

Do Deep Generative Models Know What They Don't Know?

Balaji Lakshminarayanan, Eric Nalisnick, Akihiro Matsukawa, Yee
Whye Teh, Dilan Gorur

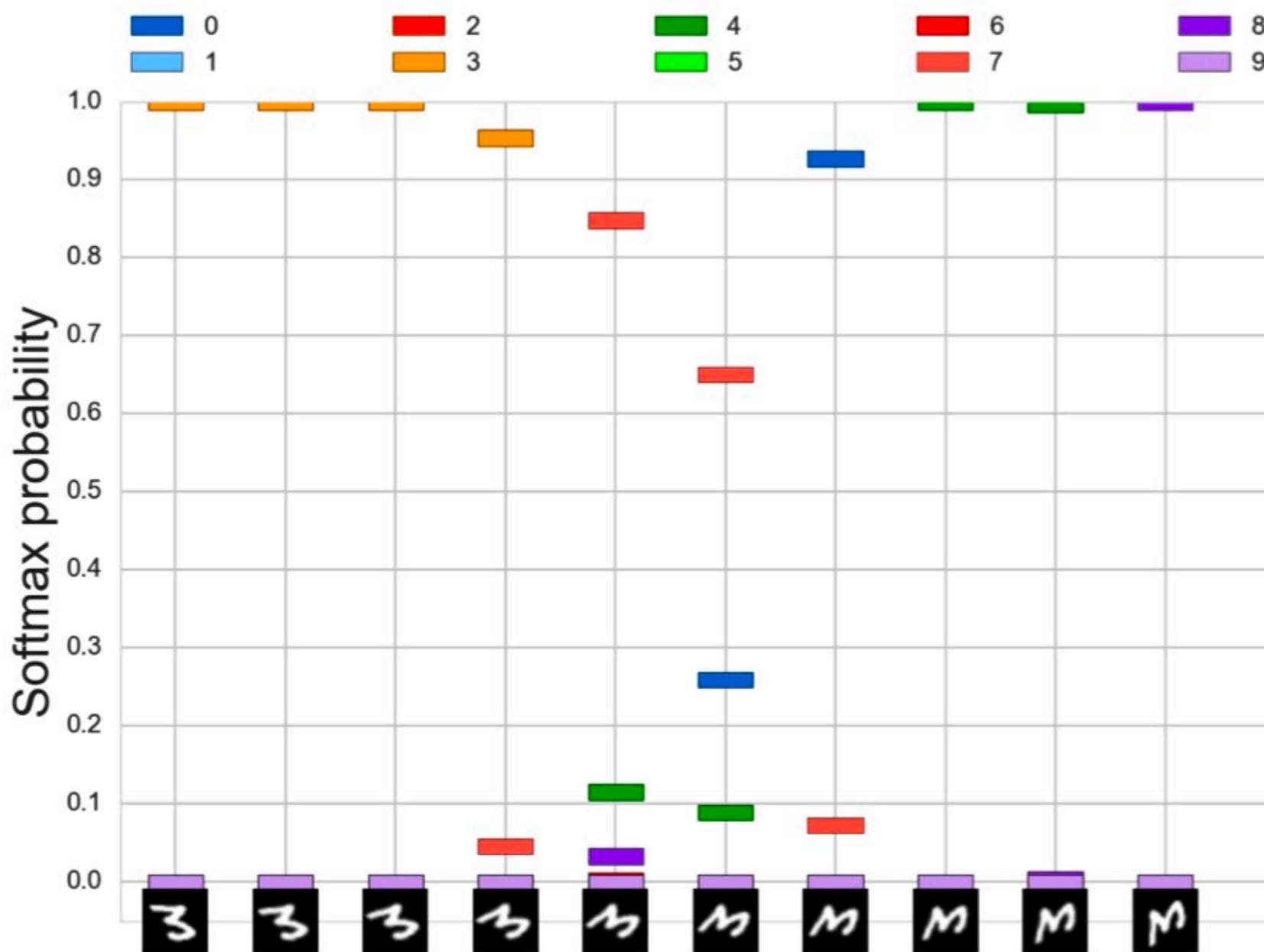


DeepMind

(не) Уверенность модели. Зачем ее оценивать?

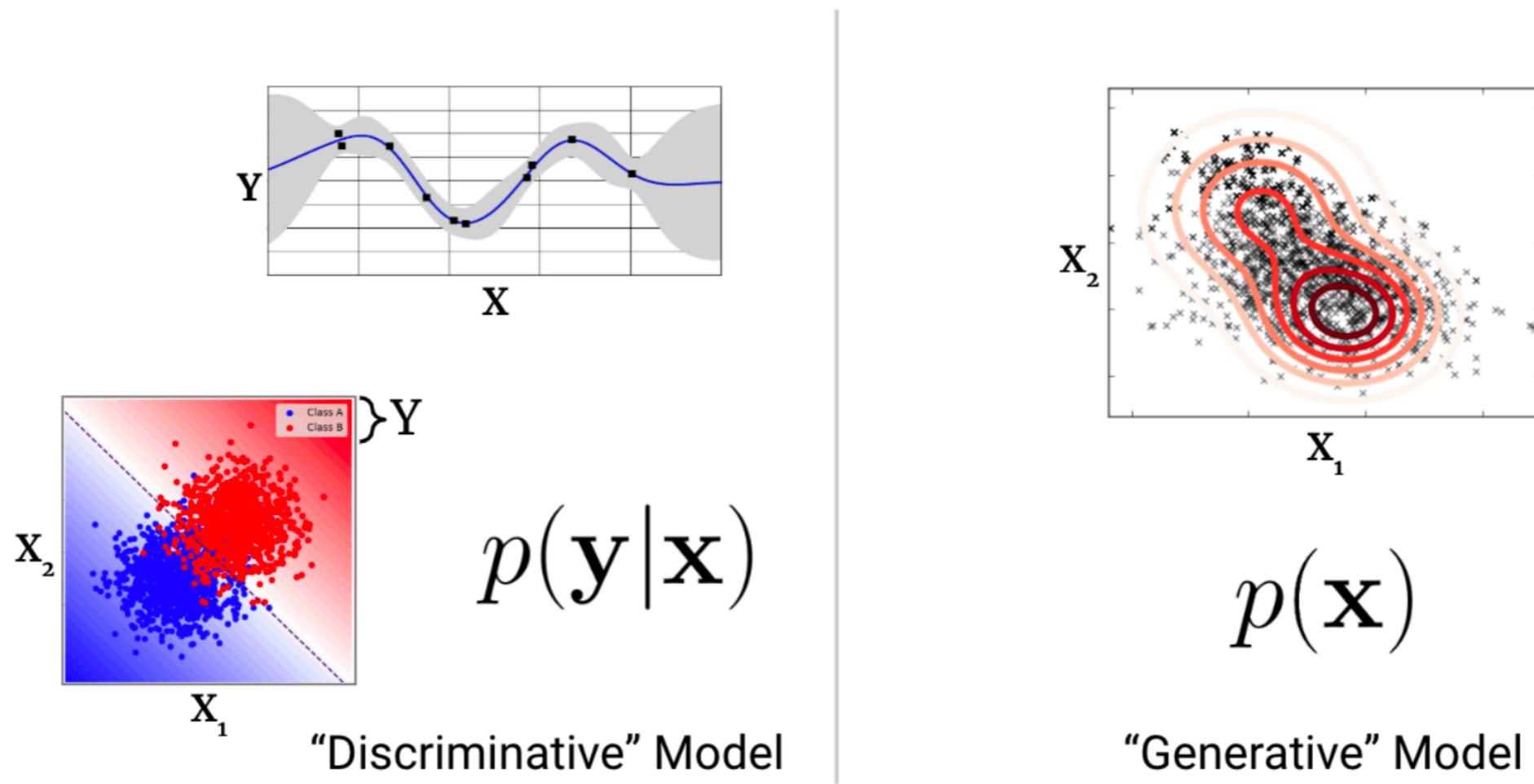
- Модели, обученные предсказывать $p(y | x)$ обычно хорошо работают для $x \sim p_{train}(X)$
- Для x из иного распределения результаты могут быть непредсказуемы
- Зная, что модель не уверена в предсказаниях мы могли бы не доверять им

(не) Уверенность модели. Зачем ее оценивать?



