

ECE 792 HW 3
Kriti Singh(ksingh23)

1.

1) For all the convolution layers ^{inputs and} outputs let's represent them as $h_{in}^i = W^i x + B^i$ _(input)
 $h_{out}^i = \text{maxpool/sigma}(h_{in}^i)$

Calculating the gradient of Loss function at the softmax output layer w.r.t input ^{at that} _{layer} to the softmax

$$\Rightarrow \frac{\partial L}{\partial h_{3in}} = \hat{y} - y$$

Calculating gradient of L at the fully connected layers.

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial h_{3in}} \times \frac{\partial h_{3in}}{\partial W^3}$$

h_{3in} is nothing but the output of the fully connected layers

$$= (\hat{y} - y) \times \frac{\partial (h_2^{out} W_3 + B_3)}{\partial W_3}$$

$$= (\hat{y} - y) (h_2^{out})^T$$

Now ~~find~~ output of 2nd convolution layer
 h_{2out} is ~~not~~

$$\text{Maxpool}[\underbrace{\text{sigma}(W_2^T x + B_2)}_{(h_{2in})}]$$

Calculating gradient ~~of~~ of L w.r.t 2nd
 layer weights,

$$\frac{\partial L}{\partial W_2} = (\hat{y} - y) \frac{\partial h_{3in}}{\partial W_2}$$

$$= (\hat{y} - y) \frac{\partial (h_{2out} W_3 + B_3)}{\partial W_2}$$

$$= (\hat{y} - y) \frac{\partial [\text{Maxpool}(\text{sigma}(h_{2in} \times W_2))]}{\partial W_2} W_3$$

$$= (\hat{y} - y) \frac{\partial}{\partial W_2} \left[\text{Maxpool}(\underbrace{\text{sigma}(h_{2out} \times W_2 + B_2)}_a) \right] W_3$$

$$= (\hat{y} - y) W_3^T \cdot \frac{\partial \text{Maxpool}(a)}{\partial W_2} \cdot \underbrace{\frac{\partial \text{sigma}(b)}{\partial W_2}}_b \times \frac{\partial (h_{2out} W_2)}{\partial W_2}$$

Here 'o' is element wise product

$$\frac{\partial L}{\partial W_2} = (h_{out})^T \cdot (\hat{y} - y) \cdot W_3^T \cdot \frac{\partial \maxpool(a)}{\partial W_2} \cdot \frac{\partial \sigma(b)}{\partial W_2}$$

↳ let us call this (Z_2)

gradient of loss wrt layer 1

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial h_{3in}} \times \frac{\partial h_{3in}}{\partial W_1}$$

~~scribbles~~

$$= (\hat{y} - y) \cdot W_3^T \cdot \frac{\partial [\max(\sigma(h_{out} W_2 + B_2))]}{\partial W_2} \times \frac{\partial (h_{out} W_2 + B_2)}{\partial W_1}$$

$$= Z_2 \times \frac{\partial [\max(\sigma(h_{out} W_1 + B_1))]}{\partial W_1} \cdot W_2$$

$$= Z_2 \times h_{out} \times \frac{\partial \max(a)}{\partial W_1} \times \frac{\partial \sigma(b)}{\partial W_1}$$

$$a = \sigma(h_{out} W_1 + B_1) \quad b = h_{out} W_1 + B_1$$

Here h_{out} is nothing but the ~~my~~ input given at the ~~0th~~ zeroth layer that is the 28x28 image which we consider x .

Gradient w.r.t Biases \rightarrow

$$\begin{aligned}\frac{\partial L}{\partial B_3} &= (\hat{y} - y) \times \frac{\partial h_{3in}}{\partial B_3} \\ &= (\hat{y} - y) \times \frac{\partial h_{2out} w_3 + B_3}{\partial B_3} \\ &= (\hat{y} - y)\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial B_2} &= (\hat{y} - y) \frac{\partial h_{3in}}{\partial B_2} \\ &= (\hat{y} - y) \frac{\partial (h_{2out} w_3 + B_3)}{\partial B_2} \\ &= (\hat{y} - y) \times \frac{\partial \left[(\max(\text{sigma}(h_{2in})) w_3) + B_3 \right]}{\partial B_2}\end{aligned}$$

