

Normalizing Flows (wtd.)

defined coupling layers:

- fast forward computation:

$$z_{\leq \tilde{D}}^{(e)} = z_{\leq \tilde{D}}^{(e-1)}$$

element-wise
↓

$$z_{> \tilde{D}}^{(e)} = s_e(z_{\leq \tilde{D}}^{(e-1)}) \cdot z_{> \tilde{D}}^{(e-1)} + t_e(z_{> \tilde{D}}^{(e-1)}) \quad \tilde{D} = \left\lfloor \frac{D}{2} \right\rfloor$$

needed for training & for inference ($\hat{=}$ compute $\hat{p}(x=x)$)

- fast backward computation:

$$z_{\leq \tilde{D}}^{(e-1)} = z_{\leq \tilde{D}}^{(e)}$$

$$z_{> \tilde{D}}^{(e-1)} = (z_{> \tilde{D}}^{(e)} - t_e(z_{\leq \tilde{D}}^{(e)})) / s_e(z_{\leq \tilde{D}}^{(e)})$$

needed for data generation

important: the nested functions $s_e(\cdot)$ and $t_e(\cdot)$ are always executed forward

\Rightarrow constructed an invertible layer, using non-invertible networks $s_e(\cdot), t_e(\cdot)$

- cheap Jacobian determinant

$$\det(f'(z^{(e-1)})) = \prod_{j=\tilde{D}+1}^D s_{e,j}(z_{\leq \tilde{D}}^{(e-1)})$$

needed for training & inference

- to build a deep network from coupling layers, we must ensure that the ship connections are not always the same dimensions \Rightarrow introduce orthonormal matrices Q_e between layers

$$f(x) = f_L \circ Q_{L-1} \circ f_{L-1} \circ Q_{L-2} \circ \dots \circ Q_1 \circ f_1$$

orthonormal is nice, because:

- easily inverted $Q^{-1} = Q^T$
- jacobian determinant $|\det Q| = 1$

two basic choices:

- $Q = \text{random-shuffle-columns}(\mathbb{I}) \rightarrow \text{permutation matrix}$
- Variants: for d odd: $Q = \begin{pmatrix} 0 & \mathbb{I}_{\tilde{D}} \\ \mathbb{I}_{\tilde{D}-D} & 0 \end{pmatrix}$ exchanges first & second half of $\mathbb{R}^{(e-1)}$

\Rightarrow every dimension has been transformed after two couplings

- Q is a random rotation

$A = \text{random-normal}(D, D)$

QR decomposition $A = QR$, use Q

Learning of Q
does not improve
models in practice

for coupling layers to be invertible, $s_{ej}(z_{\leq \tilde{D}}^{(l-1)}) \neq 0$
without loss of generality, we require $s_{ej}(z_{\leq \tilde{D}}^{(l-1)}) > 0$

\Rightarrow actually learn $\tilde{s}_{ej}(z_{\leq \tilde{D}}^{(l-1)})$ unconstrained, $s_{ej}(z_{\leq \tilde{D}}^{(l-1)}) = \exp(\tilde{s}_{ej}(z_{\leq \tilde{D}}^{(l-1)}))$

in practice, for numerical stability, we constrain $\frac{1}{\epsilon} \leq s < \epsilon$

for example

$$s_{ej}(z_{\leq \tilde{D}}^{(l-1)}) = \exp(a \cdot \tanh(\tilde{s}_{ej}(z_{\leq \tilde{D}}^{(l-1)})))$$

$$\Rightarrow \log \det f'_e(z^{(l-1)}) = \sum_{j=\tilde{D}+1}^D \tilde{s}_{ej}(z_{\leq \tilde{D}}^{(l-1)})$$

training of a NF: we want to minimize forward KL: $KL[p^*(x) \parallel \hat{p}(x)] \approx 0$

$$KL[p^*(x) \parallel p(x)] = \int p^*(x) \log \frac{p^*(x)}{p(x)} dx = \int p^*(x) (-\log p(x)) dx - \underbrace{\int p^*(x) (-\log p^*(x)) dx}_{H(p^*) \text{ entropy of data, independent of model} \Rightarrow \text{drop}}$$

