

# BackPropagation

There will be some functions that start with the word "grader" ex: grader\_sigmoid(), grader\_forwardprop(), grader\_backprop() etc, you should not change those function definition.

Every Grader function has to return True.

## Loading data

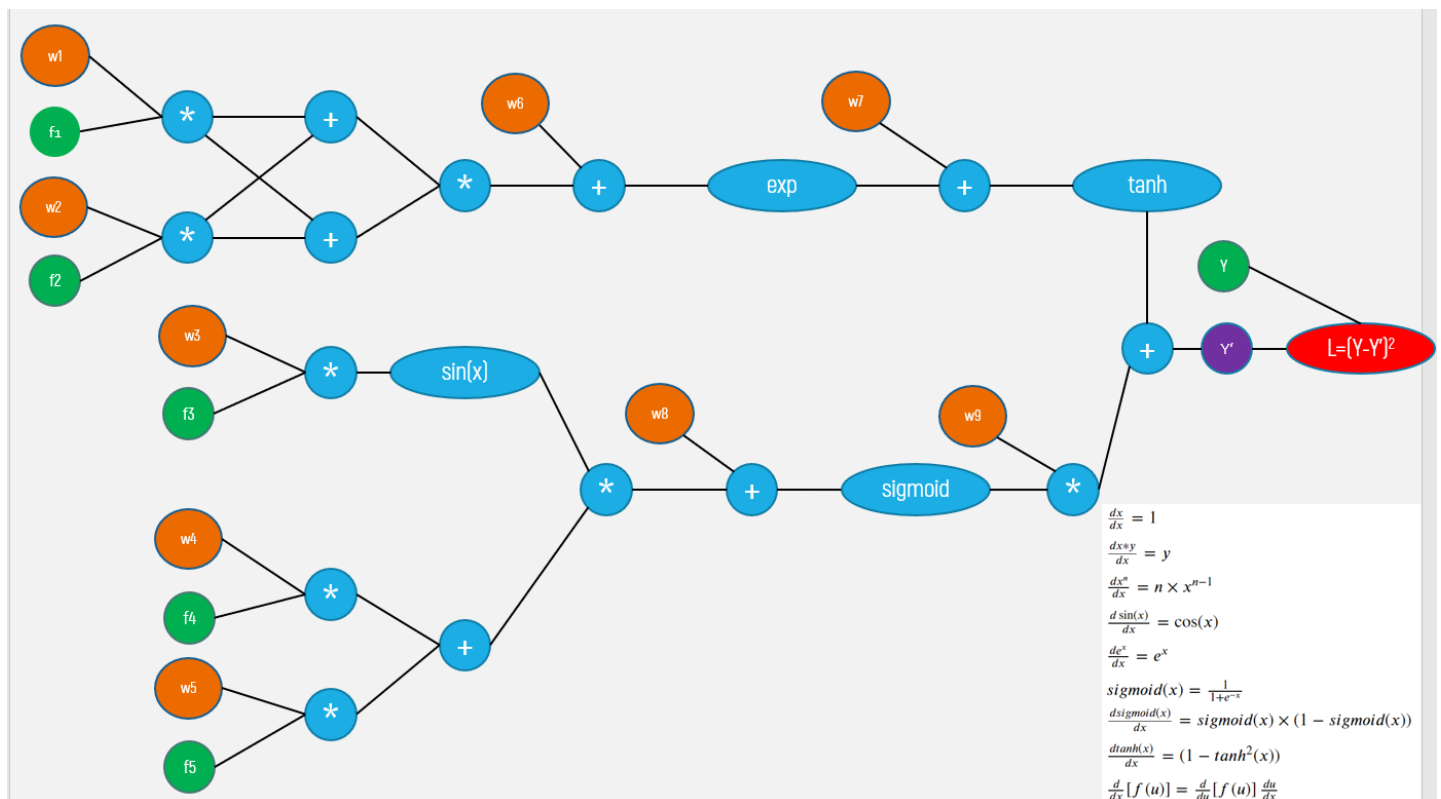
In [1]:

```
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)
```

