

Backpropagation

In this assignment, you will implement Backpropagation from scratch. You will then verify the correctness of the your implementation using a "grader" function/cell (provided by us) which will match your implementation.

The grader function would help you validate the correctness of your code.

Please submit the final Colab notebook in the classroom ONLY after you have verified your code using the grader function/cell.

Loading data

```
In [5]: import pickle
import numpy as np
from tqdm import tqdm
import matplotlib.pyplot as plt

with open('data.pkl', 'rb') as f:
    data = pickle.load(f)
print(data.shape)
x = data[:, :5]
y = data[:, -1]
print(x.shape, y.shape)

(500, 6)
(500, 5) (500,)
```

Check this video for better understanding of the computational graphs and back propagation

In [6]: From IPython.display import YouTubeVideo
YouTubeVideo("1840v96mo", width="1000", height="500")

Out [6]:

Computational graph

- If you observe the graph, we are having input features [f1, f2, f3, f4, f5] and 9 weights [w1, w2, w3, w4, w5, w6, w7, w8, w9].
- The final output of this graph is a value L, which is computed as $(Y-Y')^2$

Task 1: Implementing Forward propagation, Backpropagation and Gradient checking

Task 1.1

Forward propagation

- Forward propagation(Write your code in `def forward_propagation()`)

For easy debugging, we will break the computational graph into 3 parts.

Part 1<button>

Part 2<button>

Part 3<button>

```
In [2]: def sigmoid(z):
'''In this function, we will compute the sigmoid(z)'''
# we can use this function to forward and backward propagation
# write the code to compute the sigmoid value of z and return that value
return 1/(1+np.exp(-z))

In [3]: def grader_forward(data):
'''If you have written the code correctly then the grader function will output true'''
val=sigmoid(z)
assert(val==0.88079779778823)
return True

Out [3]: True

In [4]: def forward_propagation(x, y, w):
'''In this function, we will compute the forward propagation'''
# x: input data point, note that in this assignment you are having 5-d data points
# y: output variable
# w: weight array, its of length 9, w[0] corresponds to w1 in graph, w[1] corresponds to w2 in graph, ..., w[8] corresponds to w9 in graph.
# you have to return the following variables
# exp: part1 (compute the forward propagation until exp and then store the values in exp)
# tanh: part2 (compute the forward propagation until tanh and then store the values in tanh)
# sig = part3 (compute the forward propagation until sigmoid and then store the values in sig)
# we are computing one of the values for better understanding

val_1 = (w[0]*x[0]+w[1]*x[1]) + (w[8]*x[8])*(w[1]*x[1]) + w[5]
part_1 = np.exp(val_1)
part_2=tanh(part_1 + w[6])
val_3=(sin(w[2]*x[2]))*(w[3]*x[3])*(w[4]*x[4])) + w[7]
part_3=sigmoid(val_3)

y_pred = part_3*w[8] + part_2
loss=(y-y_pred)**2
dy_pred=-2*(y-y_pred)

# after computing part1, part2 and part3 compute the value of 'y' from the main Computational graph using required equations
# write code to compute the value of L (y-y')^2 and store it in variable loss
# compute derivative of L w.r to y' and store it in dy_pred
# Create a dictionary to store all the intermediate values i.e. dy_pred, loss, exp, tanh, sigmoid
# we will be using the dictionary to find values in backpropagation, you can add other keys in dictionary as well

forward_dict={}
forward_dict['exp']=part_1
forward_dict['tanh']=part_2
forward_dict['loss']=loss
forward_dict['dy_pred']=dy_pred

return forward_dict

In [5]: def grader_forwardprop(data):
d1 = (data['dy_pred']==-1.9285278284819143)
loss=(data['loss']==0.9280048982072913)
part1=(data['exp']==1.272987848973585)
part2=(data['tanh']==0.8437934192562146)
part3=(data['sigmoid']==0.8789793074319722)
assert(d1 and loss and part1 and part2 and part3)
w=np.ones(9)*0.1
d1=forward_propagation(X[0], y[0], w)
grader_forwardprop(d1)

Out [5]: True
```

