

# Машинное обучение и нейросетевые модели

## *Лекция 11. Вариационный автокодировщик*

Лектор: Кравченя Павел Дмитриевич

Волгоград 2025

---

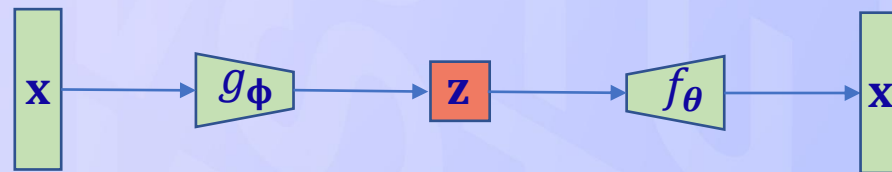
## План лекции

1. Понятие автокодировщиков. Типы автокодировщиков, их свойства.
  2. Вариационный автокодировщик, его архитектура и особенности.
  3. Семейство распределений для VAE. Функция ошибки VAE. ELBO для VAE.
  4. Явный вид функции ошибки VAE для нормального распределения данных.
  5. Процесс обучения VAE и генерация новых объектов. Примеры VAE.
  6. Распутанность представлений VAE. Модель  $\beta$ -VAE.
  7. Коллапс апостериорного распределения в моделях VAE.
  8. Структура латентного пространства автокодировщиков.
  9. Условный вариационный автокодировщик, его обучение и генерация.
  10. Особенности изображений, генерируемых VAE. Модель VAE-GAN.
-

- Автокодировщик представляет собой нейронную сеть, которая на выходе формирует сигнал, эквивалентный входному.
- Он проектируется таким образом, чтобы в процессе преобразований сигнала получить из него дополнительную информацию, например, сжатое представление сигнала в пространстве более *низкой размерности*.
- Автокодировщик состоит из двух частей:
  - Кодировщик (encoder) – преобразует входной сигнал  $\mathbf{x} \in \mathcal{X}$  в латентный вектор  $\mathbf{z} \in \mathcal{Z}$ :  $\mathbf{z} = g_{\Phi}(\mathbf{x})$ , причём, как правило,  $|\mathcal{Z}| \ll |\mathcal{X}|$ .
  - Декодировщик (decoder) – восстанавливает первоначальный сигнал  $\mathbf{x}' \in \mathcal{X}$ ,  $\mathbf{x}' \approx \mathbf{x}$ :  $\mathbf{x}' = f_{\Theta}(\mathbf{z})$ , минимизируя при этом ошибку восстановления, например:  $\mathcal{L}_{AE}(\Theta, \Phi) = \text{MSE}(\mathbf{x}, \mathbf{x}')$ .

- Автокодировщики подразделяются на несколько основных типов:

➤ Понижающие (undercomplete).



➤ Повышающие (overcomplete).

➤ Шумоподавляющие (denoising):  $\mathcal{L}_{\text{DAE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{m} \sum_{i=1}^m \left[ \mathbf{x}_i - f_{\boldsymbol{\theta}} \left( g_{\boldsymbol{\phi}}(\tilde{\mathbf{x}}_i) \right) \right]^2.$

➤ Разреженные (sparse):

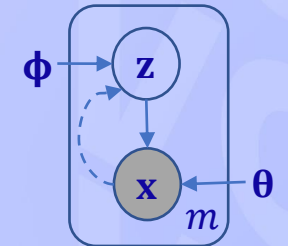
$$\mathcal{L}_{\text{SAE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) + \beta \sum_{l=1}^L \sum_{j=1}^{s_l} \text{KL} \left( \mathcal{B}_{\rho} \| \mathcal{B}_{\hat{\rho}_j^{(l)}} \right), \quad \hat{\rho}_j^{(l)} = \frac{1}{m} \sum_{i=1}^m \left[ a_j^{(l)}(\mathbf{x}_i) \right].$$

➤ Сжимающие (contrastive):

$$\mathcal{L}_{\text{SAE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) + \beta \|J_g(\mathbf{x})\|_F^2 = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) + \beta \sum_{i,j} \left[ \frac{\partial \left( g_{\boldsymbol{\phi}}(\mathbf{x}) \right)_j}{\partial x_i} \right].$$



- Вариационный автокодировщик (variational autoencoder, VAE) является представителем глубоких порождающих моделей, основанный на оценке латентных переменных с помощью вариационного инференса.
- Принцип работы VAE заключается в следующем:
  - Каждому объекту (изображению)  $\mathbf{x}_i \in \mathcal{X}$  соответствует локальная латентная переменная  $\mathbf{z}_i \in \mathcal{Z}$ .
  - $\mathbf{z} \sim q(\mathbf{z})$  отображается на  $\mathbf{x}' \in \mathcal{X}$  с помощью декодировщика:  $\mathbf{x}' = f_{\theta}(\mathbf{z})$ .
  - Переменные  $\mathbf{z}_i \in \mathcal{Z}$  обучаются с использованием амортизированного вариационного инференса, в котором кодировщик играет роль *inference network*:  $q(\mathbf{z} | \mathbf{x}) = q_{\Phi}(\mathbf{z} | g_{\Phi}(\mathbf{x}))$ .
  - Обучение VAE сводится к оптимизации:  $\theta^*, \Phi^* = \arg \max_{\theta, \Phi} \text{ELBO}(q(\mathbf{z} | \mathbf{x}))$ .



