# BackPropagation

There will be some functions that start with the word "grader" ex: grader_sigmoid(), grader_forwa
not change those function definition.

Every Grader function has to return True.

# Loading data

```
1 from google.colab import files
2 uploaded = files.upload()
```
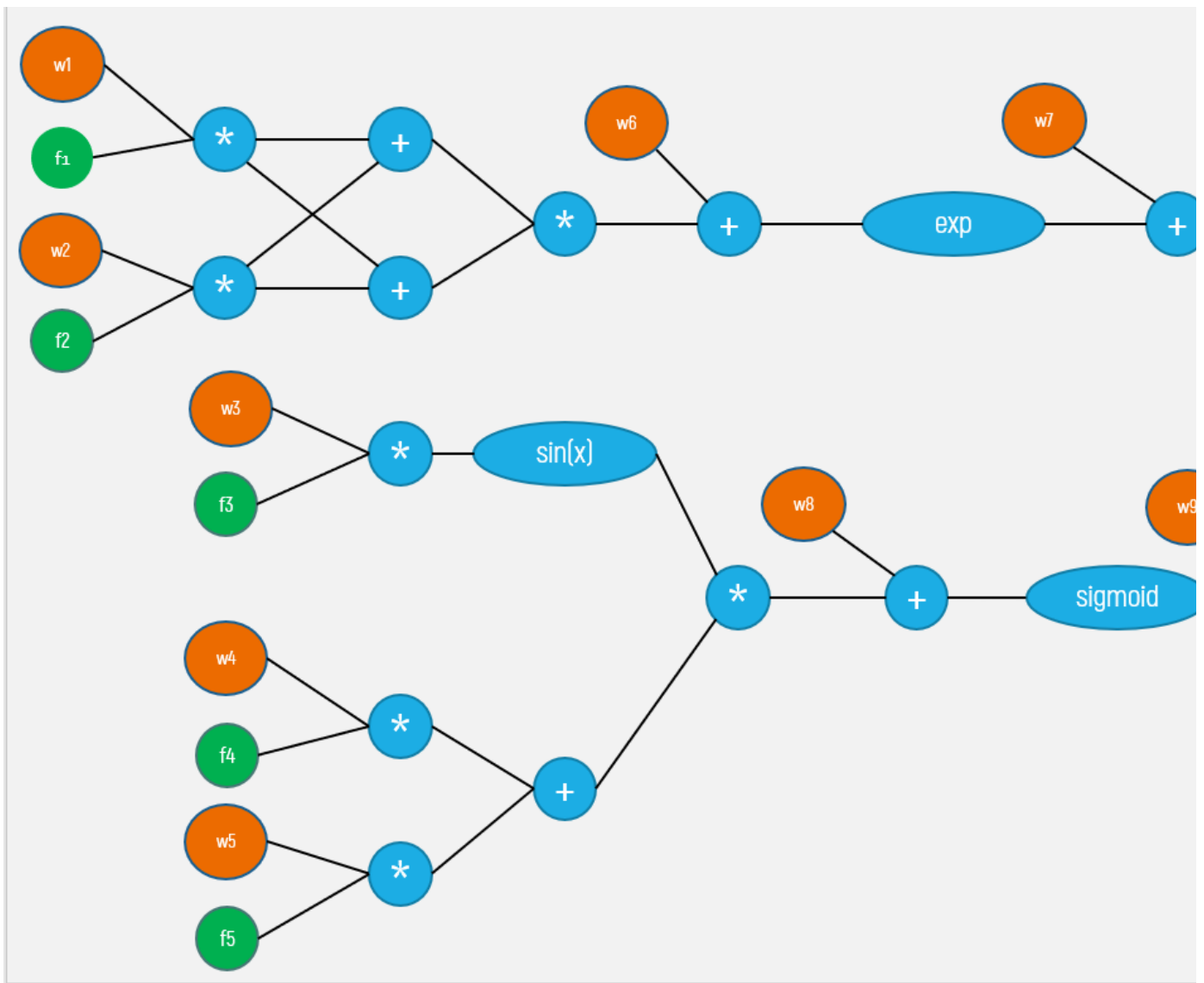
Choose Files  data.pkl
- **data.pkl**(n/a) - 24449 bytes, last modified: 5/24/2020 - 100% done
  Saving data.pkl to data (1).pkl

```
 1 import pickle
 2 import numpy as np
 3 import math
 4 from tqdm import tqdm
 5 import matplotlib.pyplot as plt
 6
 7 with open('data.pkl', 'rb') as f:
 8     data = pickle.load(f)
 9 print(data.shape)
10 X = data[:, :5]
11 y = data[:, -1]
12 print(X.shape, y.shape)
```

```
(506, 6)
(506, 5) (506,)
```

# Computational graph

- **If you observe the graph, we are having input features [f1, f2, f3, f4, f5] and 9 weights [w1, w**

- **The final output of this graph is a value L which is computed as (Y-Y')^2**

## Task 1: Implementing backpropagation and Gradient checkir

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

```
1 from IPython.display import YouTubeVideo
2 YouTubeVideo('i94OvYb6noo',width="1000",height="500")
```

