

# Implement SGD Classifier with Logloss and L2 regularization Using SGD without using sklearn

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

Every Grader function has to return True.

## Importing packages

In [1]:

```
import numpy as np
import pandas as pd
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn import linear_model
```

## Creating custom dataset

In [2]:

```
# please don't change random_state
X, y = make_classification(n_samples=50000, n_features=15, n_informative=10, n_redundant=5,
                          n_classes=2, weights=[0.7], class_sep=0.7, random_state=15)
# make_classification is used to create custom dataset
# Please check this link (https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make\_classification.html) for more details
```

In [3]:

```
X.shape, y.shape
```

Out[3]:

```
((50000, 15), (50000,))
```

## Splitting data into train and test

In [4]:

```
#please don't change random state
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=15)
```

In [5]:

```
X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

Out[5]:

```
((37500, 15), (37500,)), (12500, 15), (12500,))
```

## SGD classifier

In [6]:

```
# alpha : float
# Constant that multiplies the regularization term.

# eta0 : double
# The initial learning rate for the 'constant', 'invscaling' or 'adaptive' schedules.
```

```
clf = linear_model.SGDClassifier(eta0=0.0001, alpha=0.0001, loss='log', random_state=15, penalty='l2', tol=1e-3, verbose=2, learning_rate='constant')
clf
# Please check this documentation (https://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier.html)
```

Out[6]:

```
SGDClassifier(alpha=0.0001, average=False, class_weight=None,
              early_stopping=False, epsilon=0.1, eta0=0.0001,
              fit_intercept=True, l1_ratio=0.15, learning_rate='constant',
              loss='log', max_iter=1000, n_iter_no_change=5, n_jobs=None,
              penalty='l2', power_t=0.5, random_state=15, shuffle=True,
              tol=0.001, validation_fraction=0.1, verbose=2, warm_start=False)
```

In [7]:

```
clf.fit(X=X_train, y=y_train) # fitting our model
```

```
-- Epoch 1
Norm: 0.77, NNZs: 15, Bias: -0.316653, T: 37500, Avg. loss: 0.455552
Total training time: 0.01 seconds.
-- Epoch 2
Norm: 0.91, NNZs: 15, Bias: -0.472747, T: 75000, Avg. loss: 0.394686
Total training time: 0.02 seconds.
-- Epoch 3
Norm: 0.98, NNZs: 15, Bias: -0.580082, T: 112500, Avg. loss: 0.385711
Total training time: 0.03 seconds.
-- Epoch 4
Norm: 1.02, NNZs: 15, Bias: -0.658292, T: 150000, Avg. loss: 0.382083
Total training time: 0.03 seconds.
-- Epoch 5
Norm: 1.04, NNZs: 15, Bias: -0.719528, T: 187500, Avg. loss: 0.380486
Total training time: 0.04 seconds.
-- Epoch 6
Norm: 1.05, NNZs: 15, Bias: -0.763409, T: 225000, Avg. loss: 0.379578
Total training time: 0.05 seconds.
-- Epoch 7
Norm: 1.06, NNZs: 15, Bias: -0.795106, T: 262500, Avg. loss: 0.379150
Total training time: 0.06 seconds.
-- Epoch 8
Norm: 1.06, NNZs: 15, Bias: -0.819925, T: 300000, Avg. loss: 0.378856
Total training time: 0.07 seconds.
-- Epoch 9
Norm: 1.07, NNZs: 15, Bias: -0.837805, T: 337500, Avg. loss: 0.378585
Total training time: 0.08 seconds.
-- Epoch 10
Norm: 1.08, NNZs: 15, Bias: -0.853138, T: 375000, Avg. loss: 0.378630
Total training time: 0.08 seconds.
Convergence after 10 epochs took 0.08 seconds
```

Out[7]:

```
SGDClassifier(alpha=0.0001, average=False, class_weight=None,
              early_stopping=False, epsilon=0.1, eta0=0.0001,
              fit_intercept=True, l1_ratio=0.15, learning_rate='constant',
              loss='log', max_iter=1000, n_iter_no_change=5, n_jobs=None,
              penalty='l2', power_t=0.5, random_state=15, shuffle=True,
              tol=0.001, validation_fraction=0.1, verbose=2, warm_start=False)
```

In [8]:

```
clf.coef_, clf.coef_.shape, clf.intercept_
#clf.coef_ will return the weights
#clf.coef_.shape will return the shape of weights
#clf.intercept_ will return the intercept term
```

Out[8]:

```
(array([[ -0.42336692,  0.18547565, -0.14859036,  0.34144407, -0.2081867 ,
          0.56016579, -0.45242483, -0.09408813,  0.2092732 ,  0.18084126,
          0.19705191,  0.00421916, -0.0796037 ,  0.33852802,  0.02266721]]),
 (1, 15)).
```

