

learn conditional densities $p(x|Y)$ with NFs

- e.g. $Y = \text{digit label}$ $x = \text{MNIST}$ $x \sim p(x) \Rightarrow$ sample any digit
- $x \sim p(x|Y=2) \Rightarrow$ sample only "2"s

often needed in practice we know (or measured) Y , but want to know X

typical setup of supervised learning
 traditional networks point estimates $\hat{x} = r(Y)$ (ideally $\hat{x} = \underset{x}{\text{argmax}} p(x|Y)$)
 conditional normalizing flows = distribution of $x \hat{=}$ estimate uncertainty of x

a autoregressive function is easy to generate for conditioning

$$z = f(x; Y) = \begin{pmatrix} f_1(x_1; Y) \\ f_2(x_2; x_1, Y) \\ \vdots \\ f_D(x_D; x_{1:D-1}, Y) \end{pmatrix}$$

idea: do this for all couplings

$$f_j^{(e)}(z_j; z_{1:D}^{(e-1)} | Y) = S_j^{(e)}(z_{1:D}^{(e-1)} | Y) \cdot z_j^{(e-1)} + t_j^{(e)}(z_{1:D}^{(e-1)} | Y)$$

\Rightarrow add Y as a new input to all nested networks $S_j^{(e)}, t_j^{(e)}$

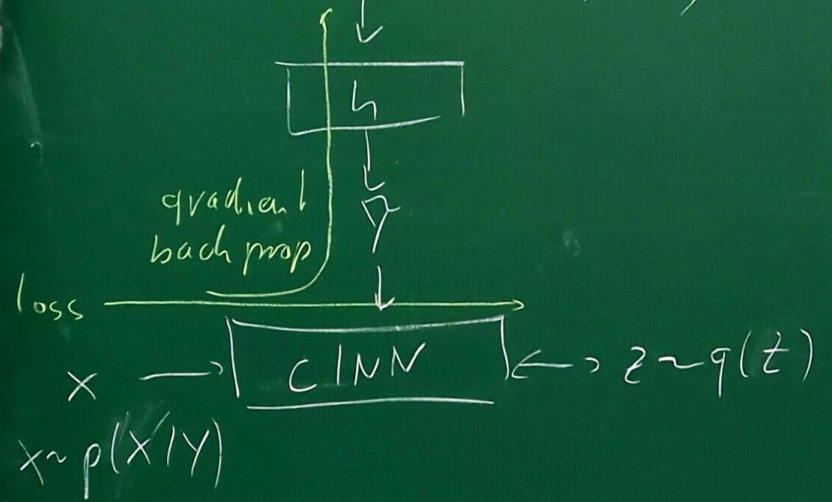
- works because Y is known for both forward and reverse network execution

if Y is complicated (e.g. high dimensional, image), processing Y again in each coupling layer is wasteful

\Rightarrow shared preprocessing network $\tilde{Y} = h(Y)$ "feature detector"

$$Y \leftarrow Y \sim p^*(Y) \text{ or } \sim p^*(Y|X)$$

"summary network"



tricks to define $h(Y)$

- use architecture of an existing regression network

$h(Y)$ is $r(Y)$ minus final layer(s)

- use a foundational model $\phi(Y)$

$$h(Y) = \tilde{h}(\phi(Y))$$

\uparrow
small

feature detector pre-trained by the big guys on massive data

e.g. images: CLIP, MOCO ($> 10^9$ training images)

train $f(x, h(Y))$ jointly