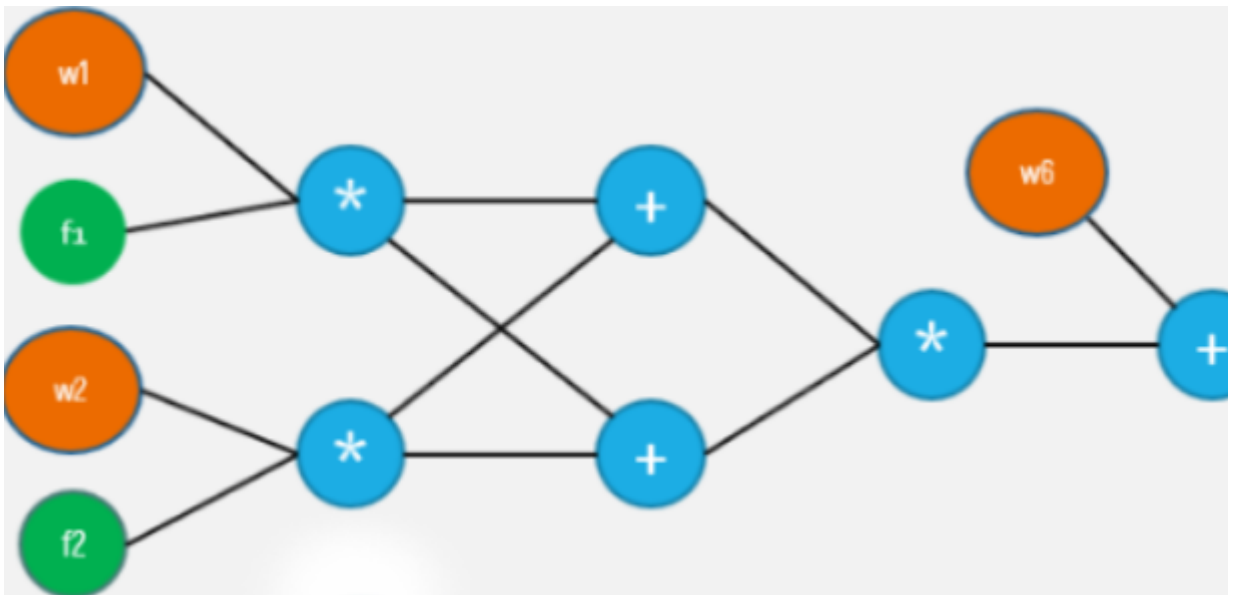## CS231n Winter 2016: Lecture 4: Backpropagation, Neural Networks 1
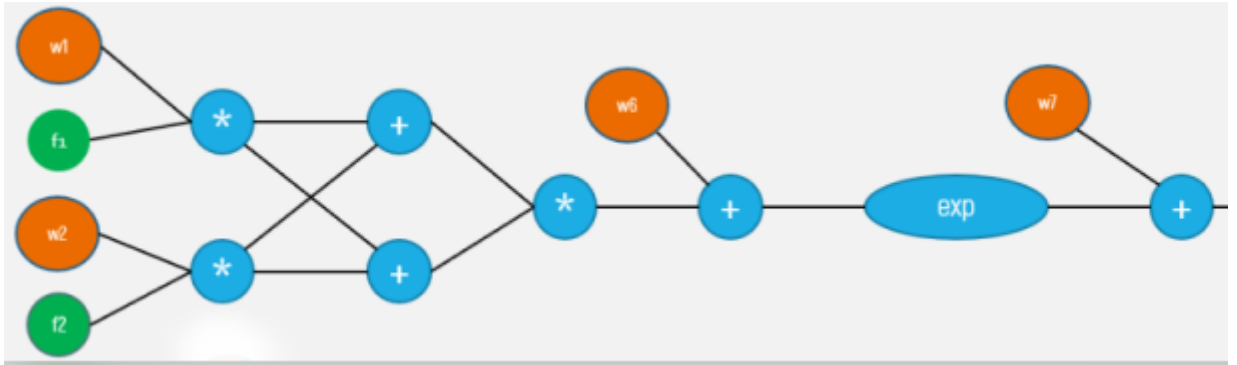
▶

