{E}_{p^*}[\log p(\mathbf{x}; \theta)] &\approx \frac{1}{2} \text{Tr} \{ \nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \theta) (\Sigma_q - \Sigma_{p^*}) \} \\ &= \frac{1}{2} \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p_z(f(\mathbf{x}_0; \phi)) + \nabla_{\mathbf{x}_0}^2 \log \left| \frac{\partial \mathbf{f}_\phi}{\partial \mathbf{x}_0} \right| \right] (\Sigma_q - \Sigma_{p^*}) \right\}, \end{aligned}$$

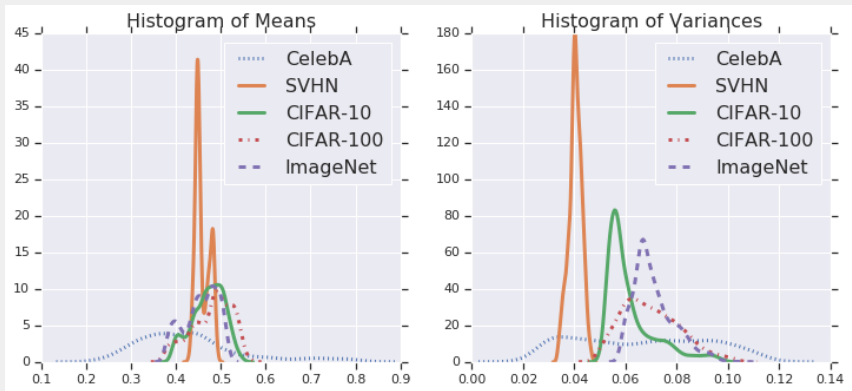


Рис.: Histogram of per-dimension means and variances (empirical).

- Пусть U_k ядро $C \times C$, где C – количество каналов, а k – номер потока.
- $\partial f_{h,w,c} / \partial x_{h,w,c} = \prod_k \sum_{j=1}^C u_{k,c,j}$
- $\partial^2 f_{h,w,c} / \partial x_{h,w,c}^2 = 0$

$$\begin{aligned} & \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \theta) \right] (\Sigma_q - \Sigma_{p^*}) \right\} \\ &= \frac{\partial^2}{\partial \mathbf{z}^2} \log p(\mathbf{z}; \psi) \sum_{c=1}^C \left(\prod_{k=1}^K \sum_{j=1}^C u_{k,c,j} \right)^2 \sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2) \end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_{\text{SVHN}}[\log p(\mathbf{x}; \boldsymbol{\theta})] - \mathbb{E}_{\text{CIFAR-10}}[\log p(\mathbf{x}; \boldsymbol{\theta})] \\
& \approx \frac{-1}{2\sigma_{\psi}^2} [\alpha_1^2(49.6 - 61.9) + \alpha_2^2(52.7 - 59.2) + \alpha_3^2(53.6 - 68.1)] \\
& = \frac{1}{2\sigma_{\psi}^2} [\alpha_1^2 \cdot 12.3 + \alpha_2^2 \cdot 6.5 + \alpha_3^2 \cdot 14.5] \geq 0
\end{aligned}$$

$$\text{где } \alpha_c = \prod_{k=1}^K \sum_{j=1}^C u_{k,c,j}$$

- Существует проблема с OOD
- ВОЗМОЖНО для FLOW объясняется разницей в дисперсии распределений
- Пока эта проблема существует, нужно думать прежде чем искать аномалии такими моделями.

- Из каких элементов состоит один слой Flow. Выписать, что происходит в каждом из них.
- Формулировка проблемы поставленной статьей.
Перечислить возможные причины ее возникновения.
- Формула замены переменных для Glow. На аккие части благодаря этому разбивается поиск оптимального параметра.