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

## Computational graph



- If you observe the graph, we are having input features  $[f_1, f_2, f_3, f_4, f_5]$  and 9 weights  $[w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8, w_9]$ .
- The final output of this graph is a value  $L$  which is computed as  $(Y - Y')^2$

# Task 1: Implementing backpropagation and Gradient checking

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

In [2]:

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

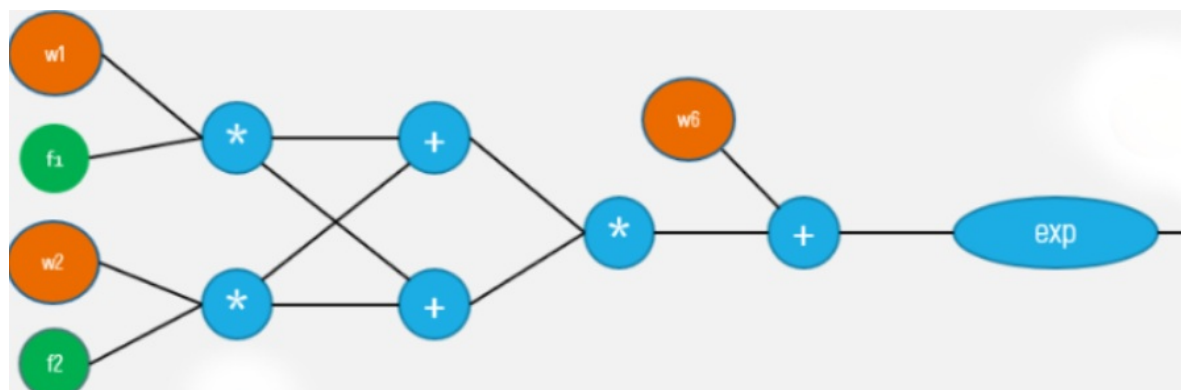
Out[2]:

- Write two functions

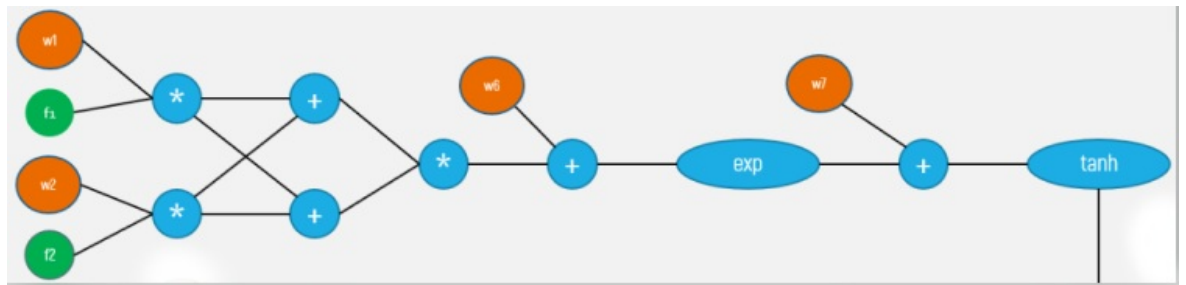
- Forward propagation (Write your code in `def forward_propagation()`)