- **Write two functions**

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

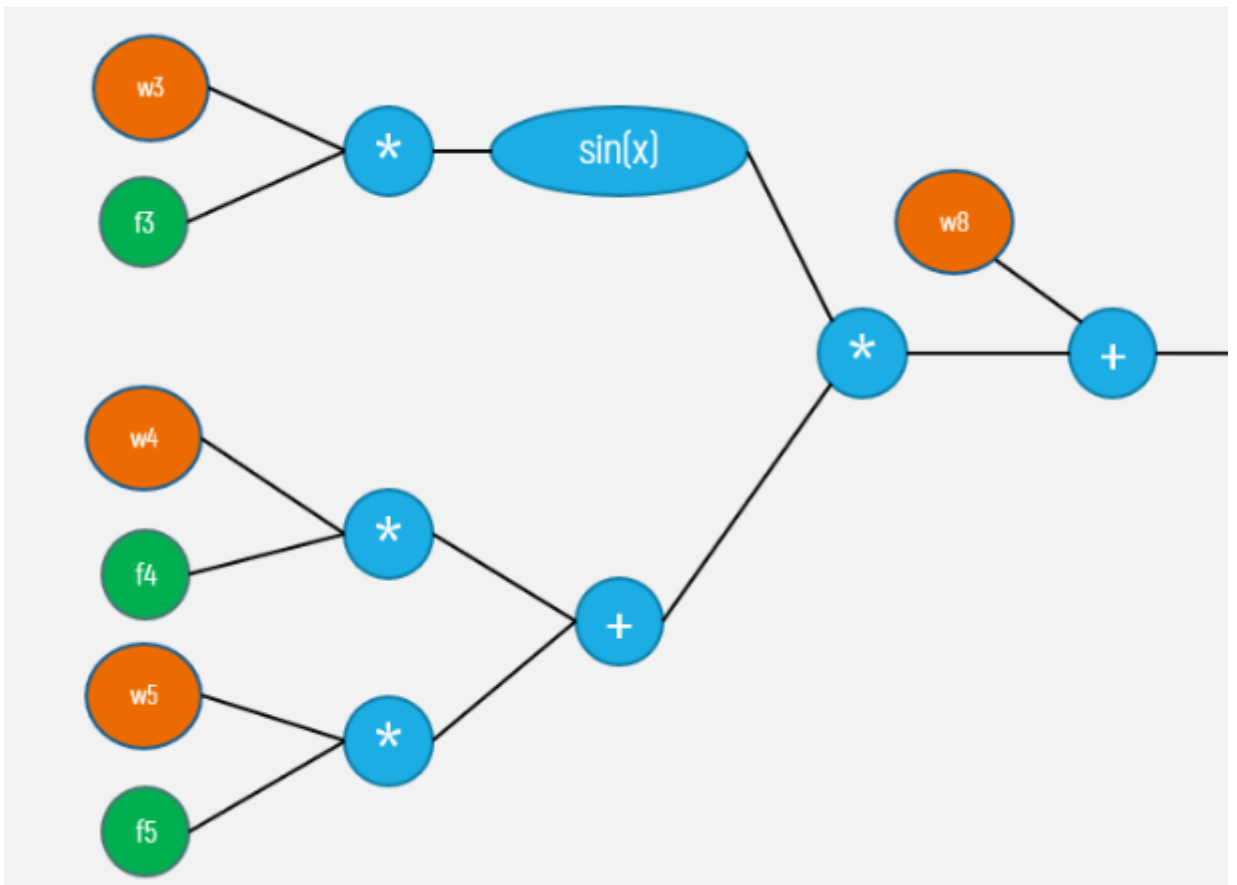    For easy debugging, we will break the computational graph into 3 parts.

    **Part 1**

**Part 2**



**Part 3**



```
def forward_propagation(X, y, W):

  # X: input data point, note that in this assignment you are having 5
  # y: output varible
  # W: weight array, its of length 9, W[0] corresponds to w1 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 sto
  # tanh =part2(compute the forward propagation until tanh and then st
  # sig = part3(compute the forward propagation until sigmoid and ther
  # now compute remaining values from computional graph and get y'
```

```
        # write code to compute the value of L=(y-y')^2
        # compute derivative of L  w.r.to Y' and store it in dl
        # Create a dictionary to store all the intermediate values
        # store L, exp,tanh,sig,dl variables

        return (dictionary, which you might need to use for back propagatior
```

- **Backward propagation**(Write your code in def backward_propagation())

```
    def backward_propagation(L, W,dictionary):

      # L: the loss we calculated for the current point
      # dictionary: the outputs of the forward_propagation() function
      # write code to compute the gradients of each weight [w1,w2,w3,...,w
      # Hint: you can use dict type to store the required variables
      # return dW, dW is a dictionary with gradients of all the weights

      return dW
```

# Gradient clipping

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

we know that the derivative of any function is

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

- The definition above can be used as a numerical approximation of the derivative. Taking an $\epsilon$ 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 backpropagatic basis of **gradient checking!**

# Gradient checking example

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 $w$

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

let calculate the aproximate gradient of $w_1$ as mentinoned in the above formula and considering $\epsilon$

$$dw_1^{approx} \quad = \quad \frac{f(w1+\epsilon,w2,x1,x2)-f(w1-\epsilon,w2,x1,x2)}{2\epsilon}$$
$$= \quad \frac{((1+0.0001)^2.3+2.4)-((1-0.0001)^2.3+2.4}{2\epsilon}$$
$$= \quad \frac{(1.00020001.3+2.4)-(0.99980001.3+2.4)}{2*0.0001}$$
$$= \quad \frac{(11.00060003)-(10.99940003)}{0.0002}$$
$$= \quad 5.99999999999$$

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 v
one of the vectors is very small. As a value for epsilon, we usually opt for 1e-7. Therefore, if gradie
means that backpropagation was implemented correctly. Otherwise, there is potentially a mistake
1e-3, then you are sure that the code is not correct.

in our example: $gradient\_check = \frac{(6-5.999999999994898)}{(6+5.999999999994898)} = 4.2514140356330737e^{-13}$

you can mathamatically derive the same thing like this

$$dw_1^{approx} \quad = \quad \frac{f(w1+\epsilon,w2,x1,x2)-f(w1-\epsilon,w2,x1,x2)}{2\epsilon}$$
$$= \quad \frac{((w_1+\epsilon)^2.x_1+w_2.x_2)-((w_1-\epsilon)^2.x_1+w_2.x_2}{2\epsilon}$$
$$= \quad \frac{4.\epsilon.w_1.x_1}{2\epsilon}$$
$$= \quad 2.w_1.x_1$$

# Implement Gradient checking

(Write your code in def gradient_checking())

