

# Генеративные модели

Лекция 2: Нормализующие потоки

# Постановка задачи

Пусть у нас есть набор данных  $\{\mathbf{x}_i\}_{i=1}^n$ ,

$\mathbf{x}_i \in \mathbb{R}^m$  — вектор, представляющий объект

$p_{data}(\mathbf{x})$  описывает как объекты распределены в пространстве  $\mathbb{R}^m$

Цель: найти неизвестное распределение  $p_{data}(\mathbf{x})$

Зная  $p_{data}(\mathbf{x})$  мы сможем оценивать вероятность новых объектов и генерировать новые

# Проблема высокой размерности

Объекты находятся в пространстве высокой размерности

Прямая оценка  $p_{data}(\mathbf{x})$  в таких пространствах затруднена

Найти истинное распределение  $p_{data}(\mathbf{x})$  почти невозможно

Будем искать аппроксимацию внутри некоторого заранее заданного *параметрического семейства распределений*  $p_{\theta}(\mathbf{x})$

$\theta$  — параметры модели

Ищем  $\theta^*$ , которые максимизируют правдоподобие:

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \log p_{\theta}(\mathbf{x}_i)$$

# Связь правдоподобия и $KL$ -дивергенции

Наша задача — найти такие параметры  $\theta$ , при которых  $p_{\theta}(\mathbf{x})$  будет максимально близко к  $p_{data}(\mathbf{x})$ :

$$p_{\theta}(\mathbf{x}) \approx p_{data}(\mathbf{x})$$

Будем искать такие параметры  $\theta$ , которые минимизируют дивергенцию между истинным и модельным распределениями:

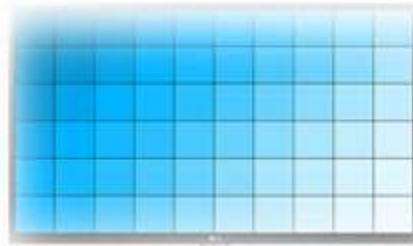
$$\min_{\theta} D(p_{data}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x}))$$

Минимизация прямой  $KL$  — дивергенции эквивалента максимизации правдоподобия

# Обучение авторегрессионных моделей

Объекты в реальном мире – многомерные случайные величины

Прямое моделирование  $p_{\theta}(\mathbf{x}) = p_{\theta}(x_1, x_2, \dots, x_D)$  – сложная задача из-за проклятия размерности



# Разложение многомерной плотности

Разложим совместную вероятность на произведение условных:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_1, x_2, \dots, x_D) = p_{\theta}(x_1) \cdot p_{\theta}(x_2|x_1) \cdot p_{\theta}(x_3|x_1, x_2) \cdot \dots \cdot p_{\theta}(x_D|x_1, x_2, \dots, x_{D-1})$$

$$= \prod_j^D p_{\theta}(x_i|\mathbf{x}_{1:j-1})$$

$\mathbf{x}_{1:j-1}$  — вектор всех компонент объекта  $\mathbf{x}$  до позиции  $j$

Теперь задача максимизации правдоподобия принимает вид:

$$\theta^* = \underset{\theta}{argmax} \sum_{i=1}^n \sum_j^D \log p_{\theta}(x_{ij}|\mathbf{x}_{i,1:j-1})$$

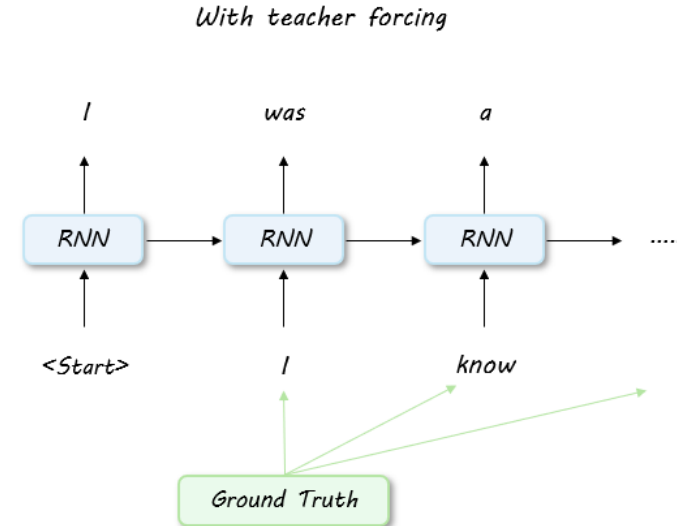
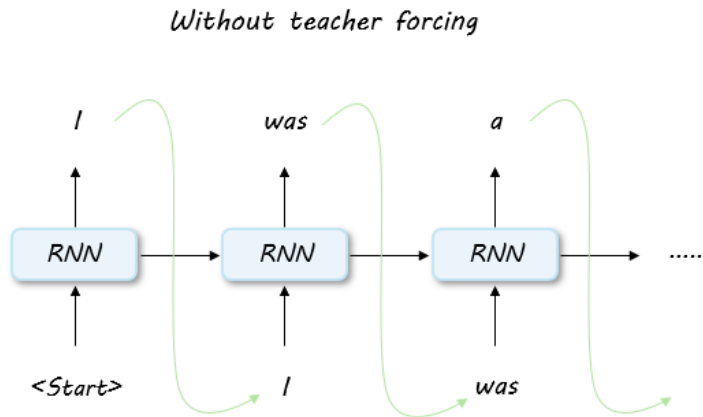


# Обучение авторегрессионных моделей

В начале обучения модель часто совершает ошибки, которые с каждым шагом накапливаются

## **Решение:**

Принудительно подавать на вход модели правильные данные (*teacher forcing*)



