ПОТОКИ В ГЛУБОКОМ ОБУЧЕНИИ

Сергей Николенко СПбГУ— Санкт-Петербург 27 ноября 2025 г.

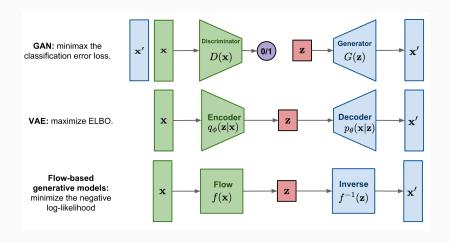




Random facts:

- 27 ноября 1095 г. Урбан II на Клермонском соборе по просьбе византийского императора Алексея I провозгласил Первый крестовый поход
- 27 ноября 1838 г. произошла битва при Сан-Хуан-де-Улуа в ходе Кондитерской войны между Мексикой и Францией; война началась по заявлению французского кондитера, якобы ограбленного мексиканскими мародёрами, Франция послала флот, чтобы вернуть долг, при бомбардировке Сан-Хуан-де-Улуа погибли более 60 человек, и президенту Мексики пришлось пообещать выплатить компенсацию Франции; впрочем, в итоге так и не выплатили
- 27 ноября 1895 г. Альфред Нобель подписал завещание, по которому большая часть его состояния поступала в фонд Нобелевской премии
- 27 ноября— день авиакатастроф: в 1962 г. Boeing 707 врезался в гору при заходе на посадку под Лимой (97 погибших), в 1970 Douglas DC-8 выкатился со взлётной полосы и загорелся (47 погибших), в 1983 Boeing 747 под Мадридом зацепил несколько холмов и разрушился (181 погибших), а в 1989 г. в окрестностях Боготы террористы по указанию Пабло Эскобара взорвали Boeing 727 (110 погибших)
- 27 ноября 1992 г. была создана Высшая школа экономики

Идея



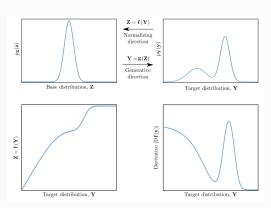
- И ещё один подход к порождающим моделям, совсем другой, но тоже очень интересный
- Нормализующие потоки (normalizing flows) это метод построения сложных распределений из маленьких кусочков, выражаемых обратимыми преобразованиями (Rezende, Mohamed, 2015)
- \cdot Если $\mathbf{z}\in\mathbb{R}^d$, $\mathbf{z}\sim q(\mathbf{z})$, и $f:\mathbb{R}^d\to\mathbb{R}^d$ обратимая функция, то у случайной величины $y=f(\mathbf{z})$ будет плотность

$$q_y(\mathbf{z}) = q(\mathbf{z}) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{z}} \right| = q(\mathbf{z}) \left| \det \frac{\partial f}{\partial \mathbf{z}} \right|^{-1}.$$

4

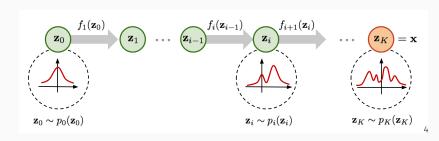
• Если $\mathbf{z}\in\mathbb{R}^d$, $\mathbf{z}\sim q(\mathbf{z})$, и $f:\mathbb{R}^d\to\mathbb{R}^d$ – обратимая функция, то у случайной величины $y=f(\mathbf{z})$ будет плотность

$$q_y(\mathbf{z}) = q(\mathbf{z}) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{z}} \right| = q(\mathbf{z}) \left| \det \frac{\partial f}{\partial \mathbf{z}} \right|^{-1}.$$

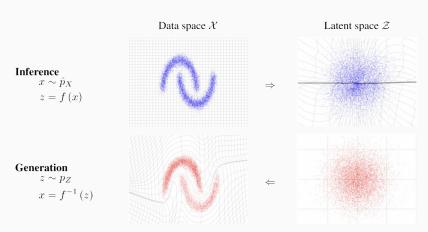


• И если сразу несколько взять таких функций подряд, то в общем тоже несложно посчитать, что получится:

$$\begin{split} \mathbf{z}_K &= f_K \circ \dots \circ f_1(\mathbf{z}_0), \quad \mathbf{z}_0 \sim q_0(\mathbf{z}_0), \\ \mathbf{z}_K \sim q_K(\mathbf{z}) &= q_0(\mathbf{z}_0) \prod_{k=1}^K \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|^{-1}, \\ \log q_K(\mathbf{z}) &= \log q_0(\mathbf{z}_0) - \sum_{k=1}^K \log \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|. \end{split}$$



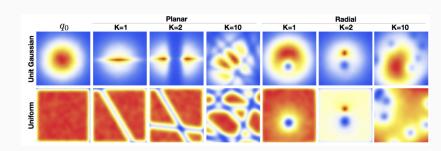
• Так можно простое распределение (гауссиан, например) превратить в очень сложное:



• Но мы ограничены преобразованиями, у которых легко подсчитать якобианы. Скоро увидим, насколько это серьёзное ограничение.