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

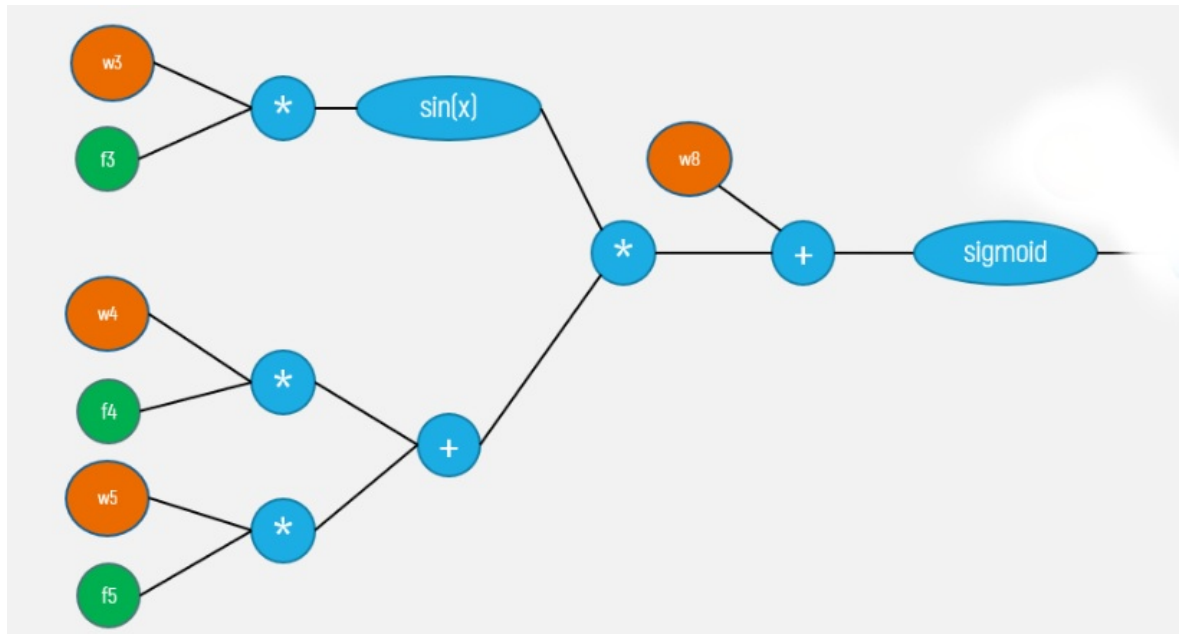
**Part 1**



**Part 2**



### Part 3



```

def forward_propagation(X, y, W):

    # X: input data point, note that in this assignment you are having 5-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)
    # now compute remaining values from computational 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 propagation)
  
```

- **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,..
    .,w9]
    # 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

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

- 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  $f$  w.r.t  $w_1$  is

$$\begin{aligned} \frac{df}{dw_1} &= dw_1 = 2 \cdot w_1 \cdot x_1 \\ &= 2 \cdot 1 \cdot 3 \\ &= 6 \end{aligned}$$

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

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

$$\begin{aligned}
& \frac{(1.00020001.3+2.4)-(0.99980001.3+2.4)}{2*0.0001} \\
&= \frac{(11.00060003)-(10.99940003)}{0.0002} \\
&= 5.999999999999
\end{aligned}$$

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

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.

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

you can mathamatically derive the same thing like this

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

## Implement Gradient checking

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

### Algorithm

```

W = initilize_randomly
def gradient_checking(data_point, W):

    # compute the L value using forward_propagation()
    # compute the gradients of W using backward_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 weights
        # subtract a small value to weight wi, and then find the values of L with
        h 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 backward_propagation() with the apr
oximation gradients of weights 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

## Task 2 : Optimizers

- As a part of this task, you will be implementing 3 type of optimizers(methods to update weight)
- 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

In [10]:

```

from IPython.display import YouTubeVideo
YouTubeVideo('gYpoJm1gyXA',width="1000",height="500")

```

### Algorithm

```

for each epoch(1-100):
    for each data point in your data:
        using the functions forward_propagation() and backward_propagation() c
ompute the gradients of weights
        update the weigts with help of gradients ex: w1 = w1-learning_rate*dw
1

```

## 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 .

## Task 1

### Forward propagation

## Forward propagation

In [8]:

```
def sigmoid(z):  
    '''In this function, we will compute the sigmoid(z)'''  
    # we can use this function in forward and backward propagation  
    return 1/(1 + np.exp(-z))
```

In [29]:

```
def forward_propagation(x, y, w):  
    '''In this function, we will compute the forward propagation '''  
  
    ## Part1  
    q1 = w[0]*x[0]  
    q2 = w[1]*x[1]  
    q3 = q1 + q2  
    q4 = q3 * q3  
    q5 = q4 + w[5]  
    exp= np.exp(q5)  
  
    ## Part2  
    q7 = exp + w[6]  
    tanh = np.tanh(q7)  
  
    ## Part3  
    q9 = w[2]*x[2]  
    q10 = np.sin(q9)  
    q11 = w[3]*x[3]  
    q12 = w[4]*x[4]  
    q13 = q11 + q12  
    q14 = q13 * q10  
    q15 = q14 + w[7]  
  
