```
function/cell (provided by us) which will match your implmentation.
           The grader fucntion 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 [1]: import pickle
           import math
           import numpy as np
           from tgdm import tgdm
           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)
           (506, 6)
           (506, 5) (506,)
           Check this video for better understanding of the computational graphs and back propagation
  In [ ]: from IPython.display import YouTubeVideo
           YouTubeVideo('i940vYb6noo', width="1000", height="500")
  Out[ ]:
                   CS231n Winter 2016: Lecture 4: Backpropagation, Neural Networks 1
                                                                                          "local gradient"
                                                                                     \frac{\partial z}{\partial x}
                                                                                                                     |z|
                                                                                                                     \frac{\partial L}{\partial z}
                                                                                                                 gradients
                                                                                                                         13 Jan 2016
                                                                 Li & Andrej Karpathy & Justin Johnson
                                                                                                      Lecture 4 - 27
            Watch on  YouTube
           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</b>
              Part 2</b>
              Part 3</b>
In [89]: def sigmoid(z):
               ""In this function, we will compute the sigmoid(z)""
               # we can use this function in 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 [90]: def grader_sigmoid(z):
             #if you have written the code correctly then the grader function will output true
             val=sigmoid(z)
             assert(val==0.8807970779778823)
             return True
           grader_sigmoid(2)
 Out[90]: True
 In [91]: 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 varible
                   # 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
                   exp = np.exp((np.dot(w[0],x[0])+np.dot(w[1],x[1]))**2 + w[5])
                   tan_h = np.tanh(exp + w[6])
                   g = ((np.sin(np.dot(w[2],x[2]))) * (np.dot(w[3],x[3])+ np.dot(w[4],x[4])))
                   sig = sigmoid(g + w[7])
                   y_pred = tan_h + np.dot(sig , w[8])
                   L = (y - y_pred)**2
                   dL = 2*(y\_pred - y)
                   # after computing part1, part2 and part3 compute the value of y' from the main Computational graph using requ
           ired 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']= exp
                   forward_dict['sigmoid'] = sig
                   forward_dict['tanh'] = tan_h
                   forward_dict['loss'] = L
                   forward_dict['dy_pred'] = dL
                   return forward_dict
In [92]: def grader_forwardprop(data):
               dl = (data['dy_pred']==-1.9285278284819143)
               loss=(data['loss']==0.9298048963072919)
               part1=(data['exp']==1.1272967040973583)
               part2=(data['tanh']==0.8417934192562146)
               part3=(data['sigmoid']==0.5279179387419721)
               assert(dl and loss and part1 and part2 and part3)
               return True
           w=np.ones(9)*0.1
           d1=forward_propagation(X[0],y[0],w)
           grader_forwardprop(d1)
 Out[92]: True
           Task 1.2
           Backward propagation
 In [93]: fp= forward_propagation(X[0],y[0],w)
           fp
 Out[93]: {'exp': 1.1272967040973583,
            'sigmoid': 0.5279179387419721,
            'tanh': 0.8417934192562146,
            'loss': 0.9298048963072919,
            'dy_pred': -1.9285278284819143}
In [94]: def backward_propagation(x,y,w,forward_dict):
               '''In this function, we will compute the backward propagation '''
               # forward_dict: the outputs of the forward_propagation() function
               # write code to compute the gradients of each weight [w1,w2,w3,...,w9]
               # Hint: you can use dict type to store the required variables
               # dw1 = # in dw1 compute derivative of L w.r.to w1
               \# 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
               dw9 = fp['dy_pred'] * fp['sigmoid']
               dw7 = fp['dy\_pred'] * (1 - (fp['tanh'])**2)
               dw8 = fp['dy_pred'] * w[8] *fp['sigmoid']*(1-fp['sigmoid'])
               dw6 = dw7 * fp['exp']
               dw1 = dw6 * (2 * (np.dot(x[0],w[0])+np.dot(x[1],w[1])) * x[0])
               dw2 = dw6 * (2 * (np.dot(x[0],w[0])+np.dot(x[1],w[1])) * x[1])
               dw3 = dw8 * (np.dot(x[3],w[3]) + np.dot(x[4],w[4])) * (math.cos(np.dot(x[2],w[2]))) *x[2]
               dw4 = dw8 * (np.sin(np.dot(x[2],w[2]))) * x[3]