**Algorithm**

```python
W = initilize_randomly
def gradient_checking(data_point, W):
    # compute the L value using forward_propagation()
    # compute the gradients of W using backword_propagation()
    approx_gradients = []
    for each wi weight value in W:
        # add a small value to weight wi, and then find the values of L with the updated
        # subtract a small value to weight wi, and then find the values of L with the upd
        # compute the approximation gradients of weight wi
        approx_gradients.append(approximation gradients of weight wi)
    # compare the gradient of weights W from backword_propagation() with the aproximation
  gradient_check formula
    return gradient_check
```

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

# Task 2 : Optimizers

- As a part of this task, you will be implementing 3 type of optimizers(methods to update weig
- Use the same computational graph that was mentioned above to do this task
- Initilze the 9 weights from normal distribution with mean=0 and std=0.01

**Check below video and [this](#) blog**

```python
1 from IPython.display import YouTubeVideo
2 YouTubeVideo('gYpoJMlgyXA',width="1000",height="500")
```

---

## CS231n Winter 2016: Lecture 5: Neural Networks Part 2

**Algorithm**

```
for each epoch(1-100):
    for each data point in your data:
        using the functions forward_propagation() and backword_propagation() compute
        update the weigts with help of gradients  ex: w1 = w1-learning_rate*dw1
```

# 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 g returning False .Recheck your logic for that variable .**

# Task 1

# Forward propagation

```
1 def sigmoid(z):
2     '''In this function, we will compute the sigmoid(z)'''
3     value=1/(1+np.exp(-z))
4     return value
5
6 def forward_propagation(x, y, w):
7   dict={}
8   exp=np.exp(((((w[0]*x[0])+(w[1]*x[1]))*((w[0]*x[0])+(w[1]*x[1])))+w[5])
9   tanh=np.tanh(exp+w[6])
```

```
10    part3=((np.sin(w[2]*x[2]))*((w[3]*x[3])+(w[4]*x[4])))+w[7]
11    sig=sigmoid(part3)
12    y_hat=tanh+sig*w[8]
13    dl=-2*(y-y_hat)
14    loss=pow(y-(tanh+sig*w[8]), 2)
15    dict["exp"]=exp
16    dict['tanh']=tanh
17    dict["sigmoid"]=sig
18    dict["dl"]=dl
19    dict["loss"]=loss
20    dict["y_hat"]=y_hat
21    return dict
22
```

## Grader function - 1

```
1 def grader_sigmoid(z):
2   val=sigmoid(z)
3   assert(val==0.8807970779778823)
4   return True
5 grader_sigmoid(2)
```

> True

## Grader function - 2

```
1 def grader_forwardprop(data):
2     dl = (np.round(data['dl'],4)==-1.9285)
3     loss=(np.round(data['loss'],4)==0.9298)
4     part1=(np.round(data['exp'],4)==1.1273)
5     part2=(np.round(data['tanh'],4)==0.8418)
6     part3=(np.round(data['sigmoid'],4)==0.5279)
7     assert(dl and loss and part1 and part2 and part3)
8     return True
9 w=np.ones(9)*0.1
10 d1=forward_propagation(X[0],y[0],w)
11 grader_forwardprop(d1)
```

> True

## Backward propagation

```
1 def backward_propagation(x,w,dict):
2   der={}
3   dw1=dict['dl']*(1-(math.pow(dict['tanh'],2)))*dict["exp"]*2*((w[0]*x[0])+(w[1]*x[1]))
4   dw2=dict['dl']*(1-(math.pow(dict['tanh'],2)))*dict["exp"]*2*((w[0]*x[0])+(w[1]*x[1]))
5   dw3=dict['dl']*(dict['sigmoid']*(1-dict['sigmoid']))*w[8]*((w[3]*x[3])+(w[4]*x[4]))*m
6   dw4=dict['dl']*(dict['sigmoid']*(1-dict['sigmoid']))*w[8]*math.sin(x[2]*w[2])*x[3]
7   dw5=dict['dl']*(dict['sigmoid']*(1-dict['sigmoid']))*w[8]*math.sin(x[2]*w[2])*x[4]
8   dw6 =dict['dl']*(1-(math.pow(dict['tanh'],2)))*dict["exp"]
9   dw7 = dict['dl']*(1-(math.pow(dict['tanh'],2)))
```

```
10   dw8 = dict['dl']*(dict['sigmoid']*(1-dict['sigmoid']))*w[8]
11   dw9 = dict['dl']*dict['sigmoid']
12   der['dw1']=dw1
13   der['dw2']=dw2
14   der['dw3']=dw3
15   der['dw4']=dw4
16   der['dw5']=dw5
17   der['dw6']=dw6
18   der['dw7']=dw7
19   der['dw8']=dw8
20   der['dw9']=dw9
21   return der
```

## Grader function - 3

```
 1 def grader_backprop(data):
 2     dw1=(np.round(data['dw1'],4)==-0.2297)
 3     dw2=(np.round(data['dw2'],4)==-0.0214)
 4     dw3=(np.round(data['dw3'],4)==-0.0056)
 5     dw4=(np.round(data['dw4'],4)==-0.0047)
 6     dw5=(np.round(data['dw5'],4)==-0.001)
 7     dw6=(np.round(data['dw6'],4)==-0.6335)
 8     dw7=(np.round(data['dw7'],4)==-0.5619)
 9     dw8=(np.round(data['dw8'],4)==-0.0481)
10     dw9=(np.round(data['dw9'],4)==-1.0181)
11     assert(dw1 and dw2 and dw3 and dw4 and dw5 and dw6 and dw7 and dw8 and dw9)
12     return True
13 w=np.ones(9)*0.1
14 d1=forward_propagation(X[0],y[0],w)
15 d1=backward_propagation(X[0],w,d1)
16 grader_backprop(d1)
```

