

Neural Networks Report #8

Introduction

The final homework is about the transformer model this time more inclined to the calculation part, nonlinear independent component analysis and variational autoencoder. A transformer is a deep learning model that adopts the mechanism of attention, weighing the influence of different parts of the input data and uses large number of adaptive elements for this. Independent component analysis is a computational method for separating a multivariate signal into additive subcomponents. VAEs are directed probabilistic graphical models whose posterior is approximated by a neural network, forming an autoencoder-like architecture, in other words autoencoder that employs probabilistic model.

Problem1

a)

$$N) X = (X_1, \dots, X_{N_{voc}})$$

$$X: N_{mod} \times N_{voc}$$

$$\therefore (N_{query} = N_v = N_{head}) \quad \left. \begin{array}{l} \text{given} \\ \text{to model} \end{array} \right\}$$

$$N_{mod} = N_v \times N_h$$

$$Q, K, V: N_{mod} \times N_{head}$$

$$W_h^j: N_{mod} \times N_{q/v}$$

$$A_h = \text{soft}[Q^T W_h^Q (W_h^K)^T K / \sqrt{N_q}]: N_{pos} \times N_{pos}$$

$$B = A_h V_h^T: N_{val} \times N_{pos}$$

$$C = W^O B \text{ where } W^O: N_{mod} \times (N_{val} \times N_{head})$$

$$C: N_{mod} \times N_{pos}, y = \text{MLP}(C): N_{mod} \times N_{hd}$$

As a result we need to train

$$W_h^Q, W_h^K, W_h^V \text{ and } W^O$$

So the total number is

$$N_t = 3 \cdot N_{mod} \times N_{qvk} + N_{mod} \times N_{val} \times N_{head}$$

$$N_1 = N_{val} \cdot N_h (3N_{qvk} + N_{val} \cdot N_{head})$$

where $N_{qvk} = N_{que} = N_{key} = N_{val}$

as they should be equal by constraint

as we also have hidden layer

there is other matrix $W: N_{hd} \times N_{mod}$

As there are N_{enc} blocks we have

$$N_1 = N_{enc} N_{val} \cdot N_h (3N_{qvk} + N_{val} \cdot N_{head} + N_{hd})$$

b)

b) for input it is evident that

$$X: N_{\text{mod}} \times N_{\text{voc}} = N_{\text{qvk}} \times N_{\text{head}} \times N_{\text{voc}}$$

Similarly the Output size should be same as input as it is only the shifted version. as we also have encoder/decoder boxes

$$N_2 = (N_{\text{enc}} + N_{\text{dec}}) \times N_{\text{qvk}} \times N_{\text{head}} \times N_{\text{voc}}$$

c)

c) The first module has same step

number, so $N_3 = N_1 \frac{N_{dec}}{N_{enc}}$
 the masked multihop fills the
 values of dimension $N_{mod} = N_{hid}$

which equals to $N_{mod} \times N_{head}$

Therefore

$$N_4 =$$

$$N_{mul} N_{voc} \cdot N_h (3N_{gvr} + N_{voc} \cdot N_{head} + N_{hid})$$

$$\text{where } N_{mul} = \frac{N_{hid}}{N_{head}}$$

$$N = \sum_i N_i$$

$$\begin{aligned} &= (N_{enc} + N_{dec}) \times N_{gvr} \times N_{head} \times N_{voc} \\ &\quad + N_{voc} \cdot N_h (3N_{gvr} + N_{voc} \cdot N_{head} + N_{hid}) \\ &\quad \times [N_{enc} + N_{dec} + \frac{N_{hid}}{N_{head}}] \end{aligned}$$

d)

d) With numerical values

$$N_1 = 6 \cdot 64 \cdot 8 (3 \cdot 64 + 64 \cdot 8 \cdot 1064)$$

$$N_1 = 5431296$$

$$N_2 = (6+6) \cdot 64 \cdot 8 \cdot 50000$$

$$N_2 = 307200000$$

$$N_3 = N_1 = 5431296$$

$$N_4 = 115864648$$

$$\text{so } N = 423064648$$

Problem2

a)

