

< 탄생배경 >

Reviewer: 김정하

이 논문은 특정한 데이터의 분포의 특징을 잘 잡아내어 다른 일반적인 혹은 간단한 분포로 바꾸어낼 수 있는 모델을 만드는 것이다.

x 가 원데이터고 h 가 우리가 해석가능한 분포 $p_H(\cdot)$ 를 가진다고 하자

$x \xrightarrow{f} h$ 같이 f 는 $h=f(x)$ 로 x 를 h 로 mapping 시켜주는 함수다. $p_H(h) = \prod p_{H_i}(h_i)$ 로 나타낼 수

있다. 이는 분포를 factorize 할 수 있다는 뜻이다. (MLP와 비슷한 느낌이라고 생각하면 된다) 중요한 점은 f 는 invertible 하고

h 의 차원은 x 와 같다. 'change of variable rule'을 따르면

① ... $p_x(x) = p_H(f(x)) \left| \det \frac{df(x)}{dx} \right|$ 를 얻을 수 있다. 후에 보겠지만 f 는 $|\det|$ 를 구하기 쉽고

쉽게 inverse 되어야 한다. 그리고 이 논문의 핵심 idea는

x 를 두개로 쪼개는 거다. (x_1, x_2) 로 쪼개고

$\begin{cases} y_1 = x_1 \\ y_2 = x_2 + m(x_1) \end{cases}$ 로 나타내는 것이다. m 은 아무런 함수여도 된다 (뉴럴 네트워크도 된다) 이렇게 나타내면 역으로도 쉽게 쓸 수 있다. (아래 같이)

$$\begin{cases} x_1 = y_1 \\ x_2 = y_2 - m(y_1) \end{cases}$$

< content >

① 식을 \log 를 취하면 (이러는 이유는 최댓값을 얻기 위해서..)

$$\log(p_x(x)) = \log(p_H(f(x))) + \log\left(\left|\det\left(\frac{df(x)}{dx}\right)\right|\right)$$

p_H 는 우리가 지정한 사전확률이다. $p_H(h) = \prod p_{H_i}(h_i)$ 라면

$$\log(p_x(x)) = \sum_{i=1}^D \left(\log(p_{H_i}(f_i(x))) + \log\left(\left|\det\left(\frac{df(x)}{dx}\right)\right|\right) \right) \text{로 나타낼 수 있다.}$$

데이터를 단순히 축소해서 함수에 넣어볼 수 있지만 자코비안 항 때문에

계산하는 것을 지양하고 곱으로 지역을 퍼뜨린다.

① triangular structure

$\det(\cdot)$ 부호의 tractable하고 계산이 쉬도록 하는 종류를 찾아야 하는데

이 눈썰미가 제대한 건 triangular matrices다. 왜냐하면

$\det(\cdot)$ 가 diagonal element를 곱하는 걸로 간단히 구할 수

있고 inverting 하는 것도 계산이 간단하다. (square matrix

$M = LU$ 로 구할 수 있다. L 은 lower triangular, U 은

upper triangular)

② Coupling layer

이 눈썰미에서 가장 기본이 되는 layer다.

$x \sim R^D$ 이고 이를 (x_{I_1}, x_{I_2}) 로 쪼갤 수 있다.

($I_1 \rightarrow 1 \sim d$ 까지 $I_2 \rightarrow d \sim D$ 까지를 나타냄.)

$\begin{cases} y_{I_1} = x_{I_1} \end{cases}$ 로 나타낼 수 있다.

$\begin{cases} y_{I_2} = g(x_{I_2}; m(x_{I_1})) \end{cases}$ g 는 가역 invertible 한 함수다.

($R^{D-d} \times m(R^d) \rightarrow R^{D-d}$ 를 만족하면 된다)

그러면 $\frac{dy}{dx} = \begin{bmatrix} I_d & 0 \\ \frac{dy_{I_2}}{dx_{I_1}} & \frac{dy_{I_2}}{dx_{I_2}} \end{bmatrix}$ 인 triangular

matrix를 얻을 수 있다. 역으로 나타낼 수 있다.

$\begin{cases} x_{I_1} = y_{I_1} \end{cases}$

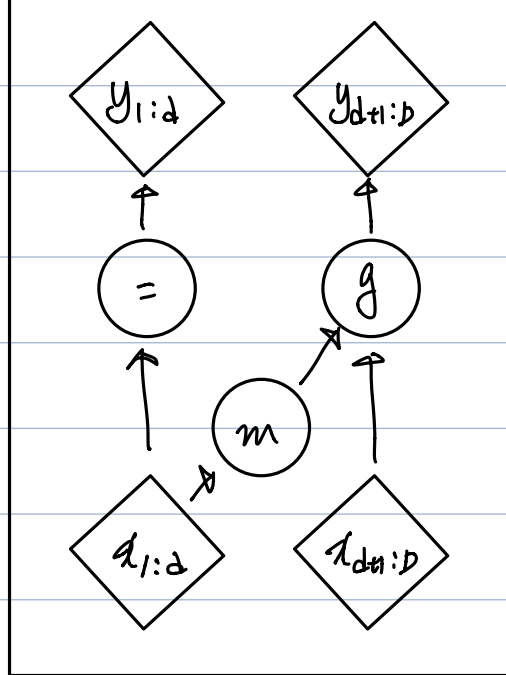
$\begin{cases} x_{I_2} = g^{-1}(y_{I_2}; m(y_{I_1})) \end{cases}$

이 layer를 복잡하게 연결만 하면 된다. '복잡하게'라는

게 인풋을 나누는 전체 각 layer마다 나누는 구간을

다르게 설정하면 그만이다. (계속 똑같이 나누면

재미는 사라지 않기 때문)



그림으로 나타내면 왼쪽과 같다.

③ Allowing rescaling

각 layer에서 자코비안 행이 있다. 그 행성도 결국 자코비안 행이 필요한 건데 이는 volume preserving의 특성을 가지고 있다. 근데 이를 깨볼 수 있는데 대각행렬 S 를 추가하는 것이다. 이를 통해 어떤 차이를 가정치를 더 주고 더 주고 할 수 있다. 또 식도 다음과 같이 바뀐다.

$$\log(p_X(x)) = \sum_{d=1}^D [\log(p_{H_i}(f_i(x))) + \log(|S_i|)]$$