

q1

Sunday, February 28, 2021 10:57 AM

Question 1:

*Note: For these questions I wrote a Python script for gradient descent with GeLU/momentum. to check my results. They are located in the files "gelu.py", "momentum.py".

1a)**Calculation by Hand**

HW1 #2 Note: bin=0.5

GD

GeLU(x) $\approx x \delta(1.702x) + x \frac{1}{1+e^{-1.702x}}$

③ Initial guess: $x_0 = 0$

a) GD equation

$x_1 \leftarrow x_0 - \eta \text{ grad}, \quad \eta = 0.1$

Find gradient

$$\frac{\partial \text{GeLU}(x)}{\partial x} = [1 \cdot \delta(1.702x)] + [\delta(1.702x)(1 - \delta(1.702x)) 1.702 \cdot x]$$

$$= \delta(1.702x) + \delta(1.702x)[1.702x \cdot (1 - \delta(1.702x))]$$

$$= \delta(1.702x)[1 + 1.702x(1 - \delta(1.702x))]$$

Thus $\frac{\partial \text{GeLU}(x)}{\partial x} \Big|_{x=0} = \delta(0)[1 + 0] = 0.5 \cdot 1 = 0.5$

- Find x_1 : $x_1 \leftarrow x_0 - \eta \text{ grad}$
 $x_1 \leftarrow 0 - (0.1 \cdot 0.5) = \boxed{-0.05}$, $\delta(1.702 \cdot -0.05) \cdot \delta(-0.05) = 0.478737$
- Find x_2 : $x_2 \leftarrow x_1 - \eta \text{ grad}$
 $x_2 \leftarrow -0.05 - (0.1 \cdot \delta(1.702 \cdot -0.05)[1 + 1.702 \cdot -0.05(1 - \delta(1.702 \cdot -0.05))])$
 $x_2 \leftarrow \boxed{-0.09575}$, $\delta(1.702 \cdot -0.09575) = 0.45935$
- Find x_3 : $x_3 \leftarrow x_2 - \eta \text{ grad}$
 $x_3 \leftarrow -0.09575 - (0.1 \cdot 0.45935) = \boxed{-0.137637}$

So $\text{GeLU}(x_1) = -0.05 \cdot \delta(1.702 \cdot -0.05) = \boxed{-0.0239}$
 $\text{GeLU}(x_2) = \boxed{-0.04398}$
 $\text{GeLU}(x_3) = -0.137637 \cdot \delta(1.702 \cdot -0.137637) = \boxed{-0.0608}$

Output of Code

```
Using lr: 0.1, x_0:0, num_iters: 3
x_1: -0.05
GeLU(x_1): -0.02393689150943369
x_2: -0.09575013021851926
GeLU(x_2): -0.04398265469655064
x_3: -0.13763771842581057
GeLU(x_3): -0.06079478852684933
```

1b)**Calculation by Hand**

b) $\eta = 1$

- $x_1 \leftarrow 0 - (1 \cdot \text{grad}) = 0 - [\delta(0)[1 + 0]] = \boxed{-0.5}$
- $x_2 \leftarrow -0.5 - (1 \cdot \delta(1.702 \cdot -0.5)[1 + 1.702 \cdot -0.5 \cdot (1 - \delta(1.702 \cdot -0.5))]) = \boxed{-0.6206}$
- $x_3 \leftarrow -0.6206 - (1 \cdot \delta(1.702 \cdot -0.6206)[1 + 1.702 \cdot -0.6206(1 - \delta(1.702 \cdot -0.6206))]) = \boxed{-0.6764}$

$$\text{GeLU}(x_1) = -0.5 \delta(1.702 \cdot -0.5) = \boxed{-0.1495}$$

$$\text{GeLU}(x_2) = -0.6206 \delta(1.702 \cdot -0.6206) = \boxed{-0.1601}$$

$$\text{GeLU}(x_3) = -0.6764 \delta(1.702 \cdot -0.6764) = \boxed{-0.1625}$$

A larger step size η means a steeper gradient, so it takes bigger steps for each iteration of GD \Rightarrow finds optimum faster.

Output of Code

```
Using lr: 1, x_0:0, num_iters: 3
x_1: -0.5
GeLU(x_1): -0.1496115633936199
x_2: -0.6207780880345858
GeLU(x_2): -0.1601399925974373
x_3: -0.676497308917168
GeLU(x_3): -0.16251748448720163
```

A larger step size means that a bigger step is taken for each iteration of gradient descent. This means it will move towards the minimum faster, which can be seen by comparing $\text{GeLU}(x_3)$. From part a), $\text{GeLU}(x_3) = -0.0608$, which is less optimal than $\text{GeLU}(x_3) = -0.1625$ when the learning rate is 1. However, a higher learning rate poses the risk of overshooting near the minimum.

1c)**GeLU, x_0 = -3 (No momentum)**

```
Using lr: 0.1, x_0:-3, num_iters: 3
x_1: -2.997545167609435
GeLU(x_1): -0.01813166529499688
x_2: -2.995082708753491
GeLU(x_2): -0.018192396733478253
x_3: -2.9926125791343785
GeLU(x_3): -0.018253507384506158
```

GeLU, x_0 = -3 (With Momentum)

GD with Momentum reaches a more optimal point/extremum (-0.0184 < -0.01825) in the same number of iterations, so it is generally faster.

GD with Momentum

```
lr: 1, x_0: -3, num_iters: 3, B: 0.9
x_1: -2.997545167609435
GeLU(x_1): -0.01813166529499688
x_2: -2.9928733596019823
GeLU(x_2): -0.01824704671202677
x_3: -2.986191702877126
GeLU(x_3): -0.01841325299731993
```

q1_code

Sunday, February 28, 2021 11:04 PM

Gradient Descent code for Q1 (gelu.py)

```
import numpy as np
import math

# Gradient Descent, GeLU Activation
def sigmoid(x):
    return 1 / (1 + math.exp(-x))

def gelu(x):
    return x * sigmoid(1.702*x)

# gradient
def grad(x):
    return sigmoid(1.702*x) * (1+1.702*x*(1-sigmoid(1.702*x)))

# Set up variables and run GD
def main():
    lr = 0.1                                # learning rate
    x = {0:-3,1:0,2:0,3:0}                   # store x[i]
    gelu_dict = {0:0, 1:0, 2:0, 3:0}          # store gelu(x[i])
    num_iters = 3

    for i in range(num_iters):
        x[i+1] = x[i] - lr * (grad(x[i]))    # update x[i+1]
        gelu_dict[i+1] = gelu(x[i+1])          # get gelu of new x[i+1]

if __name__ == "__main__":
    main()
```

GD with Momentum code (momentum.py)

```
# Set up variables and run perceptron
def main():
    B = 0.9                                # Set up variables
    lr = 1
    x = {0:-3,1:0,2:0,3:0}                  # Store x0, x1, x2, x3
    v = {0:grad(x[0]), 0:0, 0:0, 0:0}       # Store v0, v1, v2, v3
    gelu_dict = {0:0, 1:0, 2:0, 3:0}          # Store GeLU of x0, x1, x2, x3
    num_iters = 3

    print("GD with Momentum")
    print(f"lr: {lr}, x_0: {x[0]}, num_iters: {num_iters}, B: {B}")

    for i in range(3):
        v[i+1] = (B*v[i]) + ((1-B)*grad(x[i])) # update v[i+1]
        x[i+1] = x[i] - (lr * v[i+1])           # update x[i+1]
        gelu_dict[i+1] = gelu(x[i+1])            # get gelu of new x[i+1]
        print(f"GeLU(x_{i+1}): ", gelu_dict[i+1])

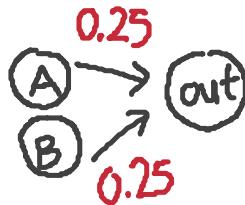
if __name__ == "__main__":
    main()
```

q2

Sunday, February 28, 2021 11:28 AM

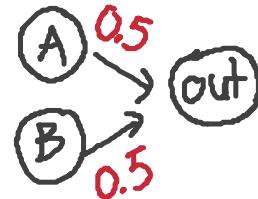
Note: all hidden nodes/output nodes have sigmoid activations, where output is 1 if the input ≥ 0.5 .

a) AND gate



Both A, B must be 1 for the output to be 1
 $0.25*1 + 0.25*1 \geq 0.5 \rightarrow \text{output} = 1$

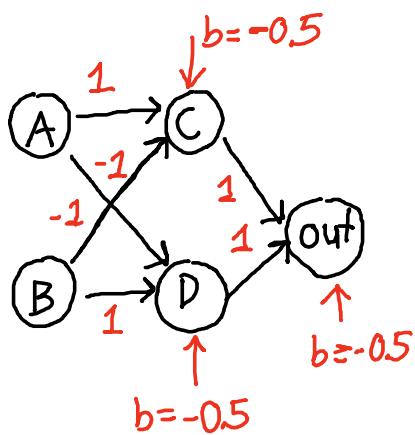
b) OR gate



Only one of A, B needs to be 1 for the output to be 1
 $0.25*1 + 0.25*0 \geq 0.5 \rightarrow \text{output} = 1$

c) XOR gate

Exactly one of A,B must be 1 for the output to be 1



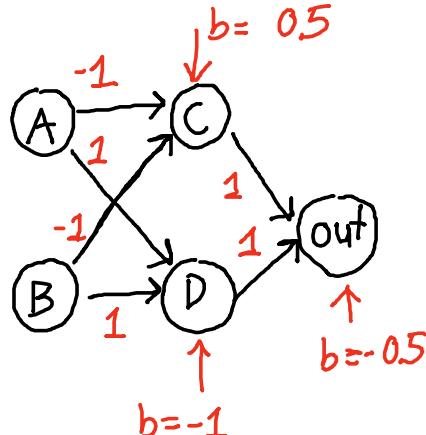
$A=0, B=0$	$A=0, B=1$
$C = \sigma(0+0-0.5) = 0$	$C = \sigma(0-1-0.5) = 0$
$D = \sigma(0+0-0.5)=0$	$D = \sigma(0+1-0.5)=1$
$\text{Out} = \sigma(0+0-0.5)=0$	$\text{Out} = \sigma(0+1-0.5)=1$

$A=1, B=0$	$A=1, B=1$
$C = \sigma(1+0-0.5) = 1$	$C = \sigma(1-1-0.5) = 0$
$D = \sigma(-1+0-0.5)=0$	$D = \sigma(-1+1-0.5)=0$
$\text{Out} = \sigma(1+0-0.5)=1$	$\text{Out} = \sigma(0+0-0.5)=0$

d)

XNOR gate

A,B must both be 1 or both be 0 for the output to be 1



$A=0, B=0$	$A=0, B=1$
$C = \sigma(0+0+0.5) = 1$	$C = \sigma(0-1+0.5) = 0$
$D = \sigma(0+0-1)=0$	$D = \sigma(0+1-1)=0$
$\text{Out} = \sigma(1+0-0.5)=1$	$\text{Out} = \sigma(0+0-0.5)=0$

$A=1, B=0$	$A=1, B=1$
$C = \sigma(-1+0+0.5) = 0$	$C = \sigma(-1-1+0.5) = 0$
$D = \sigma(1+0-1)=0$	$D = \sigma(1+1-1)=1$
$\text{Out} = \sigma(0+0-0.5)=0$	$\text{Out} = \sigma(0+1-0.5)=1$

q3

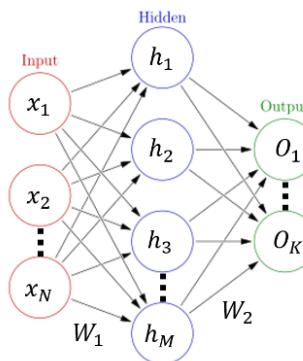
Sunday, February 28, 2021 11:54 AM

$$f_1 = xW_1 + b_1$$

$$a = \sigma(f_1)$$

$$f_2 = aW_2 + b_2$$

$$o = S(f_2)$$



$$\begin{aligned} x &= [1 \times N] & f_2 &= [1 \times K] \\ h &= [1 \times M] \\ W_1 &= [N \times M], \quad W_2 = [M \times K] \\ o &= [1 \times K] \\ f_1 &= [1 \times M] \\ a &= [1 \times M] \end{aligned}$$

a)

$$\begin{aligned} \frac{\partial E(o)}{\partial f_2} &= \left(\frac{\partial E(o)}{\partial o} \right) \left(\frac{\partial o}{\partial f_2} \right) \\ &= \left[-y_1/o_1 \quad -y_2/o_2 \dots \right] \begin{bmatrix} o_1(1-o_1) & \dots & -o_1 o_K \\ -o_2 o_1 & o_2(1-o_2) \dots & -o_2 o_K \\ \vdots & & \vdots \\ -o_K o_1 & \dots & o_K(1-o_K) \end{bmatrix} \\ &= \left[-y_1 + o_1(\sum_k y_k) \quad -y_2 + o_2(\sum_k y_k) \dots -y_K + o_K(\sum_k y_k) \right] \\ &= \boxed{\left[o_1 - y_1 \quad o_2 - y_2 \quad o_3 - y_3 \dots \quad o_K - y_K \right]} \end{aligned}$$

$$\begin{aligned} E(o) &= -\sum y_i \log o_i \\ &= -(y_1 \log o_1 + \dots + y_K \log o_K) \end{aligned}$$