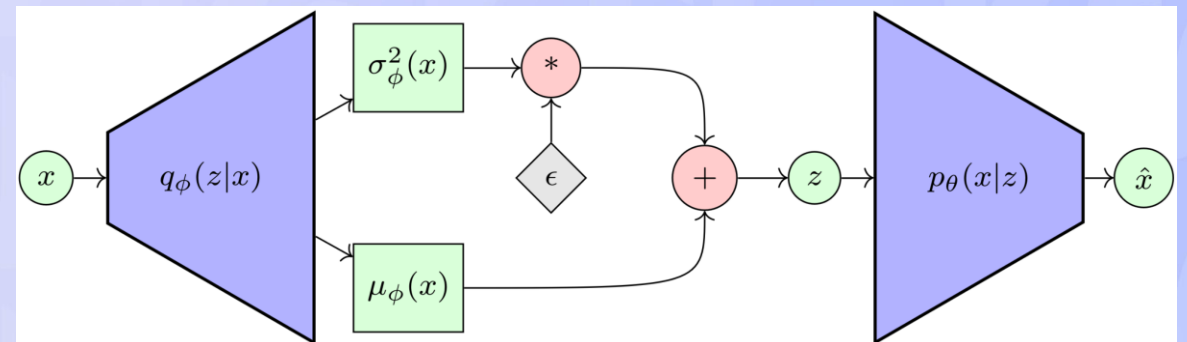
- Вариационное семейство для VAE определяется следующим образом:

$$q_{\phi}(\mathbf{z} | \mathbf{x}) = \mathcal{N} \left( \mathbf{h}_{\phi}^{(1)}(\mathbf{x}), \text{diag} \left( \exp \left( \mathbf{h}_{\phi}^{(2)}(\mathbf{x}) \right) \right) \right).$$

- Параметры вариационного распределения  $\boldsymbol{\mu} = \mathbf{h}_{\phi}^{(1)}(\mathbf{x})$  и  $\log \boldsymbol{\sigma}^2 = \mathbf{h}_{\phi}^{(2)}(\mathbf{x})$  определяются нейросетями, которые, как правило, являются составными частями кодировщика:  $(\mathbf{h}_{\phi}^{(1)}, \mathbf{h}_{\phi}^{(2)}) = g_{\phi}(\mathbf{x})$ .

- Для оценки градиента ELBO к нормальному вариационному распределению применяется трюк репараметризации:

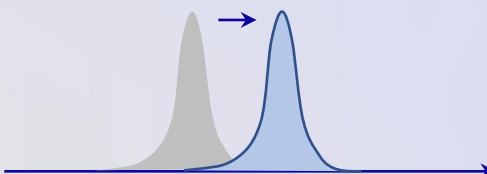
$$\mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu},$$
$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$



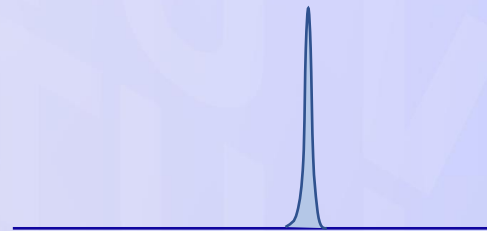
- Для VAE существует модификация, позволяющая снизить дисперсию эстиматора градиента ELBO. Перепишем ELBO следующим образом:

$$\begin{aligned} \text{ELBO}(\boldsymbol{\phi}, \boldsymbol{\theta}) &= \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i)] = \\ &= \sum_{i=1}^m \int q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i) [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i)] d\mathbf{z}_i = \\ &= \sum_{i=1}^m \int q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i) [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i | \mathbf{z}_i) + \log p(\mathbf{z}_i) - \log q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i)] d\mathbf{z}_i = \\ &= \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i | \mathbf{z}_i)] + \sum_{i=1}^m \int q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i) [\log p(\mathbf{z}_i) - \log q_{\boldsymbol{\phi}}(\mathbf{z}_i | \mathbf{x}_i)] d\mathbf{z}_i = \end{aligned}$$

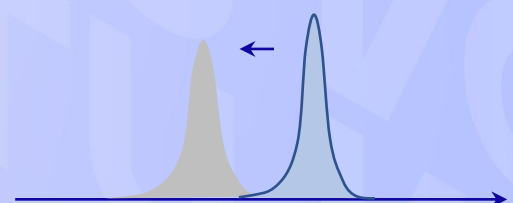
$$\begin{aligned} &= \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\Theta}(\mathbf{x}_i | \mathbf{z}_i)] + \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} \left[ \log \frac{p(\mathbf{z}_i)}{q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} \right] = \\ &= \sum_{i=1}^m \underbrace{\left\{ \mathbb{E}_{\mathbf{z}_i \sim q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\Theta}(\mathbf{x}_i | \mathbf{z}_i)] \right\}}_{\text{reconstruction loss}} - \underbrace{\text{KL}(q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) || p(\mathbf{z}_i))}_{\text{regularization term}} = \end{aligned}$$



Вариационное распределение  
стремится отличаться от  
априорного, чтобы описать  
особенности в данных



Если нет регуляризации,  
вариационное распределение  
принимает узкую форму



Регуляризация ищет баланс,  
не позволяя вариационному  
распределению сильно  
отклониться от априорного



- В общем случае, вариационный автокодировщик описывает параметры произвольных распределений.
- Но если априорное распределение латентных переменных  $p(\mathbf{z})$  считать стандартным нормальным  $\mathcal{N}(\mathbf{0}, \mathbb{I})$ , то дивергенция  $\text{KL}(q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) || p(\mathbf{z}_i))$  может быть вычислена аналитически.
- Поскольку, часть эстиматора вычисляется аналитически, то отпадает необходимость в её оценке Монте-Карло, а это приводит к уменьшению дисперсии эстиматора.
- Если наблюдаемая переменная непрерывна и с достаточной степенью точности подчиняется нормальному распределению, то и ошибка реконструкции тоже может быть преобразована для упрощения расчетов.

- Пусть  $p_{\theta}(\mathbf{x}_i | \mathbf{z}_i) \sim \mathcal{N}(f_{\theta}(\mathbf{z}_i), \sigma_{\text{dc}}\mathbb{I})$ ,  $\sigma_{\text{dc}}$  – некоторый гиперпараметр. Тогда:

$$\begin{aligned}\text{ELBO}(\Phi, \Theta) &= \sum_{i=1}^m \left\{ \mathbb{E}_{\mathbf{z}_i \sim q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\Theta}(\mathbf{x}_i | \mathbf{z}_i)] - \text{KL}(q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) \parallel p(\mathbf{z}_i)) \right\} = \\ &= \sum_{i=1}^m \left\{ \mathbb{E}_{\mathbf{z}_i \sim q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} \left[ \log \frac{1}{\sqrt{2\pi}\sigma_{\text{dc}}} - \frac{\|\mathbf{x}_i - f_{\Theta}(\mathbf{z}_i)\|_2^2}{2\sigma_{\text{dc}}^2} \right] - \text{KL}(q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) \parallel p(\mathbf{z}_i)) \right\}.\end{aligned}$$

- Таким образом, модель пытается *восстановить входной сигнал*, выполняя минимизацию  $\text{MSE}(\mathbf{x}, \mathbf{x}')$  (для задач регрессии).
- При этом модель также пытается поддерживать одинаковыми распределения  $q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)$  и  $p(\mathbf{z}_i)$ .
- Если  $p(\mathbf{z}) \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$ , то и вариационное распределение  $q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) \rightarrow \mathcal{N}(\mathbf{0}, \mathbb{I})$ .

