

Генеративные модели

Лекция 5: Векторная квантизация, *GANs*

Повторение

Проблемы вариационного вывода

Цель:

Хотим аппроксимировать сложное апостериорное распределение $p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$ с помощью более простого $q(\mathbf{z})$

Проблема 1:

Пространство всех возможных функций $q(\mathbf{z})$ бесконечно

Проблема 2:

Для датасета из N объектов $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, нам нужно решать N отдельных задач оптимизации, чтобы найти $\{q_1(\mathbf{z}_1), \dots, q_n(\mathbf{z}_n)\}$

Амортизированный вариационный вывод

Решение 1, параметризация:

Ограничим поиск $q(\mathbf{z})$ по конкретным семействам распределений:

$$q(\mathbf{z}) \rightarrow q_{\phi}(\mathbf{z})$$

Решение 2, амортизация (Amortized Variational Inference):

Вместо отдельных параметров ϕ для каждого \mathbf{x} обучим единую **нейросеть–кодировщик**, которая по \mathbf{x} будет предсказывать параметры для q :

$$q_{\phi}(\mathbf{z}) \rightarrow q_{\phi}(\mathbf{z}|\mathbf{x})$$

Амортизированный вариационный вывод

Введя эти 2 ограничения мы:

1. Ушли от необходимости считать вариационное распределение для каждого объекта \mathbf{x}
2. Введя параметры ϕ , мы признаем, что наше решение может быть не точным и KL -дивергенция между $q_\phi(\mathbf{z}|\mathbf{x})$ и $p(\mathbf{z}|\mathbf{x}, \theta)$ может быть больше нуля

Теперь наша задача сводится к нахождению оптимальных ϕ и θ , которые максимизируют $ELBO$

○ Е-шаг:

$$\phi_k = \phi_{k-1} + \eta \nabla_{\phi} \mathcal{L}_{\phi, \theta_{k-1}}(\mathbf{x}) \Big|_{\phi = \phi_{k-1}}$$

○ М-шаг:

$$\theta_k = \theta_{k-1} + \eta \nabla_{\theta} \mathcal{L}_{\phi_k, \theta}(\mathbf{x}) \Big|_{\theta = \theta_{k-1}}$$

Градиент на М-шаге

Считаем градиент по параметрам θ :

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}, \theta)] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

KL -дивергенция не зависит от θ , так как q параметризуется ϕ , а априорное распределение $p(\mathbf{z})$ обычно фиксировано:

$$\nabla_{\theta} \mathcal{L}_{\phi, \theta}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}, \theta)] = \nabla_{\theta} \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}|\mathbf{z}, \theta) d\mathbf{z} = \int q_{\phi}(\mathbf{z}|\mathbf{x}) \nabla_{\theta} \log p(\mathbf{x}|\mathbf{z}, \theta) d\mathbf{z}$$

Для оценки матожидания будем использовать метод Монте-Карло:

$$\nabla_{\theta} \mathcal{L}_{\phi, \theta}(\mathbf{x}) \approx \nabla_{\theta} \log p(\mathbf{x}|\mathbf{z}^*, \theta), \quad \mathbf{z}^* \sim q_{\phi}(\mathbf{z}|\mathbf{x})$$

Поскольку $q_{\phi}(\mathbf{z}|\mathbf{x})$ обусловлено на конкретный \mathbf{x} , то его вероятностная масса будет сосредоточена в области пространства \mathbf{z} , которая наиболее вероятна для этого \mathbf{x} , а значит дисперсия будет сильно ниже, чем при наивном подходе

Reparameterization trick

Считаем градиент по параметрам ϕ при фиксированных θ :

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}, \theta)] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

Будем использовать **трюк репараметризации (reparameterization trick)** – считать матожидание по распределению, которое не зависит от ϕ

Предположим, что мы можем выразить \mathbf{z} из распределения $q_{\phi}(\mathbf{z}|\mathbf{x})$ как детерминированную функцию \mathbf{g} от некоторой независимой случайной величины ϵ из распределения $p(\epsilon)$:

$$\mathbf{z} = \mathbf{g}_{\phi}(\mathbf{x}, \epsilon) \quad \epsilon \sim p(\epsilon)$$

Пусть $p(\epsilon) = \mathcal{N}(\mathbf{0}, I)$ и функция репараметризации имеет вид:

$$\mathbf{z} = \mathbf{g}_{\phi}(\mathbf{x}, \epsilon) = \boldsymbol{\mu}_{\phi}(\mathbf{x}) + \boldsymbol{\sigma}_{\phi}(\mathbf{x}) \cdot \epsilon$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{\phi}(\mathbf{x}), \boldsymbol{\sigma}_{\phi}(\mathbf{x}))$$

Градиент на E-шаге

Считаем градиент по параметрам ϕ при фиксированных θ :

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}, \theta)] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

Используем *reparameterization trick*, чтобы внести оператор градиента внутрь интеграла:

$$\nabla_{\phi} \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}|\mathbf{z}, \theta) d\mathbf{z} = \int p(\epsilon) \nabla_{\phi} \log p(\mathbf{x}|\mathbf{g}_{\phi}(\mathbf{x}, \epsilon), \theta) d\epsilon$$