$$\begin{aligned} \frac{\partial E(o)}{\partial o} &= \left[\frac{-y_1}{o_1} \quad \frac{-y_2}{o_2} \dots \quad \frac{-y_K}{o_K} \right] \\ \frac{\partial o}{\partial f_2} &= \underbrace{\begin{bmatrix} o_1(1-o_1) & -o_1 o_2 & \dots & -o_1 o_K \\ -o_2 o_1 & o_2(1-o_2) & \dots & -o_2 o_K \\ \vdots & & \ddots & \vdots \\ -o_K o_1 & \dots & o_K(1-o_K) & \end{bmatrix}}_{\text{Derivative of softmax}} \end{aligned}$$

$$b) \quad \frac{\partial E(o)}{\partial x} = \frac{\partial E(o)}{\partial f_2} \frac{\partial f_2}{\partial a} \frac{\partial a}{\partial f_1} \frac{\partial f_1}{\partial x}$$

$$\text{Since } f_2 = aW_2 + b_2, \quad \frac{\partial f_2}{\partial a} = W_2^T$$

$$\text{Since } f_1 = xW_1 + b_1, \quad \frac{\partial f_1}{\partial x} = W_1^T$$

$$\frac{\partial a}{\partial f_1} = \frac{\partial (\sigma(f_1))}{\partial f_1} = \begin{matrix} M \\ \begin{bmatrix} a_1(1-a_1) & \dots & \dots & 0 \\ 0 & a_2(1-a_2) & \dots & \vdots \\ \vdots & & \ddots & a_M(1-a_M) \\ 0 & \dots & \dots & \end{bmatrix} \end{matrix}$$

$$\text{So} \quad \frac{\partial E(o)}{\partial x} = \boxed{\left[o_1 - y_1 \quad o_2 - y_2 \dots o_K - y_K \right] W_2^T \begin{bmatrix} a_1(1-a_1) & 0 & \dots & 0 \\ 0 & a_2(1-a_2) & \dots & \vdots \\ \vdots & & \ddots & a_M(1-a_M) \end{bmatrix} W_1^T}$$

q4

Sunday, February 28, 2021 12:36 PM

④ CNN

	4×4	$f_1 = 3 \times 3$	$f_2 = 2 \times 2$
	$\begin{array}{ c c c c } \hline 3 & 5 & 2 & 3 \\ \hline 9 & 1 & 8 & 4 \\ \hline 6 & 4 & 3 & 7 \\ \hline 7 & 0 & 2 & 4 \\ \hline \end{array}$	*	$\begin{array}{ c c c } \hline -1 & 0.5 & -2 \\ \hline 2 & 0 & 1 \\ \hline 0 & 1 & 1.5 \\ \hline \end{array}$

a) convolve F with f to get F'

$$n=4$$

$$f=3$$

want output size: $\left(\frac{n+2p-f}{s} + 1\right)$

$$\left(\frac{4+2p-3}{s} + 1\right) = 4 \Rightarrow p=1$$

zero pad image

$$\begin{array}{|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 3 & 5 & 2 & 3 & 0 \\ \hline 0 & 9 & 1 & 8 & 4 & 0 \\ \hline 0 & 6 & 4 & 3 & 7 & 0 \\ \hline 0 & 7 & 0 & 2 & 4 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline -1 & 0.5 & -2 \\ \hline 2 & 0 & 1 \\ \hline 0 & 1 & 1.5 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 15.5 & 21 & 27 & 8 \\ \hline 4.5 & 30 & 9.5 & 22.5 \\ \hline 13.5 & -6.5 & 18 & 4 \\ \hline -5 & 6 & -2.5 & 4.5 \\ \hline \end{array} = F'$$

b) Filter 2, $p=1$.

$$n=4$$

$$f=2$$

$$\left(\frac{n+2p-f}{s} + 1\right) = 4$$

$$\left(\frac{4+2-2}{s} + 1\right) = 4$$

$$\frac{4+2-2}{s} = 3$$

$$3s=4 \rightarrow s=\underline{4/3}$$

No, since the stride can't be fractional.

c) Average pool

$$\begin{array}{|c|c|} \hline \frac{15.5+21+4.5+30}{4} & \frac{27+8+9.5+22.5}{4} \\ \hline \frac{13.5-6.5-5+6}{4} & \frac{18+4-12.5+4.5}{4} \\ \hline \end{array} = \begin{array}{|c|c|} \hline 11.75 & 16.75 \\ \hline 2 & 3.5 \\ \hline \end{array}$$

d)

Use filter with uniform weights 0.25 to "average" across 4 cells

$$f = \begin{array}{|c|c|} \hline \frac{1}{4} & \frac{1}{4} \\ \hline \frac{1}{4} & \frac{1}{4} \\ \hline \end{array}$$

no padding, stride: 2

$$2 \times 2 \quad \text{output size} = \left(\frac{n+2p-f}{s} + 1\right) = \left(\frac{4+0-2}{2} + 1\right) = 2 \quad \checkmark$$