               dw5 = dw8 * (np.sin(np.dot(x[2],w[2]))) * x[4]
               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 [95]: def grader_backprop(data):
               dw1=(np.round(data['dw1'],6)==-0.229733)
               dw2=(np.round(data['dw2'],6)==-0.021408)
               dw3=(np.round(data['dw3'],6)==-0.005625)
               dw4=(np.round(data['dw4'],6)==-0.004658)
               dw5=(np.round(data['dw5'],6)==-0.001008)
               dw6=(np.round(data['dw6'],6)==-0.633475)
               dw7=(np.round(data['dw7'],6)==-0.561942)
               dw8=(np.round(data['dw8'],6)==-0.048063)
               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
           forward_dict=forward_propagation(X[0],y[0],w)
           backward_dict=backward_propagation(X[0],y[0],w,forward_dict)
           grader_backprop(backward_dict)
 Out[95]: True
 In [96]: | bp = backward_propagation(X[0], y[0], w, forward_dict)
 Out[96]: {'dw1': -0.22973323498702003,
            'dw2': -0.021407614717752925,
            'dw3': -0.005625405580266319,
            'dw4': -0.004657941222712423,
            'dw5': -0.0010077228498574246,
            'dw6': -0.6334751873437471,
            'dw7': -0.561941842854033,
            'dw8': -0.04806288407316516,
            'dw9': -1.0181044360187037}
           Task 1.3
           Gradient clipping
           Check this <u>blog link</u> for more details on Gradient clipping
           we know that the derivative of any function is
                                                             \lim_{\epsilon	o 0}rac{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</font>
           lets understand the concept with a simple example: f(w1,w2,x1,x2)=w_1^2.\,x_1+w_2.\,x_2
           from the above function , lets assume w_1=1 , w_2=2 , x_1=3 , x_2=4 the gradient of f w.r.t w_1 is
                                                             rac{df}{dw_1}=dw_1 \quad = \quad 2.w_1.\,x_1
                                                                       = 2.1.3
           let calculate the aproximate gradient of w_1 as mentinoned in the above formula and considering \epsilon=0.0001
                                                             = \frac{((1+0.0001)^2.3+2.4)-((1-0.0001)^2.3+2.4)}{2\epsilon}
                                                             = \frac{(1.00020001.3+2.4) - (0.99980001.3+2.4)}{2*0.0001}
                                                              Then, we apply the following formula for gradient check: gradient\_check = \frac{\|(dW - dW^{approx})\|_2}{\|(dW)\|_2 + \|(dW^{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: \textit{gradient\_check} = \frac{(6-5.99999999994898)}{(6+5.99999999994898)} = 4.2514140356330737e^{-13}
           you can mathamatically derive the same thing like this
                                                   egin{array}{lll} dw_1^{approx} & = & rac{f(w1+\epsilon,w2,x1,x2)-f(w1-\epsilon,w2,x1,x2)}{2\epsilon} \ & = & rac{((w_1+\epsilon)^2.x_1+w_2.x_2)-((w_1-\epsilon)^2.x_1+w_2.x_2)}{2\epsilon} \ & = & rac{4.\epsilon.w_1.x_1}{2\epsilon} \end{array}
                                                              = 2.w_1.x_1
           Implement Gradient checking
           (Write your code in def gradient_checking())
           Algorithm
              W = initilize_randomly
              def gradient_checking(data_point, W):
                  # compute the L value using forward_propagation()
                  # compute the gradients of W using backword_propagation()</font>
                  approx_gradients = []
                  for each wi weight value in W:<font color='grey'>
                      # add a small value to weight wi, and then find the values of L with the updated weights
                      # subtract a small value to weight wi, and then find the values of L with the updated weights
                      # compute the approximation gradients of weight wi</font>
                      approx_gradients.append(approximation gradients of weight wi)<font color='grey'>
                  # compare the gradient of weights W from backword_propagation() with the aproximation gradients of weig
                  hts with <br> gradient_check formula</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 [97]: def gradient_checking(x,y,w,eps):
               # compute the dict value using forward_propagation()
               # compute the actual gradients of W using backword_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
               # you can use reverse function if the values are in reverse order
               approx_gradients_list=[]
               #now we have to write code for approx gradients, here you have to make sure that you update only one weight at a
               #write your code here and append the approximate gradient value for each weight in approx_gradients_list
               for i in range(len(w)):
                   w_plus_eps = w[i] + eps
                   w_min_{eps} = w[i] - eps
                   w[i] = w_plus_eps
                   y_plus_eps = forward_propagation(x,y,w)['loss']