For scalar u & X :

a) $u = f(x)$ we have $p(u)$

$$\begin{aligned} F_U(u) &= P(U \leq u) = P(f(X) \leq u) \\ &= P(X \leq f^{-1}(u)) = F_X(f^{-1}(u)) \end{aligned}$$

$$F(X) = \int_0^x p(x) dx \rightarrow p(x) = \frac{dF(x)}{dx}$$

$$p(u) = \frac{dF_X(f^{-1}(u))}{du}$$

$$p(x) = \frac{dF_X(x)}{dx}$$

$$\text{let } x = f^{-1}(u)$$

$$\text{then } dx = \frac{du}{f'(f^{-1}(u))} = \frac{du}{f'(x)}$$

$$\text{as } x = f^{-1}(u)$$

$$p(u) = \frac{dF_X(x)}{dx f'(x)} = \frac{p(x)}{f'(x)} = \frac{p(x)}{df/dx}$$

so generally we can say

$$p(u) = p(f^{-1}(y)) \left| \frac{d}{dy} (f^{-1}(y)) \right|$$

We know that vector is

$$f_x(x) = f_{x_1 \dots x_n}(x_1, \dots, x_n)$$

if we follow same rules but using vector definitions we get

$$p(u) = p(f^{-1}(u)) \left| \det \left[\frac{df^{-1}(\vec{z})}{d\vec{z}} \right] \right|$$

$$\text{or } p(u) = p(x) \det \left[\left(\frac{df(x)}{dx} \right)^{-1} \right] \quad z=y$$

where $\frac{df}{dx}$ is jacobian matrix

$$J(f) = \frac{df}{dx} = \begin{bmatrix} \frac{du_1}{dx_1} & \frac{du_2}{dx_1} & \dots \\ \frac{du_2}{dx_1} & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

So
$$p(u) = \frac{p(x)}{\det \left(\frac{df}{dx} \right)}$$

b)

$$b) I(u) = E[\log p_u(f(Wx))] - \sum_j E[\log p_j(u_j)]$$

$$p_u[f(Wx)] = \frac{p(Wx)}{\det(f'(Wx))} = \frac{p(x)/\det(W)}{\det(J(f))}$$

$$\begin{aligned} I(u) &= E[\log p_u(Wx)] - \log |\det(W)| \\ &\quad - \sum_j E[\log p_j(u_j)] - \log(\det(J(f))) \\ &= E[\log p_u(x)] - \log |\det(W)| \\ &\quad - E[\log(\prod_j p_j(u_j))] - \log[\det(J(f))] \end{aligned}$$

$$\text{since } \log a + \log b = \log ab \quad \& \quad E\left(\sum_{i=1}^n x_i\right) = \sum_{i=1}^n E(x_i)$$

$$\begin{aligned} \frac{dI}{dw_{ik}} &= -\frac{1}{\det(W)} \frac{d}{dw_{ik}} \det(W) \\ &\quad - \sum_j \frac{1}{p_j(u_j)} \frac{dp_j(u_j)}{du_j} \frac{du_j}{dw_{ik}} \\ &\quad - \frac{1}{\det(J(f))} \cdot \frac{d}{dw_{ik}} \det(J(f(Wx))) \\ &= -(W^{-T})_{ik} + \phi(u_i) x_k - (J(f(Wx)))_{ik} \end{aligned}$$

$$\text{where } \phi(u_i) = \frac{1}{p_i(u)} \frac{dp_i(u)}{du}$$

$$\Delta W = -\eta \frac{\partial I(u)}{\partial w_{ik}} = \eta (W^T \phi(u) x^T - J(f(Wx))^T) \\ = \eta (I - \phi(u) u^T) W^T - J(f(Wx))^T$$

making it computationally efficient, we can multiply it by $W^T W$

$$\Delta W = \eta (I - \phi(u) u^T) W - J(f(Wx))^T W^T W$$

c)

c) If we have one hidden layer perceptron
 we can say the function is sigmoid one