    ## sigmoid  
    sig_moid = sigmoid(q15)  
    q16 = sig_moid*w[8]  
  
    ## tanh  
    y_hat = q16 + tanh  
  
    ## loss-function  
    L = np.square(y-y_hat)  
  
    ## derivative of loss-function  
    dl = (-2)*(y-y_hat)  
  
    my_dict = {  
        'exp'      : exp,  
        'sigmoid'  : sig_moid,  
        'tanh'     : tanh,  
        'loss'     : L,  
        'dy_pr'    : dl,  
        'sin'      : q10,  
        'cos'      : np.cos(q9)  
    }  
    return my_dict
```

## Grader function - 1

In [30]:

```
def grader_sigmoid(z):  
    val=sigmoid(z)  
    assert (val==0.8807970779778823)  
    return True  
grader_sigmoid(2)
```

Out[30]:

True

## Grader function - 2

In [36]:

```
def grader_forwardprop(data):
    d1 = (data['dy_pr']==-1.9285278284819143)
    loss=(data['loss']==0.9298048963072919)
    part1=(data['exp']==1.1272967040973583)
    part2=(data['tanh']==0.8417934192562146)
    part3=(data['sigmoid']==0.5279179387419721)
    assert(d1 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[36]:

True

In [37]:

```
print (d1)
```

```
{'exp': 1.1272967040973583, 'sigmoid': 0.5279179387419721, 'tanh': 0.8417934192562146, 'loss': 0.9298048963072919, 'dy_pr': -1.9285278284819143, 'sin': -0.14538296400984968, 'cos': 0.9893754564247643}
```

## Backward propagation

In [26]:

```
def backward_propagation(L,W,my_dict):
    '''In this function, we will compute the backward propagation'''
    # Hint: you can use dict type to store the required variables
    # dw1 = # in dw1 compute derivative of L w.r.to w1
    dw1 = my_dict['dy_pr'] * (1-np.square(my_dict['tanh'])) * my_dict['exp'] * 2*((W[0]*
L[0]+W[1]*L[1])*L[0])

    # dw2 = # in dw2 compute derivative of L w.r.to w2
    dw2 = my_dict['dy_pr'] * (1-np.square(my_dict['tanh'])) * my_dict['exp'] * 2*((W[1]*
L[1]+W[0]*L[0])*L[1])

    # dw3 = # in dw3 compute derivative of L w.r.to w3
    dw3 = my_dict['dy_pr'] * W[8] * my_dict['sigmoid'] * (1-my_dict['sigmoid']) * (L[3]*
W[3]+L[4]*W[4]) * L[2] * my_dict['cos']

    # dw4 = # in dw4 compute derivative of L w.r.to w4
    dw4 = my_dict['dy_pr'] * W[8] * my_dict['sigmoid'] * (1-my_dict['sigmoid']) * L[3] *
my_dict['sin']

    # dw5 = # in dw5 compute derivative of L w.r.to w5
    dw5 = my_dict['dy_pr'] * W[8] * my_dict['sigmoid'] * (1-my_dict['sigmoid']) * L[4] *
my_dict['sin']

    # dw6 = # in dw6 compute derivative of L w.r.to w6
    dw6 = my_dict['dy_pr'] * (1-np.square(my_dict['tanh'])) * my_dict['exp']

    # dw7 = # in dw7 compute derivative of L w.r.to w7
    dw7 = my_dict['dy_pr'] * (1-np.square(my_dict['tanh']))

    # dw8 = # in dw8 compute derivative of L w.r.to w8
    dw8 = my_dict['dy_pr'] * W[8] * my_dict['sigmoid'] * (1-my_dict['sigmoid'])

    # dw9 = # in dw9 compute derivative of L w.r.to w9
    dw9 = my_dict['sigmoid'] * my_dict['dy_pr']
```



```

dW={
    'dw1':dw1,
    'dw2':dw2,
    'dw3':dw3,
    'dw4':dw4,
    'dw5':dw5,
    'dw6':dw6,
    'dw7':dw7,
    'dw8':dw8,
    'dw9':dw9
}

return dW