Поворот MNIST приводит к тому, что нейронная сеть (LeNet + weight decay) предсказывает неверному классу высокую вероятность

Дискриминативные и Генеративные и модели

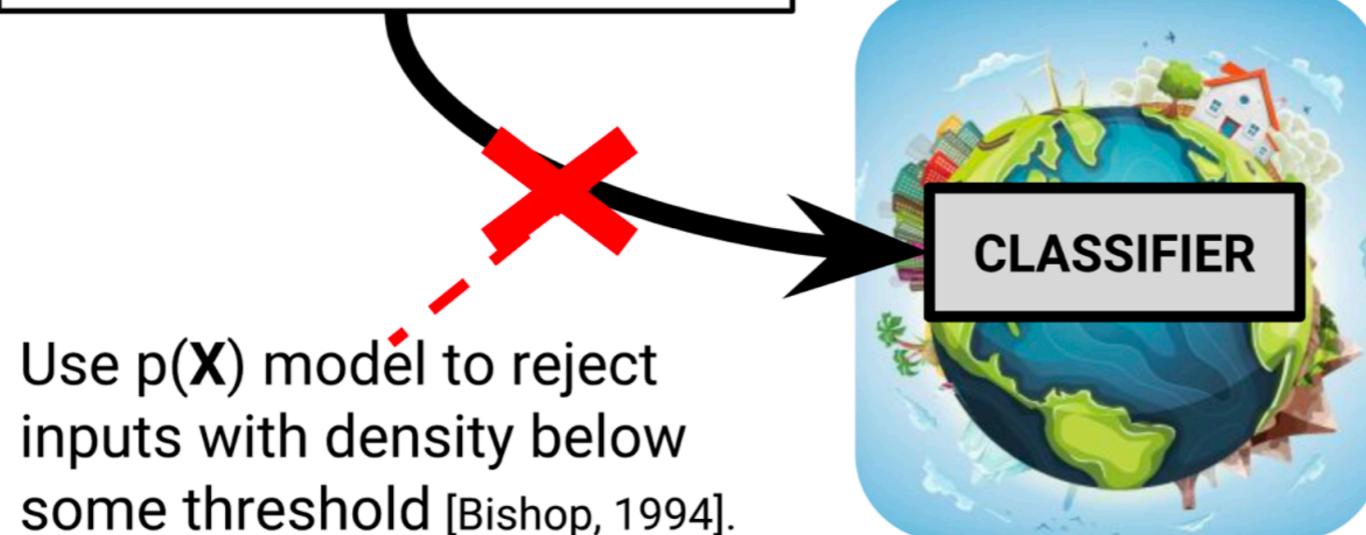


- Знаем: $p(y|x)$ ошибается на OOD (out-of-distribution)
- Идея: будем использовать генеративную модель, чтобы отличать \mathbf{x} из OOD

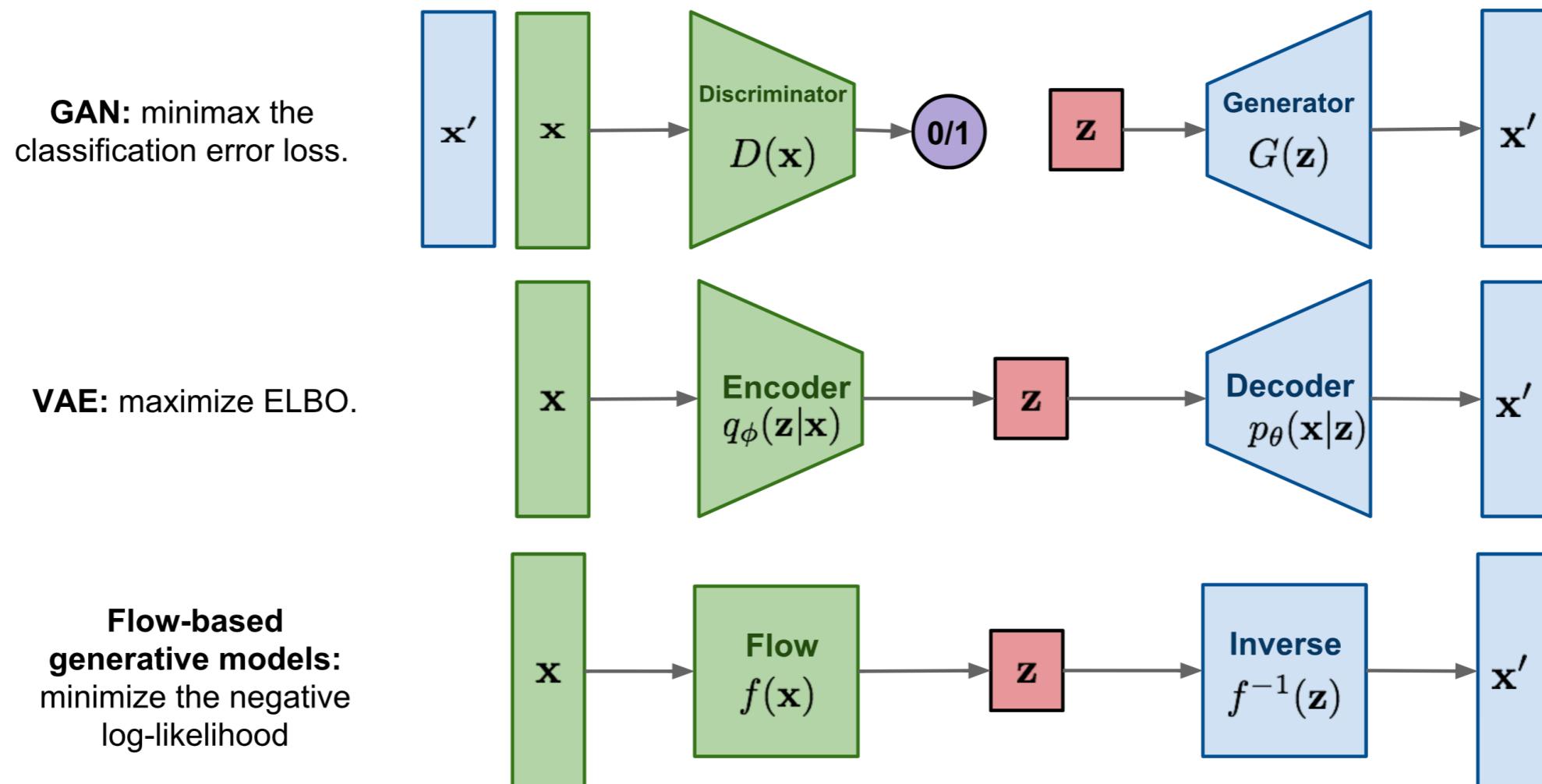
Inputs Unlike Training Data



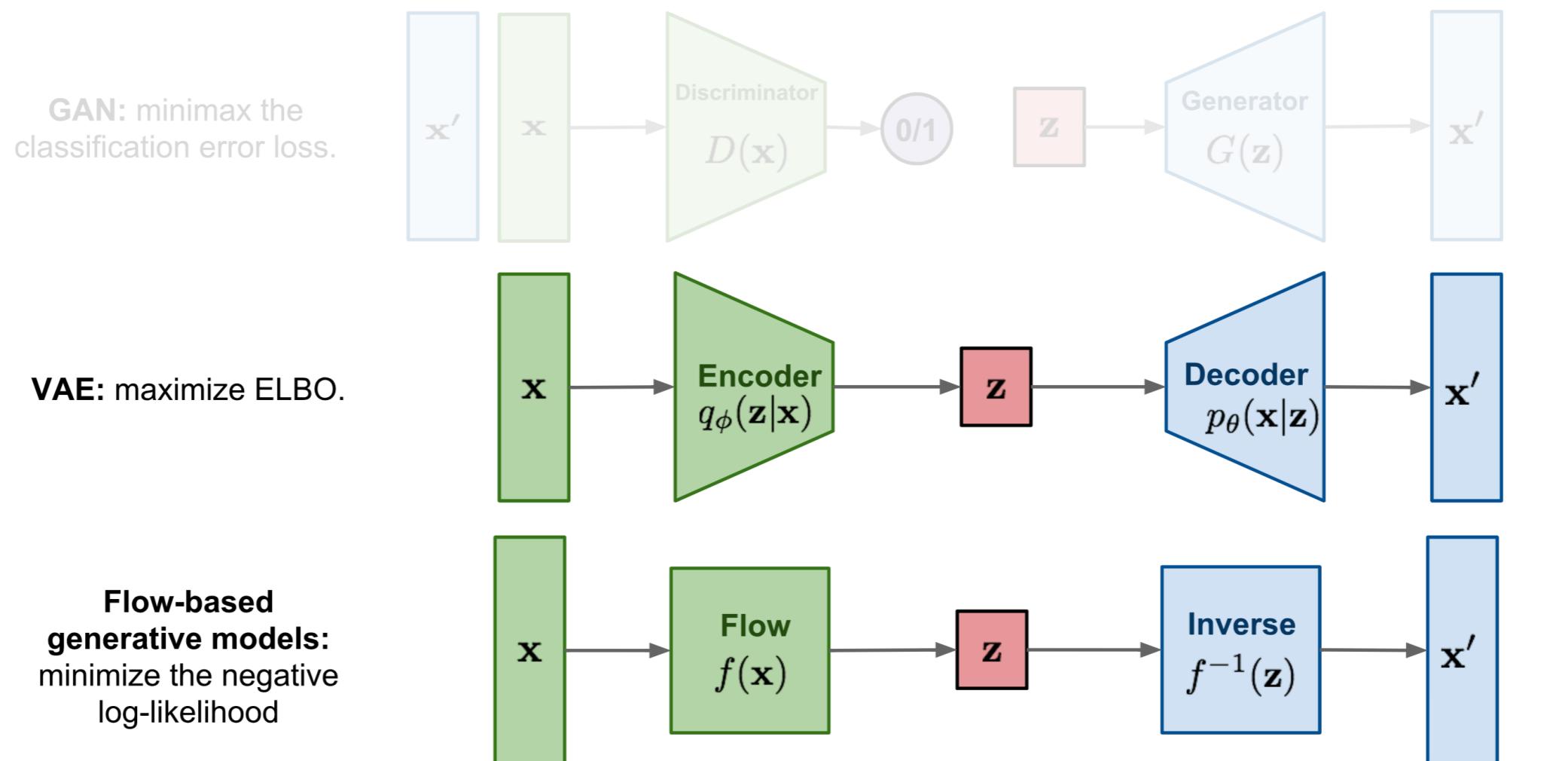
if $p(\mathbf{x}^*; \phi) < \tau$,
then reject \mathbf{x}^*