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad \frac{d\sigma}{dx} = \sigma(x)(1-\sigma(x))$$

so let's use it in our perceptron

Let $y = Wx$. $m \times n$

$$J(\sigma(y)) = \begin{bmatrix} \frac{\sigma(y_1)}{dx_1} & \frac{\sigma(y_2)}{dx_1} & \dots & \frac{\sigma(y_m)}{dx_n} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

$$= \begin{bmatrix} \sigma(y_1)(1-\sigma(y_1)) & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & \sigma(y_m)(1-\sigma(y_m)) & 0 \end{bmatrix}$$

$$J(\dots)^T = J(\dots)$$

as J is diagonal matrix

$$D^{-1} = \frac{1}{|D|} \text{adj}(D)$$

$$\det(J^T) = \prod_i \sigma(y_i)(1-\sigma(y_i))$$

let's make further simplifications &

now we have 2x2 matrices

$$J^{-T} = \frac{1}{\sigma(y_1)\sigma(y_2)(1-\sigma(y_1))(1-\sigma(y_2))} \begin{bmatrix} \sigma_2(1-\sigma_2) & 0 \\ 0 & \sigma_1(1-\sigma_1) \end{bmatrix}$$

$$= \begin{bmatrix} [\sigma_1(1-\sigma_1)]^{-1} & 0 \\ 0 & [\sigma_2(1-\sigma_2)]^{-1} \end{bmatrix}$$

so

$$\Delta W = \eta (I - \phi(u) u^T - \begin{bmatrix} [\sigma_1(1-\sigma_1)]^{-1} & 0 \\ 0 & [\sigma_2(1-\sigma_2)]^{-1} \end{bmatrix} W^T) W$$

$$\sigma_i = \sigma_i(Wx) = u_i$$

& we get

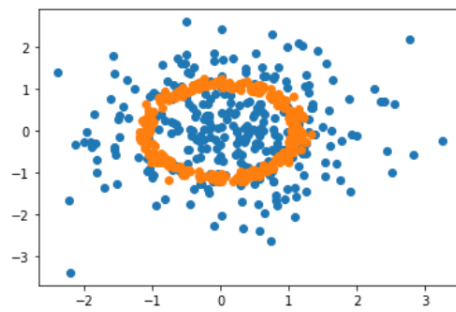
$$W[n+1] = W[n] + \Delta W$$

$$\Delta W = \begin{bmatrix} \Delta w_1 \\ \Delta w_2 \end{bmatrix} = \eta \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} \frac{1}{\sigma_1(u)} \frac{\partial \phi(u)}{\partial u_1} & \frac{1}{\sigma_2(u)} \frac{\partial \phi(u)}{\partial u_2} \\ \frac{1}{\sigma_2(u)} \frac{\partial \phi(u)}{\partial u_2} & \frac{1}{\sigma_1(u)} \frac{\partial \phi(u)}{\partial u_1} \end{bmatrix} \begin{bmatrix} u_1 & u_2 \end{bmatrix} - \begin{bmatrix} \sigma_1(Wx)(1-\sigma_1(Wx)) & 0 \\ 0 & \sigma_2(Wx)(1-\sigma_2(Wx)) \end{bmatrix} W \right)$$

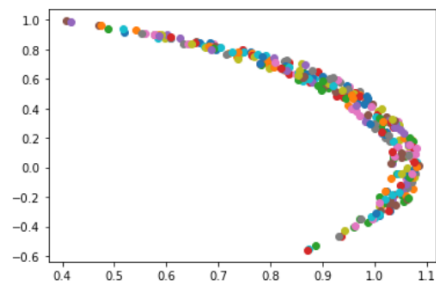
Problem3

a) The code is in the notebook provided, it just generates the gaussian and then maps it into the result.

z: blue
y: orange



b) In this task I used previous code, just modified the perceptron to give two output vector function. However, as previously, it did not work so well and as a result the output is somehow disturbed.



