

Normalizing Flow

Zarys Problemu:

- Mamy pewien rozkład $p_x(x)$, który jest "trudny" np 
- Mamy próbki z tego rozkładu x , ale niekoniecznie potrafimy obliczać $p_x(x)$
- możemy postawić dwa problemy: \rightarrow jak obliczać $p_x(x)$?
 \rightarrow jak generować próbki z tego rozkładu?

Idea normalizing Flow:

- zacznijmy od prostego rozkładu i "przetransformujmy" go w nasz żądany w sposób odwracalny
 $z \xrightleftharpoons[F^{-1}]{F} x$, wtedy $p_x(x) = p_z(z) \left| \det \left(\frac{\partial F}{\partial x} \right) \right|$ [F - musi być wystarczająco "dobre"]
 $= p_z(z) \left| \det \left(\frac{\partial F}{\partial x} \right) \right|^{-1}$ [F, F^{-1} muszą być różniczkowalne]
- z zmiennymi zmiennymi

\Rightarrow Zmiana problemu na szukanie odpowiedniej f -gi: F

zauważmy, że F może być złożeniem $F = f_k \circ f_{k-1} \circ \dots \circ f_3 \circ f_2 \circ f_1$ i jeżeli każda z f_i :

będzie spełniać nasze założenia to F też, ponieważ mamy

$$(f_2 \circ f_1)^{-1} = f_1^{-1} \circ f_2^{-1} \quad \text{oraz} \quad \det(J_{f_2 \circ f_1}(u)) = \det(J_{f_2}(f_1(u))) \cdot \det(J_{f_1}(u))$$

Intuicja ze przekształceniem i jacobianem:

F - "rozciąganie" przestrzeni \mathbb{R}^d , tak żeby zamienić $p_x(x)$ w $p_z(z)$, $|\det(J_F(\cdot))|$ odpowiada za zmiany objętości w "każdej" otoczeniu. ze względu na F . Biorąc mechanicznie małe otoczenie du i dx , ponieważ nasze prawdopodobieństwo się nie zmienia to zmiana objętości wpłynie na gęstość: $V \uparrow \Rightarrow g \downarrow$, $V \downarrow \Rightarrow g \uparrow$

Ciepło chcemy od f_i :

- wystarczająco "elastyczne" żeby przekształcić prosty rozkład w złożony
- f_i oraz f_i^{-1} muszą istnieć i być w miarę łatwe do obliczenia
- wyznacznik jacobianu $f_i^{-1}(f_i)$ musi być w miarę prosty do obliczenia

Żeby znaleźć funkcje spełniające ostatnią własność czysto konstanty, własności wyznacznika np.

Wyznacznik macierzy trójkątnej jest równy iloczynowi przekątnych

Uczenie: maksymalizujemy log-wiarygodność $p_x(x) = \log(p_u(u)) + \sum_{i=1}^k \log \left(\left| \det \frac{\partial f_i(z)}{\partial z} \right|^{-1} \right)$

Lemat o wyznaczniku macierzy: jeżeli A jest $d \times d$ i odwracalna, a V, W są macierzami $d \times m$ to $\det(A + VW^T) = \det(I + W^T A^{-1} V) \det A$

jeżeli musimy łatwo policzyć A^{-1} i $\det(A)$, oraz $m < d$ to wyznaczenie prawej strony jest obliczeniowo lepsze niż lewej $O(d^3 + d^2 m)$ [lewa strona], a dla np. diagonalnej A prawa strona to $O(m^3 + dm^2)$

Rozwar flow

$f(z) = z + v \sigma(w^T z + b)$ [sieć neuronowa z jednym neuronem ukrytym]

$v, w \in \mathbb{R}^d$, $b \in \mathbb{R}$ są parametrami flowu (do nauki) a σ jest różniczkowalną

funkcją aktywacji np. tangens hiperboliczny

Możemy interpretować ten flow jako "rozciąganie/skrywanie" przestrzeni w kierunku prostopadłym do hiperpłaszczyzny $w^T z + b = 0$