Типы генеративных моделей

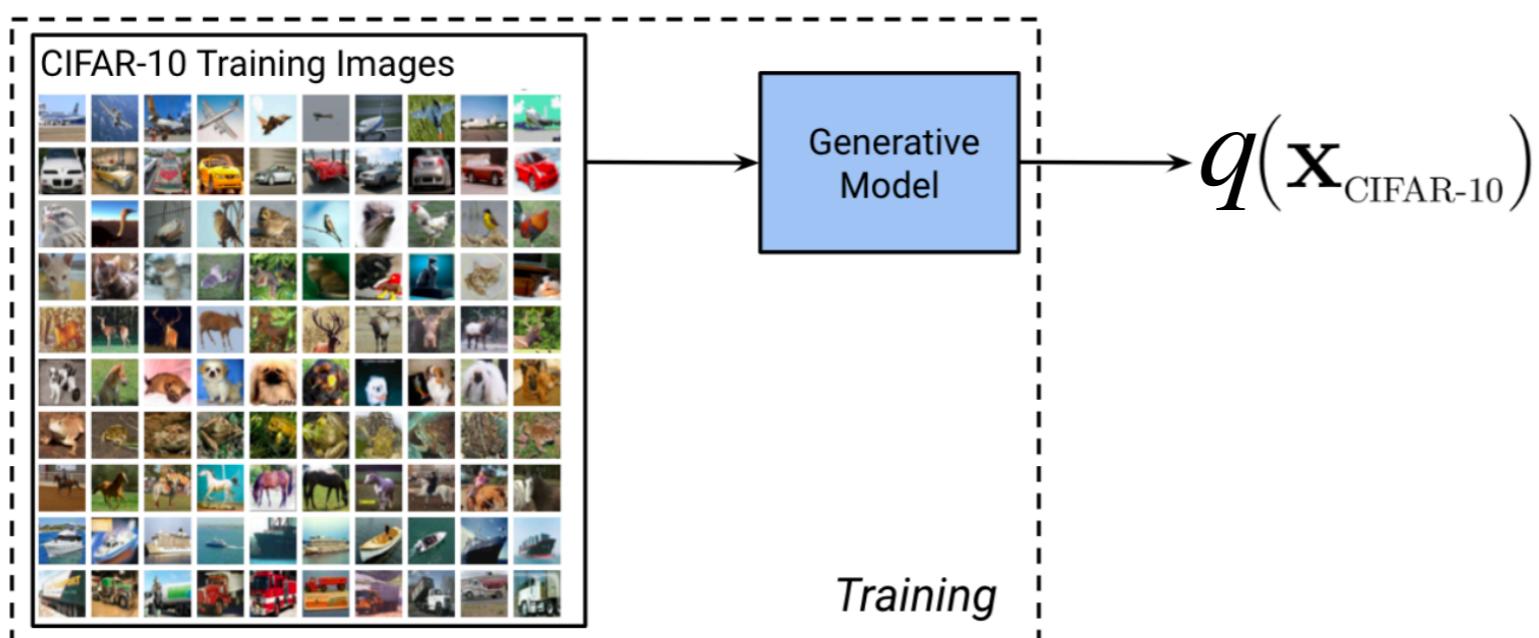


Типы генеративных моделей



- Можем вычислить $p(x)$
 - Точно для Flow-based
 - Нижнюю оценку для VAE

Оцениваем лог-правдоподобие данных: bits-per-dimension



- Обычно изображения из дискретного распределения, а модель приближает непрерывное
- Из-за этого $-\log(q(x))$ показывает неоправданно высокие значения
- Добавим к входным данным случайны шум

Оцениваем лог-правдоподобие данных: bits-per-dimension

- Пусть $p(x)$ — истинное распределение, $q(x)$ — приближённое
- Добавим к входным данным случайны шум: $y = x + u$

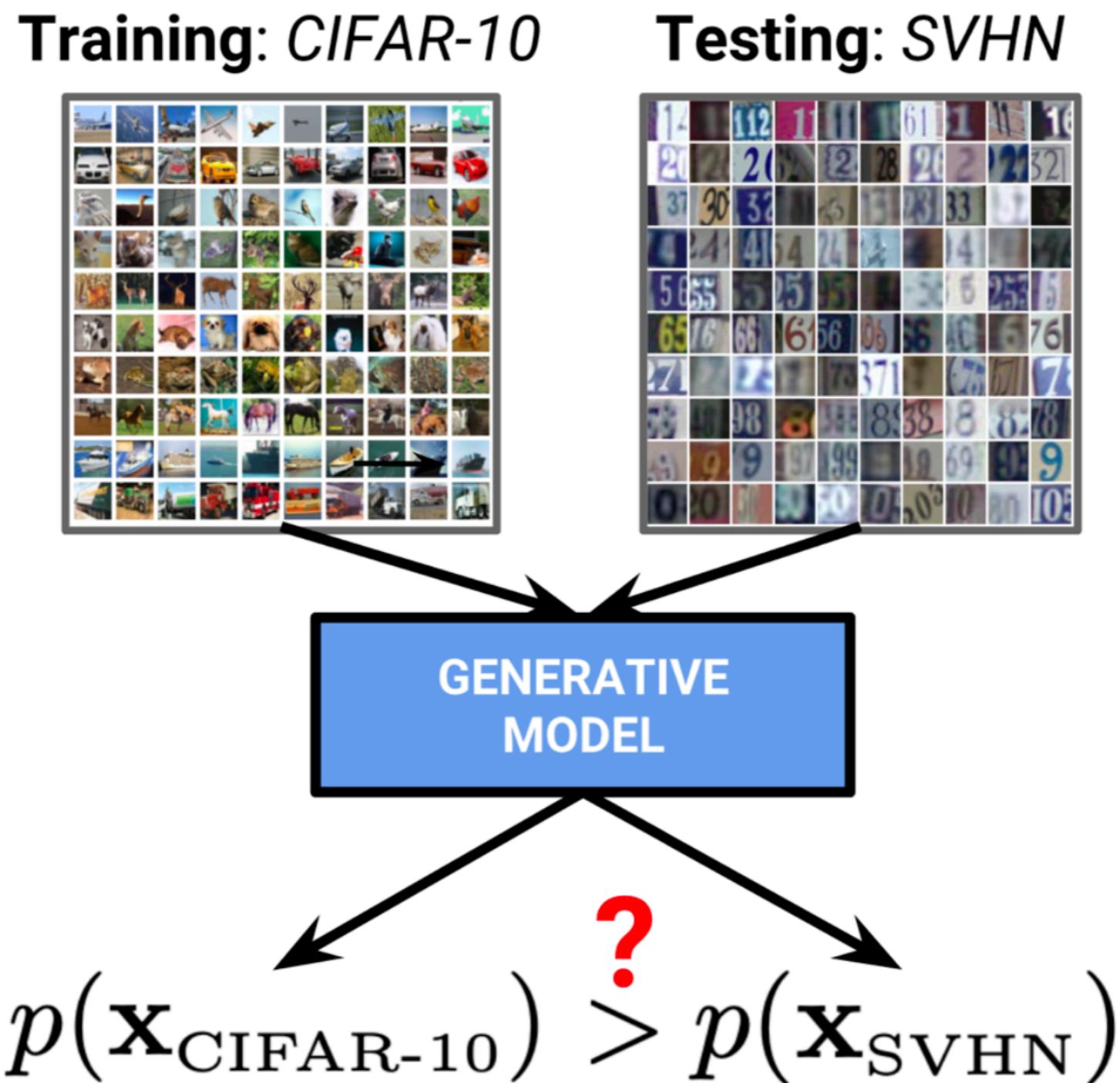
$$\underline{\int p(\mathbf{y}) \log q(\mathbf{y}) d\mathbf{y}} = \sum_{\mathbf{x}} P(\mathbf{x}) \int_{[0,1]^D} \log q(\mathbf{x} + \mathbf{u}) d\mathbf{u}$$

Усредненное лог
правдоподобие

$$\leq \sum_{\mathbf{x}} P(\mathbf{x}) \log \int_{[0,1]^D} q(\mathbf{x} + \mathbf{u}) d\mathbf{u} \\ = \sum_{\mathbf{x}} P(\mathbf{x}) \log Q(\mathbf{x}),$$