⤷     True

## ▾ Implement gradient checking

```
 1 w=np.ones(9)*0.1
```

```
 1 def gradient_checking(data_point,y, w,epsilon=1e-7):
 2   weight=w
 3   loss=forward_propagation(data_point,y,w)
 4   grad=backward_propagation(data_point,w,loss)
 5   approx_gradients = []
 6   grad_check=[]
 7   for i in range(len(w)):
 8     w=np.ones(9)*0.1
 9     w[i]=weight[i]+epsilon
10     l1=forward_propagation(data_point,y,w)
11     loss1=l1['loss']
12     w[i]=weight[i]-epsilon
13     l2=forward_propagation(data_point,y,w)
```

```
14      loss2=l2["loss"]
15      aprox=(loss1-loss2)/(2*epsilon)
16      approx_gradients.append(aprox)
17      numerator = np.linalg.norm(grad[(str("dw")+str(i+1))] - aprox)
18      denominator = np.linalg.norm(grad[(str("dw")+str(i+1))]) + np.linalg.norm(aprox)
19      difference = numerator / denominator
20      grad_check.append(difference)
21   return grad_check,approx_gradients
```

```
1 print("ACTUAL GRADIENT USING BACK PROPAGATION")
2 loss=forward_propagation(X[0],y[0],w)
3 backward_propagation(X[0],w,loss)
```

```
ACTUAL GRADIENT USING BACK PROPAGATION
{'dw1': -0.22973323498702,
 'dw2': -0.02140761471775293,
 'dw3': -0.00562540558026632,
 'dw4': -0.004657941222712424,
 'dw5': -0.0010077228498574248,
 'dw6': -0.6334751873437471,
 'dw7': -0.561941842854033,
 'dw8': -0.04806288407316517,
 'dw9': -1.0181044360187037}
```

```
1 diff,aprox_grad=gradient_checking(X[0],y[0], w,epsilon=1e-7)
2 print("APPROXIMATE GRADIENT AFTER GRADIENT CHECKING")
3 print(aprox_grad)
4 print("DIFFERENCE")
5 print(diff)
```

```
APPROXIMATE GRADIENT AFTER GRADIENT CHECKING
[-0.22973323521302547, -0.021407614569923794, -0.005625406251930087, -0.0046579412549
DIFFERENCE
[4.918867403857781e-10, 3.4527231737366113e-09, 5.96991379072821e-08, 3.4579582758165
```

- we can understand back propagation gradient calculation is working well because difference

# Task 2: Optimizers

## Algorithm with Vanilla update of weights

```
1 from sklearn.metrics import mean_squared_error
2 rate=.001
3 mu, sigma = 0, 0.01 # mean and standard deviation
4 w = np.random.normal(mu, sigma, 9) # weight intialization
5 weight=w
```

```
1 loss_value_vanilla=[]
2 epoc_vanilla=[]
```

```
 2  epoc_vanilla []
 3  for epoch in range(100):
 4    epoc_vanilla.append(epoch)
 5    y_pred=[]
 6    for point in range(len(data)):
 7      forward=forward_propagation(X[point], y[point], w)
 8      y_pred.append(forward['y_hat'])
 9      backward=backward_propagation(X[point],w,forward)
10      w[0]=w[0]-rate*backward["dw1"]
11      w[1]=w[1]-rate*backward["dw2"]
12      w[2]=w[2]-rate*backward["dw3"]
13      w[3]=w[3]-rate*backward["dw4"]
14      w[4]=w[4]-rate*backward["dw5"]
15      w[5]=w[5]-rate*backward["dw6"]
16      w[6]=w[6]-rate*backward["dw7"]
17      w[7]=w[7]-rate*backward["dw8"]
18      w[8]=w[8]-rate*backward["dw9"]
19    loss=mean_squared_error(y,y_pred)
20    loss_value_vanilla.append(loss)
```
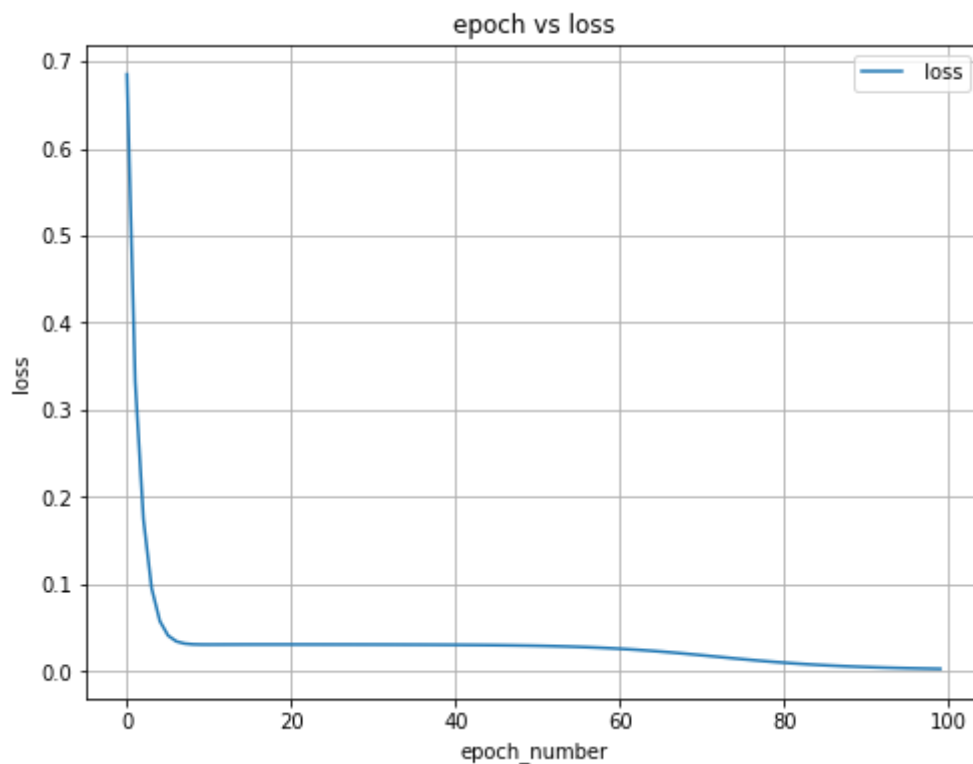
## Plot between epochs and loss

```
1  %matplotlib inline
2  import matplotlib.pyplot as plt
3  plt.figure(figsize=(8,6))
4  plt.grid()
5  plt.plot(epoc_vanilla,loss_value_vanilla, label=' loss')
6  plt.title("epoch vs loss")
7  plt.xlabel("epoch_number")
8  plt.ylabel("loss")
9  plt.legend()
```

```
<matplotlib.legend.Legend at 0x7ff9cd7a5be0>
```