- Поскольку KL-дивергенцию можно рассчитать вручную, то для расчета функции ошибки остаётся вычислить матожидание  $\log p(\mathbf{x} | \mathbf{z})$ , которое можно оценить с помощью эстиматора Монте-Карло.

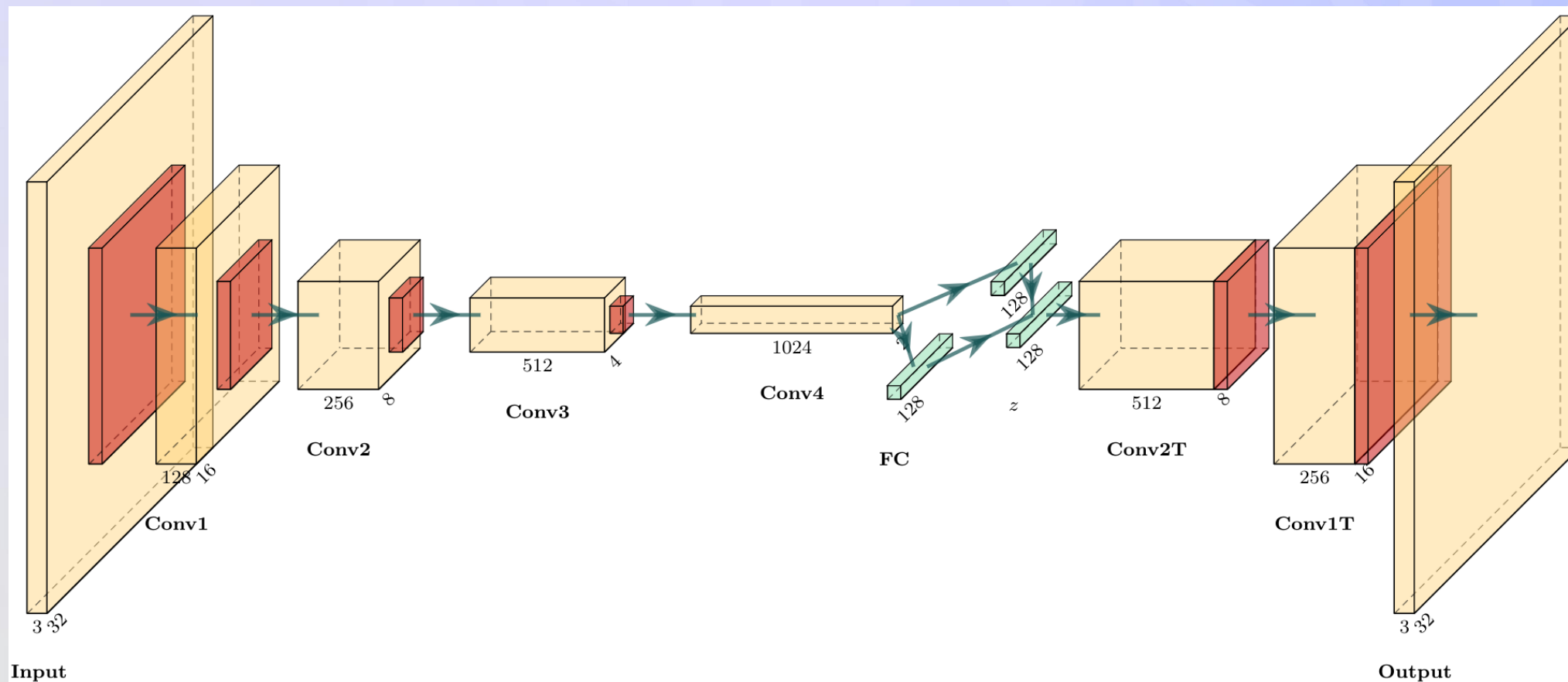
- В таком случае, функция ошибки примет вид:

$$\mathcal{L}_{\theta, \phi}(\mathbf{X}) = \frac{1}{m} \sum_{i=1}^m \left[ \frac{1}{K} \sum_{j=1}^K \log p_{\theta}(\mathbf{x}_i | \mathbf{z}_j) + \frac{1}{2} \sum_{k=1}^M \left( 1 + \ln \sigma_k^2(\mathbf{x}_i) - \mu_k^2(\mathbf{x}_i) - \sigma_k^2(\mathbf{x}_i) \right) \right].$$

Здесь  $M$  – размерность векторов  $\boldsymbol{\mu}$  и  $\boldsymbol{\sigma}$ ,  $(\boldsymbol{\mu}, \boldsymbol{\sigma}) = g_{\phi}(\mathbf{X}, \epsilon)$ .

- При достаточно большом размере мини-батча количество семплов можно взять небольшим (например, один семпл на 100 элементов в батче).
- Поскольку применен трюк репараметризации, можно считать градиент.

Обучение VAE	Генерация новых объектов с помощью VAE
<p>Пусть дан датасет <math>\mathcal{X} = \{x_1, x_2, \dots, x_n\}</math>.</p> <p>Инициализировать <math>\theta</math> и <math>\phi</math>.</p> <p><b>while</b> not converge:</p> <ul style="list-style-type: none"> <li>➤ Сформировать случайный минибатч: <math>X^m = \{x_1, x_2, \dots, x_m\} \subset \mathcal{X}</math>;</li> <li>➤ Сформировать <math>\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I})</math>;</li> <li>➤ Вычислить функцию ошибки:</li> </ul> $\mathcal{L}_{\theta, \phi}(X) = \frac{1}{m} \sum_{i=1}^m \left[ \frac{1}{K} \sum_{j=1}^K \log p_{\theta}(\mathbf{x}_i   \mathbf{z}_j) + \frac{1}{2} \sum_{k=1}^M \left( 1 + \ln \sigma_k^2(\mathbf{x}_i) - \mu_k^2(\mathbf{x}_i) - \sigma_k^2(\mathbf{x}_i) \right) \right]$ <ul style="list-style-type: none"> <li>➤ Выполнить <i>backprop()</i>;</li> <li>➤ Вычислить <math>\theta, \phi \leftarrow \theta, \phi - \nabla_{\theta, \phi} \mathcal{L}_{\theta, \phi}(X)</math>.</li> </ul>	<p>Генерация <u>новых объектов</u> выполняется с помощью декодировщика:</p> $\mathbf{x} = f_{\theta^*}(\mathbf{z}), \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbb{I}).$ <p>Кодировщик для генерации <u>новых объектов не требуется</u>.</p> <p>Однако, иногда требуется оценить вероятность генерации моделью определенного объекта:</p> $p(\mathbf{x}) = \int p(\mathbf{x}   \mathbf{z}) p(\mathbf{z}) d\mathbf{z}.$ <p>Для <u>оценки интеграла</u> требуется получить семплы <math>\mathbf{z}_i</math>. Эффективнее это сделать из распределения <math>q(\mathbf{z}   \mathbf{x})</math>, задаваемого кодировщиком, чем из <math>\mathcal{N}(\mathbf{0}, \mathbb{I})</math>.</p>





- Если каждая латентная переменная чувствительна только к одному порождающему фактору и относительно инвариантна к другим, то говорят, что представление распутано (*disentangled*).
- Одним из преимуществ *распутанного представления* является хорошая интерпретируемость и простота обобщения. Например, модель, обученная на фотографиях человеческих лиц, может определять цвет кожи, цвет и длину волос, эмоции, и т.д.
- Модель  $\beta$ -VAE представляет собой модификацию вариационного автокодировщика, в которой особое внимание уделяется обнаружению распутанных скрытых факторов.

- В  $\beta$ -VAE требуется *максимизировать* вероятность генерации реальных данных, сохраняя при этом *небольшое расстояние* между априорным и вариационным распределениями латентных переменных:

$$\max_{\theta, \phi} \mathbb{E}_{\mathbf{x} \sim \mathcal{X}} \left[ \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})} \log p_{\theta}(\mathbf{x} | \mathbf{z}) \right], \quad \text{при условии } \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z})) < \epsilon.$$

- Данное выражение может быть записано в виде Лагранжиана с множителями Лагранжа  $\beta$ :

$$\begin{aligned} \mathcal{F}(\theta, \phi, \beta) &= \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | \mathbf{z}_i)] - \beta \left[ \text{KL}(q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) || p(\mathbf{z}_i)) - \epsilon \right] = \\ &= \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i | \mathbf{z}_i)] - \beta \text{KL}(q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) || p(\mathbf{z}_i)) + \beta \epsilon \geq \end{aligned}$$

$$\geq \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\Theta}(\mathbf{x}_i | \mathbf{z}_i)] - \beta \text{KL} \left( q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) \parallel p(\mathbf{z}_i) \right), \quad \text{поскольку } \beta, \epsilon \geq 0.$$

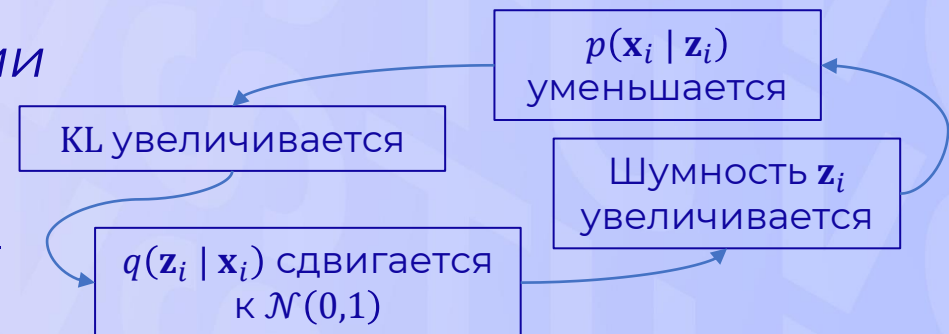
- Теперь можно ввести функцию ошибки для  $\beta$ -VAE:

$$\mathcal{L}_{\beta VAE} = - \sum_{i=1}^m \mathbb{E}_{\mathbf{z}_i \sim q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\Theta}(\mathbf{x}_i | \mathbf{z}_i)] + \beta \text{KL} \left( q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) \parallel p(\mathbf{z}_i) \right).$$

- Отрицательное значение функции ошибки является нижней границей лагранжиана. Минимизация функции ошибки приводит к максимизации лагранжиана и решению оптимизационной задачи.
- При  $\beta = 1$   $\beta$ -VAE функционирует как вариационный автокодировщик. При  $\beta > 1$  увеличиваются ограничения, которые накладываются на латентные переменные.

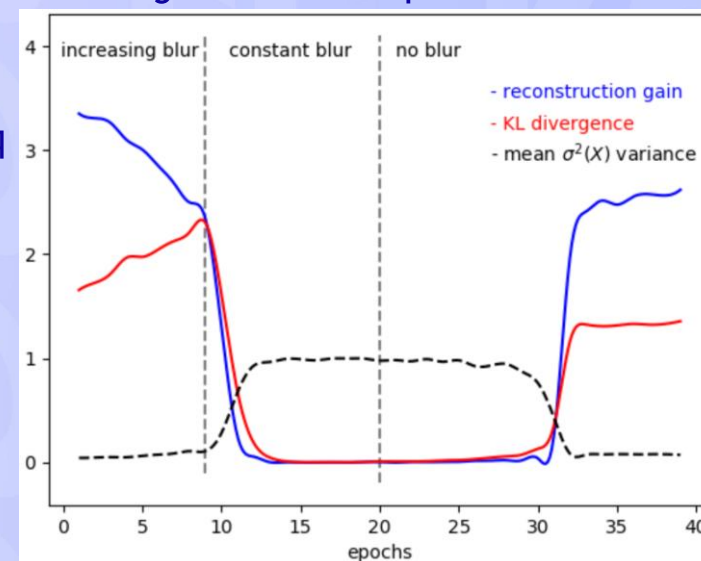
- Введённые условия ограничивают информационную ёмкость для  $z$ , что в сочетании с необходимостью максимизировать логарифмическую вероятность обучающих данных  $x$  должно стимулировать модель к изучению наиболее эффективного представления.
- Модели приходится искать компромисс между качеством реконструкции и KL-дивергенцией.
- Так как *различные* факторы по-разному влияют на потери модели при реконструкции, то ей выгодно выполнить распутывание признаков, поскольку в этом случае она может напрямую ранжировать важность (с помощью KL-дивергенции) для каждого из них.

- Если во время обучения VAE вклад ошибки восстановления для некоторой латентной переменной мал по сравнению со слагаемым Кульбака-Лейблера, то распределение этой переменной из-за будет сдвигаться из-за регуляризации к априорному распределению  $p(\mathbf{z}_i)$ , не учитывающему  $\mathbf{x}$ .
- В результате, латентная переменная станет ещё более шумной, и декодировщик перестанет учитывать её при создании новых образцов.
- А это приведет к постоянной генерации некоторых «усредненных» объектов.
- Данный эффект называется коллапсом апостериорного распределения.



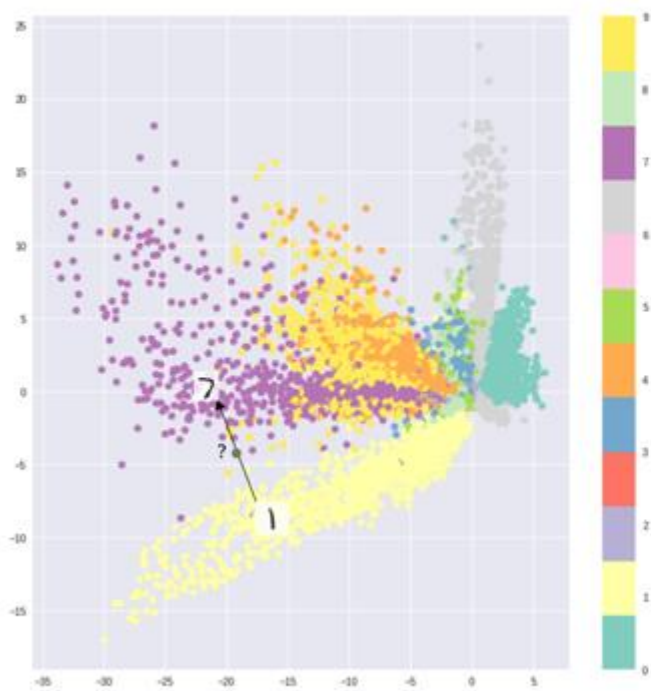


- Если *добавить* к латентной переменной искусственный шум, можно воспроизвести эффект коллапса апостериорного распределения.
- Вклад латентной переменной в реконструкцию можно вычислить как разницу между ошибками реконструкции с учётом и без учета переменной. Он называется reconstruction gain.
- Когда *reconstruction gain* переменной становится меньше, чем *KL-дивергенция*, сама переменная игнорируется сетью. Дисперсия переменной при этом увеличивается.
- Если *искусственный шум удалить*, латентная переменная вновь активируется.

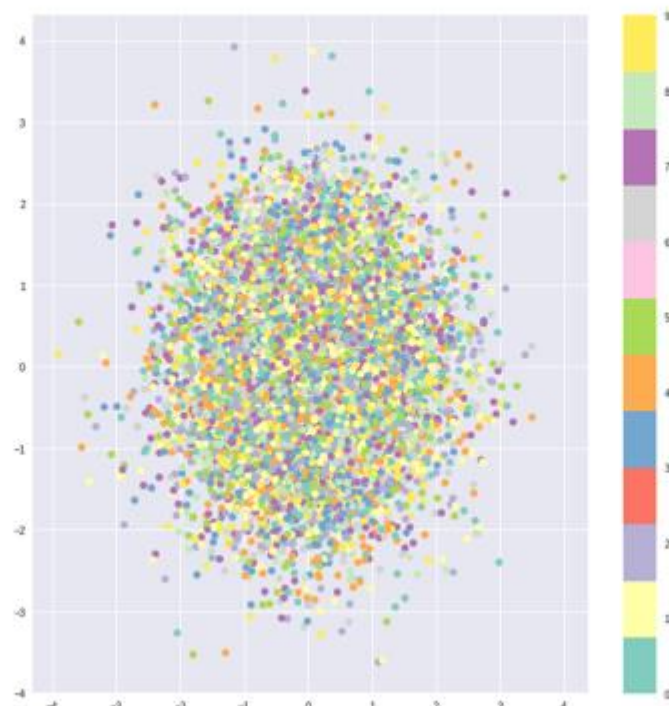


- Стандартные автокодировщики хорошо работают в задачах восстановления данных, однако, они не подходят в качестве генеративной модели, поскольку выбор случайного входного сигнала  $\mathbf{z}'$  для декодера не обязательно приведет к тому, что декодер создаст приемлемое изображение.
- Близкие значения  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ ,  $\mathbf{x}_1 \approx \mathbf{x}_2$ , не обязательно будут отображаться в близкие значения векторов латентного пространства:  $g(\mathbf{x}_1) \neq g(\mathbf{x}_2)$ .
- Близкие значения  $\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{Z}$ ,  $\mathbf{z}_1 \approx \mathbf{z}_2$ , не обязательно будут отображаться в близкие значения векторов данных:  $f(\mathbf{z}_1) \neq f(\mathbf{z}_2)$ . Более того, если  $\mathbf{x} \in \mathcal{X}$ ,  $\mathbf{z} \in \mathcal{Z}$ ,  $\mathbf{x} = f(\mathbf{z})$ , то не обязательно, чтобы при  $\mathbf{z}' \approx \mathbf{z}$  было справедливо:  $f(\mathbf{z}') \in \mathcal{X}$ .