При генерации модель использует собственные предсказания

# План

- Идея нормализующих потоков

- Авторегрессионные потоки

- Линейные потоки

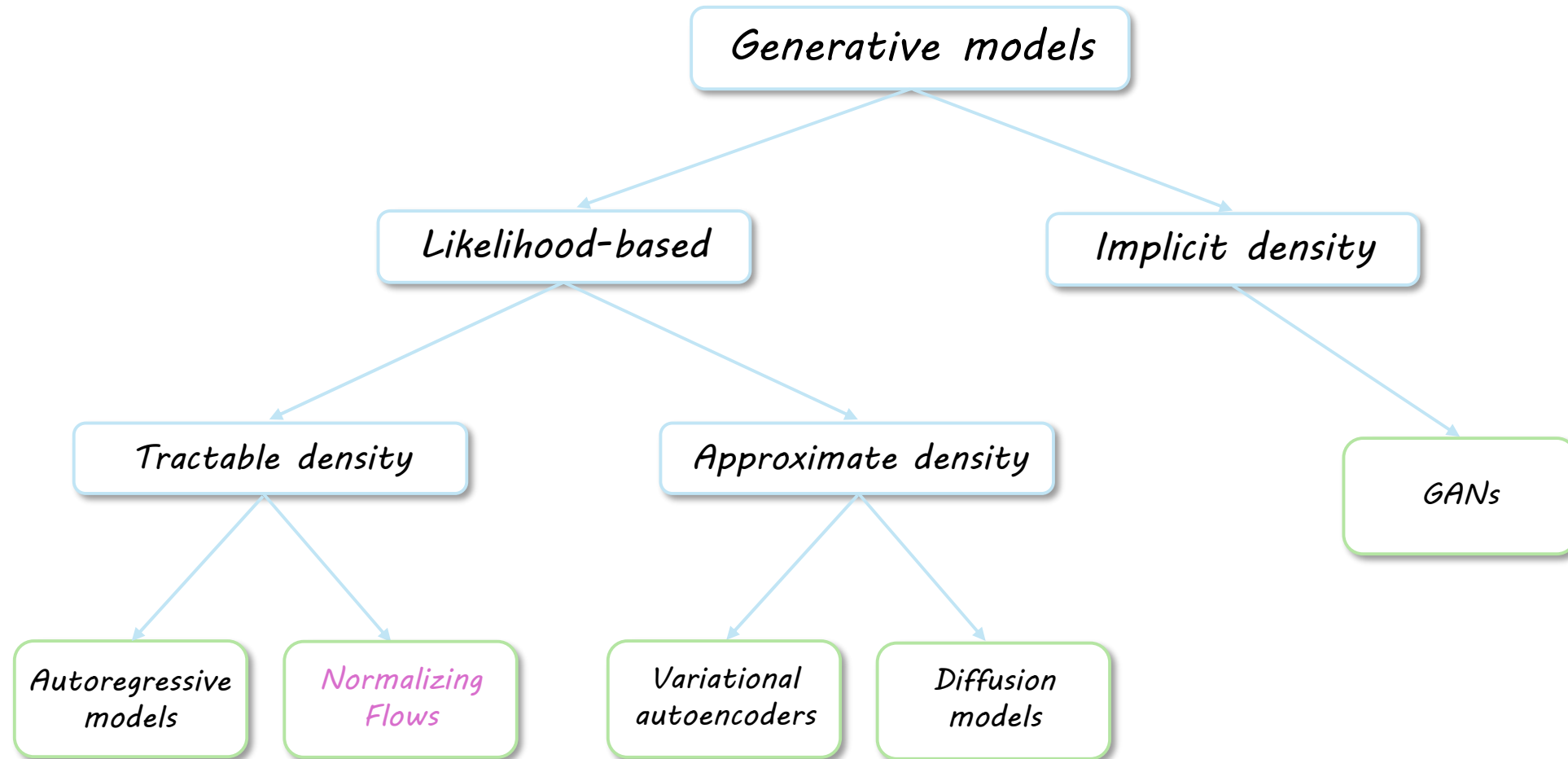
- *NICE* и *RealNVP*

- *Glow*



# Нормализующие потоки

# Зоопарк генеративных моделей



# Разложение многомерной плотности

## *Авторегрессионные модели:*

- ⊕ Точно вычисляют правдоподобие
- ⊖ Медленная генерация

## *Вариационные автокодировщики:*

- ⊕ Быстрая генерация из латентного пространства
- ⊖ Не могут точно вычислить правдоподобие

## *Нормализующие потоки:*

- ⊕ Точно вычисляют правдоподобие
- ⊕ Используют латентное пространство для быстрой генерации

Нормализующие потоки берут лучшее от двух миров

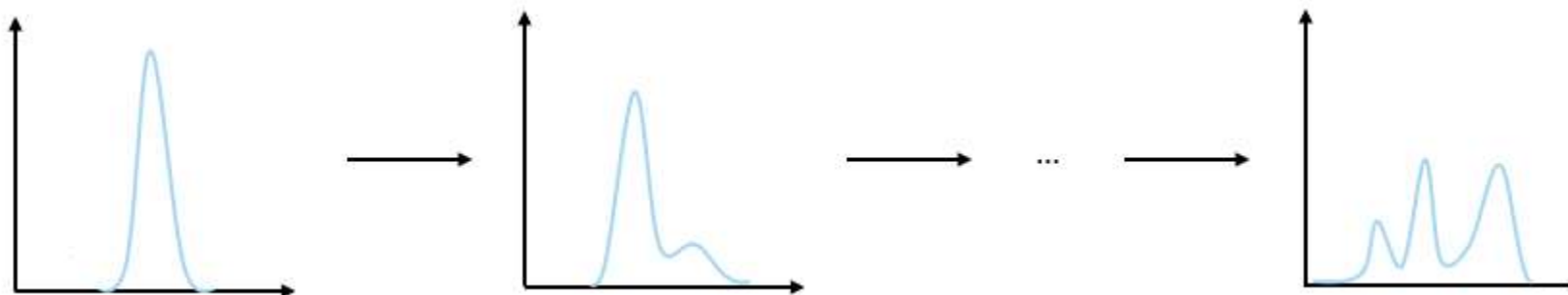
# Идея нормализующих потоков

*Идея:*

Возьмём простое, хорошо изученное распределение:

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

Шаг за шагом будем усложнять это распределение, пока оно не станет похоже на наши данные  $\mathbf{x}$



