

Normalizing Flows (contd.)

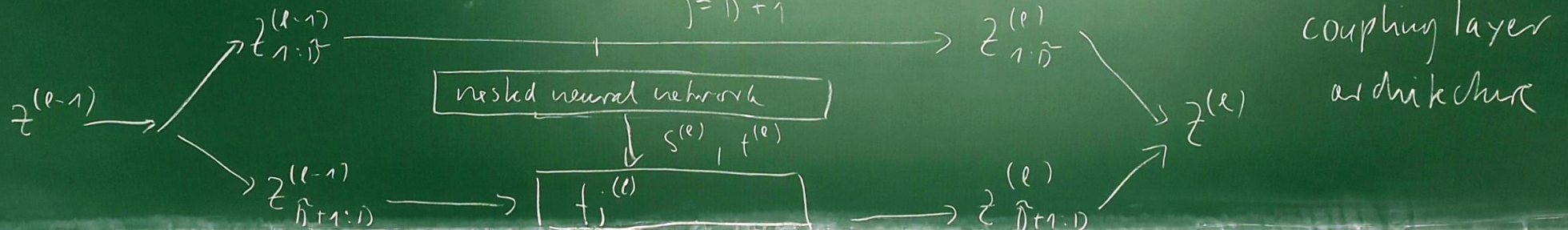
• recap: Real NVP, invertible neural network (INN)

– coupling layers with affine (⇒ easily invertible) transformations, with autoregressive form ⇒ Jacobian determinant easy

layer l :
$$z_j^{(l)} = z_j^{(l-1)} \quad \text{for } j=1, \dots, \tilde{D} \quad \tilde{D} = \lfloor \frac{D}{2} \rfloor$$

$z^{(l)} = f^{(l)}(z^{(l-1)})$
$$z_j^{(l)} = \exp\left(\tilde{s}_j^{(l)}(z_{1:\tilde{D}}^{(l-1)})\right) \cdot z_j^{(l-1)} + t_j^{(l)}(z_{1:\tilde{D}}^{(l-1)}) \quad \text{for } j = \tilde{D}+1, \dots, D$$

def
$$\left(\frac{\partial f^{(l)}(z^{(l-1)})}{\partial z^{(l-1)}} \right) = \begin{pmatrix} \tilde{D} & & \\ & \exp(\tilde{s}_j^{(l)}(z_{1:\tilde{D}}^{(l-1)})) & \\ & & 1 \end{pmatrix} \quad j = \tilde{D}+1$$



- turn this into a deep network, many coupling layers

naive: $z_{1:D}$ are always skip connections $p_E(z_{1:D}) = p^*(x_{1:D}) \neq q(z_{1:D})$

⇒ idea: change dimensions in skip part

after each coupling, add an (1) index permutation (special case of (2))

(2) orthonormal transformation Q

orthonormal is good because $|\det(Q)| = 1$, easily invertible: $Q^{-1} = Q^T$

- experimentally, it turned out that learning Q is not necessary, random fixed matrices are sufficient (may change with recent new learning alg. for orthonormal matrices)

how to create Q ?

- create a random gaussian matrix A

$$A_{ij} = N(0, 1)$$

- $Q, R = QR$ decomposition(A)

final architecture: $f = f^{(L)} \circ Q^{(L-1)} \circ f^{(L-1)} \circ \dots \circ Q^{(1)} \circ f^{(1)}$

$f^{(l)}$ coupling layer

$Q^{(l)}$ random orthonorm.

• training algorithm: minimize negative log-likelihood of data

$$\hat{f} = \underset{f}{\operatorname{argmin}} \mathbb{E}_{x \sim p^*(x)} \left[-\log p_f(x) \right] \quad p_f(x) = q(z=f(x)) \det(J_f)$$

$$\approx \underset{f}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N -\log p_f(x_i) = \underset{f}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N -\log q(f(x_i)) - \log \det(J_f(x_i))$$

$$\text{if } q(z) = \mathcal{N}(0, \Pi) \Rightarrow -\log q(z) = \frac{z^2}{2} + \log(z\pi)^{D/2}$$

$$\text{if } \det(J_f) \text{ is } \prod_{\ell=1}^L \prod_{j=\tilde{D}+1}^D s_j^{(\ell)} \underbrace{\left(z_{1:\tilde{D}}^{(\ell-1)} \right)}_{\exp(\xi^{\ell-1})} \Rightarrow \log \det(J_f) = \sum_{\ell=1}^L \sum_{j=\tilde{D}+1}^D \tilde{s}_j^{(\ell)} \left(z_{1:\tilde{D}}^{(\ell-1)} \right)$$