```

### Grader function - 3

In [38]:

```

def grader_backprop(data):
    dw1=(data['dw1']==-0.22973323498702003)
    dw2=(data['dw2']==-0.021407614717752925)
    dw3=(data['dw3']==-0.005625405580266319)
    dw4=(data['dw4']==-0.004657941222712423)
    dw5=(data['dw5']==-0.0010077228498574246)
    dw6=(data['dw6']==-0.6334751873437471)
    dw7=(data['dw7']==-0.561941842854033)
    dw8=(data['dw8']==-0.04806288407316516)
    dw9=(data['dw9']==-1.0181044360187037)
    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)
d1=backward_propagation(X[0],w,d1)
grader_backprop(d1)

```

Out[38]:

True

## Implement gradient checking

In [44]:

```

W = np.random.rand(9)
e = 0.0001
def gradient_checking(data_point, W):

    # compute the L value using forward_propagation()
    dl = forward_propagation(data_point,y[0],W)

    # compute the gradients of W using backward_propagation()
    grads_w = backward_propagation(data_point,W,dl)
    grads_w = list(grads_w.values())

    approx_gradients = []

    for i in range(0,len(W)):

        # add a small value to weight wi
        w_add = W.copy()
        w_add[i] += e
        # then find the values of L with the updated weights
        L_add = forward_propagation(data_point,y[0],w_add)
        grads_w_add = backward_propagation(data_point,w_add,L_add)
        L_add = L_add['loss']

        # subtract a small value to weight wi
        w_sub = W.copy()
        w_sub[i] -= e

```

```

# then find the values of L with the updated weights
L_sub = forward_propagation(data_point,y[0],w_sub)
grads_w_sub = backward_propagation(data_point,w_sub,L_sub)
L_sub = L_sub['loss']

# compute the approximation gradients of weight wi
approx_grads = (L_add - L_sub) / (2*e)
approx_gradients.append(approx_grads)

# compare the gradient of weights W from backward_propagation() with the approximation
gradients of weights with gradient_check formula
gradient_check = []
for i in range(0,len(W)):
    num_term = np.linalg.norm(grads_w[i] - approx_gradients[i] )
    deno_term = np.linalg.norm(grads_w[i]) + np.linalg.norm(approx_gradients[i])
    diff = num_term / deno_term
    gradient_check.append(diff)
return gradient_check

```

In [45]:

```
g = gradient_checking(X[0],W)
```

In [46]:

```
g
```

Out[46]:

```

[3.874929335698234e-08,
 4.3157001563962967e-10,
 5.1361566763371076e-11,
 4.154145689622569e-11,
 1.2818120768150566e-11,
 2.680286138793292e-08,
 3.6318330152259967e-09,
 1.1754813341143227e-10,
 7.762460940103805e-13]

```

## Task 2: Optimizers

### Algorithm with Vanilla update of weights

In [47]:

```
W = list(np.random.normal(0.0, 0.01, 9))
W
```

Out[47]:

```

[-0.0025199913735369835,
 -0.0011575076197623772,
 -0.002123680793884558,
 -0.0020039254684206934,
 0.00797482240391943,
 0.007205892985395164,
 -0.008145617506374811,
 0.0016209220498196748,
 0.0039333375530683129]

```

In [50]:

```

W_van      = W
loss_van   = []

for epoch in range(100):
    for i,j in zip(X,y):
        dl = forward_propagation(i,j,W_van)
        loss_v = dl['loss']

```

```

dw      = backward_propagation(i,W_van,dl)
dw      = list(dw.values())
dw      = [i * 0.01 for i in dw]
W_van  = np.subtract(W_van,dw)

```

```
loss_van.append(loss_v)
```

### Plot between epochs and loss

In [51]:

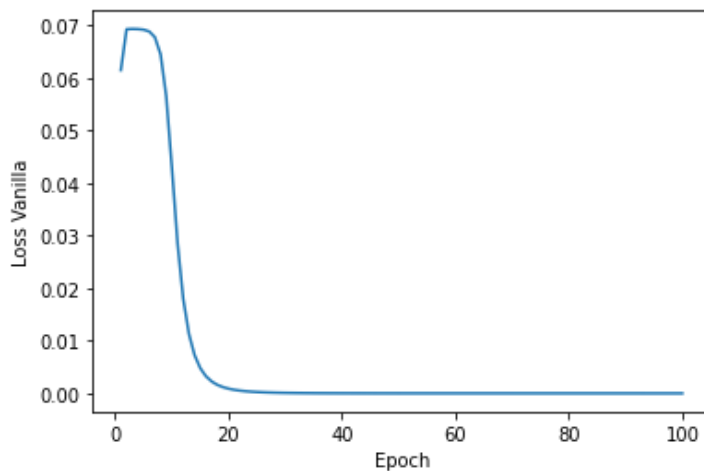
```

import matplotlib.pyplot as plt
epoch=list(range(1,101))
plt.plot(epoch,loss_van)
plt.xlabel('Epoch')
plt.ylabel('Loss Vanilla')

```

Out[51]:

Text(0, 0.5, 'Loss Vanilla')



### Algorithm with Momentum update of weights

In [64]:

```

w_mom = W
v = list(np.zeros(9))
loss_mom=[]

for epoch in range(100):
    for i,j in zip(X,y):
        dl = forward_propagation(i,j,w_mom)
        loss_m = dl['loss']

        dw = backward_propagation(i,w_mom,dl)
        dw = list(dw.values())

        for i in range(len(dw)):
            dw[i] = v[i] = 0.9*v[i] - 0.01*dw[i]

        w_mom = np.add(w_mom,v)
    loss_mom.append(loss_m)

```

### Plot between epochs and loss

In [65]:

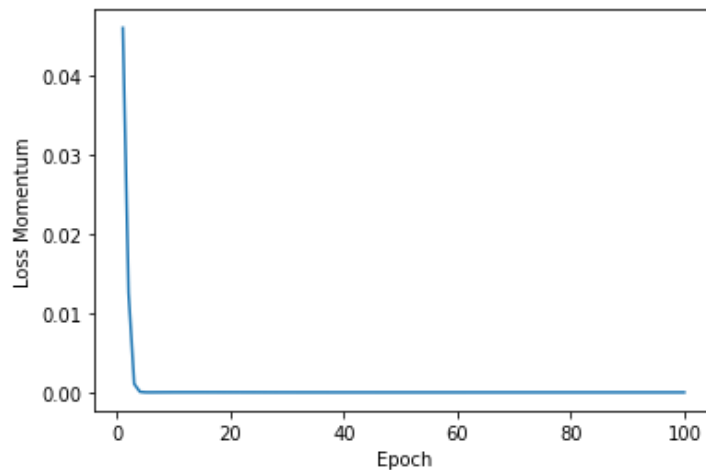
```

epoch=list(range(1,101))
plt.plot(epoch,loss_mom)
plt.xlabel('Epoch')
plt.ylabel('Loss Momentum')

```

Out [65]:

```
Text(0, 0.5, 'Loss Momentum')
```



## Algorithm with Adam update of weights

In [66]:

```
w_adam = W
v_adam = list(np.zeros(9))
m_adam = list(np.zeros(9))
loss_adam = []

for epoch in range(100):
    for i,j in zip(X,y):
        dl = forward_propagation(i,j,w_adam)
        loss_a = dl['loss']

    dw = backward_propagation(i,w_adam,dl)
    dw = list(dw.values())

    for i in range(len(dw)):
        m_adam[i] = 0.9*m_adam[i] + (1-0.9)*dw[i]
        v_adam[i] = 0.999*v_adam[i] + (1-0.999)*(dw[i]**2)
        w_adam[i] += ((-0.1)* m_adam[i]) / (np.sqrt(v_adam[i]) + 1e-8)

    loss_adam.append(loss_a)
```