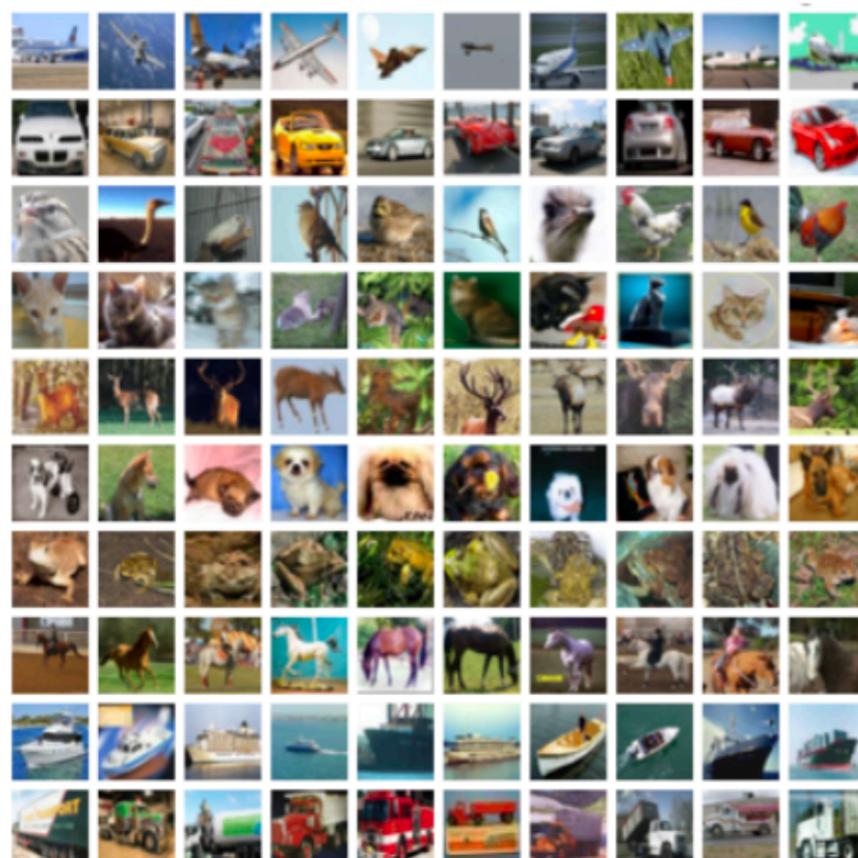
Отрицание этого
выражение — BPD

Обучаем на CIFAR тестируем на SVHN



Обучаем на CIFAR тестируем на SVHN

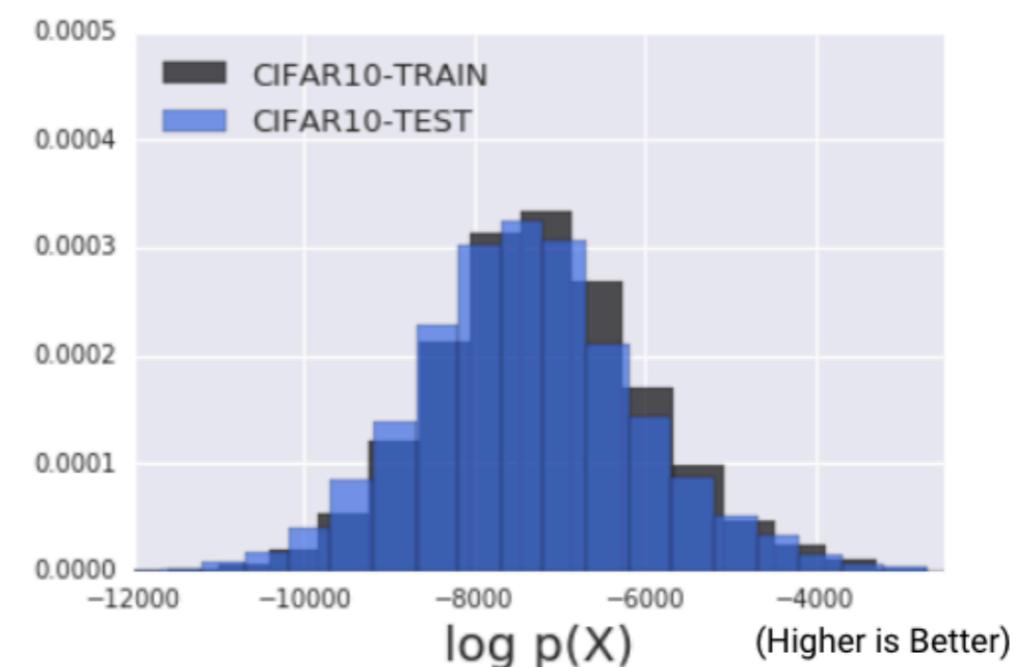
CIFAR-10 Training Images



Bits Per Dimension
(NLL / # dims / log 2)

CIFAR10-Train	3.386
CIFAR10-Test	3.464

(Lower is Better)



(Higher is Better)

Обучаем на CIFAR тестируем на SVHN

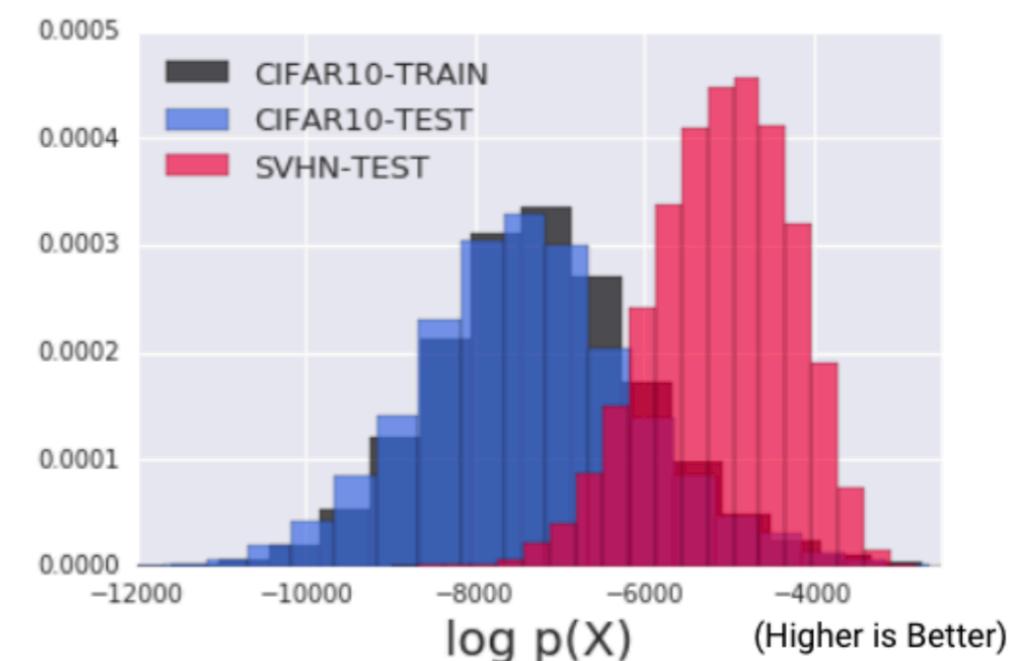
SVHN Test Images



Bits Per Dimension
(NLL / # dims / log 2)

	Bits Per Dimension (NLL / # dims / log 2)
CIFAR10-Train	3.386
CIFAR10-Test	3.464
SVHN-Test	2.389

(Lower is Better)