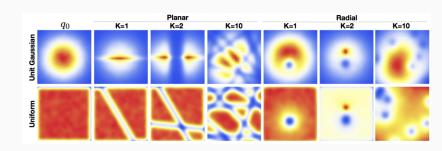
 Например, планарный поток (линейный) «разрезает» пространство гиперплоскостями:

$$\begin{split} f(\mathbf{z}) &= \mathbf{z} + \mathbf{u} h(\mathbf{w}^{\top} \mathbf{z} + b), \\ \left| \det \frac{\partial f}{\partial \mathbf{z}} \right| &= \left| 1 + \mathbf{u}^{\top} \left(h'(\mathbf{w}^{\top} \mathbf{z} + b) \mathbf{w} \right) \right|. \end{split}$$



· Радиальный поток «разрезает» пространство сферами:

$$\begin{split} f(\mathbf{z}) &= \mathbf{z} + \beta h(\alpha, r) \left(\mathbf{z} - \mu \right), \\ \text{где } r &= \left\| \mathbf{z} - \mu \right\|_2, \ h(\alpha, r) = \frac{1}{\alpha + r}. \end{split}$$



- А как сделать что-то поумнее? Нам надо уметь считать якобиан, т.е. определитель $\frac{\partial f}{\partial x}$.
- Но если матрица треугольная, то определитель от неё несложно посчитать!
- Это в точности идея авторегрессионных потоков (autoregressive flows):

$$y_i = f(z_1, \dots, z_i), \quad \det \frac{\partial f}{\partial \mathbf{z}} = \prod_{i=1}^d \frac{\partial y_i}{\partial z_i}.$$

4

• Real Non-Volume Preserving Flows (R-NVP) – не слишком выразительные, но очень эффективные: для некоторых функция $\mu,\sigma:\mathbb{R}^k \to \mathbb{R}^{d-k}$ определим

$$\mathbf{y}_{1:k} = \mathbf{z}_{1:k}, \quad \mathbf{y}_{k+1:d} = \mathbf{z}_{k+1:d} \cdot \sigma(\mathbf{z}_{1:k}) + \mu(\mathbf{z}_{1:k}).$$

• Т.е. копируем первые k размерностей, потом для остальных применяем линейное преобразование; получается сначала единичная матрица, потом нижнетреугольная:

$$\det \frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \prod_{i=1}^{d-k} \sigma_i(\mathbf{z}_{1:k}).$$

• Очень легко посчитать (всё параллельно), легко обратить даже для необратимых μ и σ , и можно использовать как параметризацию для VAE, а также как просто правдоподобие.

• Можно обобщить:

$$y_1 = \mu_1 + \sigma_1 z_1, \quad y_i = \mu(\mathbf{y}_{1:i-1}) + \sigma(\mathbf{y}_{1:i-1}) z_i.$$

- Более выразительная модель, якобиан такой же простой, $\prod_i \sigma(\mathbf{y}_{1:i-1})$, тоже можно легко обратить.
- Но теперь параллельно не вычисляется (по крайней мере в общем виде), преобразование по сути последовательное.

Masked Autoregressive Flow (MAF):

$$\begin{split} p\left(\mathbf{x}\right) &= \prod_{i=1}^{D} p\left(x_{i} | \mathbf{x}_{1:i-1}\right), \\ x_{i} &\sim p\left(x_{i} | \mathbf{x}_{1:i-1}\right) = z_{i} \odot \sigma_{i}\left(\mathbf{x}_{1:i-1}\right) + \mu_{i}\left(\mathbf{x}_{1:i-1}\right). \end{split}$$

- Функции μ и σ могут быть абсолютно произвольными! Например, параметризованными сложными нейросетями.
- Быстро работает для оценки плотности (обучения), но для сэмплирования (порождения) медленно, то есть авторегрессивно.