                   w[i] = w_min_{eps}
                   y_min_eps = forward_propagation(x,y,w)['loss']
                   app_grad = (y_plus_eps - y_min_eps) / (2*eps)
                   approx_gradients_list.append(app_grad)
               #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 [98]: def grader_grad_check(value):
               print(value)
               assert(np.all(value <= 10**-3))</pre>
               return True
           w=[0.00271756, 0.01260512, 0.00167639, -0.00207756, 0.00720768,
              0.00114524, 0.00684168, 0.02242521, 0.01296444]
           eps=10**-7
           value= gradient_checking(X[0],y[0],w,eps)
           grader_grad_check(value)
           [-0.17556562 -0.17555313 -0.06174431 -0.06150066 -0.06168482 -0.1755656
            -0.23240444 -0.06149022 -0.03845871]
 Out[98]: 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
              otherwiswe you can face vanishing gradient and exploding gradients problem.
           Check below video for reference purpose
  In [ ]: from IPython.display import YouTubeVideo
           YouTubeVideo('gYpoJMlgyXA', width="1000", height="500")
           Algorithm
                   for each epoch(1-20):
                       for each data point in your data:
                           using the functions forward_propagation() and backword_propagation() compute the gradients of weig
              hts
                           update the weigts with help of gradients
           Implement below tasks</b>

    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 .Recheck your logic for that variable .
           2.1 Algorithm with Vanilla update of weights
In [105]: mu, sigma = 0 , 0.01
           learning_rate = 0.001
           w1 = np.random.normal(mu, sigma, 9)
           loss=[]
           for i in range(0,20):
               11=[]
               for j in range(len(X)):
                   fp = forward_propagation(X[j],y[j],w1)
                   bp = backward_propagation(X[j],y[j],w1,fp)
                   grad = np.asarray(list(bp.values()))
                   w1 = w1 - (learning_rate*grad)
                   11.append(fp['loss'])
               loss.append(np.mean(11))
           print('Vanilla update of weights:\n',loss)
           Vanilla update of weights:
            [0.6879032566291252,\ 0.3376373168892776,\ 0.17905353093593102,\ 0.09814754386653339,\ 0.05919911730920271,\ 0.0419305265
           5699068, 0.03484713079323609, 0.032149666568279535, 0.03120717764387713, 0.030919202908387903, 0.03085424307921121,
           0.030853234557867086, 0.03086200941485578, 0.0308645625092576, 0.03085771434190855, 0.03084196797197829, 0.0308184444
           2642217, 0.030787994622622335, 0.030751035903271235, 0.030707596141364434]
In [106]: plt.figure(figsize=(10,10))
           epoch = np.arange(0,20)
           plt.plot(epoch, loss, label='Vanilla')
           plt.scatter(epoch,loss,label = 'Vanilla points')
           plt.legend()
           plt.xlabel('Epoch')
           plt.ylabel('Loss')
           plt.title('Vanilla-Loss Plot')
           plt.grid()
           plt.show()
                                              Vanilla-Loss Plot
                                                                           --- Vanilla
             0.7

    Vanilla points

             0.6
             0.5
             0.3
             0.2
             0.1
             0.0
                                                             12.5
                                                                             17.5
                   0.0
                           2.5
                                    5.0
                                            7.5
                                                    10.0
                                                                     15.0
                                                  Epoch
           2.2 Algorithm with Momentum update of weights
           Here Gamma referes to the momentum coefficient, eta is leaning rate and v_t is moving average of our gradients at timestep t
In [121]: mu, sigma = 0 , 0.01
           learning_rate = 0.001
           w1 = np.random.normal(mu, sigma, 9)
           m_loss=[]
           gamma = 0.9
           v_t=0
           for i in range(0,20):
               11=[]
               for j in range(len(X)):
                   fp = forward_propagation(X[j], y[j], w1)
                   bp = backward_propagation(X[j],y[j],w1,fp)
                   grad = np.asarray(list(bp.values()))
                   v_t = (gamma*v_t) + (1- gamma) *grad
                   w1 = w1 - (learning_rate*v_t)
                   11.append(fp['loss'])
               m_loss.append(np.mean(l1))
           print('Momentum update of weights:\n',m_loss)
           Momentum update of weights:
            [0.7079929625815475, 0.3371866243781629, 0.17591009273370498, 0.09550673109225633, 0.05766246566146238, 0.0412718768
           5188669, 0.034716503650040344, 0.03229408777820288, 0.03148395432767918, 0.031258064359498676, 0.031223807550487464,
           0.031240892148264428, 0.03126365024689717, 0.03127978712891493, 0.031287738155729584, 0.031288758714026976, 0.0312843