Обучаем на CIFAR тестируем на SVHN

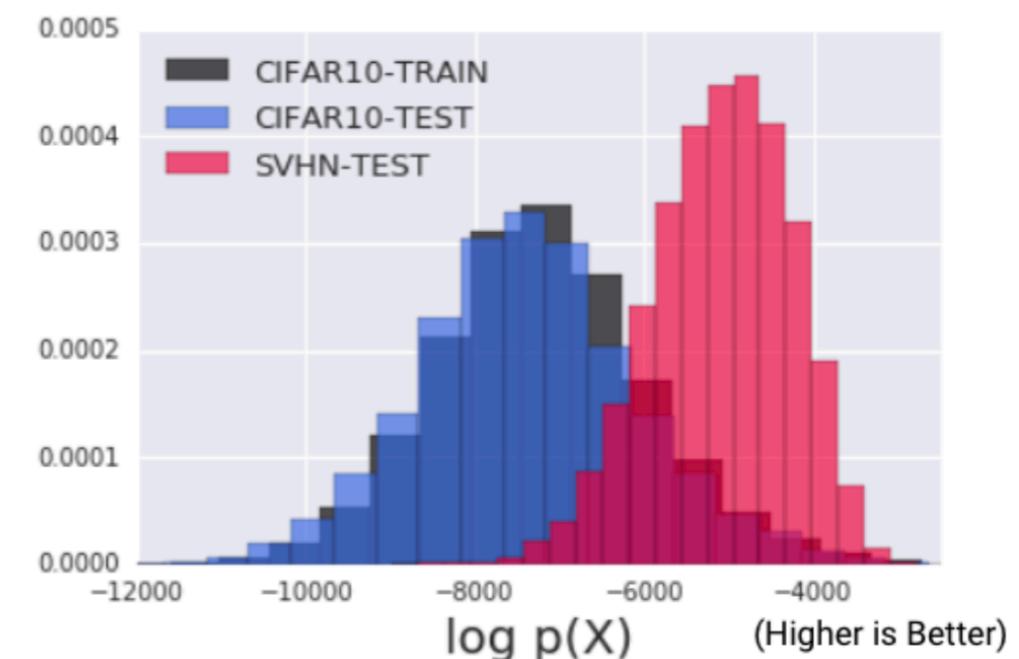
SVHN Test Images



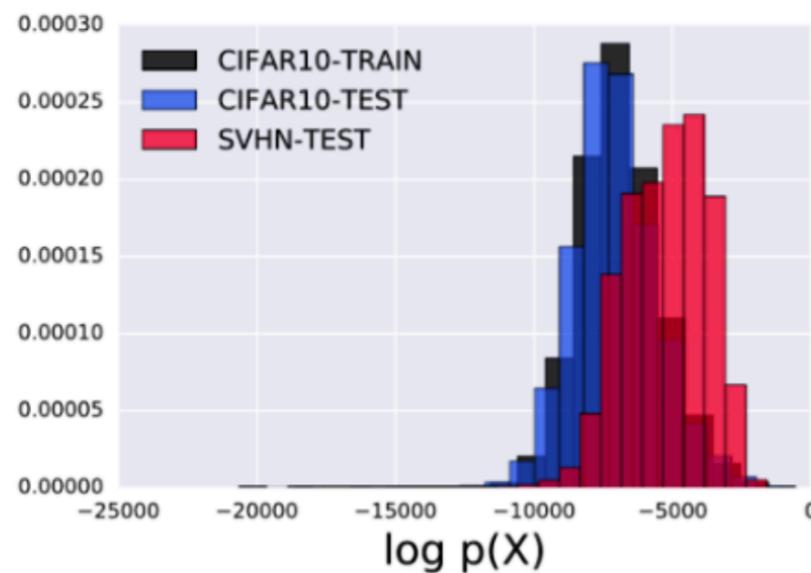
Bits Per Dimension
(NLL / # dims / log 2)

	Bits Per Dimension (NLL / # dims / log 2)
CIFAR10-Train	3.386
CIFAR10-Test	3.464
SVHN-Test	2.389

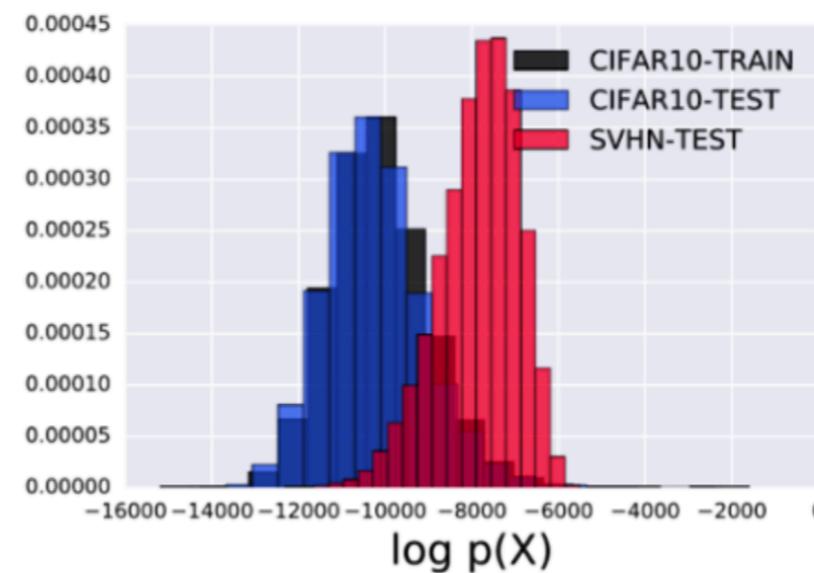
(Lower is Better)