## Algorithm with momentum update of weights

```
1 rate=.001
2 m=np.zeros(9)
3 b=.9
4 mu, sigma = 0, 0.01 # mean and standard deviation
5 w = np.random.normal(mu, sigma, 9) # weight intialization
```

```
1 loss_value_momentum=[]
2 epoc_momentum=[]
3 for epoch in range(100):
4   epoc_momentum.append(epoch)
5   y_pred=[]
6   for point in range(len(data)):
7     forward=forward_propagation(X[point], y[point], w)
8     y_pred.append(forward['y_hat'])
9     backward=backward_propagation(X[point],w,forward)
10    m[0]=b*m[0]+(1-b)*backward["dw1"]
11    w[0]=w[0]-rate*m[0]
12    m[1]=b*m[1]+(1-b)*backward["dw2"]
13    w[1]=w[1]-rate*m[1]
14    m[2]=b*m[2]+(1-b)*backward["dw3"]
15    w[2]=w[2]-rate*m[2]
16    m[3]=b*m[3]+(1-b)*backward["dw4"]
17    w[3]=w[3]-rate*m[3]
18    m[4]=b*m[4]+(1-b)*backward["dw5"]
19    w[4]=w[4]-rate*m[4]
20    m[5]=b*m[5]+(1-b)*backward["dw6"]
21    w[5]=w[5]-rate*m[5]
22    m[6]=b*m[6]+(1-b)*backward["dw7"]
23    w[6]=w[6]-rate*m[6]
24    m[7]=b*m[7]+(1-b)*backward["dw8"]
25    w[7]=w[7]-rate*m[7]
26    m[8]=b*m[8]+(1-b)*backward["dw9"]
27    w[8]=w[8]-rate*m[8]
28   loss=mean_squared_error(y,y_pred)
29   loss_value_momentum.append(loss)
30
```
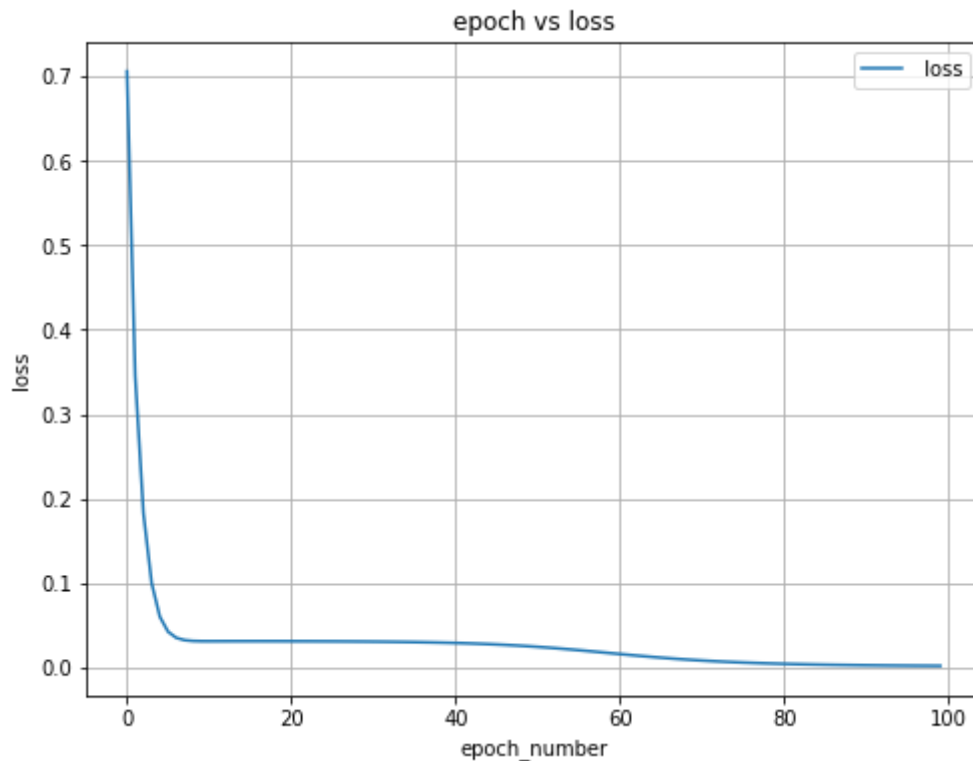
## Plot between epochs and loss

```
1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 plt.figure(figsize=(8,6))
4 plt.grid()
5 plt.plot(epoc_momentum,loss_value_momentum, label=' loss')
6 plt.title("epoch vs loss")
7 plt.xlabel("epoch_number")
8 plt.ylabel("loss")
9 plt.legend()
```

```
<matplotlib.legend.Legend at 0x7ff9cd1db630>
```