The behavior is similar to the noisy one, and every time it gives different output. So the perceptron does not work properly. Btw, the hidden layer number is set to be 5.

c)

$$\begin{aligned}
 c) D_{\text{KL}} &= \int p_x(x) \log \left(\frac{p_x(x)}{q_x(x)} \right) dx \\
 &= E [\log(p_x(x)) - \log(q_x(x))]
 \end{aligned}$$

We have $N(\mu_0, \Sigma_0)$ & $N(\mu_1, \Sigma_1)$

$$N(\mu, \Sigma) = \frac{1}{\sqrt{2\pi} |\Sigma|} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$\begin{aligned}
 D &= E \left[-\frac{1}{2} \log 2\pi |\Sigma_0| - \frac{1}{2} (x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0) \right. \\
 &\quad \left. + \frac{1}{2} \log 2\pi |\Sigma_1| + \frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) \right]
 \end{aligned}$$

$$\begin{aligned}
 &= E \left[\frac{1}{2} \log \frac{2\pi |\Sigma_1|}{2\pi |\Sigma_0|} \right] - \frac{E}{2} \left[(x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0) \right] \\
 &\quad + \frac{E}{2} \left[(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) \right]
 \end{aligned}$$

$(x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0) \in \mathbb{R}$ so we can write it as

$$\text{trace: } \text{tr}[(x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0)]$$

Using a theorem $\text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$

we get a term

$$\frac{1}{2} E [\text{tr} \{ (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) \}]$$

as expectation is linear function we can interchange:

$$= \frac{1}{2} \text{tr} \{ E [(x - \mu_0) (x - \mu_0)^T \Sigma_0^{-1}] \}$$

$$= \frac{1}{2} \text{tr} \{ \underbrace{E [(x - \mu_0) (x - \mu_0)^T]}_{\Sigma_0} \Sigma_0^{-1} \}$$

$$= \frac{1}{2} \text{tr} \{ \Sigma_0 \Sigma_0^{-1} \}$$

$$= \frac{1}{2} \text{tr} \{ I_K \} = \frac{K}{2} \quad (1)$$

Using similar operations for the third term.

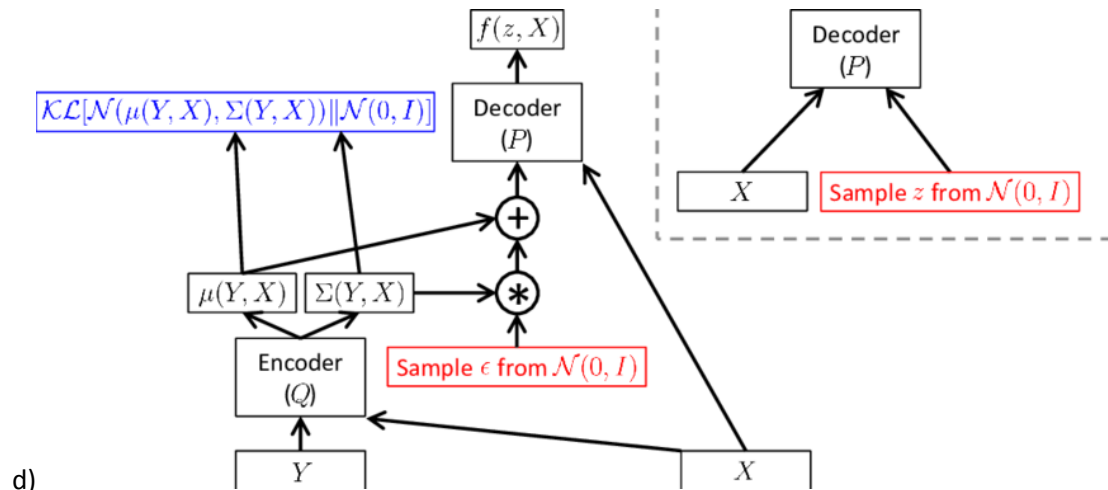
$$E [(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_0)] = (\mu_0 - \mu_1)^T \Sigma_1^{-1} (\mu_0 - \mu_1) + \text{tr} \{ \Sigma_1^{-1} \Sigma_0 \} \quad (3)$$

