

# C1\_W1\_Assignment

January 23, 2022

## 1 Assignment 1: Logistic Regression

Welcome to week one of this specialization. You will learn about logistic regression. Concretely, you will be implementing logistic regression for sentiment analysis on tweets. Given a tweet, you will decide if it has a positive sentiment or a negative one. Specifically you will:

- Learn how to extract features for logistic regression given some text
- Implement logistic regression from scratch
- Apply logistic regression on a natural language processing task
- Test using your logistic regression
- Perform error analysis

We will be using a data set of tweets. Hopefully you will get more than 99% accuracy. Run the cell below to load in the packages.

### 1.1 Import functions and data

```
[1]: # run this cell to import nltk
import nltk
from os import getcwd
import w1_unittest

nltk.download('twitter_samples')
nltk.download('stopwords')
```

```
[nltk_data] Downloading package twitter_samples to
[nltk_data] /home/jovyan/nltk_data...
[nltk_data] Package twitter_samples is already up-to-date!
[nltk_data] Downloading package stopwords to /home/jovyan/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

```
[1]: True
```

#### 1.1.1 Imported functions

Download the data needed for this assignment. Check out the [documentation for the twitter\\_samples dataset](#).

- `twitter_samples`: if you're running this notebook on your local computer, you will need to download it using:

```
nltk.download('twitter_samples')
```

- `stopwords`: if you're running this notebook on your local computer, you will need to download it using:

```
nltk.download('stopwords')
```

### Import some helper functions that we provided in the `utils.py` file:

- `process_tweet`: cleans the text, tokenizes it into separate words, removes stopwords, and converts words to stems.
- `build_freqs`: this counts how often a word in the 'corpus' (the entire set of tweets) was associated with a positive label '1' or a negative label '0', then builds the 'freqs' dictionary, where each key is the (word,label) tuple, and the value is the count of its frequency within the corpus of tweets.

```
[2]: filePath = f"{getcwd()}/../tmp2/"
nltk.data.path.append(filePath)
```

```
[3]: import numpy as np
import pandas as pd
from nltk.corpus import twitter_samples

from utils import process_tweet, build_freqs
```

#### 1.1.2 Prepare the data

- The `twitter_samples` contains subsets of five thousand `positive_tweets`, five thousand `negative_tweets`, and the full set of 10,000 tweets.
  - If you used all three datasets, we would introduce duplicates of the positive tweets and negative tweets.
  - You will select just the five thousand positive tweets and five thousand negative tweets.

```
[4]: # select the set of positive and negative tweets
all_positive_tweets = twitter_samples.strings('positive_tweets.json')
all_negative_tweets = twitter_samples.strings('negative_tweets.json')
```

- Train test split: 20% will be in the test set, and 80% in the training set.

```
[5]: # split the data into two pieces, one for training and one for testing
      ↪ (validation set)
test_pos = all_positive_tweets[4000:]
train_pos = all_positive_tweets[:4000]
test_neg = all_negative_tweets[4000:]
```

```
train_neg = all_negative_tweets[:4000]

train_x = train_pos + train_neg
test_x = test_pos + test_neg
```

- Create the numpy array of positive labels and negative labels.

```
[6]: # combine positive and negative labels
train_y = np.append(np.ones((len(train_pos), 1)), np.zeros((len(train_neg), 1)), axis=0)
test_y = np.append(np.ones((len(test_pos), 1)), np.zeros((len(test_neg), 1)), axis=0)
```

```
[7]: # Print the shape train and test sets
print("train_y.shape = " + str(train_y.shape))
print("test_y.shape = " + str(test_y.shape))
```

```
train_y.shape = (8000, 1)
test_y.shape = (2000, 1)
```

- Create the frequency dictionary using the imported build\_freqs function.
  - We highly recommend that you open utils.py and read the build\_freqs function to understand what it is doing.
  - To view the file directory, go to the menu and click File->Open.

```
for y,tweet in zip(ys, tweets):
    for word in process_tweet(tweet):
        pair = (word, y)
        if pair in freqs:
            freqs[pair] += 1
        else:
            freqs[pair] = 1
```

- Notice how the outer for loop goes through each tweet, and the inner for loop steps through each word in a tweet.
- The 'freqs' dictionary is the frequency dictionary that's being built.
- The key is the tuple (word, label), such as ("happy",1) or ("happy",0). The value stored for each key is the count of how many times the word "happy" was associated with a positive label, or how many times "happy" was associated with a negative label.

```
[8]: # create frequency dictionary
freqs = build_freqs(train_x, train_y)

# check the output
print("type(freqs) = " + str(type(freqs)))
print("len(freqs) = " + str(len(freqs.keys())))
```

```
type(freqs) = <class 'dict'>
len(freqs) = 11436
```

### Expected output

```
type(freqs) = <class 'dict'>
len(freqs) = 11436
```

#### 1.1.3 Process tweet

The given function ‘process\_tweet’ tokenizes the tweet into individual words, removes stop words and applies stemming.

```
[9]: # test the function below
print('This is an example of a positive tweet: \n', train_x[0])
print('\nThis is an example of the processed version of the tweet: \n',
      ↪process_tweet(train_x[0]))
```

This is an example of a positive tweet:

```
#FollowFriday @France_Inte @PKuchly57 @Milipol_Paris for being top engaged
members in my community this week :)
```

This is an example of the processed version of the tweet:

```
['followfriday', 'top', 'engag', 'member', 'commun', 'week', ':)']
```

### Expected output

This is an example of a positive tweet:

```
#FollowFriday @France_Inte @PKuchly57 @Milipol_Paris for being top engaged members in my commu
```

This is an example of the processes version:

```
['followfriday', 'top', 'engag', 'member', 'commun', 'week', ':)']
```

## 2 Part 1: Logistic regression

### 2.0.1 Part 1.1: Sigmoid

You will learn to use logistic regression for text classification. \* The sigmoid function is defined as:

$$h(z) = \frac{1}{1 + \exp^{-z}} \quad (1)$$

It maps the input ‘z’ to a value that ranges between 0 and 1, and so it can be treated as a probability.

Figure 1

#### Instructions: Implement the sigmoid function

- You will want this function to work if z is a scalar as well as if it is an array.

Hints

numpy.exp

```
[10]: # UNQ_C1 GRADED FUNCTION: sigmoid
def sigmoid(z):
    '''
    Input:
        z: is the input (can be a scalar or an array)
    Output:
        h: the sigmoid of z
    '''

    ### START CODE HERE ###
    # calculate the sigmoid of z
    h = None
    ### END CODE HERE ###
    h = 1 / (1 + np.exp(-1 * z))

    return h
```

```
[11]: # Testing your function
if (sigmoid(0) == 0.5):
    print('SUCCESS!')
else:
    print('Oops!')

if (sigmoid(4.92) == 0.9927537604041685):
    print('CORRECT!')
else:
    print('Oops again!')
```

SUCCESS!

CORRECT!

```
[12]: # Test your function
w1_unittest.test_sigmoid(sigmoid)
```

All tests passed

## 2.0.2 Logistic regression: regression and a sigmoid

Logistic regression takes a regular linear regression, and applies a sigmoid to the output of the linear regression.

Regression:

$$z = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \theta_N x_N$$

Note that the  $\theta$  values are “weights”. If you took the deep learning specialization, we referred to the weights with the ‘w’ vector. In this course, we’re using a different variable  $\theta$  to refer to the weights.

Logistic regression

$$h(z) = \frac{1}{1 + \exp^{-z}}$$

$$z = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \theta_N x_N$$

We will refer to ‘z’ as the ‘logits’.

### 2.0.3 Part 1.2 Cost function and Gradient

The cost function used for logistic regression is the average of the log loss across all training examples:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h(z(\theta)^{(i)})) + (1 - y^{(i)}) \log(1 - h(z(\theta)^{(i)})) \quad (5)$$

\*  $m$  is the number of training examples \*  $y^{(i)}$  is the actual label of training example ‘i’. \*  $h(z^{(i)})$  is the model’s prediction for the training example ‘i’.

The loss function for a single training example is

$$Loss = -1 \times \left( y^{(i)} \log(h(z(\theta)^{(i)})) + (1 - y^{(i)}) \log(1 - h(z(\theta)^{(i)})) \right)$$

- All the  $h$  values are between 0 and 1, so the logs will be negative. That is the reason for the factor of -1 applied to the sum of the two loss terms.
- Note that when the model predicts 1 ( $h(z(\theta)) = 1$ ) and the label ‘y’ is also 1, the loss for that training example is 0.
- Similarly, when the model predicts 0 ( $h(z(\theta)) = 0$ ) and the actual label is also 0, the loss for that training example is 0.
- However, when the model prediction is close to 1 ( $h(z(\theta)) = 0.9999$ ) and the label is 0, the second term of the log loss becomes a large negative number, which is then multiplied by the overall factor of -1 to convert it to a positive loss value.  $-1 \times (1 - 0) \times \log(1 - 0.9999) \approx 9.2$ . The closer the model prediction gets to 1, the larger the loss.

```
[13]: # verify that when the model predicts close to 1, but the actual label is 0,
      ↪ the loss is a large positive value
      -1 * (1 - 0) * np.log(1 - 0.9999) # loss is about 9.2
```

```
[13]: 9.210340371976294
```

- Likewise, if the model predicts close to 0 ( $h(z) = 0.0001$ ) but the actual label is 1, the first term in the loss function becomes a large number:  $-1 \times \log(0.0001) \approx 9.2$ . The closer the prediction is to zero, the larger the loss.

```
[14]: # verify that when the model predicts close to 0 but the actual label is 1, the
      ↪ loss is a large positive value
```

```
-1 * np.log(0.0001) # loss is about 9.2
```

[14]: 9.210340371976182

**Update the weights** To update your weight vector  $\theta$ , you will apply gradient descent to iteratively improve your model's predictions.

The gradient of the cost function  $J$  with respect to one of the weights  $\theta_j$  is:

$$\nabla_{\theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h^{(i)} - y^{(i)}) x_j^{(i)} \quad (5)$$

\* 'i' is the index across all 'm' training examples. \* 'j' is the index of the weight  $\theta_j$ , so  $x_j^{(i)}$  is the feature associated with weight  $\theta_j$

- To update the weight  $\theta_j$ , we adjust it by subtracting a fraction of the gradient determined by  $\alpha$ :

$$\theta_j = \theta_j - \alpha \times \nabla_{\theta_j} J(\theta)$$

- The learning rate  $\alpha$  is a value that we choose to control how big a single update will be.

## 2.1 Instructions: Implement gradient descent function

- The number of iterations 'num\_iters' is the number of times that you'll use the entire training set.
- For each iteration, you'll calculate the cost function using all training examples (there are 'm' training examples), and for all features.
- Instead of updating a single weight  $\theta_i$  at a time, we can update all the weights in the column vector:

$$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{pmatrix}$$

- $\theta$  has dimensions  $(n+1, 1)$ , where 'n' is the number of features, and there is one more element for the bias term  $\theta_0$  (note that the corresponding feature value  $\mathbf{x}_0$  is 1).
- The 'logits', 'z', are calculated by multiplying the feature matrix 'x' with the weight vector 'theta'.  $z = \mathbf{x}\theta$ 
  - $\mathbf{x}$  has dimensions  $(m, n+1)$
  - $\theta$  has dimensions  $(n+1, 1)$
  - $\mathbf{z}$  has dimensions  $(m, 1)$
- The prediction 'h', is calculated by applying the sigmoid to each element in 'z':  $h(z) = \text{sigmoid}(z)$ , and has dimensions  $(m, 1)$ .

- The cost function  $J$  is calculated by taking the dot product of the vectors 'y' and 'log(h)'. Since both 'y' and 'h' are column vectors (m,1), transpose the vector to the left, so that matrix multiplication of a row vector with column vector performs the dot product.

$$J = \frac{-1}{m} \times \left( \mathbf{y}^T \cdot \log(\mathbf{h}) + (\mathbf{1} - \mathbf{y})^T \cdot \log(\mathbf{1} - \mathbf{h}) \right)$$

- The update of theta is also vectorized. Because the dimensions of  $\mathbf{x}$  are (m, n+1), and both  $\mathbf{h}$  and  $\mathbf{y}$  are (m, 1), we need to transpose the  $\mathbf{x}$  and place it on the left in order to perform matrix multiplication, which then yields the (n+1, 1) answer we need:

$$\theta = \theta - \frac{\alpha}{m} \times (\mathbf{x}^T \cdot (\mathbf{h} - \mathbf{y}))$$

Hints

use numpy.dot for matrix multiplication.

To ensure that the fraction -1/m is a decimal value, cast either the numerator or denominator (or both), like `float(1)`, or write `1.` for the float version of 1.

```
[15]: # UNQ_C2 GRADED FUNCTION: gradientDescent
def gradientDescent(x, y, theta, alpha, num_iters):
    '''
    Input:
        x: matrix of features which is (m,n+1)
        y: corresponding labels of the input matrix x, dimensions (m,1)
        theta: weight vector of dimension (n+1,1)
        alpha: learning rate
        num_iters: number of iterations you want to train your model for
    Output:
        J: the final cost
        theta: your final weight vector
    Hint: you might want to print the cost to make sure that it is going down.
    '''
    ### START CODE HERE ###
    # get 'm', the number of rows in matrix x
    m = x.shape[0]

    for i in range(0, num_iters):

        # get z, the dot product of x and theta
        z = np.dot(x, theta)

        # get the sigmoid of z
        h = sigmoid(z)

        # calculate the cost function
        J = -1./m * (np.dot(y.transpose(), np.log(h)) + np.dot((1 - y).
→transpose(), np.log(1 - h)))
```



```

    # update the weights theta
    theta = theta - (alpha / m) * (np.dot(x.transpose(), (h - y)))

    ### END CODE HERE ###
    J = float(J)
    return J, theta

```

```

[16]: # Check the function
      # Construct a synthetic test case using numpy PRNG functions
      np.random.seed(1)
      # X input is 10 x 3 with ones for the bias terms
      tmp_X = np.append(np.ones((10, 1)), np.random.rand(10, 2) * 2000, axis=1)
      # Y Labels are 10 x 1
      tmp_Y = (np.random.rand(10, 1) > 0.35).astype(float)

      # Apply gradient descent
      tmp_J, tmp_theta = gradientDescent(tmp_X, tmp_Y, np.zeros((3, 1)), 1e-8, 700)
      print(f"The cost after training is {tmp_J:.8f}.")
      print(f"The resulting vector of weights is {[round(t, 8) for t in np.
      ↳squeeze(tmp_theta)]}")

```

The cost after training is 0.67094970.

The resulting vector of weights is [4.1e-07, 0.00035658, 7.309e-05]

### Expected output

The cost after training is 0.67094970.

The resulting vector of weights is [4.1e-07, 0.00035658, 7.309e-05]

```

[17]: # Test your function
      w1_unittest.test_gradientDescent(gradientDescent)

```

All tests passed

## 2.2 Part 2: Extracting the features

- Given a list of tweets, extract the features and store them in a matrix. You will extract two features.
  - The first feature is the number of positive words in a tweet.
  - The second feature is the number of negative words in a tweet.
- Then train your logistic regression classifier on these features.
- Test the classifier on a validation set.

### 2.2.1 Instructions: Implement the `extract_features` function.

- This function takes in a single tweet.

- Process the tweet using the imported `process_tweet` function and save the list of tweet words.
- Loop through each word in the list of processed words
  - For each word, check the 'freqs' dictionary for the count when that word has a positive '1' label. (Check for the key (word, 1.0))
  - Do the same for the count for when the word is associated with the negative label '0'. (Check for the key (word, 0.0).)

#### Hints

Make sure you handle cases when the (word, label) key is not found in the dictionary.

Search the web for hints about using the 'get' function of a Python dictionary. Here is an example

```
[18]: # UNQ_C3 GRADED FUNCTION: extract_features
def extract_features(tweet, freqs, process_tweet=process_tweet):
    """
    Input:
        tweet: a list of words for one tweet
        freqs: a dictionary corresponding to the frequencies of each tuple_
    ↪ (word, label)
    Output:
        x: a feature vector of dimension (1,3)
    """
    # process_tweet tokenizes, stems, and removes stopwords
    word_l = process_tweet(tweet)

    # 3 elements in the form of a 1 x 3 vector
    x = np.zeros((1, 3))

    # bias term is set to 1
    x[0,0] = 1

    ### START CODE HERE ###

    # loop through each word in the list of words
    for word in word_l:

        # increment the word count for the positive label 1
        x[0,1] += freqs.get((word, 1.0), 0)

        # increment the word count for the negative label 0
        x[0,2] += freqs.get((word, 0.0), 0)

    ### END CODE HERE ###
    assert(x.shape == (1, 3))
    return x
```

```
[19]: # Check your function
# test 1
# test on training data
tmp1 = extract_features(train_x[0], freqs)
print(tmp1)
```

```
[[1.000e+00 3.133e+03 6.100e+01]]
```

### Expected output

```
[[1.000e+00 3.133e+03 6.100e+01]]
```

```
[20]: # test 2:
# check for when the words are not in the freqs dictionary
tmp2 = extract_features('blorb bleeeeb bloooob', freqs)
print(tmp2)
```

```
[[1. 0. 0.]]
```

### Expected output

```
[[1. 0. 0.]]
```

```
[21]: # Test your function
w1_unittest.test_extract_features(extract_features, freqs)
```

```
All tests passed
```

## 2.3 Part 3: Training Your Model

To train the model: \* Stack the features for all training examples into a matrix X. \* Call `gradientDescent`, which you've implemented above.

This section is given to you. Please read it for understanding and run the cell.

```
[22]: # collect the features 'x' and stack them into a matrix 'X'
X = np.zeros((len(train_x), 3))
for i in range(len(train_x)):
    X[i, :] = extract_features(train_x[i], freqs)

# training labels corresponding to X
Y = train_y

# Apply gradient descent
J, theta = gradientDescent(X, Y, np.zeros((3, 1)), 1e-9, 1500)
print(f"The cost after training is {J:.8f}.")
print(f"The resulting vector of weights is {[round(t, 8) for t in np.
    ↳squeeze(theta)]}")
```

The cost after training is 0.22522315.

The resulting vector of weights is [6e-08, 0.00053818, -0.0005583]

#### Expected Output:

The cost after training is 0.22522315.

The resulting vector of weights is [6e-08, 0.00053818, -0.0005583]

### 3 Part 4: Test your logistic regression

It is time for you to test your logistic regression function on some new input that your model has not seen before.

**Instructions:** Write `predict_tweet` Predict whether a tweet is positive or negative.

- Given a tweet, process it, then extract the features.
- Apply the model's learned weights on the features to get the logits.
- Apply the sigmoid to the logits to get the prediction (a value between 0 and 1).

$$y_{pred} = \text{sigmoid}(\mathbf{x} \cdot \theta)$$

```
[23]: # UNQ_C4 GRADED FUNCTION: predict_tweet
def predict_tweet(tweet, freqs, theta):
    '''
    Input:
        tweet: a string
        freqs: a dictionary corresponding to the frequencies of each tuple_
        ↪ (word, label)
        theta: (3,1) vector of weights
    Output:
        y_pred: the probability of a tweet being positive or negative
    '''
    ### START CODE HERE ###

    # extract the features of the tweet and store it into x
    x = extract_features(tweet, freqs)

    # make the prediction using x and theta
    y_pred = sigmoid(np.dot(x, theta))

    ### END CODE HERE ###

    return y_pred
```

```
[24]: # Run this cell to test your function
```

```
for tweet in ['I am happy', 'I am bad', 'this movie should have been great.',
             'great', 'great great', 'great great great', 'great great great great']:
    print( '%s -> %f' % (tweet, predict_tweet(tweet, freqs, theta)))
```

```
I am happy -> 0.519275
I am bad -> 0.494347
this movie should have been great. -> 0.515979
great -> 0.516065
great great -> 0.532096
great great great -> 0.548062
great great great great -> 0.563929
```

#### Expected Output:

```
I am happy -> 0.519275
I am bad -> 0.494347
this movie should have been great. -> 0.515979
great -> 0.516065
great great -> 0.532096
great great great -> 0.548062
great great great great -> 0.563929
```

```
[25]: # Feel free to check the sentiment of your own tweet below
my_tweet = 'I am learning :)'
predict_tweet(my_tweet, freqs, theta)
```

```
[25]: array([[0.83110307]])
```

```
[26]: # Test your function
w1_unittest.test_predict_tweet(predict_tweet, freqs, theta)
```

```
All tests passed
```

### 3.1 Check performance using the test set

After training your model using the training set above, check how your model might perform on real, unseen data, by testing it against the test set.

#### Instructions: Implement `test_logistic_regression`

- Given the test data and the weights of your trained model, calculate the accuracy of your logistic regression model.
- Use your ‘predict\_tweet’ function to make predictions on each tweet in the test set.
- If the prediction is  $> 0.5$ , set the model’s classification ‘y\_hat’ to 1, otherwise set the model’s classification ‘y\_hat’ to 0.
- A prediction is accurate when the y\_hat equals the test\_y. Sum up all the instances when they are equal and divide by m.

Hints

Use `np.asarray()` to convert a list to a numpy array

Use `numpy.squeeze()` to make an `(m,1)` dimensional array into an `(m,)` array

```
[27]: # UNQ_C5 GRADED FUNCTION: test_logistic_regression
def test_logistic_regression(test_x, test_y, freqs, theta,
    predict_tweet=predict_tweet):
    """
    Input:
        test_x: a list of tweets
        test_y: (m, 1) vector with the corresponding labels for the list of
    tweets
        freqs: a dictionary with the frequency of each pair (or tuple)
        theta: weight vector of dimension (3, 1)
    Output:
        accuracy: (# of tweets classified correctly) / (total # of tweets)
    """

    ### START CODE HERE ###

    # the list for storing predictions
    y_hat = []

    for tweet in test_x:
        # get the label prediction for the tweet
        y_pred = predict_tweet(tweet, freqs, theta)

        if y_pred > 0.5:
            # append 1.0 to the list
            y_hat.append(1.0)
        else:
            # append 0 to the list
            y_hat.append(0.0)

    # With the above implementation, y_hat is a list, but test_y is (m,1) array
    # convert both to one-dimensional arrays in order to compare them using the
    '==' operator
    accuracy = (y_hat==np.squeeze(test_y)).sum()/len(test_x)

    ### END CODE HERE ###

    return accuracy
```

```
[28]: tmp_accuracy = test_logistic_regression(test_x, test_y, freqs, theta)
print(f"Logistic regression model's accuracy = {tmp_accuracy:.4f}")
```

Logistic regression model's accuracy = 0.9950

**Expected Output:** 0.9950

Pretty good!

```
[29]: # Test your function
w1_unittest.unittest_test_logistic_regression(test_logistic_regression, freqs,
→theta)
```

All tests passed

## 4 Part 5: Error Analysis

In this part you will see some tweets that your model misclassified. Why do you think the misclassifications happened? Specifically what kind of tweets does your model misclassify?

```
[30]: # Some error analysis done for you
print('Label Predicted Tweet')
for x,y in zip(test_x,test_y):
    y_hat = predict_tweet(x, freqs, theta)

    if np.abs(y - (y_hat > 0.5)) > 0:
        print('THE TWEET IS:', x)
        print('THE PROCESSED TWEET IS:', process_tweet(x))
        print('%d\t0.8f\t%s' % (y, y_hat, ' '.join(process_tweet(x)).
→encode('ascii', 'ignore')))
```

Label Predicted Tweet

THE TWEET IS: @MarkBreech Not sure it would be good thing 4 my bottom daring 2 say 2 Miss B but Im gonna be so stubborn on mouth soaping ! #NotHavingit :p

THE PROCESSED TWEET IS: ['sure', 'would', 'good', 'thing', '4', 'bottom', 'dare', '2', 'say', '2', 'miss', 'b', 'im', 'gonna', 'stubborn', 'mouth', 'soap', 'nothavingit', ':p']

1 0.48901497 b'sure would good thing 4 bottom dare 2 say 2 miss b im gonna stubborn mouth soap nothavingit :p'

THE TWEET IS: I'm playing Brain Dots : ) #BrainDots

http://t.co/UGQz0x0huu

THE PROCESSED TWEET IS: ["i'm", 'play', 'brain', 'dot', 'braindot']

1 0.48418949 b"i'm play brain dot braindot"

THE TWEET IS: I'm playing Brain Dots : ) #BrainDots http://t.co/aOKldo3GMj

http://t.co/xWCM9qyRG5

THE PROCESSED TWEET IS: ["i'm", 'play', 'brain', 'dot', 'braindot']

1 0.48418949 b"i'm play brain dot braindot"

THE TWEET IS: I'm playing Brain Dots : ) #BrainDots http://t.co/R2JB08iNww

http://t.co/ow5BBwdEMY

THE PROCESSED TWEET IS: ["i'm", 'play', 'brain', 'dot', 'braindot']

1 0.48418949 b"i'm play brain dot braindot"

THE TWEET IS: off to the park to get some sunlight : )

THE PROCESSED TWEET IS: ['park', 'get', 'sunlight']

```

1      0.49636374      b'park get sunlight'
THE TWEET IS: @msarosh Uff Itna Miss karhy thy ap :p
THE PROCESSED TWEET IS: ['uff', 'itna', 'miss', 'karhi', 'thi', 'ap', ':p']
1      0.48237069      b'uff itna miss karhi thi ap :p'
THE TWEET IS: @phenomyoutube u probs had more fun with david than me : (
THE PROCESSED TWEET IS: ['u', 'prob', 'fun', 'david']
0      0.50988239      b'u prob fun david'
THE TWEET IS: pats jay : (
THE PROCESSED TWEET IS: ['pat', 'jay']
0      0.50040365      b'pat jay'
THE TWEET IS: my beloved grandmother : ( https://t.co/wt4oXq5xCf
THE PROCESSED TWEET IS: ['belov', 'grandmoth']
0      0.50000002      b'belov grandmoth'
THE TWEET IS: Sr. Financial Analyst - Expedia, Inc.: (#Bellevue, WA)
http://t.co/ktnMhvwCI #Finance #ExpediaJobs #Job #Jobs #Hiring
THE PROCESSED TWEET IS: ['sr', 'financi', 'analyst', 'expedia', 'inc',
'bellevu', 'wa', 'financ', 'expediajob', 'job', 'job', 'hire']
0      0.50648681      b'sr financi analyst expedia inc bellevu wa financ
expediajob job job hire'

```

Later in this specialization, we will see how we can use deeplearning to improve the prediction performance.

## 5 Part 6: Predict with your own tweet

```

[31]: # Feel free to change the tweet below
my_tweet = 'This is a ridiculously bright movie. The plot was terrible and I_
↳was sad until the ending!'
print(process_tweet(my_tweet))
y_hat = predict_tweet(my_tweet, freqs, theta)
print(y_hat)
if y_hat > 0.5:
    print('Positive sentiment')
else:
    print('Negative sentiment')

```

```

['ridicul', 'bright', 'movi', 'plot', 'terribl', 'sad', 'end']
[[0.48125423]]
Negative sentiment

```



# utils

January 23, 2022

```
[ ]: import re
import string
import numpy as np
```

```
[ ]: from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.tokenize import TweetTokenizer
```

```
[ ]: def process_tweet(tweet):
    """Process tweet function.
    Input:
        tweet: a string containing a tweet
    Output:
        tweets_clean: a list of words containing the processed tweet

    """
    stemmer = PorterStemmer()
    stopwords_english = stopwords.words('english')
    # remove stock market tickers like $GE
    tweet = re.sub(r'\$\w*', '', tweet)
    # remove old style retweet text "RT"
    tweet = re.sub(r'^RT[\s]+', '', tweet)
    # remove hyperlinks
    tweet = re.sub(r'https?:\/\/[^\s\n\r]+', '', tweet)
    # remove hashtags
    # only removing the hash # sign from the word
    tweet = re.sub(r'#', '', tweet)
    # tokenize tweets
    tokenizer = TweetTokenizer(preserve_case=False, strip_handles=True,
                               reduce_len=True)
    tweet_tokens = tokenizer.tokenize(tweet)

    tweets_clean = []
    for word in tweet_tokens:
        if (word not in stopwords_english and # remove stopwords
            word not in string.punctuation): # remove punctuation
            # tweets_clean.append(word)
```

```

        stem_word = stemmer.stem(word) # stemming word
        tweets_clean.append(stem_word)

    return tweets_clean

```

```

[ ]: def build_freqs(tweets, ys):
    """Build frequencies.
    Input:
        tweets: a list of tweets
        ys: an m x 1 array with the sentiment label of each tweet
            (either 0 or 1)
    Output:
        freqs: a dictionary mapping each (word, sentiment) pair to its
            frequency
    """
    # Convert np array to list since zip needs an iterable.
    # The squeeze is necessary or the list ends up with one element.
    # Also note that this is just a NOP if ys is already a list.
    yslist = np.squeeze(ys).tolist()

    # Start with an empty dictionary and populate it by looping over all tweets
    # and over all processed words in each tweet.
    freqs = {}
    for y, tweet in zip(yslist, tweets):
        for word in process_tweet(tweet):
            pair = (word, y)
            if pair in freqs:
                freqs[pair] += 1
            else:
                freqs[pair] = 1

    return freqs

```

## w1\_unittest

January 23, 2022

```
[ ]: # import nltk
import numpy as np

[ ]: def test_sigmoid(target):
    successful_cases = 0
    failed_cases = []

    test_cases = [
        {"name": "default_check", "input": {"z": 0}, "expected": 0.5},
        {
            "name": "positive_check",
            "input": {"z": 4.92},
            "expected": 0.9927537604041685,
        },
        {"name": "negative_check", "input": {"z": -1}, "expected": 0.
↪2689414213699951},
        {
            "name": "larger_neg_check",
            "input": {"z": -20},
            "expected": 2.0611536181902037e-09,
        },
    ]

    for test_case in test_cases:
        result = target(**test_case["input"])

        try:
            assert np.isclose(result, test_case["expected"])
            successful_cases += 1
        except:
            failed_cases.append(
                {
                    "name": test_case["name"],
                    "expected": test_case["expected"],
                    "got": result,
                }
            )
```

```

        print(
            f"Wrong output from sigmoid function. \n\tExpected:␣
→{failed_cases[-1].get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

    if len(failed_cases) == 0:
        print("\033[92m All tests passed")
    else:
        print("\033[92m", successful_cases, " Tests passed")
        print("\033[91m", len(failed_cases), " Tests failed")

    # return failed_cases, len(failed_cases) + successful_cases

```

```

[ ]: def test_gradientDescent(target):
    successful_cases = 0
    failed_cases = []

    test_cases = [
        {
            "name": "default_check",
            "input": {
                "random_seed": 1,
                "input_dict": {
                    "x": np.array(
                        [
                            [1.00000000e00, 8.34044009e02, 1.44064899e03],
                            [1.00000000e00, 2.28749635e-01, 6.04665145e02],
                            [1.00000000e00, 2.93511782e02, 1.84677190e02],
                            [1.00000000e00, 3.72520423e02, 6.91121454e02],
                            [1.00000000e00, 7.93534948e02, 1.07763347e03],
                            [1.00000000e00, 8.38389029e02, 1.37043900e03],
                            [1.00000000e00, 4.08904499e02, 1.75623487e03],
                            [1.00000000e00, 5.47751864e01, 1.34093502e03],
                            [1.00000000e00, 8.34609605e02, 1.11737966e03],
                            [1.00000000e00, 2.80773877e02, 3.96202978e02],
                        ]
                    ),
                    "y": np.array(
                        [
                            [1.0],
                            [1.0],
                            [0.0],
                            [1.0],
                            [1.0],
                            [1.0],
                            [0.0],
                            [0.0],
                        ]
                    )
                }
            }
        }
    ]

```

```

        [0.0],
        [1.0],
    ]
),
    "theta": np.zeros((3, 1)),
    "alpha": 1e-8,
    "num_iters": 700,
},
},
"expected": {
    "J": 0.6709497038162118,
    "theta": np.array(
        [[4.10713435e-07], [3.56584699e-04], [7.30888526e-05]]
    ),
},
},
{
    "name": "larger_check",
    "input": {
        "random_seed": 2,
        "input_dict": {
            "x": np.array(
                [
                    [1.0, 435.99490214, 25.92623183, 549.66247788],
                    [1.0, 435.32239262, 420.36780209, 330.334821],
                    [1.0, 204.64863404, 619.27096635, 299.65467367],
                    [1.0, 266.8272751, 621.13383277, 529.14209428],
                    [1.0, 134.57994534, 513.57812127, 184.43986565],
                    [1.0, 785.33514782, 853.97529264, 494.23683738],
                    [1.0, 846.56148536, 79.64547701, 505.24609012],
                    [1.0, 65.28650439, 428.1223276, 96.53091566],
                    [1.0, 127.1599717, 596.74530898, 226.0120006],
                    [1.0, 106.94568431, 220.30620707, 349.826285],
                    [1.0, 467.78748458, 201.74322626, 640.40672521],
                    [1.0, 483.06983555, 505.23672002, 386.89265112],
                    [1.0, 793.63745444, 580.00417888, 162.2985985],
                    [1.0, 700.75234661, 964.55108009, 500.00836117],
                    [1.0, 889.52006395, 341.61365267, 567.14412763],
                    [1.0, 427.5459633, 436.74726303, 776.559185],
                    [1.0, 535.6041735, 953.74222694, 544.20816015],
                    [1.0, 82.09492228, 366.34240168, 850.850504],
                    [1.0, 406.27504305, 27.20236589, 247.177239],
                    [1.0, 67.14437074, 993.85201142, 970.58031338],
                ]
            ),
            "y": np.array(
                [

```

```

        [1.0],
        [1.0],
        [1.0],
        [0.0],
        [0.0],
        [1.0],
        [0.0],
        [0.0],
        [1.0],
        [0.0],
        [1.0],
        [0.0],
        [0.0],
        [0.0],
        [1.0],
        [1.0],
        [0.0],
        [0.0],
        [1.0],
        [0.0],
    ]
),
    "theta": np.zeros((4, 1)),
    "alpha": 1e-4,
    "num_iters": 30,
},
},
"expected": {
    "J": 6.5044107216556135,
    "theta": np.array(
        [
            [9.45211976e-05],
            [2.40577958e-02],
            [-1.77876847e-02],
            [1.35674845e-02],
        ]
    ),
},
},
]

for test_case in test_cases:
    # Setting the random seed for reproducibility
    result_J, result_theta = target(**test_case["input"]["input_dict"])

    try:
        assert isinstance(result_J, float)

```

```

        successful_cases += 1
    except:
        failed_cases.append(
            {"name": test_case["name"], "expected": float, "got":
→type(result_J),}
        )
        print(
            f"Wrong output type for loss function. \n\tExpected:
→{failed_cases[-1].get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

    try:
        assert np.isclose(result_J, test_case["expected"]["J"])
        successful_cases += 1
    except:
        failed_cases.append(
            {
                "name": test_case["name"],
                "expected": test_case["expected"]["J"],
                "got": result_J,
            }
        )
        print(
            f"Wrong output for the loss function. Check how you are
→implementing the matrix multiplications. \n\tExpected: {failed_cases[-1].
→get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

    try:
        assert result_theta.shape ==
→test_case["input"]["input_dict"]["theta"].shape
        successful_cases += 1
    except:
        failed_cases.append(
            {
                "name": test_case["name"],
                "expected": test_case["input"]["input_dict"]["theta"].shape,
                "got": result_theta.shape,
            }
        )
        print(
            f"Wrong shape for weights matrix theta. \n\tExpected:
→{failed_cases[-1].get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

    try:
        assert np.allclose(

```

```

        np.squeeze(result_theta), np.
→squeeze(test_case["expected"]["theta"]),
    )
    successful_cases += 1
except:
    failed_cases.append(
        {
            "name": test_case["name"],
            "expected": test_case["expected"]["theta"],
            "got": result_theta,
        }
    )
    print(
        f"Wrong values for weight's matrix theta. Check how you are_
→updating the matrix of weights. \n\tExpected: {failed_cases[-1].
→get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
    )

if len(failed_cases) == 0:
    print("\033[92m All tests passed")
else:
    print("\033[92m", successful_cases, " Tests passed")
    print("\033[91m", len(failed_cases), " Tests failed")

```

[ ]:

```

[ ]: def test_extract_features(target, freqs):
    successful_cases = 0
    failed_cases = []

    test_cases = [
        {
            "name": "default_check",
            "input": {
                "tweet": "#FollowFriday @France_Inte @PKuchly57 @Milipol_Paris_
→for being top engaged members in my community this week :)",
                "freqs": freqs,
            },
            "expected": np.array(
                [[1.00e00, 3.133e03, 6.10e01]]
            ),
        },
        {
            "name": "unk_words_check",
            "input": {"tweet": "blorb bleeeeb bloooob", "freqs": freqs},
            "expected": np.array([[1.0, 0.0, 0.0]]),
        },
    ]

```



```

    {
        "name": "good_words_check",
        "input": {"tweet": "Hello world! All's good!", "freqs": freqs},
        "expected": np.array([[1.0, 263.0, 106.0]]),
    },
    {
        "name": "bad_words_check",
        "input": {"tweet": "It is so sad!", "freqs": freqs},
        "expected": np.array([[1.0, 5.0, 100.0]]),
    },
]

for test_case in test_cases:
    result = target(**test_case["input"])

    try:
        assert result.shape == test_case["expected"].shape
        successful_cases += 1
    except:
        failed_cases.append(
            {
                "name": test_case["name"],
                "expected": test_case["expected"].shape,
                "got": result.shape,
            }
        )
        print(
            f"Wrong output shape. \n\tExpected: {failed_cases[-1].get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

    try:
        assert np.allclose(result, test_case["expected"])
        successful_cases += 1
    except:
        failed_cases.append(
            {
                "name": test_case["name"],
                "expected": test_case["expected"],
                "got": result,
            }
        )
        print(
            f"Wrong output values. Check how you are computing the positive or negative word count. \n\tExpected: {failed_cases[-1].get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

```

```

if len(failed_cases) == 0:
    print("\033[92m All tests passed")
else:
    print("\033[92m", successful_cases, " Tests passed")
    print("\033[91m", len(failed_cases), " Tests failed")

```

```

[ ]: def test_predict_tweet(target, freqs, theta):
    successful_cases = 0
    failed_cases = []

    test_cases = [
        {
            "name": "default_check1",
            "input": {"tweet": "I am happy", "freqs": freqs, "theta": theta},
            "expected": np.array([[0.5192746]]),
        },
        {
            "name": "default_check2",
            "input": {"tweet": "I am bad", "freqs": freqs, "theta": theta},
            "expected": np.array([[0.49434685]]),
        },
        {
            "name": "default_check3",
            "input": {
                "tweet": "this movie should have been great",
                "freqs": freqs,
                "theta": theta,
            },
            "expected": np.array([[0.5159792]]),
        },
        {
            "name": "default_check5",
            "input": {"tweet": "It is a good day", "freqs": freqs, "theta":
→theta},
            "expected": np.array([[0.52320595]]),
        },
        {
            "name": "default_check6",
            "input": {"tweet": "It is a bad bad day", "freqs": freqs, "theta":
→theta},
            "expected": np.array([[0.49780224]]),
        },
        {
            "name": "default_check7",
            "input": {
                "tweet": "It is a good day",

```

```

        "freqs": freqs,
        "theta": np.array([[5.0000e-04], [-3.4e-02], [3.2e-02]]),
    },
    "expected": np.array([[0.00147813]]),
},
{
    "name": "default_check8",
    "input": {
        "tweet": "It is a bad bad day",
        "freqs": freqs,
        "theta": np.array([[5.0000e-04], [-3.4e-02], [3.2e-02]]),
    },
    "expected": np.array([[0.45673348]]),
},
{
    "name": "default_check9",
    "input": {
        "tweet": "this movie should have been great",
        "freqs": freqs,
        "theta": np.array([[5.0000e-04], [-3.4e-02], [3.2e-02]]),
    },
    "expected": np.array([[0.01561938]]),
},
]

for test_case in test_cases:
    result = target(**test_case["input"])

    try:
        assert result.shape == test_case["expected"].shape
        successful_cases += 1
    except:
        failed_cases.append(
            {
                "name": test_case["name"],
                "expected": test_case["expected"].shape,
                "got": result.shape,
            }
        )
        print(
            f"Wrong output shape. \n\tExpected: {failed_cases[-1].get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

    try:
        assert np.allclose(result, test_case["expected"])
        successful_cases += 1

```

```

except:
    failed_cases.append(
        {
            "name": test_case["name"],
            "expected": test_case["expected"],
            "got": result,
        }
    )
    print(
        f"Wrong predicted values. \n\tExpected: {failed_cases[-1].
        ↳get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
    )

if len(failed_cases) == 0:
    print("\033[92m All tests passed")
else:
    print("\033[92m", successful_cases, " Tests passed")
    print("\033[91m", len(failed_cases), " Tests failed")

```

[ ]:

```

[ ]: def unittest_test_logistic_regression(target, freqs, theta):
    successful_cases = 0
    failed_cases = []

    test_cases = [
        {
            "name": "default_check1",
            "input": {
                "test_x": [
                    "Bro:U wan cut hair anot,ur hair long Liao bo\nMe:since ord_
                    ↳liao,take it easy lor treat as save $ leave it longer :) \nBro:LOL Sibe_
                    ↳xialan",
                    "@heyclaireee is back! thnx God!!! i'm so happy :-)",
                    "@BBCRadio3 thought it was my ears which were_
                    ↳malfunctioning, thank goodness you cleared that one up with an apology :-)",
                    "@HumayAG 'Stuck in the centre right with you. Clowns to_
                    ↳the right, jokers to the left...' :) @orgasticpotency @ahmedshaheed_
                    ↳@AhmedSaeedGahaa",
                    "Happy Friday :-) http://t.co/iymP1lWXYF",
                    "I wanna change my avi but uSanele :(",
                    "MY PUPPY BROKE HER FOOT :(",
                    "where's all the jaebum baby pictures :(",
                    "But but Mr Ahmad Maslan cooks too :( https://t.co/
                    ↳ArCiD31Zv6",

```

```

        "@eawoman As a Hull supporter I am expecting a misserable_
↪few weeks :-(",
    ],
    "test_y": np.array(
        [
            [1.0],
            [1.0],
            [1.0],
            [1.0],
            [1.0],
            [0.0],
            [0.0],
            [0.0],
            [0.0],
            [0.0],
            [0.0],
        ]
    ),
    "freqs": freqs,
    "theta": theta,
},
"expected": 1.0,
},
{
    "name": "default_check1",
    "input": {
        "test_x": [
            "Bro:U wan cut hair anot,ur hair long Liao bo\nMe:since ord_
↪liao,take it easy lor treat as save $ leave it longer :)\nBro:LOL Sibe_
↪xialan",
            "@heyclaireeee is back! thnx God!!! i'm so happy :)",
            "@BBCRadio3 thought it was my ears which were_
↪malfunctioning, thank goodness you cleared that one up with an apology :-)",
            "@HumayAG 'Stuck in the centre right with you. Clowns to_
↪the right, jokers to the left...' :) @orgasticpotency @ahmedshaheed_
↪@AhmedSaeedGahaa",
            "Happy Friday :-) http://t.co/iymPILWXYF",
            "I wanna change my avi but uSanele :(",
            "MY PUPPY BROKE HER FOOT :(",
            "where's all the jaebum baby pictures :(",
            "But but Mr Ahmad Maslan cooks too :( https://t.co/
↪ArCiD31Zv6",
            "@eawoman As a Hull supporter I am expecting a misserable_
↪few weeks :-(",
        ],
        "test_y": np.array(
            [

```

```

        [1.0],
        [1.0],
        [1.0],
        [1.0],
        [1.0],
        [0.0],
        [0.0],
        [0.0],
        [0.0],
        [0.0],
    ]
),
    "freqs": freqs,
    "theta": np.array([[5.0000e-04], [-3.4e-02], [3.2e-02]]),
},
    "expected": 0.0,
},
]

for test_case in test_cases:
    result = target(**test_case["input"])

    try:
        assert isinstance(result, np.float64)
        successful_cases += 1
    except:
        failed_cases.append(
            {
                "name": test_case["name"],
                "expected": np.float64,
                "got": type(result),
            }
        )
        print(
            f"Wrong output type. \n\tExpected: {failed_cases[-1].get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
        )

    try:
        assert np.isclose(result, test_case["expected"])
        successful_cases += 1
    except:
        failed_cases.append(
            {
                "name": test_case["name"],
                "expected": test_case["expected"],
                "got": result,
            }
        )

```

```

        }
    )
    print(
        f"Wrong accuracy value. \n\tExpected: {failed_cases[-1].
→get('expected')}. \n\tGot: {failed_cases[-1].get('got')}."
    )

if len(failed_cases) == 0:
    print("\033[92m All tests passed")
else:
    print("\033[92m", successful_cases, " Tests passed")
    print("\033[91m", len(failed_cases), " Tests failed")

```