ВАРИАЦИОННЫЕ АВТОКОДИРОВЩИКИ

Сергей Николенко СПбГУ— Санкт-Петербург 6 ноября 2025 г.





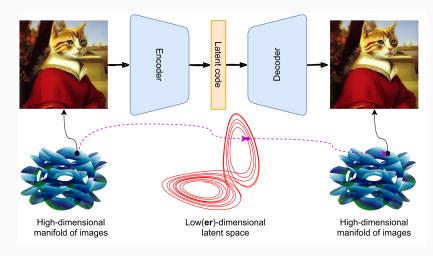
Random facts:

- 6 ноября 447 г. мощное землетрясение разрушило большие участки стен Константинополя, в том числе 57 башен
- 6 ноября 1736 г. Михаил Ломоносов был зачислен в Марбургский университет в Германии
- 6 ноября 1928 г. американский полковник Джейкоб Шик запатентовал электробритву
- 6 ноября 1947 г. на NBC впервые вышла Meet the Press, самая долгоживущая телепрограмма в истории (уже 77 сезонов и почти 5000 выпусков)
- 6 ноября 1957 г. на Марсовом поле в Ленинграде зажгли первый в стране Вечный огонь, а 6 ноября 1968 г. на Кутузовском проспекте у Поклонной горы была открыта Триумфальная арка
- 6 ноября 1975 г. группа «Sex Pistols» дала свой первый концерт в художественной школе Святого Мартина, а 6 ноября 1995 г. был выпущен последний альбом группы «Queen» с участием Фредди Меркьюри

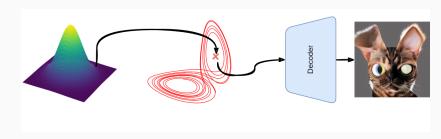
АВТОКОДИРОВЩИКИ

Вариационные

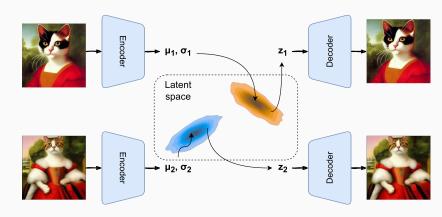
• Главная задача — как всегда: как сделать из обычного автокодировщика порождающую модель?



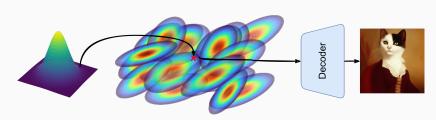
• Даже в пространстве низкой размерности нелегко выбрать хорошую точку:



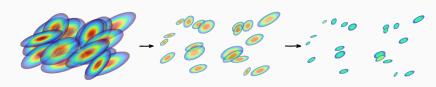
• Главная идея вариационного автокодировщика (Kingma et al., 2013) — пусть из каждой точки $\mathbf x$ получается не точка $\mathbf z$, а целое распределение $p(\mathbf x|\mathbf z)$:



- Интуиция в том, что декодеру приходится быть устойчивым к небольшим изменениям в **z**; некий более идейный аналог шумоподавляющего автокодировщика
- В идеале мы покроем пространство латентных кодов достаточно «широкими» распределениями:

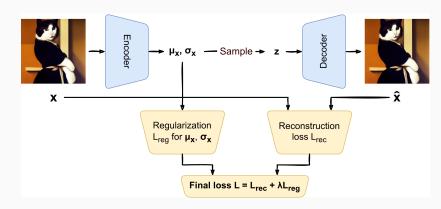


• Однако есть проблема: если просто обучать автокодировщик, ему будет легче декодировать «узкие» распределения, и мы в итоге придём к тем же точкам с нулевой дисперсией

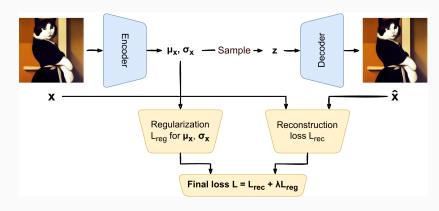


• Как решить эту проблему?

• Нужно добавить некое ограничение на $p\left(\mathbf{x}|\mathbf{z}\right)$, то есть регуляризовать его:

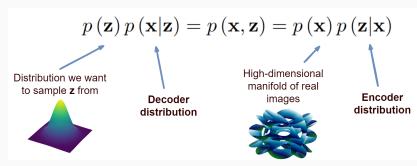


- Теперь остались два вопроса:
 - какие выбрать функцию ошибки и регуляризатор?
 - как протаскивать градиенты через сэмплирование?
- Давайте на оба этих вопроса ответим...