It is done using a theorem (380) in the referenced book

Combining, we get

$$D[N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)] =$$

$$\frac{1}{2} (\text{tr}(\Sigma_1^{-1} \Sigma_0) + \mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - k \\ + \log\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) \quad \text{[red mark]}$$



As always there are two parts, encoder and decoder. Encoder wants to encode Y into its hidden representation while decoder wants to find this hidden representation to decode into the original image. So, both encoder Q and decoder P are neural networks that we aim to train to. Q takes Y and produces some distribution which is gaussian for simplicity, thus Q needs to find parameters average and variance of the input and at the same time wants it to be close to $\mathcal{N}(0, I)$ that's why its divergence with the encoded part is measured. After that, the samples are transferred to the decoder that wants to use these hidden features to decode the input. As the transferred sample are hard to train, the reparameterization trick is used instead, instead of samples, the parameters are given to P and using the normal function this data is reproduced again. Then feed forward part for the decoder occurs and result is given. The conditional variational encoder helps to control the learning process by introducing new conditional variable X to both Q and P , after that the distribution produced depends on the condition of X so we have $P(z|X)$ instead of $P(z)$, $Q(z|Y, X)$ instead of $Q(z|Y)$ etc. In other words, let's say, given an input X (label of the image) we want our generative model to produce output Y (image). So, the process of VAE will be modified as the following: given observation y , z is drawn from the prior distribution $P(z|y)$, and the output x is generated from the distribution $P(x|y, z)$. X is like the label of the image which makes the learning resemble supervising one.

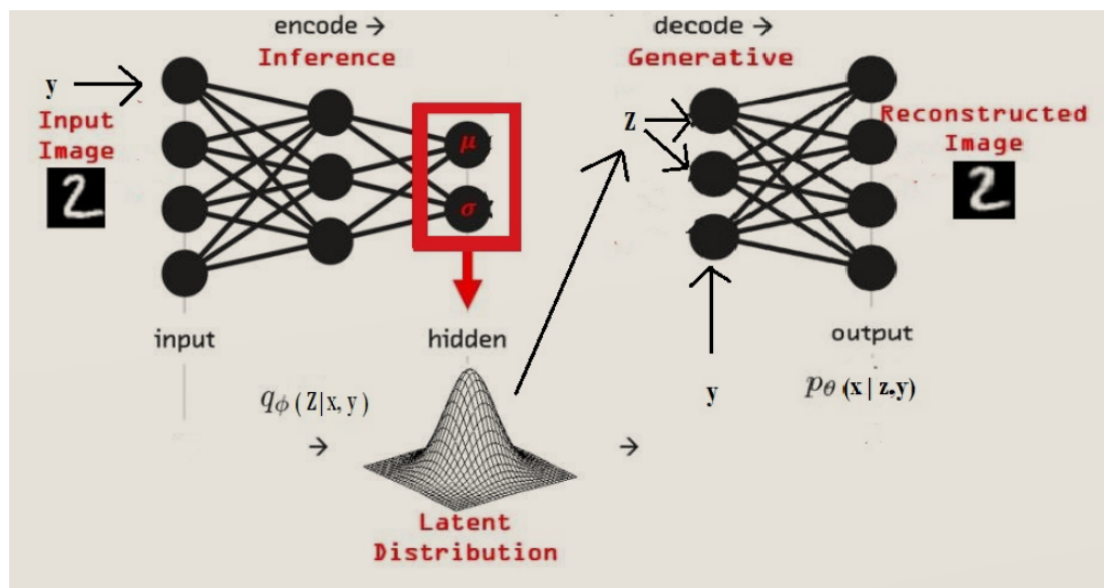
To make the back propagation we need to define the objective function that we would like to minimize/maximize. This function will help to derive the loss function which will be used in the next section for deriving the learning rules. First of all, we want the input and the decoded part to be the same, i.e. we would like to minimize the information is lost during the decoding process. It can be achieved by using KL divergence between $Q(z|Y, X)$ and $P(z|Y, X)$. Another thing is we would like to maximize the distribution of the resultant pdf. The greater the distribution of data, the more information, the larger is capability of the model is for the decoder to decode different data. As a result, objective function is