$$= (\hat{y} - y) \times \frac{\partial \maxpool(h_{1, out} \times W_2 + B_2)}{\partial B_2} \times \frac{\partial \sigma(h_{1, out} \times W_2 + B_2)}{\partial W_2}$$

$$\times \frac{\partial (h_{1, out} \times W_2 + B_2)}{\partial W_2} \text{ [chain rule]} \times W_3^T$$

$$\boxed{\frac{\partial L}{\partial B_2}} = (\hat{y} - y) \times \frac{\partial \maxpool(a)}{\partial B_2} \times \frac{\partial \sigma(b)}{\partial B_2} \times W_3^T$$

$$a = \sigma(h_{1, out} \times W_2 + B_2)$$

$$b = [h_{1, out} \times W_2 + B_2]$$

$$\frac{\partial L}{\partial B_1} = (\hat{y} - y) \times \frac{\partial h_{3, in}}{\partial B_1}$$

$$= \cancel{(\hat{y} - y)} \times \cancel{W_3^T} \times \cancel{W_2}$$

$$= \frac{\partial L}{\partial B_2} \times W_2^T \times \frac{\partial \maxpool(\sigma(h_{1, out} \times W_1 + B_1))}{\partial B_1}$$

2. We have, at the encoder size,

$$h(x) = a(wx + b)$$

$$\hat{x} = g(\hat{a}(w^*h(x) + c))$$

$$\text{Loss function} \rightarrow L(x, \hat{x}) = ||x - \hat{x}||_2^2$$

Now calculating the gradient of L w.r.t weight parameters w and b ,

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{x}} \times \frac{\partial \hat{x}}{\partial h} \times \frac{\partial h}{\partial w}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{x}} \times \frac{\partial \hat{x}}{\partial h} \times \frac{\partial h}{\partial b}$$

$$\frac{\partial L}{\partial \hat{x}} = 2(\hat{x} - x)$$

$$\begin{aligned} \frac{\partial \hat{x}}{\partial h} &= \frac{\partial \hat{a}(w^*h(x) + c)}{\partial (w^*h + c)} + \frac{\partial (w^*h + c)}{\partial h} \\ &= \hat{a}'(w^*h + c) \times w^* \end{aligned}$$

$$\frac{\partial h}{\partial w} = \frac{\partial a(wx+b)}{\partial (wx+b)} \times \frac{\partial (wx+b)}{\partial w}$$

$$= a'(wx+b) \times x^T$$

we know $w^* = w^T$

$$\frac{\partial L}{\partial w} = x^T \times a'(wx+b) \times \hat{a}'(w^* \hat{x} + c) \times 2(\hat{x} - x)$$

$$\frac{\partial L}{\partial b} = a'(wx+b) \times \hat{a}'(w^* \hat{x} + c) \times 2(\hat{x} - x)$$

The gradient of L w.r.t w and b are products of gradients of encoder and decoder.

So we can say gradient ~~w.r.t~~ of L w.r.t w and w^* is sum of ~~encoder~~ gradients of encoder and decoder.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial w^*} = x^T \cdot a'(wx+b) + \hat{a}'(w^* \hat{x} + c) \times 2(\hat{x} - x)$$

3). KL Divergence between two gaussian distr,

$$D_{KL}[N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)] = \frac{1}{2} (\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T$$

$$(\Sigma_1^{-1})(\mu_1 - \mu_0) - k + \log\left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)}\right)$$

where k is dimension of Gaussian distr.

To show this, we see the definition of KL divergence,

$$D_{KL}[N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)] = \int p(x) \log \left[\frac{p(x)}{q(x)} \right] dx$$

where $p(x)$ is true distr of $N(\mu_0, \Sigma_0)$ and $q(x)$ is approximate distr $N(\mu_1, \Sigma_1)$

$$D_{KL}[N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)] = - \int p(x) \log \left[\frac{q(x)}{p(x)} \right] dx$$

$$\text{we know, } q(x) = N(x | \mu_1, \Sigma_1) = \frac{1}{(2\pi)^{k/2} (\Sigma_1)^{1/2}} e^{-\frac{1}{2} \frac{(x - \mu_1)^T (x - \mu_1)}{\Sigma_1}}$$

substituting this expression for $q(x)$,

$$D_{KL}[N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)] = - \int p(x) \log \left[\frac{1}{(2\pi)^{k/2} (\Sigma_1)^{1/2}} e^{-\frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)} \right] \frac{1}{(2\pi)^{k/2} (\Sigma_0)^{1/2}} e^{-\frac{1}{2} (x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0)} dx$$