Task 1.2

Backward propagation

```
In [6]: def backward_propagation(x,y,w,forward_dict):
'''In this function, we will compute the backward propagation'''
# Hint: you can use dict type to store the required variables
# dw = # in dw compute derivative of L w.r to w
# dw2 = # in dw2 compute derivative of L w.r to w2
# dw3 = # in dw3 compute derivative of L w.r to w3
# dw4 = # in dw4 compute derivative of L w.r to w4
# dw5 = # in dw5 compute derivative of L w.r to w5
# dw6 = # in dw6 compute derivative of L w.r to w6
# dw7 = # in dw7 compute derivative of L w.r to w7
# dw8 = # in dw8 compute derivative of L w.r to w8
# dw9 = # in dw9 compute derivative of L w.r to w9

dw1=forward_dict['dy_pred']*(1-(pow(forward_dict['tanh'],2)))*forward_dict['exp']**2*((w[0]*x[0])+(w[1]*x[1]))*x[9]
dw2=forward_dict['dy_pred']*(1-(pow(forward_dict['tanh'],2)))*forward_dict['exp']**2*((w[0]*x[0])+(w[1]*x[1]))*x[1]
dw3=forward_dict['dy_pred']*(forward_dict['sigmoid']*(1-forward_dict['sigmoid']))*(w[8]*(w[3]*x[3])+(w[1]*x[1]))*x[2]**w[2]**x[2]
dw4=forward_dict['dy_pred']*(forward_dict['sigmoid']*(1-forward_dict['sigmoid']))*(w[8]*(w[3]*x[3])+(w[1]*x[1]))*x[3]
dw5=forward_dict['dy_pred']*(1-(forward_dict['tanh']**2))*forward_dict['exp']
dw6=forward_dict['dy_pred']*(1-(forward_dict['sigmoid']*(1-forward_dict['sigmoid']))*(w[8]
dw7=forward_dict['dy_pred']*(1-(forward_dict['sigmoid']*(1-forward_dict['sigmoid']))*(w[8]
dw8=forward_dict['dy_pred']*(1-(forward_dict['sigmoid']*(1-forward_dict['sigmoid']))*(w[8]
dw9=forward_dict['dy_pred']*(1-(forward_dict['sigmoid']*(1-forward_dict['sigmoid']))*(w[8]

backward_dict={}
# store the variables dw1, dw2 etc. in a dict as backward_dict['dw1']= dw1, backward_dict['dw2']= dw2...
backward_dict['dw1']=dw1
backward_dict['dw2']=dw2
backward_dict['dw3']=dw3
backward_dict['dw4']=dw4
backward_dict['dw5']=dw5
backward_dict['dw6']=dw6
backward_dict['dw7']=dw7
backward_dict['dw8']=dw8
backward_dict['dw9']=dw9

return backward_dict

In [7]: def grader_backwardprop(data):
dw1=np.round(data['dw1'],6)==-0.2297933
dw2=np.round(data['dw2'],6)==-0.021408
dw3=np.round(data['dw3'],6)==-0.066625
dw4=np.round(data['dw4'],6)==-0.094588
dw5=np.round(data['dw5'],6)==-0.043888
dw6=np.round(data['dw6'],6)==0.633475
dw7=np.round(data['dw7'],6)==-0.56342
dw8=np.round(data['dw8'],6)==-0.048093
dw9=np.round(data['dw9'],6)==-1.018104
assert(dw1 and dw2 and dw3 and dw4 and dw5 and dw6 and dw7 and dw8 and dw9)
return True
w=np.ones(9)*0.1
d1=forward_propagation(X[0], y[0], w)
backward_dict=backward_propagation(X[0], y[0], w, forward_dict)
grader_backwardprop(backward_dict)

Out [7]: True
```

Task 1.3

Gradient clipping

Check this [blog link](#) for more details on Gradient clipping

we know that the derivative of any function is

$$\lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$$

- The definition above can be used as a numerical approximation of the derivative. Taking an epsilon small enough, the calculated approximation will have an error in the range of epsilon squared.
- In other words, if epsilon is 0.001, the approximation will be off by 0.00001.

Therefore, we can use this to approximate the gradient, and in turn make sure that backpropagation is implemented properly. This forms the basis of gradient checking!

Gradient checking example

lets understand the concept with a simple example: $f(w_1, w_2, x_1, x_2) = w_1^2 \cdot x_1 + w_2 \cdot x_2$

from the above function , lets assume $w_1 = 1, w_2 = 2, x_1 = 3, x_2 = 4$ the gradient of w_1 wrt w_1 is

$$\frac{df}{dw_1} = dw_1 = \frac{2 \cdot w_1 \cdot x_1}{2.1.3} = 2.1.3 = 6$$

let calculate the approximate gradient of w_1 as mentioned in the above formula and considering $\epsilon = 0.0001$

$$\begin{aligned} dw_1^{approx} &= \frac{f(w_1+\epsilon, w_2, x_1, x_2) - f(w_1-\epsilon, w_2, x_1, x_2)}{2\epsilon} \\ &= \frac{((1+0.0001)^2 \cdot 3 \cdot 4) - ((1-0.0001)^2 \cdot 3 \cdot 4)}{2 \cdot 0.0001} \\ &= \frac{(1.00020001 \cdot 3 \cdot 4) - (0.99980001 \cdot 3 \cdot 4)}{0.0002} \\ &= \frac{12.00600012 - 11.99400008}{0.0002} \\ &= 5.9999999999 \end{aligned}$$