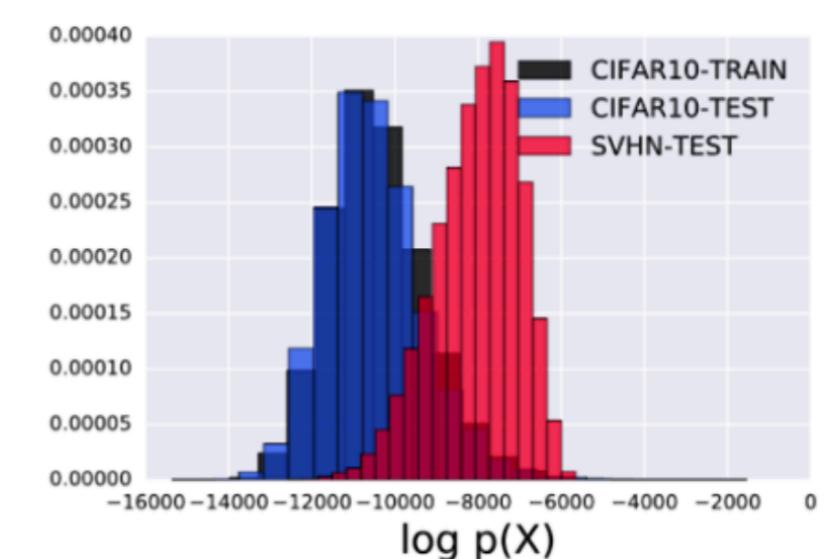
Воспроизводимость на других моделях



(a) PixelCNN

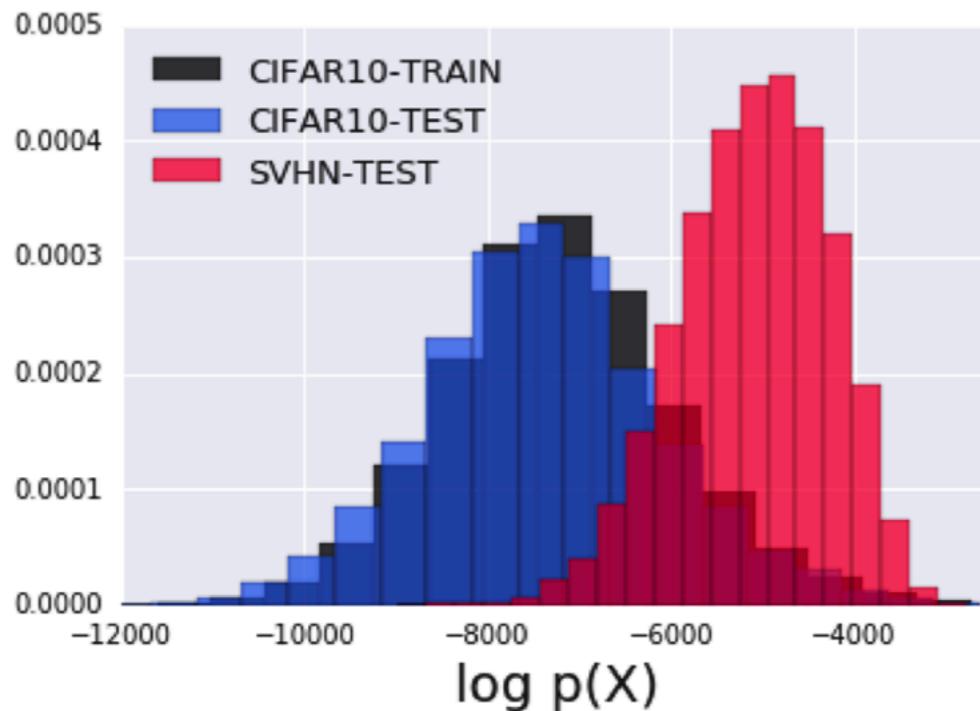


(b) VAE with RNVP as encoder

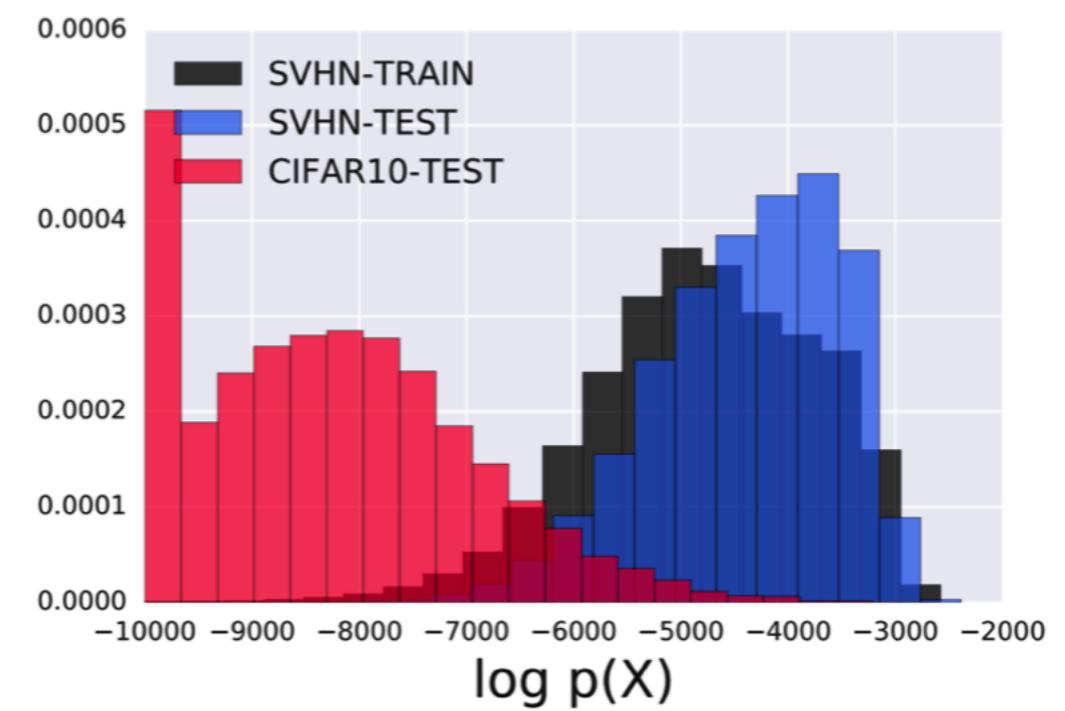


(c) VAE conv-categorical likelihood

Асимметричность по данным в обучении

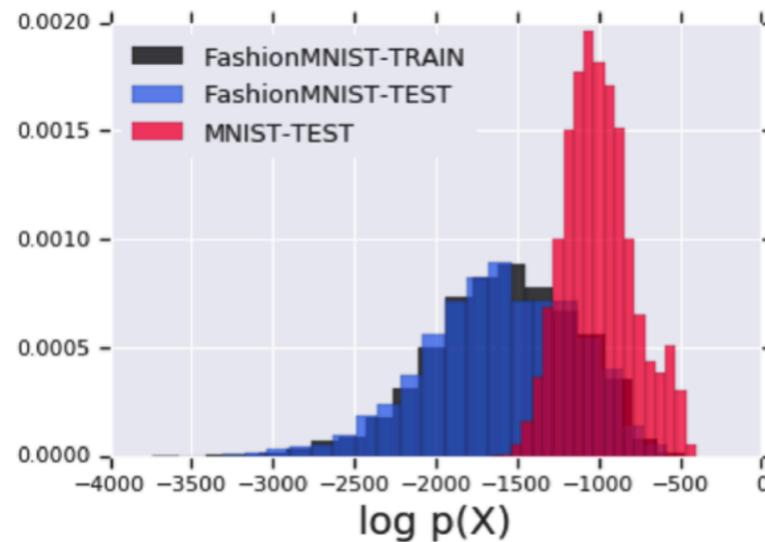


CIFAR-10 vs SVHN

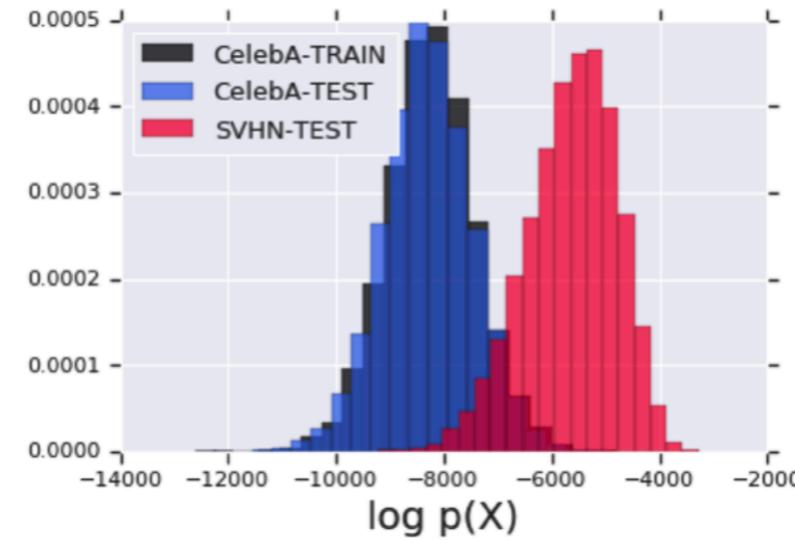


SVHN vs CIFAR-10

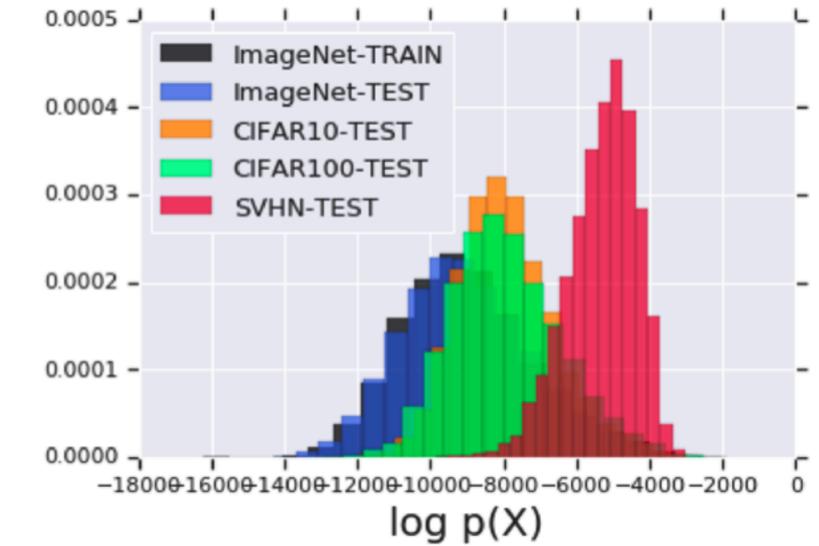
Другие пары наборов данных



FashionMNIST vs MNIST



CelebA vs SVHN

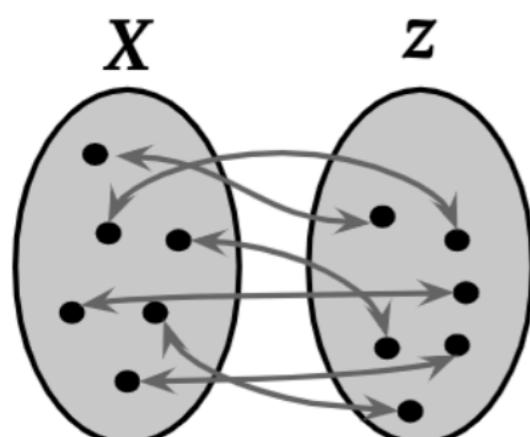


**ImageNet vs CIFAR-10
vs SVHN**

Коротко про потоки

Define Z by a transformation of another variable X : $Z = f(X)$

$f(x)$ is a *bijection*
(invertible 1:1 mapping)



$$x = f^{-1}(z) \quad z = f(x)$$

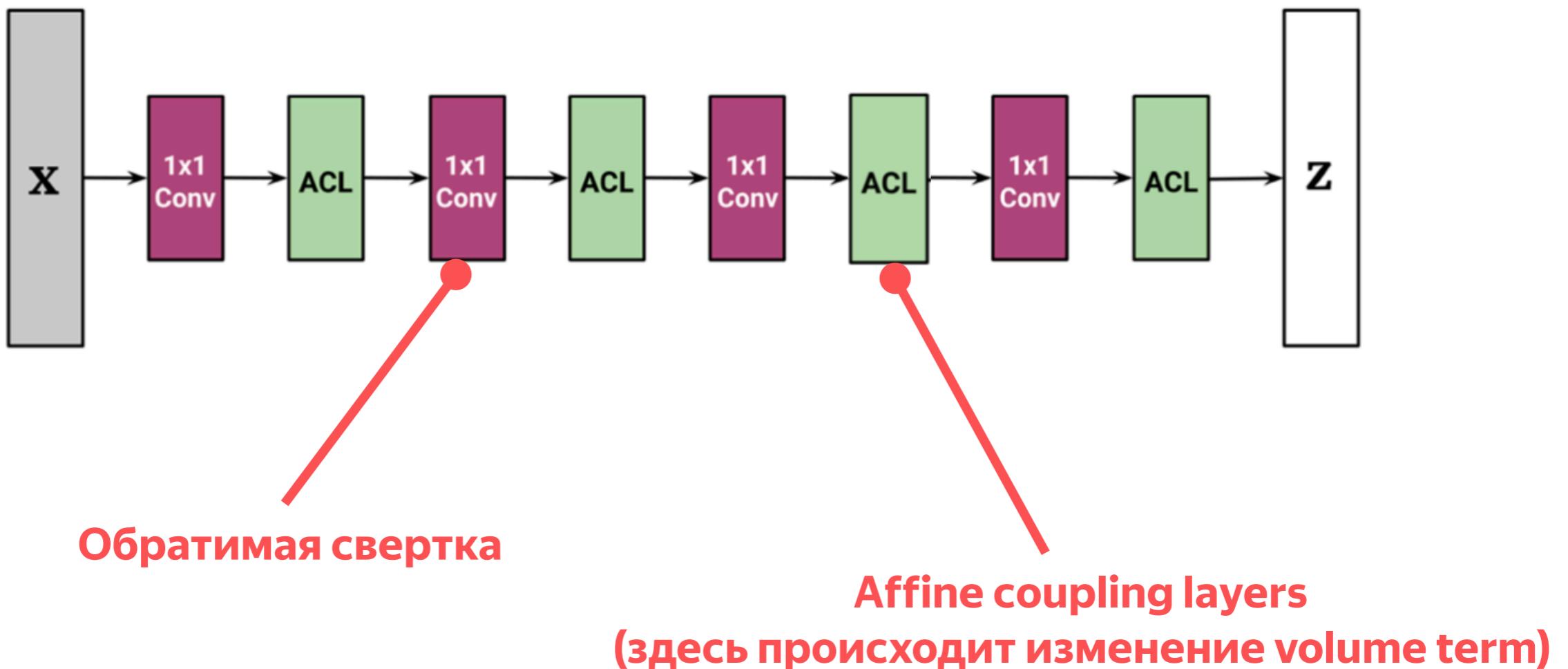
Change of Variables Formula ($X \rightarrow Z$):