```
array([-0.8531383]))
```

```
# This is formatted as code
```

## Implement Logistic Regression with L2 regularization Using SGD: without using sklearn

1. We will be giving you some functions, please write code in that functions only.
2. After every function, we will be giving you expected output, please make sure that you get that output.

- Initialize the weight\_vector and intercept term to zeros (Write your code in `def initialize_weights()`)
- Create a loss function (Write your code in `def logloss()`)

$\text{log loss} = -1 * \frac{1}{n} \sum_{\text{for each } Y_t, Y_{\text{pred}}} (Y_t \log_{10}(Y_{\text{pred}}) + (1 - Y_t) \log_{10}(1 - Y_{\text{pred}}))$

- for each epoch:
  - for each batch of data points in train: (keep batch size=1)
    - calculate the gradient of loss function w.r.t each weight in weight vector (write your code in `def gradient_dw()`)  
$$\frac{dw^{\{t\}}}{dt} = x_n(y_n - \sigma((w^{\{t\}})^T x_n + b^{\{t\}})) - \frac{\lambda}{N} w^{\{t\}}$$
    - Calculate the gradient of the intercept (write your code in `def gradient_db()`) [check this](#)  
$$\frac{db^{\{t\}}}{dt} = y_n - \sigma((w^{\{t\}})^T x_n + b^{\{t\}})$$
    - Update weights and intercept (check the equation number 32 in the above mentioned [pdf](#)):  
$$w^{\{t+1\}} \leftarrow w^{\{t\}} + \alpha \frac{dw^{\{t\}}}{dt}$$
  
$$b^{\{t+1\}} \leftarrow b^{\{t\}} + \alpha \frac{db^{\{t\}}}{dt}$$
  - calculate the log loss for train and test with the updated weights (you can check the python assignment 10th question)
  - And if you wish, you can compare the previous loss and the current loss, if it is not updating, then you can stop the training
  - append this loss in the list ( this will be used to see how loss is changing for each epoch after the training is over )

### Initialize weights

In [9]:

```
def initialize_weights(dim):
    ''' In this function, we will initialize our weights and bias'''
    #initialize the weights to zeros array of (1,dim) dimensions
    #you use zeros_like function to initialize zero, check this link
    https://docs.scipy.org/doc/numpy/reference/generated/numpy.zeros_like.html
    #initialize bias to zero

    w = np.zeros_like(dim)
    b = 0

    return w,b
```

In [10]:

```
dim=X_train[0]
w,b = initialize_weights(dim)
print('w =', (w))
print('b =',str(b))
```

```
w = [0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
b = 0
```

### Grader function - 1

In [11]:

```
dim=X_train[0]
w,b = initialize_weights(dim)
```

```
def grader_weights(w,b):
    assert((len(w)==len(dim)) and b==0 and np.sum(w)==0.0)
    return True
grader_weights(w,b)
```

Out[11]:

True

### Compute sigmoid

$\text{sigmoid}(z) = 1/(1+\exp(-z))$

In [12]:

```
def sigmoid(z):
    ''' In this function, we will return sigmoid of z'''
    # compute sigmoid(z) and return

    sigma = (1/(1+np.exp(-z)))

    return sigma
```

### Grader function - 2

In [13]:

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

Out[13]:

True

### Compute loss

$\text{log loss} = -\frac{1}{n} \sum_{\text{for each } Y_t, Y_{\text{pred}}} (Y_t \log_{10}(Y_{\text{pred}}) + (1-Y_t) \log_{10}(1-Y_{\text{pred}}))$

In [14]:

```
def logloss(y_true,y_pred):
    '''In this function, we will compute log loss'''

    n = len(y_true)
    x = np.log10(y_pred)
    x1 = np.log10(np.ones_like(y_pred) - y_pred)
    loss = 0
    for j in range(n):
        loss = loss + ((y_true[j]*x[j]) + ((1-y_true[j])*x1[j]))
    loss*=(-1/n)
    return loss
```

### Grader function - 3

In [15]:

```
def grader_logloss(true,pred):
    loss=logloss(true,pred)
    assert(loss==0.07644900402910389)
    return True
true=[1,1,0,1,0]
pred=[0.9,0.8,0.1,0.8,0.2]
grader_logloss(true,pred)
```

Out[15]:

True

#### Compute gradient w.r.to 'w'

$$\frac{dw}{dt} = x_n(y_n - \sigma((w^{(t)})^T x_n + b^{(t)})) - \frac{\lambda}{N} w^{(t)}$$

In [16]:

```
def gradient_dw(x,y,w,b,alpha,N):  
    '''In this function, we will compute the gradient w.r.to w'''  
  
    dw = (x*(y - sigmoid(np.matmul(w,x) + b))) - ((alpha/N)*w)  
  
    return dw
```

#### Grader function - 4

In [17]:

```
def grader_dw(x,y,w,b,alpha,N):  
    grad_dw=gradient_dw(x,y,w,b,alpha,N)  
    assert(np.sum(grad_dw)==2.613689585)  
    return True  
grad_x=np.array([-2.07864835,  3.31604252, -0.79104357, -3.87045546, -1.14783286,  
                 -2.81434437, -0.86771071, -0.04073287,  0.84827878,  1.99451725,  
                 3.67152472,  0.01451875,  2.01062888,  0.07373904, -5.54586092])  
grad_y=0  
grad_w,grad_b=initialize_weights(grad_x)  
alpha=0.0001  
N=len(X_train)  
grader_dw(grad_x,grad_y,grad_w,grad_b,alpha,N)
```

Out [17]:

True

#### Compute gradient w.r.to 'b'

$$\frac{db}{dt} = y_n - \sigma((w^{(t)})^T x_n + b^{(t)})$$

In [18]:

```
def gradient_db(x,y,w,b):  
    '''In this function, we will compute gradient w.r.to b'''  
  
    db = (y - sigmoid(np.matmul(w,x) + b))  
  
    return db
```

#### Grader function - 5

In [19]:

```
def grader_db(x,y,w,b):  
    grad_db=gradient_db(x,y,w,b)  
    assert(grad_db== -0.5)  
    return True  
grad_x=np.array([-2.07864835,  3.31604252, -0.79104357, -3.87045546, -1.14783286,  
                 -2.81434437, -0.86771071, -0.04073287,  0.84827878,  1.99451725,  
                 3.67152472,  0.01451875,  2.01062888,  0.07373904, -5.54586092])  
grad_y=0  
grad_w,grad_b=initialize_weights(grad_x)  
alpha=0.0001  
N=len(X_train)  
grader_db(grad_x,grad_y,grad_w,grad_b)
```

Out [19]:

True

## Implementing logistic regression

In [20]:

```
def train(X_train,y_train,X_test,y_test,epochs,alpha,eta0):
    ''' In this function, we will implement logistic regression'''
    #Here eta0 is learning rate
    #implement the code as follows
    # initialize the weights (call the initialize_weights(X_train[0]) function)
    # for every epoch
        # for every data point(X_train,y_train)
            #compute gradient w.r.to w (call the gradient_dw() function)
            #compute gradient w.r.to b (call the gradient_db() function)
            #update w, b
        # predict the output of x_train[for all data points in X_train] using w,b
        #compute the loss between predicted and actual values (call the loss function)
        # store all the train loss values in a list
        # predict the output of x_test[for all data points in X_test] using w,b
        #compute the loss between predicted and actual values (call the loss function)
        # store all the test loss values in a list
        # you can also compare previous loss and current loss, if loss is not updating then stop the process and return w,b

    w,b = initialize_weights(X_train[0])
    N = len(X_train)

    loss_train = []
    loss_test = []
    for i in range(epochs):
        w_prev,b_prev = w,b
        y_pred_train = []
        y_pred_test = []
        for j in range(len(X_train)):
            dw = gradient_dw(X_train[j,:], y_train[j], w, b, alpha, N)
            db = gradient_db(X_train[j,:], y_train[j], w, b)

            w = w + (eta0 * dw)
            b = b + (eta0 * db)
        for j in range(len(X_train)):
            y_pred_train.append(sigmoid(np.matmul(w,X_train[j,:] + b))

        present=logloss(y_train, y_pred_train)
        loss_train.append(present)
        print(present)

        for j in range(len(X_test)):
            y_pred_test.append(sigmoid(np.matmul(w,X_test[j,:] + b))
        loss_test.append(logloss(y_test, y_pred_test))

        if (i!=0 and loss_train[i-1]<present):
            return w_prev,b_prev,loss_train[:i],loss_test[:i]

    return w,b,loss_train,loss_test
```

In [21]:

```
alpha=0.0001
eta0=0.0001
N=len(X_train)
epochs=1000
w,b,train_loss,test_loss = train(X_train,y_train,X_test,y_test,epochs,alpha,eta0)
```