Then, we apply the following formula for gradient check: $gradient_check = \frac{||dw_1 - dw_1^{approx}||_2}{||dw_1||_2 + ||dw_1^{approx}||_2}$

The equation above is basically the Euclidean distance normalized by the sum of the norm of the vectors. We use normalization in case that one of the vectors is very small. As a value for epsilon, we usually opt for $1e-7$. Therefore, if gradient check return a value less than $1e-7$, then it means that backpropagation was implemented correctly. Otherwise, there is potentially a mistake in your implementation. If the value exceeds $1e-3$, then you are sure that the code is not correct.

In our example: $gradient_check = \frac{(0.10000000000000009)}{(0.10000000000000009)} = 4.251414035630797e-13$

you can mathematically derive the same thing like this

$$\begin{aligned} dw_1^{approx} &= \frac{f(w_1+\epsilon, w_2, x_1, x_2) - f(w_1-\epsilon, w_2, x_1, x_2)}{2\epsilon} \\ &= \frac{(w_1+\epsilon)^2 \cdot x_1 + w_2 \cdot x_2 - (w_1-\epsilon)^2 \cdot x_1 + w_2 \cdot x_2}{2\epsilon} \\ &= \frac{4 \cdot w_1 \cdot x_1}{2\epsilon} \\ &= 2 \cdot w_1 \cdot x_1 \end{aligned}$$

Implement Gradient checking

(Write your code in `def gradient_checking()`)

```
Algorithm

W = initialize randomly
def gradient_checking(data_point, W):

# compute the L value using forward_propagation()
# compute the gradients of W using backward_propagation()</font>
approx_gradients = {}
for each w_i weight value in W:
    # add a small value to weight w_i, and then find the values of L with the updated weights
    # subtract a small value to weight w_i, and then find the values of L with the updated weights
    # compute the approximation gradients of weight w_i</font>
    # compute the approximation gradients of weight w_i</font>
    # compare the gradient of weights W from backward_propagation() with the approximation gradients of weights with <br> gradient check formulae</font>
    return gradient_check</font>

NOTE: you can do sanity check by checking all the return values of gradient checking(), they have to be zero. if not you have bug in your code

In [8]: def gradient_checking(x,y,w,eps):
# compute the dict value using forward_propagation()
# compute the actual gradients of W using backward_propagation()
forward_dict=forward_propagation(x,y,w)
backward_dict=backward_propagation(x,y,w,forward_dict)

#we are storing the original gradients for the given datapoints in a list
original_gradients_list=list(backward_dict.values())
# make sure that the order is correct i.e. first element in the list corresponds to dw1 , second element is dw2 etc.
# you can use reverse function if the values are in reverse order
w=np.add(w,eps)
w2=np.subtract(w,eps)
backward_dict_1=forward_propagation(x,y,w)
forward_dict_2=forward_propagation(x,y,w2)

approx_gradients_list=[forward_dict_1['loss']-forward_dict_2['loss']]/(2*eps)
#now we have to write code for approx gradients, here you have to make sure that you update only one weight at a time
#write your code here and append the approximate gradient value for each weight in approx_gradients_list

#performing gradient check operation
original_gradients_list=np.array(original_gradients_list)
approx_gradients_list=np.array(approx_gradients_list)
gradient_check_value=(original_gradients_list-approx_gradients_list)/(original_gradients_list+approx_gradients_list)
return gradient_check_value

In [9]: def grader_grad_check(value):
print(value)
assert(np.all(value == 10**(-3)))
return True

w1 = 0.00271756, 0.01260812, 0.00167639, 0.00207756, 0.00729768,
0.00140254, 0.00064168, 0.02422523, 0.01295444]
eps=1e-7
value= gradient_checking(X[0], y[0], w, eps)
grader_grad_check(value)

[-0.86437591 -0.99851478 -1.00000496 -0.99998429 -0.99999966 -0.23635657
-0.23690899 -0.996408 -0.1416402 ]
True

Out [9]: True
```

Task 2 : Optimizers

- As a part of this task, you will be implementing 2 optimizers(methods to update weight)
- Use the same computational graph that was mentioned above to do this task
- The weights have been initialized from normal distribution with mean=0 and std=0.01. The initialization of weights is very important otherwise you can face vanishing gradient and exploding gradients problem.