           98247743476, 0.031275839416183546, 0.03126382783996774, 0.0312487537651831]
In [122]: plt.figure(figsize=(10,10))
           epoch = np.arange(0,20)
           plt.plot(epoch, m_loss, label='Momentum')
           plt.scatter(epoch,m_loss,label = 'Momentum points')
           plt.legend()
           plt.xlabel('Epoch')
           plt.ylabel('Loss')
           plt.title('Momentum-Loss Plot')
           plt.grid()
           plt.show()
                                            Momentum-Loss Plot
                                                                        — Momentum

    Momentum points

             0.7
             0.5
             0.3
             0.2
             0.1
                           2.5
                                    5.0
                                            7.5
                                                    10.0
                                                             12.5
                                                                             17.5
                   0.0
                                                  Epoch
           2.3 Algorithm with Adam update of weights
In [118]: mu, sigma = 0 , 0.01
           learning_rate = 0.001
           w1 = np.random.normal(mu, sigma, 9)
           a_loss=[]
           beta1 = 0.9
           beta2=0.99
           mt=0
           vt=0
           mt1=0
           vt1=0
           eps=1e-8
           for i in range(0,20):
               11=[]
               for j in range(len(X)):
                   fp = forward_propagation(X[j],y[j],w1)
                   bp = backward_propagation(X[j],y[j],w1,fp)
                   grad = np.asarray(list(bp.values()))
                   mt = (beta1*mt) + (1-beta1)* grad
                   vt = (beta2 * vt) + (1 - beta2)* ((grad)**2)
                   mt1 = mt/(1-beta1)
                   vt1 = vt/(1-beta2)
                   w1 = w1 - (learning_rate/np.sqrt(vt1+eps))*mt1
                   11.append(fp['loss'])
               a_loss.append(np.mean(l1))
           print('Adam update of weights:\n',a_loss)
           Adam update of weights:
            [0.4692627165061932, 0.07568400257146835, 0.005112652936817799, 0.0016766618250525563, 0.0007510070677249286, 0.0003
           224273765460544, 0.00012577925851157731, 5.371464941987592e-05, 2.4051875768510914e-05, 8.583552512094133e-06, 2.2078
           599323298643e-06, 8.124108569051007e-07, 9.747133697596357e-07, 1.1320714910822964e-06, 1.0310456539732363e-06, 1.093
           5732306982432e-06, 9.655087812129749e-07, 1.0599517480542574e-06, 1.0074276094307395e-06, 1.0393072682933578e-06]
In [119]: plt.figure(figsize=(10,10))
           epoch = np.arange(0,20)
           plt.plot(epoch, a_loss, label='Adam')
           plt.scatter(epoch, a_loss, label = 'Adam points')
           plt.legend()
           plt.xlabel('Epoch')
           plt.ylabel('Loss')
           plt.title('Adam-Loss Plot')
           plt.grid()
           plt.show()
                                               Adam-Loss Plot
                                                                           — Adam

    Adam points

             0.4
             0.3
             0.2
             0.1
                   0.0
                           2.5
                                            7.5
                                                    10.0
                                                             12.5
                                                                     15.0
                                                                             17.5
                                                                                      20.0
                                                  Epoch
           Comparision plot between epochs and loss with different optimizers. Make sure that loss is conerging with increaing epochs
In [123]: #plot the graph between loss vs epochs for all 3 optimizers.
           plt.figure(figsize=(10,10))
           epoch = np.arange(0,20)
           plt.plot(epoch, a_loss, label='Adam')
           plt.plot(epoch, m_loss, label='Momentum')
           plt.plot(epoch, loss, label='Vanilla')
           plt.legend()
           plt.xlabel('Epoch')
           plt.ylabel('Loss')
           plt.title('Vanilla-Momentum-Adam-Loss Plot')
           plt.grid()
           plt.show()
                                       Vanilla-Momentum-Adam-Loss Plot
                                                                            --- Adam
                                                                               Momentum
             0.7
                                                                            — Vanilla
             0.5
```

0.4

0.3

0.2

0.1

0.0

10.0

Epoch

You can go through the following blog to understand the implementation of other optimizers .

[Gradients update blog](https://cs231n.github.io/neural-networks-3/)

12.5

17.5

Backpropagation

In this assignment, you will implement Backpropagation from scratch. You will then verify the correctness of the your implementation using a "grader"