minimal NLL loss: $= \mathbb{E}_{x \sim p^*(x)} [-\log p(x)] + \text{const}$

$$\approx \frac{1}{N} \sum_{i=1}^N -\log p(x_i) \quad \begin{matrix} \text{change of} \\ \text{variables} \end{matrix} \quad \frac{1}{N} \sum_{i=1}^N -\log \underbrace{q(f(x_i))}_{\text{desired code distribution } = \mathcal{N}(0, I)} - \log \det f'(x_i)$$

if $q(z) = \mathcal{N}(0, I) \Rightarrow -\log q(f(x_i)) = \underbrace{-\log \text{normalization}}_{\text{const.} \Rightarrow \text{drop}} + \frac{f(x_i)^2}{2}$

if $f(x_i)$ is a coupling network: $\log \det f'(x_i) = \sum_{l=1}^L \sum_{j=\tilde{n}+1}^D \tilde{S}_{ej} \left(z_{\leq \tilde{n}}^{(l-1)} \right)$

training objective of NF with $q(z) = N(0, I)$

$$\hat{f}(x) = \arg \min_f \frac{1}{N} \sum_{i=1}^N \left(\frac{\|f(x_i)\|^2}{2} - \sum_{l=1}^L \sum_{j=\tilde{D}+1}^D \tilde{s}_{ej}(\tilde{z}_{i, \leq \tilde{D}}^{(l-1)}) \right)$$

alg: ① initialize s_e, t_e randomly as in ordinary neural networks, init Q_e randomly (not trainable)

② for $t=1, \dots, T$

do a gradient step of loss to update s_e, t_e (ordinary back-prop or ADAM)

known as "Real NVP" network (NVP = non-volume-preserving = $\det f'(x) \neq 1$)

[Dinh et al. 2017]

implemented in various libraries, e.g. FrEIA, BayesFlow [org]

universal, i.e.

[Draxler et al. 2022-23]

other

• simplification: volume-preserving NFs: $\det f'(x) = 1$ for all x

- NICE architecture [Dinh et al. 2015]

$$z_{>\tilde{n}}^{(e)} = z_{>\tilde{n}}^{(e-1)} + t_e (z_{\leq \tilde{n}}^{(e-1)})$$

$$\Rightarrow f'_e(z^{(e-1)}) = \mathbb{I}, \det f'_e(z^{(e-1)}) = 1$$

- GIN architecture [Sorensen et al. 2020]
"General incompressible NF"

$$z_{>\tilde{n}}^{(e)} = \exp(\tilde{s}_e(z_{\leq \tilde{n}}^{(e-1)})) \cdot z_{>\tilde{n}}^{(e-1)} + t_e (z_{\leq \tilde{n}}^{(e-1)})$$

$$\tilde{s}_e(z_{\leq \tilde{n}}^{(e-1)}) = \tilde{s}_e(z_{\leq \tilde{n}}^{(e-1)}) - \text{median}_{j=\tilde{n}+1 \dots D} \tilde{s}_e(z_{\leq \tilde{n}}^{(e-1)})$$

$$\log \det f'_e(z^{(e-1)}) = \sum_{j=\tilde{n}+1}^D \tilde{s}_{e,j}(z_{\leq \tilde{n}}^{(e-1)}) = 0 \Rightarrow \det f'(z) = 1$$

- advantage $\hat{p}(x) = q(z=f(x))$,
e.g. simplifies power scaling $\hat{p}(x)^p$

disadvantage: volume-preserving NFs are not universal
e.g. cannot change # modes

• if X is an image: GLOW architecture [Kingma & Dhariwal 2018]

- $s_0(\cdot)$ and $t_0(\cdot)$ are convolutional networks
- active and passive part of $z^{(t+1)}$ is selected among the channels (not locations)
- Wavelet Flow decompose the image hierarchically using wavelets (similar to Fourier transform)

next time: how to learn conditional distributions $p(X|Y)$

latent properties \rightarrow observations