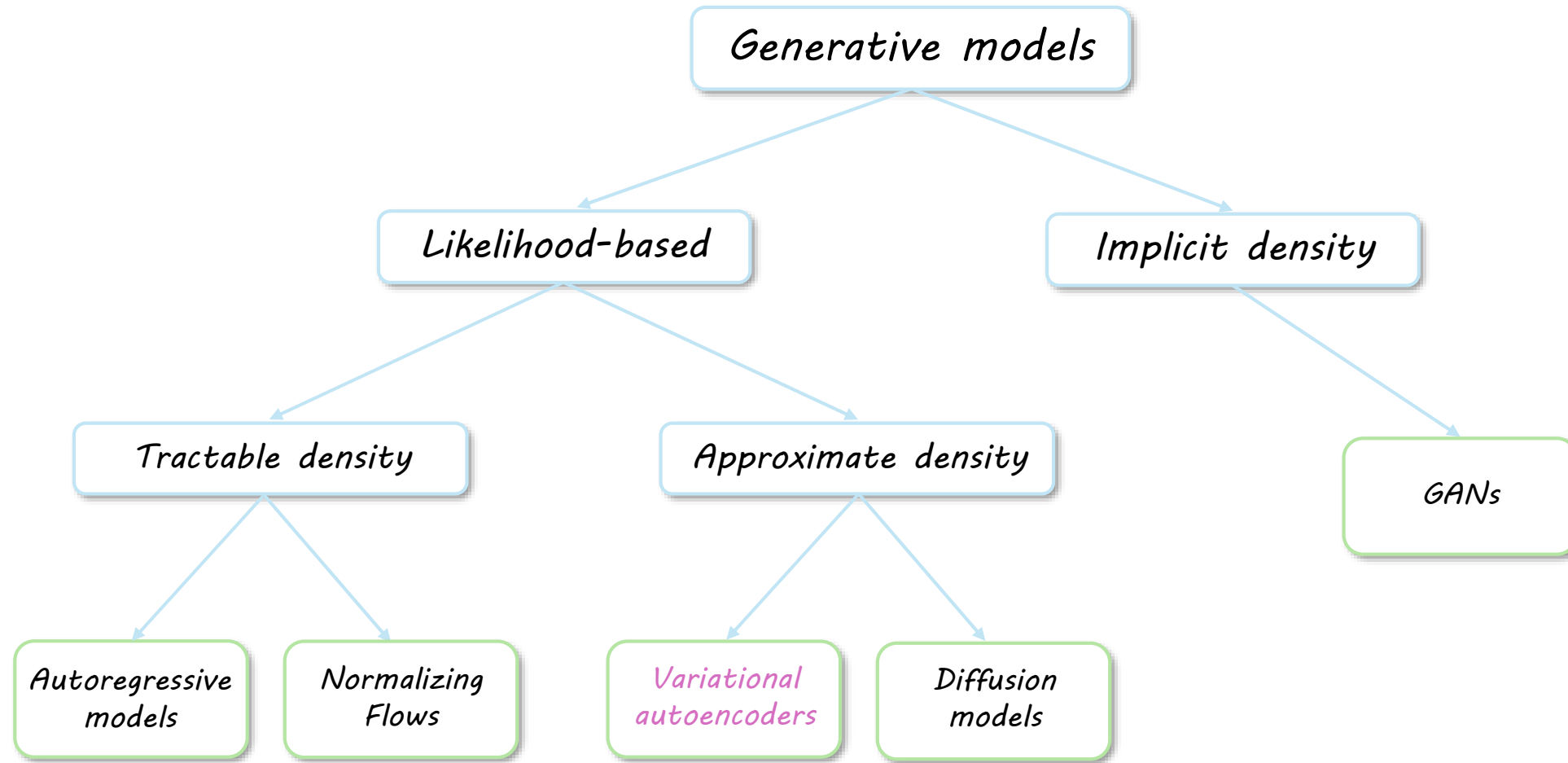
Оценим этот интеграл с помощью метода Монте-Карло:

$$\int p(\epsilon) \nabla_{\phi} \log p(\mathbf{x}|\mathbf{g}_{\phi}(\mathbf{x}, \epsilon), \theta) d\epsilon \approx \nabla_{\phi} \log p(\mathbf{x}|\mu_{\phi}(\mathbf{x}) + \sigma_{\phi}(\mathbf{x}) \cdot \epsilon^*, \theta), \quad \epsilon^* \sim \mathcal{N}(\mathbf{0}, I)$$

Для KL –дивергенции между двумя нормальными распределениями существует аналитическая формула:

$$\nabla_{\phi} KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) = \nabla_{\phi} KL(\mathcal{N}(\mu_{\phi}(\mathbf{x}), \sigma_{\phi}(\mathbf{x}))||\mathcal{N}(\mathbf{0}, I))$$

Зоопарк генеративных моделей

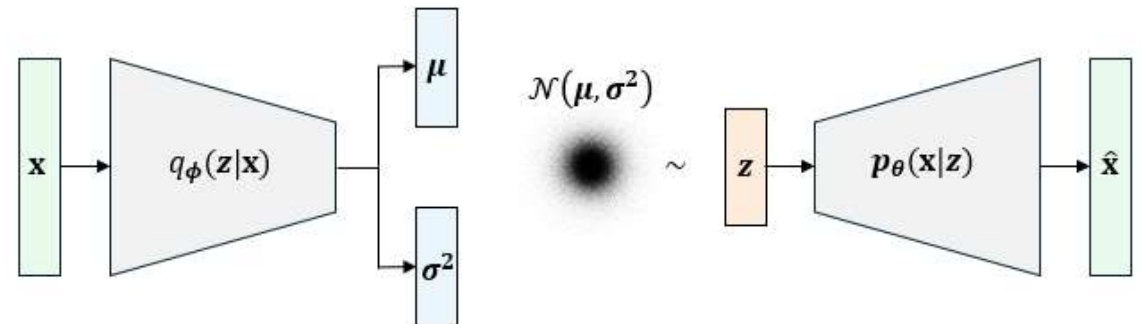


Variational Autoencoder

- **VAE** состоит из двух нейросетей: **кодировщика** и **декодировщика**

Общая идея:

- Одному объекту соответствует не одна точка, а целое распределение $q_{\phi}(\mathbf{z}|\mathbf{x})$
- Кодировщик $q_{\phi}(\mathbf{z}|\mathbf{x})$ выдает не один вектор, а два – вектор средних μ и дисперсий σ^2
- Из этого распределения берется случайная точка \mathbf{z} , которая затем подается в декодер для восстановления



Обучение VAE

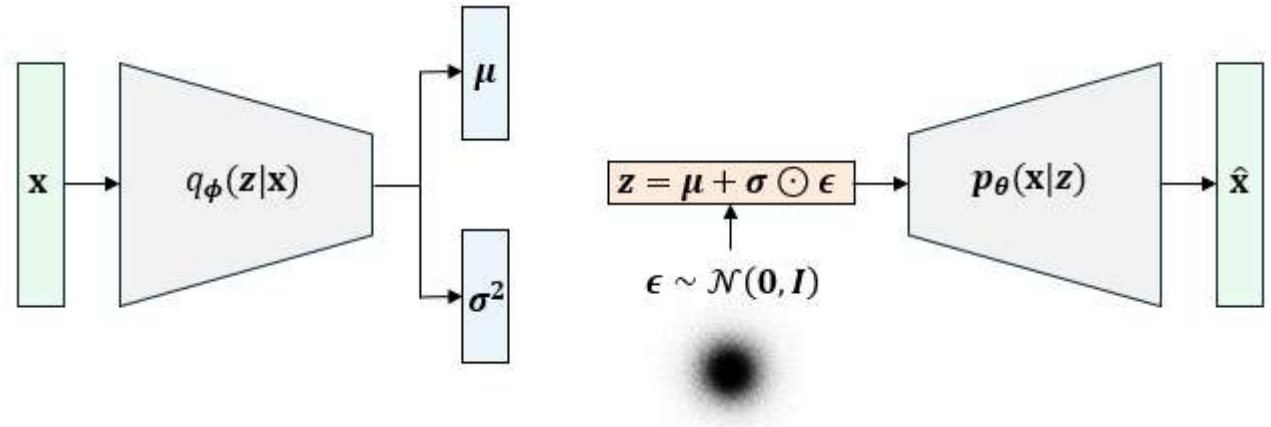
Алгоритм:

- Берем случайный объект \mathbf{x}
- Пропускаем его через кодировщик $q_{\phi}(\mathbf{z}|\mathbf{x})$, получаем $\boldsymbol{\mu}_{\phi}(\mathbf{x})$ и $\boldsymbol{\sigma}_{\phi}(\mathbf{x})$
- Вычисляем \mathbf{z}^* с помощью репараметризации:

Берем $\boldsymbol{\epsilon}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

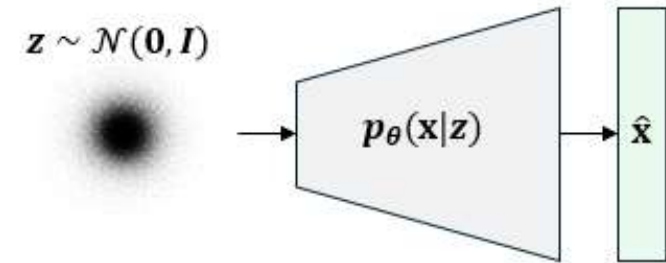
Считаем $\mathbf{z}^* = \boldsymbol{\mu}_{\phi}(\mathbf{x}) + \boldsymbol{\sigma}_{\phi}(\mathbf{x}) \cdot \boldsymbol{\epsilon}^*$

- Пропускаем \mathbf{z}^* через декодер и получаем восстановленный $\hat{\mathbf{x}}$
- Считаем лосс, обновляем параметры ϕ и θ



Генерация VAE

- Кодировщик $q_{\phi}(\mathbf{z}|\mathbf{x})$ больше не нужен
- Берем случайный вектор \mathbf{z}^* из априорного распределения $\mathcal{N}(\mathbf{0}, \mathbf{I})$
- Подаем \mathbf{z}^* на вход обученному декодеру



План

- Discrete VAE

- Vector Quantization

- Generative Adversarial Networks

Discrete VAE

Дискретный VAE

- Многие данные по своей природе дискретны
- Трансформеры работают именно с дискретными токенами

Пусть c — новая скрытая переменная является одной из K классов:

$$c \sim \text{Categorical}(\boldsymbol{\pi})$$

$\boldsymbol{\pi}$ — вектор вероятностей:

$$\boldsymbol{\pi} = (\pi_1, \dots, \pi_K), \quad \pi_k = P(c = k), \quad \sum_{k=1}^K \pi_k = 1$$

Выберем априорное распределение $p(c)$ самым простым – равномерным:

$$p(c) = \frac{1}{K}$$

Дискретный VAE

- Кодировщик $q_{\phi}(c|\mathbf{x})$ теперь будет предсказывать вектор вероятностей $\boldsymbol{\pi}$

Преобразуем *ELBO*:

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(c|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|c)] - KL(q_{\phi}(c|\mathbf{x})||p(c))$$

Распишем *KL* – дивергенцию:

$$KL(q_{\phi}(c|\mathbf{x})||p(c)) = \sum_{k=1}^K q_{\phi}(k|\mathbf{x}) \log \frac{q_{\phi}(k|\mathbf{x})}{p(k)} =$$

$$\sum_{k=1}^K q_{\phi}(k|\mathbf{x}) \log q_{\phi}(k|\mathbf{x}) - \sum_{k=1}^K q_{\phi}(k|\mathbf{x}) \log p(k) = -H(q_{\phi}(c|\mathbf{x})) + \log K$$

Дискретный VAE

Итоговая формула *ELBO* для дискретного VAE:

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(c|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|c)] + H(q_{\phi}(c|\mathbf{x})) - \log K$$

- Чтобы посчитать матожидание $\mathbb{E}_{q_{\phi}(c|\mathbf{x})}$, нам нужно сэмплировать c из распределения, которое выдал кодировщик
- Эта операция недифференцируема
- Трюк репараметризации больше не работает, нужно искать аналог для дискретных распределений

Vector Quantization

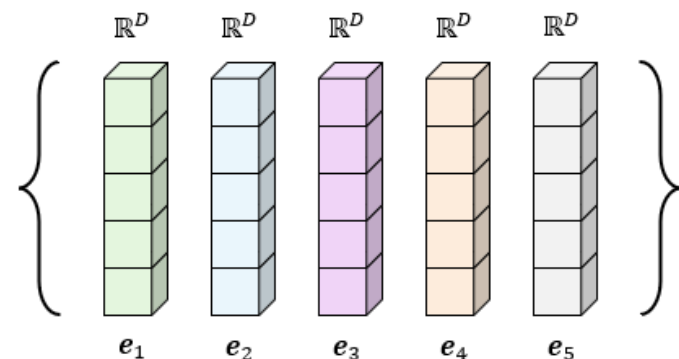
Векторное квантование

Идея **векторной квантизации** (*Vector Quantization, VQ*):

- Будем аппроксимировать непрерывное пространство конечным набором векторов

- Создадим кодовую книгу (**codebook**) – словарь из K обучаемых векторов $\{\mathbf{e}_k\}_{k=1}^K$

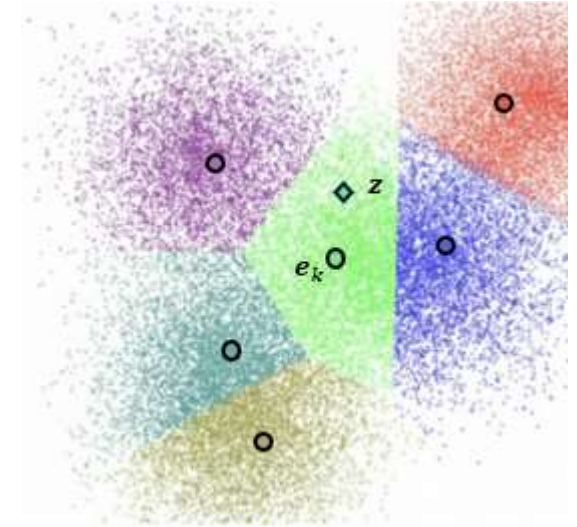
- Каждый \mathbf{e}_k – вектор в пространстве \mathbb{R}^D :
 K – размер словаря
 D – размерность каждого вектора



Векторное квантование

Для вектора $\mathbf{z}_e \in \mathbb{R}^D$ находим ближайший к нему вектор из кодовой книги и получаем квантованный вектор \mathbf{z}_q :

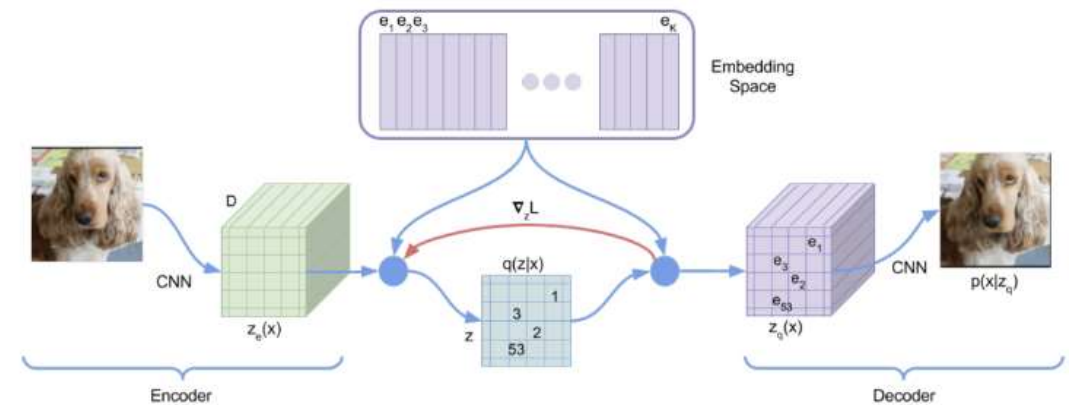
$$\mathbf{z}_q = q(\mathbf{z}_e) = \mathbf{e}_{k^*}, \quad k^* = \arg \min_k \|\mathbf{z}_e - \mathbf{e}_k\|$$



Это детерминированная операция, похожая на поиск ближайшего соседа в *K – Means*

Векторное квантование

- **Encoder**: принимает на вход объект x , сжимает его и выдает непрерывный вектор z_e (без μ и σ)
- **VQ**: берёт вектор z_e и находит ближайший к нему вектор e_{k^*} из кодовой книги
- **Decoder**: на вход декодеру приходит уже квантованный вектор $z_q = e_{k^*}$, из которого восстанавливается исходный объект



Векторное квантование

Особенностью $VQ - VAE$ является использование детерминированного вариационного апостериорного распределения:

$$q_{\phi}(c = k^* | \mathbf{x}) = \begin{cases} 1 & \text{для } k^* = \arg \min_k ||\mathbf{z}_e - \mathbf{e}_k|| \\ 0 & \text{иначе} \end{cases}$$

Такой выбор значительно упрощает $ELBO$:

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(c|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|c)] + \underbrace{H(q_{\phi}(c|\mathbf{x}))}_0 - \log K$$

Оптимизация сводится к максимизации члена реконструкции:

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(c|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|c)] - \log K = \log p_{\theta}(\mathbf{x}|\mathbf{z}_q) - \log K$$

Straight – Through Estimator

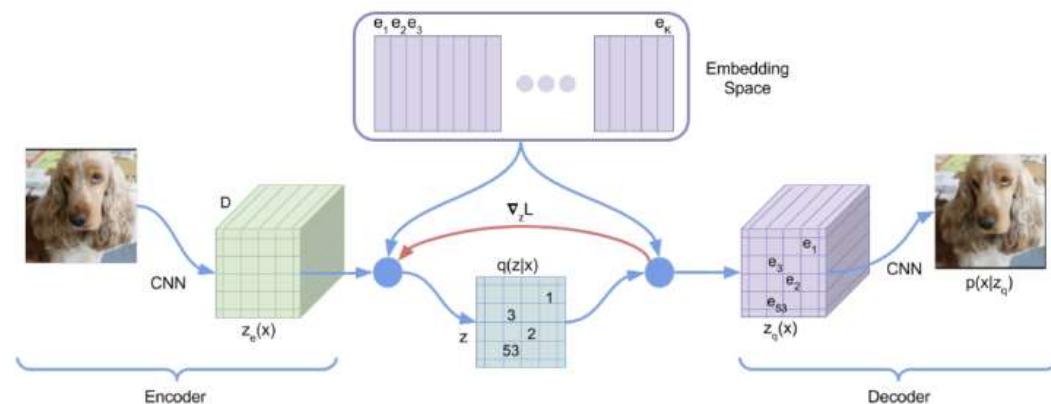
Операция квантования $k^* = \arg \min_k ||\mathbf{z}_e - \mathbf{e}_k||$ является недифференцируемой

Авторы $VQ - VAE$ предложили **прямое оценивание градиента (Straight – Through Estimator, STE)**

Механизм **STE**:

- На этапе обратного распространения градиент просто копируется с выхода недифференцируемого блока на его вход
- В $VQ - VAE$ градиент по \mathbf{z}_q напрямую передается \mathbf{z}_e :

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}_e} \approx \frac{\partial \mathcal{L}}{\partial \mathbf{z}_q}$$



Straight – Through Estimator

Математически это можно представить как аппроксимацию в *chain rule*:

$$\frac{\partial \log p_{\theta}(\mathbf{x}|\mathbf{z}_q)}{\partial \phi} = \frac{\partial \log p_{\theta}(\mathbf{x}|\mathbf{z}_q)}{\partial \mathbf{z}_q} \cdot \frac{\partial \mathbf{z}_q}{\partial \mathbf{z}_e} \cdot \frac{\partial \mathbf{z}_e}{\partial \phi} \approx \frac{\partial \log p_{\theta}(\mathbf{x}|\mathbf{z}_q)}{\partial \mathbf{z}_q} \cdot \frac{\partial \mathbf{z}_e}{\partial \phi}$$

- Выход кодировщика \mathbf{z}_e и квантованный вектор \mathbf{z}_q находятся рядом в одном и том же пространстве \mathbb{R}^D
- Градиент от декодера $\nabla_{\mathbf{z}_q} \mathcal{L}$ показывает, куда сдвинуть \mathbf{z}_q , чтобы улучшить реконструкцию
- Кодировщик применяет это к \mathbf{z}_e и подталкивает модель выбирать более подходящий вектор \mathbf{e}_k

Хотя такая оценка является смещенной, она часто указывает в направлении истинного градиента и хорошо работает на практике

VQ – VAE

Нам нужно обучать кодировщик ϕ , декодер θ и кодовую книгу $\{\mathbf{e}_k\}_{k=1}^K$

В VQ – VAE используется составная функция потерь:

$$\mathcal{L} = \underbrace{\log p_{\theta}(\mathbf{x}|\mathbf{z}_q)}_{\text{Reconstruction loss}} + \underbrace{\|sg[\mathbf{z}_e] - \mathbf{e}_{k^*}\|^2}_{\text{VQ – Loss}} + \underbrace{\beta \|\mathbf{z}_e - sg[\mathbf{e}_{k^*}]\|^2}_{\text{Commitment Loss}}$$

sg – оператор **stop – gradient**

VQ – Loss:

- Обновляет только *codebook*
- Заставляет выбранный вектор \mathbf{e}_{k^*} сдвинуться ближе к \mathbf{z}_e

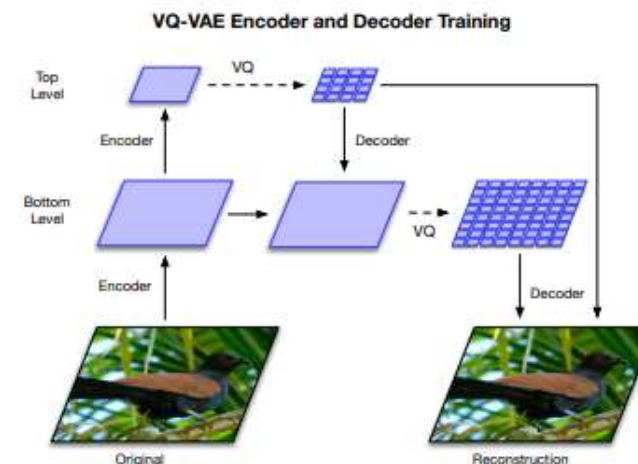
Commitment – Loss:

- Обновляет только кодировщик
- Заставляет выход кодировщика \mathbf{z}_e сдвинуться ближе к \mathbf{e}_{k^*}

VQ – VAE2

- $VQ - VAE$ использует один *codebook*
- Сразу хотим, чтобы модель умела улавливать и локальные, и глобальные детали

- **VQ – VAE2** использует несколько иерархических *codebooks* (верхний и нижний)
- Кодировщик выдает несколько карт признаков \mathbf{z}_e с разным пространственным разрешением
- Карты признаков верхнего и нижнего уровня квантуются с помощью своих *codebooks*



Generative Adversarial Networks

Likelihood based models

- Мы хотим научиться воспроизводить объекты из истинного распределения $p_{data}(\mathbf{x})$
- Обычно мы просто максимизируем правдоподобие $p_{\theta}(\mathbf{x})$

Модель 1:

$$p_1(\mathbf{x}) = \sum_{i=1}^n \mathcal{N}(\mathbf{x}|\mathbf{x}_i, \epsilon I)$$

- Строим узкий гауссовский пик каждого \mathbf{x}_i
- Для малого ϵ при генерации модель выдаст образец, почти идентичный одному из обучающих
- Если взять тестовый образец, то его правдоподобие будет низким

Модель 2:

$$p_2(\mathbf{x}) = 0.01p(\mathbf{x}) + 0.99p_{noise}(\mathbf{x})$$

- Модель с вероятностью 0.01 использует хорошую модель $p(\mathbf{x})$, и в 0.99 случаях генерируем шум
- Почти все сгенерированные объекты будут представлять шум
- Правдоподобие такой модели:

$$\begin{aligned} \log [0.01p(\mathbf{x}) + 0.99p_{noise}(\mathbf{x})] &\geq \\ \log [0.01p(\mathbf{x})] &= \log p(\mathbf{x}) - \log 100 \end{aligned}$$

Likelihood based models

- Способность генерировать качественные объекты не гарантирует высокого правдоподобия модели
 - Высокое правдоподобие не гарантирует, что модель будет генерировать качественные образцы
- Правдоподобие – не идеальная мера качества для генеративных моделей
 - Правдоподобие часто может быть невычислимым
- Это заставляет задуматься о новой идее, свободной от правдоподобия – ***идее состязательного обучения***

Likelihood free – learning

- Мы отказываемся от прямой аппроксимации плотности $p_{data}(\mathbf{x})$
- Вместо этого будем учить функцию генератор $p_{\theta}(\mathbf{x})$ создавать образцы, которые выглядят как будто они взяты из $p_{data}(\mathbf{x})$

У нас будет 2 набора данных:

- **Реальные образцы (real samples)** $\{\mathbf{x}_i\}_{i=1}^{n_1}$ из распределения $p_{data}(\mathbf{x})$
- **Сгенерированные образцы (fake samples)** $\{\mathbf{x}_i\}_{i=1}^{n_2}$ из распределения $p_{\theta}(\mathbf{x})$

Вводим **классификатор (Discriminator)**, который будет учиться различать эти 2 набора:

- $p(y = 1|\mathbf{x})$ – вероятность того, что образец \mathbf{x} является реальным
- $p(y = 0|\mathbf{x})$ – вероятность того, что образец \mathbf{x} является сгенерированным

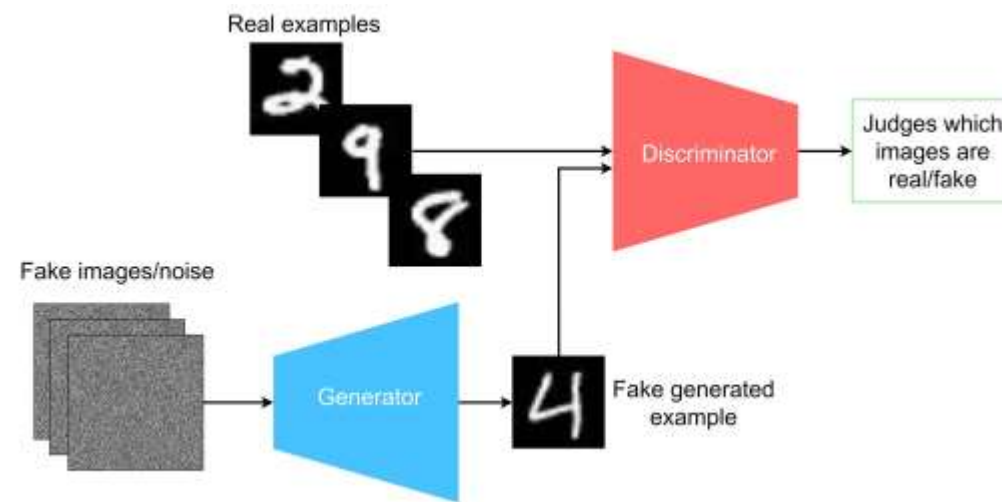
Likelihood free – learning

- Мы считаем, что генератор идеально изучил истинное распределение $p_{\theta}(\mathbf{x}) = p_{data}(\mathbf{x})$, когда лучший из возможных дискриминаторов не может отличить реальные образцы от сгенерированных
- В этой точке дискриминатор будет для любого \mathbf{x} выдавать вероятность 0.5
- Наша цель – достичь этого равновесия

Generative Adversarial Networks

Обучение GANs реализуется через состязание двух нейросетей:

- **Генератор G** : создаёт реалистичные данные $x = G(z)$ из вектора $z \sim p(z)$
- **Дискриминатор D** : бинарный классификатор, который получает на вход объект x и выдает вероятность $D(x)$, что x — реальный



Generator

По аналогии моделей со скрытыми переменными \mathbf{z} — это латентный вектор, $p(\mathbf{z})$ — априорное распределение

Мы можем думать о генераторе, как о детерминированном декодере:

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z}) \cdot p(\mathbf{z})$$

где $p_{\theta}(\mathbf{x}|\mathbf{z}) = \delta(\mathbf{x} - G_{\theta}(\mathbf{z}))$

δ — дельта-функция Дирака, для одного \mathbf{z} на выходе вся вероятность сосредоточена в одной точке $\mathbf{x} = G_{\theta}(\mathbf{z})$

Функция потерь

- Дискриминатор решает задачу бинарной классификации с метками $y = 1$ (реальный) и $y = 0$ (сгенерированный) и выдает вероятность $p(y = 1|\mathbf{x})$

- Хотим обучить D так, чтобы он максимизировал вероятность (правдоподобие) присвоения правильных меток всем классам:

$$\begin{aligned} \max_{p(y|\mathbf{x})} (\mathbb{E}_{p_{data}(\mathbf{x})} \log p(y = 1|\mathbf{x}) + \mathbb{E}_{p_{\theta}(\mathbf{x})} \log p(y = 0|\mathbf{x})) = \\ = \max_D (\mathbb{E}_{p_{data}(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p_{\theta}(\mathbf{x})} \log(1 - D(\mathbf{x}))) \end{aligned}$$

- Генератор, наоборот, пытается минимизировать эту функцию:

$$\begin{aligned} \min_G \max_D (\mathbb{E}_{p_{data}(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p_{\theta}(\mathbf{x})} \log(1 - D(\mathbf{x}))) \\ \min_G \max_D (\mathbb{E}_{p_{data}(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z})))) \end{aligned}$$

Функция потерь 1

$$\min_{p(y|\mathbf{x})} (-\mathbb{E}_{p_{data}(\mathbf{x})}[\log p(y = 1|\mathbf{x})] - \mathbb{E}_{p_{\theta}(\mathbf{x})}[\log p(y = 0|\mathbf{x})])$$

- $-\mathbb{E}_{p_{data}(\mathbf{x})}[\log p(y = 1|\mathbf{x})]$ отвечает за реальные данные
- Хотим, чтобы $D(\mathbf{x}) = p(y = 1|\mathbf{x})$ был равен 1
- Если $D(\mathbf{x}) = 1$, то $\log 1 = 0$
- Если $D(\mathbf{x}) = 0.01$, то $\log 0.01 \approx -4.6$
- $-\mathbb{E}_{p_{data}(\mathbf{x})}[\log p(y = 1|\mathbf{x})]$ – это штраф за то, что дискриминатор назвал реальный объект сгенерированным

Функция потерь 2

$$\min_{p(y|\mathbf{x})} (-\mathbb{E}_{p_{data}(\mathbf{x})} [\log p(y = 1|\mathbf{x})] - \mathbb{E}_{p_{\theta}(\mathbf{x})} [\log p(y = 0|\mathbf{x})])$$

- $-\mathbb{E}_{p_{\theta}(\mathbf{x})} [\log p(y = 0|\mathbf{x})]$ отвечает за сгенерированные данные
- Хотим, чтобы $1 - D(\mathbf{x}) = p(y = 0|\mathbf{x})$ был равен 1
- Если $D(\mathbf{x}) = 0$, то $\log 1 = 0$
- Если $D(\mathbf{x}) = 0.99$, то $\log 0.01 \approx -4.6$
- $-\mathbb{E}_{p_{\theta}(\mathbf{x})} [\log p(y = 0|\mathbf{x})]$ – это штраф за то, что дискриминатор назвал сгенерированный объект реальным

Оптимальный дискриминатор

Найдем оптимальный дискриминатор D^* при любом фиксированном генераторе G

Хотим найти D , который максимизирует $V(G, D)$:

$$V(G, D) = \mathbb{E}_{p_{data}(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p_{\theta}(\mathbf{x})} \log(1 - D(\mathbf{x})) =$$
$$\int [p_{data}(\mathbf{x}) \log D(\mathbf{x}) + p_{\theta}(\mathbf{x}) \log(1 - D(\mathbf{x}))] d\mathbf{x}$$

$$\frac{dy(D)}{dD} = \frac{p_{data}(\mathbf{x})}{D(\mathbf{x})} - \frac{p_{\theta}(\mathbf{x})}{1 - D(\mathbf{x})} = 0$$

$$D^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_{\theta}(\mathbf{x})}$$

Генератор

Найдем цель генератора при условии оптимального дискриминатора D^*

Подставим D^* в нашу функцию $V(G, D)$:

$$\begin{aligned} V(G, D^*) &= \int p_{data}(\mathbf{x}) \log D(\mathbf{x}) d\mathbf{x} + \int p_{\theta}(\mathbf{x}) \log(1 - D(\mathbf{x})) \\ &= \int p_{data}(\mathbf{x}) \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_{\theta}(\mathbf{x})} d\mathbf{x} + \int p_{\theta}(\mathbf{x}) \log \left(1 - \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_{\theta}(\mathbf{x})} \right) \\ &= KL \left(p_{data}(\mathbf{x}) \parallel \frac{p_{data}(\mathbf{x}) + p_{\theta}(\mathbf{x})}{2} \right) + KL \left(p_{\theta}(\mathbf{x}) \parallel \frac{p_{data}(\mathbf{x}) + p_{\theta}(\mathbf{x})}{2} \right) - 2 \log 2 \\ &= 2JSD(p_{data}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) - 2 \log 2 \end{aligned}$$

Оптимальный GAN

Дивергенция Йенсена-Шеннона:

$$JSD(p_{data}(\mathbf{x}) || p_{\theta}(\mathbf{x})) = \frac{1}{2} \left[KL \left(p_{data}(\mathbf{x}) || \frac{p_{data}(\mathbf{x}) + p_{\theta}(\mathbf{x})}{2} \right) + KL \left(p_{\theta}(\mathbf{x}) || \frac{p_{data}(\mathbf{x}) + p_{\theta}(\mathbf{x})}{2} \right) \right]$$

Истинная цель генератора – минимизация дивергенции Йенсена-Шеннона:

$$2JSD(p_{data}(\mathbf{x}) || p_{\theta}(\mathbf{x})) - 2 \log 2 \rightarrow \min_{\theta}$$

- $JSD(p||q) = 0$ когда $p = q$
- В этой точке $D^*(\mathbf{x}) = \frac{1}{2}$, $V(G^*, D^*) = -2 \log 2$
- Теоретическое равновесие достигается, когда генератор идеально воспроизводит данные, а дискриминатор не может отличить реальные объекты от сгенерированных

Проблемы GAN

Теоретическое равновесие требует, чтобы модели G и D были достаточно мощными (любыми функциями)

В действительности G и D – нейросети с конечным набором параметров

Из-за этой ограниченной мощности настоящее теоретическое равновесие для многих GAN архитектур может не существовать

Это и является основной причиной многих проблем при обучении GAN

Обучение GAN

Теперь наша *min max* задача будет происходить над конкретными параметрами

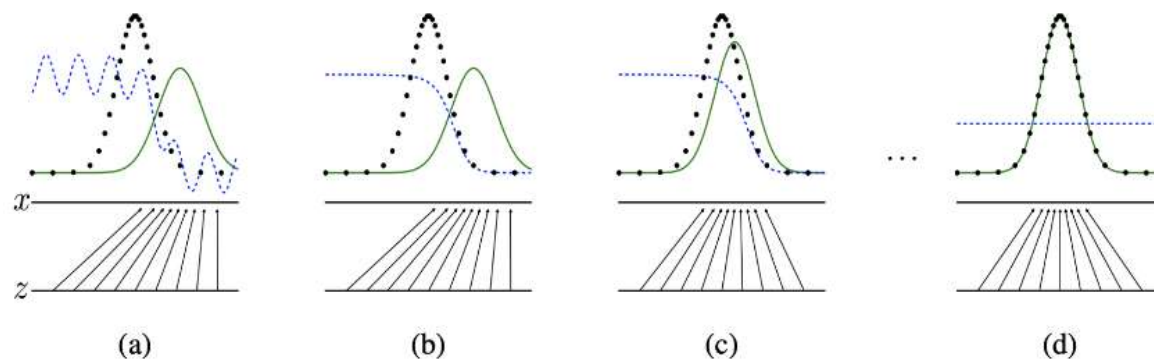
$$\min_{\theta} \max_{\phi} [\mathbb{E}_{p_{data}(\mathbf{x})} \log D_{\phi}(\mathbf{x}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D_{\phi}(G_{\theta}(\mathbf{z})))]$$

Обучаем D :

- Фиксируем веса генератора θ
- Делаем шаг градиентного спуска по ϕ , чтобы максимизировать $V(G, D)$

Обучаем G :

- Фиксируем веса дискриминатора ϕ
- Делаем шаг градиентного спуска по θ , чтобы минимизировать $V(G, D)$



Спасибо за внимание!