· Inverse Autoregressive Flow (IAF) работает по сути наоборот:

$$\tilde{x}_i \sim p\left(\tilde{x}_i \middle| \mathbf{z}_{1:i}\right) = z_i \odot \sigma_i(\mathbf{z}_{1:i-1}) + \mu_i(\mathbf{z}_{1:i-1}).$$

• Теперь ${f y}$ зависит только от ${f z}$, и можно вычислить быстро (параллельно):

$$\tilde{\mathbf{x}} = \mathbf{z} \cdot \sigma(\mathbf{z}) + \mu(\mathbf{z}), \quad \det \frac{\partial \tilde{\mathbf{x}}}{\partial \mathbf{z}} = \prod_{i=1}^{d} \sigma_i(\mathbf{z}),$$

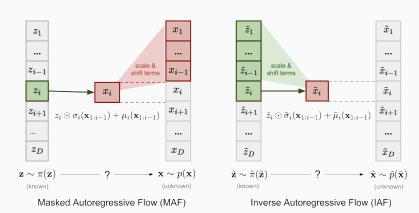
и если собрать целую последовательность $f_K \circ ... \circ f_0(\mathbf{z})$, то

$$\log q_K(\mathbf{z}_K) = \log q(\mathbf{z}) - \sum_{k=0}^K \sum_{i=1}^d \log \sigma_{k,i}.$$

• Теперь сэмплировать легко, а плотность считать (обучать) сложно, нужно авторегрессивно.

4

• Получается такая конструкция:



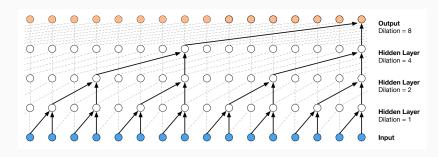
• В качестве всех этих μ и σ можно и нужно использовать нейронные сети.

- Haпpumep, PixelRNN и MADE это в точности нормализующие потоки.
- То есть МАГ и ІАГ эффективны в разных направлениях:

	Base distribution	Target distribution	Model	Data generation	Density estimation
MAF	$\mathbf{z} \sim \pi(\mathbf{z})$	$\mathbf{x} \sim p(\mathbf{x})$	$x_i = z_i \odot \sigma_i(\mathbf{x}_{1:i-1}) + \mu_i(\mathbf{x}_{1:i-1})$	Sequential; slow	One pass; fast
IAF	$ ilde{\mathbf{z}} \sim ilde{\pi}(ilde{\mathbf{z}})$	$ ilde{\mathbf{x}} \sim ilde{p}(ilde{\mathbf{x}})$	$ ilde{x}_i = ilde{z}_i \odot ilde{\sigma}_i(ilde{\mathbf{z}}_{1:i-1}) + ilde{\mu}_i(ilde{\mathbf{z}}_{1:i-1})$	One pass; fast	Sequential; slow

• Как это использовать?

• Следующий интересный шаг – Parallel WaveNet (van Oord, 2017), который ввёл метод параллельной дистилляции (parallel density distillation).



- WaveNet это как раз модель, которая быстро (параллельно) обучается, но медленно (последовательно) сэмплирует.
- Чтобы быстро сэмплировать, было бы хорошо перейти к IAF:

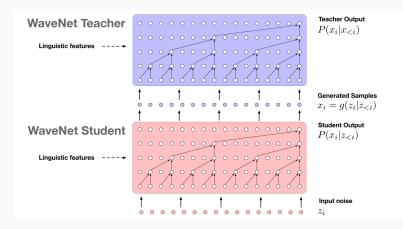
$$\log p_X(\mathbf{x}) = \log p_Z(\mathbf{z}) - \log \left| \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right|, \quad \log \left| \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right| = \sum_t \log \frac{\partial f(\mathbf{z}_{\leq t})}{\partial z_t}.$$

- Тогда сэмплирование – это просто $\mathbf{z} \sim q(\mathbf{z})$ (у WaveNet – логистическое), а дальше

$$x_t = z_t \cdot s(\mathbf{z}_{\leq t}, \theta) + \mu(\mathbf{z}_{\leq t}, \theta),$$

и для s и μ можно использовать те же свёрточные сети, что в WaveNet; можно в принципе одну, но оказывается получше, если композицию из четырёх сделать.

• А теперь главный трюк: обучать IAF очень долго, для WaveNet совсем непрактично. Поэтому parallel density distillation – давайте обучать студента $p_S(\mathbf{x})$ на основе учителя $p_T(\mathbf{x})$ с функцией ошибки $\mathrm{KL}\,(P_S\|P_T) = H(P_S,P_T) - H(P_S)$:



• Теперь и логистические распределения пригождаются:

$$\begin{split} H(P_S) &= \mathbb{E}_{\mathbf{z} \sim L(0,1)} \left[\sum_{t=1}^T -\ln p_S(x_t \mid \mathbf{z}_{< t}) \right] = \\ &= \mathbb{E}_{\mathbf{z} \sim L(0,1)} \left[\sum_{t=1}^T \ln s(\mathbf{z}_{< t}, \theta) \right] + 2T, \end{split}$$

потому что энтропия $L(\mu,s)$ равна $\ln s + 2$

• А кросс-энтропия тоже нехитро:

$$\begin{split} H(P_S, P_T) &= \int_{\mathbf{x}} p_S(\mathbf{x}) \ln p_T(\mathbf{x}) \mathrm{d}\mathbf{x} = \ldots = \\ &= \sum_{t=1}^T \mathbb{E}_{p_S(\mathbf{x}_{< t})} \left[H\left(p_S(x_t \mid \mathbf{x}_{< t}), p_T(x_t \mid \mathbf{x}_{< t}) \right) \right]. \end{split}$$

5

- Т.е. можно порождать из student \mathbf{x} , вычислять $p_T(x_t \mid \mathbf{x}_{< t})$ параллельно в teacher, а потом эффективно сэмплировать несколько x_t из $p_S(x_t \mid \mathbf{x}_{< t})$ для каждого t.
- Осталось добавить несколько дополнительных loss functions для student-сети (power loss для уровня мощности, perceptual loss на WaveNet-подобном классификаторе, contrastive loss при разных условиях).

• И получается такая же модель по качеству:

Method	Subjective 5-scale MOS			
16kHz, 8-bit μ-law, 25h data:				
LSTM-RNN parametric [27]	3.67 ± 0.098			
HMM-driven concatenative [27]	3.86 ± 0.137			
WaveNet [27]	4.21 ± 0.081			
24kHz, 16-bit linear PCM, 65h data:				
HMM-driven concatenative	4.19 ± 0.097			
Autoregressive WaveNet	4.41 ± 0.069			
Distilled WaveNet	4.41 ± 0.078			

• Но авторегрессионный WaveNet сэмплирует 172 отсчёта в секунду, а дистиллированный – >500000 отсчётов в секунду на том же железе...

• Следующий шаг – давайте обобщать:

$$\mathbf{z} \sim p_{\theta}(\mathbf{z}), \quad \mathbf{x} = g_{\theta}(\mathbf{z}), \quad \mathbf{z} = f_{\theta}(\mathbf{x}) = g_{\theta}^{-1}(\mathbf{x}).$$

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

где $\mathbf{x} = \mathbf{h}_0$, $\mathbf{z} = \mathbf{h}_K$, и остальные в промежутке

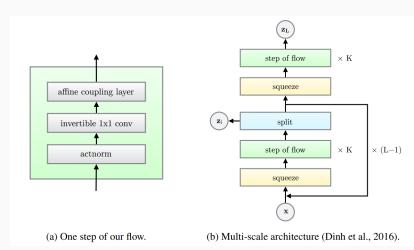
• Даже ограничимся преобразованиями, где якобиан – треугольная матрица.

6

• Но всё равно разнообразие довольно большое можно придумать:

Description	Function	Reverse Function	Log-determinant
Actnorm. See Section 3.1.	$\forall i,j: \mathbf{y}_{i,j} = \mathbf{s} \odot \mathbf{x}_{i,j} + \mathbf{b}$	$\mid \forall i, j : \mathbf{x}_{i,j} = (\mathbf{y}_{i,j} - \mathbf{b})/\mathbf{s}$	$ \mid h \cdot w \cdot \mathtt{sum}(\log \mathbf{s}) $
$\begin{aligned} &\text{Invertible } 1\times 1 \text{ convolution.} \\ &\mathbf{W}: [c\times c]. \\ &\text{See Section } 3.2. \end{aligned}$	$\forall i,j: \mathbf{y}_{i,j} = \mathbf{W}\mathbf{x}_{i,j}$	$\mid \forall i, j : \mathbf{x}_{i,j} = \mathbf{W}^{-1} \mathbf{y}_{i,j}$	$ \begin{aligned} h \cdot w \cdot \log \det(\mathbf{W}) \\ \text{or} \\ h \cdot w \cdot \text{sum}(\log \mathbf{s}) \\ \text{(see eq. (10))} \end{aligned} $
Affine coupling layer. See Section 3.3 and (Dinh et al., 2014)	$\begin{aligned} &\mathbf{x}_a, \mathbf{x}_b = \mathtt{split}(\mathbf{x}) \\ &(\log \mathbf{s}, \mathbf{t}) = \mathtt{NN}(\mathbf{x}_b) \\ &\mathbf{s} = \mathtt{exp}(\log \mathbf{s}) \\ &\mathbf{y}_a = \mathbf{s} \odot \mathbf{x}_a + \mathbf{t} \\ &\mathbf{y}_b = \mathbf{x}_b \\ &\mathbf{y} = \mathtt{concat}(\mathbf{y}_a, \mathbf{y}_b) \end{aligned}$	$ \begin{aligned} &\mathbf{y}_a, \mathbf{y}_b = \mathtt{split}(\mathbf{y}) \\ &(\log \mathbf{s}, \mathbf{t}) = \mathtt{NN}(\mathbf{y}_b) \\ &\mathbf{s} = \exp(\log \mathbf{s}) \\ &\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t})/\mathbf{s} \\ &\mathbf{x}_b = \mathbf{y}_b \\ &\mathbf{x} = \mathtt{concat}(\mathbf{x}_a, \mathbf{x}_b) \end{aligned} $	$\mathtt{sum}(\log(\mathbf{s}))$

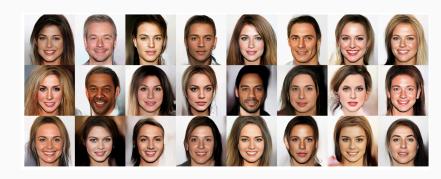
· И из этих компонентов составим модель; получится GLOW (Kingma, Dhariwal, 2018):



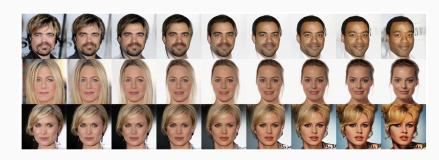
• И всё, это уже можно обучать делать неплохие, например, лица людей:



• Вполне разнообразно:



• По пространству признаков прогулки хорошие:



WAVEGLOW

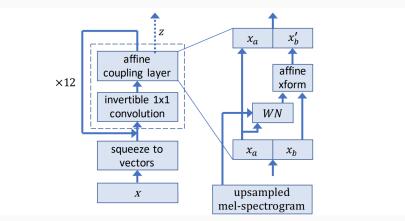
· А если совместить WaveNet и GLOW, получится



• WaveGlow генерирует речь из мел-спектрограммы (т.е. надо ещё отдельно её породить, но это нормально)

WAVEGLOW

• Смысл в целом тот же, но идея в том, что можно сделать из необратимого WN (типа WaveNet) обратимое преобразование:



• И можно добраться до 2M отсчётов/сек, около 100x real-time

- Flow++ (Ho et al., 2019) несколько улучшений для потоковых моделей:
 - предыдущие модели используют деквантизацию, т.е. добавляют шум к дискретным данным, чтобы p_{model} не коллапсировала в смесь точечных распределений; тогда получается, что мы моделируем

$$P_{\mathrm{model}}\left(\mathbf{x}\right) = \int_{[0,1]^{D}} p_{\mathrm{model}}\left(\mathbf{x} + \mathbf{u}\right) \mathrm{d}\mathbf{u};$$

- но равномерный шум не очень хорош для деквантизации: он просит нейросеть обучить равномерное распределение в кубиках вокруг набора точек, очень странно;
- · Flow++ предлагает использовать

- Flow++ (Ho et al., 2019) несколько улучшений для потоковых моделей:
 - · Flow++ предлагает добавить распределение шума $q\left(\mathbf{u} \mid \mathbf{x}\right)$ и оптимизировать вариационную оценку:

$$\begin{split} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log P_{\text{model}} \left(\mathbf{x} \right) \right] &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \int_{[0,1]^D} p_{\text{model}} \left(\mathbf{x} + \mathbf{u} \right) \mathrm{d} \mathbf{u} \right] \geq \\ &\geq \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\int_{[0,1]^D} q \left(\mathbf{u} \mid \mathbf{x} \right) \log \frac{p_{\text{model}} \left(\mathbf{x} + \mathbf{u} \right)}{q \left(\mathbf{u} \mid \mathbf{x} \right)} \mathrm{d} \mathbf{u} \right] = \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}, \mathbf{u} \sim q(\mathbf{u} \mid \mathbf{x})} \left[\log \frac{p_{\text{model}} \left(\mathbf{x} + \mathbf{u} \right)}{q \left(\mathbf{u} \mid \mathbf{x} \right)} \right]; \end{split}$$