Simplifying the above expressions,

$$D_{KL}[N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)] = \frac{1}{2} [\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T$$

$$\Sigma_1^{-1} (\mu_1 - \mu_0) - k + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right)]$$

$$\rightarrow D_{KL}[N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)] = \frac{1}{2} [\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T$$

$$\Sigma_1^{-1} (\mu_1 - \mu_0) - k + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right)]$$

$$\rightarrow D_{KL}[N(\mu_1, \Sigma_1) \| N(\mu_0, \Sigma_0)] = \frac{1}{2} [\text{tr}(\Sigma_0^{-1} \Sigma_1) + (\mu_0 - \mu_1)^T$$

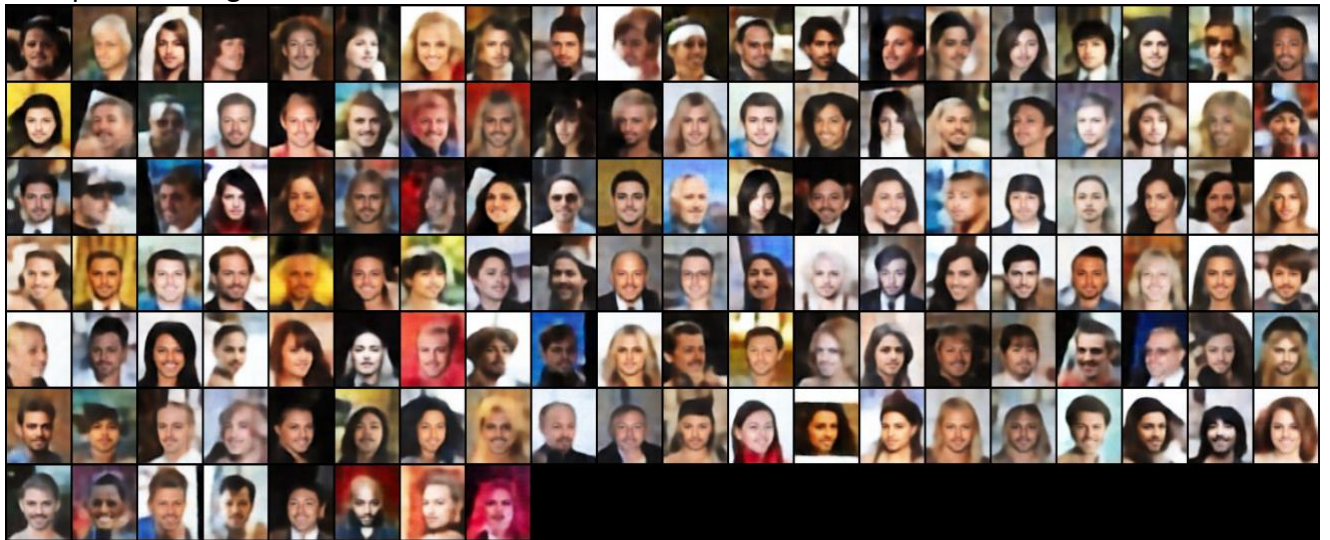
$$\Sigma_0^{-1} (\mu_0 - \mu_1) - k + \log \left(\frac{\det \Sigma_0}{\det \Sigma_1} \right)]$$

We see that the two k dimensional distⁿ is symmetric with only the subscripts of μ and Σ flipped.

4. Wrote the code for **TASK 1** and ran the model for **8 epochs** due to processing constraints. The batch size was **128** and learning rate was **0.001**.

TASK 2:

Manipulated images with Beard



Original Images



Original Images



Manipulated images for smiling faces



Manipulated Images for eyeglasses



Original Images



TASK 3:

Took 2 images and got the following interpolation for 10 samples.

