

CSE 676 : Assignment #2

Due Date: 11:59PM April 25, 2021

Q1.1 Padding size = 1

$$x = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 3 & 4 & 0 \\ 0 & 5 & 6 & 7 & 8 & 0 \\ 0 & 9 & 10 & 11 & 12 & 0 \\ 0 & 13 & 14 & 15 & 16 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

→ with padding.

as stride = 2, convolutional output is,

$$O_{ij} = (x * w)_{ij} \triangleq \sum_k \sum_l x_{k+i} \cdot l_{j+l} w_{kl}$$

assuming matrix indices start from 0.

$$\begin{aligned} O_{0,0} &= x_{0,0} \cdot w_{0,0} + x_{0,1} \cdot w_{0,1} + x_{1,0} \cdot w_{1,0} + x_{1,1} \cdot w_{1,1} \\ &= 0.4 + 0.2 + 0.5 + 1.1 \\ &= 1 \end{aligned}$$

moving filter by 2 positions to the right,

$$\begin{aligned} O_{0,1} &= x_{0,2} \cdot w_{0,0} + x_{0,3} \cdot w_{0,1} + x_{1,2} \cdot w_{1,0} + x_{1,3} \cdot w_{1,1} \\ &= 0.4 + 0.2 + 2.5 + 3 \cdot 1 \end{aligned}$$

111th

$$\begin{aligned} O_{0,2} &= x_{0,4} \cdot w_{0,0} + x_{0,5} \cdot w_{0,1} + x_{1,4} \cdot w_{1,0} + x_{1,5} \cdot w_{1,1} \\ &= 0.4 + 0.2 + 4.5 + 0.1 = 20 \end{aligned}$$

$$O_{1,0} = x_{2,0} \cdot w_{0,0} + x_{2,1} \cdot w_{0,1} + x_{3,0} \cdot w_{1,0} + x_{3,1} \cdot w_{1,1}$$

$$= 0.4 + 0.5 \cdot 2 + 0.1 + 0.5$$

$$= 1.9$$

$$O_{1,1} = x_{2,2} \cdot w_{0,0} + x_{2,3} \cdot w_{0,1} + x_{3,2} \cdot w_{1,0} + x_{3,3} \cdot w_{1,1}$$

$$= 6 \cdot 4 + 7 \cdot 2 + 10 \cdot 5 + 11 \cdot 1$$

$$= 99$$

$$O_{1,2} = x_{2,4} \cdot w_{0,0} + x_{2,5} \cdot w_{0,1} + x_{3,4} \cdot w_{1,0} + x_{3,5} \cdot w_{1,1}$$

$$= 8 \cdot 4 + 0.2 + 12 \cdot 5 + 0 \cdot 1$$

$$= 92$$

$$O_{2,0} = x_{4,0} \cdot w_{0,0} + x_{4,1} \cdot w_{0,1} + x_{5,0} \cdot w_{1,0} + x_{5,1} \cdot w_{1,1}$$

$$= 0 \cdot 4 + 13 \cdot 2 + 0 \cdot 5 + 0 \cdot 1$$

$$= 26$$

$$O_{2,1} = x_{4,2} \cdot w_{0,0} + x_{4,3} \cdot w_{0,1} + x_{5,2} \cdot w_{1,0} + x_{5,3} \cdot w_{1,1}$$

$$= 14 \cdot 4 + 15 \cdot 2 + 0 \cdot 5 + 0 \cdot 1 = 86$$

$$= 86$$

$$O_{2,2} = x_{4,4} \cdot w_{0,0} + x_{4,5} \cdot w_{0,1} + x_{5,4} \cdot w_{1,0} + x_{5,5} \cdot w_{1,1}$$

$$= 16 \cdot 4 + 0 \cdot 2 + 0 \cdot 2 + 0 \cdot 1 = 64$$

output = $\begin{bmatrix} 1 & 19 & 26 \\ 19 & 99 & 92 \\ 26 & 86 & 64 \end{bmatrix}$

$$\text{The relation} \Rightarrow O = \frac{L + 2(P) - F}{S} + 1$$

As stride value increases the value skips decreases by same ratio. Both row and column at both ends are added due to padding and filter size decreases the number by its size.

2. (ii) Given :- input = 32×32

layer 1 \rightarrow 6 filters

Each filter \rightarrow dimension of 5×5

~~Number of parameters per filter~~

$$\text{No. of parameters per filter} = (5 \times 5) + 1 = 26$$

$$\text{Total no. of parameters} = 26 \times \text{No. of filters}$$

$$= 26 \times 6$$

$$= 156$$

Dimension of o/p

$$\text{o/p dimension of each filter} = O = \frac{L + 2(P) - F}{S} + 1$$

$$O = \frac{32 + 2(0) - 5}{4} + 1$$

$$O = 28$$

Each filter will output a 28×28 matrix

6 filters $\rightarrow 28 \times 28 \times 6$ output dimension.

(2) So the LeNet Architecture has 2 convolutions layer followed by sub Sampling layer which is the max pooling layers and after flattening the output of second convolutional block it is fed to 2 fully connected layers.

We can take the following steps to reduce the no. of parameters. Due to the accuracy trade off and also poorer performance due to the model complexity they can lead to poorer performance:

→ Reducing kernel size to 4×4 .

~~Reduction of no. of parameters per filter~~

→ Change value of stride 1 to 2

→ Reducing no. of kernel : 6 to 5

According to above changes,

$$\bullet \text{No. of parameters per filter} = (4 \times 4) * 1 = 16$$

$$\therefore \text{Total no. of parameters} = 16 \times \text{No. of filters}$$
$$= 16 \times 5$$

$$\therefore \text{O/P dimension of each filter} = O = \frac{L + 2(P) - F}{S} + 1$$

$$O = \frac{32 + 2(0) - 4}{2} + 1$$

$$O = 15$$

• 15×15 matrix is the o/p of each filter.

• 5 filters therefore o/p dimension is $15 \times 15 \times 5$

• The direct effect of decrease in o/p dimensionality is the number of parameters of further layers.

- The number of parameters has gone from 156 to 85 in just the first layer by 2 small changes that were made.

Q2.1 Given: Gradients calculated for timestamp $t+1$.

Need to calculate backprogs pageite based on given values and need to calculate the gradients for time t

First let us calculate $\frac{\partial L}{\partial h_t}$ in terms of gradients from timestamp $t+1$. And then using the value of $\frac{\partial L}{\partial h_t}$ to update all the gradient value at time t .

In the LSTM architecture can be seen from the figure the Back propagation through time takes place by 2 nodes i.e c_{t+1} and h_t .

So, now solving the equations at timestamp $t+1$ to calculate the gradient value for h_t which would be required for the calculation of all other gradient at time t .

$$\Rightarrow z_t = \begin{pmatrix} i_t \\ f_t \\ \hat{o}_t \\ \hat{a}_t \end{pmatrix} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$z_t = w \times I_t$$

$$\hat{i}_t = \sigma(\hat{i}_t)$$

$$f_t = \sigma(\hat{f}_t)$$

$$\hat{o}_t = \sigma(\hat{o}_t)$$

$$g_t = \sigma(\hat{g}_t)$$

so we have $\frac{\partial L}{\partial i_{t+1}}$, $\frac{\partial L}{\partial f_{t+1}}$, $\frac{\partial L}{\partial o_{t+1}}$, $\frac{\partial L}{\partial g_{t+1}}$, $\frac{\partial L}{\partial c_{t+1}}$ and $\frac{\partial L}{\partial h_{t+1}}$ i.e

$$\delta i_{t+1}, \delta f_{t+1}, \delta o_{t+1}, \delta g_{t+1}, \delta c_{t+1} \text{ and } \delta h_{t+1}$$

To calculate δh_t , we are applying partial derivatives in equation 2 for time $t+1$

$$\delta \hat{i}_{t+1} = \delta i_{t+1} \odot i_t \odot (1 - i_{t+1})$$

$$\delta \hat{f}_{t+1} = \delta f_{t+1} \odot f_t \odot (1 - f_{t+1})$$

$$\delta \hat{o}_{t+1} = \delta o_{t+1} \odot o_t \odot (1 - o_{t+1})$$

$$\delta \hat{g}_{t+1} = \delta g_{t+1} \odot (1 - \tanh^2(\hat{g}_{t+1}))$$

So now,

$$\delta z_{t+1} = \begin{pmatrix} \delta \hat{i}_{t+1} \\ \delta \hat{f}_{t+1} \\ \delta \hat{o}_{t+1} \\ \delta \hat{g}_{t+1} \end{pmatrix}$$

we get

$$\delta I_{t+1} = W^T \times \delta z_{t+1}$$

$$I_{t+1} = \begin{pmatrix} h_t \\ x_{t+1} \end{pmatrix}$$

$$\delta I_{t+1} = \begin{pmatrix} \delta h_t \\ \delta x_{t+1} \end{pmatrix}$$

we get δh_t from above matrix which is being back propagated by further time steps.

So Secondly, considering the state of the cell at time t. The δh_t will contain 2 parts i.e the backpropagated error calculated earlier and the error difference calculated for that particular time stamp.

$$\Rightarrow \delta h_t = \Delta_t + \delta h_t (\text{from time } t+1)$$

Now solving the equations for time t

The front pass equations are given as -

$$c_t = f_t \odot c_{t+1} + i_t \odot g_t$$

$$\rightarrow h_t = b_t \odot \tanh(c_t)$$

Since we have δh_t , we can use it to calculate δc_t , δf_t , δo_t , δg_t and δi_t .

Starting with eq. 9. calculating δc_t . Using partial derivatives.

$$\text{So unknown is } \delta o_t = \frac{\partial L}{\partial o_t} = \frac{\partial L}{\partial h_t} \cdot \frac{\partial h_t}{\partial o_t}$$

$$\text{function for } o_t \text{ is } o_t = \tanh(h_t) \text{ so } \frac{\partial h_t}{\partial o_t} = 1 - \tanh^2(o_t)$$

$$\text{so } \delta o_t = \delta h_t \cdot \tanh(o_t)$$

$$\text{so } \delta o_t = \delta h_t \odot \tanh(o_t)$$

Now, calculating δc_t using eq. 9. The back propagation path to c_t comes from h_t and c_{t+1} .

$$\text{So, } \delta c_t = \frac{\partial L}{\partial c_t} = \frac{\partial L}{\partial h_t} \cdot \frac{\partial h_t}{\partial c_t} + \frac{\partial L}{\partial c_{t+1}} \cdot \frac{\partial c_{t+1}}{\partial c_t}$$

$$\delta c_t = \delta h_t \odot o_t \odot (1 - \tanh^2(o_t)) + \delta c_{t+1} \cdot f_t$$

Now, calculating gradients i.e δi_t , δf_t and δg_t as they depend on δc_t .

$$\delta i_t = \frac{\partial L}{\partial i_t} = \frac{\partial L}{\partial c_t} \cdot \frac{\partial c_t}{\partial i_t}$$

$$\delta i_t = \delta c_t \cdot g_t$$

$$\delta i_t = \delta h_t \odot g_t$$

Similarly,

$$\delta f_t = \frac{\partial L}{\partial f_t} = \frac{\partial L}{\partial c_t} \cdot \frac{\partial c_t}{\partial f_t}$$

$$\delta f_t = \delta c_t \cdot o_{ct-1}$$

$$\delta f_t = \delta h_t \odot o_{ct-1}$$

Similarly,

$$\delta g_t = \frac{\partial L}{\partial g_t} = \frac{\partial L}{\partial c_t} \cdot \frac{\partial c_t}{\partial g_t}$$

$$\delta g_t = \delta c_t \cdot i_t$$

$$\delta g_t = \delta h_t \odot i_t$$

So, to summarize the solution. We obtain the value of δh_t . The value is calculated as the sum of gradient calculated at that times stamp and the back propagated value from δI_{t+1} , which is a stack of δh_t .

$$\delta x_t$$

$$\delta o_t = \delta h_t \odot \tanh(c_t)$$

$$\delta c_t = \delta h_t \odot o_t \odot (1 - \tanh^2(c_t)) + \delta c_{t+1} \odot f_t$$

$$\delta i_t = \delta h_t \odot g_t$$

$$\delta f_t = \delta h_t \odot c_t$$

$$\delta g_t = \delta h_t \odot i_t$$

Q3. (1) The (a), (b) and (c) can be learned using the standard VAE.

(ex. spatial + temporal latents)

(2) The standard VAE is based on the assumption that the latent space variables are a Gaussian Distribution. This assumption is made to simplify the calculation of Variational Inference problem.

so, the fourth distribution is not a Gaussian distribution
so it can't be learnt

$$KL\left(\frac{q(x)}{p(z|x)}\right) = -\sum q(x) \log \left(\frac{p(z|x)}{q(z)} \right)$$

$$= -\sum q(x) \log \left(\frac{p(x,z)}{q(z) * p(x)} \right)$$

$$= -\sum q(x) \left(\log \left(\frac{p(x,z)}{q(z)} \right) + \log \left(\frac{1}{p(x)} \right) \right)$$

$$\text{so we want to minimize } -\sum q(x) \left(\log \left(\frac{p(x,z)}{q(z)} \right) + \log \left(\frac{1}{p(x)} \right) \right)$$

$$\text{so we want to maximize } \sum q(x) \log \left(\frac{p(x,z)}{q(z)} \right) - \log(p(x))$$

$$\text{so we want to maximize } \sum q(x) \log \left(\frac{p(x,z)}{q(z)} \right) + \log(p(x))$$

$$\text{so we want to maximize } \sum q(x) \log \left(\frac{p(x,z)}{q(z)} \right) + \log(p(x))$$

$$\text{so we want to maximize } \sum q(x) \log \left(\frac{p(x,z)}{q(z)} \right) + \log(p(x))$$

$$KL\left(\frac{q(x)}{p(z|x)}\right) = - \sum q(x) \left(\log \left(\frac{p(x,z)}{q(z)} \right) \right) + \log(p(x))$$

$$KL\left(\frac{q(x)}{p(z|x)}\right) = - \text{Variational Lower bound} + \log(p(x))$$

In the above equation the $\log(p(x))$ is the constant term. here we are trying to minimize the KL divergence.

Minimizing KL Divergence is same as maximizing Lower Bound.

$$\begin{aligned} \text{Lower Bound} &= \sum q(x) \left(\log \left(\frac{p(x,z)}{q(z)} \right) \right) \\ &= \sum q(x) \left(\log \left(\frac{p(x|z) \cdot p(z)}{q(z)} \right) \right) \\ &= \sum q(x) \left(\log(p(x|z)) + \log \left(\frac{p(z)}{q(z)} \right) \right) \\ &= \sum q(x) \log(p(x|z)) - D_{KL}[q(z) || p(z)] \\ &= E_{q(z)} \log(p(x|z)) - D_{KL}[q(z) || p(z)] \end{aligned}$$

Now, while optimizing we try to maximize the lower Bound. The first term corresponds to the reconstruction loss ie how close is the reconstructed output with respect to the input. Since here we assume that $q(z)$ distribution is Gaussian, the second term corresponding to the KL Divergence tries to Normalize the latent space vector as close to the Gaussian Distribution as possibly by reducing the divergence.

Q4.1. Generative Adversarial Nets (GANs)

- Method to generate samples (Based on game theory)
 - they generate new samples from the same data distribution
- training data $\sim P_{\text{data}}(x)$; generated samples $\sim P_{\text{model}}(x)$
the goal is to learn $P_{\text{model}}(x)$ such that it is close to $P_{\text{data}}(x)$
- GAN training framework: It corresponds to a minmax two-player game

- For arbitrary G and D, a unique solution exists.

- If G and D are MLP's can be trained using backprop.

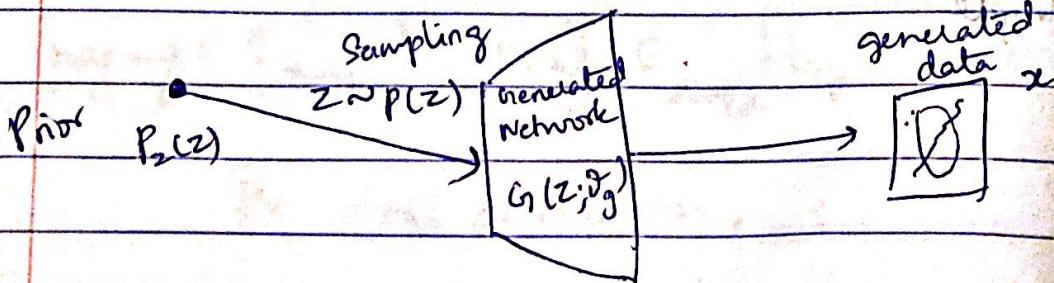
In the space of arbitrary function G and D, a unique sol. exists.
Resulting G is used to generate images.

Generator Network: To learn the generator distribution

p_g over data x , we define a prior on input noise variables $p_z(z)$.

Represent a mapping to data spaces as $G(z; \theta_g)$.

where $g \rightarrow$ differentiable generator functions.



- Generator Network produces samples
 $x = G(z; \theta_g)$
- Discriminator network attempts to distinguish samples drawn from training data and samples drawn from generator and also emits a probability value - $D(x; \theta_d)$
- We also have a second multilayer perceptron $D(x; \theta_d)$ that outputs a single scalar.

Training D & G_i : Trains to maximize the prob. of assigning the correct label to both training examples and samples from G_i

Simultaneously train G_i to minimize $\log(1 - D(G_i(z)))$

$$\min_{G_i} \max_{D} V(D, G_i) = E_{x \sim p_{\text{data}}(x)} [\log D(x)] + E_{z \sim p_{G_i}(z)} [\log(1 - D(G_i(z)))]$$

D is accurate over real data by minimizing

$$E_{x \sim p_{\text{data}}(x)} [\log D(x)]$$

G_i is trained to increase D 's prob. for a fake example.

Combined together G and D are playing a minmax game in which we optimize:

- To find best value for D consider

$$L(G, D) = \int_x [P_{\text{data}}(x) \log D(x) + P_g(x) \log (1 - D(x))] dx$$

- Setting the derivative to zero

$$D^*(x) = x^+ = \frac{A}{A+B} = \frac{P_r(x)}{P_r(x) + P_g(x)} \in [0, 1]$$

- Once G is trained to optimal, P_g is very close to

$$P_{\text{data}}$$

- When $P_g = P_{\text{data}}$, $D^*(x)$ become $1/2$

D is trained to discriminate samples from data,

$$\text{converging to } D^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}$$

After an update of G, gradient of D has guided G(z) to flow to regions that are more likely to be classified as data.

* GAN Objective: Value function of Minmax played by G, D.

$$\text{minmax } \underset{G, D}{E}_{x \sim p_x} [\log D(x)] + E_{z \sim p_z} [\log (1 - D(G(z)))]$$

p_x : data distribution,

$P_{G(z)}$: model distribution; z modeled as Gaussian.

→ We must implement the game using an iterative numerical approach.

Optimizing D to completion in inner loop of training is computationally prohibitive and results in overfitting in finite datasets.

Instead we alternate between K steps of optimizing D and one step of optimizing G .

Alternating gradient updates :-

1: fix G and perform a gradient step for maximization

$$\max_D E_{x \sim p_X} [\log D(x)] + E_{z \sim p_Z} [\log (1 - D(G(z)))]$$

2: fix D and perform a gradient step for minimization

In theory: treat $\nabla_{\theta} \log D(\theta)$ as a batch of (s, r)

$$\min_G E_{z \sim p_Z} [\log (1 - D(G(z)))]$$

In practice usually initialize D with ReLU :

$$\max_G E_{z \sim p_Z} [\log b(G(z))]$$

2. Mode collapse happens when few modes generated
Partial mode collapse, a hard problem to solve
in GAN. A complete collapse is not common but
a partial collapse happens often.

The Reason for Mode Collapse:

Objective of G : create images to fool D :

$$\nabla_{\mathbf{z}} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)})))$$

Case when G is trained without update to D .

- Generated images converge to \mathbf{x}^* that fool D the most -- most realistic from the D perspective.

• In this extreme, \mathbf{x}^* will be independent of \mathbf{z} .

$$\mathbf{x}^* = \text{argmax}_{\mathbf{x}} D(\mathbf{x})$$

Mode collapses to a single point

• Gradient associated with \mathbf{z} approaches zero $\frac{\partial J}{\partial \mathbf{z}} \approx 0$

When we restart the training in D , the most effective way to detect generated images is to detect this single mode.

- Since G ~~desensitizes~~ already desensitizes the impact of \mathbf{z} already, the gradient from D will likely push the single point around for the next most vulnerable mode.

This is not good to find. It produces such an imbalance of modes in training that it deteriorates its capability to detect others.

Now, both network are overfitted to exploit the short term opponent weakness. This turns into a cat and mouse chase and the model will not converge.

at what point the two reward are 0?

at what point the two reward is equal to zero?

at what point the two reward is different?

(x,y) = (0,0)

many signs of convergence start

as the two reward becomes more and more similar

and also when one of the reward starts to move towards the other reward then the two reward becomes more and more similar

so the two reward becomes more and more similar

from time to time both reward will diverge like

about what extent the two reward moves