epoch vs loss

## Algorithm with Adam update of weights

```
1 rate=.001
2 m=np.zeros(9)
3 v=np.zeros(9)
4 b1=.9
5 b2=0.999
6 z=1e-8
7 mu, sigma = 0, 0.01 # mean and standard deviation
8 w = np.random.normal(mu, sigma, 9) # weight intialization
```

```
1 loss_value_adam=[]
2 epoc_adam=[]
3 for epoch in range(100):
4   epoc_adam.append(epoch)
5   y_pred=[]
6   for point in range(len(data)):
7     forward=forward_propagation(X[point], y[point], w)
8     y_pred.append(forward['y_hat'])
9     backward=backward_propagation(X[point],w,forward)
10    m[0]=b1*m[0]+(1-b1)*backward["dw1"]
11    mt=m[0]/(1-b1)
12    v[0]=b2*v[0]+(1-b2)*math.pow(backward["dw1"],2)
13    vt=v[0]/(1-b2)
14    w[0]=w[0]-(rate/(math.sqrt(vt)+z))*mt
15    m[1]=b1*m[1]+(1-b1)*backward["dw2"]
16    mt=m[1]/(1-b1)
17    v[1]=b2*v[1]+(1-b2)*math.pow(backward["dw2"],2)
18    vt=v[1]/(1-b2)
19    w[1]=w[1]-(rate/(math.sqrt(vt)+z))*mt
20    m[2]=b1*m[2]+(1-b1)*backward["dw3"]
```

```
20      m[2]=b1*m[2]+(1-b1)*backward["dw3"]
21      mt=m[2]/(1-b1)
22      v[2]=b2*v[2]+(1-b2)*math.pow(backward["dw3"],2)
23      vt=v[2]/(1-b2)
24      w[2]=w[2]-(rate/(math.sqrt(vt)+z))*mt
25      m[3]=b1*m[3]+(1-b1)*backward["dw4"]
26      mt=m[3]/(1-b1)
27      v[3]=b2*v[3]+(1-b2)*math.pow(backward["dw4"],2)
28      vt=v[3]/(1-b2)
29      w[3]=w[3]-(rate/(math.sqrt(vt)+z))*mt
30      m[4]=b1*m[4]+(1-b1)*backward["dw5"]
31      mt=m[4]/(1-b1)
32      v[4]=b2*v[4]+(1-b2)*math.pow(backward["dw5"],2)
33      vt=v[4]/(1-b2)
34      w[4]=w[4]-(rate/(math.sqrt(vt)+z))*mt
35      m[5]=b1*m[5]+(1-b1)*backward["dw6"]
36      mt=m[5]/(1-b1)
37      v[5]=b2*v[5]+(1-b2)*math.pow(backward["dw6"],2)
38      vt=v[5]/(1-b2)
39      w[5]=w[5]-(rate/(math.sqrt(vt)+z))*mt
40      m[6]=b1*m[6]+(1-b1)*backward["dw7"]
41      mt=m[6]/(1-b1)
42      v[6]=b2*v[6]+(1-b2)*math.pow(backward["dw7"],2)
43      vt=v[6]/(1-b2)
44      w[6]=w[6]-(rate/(math.sqrt(vt)+z))*mt
45      m[7]=b1*m[7]+(1-b1)*backward["dw8"]
46      mt=m[7]/(1-b1)
47      v[7]=b2*v[7]+(1-b2)*math.pow(backward["dw8"],2)
48      vt=v[7]/(1-b2)
49      w[7]=w[7]-(rate/(math.sqrt(vt)+z))*mt
50      m[8]=b1*m[8]+(1-b1)*backward["dw9"]
51      mt=m[8]/(1-b1)
52      v[8]=b2*v[8]+(1-b2)*math.pow(backward["dw9"],2)
53      vt=v[8]/(1-b2)
54      w[8]=w[8]-(rate/(math.sqrt(vt)+z))*mt
55   loss=mean_squared_error(y,y_pred)
56   loss_value_adam.append(loss)
57
```