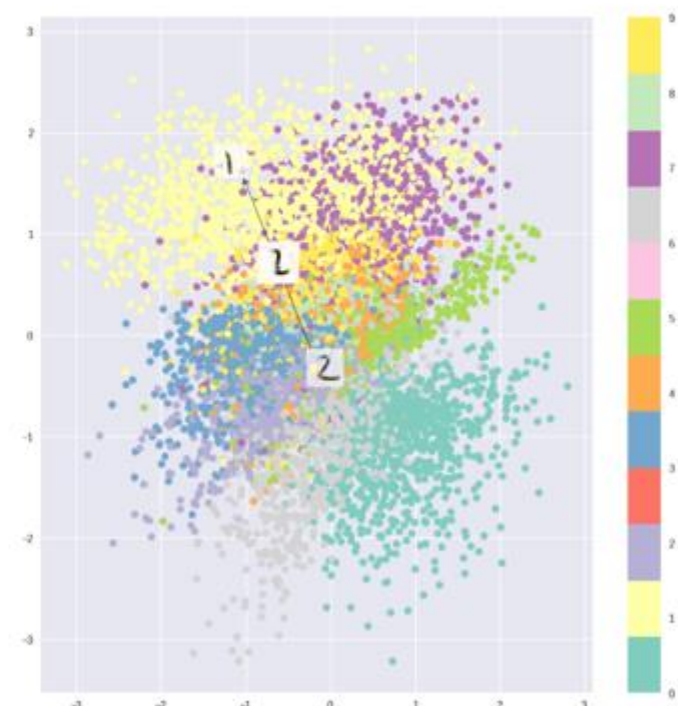
Only reconstruction loss



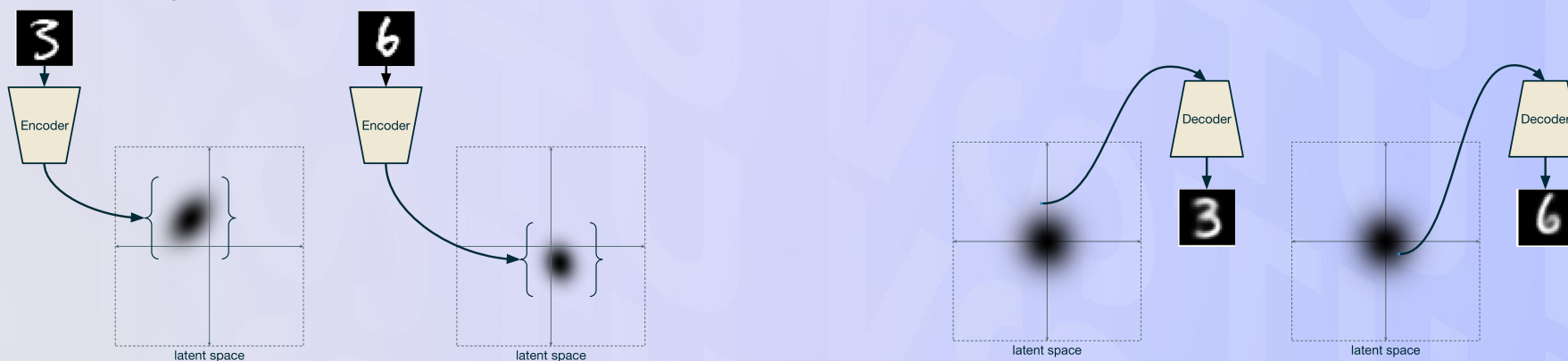
Only KL divergence



Combination



- Нужен какой-то способ *гарантировать*, что *декодер* готов преобразовывать любой *входной* сигнал в *разумное* изображение. Для этого нужно заранее определить распределение входных данных, которые *декодер* должен ожидать. В вариационном автокодировщике с этой целью обычно используется стандартное нормальное распределение  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .
- Онлайн-ресурс с VAE: <https://www.siares.com/projects/variational-autoencoder>





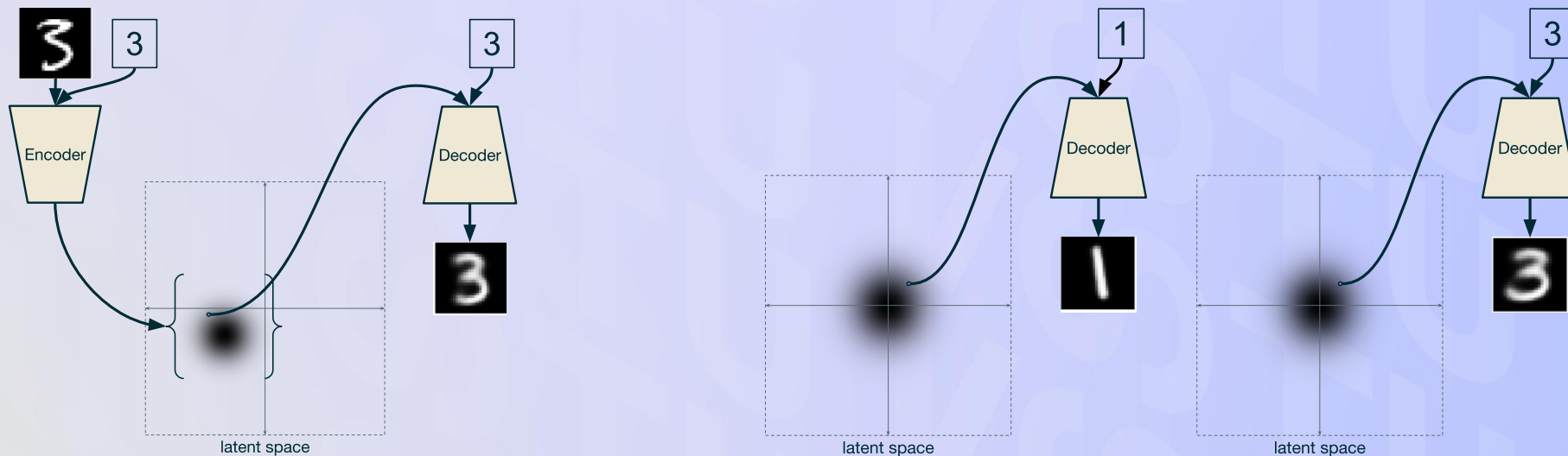
- Часто требуется сгенерировать не просто какой-либо произвольный объект из датасета, а относящийся к конкретной группе или классу.
- Для этого потребуем, чтобы распределения, участвующие при построении вариационного автокодировщика, были условными по некоторой переменной  $y \in \mathcal{Y}$ , характеризующей принадлежность объекта к требуемому классу:  $q_{\phi}(\mathbf{z} | \mathbf{x}, y)$ ,  $p_{\theta}(\mathbf{x} | \mathbf{z}, y)$ .
- В этом случае, функция ошибки для условного автокодировщика:

$$\text{ELBO}(\phi, \theta) = \sum_{i=1}^m \left\{ \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i, y_i)} [\log p_{\theta}(\mathbf{x}_i | \mathbf{z}_i, y_i)] + \text{KL}(q_{\phi}(\mathbf{z}_i | \mathbf{x}_i, y_i) || p(\mathbf{z}_i, y_i)) \right\}.$$

- Переменная  $y \in \mathcal{Y}$  может иметь любую природу (быть, например, меткой).



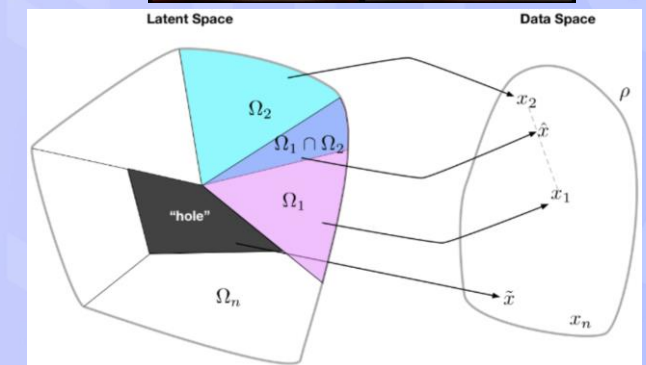
- Обучение CVAE отличается от обучения нескольких независимых VAE, поскольку веса CVAE являются общими для всех классов.
- При реализации CVAE входы кодировщика и декодировщика обычно конкатенируют с тензором  $y$ .



- Изображение *классического VAE* часто размыто (*blurry*).
- Одной из причин этого является неоднозначность отображения входных данных в латентное пространство.
- Для разных векторов данных могут наблюдаться перекрывающиеся латентные переменные.

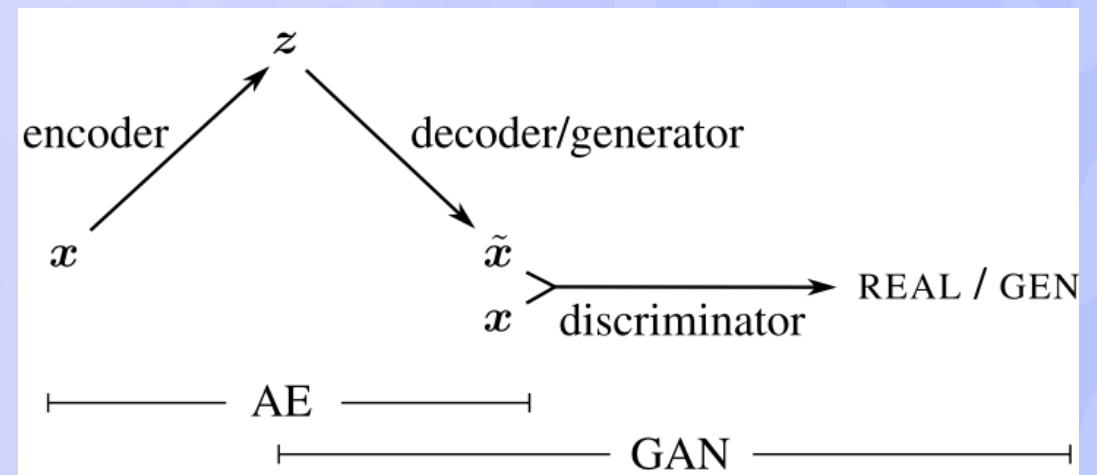
Оптимальная реконструкция для классических VAE без ограничений предполагает расчет среднего между этими векторами, что приводит к размытию.

- Также, в VAE может наблюдаться проблема «дыр»: областях из  $\mathcal{Z}$ , в которых декодер ничем не ограничен, поэтому, может выдавать произвольные выходные данные.

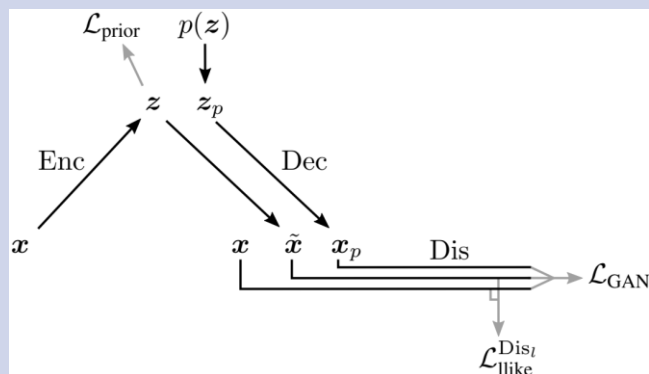


- Еще одной *причиной* появления размытости является использование *попиксельных* ошибок (например, *среднеквадратичной*) VAE при расчете ошибки восстановления.
- Поскольку архитектура VAE предполагает сжатие входного сигнала, она не сохраняет всю информацию о входном изображении. Нюансы изображения часто *теряются*, и VAE будет прогнозировать что-то среднее между изученными изображениями, чтобы минимизировать ошибку восстановления. Это приводит к тому, что изображение становится размытым.
- Есть несколько *подходов*, предполагающих замену функции ошибки на более *подходящую* в конкретных случаях, например, архитектура VAE-GAN.