$$p_z(f(X)) \left| \frac{df(X)}{dX} \right| = p(X)$$

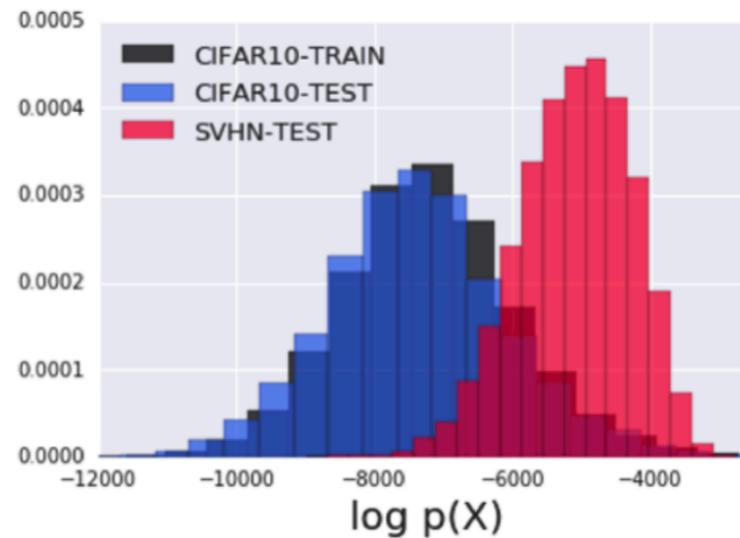
Use simple p_z distribution
(e.g. standard normal)
(Distribution term)

Use f such that the
Jacobian df/dx is easy to
compute
(Volume term)

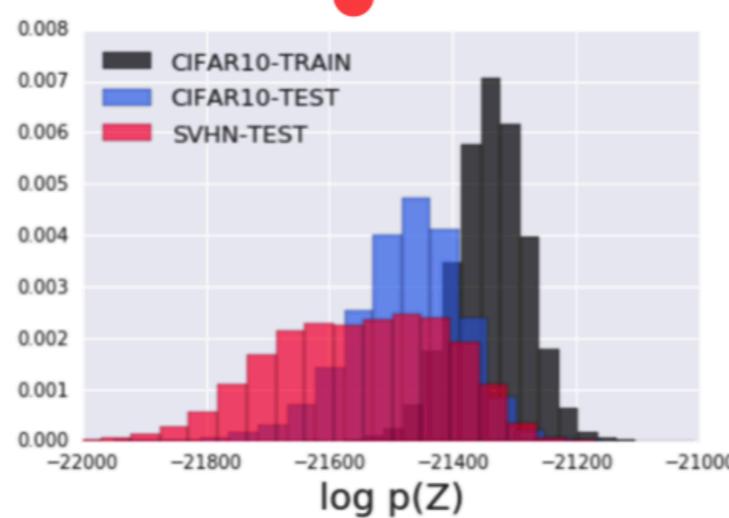
Коротко про потоки



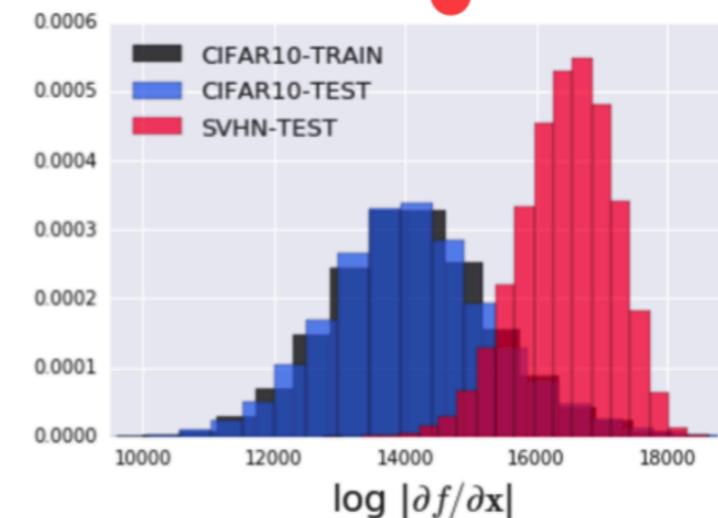
Разложение правдоподобия для flow-based моделей



CIFAR-10 vs SVHN



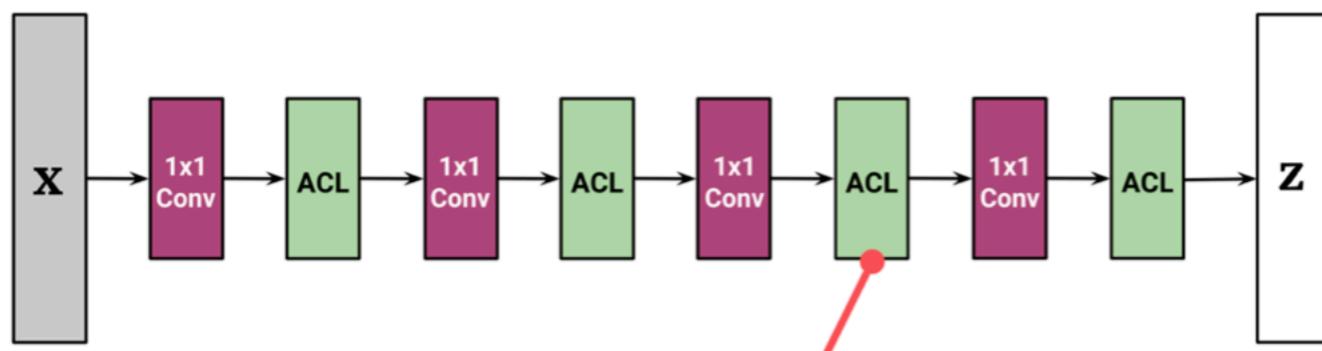
Distribution Term



Volume Term

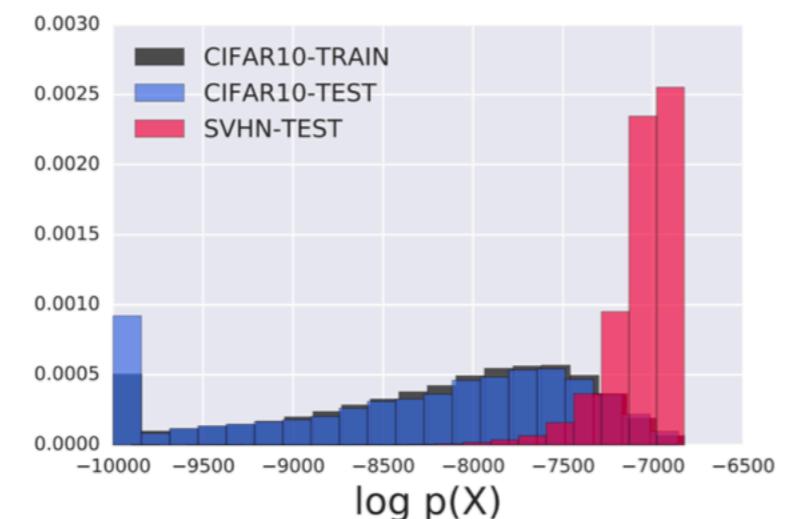
Constant-Volume GLOW

We define a sub-class we term *constant-volume* (w.r.t. input) flows.



Use only translation operations.