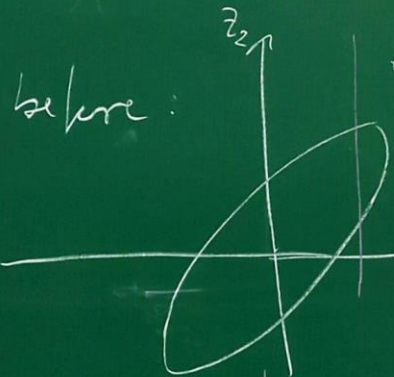
$$\hat{f} = \underset{f}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \left(\frac{f(x_i)^2}{2} - \sum_{\ell=1}^L \sum_{j=\tilde{D}+1}^D \tilde{s}_j^{(\ell)} \left(z_{i,1:\tilde{D}}^{(\ell-1)} \right) \right)$$

• Why does affine coupling work?

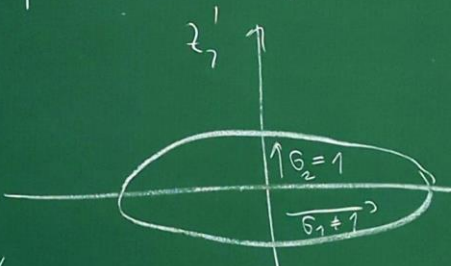
$p(z^{(e-1)}) = N(0, \Sigma) \rightarrow$ want that $p(z^{(e)})$ is closer to $N(0, \mathbb{I})$ than $p(z^{(e-1)})$

$$KL[p(z^{(e)}) \parallel N(0, \mathbb{I})] \approx \frac{1}{2} KL[p(z^{(e-1)}) \parallel N(0, \mathbb{I})]$$

holds empirically for most $p(z^{(e-1)})$, proven for $p(z^{(e-1)}) = N(0, \Sigma)$

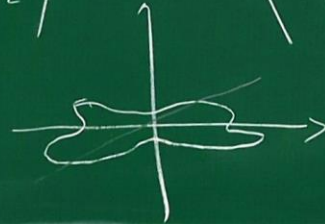
before:  ← coupling scales & translates every slice to mean=0 variance=1

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \rightsquigarrow$$

 random rotate and repeat

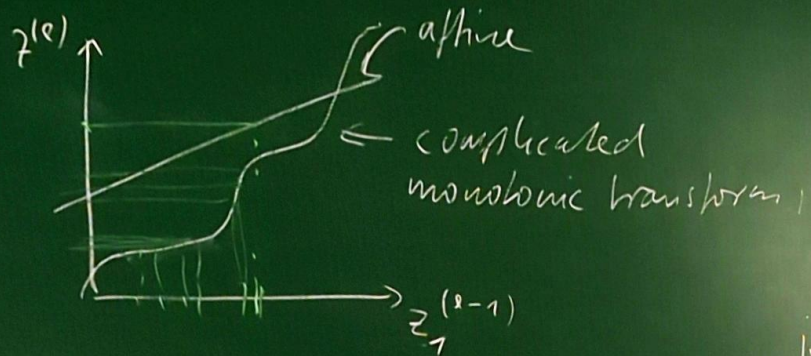
$$\begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} - (\Sigma_{12} \Sigma_{11}^{-1} \Sigma_{21}) \end{pmatrix}$$

after many random rotations enough "good slices" are encountered to get a standard normal



• spline mappings are superior to affine mappings in low dimensions ($D < 20$)

$$f_i^{(e)} = S_i^{(e)} z_j + f_j^{(e)}$$



ideally

$$\Phi^{-1} \left(\text{CDF} \left(z_j^{(e-1)} \mid z_{1:0}^{(e-1)} \right) \right)$$

inverse CDF of standard normal



Hermite splines

- location of knots $z_k^{(e-1)}$
- function at knots $z_k^{(e)}$
- derivatives at knots $f_k^{(e)}$

⇒ spline interpolates in between

affine approximates as straight line
spline as a piecewise simple fct.

→ nested networks learn these parameters & define some spline to interpolate

most popular: rational quadratic spline

"Neural spline Flows" [2019]

$$\frac{az^2 + bz + c}{dz^2 + e \cdot z + g}$$

- various libraries (e.g. FrEIA in pytorch, Bayesflow in tensorflow)
 - implement INN's
 - major variant for images: nested networks are convolutional "Glow" [Kingma & Dhariwal 2018]
 - funny face manipulation demo in the web
 - new: dict norm layers $\hat{=}$ adapt batchnorm to normalizing flows
 - better variants of normalizing flows
 - diffusion models (e.g. "stable diffusion"): implement $f(x)$ by a differential equation
 - free-form flows [Draxler et al. 2023]: autoencoder without bottleneck ($\dim(z) = \dim(x)$) \Rightarrow can be any network, no couplings or other restrictions
- later