Check below video for reference purpose

```
In [10]: from IPython.display import YouTubeVideo
YouTubeVideo("gYpoJHlgyXA", width="1000", height="500")

Out [10]:
```

Algorithm

For each epoch(1-20):
For each data point in your data:
using the functions forward_propagation() and backward_propagation() compute the gradients of weights
update the weights with help of gradients

Implement below tasks

- Task 2.1: you will be implementing the above algorithm with Vanilla update of weights
- Task 2.2: you will be implementing the above algorithm with Momentum update of weights
- Task 2.3: you will be implementing the above algorithm with Adam update of weights

Note : If you get any assertion error while running grader functions, please print the variables in grader functions and check which variable is returning False. Rereck your logic for that variable .

2.1 Algorithm with Vanilla update of weights

```
In [28]: from sklearn.metrics import mean_squared_error
vanilla_epochs=[]
rate=0.01
w=np.zeros(9)
w=np.random.normal(mu,std,9)
loss_values_vanilla=[]
epochs=20
for i in range(epochs):
    vanilla_epochs.append(i)
    y_pred=[]
    for j in range(len(data)):
        forward=forward_propagation(X[j],y[j],w)
        y_pred.append(forward['y_pred'])
        backward=backward_propagation(X[j],y[j],w,forward)
        for k in range(len(w)):
            w[k]=w[k]+rate*backward['dw'+str(k+1)]
    loss=mean_squared_error(y,y_pred)
    loss_values_vanilla.append(loss)

Out [28]: <matplotlib.legend.Legend at 0x27d036ee6bb>
```

2.2 Algorithm with Momentum update of weights

Here Gamma refers to the momentum coefficient, eta is learning rate and v_1 is moving average of our gradients at timestep t

```
In [49]: rate=0.01
w=np.zeros(9)
w=np.random.normal(mu,std,9)
b=0.9

In [50]: momentum_epochs=[]
loss_values_momentum=[]
for i in range(epochs):
    momentum_epochs.append(i)
    y_pred=[]
    for point in range(len(data)):
        forward=forward_propagation(X[point],y[point],w)
        y_pred.append(forward['y_pred'])
        backward=backward_propagation(X[point],y[point],w,forward)
        for j in range(len(w)):
            v[j]=v[j]+rate*w[j]
            w[j]=w[j]+(1-b)*backward['dw'+str(j+1)]+b*v[j]
    loss=mean_squared_error(y,y_pred)
    loss_values_momentum.append(loss)

In [51]: # Graph plot for epoch vs loss
plt.grid()
plt.plot(momentum_epochs,loss_values_momentum,label="loss")
plt.title('epochs vs loss')
plt.xlabel('epochs')
plt.ylabel('loss')
plt.legend()

Out [51]: <matplotlib.legend.Legend at 0x27da336c18b>
```

2.3 Algorithm with Adam update of weights

```
In [ ]: m = beta1*m + (1-beta1)*dw
v = beta2*v + (1-beta2)*dw**2
x = - learning_rate * m / (np.sqrt(v) + eps)

In [62]: w=np.zeros(9)
w=np.zeros(9)
rate=0.001
beta1=0.9
beta2=0.99
w=np.random.normal(mu,std,9)

In [63]: adam_epochs=[]
loss_values_adam=[]
for i in range(epochs):
    adam_epochs.append(i)
    y_pred=[]
    for point in range(len(data)):
        forward=forward_propagation(X[point],y[point],w)
        y_pred.append(forward['y_pred'])
        backward=backward_propagation(X[point],y[point],w,forward)
        for j in range(len(w)):
            m[j]=beta1*m[j] + (1-beta1) * backward['dw'+str(j+1)]
            v[j]=beta2*v[j] + (1-beta2) * (backward['dw'+str(j+1)]**2)
            w[j]=w[j]-rate*m[j]/(np.sqrt(v[j])+eps)
    loss=mean_squared_error(y,y_pred)
    loss_values_adam.append(loss)

In [67]: # Graph plot for epoch vs loss
plt.grid()
plt.plot(adam_epochs,loss_values_adam,label="loss")
plt.title('epochs vs loss')
plt.xlabel('epochs')
plt.ylabel('loss')
plt.legend()

Out [67]: <matplotlib.legend.Legend at 0x27da336b0ab>
```

Comparison plot between epochs and loss with different optimizers. Make sure that loss is converging with increasing epochs

```
In [66]: # Grid the graph between loss vs epochs for all 3 optimizers.
plt.grid()
plt.plot(momentum_epochs,loss_values_momentum,label="momentum")
plt.plot(vanilla_epochs,loss_values_vanilla,label="vanilla")
plt.plot(adam_epochs,loss_values_adam,label="adam")
plt.title('epochs vs loss')
plt.xlabel('epochs')
plt.ylabel('loss')
plt.legend()

Out [66]: <matplotlib.legend.Legend at 0x27da336d09b>
```

You can go through the following blog to understand the implementation of other optimizers .
[Graders update blog](<https://cs231n.github.io/neural-networks-3/>)

```
In [ ]:
```