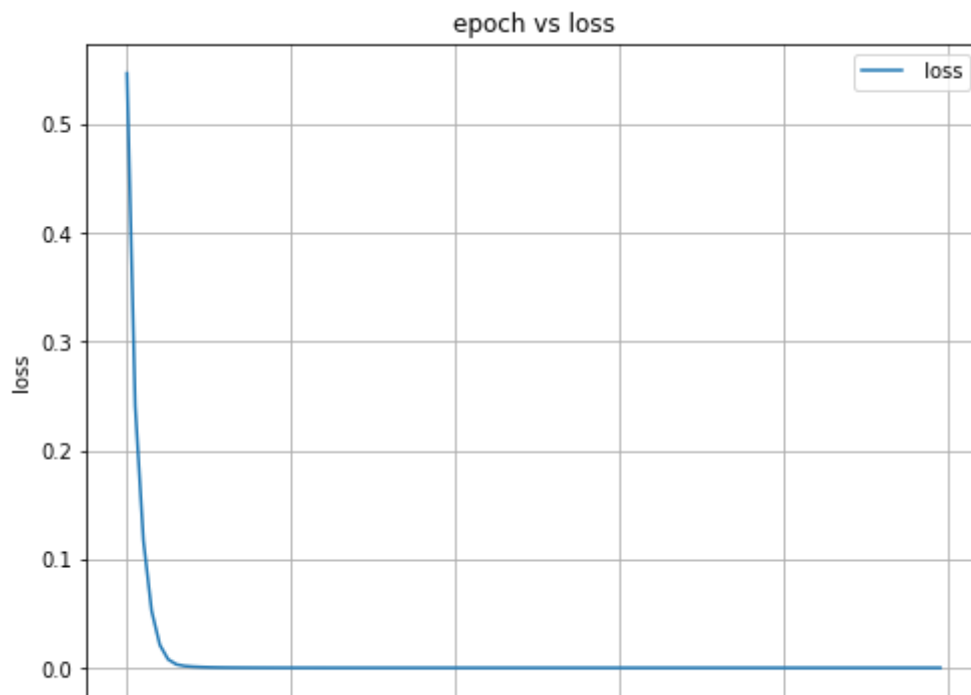## Plot between epochs and loss

```
 1
 2 %matplotlib inline
 3 import matplotlib.pyplot as plt
 4 plt.figure(figsize=(8,6))
 5 plt.grid()
 6 plt.plot(epoc_adam,loss_value_adam, label=' loss')
 7 plt.title("epoch vs loss")
 8 plt.xlabel("epoch_number")
 9 plt.ylabel("loss")
10 plt.legend()
```
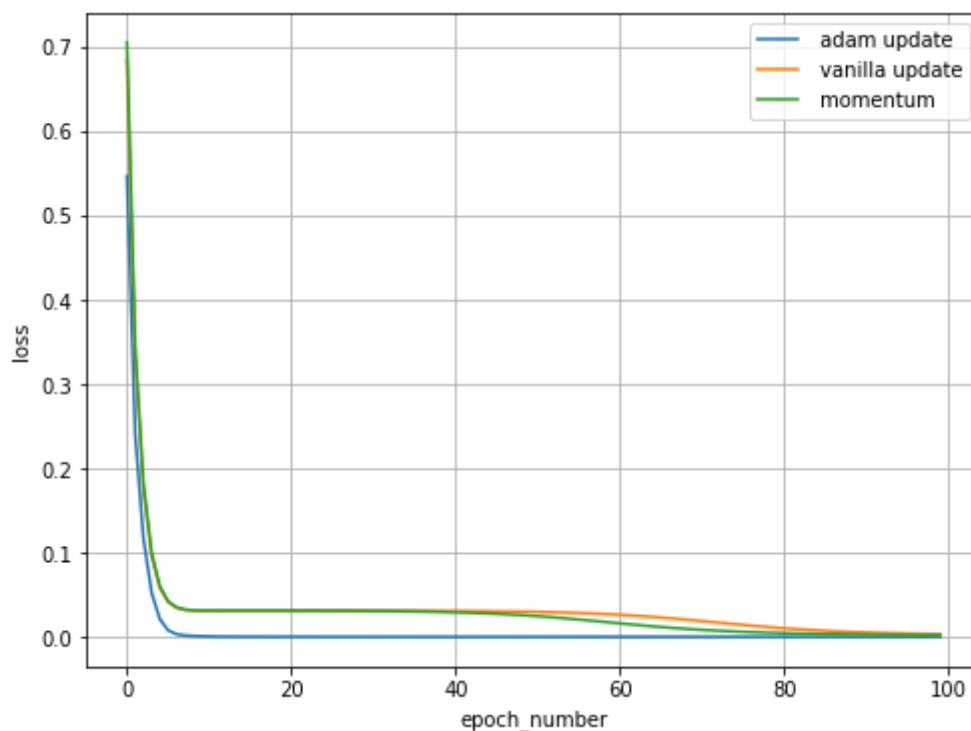
```
<matplotlib.legend.Legend at 0x7ff9ccc15438>
```



## Comparision plot between epochs and loss with different optimizers

```
1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 plt.figure(figsize=(8,6))
4 plt.grid()
5 plt.plot(epoc_adam,loss_value_adam, label=' adam update')
6 plt.plot(epoc_vanilla,loss_value_vanilla, label=' vanilla update')
7 plt.plot(epoc_momentum,loss_value_momentum, label=' momentum')
8 plt.xlabel("epoch_number")
9 plt.ylabel("loss")
10 plt.legend()
```

```
<matplotlib.legend.Legend at 0x7ff9ccb0f978>
```

- from the above graph we can understand that adam optmizer convergence faster.

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.