• и теперь $q\left(\mathbf{u}\mid\mathbf{x}\right)$ можно взять тоже как потоковую порождающую модель $\mathbf{u}=q_{\mathbf{x}}(\epsilon)$, $\epsilon\sim N\left(0,\mathbf{I}\right)$;

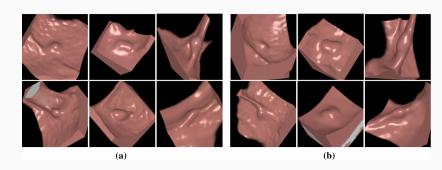
- Flow++ (Ho et al., 2019) несколько улучшений для потоковых моделей:
 - теперь

$$\begin{split} q\left(\mathbf{u}\mid\mathbf{x}\right) &= p\left(q_{\mathbf{x}}^{-1}(\mathbf{u})\right) \left|\frac{\partial q_{\mathbf{x}}^{-1}}{\partial \mathbf{u}}\right|, \quad \mathbf{M} \\ \mathbb{E}_{\mathbf{x} \sim p_{\mathrm{data}}}\left[\log P_{\mathrm{model}}\left(\mathbf{x}\right)\right] &\geq \mathbb{E}_{\mathbf{x} \sim p_{\mathrm{data}}, \epsilon}\left[\log \frac{p_{\mathrm{model}}\left(\mathbf{x} + \mathbf{u}\right)}{p\left(\epsilon\right) \left|\frac{\partial q_{\mathbf{x}}}{\partial \epsilon}\right|^{-1}}\right], \end{split}$$

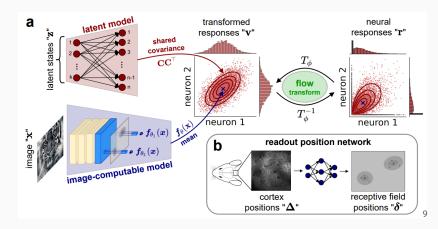
и это теперь можно максимизировать совместно по p_{model} и q.

- Flow++ (Ho et al., 2019) несколько улучшений для потоковых моделей:
 - ещё улучшения для coupling layers;
 - улучшения для архитектуры нейросети, вычисляющей параметры этого преобразования (здесь трансформеры появились);
 - \cdot и в итоге лучше предыдущих работало на 64×64 картинках, но главное идеи, они потом пошли дальше в потоковые модели.

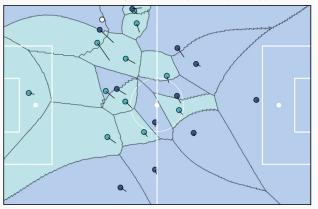
- Потоковые модели активно используют для задач, где датасеты поменьше.
- · (Uemura et al., 2020): 3D GLOW для моделирования полипов



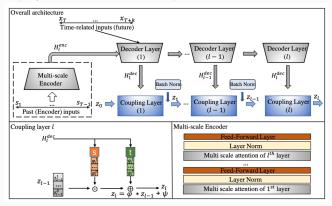
- Потоковые модели активно используют для задач, где датасеты поменьше.
- (Bashiri et al., 2021): классификация активностей нейронов (настоящих, биологических)



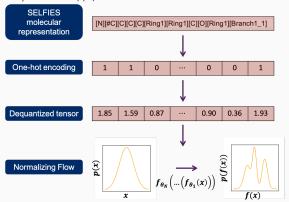
- Потоковые модели активно используют для задач, где датасеты поменьше.
- · (Fadel et al., 2021): моделирование траекторий футболистов



- Потоковые модели активно используют для задач, где датасеты поменьше.
- (Feng et al., 2022): потоки для моделирования временных рядов (ещё multi-scale attention)



- Потоковые модели активно используют для задач, где датасеты поменьше.
- (Frey et al., 2022): потоки для порождения молекул (SMILES fingerprints, как всегда)



Спасибо!

Спасибо за внимание!