Jakobien danych jest wzorem $J_f(z) = I + \sigma'(w^T z + b) v w^T$

σ' - pochodna f-ji aktywacji

kontynuując z lematu o wyznaczniku: $\det J_f(z) = 1 + \sigma'(w^T z + b) w^T v$

Żeby ten flow był odwracalny musimy uzyskać dodatnie wartości:

- σ' jest dodatnie i ograniczone

- Warunek dostateczny dla odwracalności $w^T v > - \frac{1}{\sup_t \sigma'(t)}$

Sylvester Flow (uogólnienie Planer flow na M meromorficznych)

$f(z) = z + V\sigma(W^T z + b)$, $V, W \in \mathbb{R}^{d \times m}$, $b \in \mathbb{R}^m$ to parametry do nauczenia

σ - f -y σ aktywacji, metadane po współrzędnych

Jacobian przekształcenia $J_f(z) = I + VS(z)W^T$, gdzie $S(z)$ jest macierzą diagonalną $m \times m$ o diagonalce $\sigma'(W^T z + b)$, rząd z lenetki:

$\det J_f(z) = \det(I + S(z)W^T V)$, możemy zdefiniować $V = QU$, $W = QL$

gdzie $Q \in \mathbb{R}^{d \times m}$ ma kolumny ortonormalne, $U \in \mathbb{R}^{m \times m}$ jest górnokątowa, a $L \in \mathbb{R}^{m \times m}$ dolnokątowa, ponieważ mamy $Q^T Q = I$ oraz $L^T U$ jest dolnokątowa to:

$$\det J_f(z) = \det(I + S(z) \underbrace{L^T Q^T Q U}_{I}) = \prod_{i=1}^m (1 + S_{ii}(z) L_{ii} U_{ii})$$

dobrych założenia gwarantujące odwracalność:

- σ' jest dodatnie i ograniczone
- $L_{ii} U_{ii} > -\frac{1}{\sup_t \sigma'(t)}$, $i \in \{1, \dots, d\}$

Flow sprężysty

Dzielmy z na dwie części (z^A, z^B) i niech h będzie bijekcją $\mathbb{R}^d \rightarrow \mathbb{R}^d$
definiujemy $f(z)$ w ten sposób: $f(z)^B = h(z^B, \Theta(z^A))$

$$f(z)^A = z^A$$

$\Theta(z^B)$ - dowolna funkcja

f jest odwracalny i h jest odwracalna i $f^{-1}(z)^B = h^{-1}(f(z)^B, \Theta(z^A))$
 $f^{-1}(z)^A = f(z)^A = z^A$

Jakobian f :

I	O
A	Dh

 $\Rightarrow \det(J_f) = \det(Dh)$

Najczęściej h bierze postać $h(x) = Ax + b$ i wtedy Dh jest dostrójkowane

Zauważ, że $\Theta(x^B)$ może być dowolnie skomplikowane

Przykład: Addytywne warstwa sprężysta

$$x^A = z^A \quad x^B = z^B + m(z^A) \quad \text{gdzie } m - \text{np. sieć neuronowa}$$

ale zauważ, że $\det J_f = 1$ [cyki zachowują objętość \rightarrow mocne ograniczenie]

Poprawka: Przed warstwą sprężystą dodajemy warstwę skalującą

$$x = Sz \quad S - \text{diagonalna}$$

Albo: Real NVP

$$x^A = z^A \quad x^B = S^\Theta(x^A) \odot x^B + t^\Theta(x^A) \quad [\text{tętuś odwracalna}]$$

\uparrow skalowanie \uparrow transkrypcja

Uwagi i ogólne architektury:

- jeżeli mamy warstwę, które zmieniają tylko część wejść to pomiędzy musimy dodać warstwę z permutacjami
- możemy rozwiązać złożenie kilku warstw
- możemy badać prawdopodobieństwo próbki x : $P_x(x) = P_z(F^{-1}(x)) \det(J_{F^{-1}}(x))$
- możemy generować nowe próbki: generujemy $z = P_z(z)$ i przekształcamy $x = F(z)$