# Идея нормализующих потоков

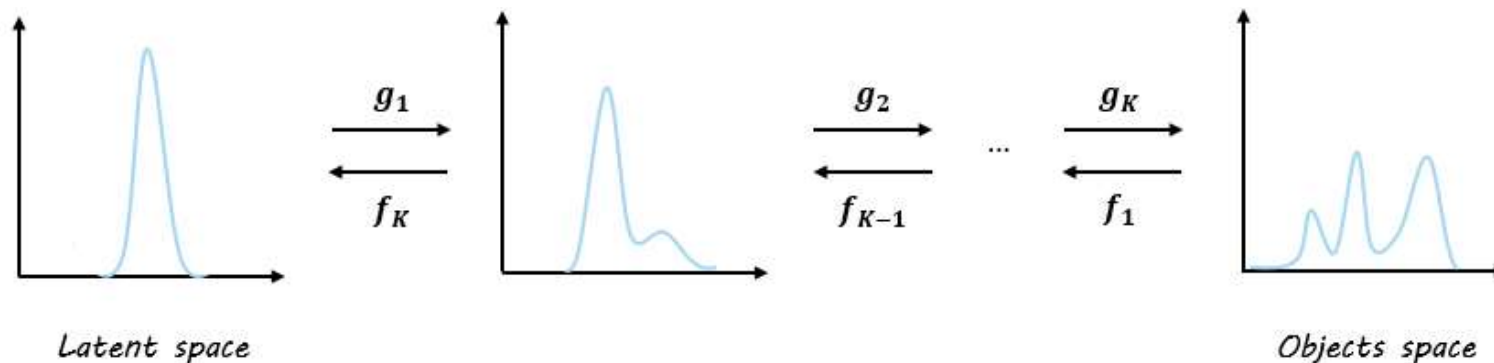
Особенность потоков – наличие двух направлений:

- **Нормализующее направление (для обучения)**  $x \rightarrow z$ :

Берём реальный объект  $x$  и с помощью нормализующей функции  $f$  превращаем его в простой  $z$

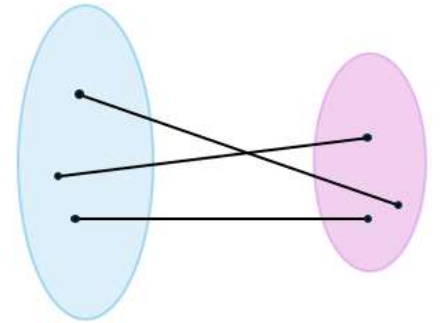
- **Генеративное направление**  $z \rightarrow x$ :

Берём простой объект  $z$  и с помощью генеративной функции  $g = f^{-1}$  создаём объект  $x$



# Идея нормализующих потоков

Функции взаимно обратные  $\rightarrow$  для каждой точки в пространстве  $\mathbf{x}$  существует ровно одна из пространства  $\mathbf{z}$



Размерности исходного и преобразованного пространств должны совпадать:

$$\dim \mathbf{x} = \dim \mathbf{z}$$

# Теорема о замене переменных

Наша цель – найти распределение  $p(\mathbf{x})$

Связать неизвестную плотность  $p(\mathbf{x})$  с известной плотностью простого распределения  $p(\mathbf{z})$  нам позволяет **формула замены переменных**:

$$p(\mathbf{x}) = p(\mathbf{z}) |\det(J_f)| = p(\mathbf{z}) \left| \det \left( \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right) \right| = p(\mathbf{f}(\mathbf{x})) \left| \det \left( \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$

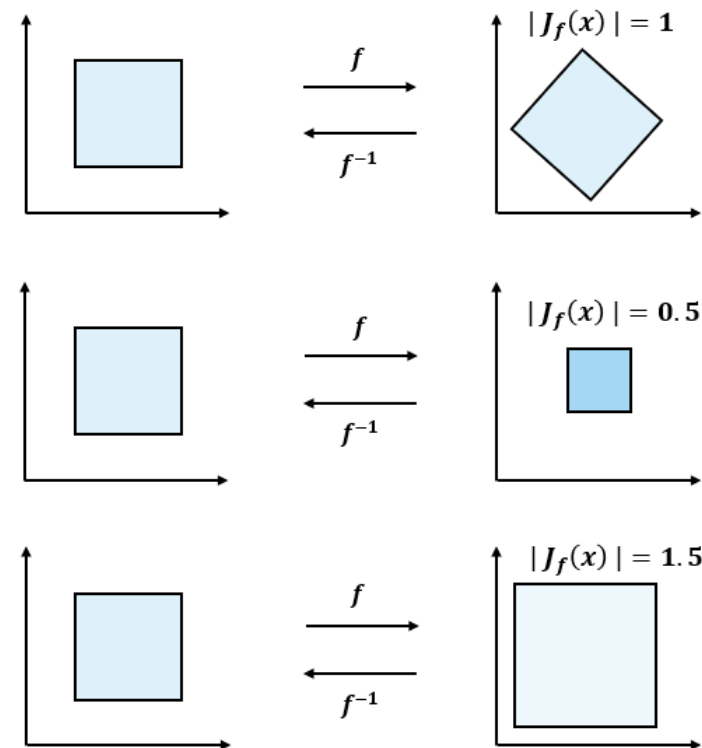
Здесь  $J_f(\mathbf{x})$  – матрица Якоби:

$$J_f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_D}{\partial x_1} & \dots & \frac{\partial f_D}{\partial x_D} \end{pmatrix}$$

# Якобиан

- Интуитивно определитель Якобиана можно понимать как коэффициент изменения объема
- Он показывает, насколько  $f$  сжимает или растягивает пространство

- Если объем не изменился  $|\det J_f(x)| = 1$ , плотность в этой точке не изменится
- Если пространство сжалось  $|\det J_f(x)| < 1$ , плотность в этой точке должна увеличиться
- Если пространство растянулось  $|\det J_f(x)| > 1$ , плотность в этой точке должна уменьшиться