CIFAR-10 vs SVHN



Анализ CV-GLOW

Mathematical characterization:

$$0 < \underbrace{\mathbb{E}_q[\log p(x; \theta)]}_{\text{Non-Training Distribution}} - \underbrace{\mathbb{E}_{p^*}[\log p(x; \theta)]}_{\text{Training Distribution}}$$

Анализ CV-GLOW

Mathematical characterization:

$$0 < \underbrace{\mathbb{E}_q[\log p(x; \theta)]}_{\text{Non-Training Distribution}} - \underbrace{\mathbb{E}_{p^*}[\log p(x; \theta)]}_{\text{Training Distribution}}$$

$$\log p(x; \theta) \approx \log p(x_0; \theta)$$

$$\begin{aligned} & + \nabla_{x_0} \log p(x_0; \theta)^T (x - x_0) \\ & + \frac{1}{2} \operatorname{Tr} \left\{ \nabla^2_{x_0} \log p(x_0; \theta) (x - x_0)(x - x_0)^T \right\} \end{aligned}$$

Анализ CV-GLOW

Mathematical characterization:

$$0 < \underbrace{\mathbb{E}_q[\log p(\mathbf{x}; \boldsymbol{\theta})]}_{\text{Non-Training Distribution}} - \underbrace{\mathbb{E}_{p^*}[\log p(\mathbf{x}; \boldsymbol{\theta})]}_{\text{Training Distribution}}$$

$$\approx \nabla_{x_0} \log p(x_0; \theta)^T (\mathbb{E}q[x] - \mathbb{E}p_*[x]) \quad \mid \text{Zero for equal means}$$

$$+ \frac{1}{2} Tr \left\{ \nabla^2_{x_0} \log p(x_0; \theta) \left(\sum q - \sum p_* \right) \right\}$$

Анализ CV-GLOW

Mathematical characterization:

$$0 < \mathbb{E}_{\underline{q}}[\log p(\mathbf{x}; \boldsymbol{\theta})] - \mathbb{E}_{\underline{p^*}}[\log p(\mathbf{x}; \boldsymbol{\theta})]$$

Non-Training Distribution Training Distribution Second Moment of Training Distribution

$$\approx \frac{1}{2} \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p_z(f(\mathbf{x}_0; \boldsymbol{\phi})) + \nabla_{\mathbf{x}_0}^2 \log \left| \frac{\partial f_{\boldsymbol{\phi}}}{\partial \mathbf{x}_0} \right| \right] (\underline{\Sigma_q} - \underline{\Sigma_{p^*}}) \right\}$$

Change-of-Variable Terms Second Moment of Non-Training Distribution

Анализ CV-GLOW

Mathematical characterization:

$$0 < \mathbb{E}_{\underline{q}}[\log p(\mathbf{x}; \boldsymbol{\theta})] - \mathbb{E}_{\underline{p^*}}[\log p(\mathbf{x}; \boldsymbol{\theta})]$$

Non-Training Distribution Training Distribution Second Moment of Training Distribution

$$\approx \frac{1}{2} \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p_z(f(\mathbf{x}_0; \boldsymbol{\phi})) + \nabla_{\mathbf{x}_0}^2 \log \left| \frac{\partial f_{\boldsymbol{\phi}}}{\partial \mathbf{x}_0} \right| \right] (\underline{\Sigma_q} - \underline{\Sigma_{p^*}}) \right\}$$

Change-of-Variable Terms Second Moment of Non-Training Distribution

The diagram illustrates the mathematical characterization of CV-GLOW. It starts with the inequality $0 < \mathbb{E}_{\underline{q}}[\log p(\mathbf{x}; \boldsymbol{\theta})] - \mathbb{E}_{\underline{p^*}}[\log p(\mathbf{x}; \boldsymbol{\theta})]$. This is approximated by a trace term involving the second derivatives of the log probability density function of the latent variable p_z and the change-of-variable Jacobian determinant. A red X is drawn over the term involving the second derivative of the log Jacobian. Below the equation, a horizontal double-headed arrow labeled "Change-of-Variable Terms" spans the two terms in the trace expression. To the left of the first term is "Non-Training Distribution" and to the right of the second term is "Training Distribution". To the right of the entire expression is "Second Moment of Training Distribution". Below the entire expression is "Second Moment of Non-Training Distribution".

Анализ CV-GLOW

Plugging in the CV-Glow transform:

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

1x1 Conv. Params

Second Moment
of Non-Training
Distribution

Second Moment
of Training
Distribution

Анализ CV-GLOW

Plugging in the CV-Glow transform:

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

1x1 Conv. Params

Second Moment
of Non-Training
Distribution

Second Moment
of Training
Distribution

Анализ CV-GLOW

Plugging in the CV-Glow transform:

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

Diagram annotations:

- Sums over channel dimensions**: Points to the summation over j in the equation.
- Product over steps in flow**: Points to the summation over k in the equation.
- 1x1 Conv. Params**: Points to the product term $u_{k,c,j}$.
- Second Moment of Non-Training Distribution**: Points to the term $\overline{\sigma_{q,h,w,c}^2}$.
- Second Moment of Training Distribution**: Points to the term $\overline{\sigma_{p^*,h,w,c}^2}$.
- Sum over spatial dimensions**: Points to the summation over h, w in the equation.

Анализ CV-GLOW

Plugging in the CV-Glow transform:

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

Second Moment
of Non-Training
Distribution

Second Moment
of Training
Distribution

$\overline{\sigma_{q,h,w,c}^2}$

$\overline{\sigma_{p^*,h,w,c}^2}$

$\frac{\partial^2}{\partial z^2} \log p(z; \psi) < 0$ for all log-concave densities (e.g. Gaussian)

Non-negative due to square

Анализ CV-GLOW

Plugging in the CV-Glow transform:

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

Second Moment of Non-Training Distribution $\sum_{h,w} (\bar{\sigma}_{q,h,w,c}^2 - \bar{\sigma}_{p^*,h,w,c}^2)$

Second Moment of Training Distribution $\sum_{h,w} (\underline{\sigma}_{q,h,w,c}^2 - \underline{\sigma}_{p^*,h,w,c}^2)$

< 0 for all log-concave densities (e.g. Gaussian)

Non-negative due to square

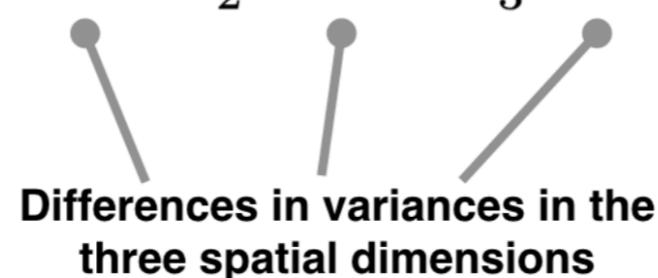
Sign boils down to difference in moments.
Speaks to asymmetric behavior.

Анализ CV-GLOW

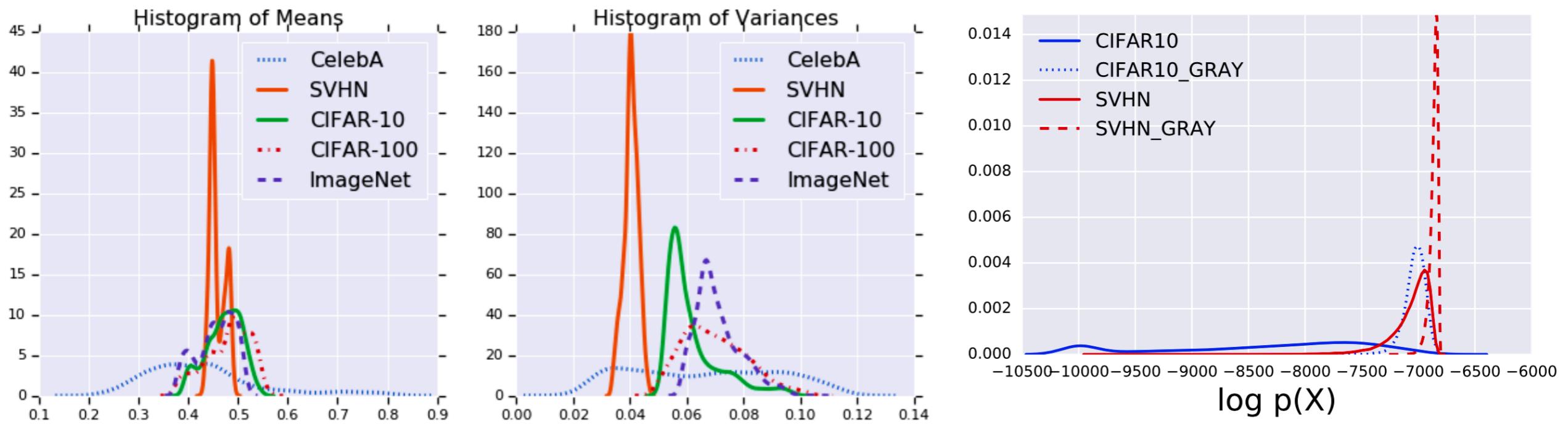
Plugging in the CIFAR-10 and SVHN statistics:

$$\mathbb{E}_{\text{SVHN}}[\log p(\mathbf{x}; \boldsymbol{\theta})] - \mathbb{E}_{\text{CIFAR10}}[\log p(\mathbf{x}; \boldsymbol{\theta})]$$

$$\approx \frac{1}{2\sigma_{\psi}^2} [\alpha_1^2 \cdot 12.3 + \alpha_2^2 \cdot 6.5 + \alpha_3^2 \cdot 14.5] \geq 0 \quad \text{where } \alpha_c = \prod_{k=1}^K \sum_{j=1}^C u_{k,c,j}$$



Анализ CV-GLOW



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

(b) Graying images increases likelihood.

Выводы

- У существующих генеративных моделей есть проблемы с выявлением OOD
- Для flow-based моделей это объясняется разницей в дисперсии распределений
- Нужно быть осторожным при использовании генеративных моделей для выявления аномалий

Вопросы

- Почему нужно прибавлять случайный шум к данным при обучении генеративных моделей? Что такое BPD?
- Какой эффект обнаружили авторы при сравнении правдоподобия данных из обучения и из другого домена? Ожидаемый ли он? Почему?
- Выписать неравенство, которым объясняется феномен для Constant-Volume GLOW. Что значит каждый член этого неравенства?