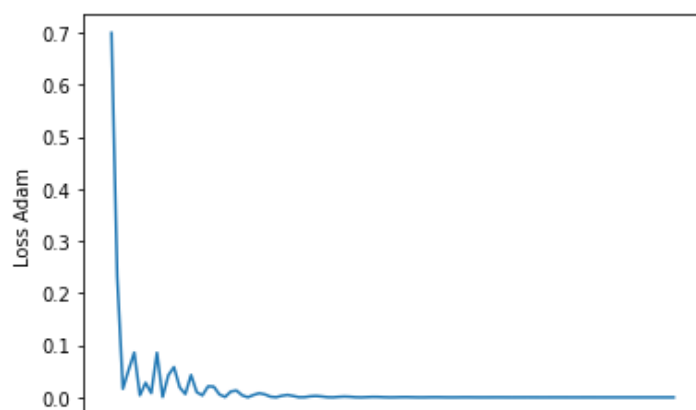
## Plot between epochs and loss

In [67]:

```
epoch=list(range(1,101))
plt.plot(epoch,loss_adam)
plt.xlabel('Epoch')
plt.ylabel('Loss Adam')
```

Out [67]:

```
Text(0, 0.5, 'Loss Adam')
```

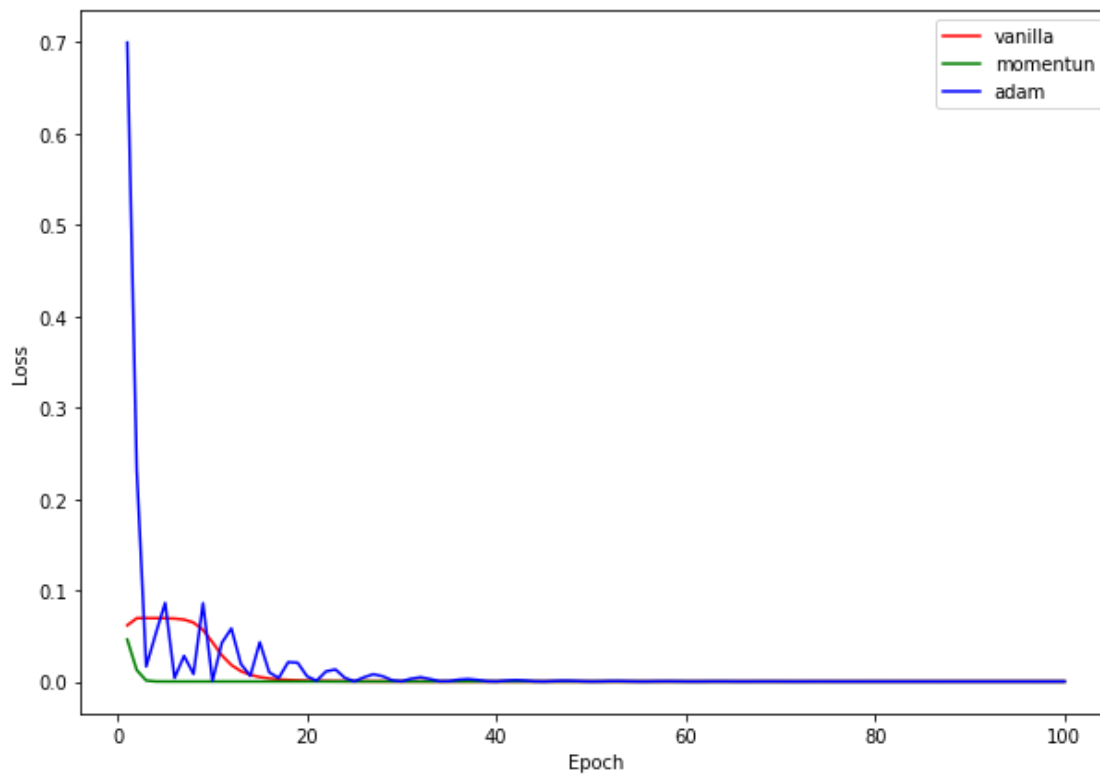


0 20 40 60 80 100  
Epoch

## Comparison plot between epochs and loss with different optimizers

In [92]:

```
plt.figure(figsize=(10, 7))  
plt.plot(epoch, loss_van, 'r')  
plt.plot(epoch, loss_mom, 'g')  
plt.plot(epoch, loss_adam, 'b')  
plt.legend(['vanilla', 'momentun', 'adam'])  
plt.xlabel('Epoch')  
plt.ylabel('Loss')  
plt.show()
```



In [ ]: