

# Binary Logistic Regression

## Dataset Description:

The CSV file contains 5172 rows, each row for each email. There are 3002 columns. The first column indicates the Email name. The name has been set with numbers, not recipients', to protect privacy. The last column has the labels for prediction: 1 for spam, and 0 for not spam. The remaining 3000 columns are the 3000 most common words in all the emails, *after excluding the non-alphabetical characters/words*.

## Dataset Preparation:

1. Select proper columns.
2. Normalize the dataset.
3. Randomly split the dataset into TRAINING(80%) and TEST(20%) sets.

## Train (update $\Theta$ ):

1. for each sample,  $X = [x_1, x_2, \dots, x_n]$  in TRAINING set:  
Concatenate 1 and turn it into  $X' = [x_1, x_2, \dots, x_n, 1]$
2. randomly initialize  $\theta = [\theta_1, \theta_2, \dots, \theta_{n+1}]$  within 0 to 1 //  $\theta_1, \theta_2, \dots, \theta_n$ :weights,  $\theta_{n+1}$ :bias
3. max\_iter = 500, lr = 0.01
4. history = list()
5. for itr in [1, max\_iter]:  
     $J = 0$  //total cost  
    for each sample,  $X'$  in TRAINING set:  
         $z = X' \cdot \theta$  //use np.dot function  
         $h = \text{sigmoid}(z)$  //sigmoid available in Python  
         $L = -y \log(h) - (1 - y) \log(1 - h)$  //h = prediction label, y = true label  
         $J += L$   
         $dw = X' \cdot (h - y)$  //dim(dw) = n + 1  
         $\theta = \theta - dw * lr$  //dim( $\theta$ ) = n + 1  
     $J /= N_{\text{train}}$  //N\_train = # of training samples  
    append J into history

## Validation:

1. correct = 0
2. For each sample in  $X'$  in the TEST set:  
     $z = X' \cdot \theta$   
     $h = \text{sigmoid}(z)$   
    if  $h \geq 0.5$ : h=1  
    else: h=0  
    if  $h == y$ : correct += 1
3.  $\text{test\_acc} = \text{correct} * 100 / N_{\text{test}}$  //N\_test = # of testing samples

## Instruction:

- Plot total cost vs iteration(History) for different learning rate,  $lr = 0.1, 0.01, 0.001$  and  $0.0001$  ( $\text{max\_iter} = 500$ )
- Make a table with 2 columns: learning rate  $lr$  and  $\text{test\_acc}$ .
- Submit a .ipynb file and a report .pdf file.
- Do not use libraries such as: SKlearn, Scikit learning for this assignment.