max E

$$E = \log P(Y|X) - D[Q(z|Y, X) || P(z|Y, X)]$$

For deriving the learning rule, the next step is make the loss term from this function and feed it to the model do make the update process.

e) The process can be explained by picture below



$$d) \log P(Y|X) - D[Q(z|Y, X) || P(z|Y, X)]$$

$$= E_{z \sim Q(\cdot|Y, X)} [\log P(Y|z, X)] - D[Q(z|Y, X)$$

$$\uparrow \\ || P(z|X)]$$

(lower bound

This is the function we would like to

maximize, the loss function needs to

be minimized, the lower bound of

objective function is then

$$\max - E[\log P(Y|z, X)] + D[Q(z|Y, X) || P(z|X)]$$

As our model wants $P(z|X)$ to be close to $N(0, I)$ we may assume $P(z) \sim N(0, 1)$

then $D[\dots]$ has a closed form solution.

Putting it together

$$E_{z \sim Q(z|X, Y)} \log P(Y|z, X) \propto \|Y - f(z, X)\|^2$$

as the performance of decoder strongly

depends on the result of neural network
To maximize the function, we would
like to minimize its lower bound so
now it can have larger range for
the following conditions. So we can
write loss as

$$L_e = \|X - f(z|x)\|^2 - \lambda D[Q(z|x) \| P(z|x)]$$

where the 1st part is pixel difference
between input & output and the second
one is regularization term that avoids
overfitting & λ is parameter $0 \leq \lambda < 1$

This is tricky as we want to minimize
 D , however when training, this term needs
to be high in order to the model to not
stick to the parameters provided & be

flexible.

So generally we can write as

$$Q(z|Y, X) = \sigma(W_n \underbrace{\sigma(\dots \sigma(W_1 x))}_{\text{nonlinear function}}) \dots$$

$$f(z, X) = \sigma(W_n \sigma(\dots \sigma(W_1 Q(z|Y, X)) \dots))$$

$$L = \|X - f(z, X)\|^2 - \lambda \mathbb{E} [Q(z, X) \log P(z, X)]$$

$$\Delta W_p = -\eta L \quad 0 < \eta < 1$$

$$\Delta W_q = -\eta L$$

$$W[n+1] = W[n] + \Delta W$$

The rules defined above cannot explain the learning process well as they are not the algorithm. So I will also provide small pseudocode below[3]:

Given a dataset of examples $Y = \{Y_1, Y_2, \dots\}$

Initialize parameters for Encoder and Decoder

Repeat till convergence:

$Y^M \leftarrow$ Random minibatch of M examples from Y under condition X

$\epsilon \leftarrow$ Sample M noise vectors from $N(0, I)$

Compute $L(Y^M, \epsilon, \theta, X)$ (i.e. run a forward pass in the neural network)

Gradient descent on L to updated Encoder and Decoder.

Reference

- [1] <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>
- [2] <https://towardsdatascience.com/understanding-conditional-variational-autoencoders-cd62b4f57bf8>
- [3] https://slazebni.cs.illinois.edu/spring17/lec12_vae.pdf