# Пример

Пусть есть одномерная случайная величина  $z \sim \mathcal{N}(0, 1)$ :

$$\pi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

Хотим получить новую переменную:

$$x = g(z) = 5z + 1$$

Найдем нормализующую функцию  $f$ , которая является обратной к  $g$ :

$$z = f(x) = g^{-1}(x) = \frac{x - 1}{5}$$

# Пример

Якобиан в одномерном случае:

$$\left| \frac{df(x)}{dx} \right| = \left| \frac{d}{dx} \left( \frac{x-1}{5} \right) \right| = \left| \frac{1}{5} \right| = \frac{1}{5}$$

Подставляем и получаем распределение для  $x$ :

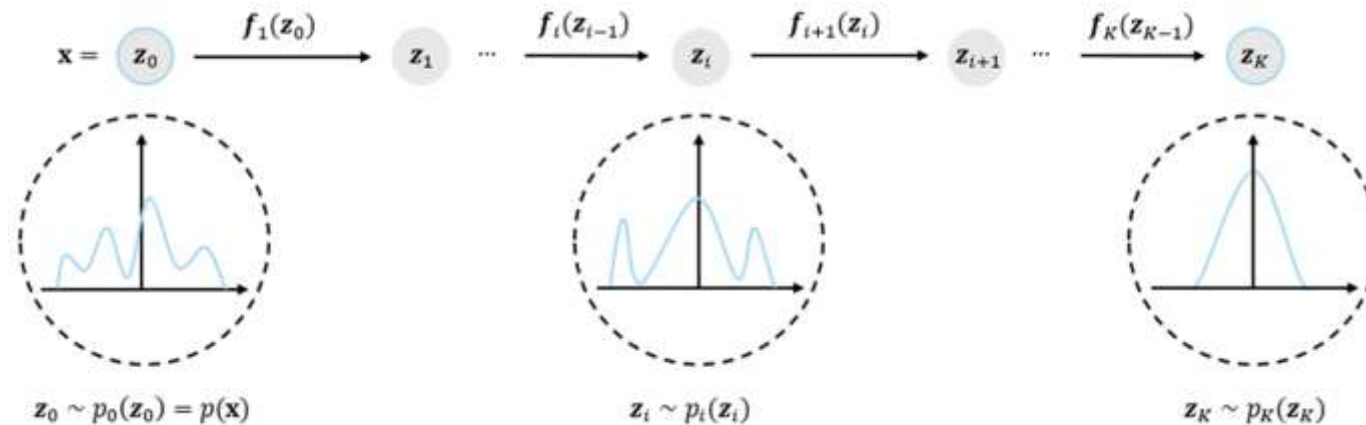
$$p(x) = \frac{1}{5} \pi \left( \frac{x-1}{5} \right) = \frac{1}{\sqrt{2\pi \cdot 5^2}} e^{-\frac{(x-1)^2}{2 \cdot 5^2}} = \mathcal{N}(1, 25)$$

# Композиция преобразований

- Одной функции  $f$  недостаточно, чтобы из  $\mathcal{N}(0,1)$  получить сложное распределение наших данных
- Потoki используют последовательность из  $K$  простых и обратимых преобразований:

$$\mathbf{x} = \mathbf{z}_0 \xrightarrow{f_1} \mathbf{z}_1 \xrightarrow{f_2} \dots \xrightarrow{f_K} \mathbf{z}_K$$

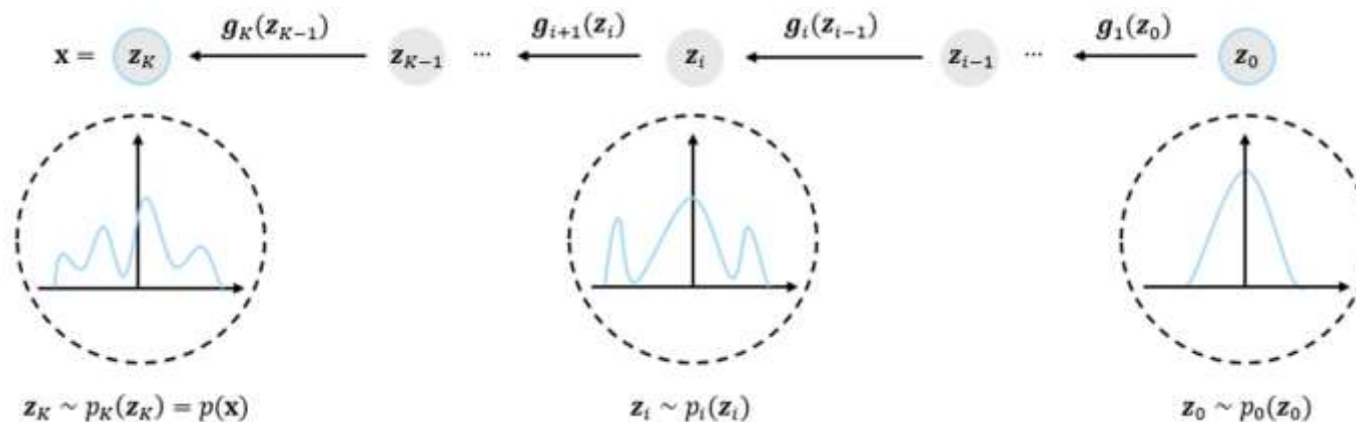
- Во время обучения мы используем нормализующее направление  $\mathbf{x} \rightarrow \mathbf{z}$ :



# Композиция преобразований

- Генеративное направление  $\mathbf{z} \rightarrow \mathbf{x}$  используется для создания новых объектов:

$$\mathbf{z}_0 \xrightarrow{g_1} \mathbf{z}_1 \xrightarrow{g_2} \dots \xrightarrow{g_K} \mathbf{z}_K = \mathbf{x}$$



# Композиция преобразований

Нужно найти определитель якобиана всей композиции

Два математических факта:

- Матрица Якоби для композиций функций равна их произведению:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}_K}{\partial \mathbf{z}_{K-1}} \cdot \frac{\partial \mathbf{f}_{K-1}}{\partial \mathbf{z}_{K-2}} \cdot \dots \cdot \frac{\partial \mathbf{f}_1}{\partial \mathbf{z}_0}$$

- Определитель произведения матриц равен произведению определителей:

$$\det(\mathbf{J}_f) = \det(\mathbf{J}_{f_K}) \cdot \det(\mathbf{J}_{f_{K-1}}) \cdot \dots \cdot \det(\mathbf{J}_{f_1})$$

Будем считать определители для каждого шага  $\mathbf{f}_k$  и перемножать их:

$$\det\left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}}\right) = \prod_{k=1}^K \det\left(\frac{\partial \mathbf{f}_k}{\partial \mathbf{z}_{k-1}}\right)$$

# Обучение нормализующих потоков

Чтобы сделать эту композицию обучаемой, представим  $f_k$  как нейросеть с параметрами  $\theta$

Плотность, которую моделирует наша модель:

$$p_{\theta}(\mathbf{x}) = p(\mathbf{f}_{\theta}(\mathbf{x})) \left| \det \left( \frac{\partial \mathbf{f}_{\theta}}{\partial \mathbf{x}} \right) \right| = p(\mathbf{z}_K) \left| \prod_{k=1}^K \det \left( \frac{\partial \mathbf{f}_{k,\theta}}{\partial \mathbf{z}_{k-1}} \right) \right|$$

$\mathbf{z}_K$  — финальный вектор в латентном пространстве

Максимизируем правдоподобие:

$$\log p_{\theta}(\mathbf{x}) = \log p(\mathbf{z}_K) + \sum_{k=1}^K \log \left| \det \left( \frac{\partial \mathbf{f}_{k,\theta}}{\partial \mathbf{z}_{k-1}} \right) \right|$$

- $\log p(\mathbf{z}_K)$  отвечает за то, чтобы выход потока соответствовал базовому распределению
- $\sum_{k=1}^K \log \left| \det \left( \frac{\partial \mathbf{f}_{k,\theta}}{\partial \mathbf{z}_{k-1}} \right) \right|$  отслеживает суммарное изменение объема после всех преобразований

# Функция потерь

Функция потерь:

$$\mathcal{L}_\theta(\mathbf{x}) = -\log p_\theta(\mathbf{x}) = -\left(\log p(\mathbf{z}_K) + \sum_{k=1}^K \log \left| \det \left( \frac{\partial \mathbf{f}_{k,\theta}}{\partial \mathbf{z}_{k-1}} \right) \right| \right)$$

Требования к преобразованиям  $\mathbf{f}_{k,\theta}$ :

- Наличие обратной функции  $\mathbf{g}_\theta = \mathbf{f}_\theta^{-1}$
- Быстрое вычисление определителя  $\det(\mathbf{J}_f)$  (в идеале за  $\mathcal{O}(D)$ )



Будем искать функции, которые удовлетворяют этим условиям

# Авторегрессионные потоки



# Авторегрессионные потоки

Авторегрессионные потоки вдохновлены идеей классических  $AR$  моделей:

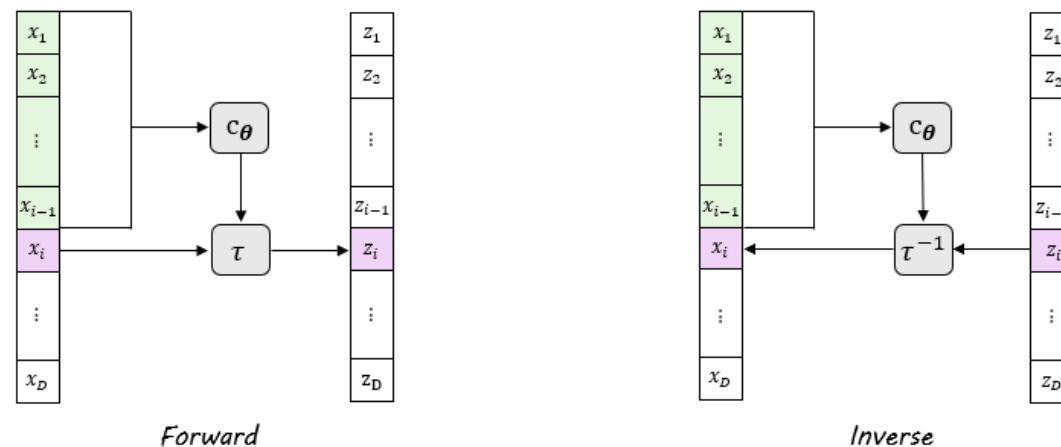
$$z_i = \tau(x_i; c_\theta(\mathbf{x}_{1:i-1}))$$

$\tau$  – **transformer**, обратимая функция, которая преобразует  $x_i$  в  $z_i$

$c_\theta(\mathbf{x}_{1:i-1})$  – **conditioner**, модель, которая генерирует параметры для  $\tau$  на основе  $\mathbf{x}_{1:i-1}$

Если  $\tau$  обратим, то и весь поток обратим:

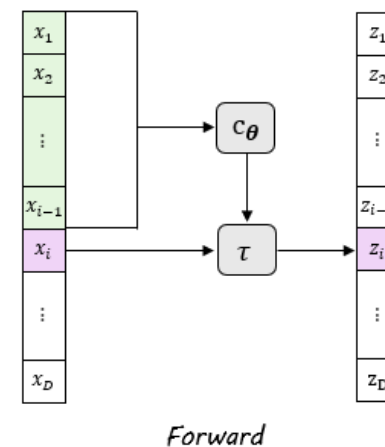
$$x_i = \tau^{-1}(z_i; c_\theta(\mathbf{x}_{1:i-1}))$$



# Авторегрессионные потоки

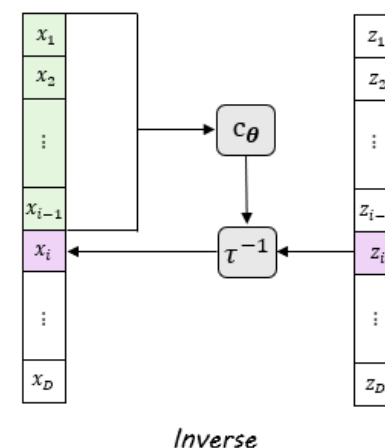
- При прямом проходе  $\mathbf{x}$  полностью известен и мы можем параллельно вычислить все  $z_i$

$$z_i = \tau(x_i; c_{\theta}(\mathbf{x}_{1:i-1}))$$



- При обратном проходе  $\mathbf{z}$  полностью известен, но для генерации  $x_i$  нам нужен контекст  $\mathbf{x}_{1:i-1}$

$$x_i = \tau^{-1}(z_i; c_{\theta}(\mathbf{x}_{1:i-1}))$$



# Авторегрессионные потоки

Якобиан имеет треугольный вид, поскольку  $\frac{\partial z_i}{\partial x_j} = 0$  при  $j > i$ :

$$J_{f_\theta}(\mathbf{x}) = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ \frac{\partial z_D}{\partial x_1} & \dots & \frac{\partial z_D}{\partial z_D} \end{pmatrix}$$

Определитель треугольной матрицы считается за  $\mathcal{O}(D)$ :

$$\log |\det J_{f_\theta}(\mathbf{x})| = \sum_{i=1}^D \log \left| \frac{\partial \tau(x_i; c_\theta(\mathbf{x}_{1:i-1}))}{\partial x_i} \right|$$

# Affine Transformers

Аффинное преобразование:

$$z_i = \alpha_i \cdot x_i + \beta_i$$

$\alpha_i$  и  $\beta_i$  — параметры масштаба и сдвига

Такой поток всегда обратим (если  $\alpha_i \neq 0$ ):

$$x_i = \frac{z_i - \beta_i}{\alpha_i}$$

Частная производная  $\frac{\partial z_i}{\partial x_j} = \alpha_i$

Логарифм определителя якобиана для всего преобразования:

$$\log |\det J_{f_\theta}(\mathbf{x})| = \sum_i^D |\alpha_i|$$

# Masked Conditioners

Хотим посчитать параметры  $c_{\theta}(\mathbf{x}_{1:i-1})$  для генерации  $z_i$ :

$$z_i = \tau(x_i; c_{\theta}(\mathbf{x}_{1:i-1}))$$

Будем использовать единую нейросеть для вычисления всех параметров

Будем использовать маски, чтобы сохранить авторегрессию в этой сети:

- В полносвязных сетях на веса накладываются бинарные маски (**MADE**)
- В свёрточных сетях используются **causal convolution** (**PixelCNN**)

Либо на обучении, либо на генерации будет медленная скорость

Есть 2 популярные архитектуры:

- **MAF** (**Masked Autoregressive Flow**) — создана для быстрого обучения
- **IAF** (**Inverse Autoregressive Flow**) — та же архитектура, но развернута для быстрой генерации

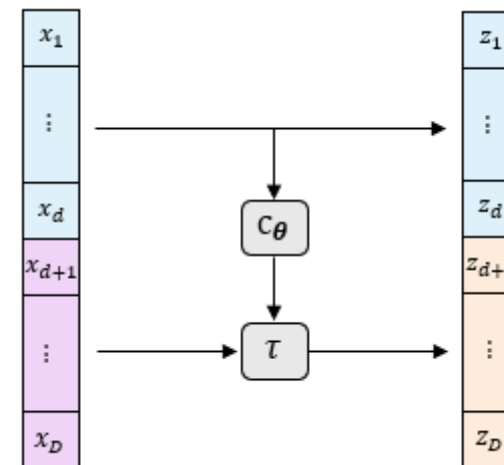
# Coupling Layers

Чтобы решить проблему вычислительной асимметрии были предложены **слои связи** (*coupling layers*)

Делим входной вектор  $\mathbf{x}$  на две части:

$$\mathbf{x} = [\mathbf{x}_{1:d}; \mathbf{x}_{d+1:D}]$$

- Первая часть  $\mathbf{x}_{1:d}$  преобразуется напрямую без использования  $c_\theta$
- Вторая часть  $\mathbf{x}_{d+1:D}$  преобразуется с помощью  $\tau$



Обычно преобразование первой части тождественное:

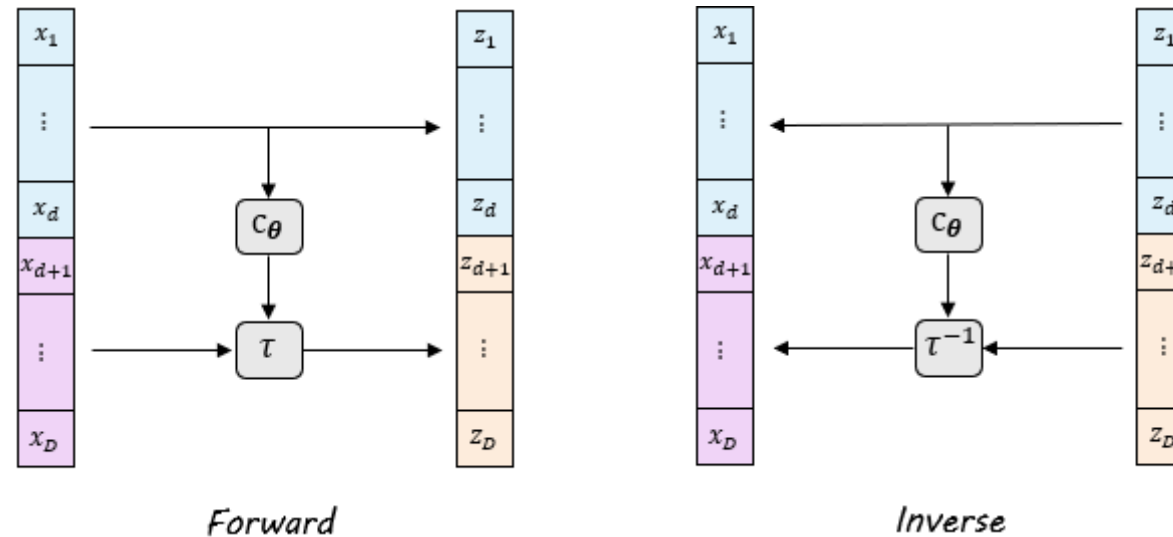
$$\begin{cases} \mathbf{z}_{1:d} = \mathbf{x}_{1:d} \\ \mathbf{z}_{d+1:D} = \tau(\mathbf{x}_{d+1:D}; c_\theta(\mathbf{x}_{1:d})) \end{cases}$$

# Свойства слоёв связи

## Симметричные вычисления:

- Прямой проход быстр, так как  $\mathbf{x}$  известен
- Обратный проход тоже быстр, так как мы знаем первую половину  $\mathbf{x}_{1:d}$ :

$$\begin{cases} \mathbf{x}_{1:d} = \mathbf{z}_{1:d} \\ \mathbf{x}_{d+1:D} = \tau^{-1}(\mathbf{z}_{d+1:D}; c_{\theta}(\mathbf{x}_{1:d})) = \tau^{-1}(\mathbf{z}_{d+1:D}; c_{\theta}(\mathbf{z}_{1:d})) \end{cases}$$



# Свойства слоёв связи

*Блочный нижнетреугольный якобиан:*

$$J_f = \begin{pmatrix} \frac{\partial \mathbf{z}_{1:d}}{\partial \mathbf{x}_{1:d}} & \frac{\partial \mathbf{z}_{1:d}}{\partial \mathbf{x}_{d+1:D}} \\ \frac{\partial \mathbf{z}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \frac{\partial \mathbf{z}_{d+1:D}}{\partial \mathbf{x}_{d+1:D}} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A} & \mathbf{D} \end{pmatrix}$$

Определитель равен произведению определителей диагонального блока:

$$\det(J_f) = \det(\mathbf{I}) \cdot \det(\mathbf{D}) = \det(\mathbf{D})$$



# Линейные потоки

# Линейные потоки

- Авторегрессионные потоки и слои связи чувствительны к порядку переменных
- Обычное *перемешивание* (*permutation*) – не обучаемая операция

## **Решение:**

- Вставить между нелинейными слоями обучаемый линейный поток:

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

- Это позволит модели самой находить оптимальное перемешивание

- Линейные потоки ограничены в выразительности
- Используются в основном для перестановок

# Линейные потоки

## **Проблема:**

- Определитель якобиана считается за  $\mathcal{O}(D^3)$

## **Решение:**

- Не будем обучать матрицу напрямую, параметризуем её через матричные разложения

Самый популярный способ – использовать ***LU-разложение***:

$$W = PLU$$

***P*** – фиксированная матрица перестановок (не обучается)

***L*** – нижнетреугольная матрица

***U*** – верхнетреугольная матрица

# Линейные потоки

- Матрица ***P*** не обучается
- Позднее были работы, которые использовали ***QR*-разложение:**

$$W = QR$$

Для изображений обычно применяется ***обратимая 1 × 1 свертка***

- Это линейный поток, который перемешивает информацию между каналами

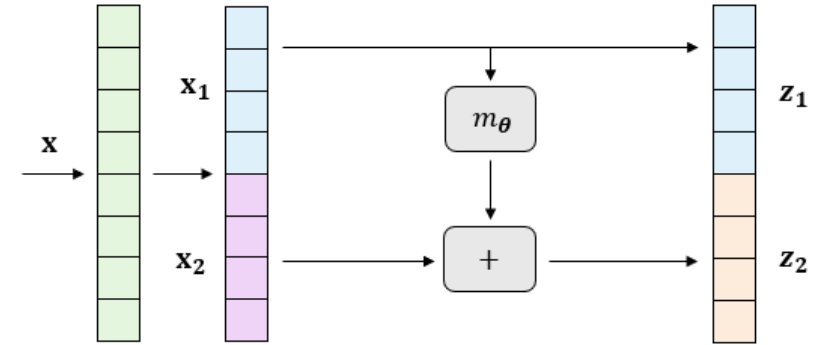
*NICE и RealNVP*

# NICE: Additive Coupling Layers

Разделяем вектор  $\mathbf{x}$  на 2 части  $\mathbf{x}_{1:d}$  и  $\mathbf{x}_{d+1:D}$ :

$$\begin{cases} \mathbf{z}_{1:d} = \mathbf{x}_{1:d} \\ \mathbf{z}_{d+1:D} = \mathbf{x}_{d+1:D} + m_{\theta}(\mathbf{x}_{1:d}) \end{cases}$$

$m_{\theta}$  — произвольная нейросеть



Такая архитектура легко обратима:

$$\begin{cases} \mathbf{x}_{1:d} = \mathbf{z}_{1:d} \\ \mathbf{x}_{d+1:D} = \mathbf{z}_{d+1:D} - m_{\theta}(\mathbf{z}_{1:d}) \end{cases}$$

Якобиан имеет блочный нижнетреугольный вид:

$$J = \begin{pmatrix} \mathbf{I}_d & \mathbf{0} \\ \frac{\partial \mathbf{z}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \mathbf{I}_{D-d} \end{pmatrix}$$

# *NICE*: Проблемы

Определитель якобиана всегда равен единице, поэтому такое преобразование сохраняет объем

Даже композиция из нескольких аддитивных слоёв все равно сохраняет объем

Авторы *NICE* предложили в конце добавить слой масштабирования:

$$\mathbf{z} = \mathbf{s} \odot \mathbf{h}$$

$\mathbf{h}$  — выход последнего связующего слоя

$\mathbf{s}$  — обучаемый вектор масштаба

$\odot$  — поэлементное умножение

Определитель якобиана:

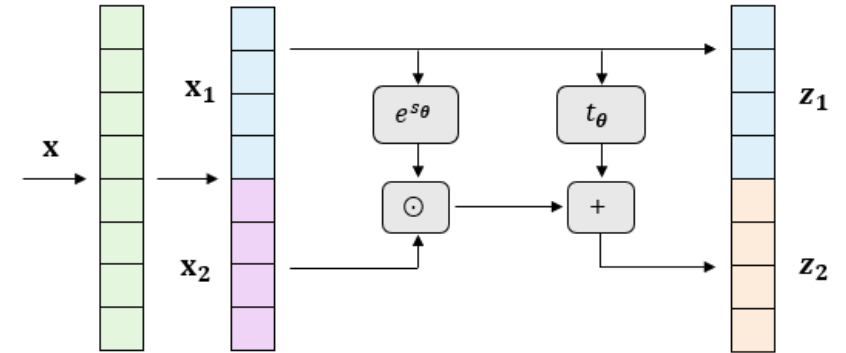
$$\det(\mathbf{J}) = \prod_{i=1}^D s_i$$

# RealNVP: Affine Coupling Layers

Разделяем вектор  $\mathbf{x}$  на 2 части  $\mathbf{x}_{1:d}$  и  $\mathbf{x}_{d+1:D}$ :

$$\begin{cases} \mathbf{z}_{1:d} = \mathbf{x}_{1:d} \\ \mathbf{z}_{d+1:D} = \mathbf{x}_{d+1:D} \odot e^{s_\theta(\mathbf{x}_{1:d})} + t_\theta(\mathbf{x}_{1:d}) \end{cases}$$

$s_\theta$ ,  $t_\theta$  — произвольные нейросети



Такая архитектура всё равно легко обратима:

$$\begin{cases} \mathbf{x}_{1:d} = \mathbf{z}_{1:d} \\ \mathbf{x}_{d+1:D} = (\mathbf{z}_{d+1:D} - t_\theta(\mathbf{z}_{1:d})) \odot e^{-s_\theta(\mathbf{z}_{1:d})} \end{cases}$$

И якобиан имеет блочный нижнетреугольный вид:

$$J = \begin{pmatrix} \mathbf{I}_d & \mathbf{0} \\ \frac{\partial \mathbf{z}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(e^{s_\theta(\mathbf{x}_{1:d})}) \end{pmatrix}$$

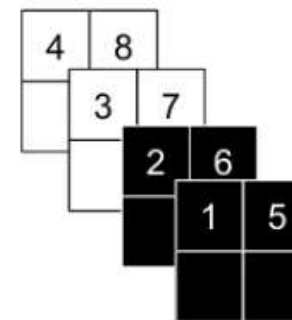


# RealNVP

- Для эффективной работы с изображениями в **RealNVP** используют иерархическую (**Multi – Scale**) архитектуру
- Это позволяет переходить от анализа мелких деталей к более крупным

Типичный блок устроен следующим образом:

- Сначала применяются несколько связующих слоёв с **шахматными масками**
- Затем сжимаем  $h \times w \times c \rightarrow \frac{h}{2} \times \frac{w}{2} \times 4c$
- Применяем связующие слои с **канальными масками**



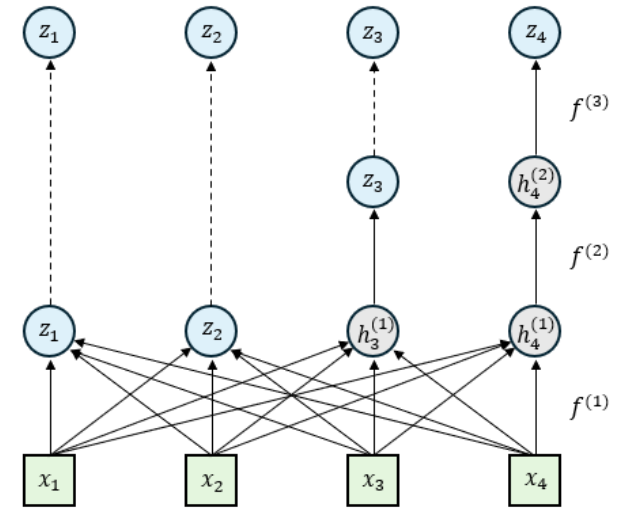
# RealNVP

## Проблема:

Прогонять весь входной вектор через все слои довольно затратно

## Решение:

После каждого блока будем выносить половину переменных из потока сразу в финальный вектор  $\mathbf{z}$



# *RealNVP: Samples*



*Glow*

# *Glow: ActNorm*

**Glow** – идейный наследник **RealNVP**, использующий **Affine Coupling Layers** и **Multi – Scale** архитектуру

Авторы решили заменить эвристики **RealNVP** на более гибкие, обучаемые аналоги

- **BatchNorm** плохо работает с малыми размерами батчей (шумные статистики)
- Чтобы не зависеть от статистик батча, авторы предложили использовать **ActNorm** (**Activation Normalization**) слой

# *Glow: Invertible $1 \times 1$ Conv*

## **Проблема:**

Фиксированные перестановки в *RealNVP* не обучаются и не гарантируют оптимальное смешивание каналов

## **Решение:**

Будем обучать обратимую свёртку  $1 \times 1$  – линейный поток, который позволит модели самой найти оптимальный способ перемешивания каналов

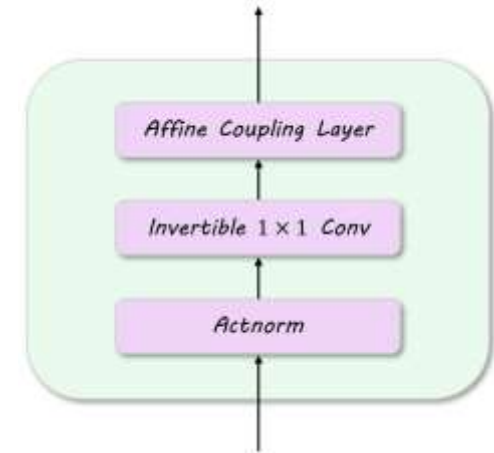
## **Трюк:**

Чтобы быстро считать определитель матрицы свёртки  $W$ , используется *LU-разложение*

# *Glow: Flow Step*

***Glow*** объединяет все 3 операции в один обратимый блок (***Flow Step***), который многократно применяется в модели:

***ActNorm***  $\rightarrow$  ***1  $\times$  1 Conv***  $\rightarrow$  ***Affine Coupling Layer***



***Multi – Scale архитектура:***

- Сжимаем  $h \times w \times c \rightarrow \frac{h}{2} \times \frac{w}{2} \times 4c$
- Делаем  $K$  шагов потока
- Половину каналов выносим сразу в финальный вектор  $\mathbf{z}$

# *Glow: Samples*





Спасибо за внимание!