- Подход строится на том, чтобы обучить вариационный автокодировщик генерировать новые изображения, но требуемую метрику определять автоматически посредством GAN.
- Объединим декодировщик VAE и генератор GAN в единую нейросеть, которая занимается задачей формирования восстановленных объектов.
- Дискриминатор GAN пытается отличить реальный объект от сгенерированного.
- Для оценки схожести объектов требуется ввести функцию ошибки на основе сравнения признаков.



## Архитектура VAE-GAN



$$p(\text{Dis}_l(\mathbf{x})) = \mathcal{N}(\text{Dis}_l(\mathbf{x}) | \text{Dis}_l(\tilde{\mathbf{x}}), \mathbb{I}),$$

$$\mathcal{L}_{\text{llike}}^{\text{Dis}_l} = -\mathbb{E}_{q(\mathbf{z} | \mathbf{x})} [\log p(\text{Dis}_l(\mathbf{x}) | \mathbf{z})],$$

$$\mathcal{L}_{\text{prior}} = \text{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z})),$$

$$\mathcal{L}_{\text{GAN}} = \log(\text{Dis}(\mathbf{x})) + \log(1 - \text{Dis}(\text{Dec}(\mathbf{z}))), \\ + \log(1 - \text{Dis}(\text{Dec}(\text{Enc}(\mathbf{x}))))$$

$$\mathcal{L} = \mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{llike}}^{\text{Dis}_l} + \mathcal{L}_{\text{GAN}}.$$

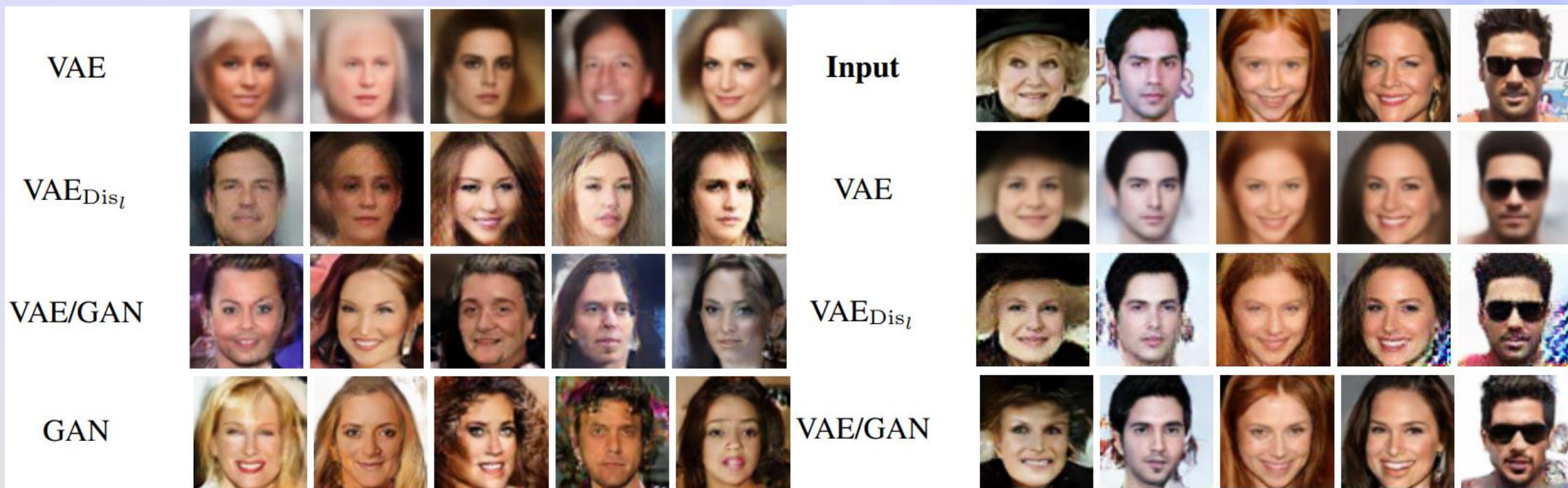
## Процесс обучения VAE-GAN

Инициализировать  $\theta_{\text{Enc}}$ ,  $\theta_{\text{Dec}}$ ,  $\theta_{\text{Dis}}$ .

**while** not converge:

- Сформировать случайный *минибатч*  $\mathbf{X}$ ;
- $\mathbf{Z} \leftarrow \text{Enc}(\mathbf{X})$ ;
- $\tilde{\mathbf{X}} \leftarrow \text{Dec}(\mathbf{Z})$ ;
- $\mathbf{Z}_p \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$ ;
- $\mathbf{X}_p \leftarrow \text{Dec}(\mathbf{Z}_p)$ ;
- Рассчитать  $\mathcal{L}_{\text{GAN}}(\mathbf{X}, \tilde{\mathbf{X}}, \mathbf{X}_p)$ ;
- $\theta_{\text{Enc}} \leftarrow -\nabla_{\theta_{\text{Enc}}} (\mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{llike}}^{\text{Dis}_l})$ ;
- $\theta_{\text{Dec}} \leftarrow -\nabla_{\theta_{\text{Dec}}} (\gamma \mathcal{L}_{\text{llike}}^{\text{Dis}_l} - \mathcal{L}_{\text{GAN}})$ ;
- $\theta_{\text{Dis}} \leftarrow -\nabla_{\theta_{\text{Dis}}} \mathcal{L}_{\text{GAN}}$ ;





Семплирование

Реконструкция

# Демонстрация практических примеров

---

## Заключение

1. Ввели понятие новой глубокой порождающей модели – VAE – и поговорили про её особенности и свойства.
2. Получили функцию ошибки вариационного автокодировщика, разобрались с алгоритмом его обучения и генерации новых образцов.
3. Поговорили про особенности латентного пространства автокодировщиков.
4. Исследовали распутанность представлений VAE и ввели архитектуру  $\beta$ -VAE.
5. Ввели понятие условного вариационного автокодировщика.
6. Поговорили про особенности и проблемы, связанные с обучением VAE, ввели понятие коллапса апостериорного распределения.
7. Теоретически и на практическом примере разобрались с принципами функционирования вариационных автокодировщиков.

# **Спасибо за внимание!**

Волгоград 2025

---