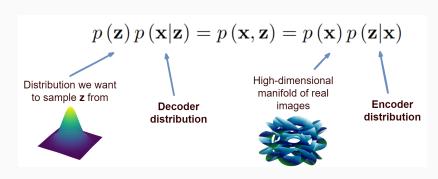


- Давайте начнём с начала мы хотим обучить кодировщик делать латентные коды ${\bf z}$ из картинок ${\bf x}$ и декодировщик, который ${\bf z}$ превращает обратно в ${\bf x}$.
- Совместное распределение можно разложить по-разному:

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



- · Здесь:
 - $\cdot p\left(\mathbf{x}\right)$ это распределение изображений,
 - $\cdot p(\mathbf{z})$ это распределение латентных кодов, т.е. простое распределение, из которого можно сэмплировать;
 - \cdot а остальные два нам нужно найти это распределение кодировщика $p(\mathbf{z}|\mathbf{x})$ и декодировщика $p(\mathbf{x}|\mathbf{z})$.

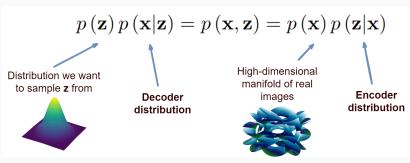


- \cdot Для порождающей модели надо обучить и $p\left(\mathbf{x}|\mathbf{z}\right)$, и $p\left(\mathbf{z}|\mathbf{x}\right)$.
- Но у нас $p(\mathbf{z})$ это простое распределение, а $p(\mathbf{x})$ очень сложное!
- Значит, придётся выбрать одно из двух:
 - · или предположим, что $p\left(\mathbf{x}|\mathbf{z}\right)$ простое, а затем будем искать сложное $p\left(\mathbf{z}|\mathbf{x}\right)$;
 - · или наоборот, предположим, что $p\left(\mathbf{z}|\mathbf{x}\right)$ простое, а затем будем искать сложное $p\left(\mathbf{x}|\mathbf{z}\right)$.

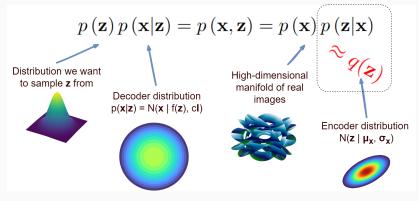
• В VAE удобнее использовать первый вариант: предположим, что

$$p\left(\mathbf{x}|\mathbf{z}\right) = N\left(\mathbf{x}|f(\mathbf{z}), c\mathbf{I}\right),$$
 где $f(\mathbf{z}) = \mathrm{Decoder}(\mathbf{z})$

• Слева получается $p\left(\mathbf{z}|\mathbf{x}\right)$ умножить на $p\left(\mathbf{z}\right)=N\left(\mathbf{z}|0,\mathbf{I}\right)$, то есть опять гауссиан.



· А справа будем делать приближение $q(\mathbf{z}) \approx p\left(\mathbf{z}|\mathbf{x}\right)$:



• А именно вариационное приближение...

• Вспомним идею вариационных приближений:

$$\begin{split} p\left(\mathbf{x}, \mathbf{z}\right) &= p\left(\mathbf{x}\right) p\left(\mathbf{z} | \mathbf{x}\right), \\ \log p\left(\mathbf{x}\right) &= \log p\left(\mathbf{x}, \mathbf{z}\right) - \log p\left(\mathbf{z} | \mathbf{x}\right), \end{split}$$

возьмём ожидание по $q\left(\mathbf{z}\right)$:

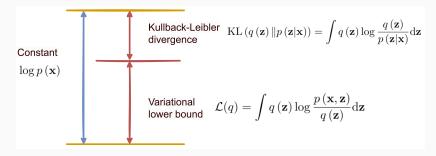
$$\mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{x}\right)\right] = \mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{x},\mathbf{z}\right)\right] - \mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{z}|\mathbf{x}\right)\right],$$

а затем сделаем стандартные преобразования:

$$\begin{split} \log p\left(\mathbf{x}\right) = & \mathbb{E}_{q(\mathbf{z})} \left[\log p\left(\mathbf{x}, \mathbf{z}\right) \right] - \mathbb{E}_{q(\mathbf{z})} \left[\log q\left(\mathbf{z}\right) \right] + \\ & + \mathbb{E}_{q(\mathbf{z})} \left[\log q\left(\mathbf{z}\right) \right] - \mathbb{E}_{q(\mathbf{z})} \left[\log p\left(\mathbf{z}|\mathbf{x}\right) \right], \\ \log p\left(\mathbf{x}\right) = & \int q\left(\mathbf{z}\right) \log \frac{p\left(\mathbf{x}, \mathbf{z}\right)}{q\left(\mathbf{z}\right)} \mathrm{d}\mathbf{z} + \int q\left(\mathbf{z}\right) \log \frac{q\left(\mathbf{z}\right)}{p\left(\mathbf{z}|\mathbf{x}\right)} \mathrm{d}\mathbf{z}. \end{split}$$