```
0.1754574844285461
0.16867157050333045
0.1663916799246292
0.16536827537403162
0.16485707459547086
0.16458820012928274
0.16444271323364384
0.16436263615826988
0.1643180694666775
```

0.16429307374132515  
0.1642789743093407  
0.16427098545835503  
0.1642664419100352  
0.16426384911424854  
0.16426236468266475  
0.16426151190514932  
0.16426102013167446  
0.16426073527505572  
0.16426056938842287  
0.16426047215122602  
0.164260414696778  
0.1642603804172862  
0.16426035972510467  
0.16426034706229276  
0.16426033919038657  
0.16426033421042477  
0.1642603310001301  
0.1642603288898659  
0.16426032747541439  
0.16426032650945105  
0.16426032583825456  
0.16426032536459445  
0.1642603250258187  
0.1642603247807538  
0.16426032460180723  
0.1642603244701486  
0.16426032437269608  
0.16426032430021448  
0.16426032424610498  
0.16426032420559736  
0.16426032417520106  
0.16426032415235473  
0.16426032413516098  
0.16426032412221028  
0.16426032411244418  
0.16426032410507813  
0.1642603240995188  
0.16426032409532287  
0.16426032409215466  
0.16426032408976068  
0.16426032408795205  
0.16426032408658778  
0.1642603240855564  
0.16426032408477767  
0.1642603240841884  
0.16426032408374264  
0.1642603240834062  
0.16426032408315108  
0.16426032408296104  
0.1642603240828153  
0.1642603240827063  
0.16426032408262436  
0.16426032408256108  
0.1642603240825134  
0.16426032408247754  
0.16426032408244948  
0.1642603240824297  
0.16426032408241492  
0.16426032408240282  
0.16426032408239452  
0.16426032408238678  
0.16426032408238248  
0.16426032408237706  
0.16426032408237606  
0.16426032408237418  
0.16426032408237204  
0.16426032408237196  
0.16426032408236904  
0.16426032408236896  
0.1642603240823696

0.16426032408236896 is the minimum loss achieved through above SGD algorithm.

# This is formatted as code

## Goal of assignment

Compare your implementation and SGDClassifier's the weights and intercept, make sure they are as close as possible i.e difference should be in terms of  $10^{-3}$

In [25]:

```
# these are the results we got after we implemented sgd and found the optimal weights and intercept
w-clf.coef_, b-clf.intercept_
```

Out[25]:

```
(array([[ -0.00642561,  0.00755958,  0.00012038, -0.00335043, -0.01309575,
         0.0097832 ,  0.00724319,  0.00418412,  0.01255635, -0.00701157,
         0.0016966 , -0.00480349, -0.00173046,  0.00056209,  0.00032076]]),
 array([-0.03911443]))
```

## Plot epoch number vs train , test loss

- epoch number on X-axis
- loss on Y-axis

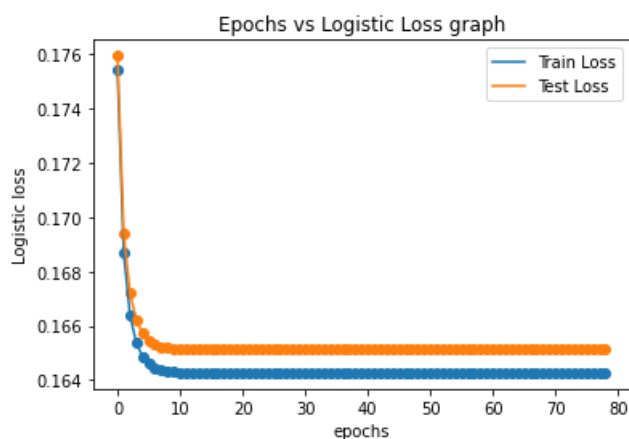
In [26]:

```
import matplotlib.pyplot as plt

epochs=[i for i in range(len(train_loss))]
plt.plot(epochs,train_loss,label= 'Train Loss')
plt.plot(epochs,test_loss,label= 'Test Loss')

plt.scatter(epochs,train_loss)
plt.scatter(epochs,test_loss)

plt.legend()
plt.title('Epochs vs Logistic Loss graph')
plt.xlabel('epochs')
plt.ylabel('Logistic loss')
plt.show()
```



We can observe from the above plots

1. In initial iterations the reduction of loss is high, but as no of epochs increases the reduction of loss is very low. This indicates that gradient of objective function is approaching towards optimal weight gradient.

In [27]:

```
def pred(w,b, X):
    N = len(X)
    predict = []
    for i in range(N):
```



```
z=np.dot(w,X[i])+b
if sigmoid(z) >= 0.5: # sigmoid(w,x,b) returns 1/(1+exp(-(dot(x,w)+b)))
    predict.append(1)
else:
    predict.append(0)
return np.array(predict)
print(1-np.sum(y_train - pred(w,b,X_train))/len(X_train))
print(1-np.sum(y_test - pred(w,b,X_test))/len(X_test))
```

0.9522133333333334

0.95