• Итого получается, что

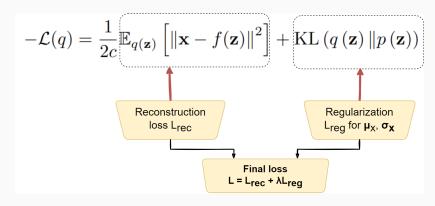
$$\log p\left(\mathbf{x}
ight) = L(q) + \mathrm{KL}\left(q\left(z
ight) \| p\left(\mathbf{z} | \mathbf{x}
ight)
ight),$$
 где $L(q) = \int q\left(\mathbf{z}
ight) \log rac{p\left(\mathbf{x}, \mathbf{z}
ight)}{q\left(\mathbf{z}
ight)} \mathrm{d}\mathbf{z}.$



• Значит, можно аппроксимировать $p\left(\mathbf{z}|\mathbf{x}\right)$ через $q\left(\mathbf{z}\right)$, максимизируя L(q). В случае VAE его можно найти:

$$\begin{split} L(q) &= \int q\left(\mathbf{z}\right) \log \frac{p\left(\mathbf{x}, \mathbf{z}\right)}{q\left(\mathbf{z}\right)} \mathrm{d}\mathbf{z} = \int q\left(\mathbf{z}\right) \log \frac{p\left(\mathbf{z}\right) p\left(\mathbf{x} | \mathbf{z}\right)}{q\left(\mathbf{z}\right)} \mathrm{d}\mathbf{z} \\ &= \int q\left(\mathbf{z}\right) \log p\left(\mathbf{x} | \mathbf{z}\right) \mathrm{d}\mathbf{z} + \int q\left(\mathbf{z}\right) \log \frac{p\left(\mathbf{z}\right)}{q\left(\mathbf{z}\right)} \mathrm{d}\mathbf{z} \\ &= \int q\left(\mathbf{z}\right) \log N\left(\mathbf{x} | f(\mathbf{z}), c\mathbf{I}\right) \mathrm{d}\mathbf{z} - \mathrm{KL}\left(q\left(\mathbf{z}\right) \| p\left(\mathbf{z}\right)\right) \\ &= -\frac{1}{2c} \mathbb{E}_{q(\mathbf{z})} \left[\left\|\mathbf{x} - f(\mathbf{z})\right\|^2 \right] - \mathrm{KL}\left(q\left(\mathbf{z}\right) \| p\left(\mathbf{z}\right)\right). \end{split}$$

• И вот получилось как раз то, чего мы ожидали интуитивно:



 \cdot Только теперь мы знаем, как выглядит $L_{
m rec}$, да и $\,$ КL-дивергенцию между двумя гауссианами можно подсчитать:

$$\mathrm{KL}\left(q\left(\mathbf{z}\right)\|p\left(\mathbf{z}\right)\right) = \frac{1}{2}\sum_{j=1}^{d}\left(\sigma_{\mathbf{x},j}^{2} + \mu_{\mathbf{x},j}^{2} - \log\sigma_{\mathbf{x},j}^{2} - 1\right).$$

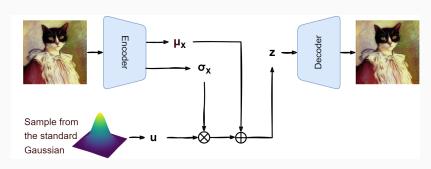
• Получается окончательная функция ошибки для вариационного автокодировщика:

$$L(q) = -\frac{1}{2c}\mathbb{E}_{q(\mathbf{z})}\left[\left\|\mathbf{x} - f(\mathbf{z})\right\|^2\right] - \frac{1}{2}\sum_{j=1}^d\left(\sigma_{\mathbf{x},j}^2 + \mu_{\mathbf{x},j}^2 - \log\sigma_{\mathbf{x},j}^2 - 1\right).$$

• Осталась только проблема с сэмплированием...

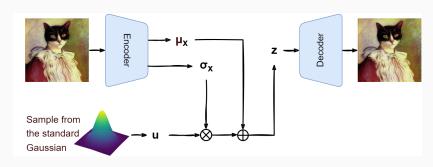
REPARAMETRIZATION TRICK

- Чтобы обучить модель, нужно пропустить градиент через **z** обратно к кодировщику, но там сэмплирование
- Reparametrization trick: давайте сначала сэмплировать, а потом преобразовывать результат!



REPARAMETRIZATION TRICK

• И теперь уже точно все компоненты на месте: сэмплируем мини-батч ${f u}$ для мини-батча ${f x}$ и используем их для обучения



Спасибо!

Спасибо за внимание!



