# [COM6513] Assignment 1: Sentiment Analysis with Logistic Regression

## Instructor: Nikos Aletras

The goal of this assignment is to develop and test a **text classification** system for **sentiment analysis**, in particular to predict the sentiment of movie reviews, i.e. positive or negative (binary classification).

For that purpose, you will implement:

- Text processing methods for extracting Bag-Of-Word features, using
  - n-grams (BOW), i.e. unigrams, bigrams and trigrams to obtain vector representations of documents where n=1,2,3 respectively. Two vector weighting schemes should be tested: (1) raw frequencies (**1 mark**); (2) tf.idf (**1 mark**).
  - character n-grams (BOCN). A character n-gram is a contiguous sequence of characters given a word, e.g. for n=2, 'coffee' is split into {'co', 'of', 'ff', 'fe', 'ee'}. Two vector weighting schemes should be tested: (1) raw frequencies (**1 mark**); (2) tf.idf (**1 mark**). **Tip: Note the large vocabulary size!**
  - a combination of the two vector spaces (n-grams and character n-grams) choosing your best performing wighting respectively (i.e. raw or tfidf). (**1 mark**) **Tip: you should merge the two representations**

- Binary Logistic Regression (LR) classifiers that will be able to accurately classify movie reviews trained with:
  - (1) BOW-count (raw frequencies)
  - (2) BOW-tfidf (tf.idf weighted)
  - (3) BOCN-count
  - (4) BOCN-tfidf
  - (5) BOW+BOCN (best performing weighting; raw or tfidf)

- The Stochastic Gradient Descent (SGD) algorithm to estimate the parameters of your Logistic Regression models. Your SGD algorithm should:
  - Minimise the Binary Cross-entropy loss function (**1 mark**)
  - Use L2 regularisation (**1 mark**)
  - Perform multiple passes (epochs) over the training data (**1 mark**)
  - Randomise the order of training data after each pass (**1 mark**)
  - Stop training if the difference between the current and previous development loss is smaller than a threshold (**1 mark**)
  - After each epoch print the training and development loss (**1 mark**)

- Discuss how did you choose hyperparameters (e.g. learning rate and regularisation strength) for each LR model? You should use a table showing model performance using different set of hyperparameter values. (**2 marks**). **\*\*Tip: Instead of using all possible combinations, you could perform a random sampling of combinations.**

- After training each LR model, plot the learning process (i.e. training and validation loss in each epoch) using a line plot. Does your model underfit, overfit or is it about right? Explain why. (**1 mark**).

- Identify and show the most important features (model interpretability) for each class (i.e. top-10 most positive and top-10 negative weights). Give the top 10 for each class and comment on whether they make sense (if they don't you might have a bug!). If you were to apply the classifier into a different domain such laptop reviews or restaurant reviews, do you think these features would generalise well? Can you propose what features the classifier could pick up as important in the new domain? (**2 marks**)

- Provide well documented and commented code describing all of your choices. In general, you are free to make decisions about text processing (e.g. punctuation, numbers, vocabulary size) and hyperparameter values. We expect to see justifications and discussion for all of your choices (**2 marks**).

- Provide efficient solutions by using Numpy arrays when possible (you can find tips in Lab 1 sheet). Executing the whole notebook with your code should not take more than 5 minutes on a any standard computer (e.g. Intel Core i5 CPU, 8 or 16GB RAM) excluding hyperparameter tuning runs (**2 marks**).

## Data

The data you will use are taken from here: [http://www.cs.cornell.edu/people/pabo/movie-review-data/ (http://www.cs.cornell.edu/people/pabo/movie-review-data/)](http://www.cs.cornell.edu/people/pabo/movie-review-data/) and you can find it in the `./data_sentiment` folder in CSV format:

- `data_sentiment/train.csv` : contains 1,400 reviews, 700 positive (label: 1) and 700 negative (label: 0) to be used for training.
- `data_sentiment/dev.csv` : contains 200 reviews, 100 positive and 100 negative to be used for hyperparameter selection and monitoring the training process.
- `data_sentiment/test.csv` : contains 400 reviews, 200 positive and 200 negative to be used for testing.

## Submission Instructions

You should submit a Jupyter Notebook file (assignment1.ipynb) and an exported PDF version (you can do it from Jupyter: `File->Download as->PDF via Latex` or you can print it as PDF using your browser).

You are advised to follow the code structure given in this notebook by completing all given funtions. You can also write any auxilliary/helper functions (and arguments for the functions) that you might need but note that you can provide a full solution without any such functions. Similarly, you can just use only the packages imported below but you are free to use any functionality from the [Python Standard Library (https://docs.python.org/2/library/index.html)](https://docs.python.org/2/library/index.html), NumPy, SciPy (excluding built-in softmax funtcions) and Pandas. You are not allowed to use any third-party library such as Scikit-learn (apart from metric functions already provided), NLTK, Spacy, Keras etc..

There is no single correct answer on what your accuracy should be, but correct implementations usually achieve F1-scores around 80% or higher. The quality of the analysis of the results is as important as the accuracy itself.

This assignment will be marked out of 20. It is worth 20% of your final grade in the module.

The deadline for this assignment is **23:59 on Mon, 14 Mar 2022** and it needs to be submitted

via Blackboard. Standard departmental penalties for lateness will be applied. We use a range of strategies to **detect unfair means (https://www.sheffield.ac.uk/ssid/unfair-means/index)**, including Turnitin which helps detect plagiarism. Use of unfair means would

```
In [1]:  1  import pandas as pd
         2  import numpy as np
         3  from collections import Counter
         4  import re
         5  import matplotlib.pyplot as plt
         6  from sklearn.metrics import accuracy_score, precision_score, recall
         7  import random
         8  import string
         9
        10  # fixing random seed for reproducibility
        11  random.seed(123)
        12  np.random.seed(123)
        13
```

## Load Raw texts and labels into arrays

First, you need to load the training, development and test sets from their corresponding CSV files (tip: you can use Pandas dataframes).

```
In [2]:  1  # read csv file with the test sets
         2  test = pd.read_csv("data_sentiment/test.csv")
         3
         4
         5  # displaying the list of column names
         6  #Column 1 = TEXT
         7  #Column 2 = LABELS
         8
         9  # creating a list of column names by
        10  # calling the columns
```

If you use Pandas you can see a sample of the data.

```
In [3]:  1  # read csv file with the train sets
         2  train = pd.read_csv("data_sentiment/train.csv")
         3
         4
         5  # displaying the list of column names
         6  #Column 0 = TEXT
         7  #Column 1 = LABELS
         8
         9
        10  # creating a list of column names by
        11  # calling the .columns
        12  train_column_names = list(train.columns)
        13
        14
        15  ##########################################################
        16
        17
        18  # read csv file with the development sets
        19  dev  = pd.read_csv("data_sentiment/dev.csv")
        20
```

```
21
22  # displaying the list of column names
23  #Column 0 = TEXT
24  #Column 1 = LABELS
25
26
27  # creating a list of column names by
28  # calling the .columns
```

The next step is to put the raw texts into Python lists and their corresponding labels into NumPy arrays:

```
 1  #put the trainning raw texts into Python lists
 2  train_text = list(train[train_column_names[0]])
 3
 4  #print the text for verification
 5  #print(train_text,"\n")
 6
 7  #put the trainning labels into a NumPy arrays
 8  train_label = train[train_column_names[1]].values
 9
10  #print the train label for verification
11  print(train_label,"\n")
12
13
14  #############################################
15
16
17  #put the testing raw texts into Python lists
18  test_text = list(test[test_column_names[0]])
19
20  #print the text for verification
21  #print(test_text,"\n")
22
23  #put the testing labels into a NumPy arrays
24  test_label = test[test_column_names[1]].values
25
26  #print the test label for verification
27  print(test_label,"\n")
28
29
30  #############################################
31
32
33  #put the development raw texts into Python lists
34  dev_text = list(dev[dev_column_names[0]])
35
36  #print the text for verification
37  #print(dev_text,"\n")
38
39  #put the development labels into a NumPy arrays
40  dev_label = dev[dev_column_names[1]].values
41
42  #print the dev label for verification
```

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

[1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1
 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0]

[1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0
 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0
```

# Vector Representations of Text

To train and test Logisitc Regression models, you first need to obtain vector representations for all documents given a vocabulary of features (unigrams, bigrams, trigrams).

## Text Pre-Processing Pipeline

To obtain a vocabulary of features, you should:

- tokenise all texts into a list of unigrams (tip: using a regular expression)
- remove stop words (using the one provided or one of your preference)
- compute bigrams, trigrams given the remaining unigrams (or character ngrams from the unigrams)
- remove ngrams appearing in less than K documents
- use the remaining to create a vocabulary of unigrams, bigrams and trigrams (or character n-grams). You can keep top N if you encounter memory issues.

```
In [5]:    1  #make a list for unigrams from each
           2  def generate_ngrams(s, n):
           3
           4      # Convert to lowercases
           5      #s = s.lower()
           6
```

```python
 7      # Replace all none alphanumeric characters with spaces
 8      s = re.sub(r'[^a-zA-Z0-9\s]', ' ', str(s))
 9
10      # Break sentence in the token, remove empty tokens
11      tokens = [token for token in s.split(" ") if token != ""]
12
13      # Use the zip function to help generate n-grams
14      # Concatentate the tokens into ngrams and return
15      ngrams = zip(*[tokens[i:] for i in range(n)])
16      return [" ".join(ngram) for ngram in ngrams]
17
18  #Generate unigrams, using the trainning set
19  train_unigram = generate_ngrams(train_text, n=1)
20
21  #Generate unigrams, using the testing set
22  test_unigram = generate_ngrams(test_text, n=1)
23
24  #Generate unigrams, using the development set
25  dev_unigrams = generate_ngrams(dev_text, n=1)
26
27
28
29
30  stop_words = ['a','in','on','at','and','or',
31                'to', 'the', 'of', 'an', 'by',
32                'as', 'is', 'was', 'were', 'been', 'be',
33                'are','for', 'this', 'that', 'these', 'those', 'you'
34             'it', 'he', 'she', 'we', 'they', 'will', 'have', 'has
35                'do', 'did', 'can', 'could', 'who', 'which', 'what',
36                'his', 'her', 'they', 'them', 'from', 'with', 'its']
37
38
39  #Remove these stop words from the list of ngrams
40  def remove_stop_word(ngram):
41
42      #use list comprehension,
43      #to only return words not inlcuded in stop_words
44      return [word for word in ngram if word not in stop_words]
45
46  #Remove stopwords on the train/test/dev unigrams
47  uni_train_no_sw = remove_stop_word(train_unigram)
48  uni_test_no_sw = remove_stop_word(test_unigram)
49  uni_dev_no_sw = remove_stop_word(dev_unigrams)
50
51  #Make Bigrams from the created unigrams
52  bi_train = generate_ngrams(uni_train_no_sw, n=2)
53  bi_test = generate_ngrams(uni_test_no_sw, n=2)
54  bi_dev = generate_ngrams(uni_dev_no_sw, n=2)
55
56  #Make Trigrams from the created unigrams
57  tri_train = generate_ngrams(uni_train_no_sw, n=3)
58  tri_test = generate_ngrams(uni_test_no_sw, n=3)
59  tri_dev = generate_ngrams(uni_dev_no_sw, n=3)
60
61
62  ###################################
63  #Remove ngrams appearing in less than K documents
64  def doc_counter(set_train, set_test, dev_test):
65
66      #use Counter type in order to count all unique words,
```

```
67      #from each train/test/dev set
68      c = Counter()
69      c.update(set_train)
70      c.update(set_test)
71      c.update(dev_test)
72
73      return c
74
75  #Initialise the set versions of the train/test/dev unigrams
76  set_uni_train = set(uni_train_no_sw)
77  set_uni_test = set(uni_test_no_sw)
78  set_uni_dev = set(uni_dev_no_sw)
79
80  #Call doc_counter for all the unigram sets
81  uni_doc_appearances = doc_counter(set_uni_train, set_uni_test, set
82
83  #Initialise the set versions of the train/test/dev bigrams
84  set_bi_train = set(bi_train)
85  set_bi_test = set(bi_test)
86  set_bi_dev = set(bi_dev)
87
88  #Call doc_counter for all the bigram sets
89  bi_doc_appearances = doc_counter(set_bi_train, set_bi_test, set_bi
90
91
92  set_tri_train = set(tri_train)
93  set_tri_test = set(tri_test)
94  set_tri_dev = set(tri_dev)
95
96  #Initialise the set versions of the train/test/dev trigrams
97  tri_doc_appearances = doc_counter(set_tri_train, set_tri_test, set
98
99
100 #Number of documents(set between 1 and 3)
101 def find_words(c, k):
102
103     #Output list variable
104     found_words =[]
105
106     #c is a Counter,
107     # go through every word contained by c
108     for words in c.keys():
109
110     #if documents appearance value is smaller than k
111     #in that case continue
112         if c[words] < k :
113             continue
114         else:
115             #Add this ngram to the list
116             found_words.append(words)
117
118     #No need to keep it as a list,
119     # arrays will help with efficiency
120     return np.array(found_words)
121
122
123 #Function that will return lists,
124 #containning only ngrams  appearing at least in K documents
125 def remove_k(set_train, set_test, set_dev ,doc_ap):
126
```

```
127    clean_train = []
128    clean_test = []
129    clean_dev = []
130
131    print("Finding words....", "\n")
132
133    #Get all ngrams that appear in at least k documents
134    found_words = find_words(doc_ap, k=3)
135
136    print(found_words)
137
138    print("Starting trainning....", "\n")
139    clean_train = [word for word in set_train if word in found_wor
140
141    print("Starting testing....", "\n")
142    clean_test = [word for word in set_test if word in found_words
143
144    print("Starting dev....", "\n")
145    clead_dev = [word for word in set_dev if word in found_words]
146
147    return np.array(clean_train), np.array(clean_test), np.array(c
148
```

**This code cell is separated from the rest due to the heavy computation needed.**

**Only run this cell once!**

**Estimated processing time with the full lists of ngrams: 12 minutes.**

**5 Minutes with the sets version**

In [6]:
```
 1
 2
 3
 4  print("Starting UNIGRAMS....", "\n")
 5
 6
 7  clean_uni_train, clean_uni_test, clean_uni_dev  = remove_k(set_uni_
 8                                                    set_u
 9                                                    set_u
10                                                    uni_doc
11
12
13  print("Starting BIGRAMS....", "\n")
14
15
16  clean_bi_train, clean_bi_test, clean_bi_dev = remove_k(set_bi_trair
17                                                 set_bi_test
18                                                 set_bi_dev,
19                                                  bi_doc_ap
20
21
```

```
22  print("Starting TRIGRAMS....", "\n")
23
24
25  clean_tri_train, clean_tri_test, clean_tri_dev = remove_k(set_tri_t
26                                                            set_tri_te
27                                                            set_tri_de
28                                                               tri_doc_
```

Starting UNIGRAMS....

Finding words....

['emotive' 'macy' 'sketched' ... 'ms' 'virtues' 'aiming']
Starting trainning....

Starting BIGRAMS....

Finding words....

['no spark' 'take back' 'one films' ... 'every movie' 'friend amazin
g'
 'known actors']
Starting trainning....

Starting TRIGRAMS....

Finding words....

['still doesn t' 'movie going experience' 'few far between' ...
 'pretty good but' 'scale 0 4' 'best thing about']
Starting trainning....

In [7]:
```
 1  #Create a vocabulary of unigrams, bigrams and trigrams
 2  vocab = set(clean_uni_train)
 3  vocab.update(clean_uni_test)
 4  vocab.update(clean_uni_dev)
 5
 6  vocab.update(clean_bi_train)
 7  vocab.update(clean_bi_test)
 8  vocab.update(clean_bi_dev)
 9
10  vocab.update(clean_tri_train)
11  vocab.update(clean_tri_test)
12  vocab.update(clean_tri_dev)
13
14
15
16  print(vocab, "\n")
17
```

```
{'emotive', 'macy', 'sketched', 'shrewd', 'director lasse', 'deep ri
sing', 'very believable', 'afraid', 'there more', 'goof', 'parallel
', 'civilization', 'physical appearance', 'ricci', 's possible', 'de
picting', 'covert', 'so ll', 'when', 'no sense all', 'no spark', 'po
lice station', 'store', 'years s', 'time just', 'gentleman', 's fell
ow', 'if s', 'taste', 'rd', 'gumption', 'bitchy', 'but comes', 'even
some', 'but overall', 'right hand', 'dr evil s', 'anytime', 'sexual
relationship', 'knocks', 'focuses', 'd probably', 'universal studios
', 'take back', 'doesn t get', 'founded', 'their friends', 'cloud',
```

## N-gram extraction from a document

You first need to implement the `extract_ngrams` function. It takes as input:

- `x_raw` : a string corresponding to the raw text of a document
- `ngram_range` : a tuple of two integers denoting the type of ngrams you want to extract,
  e.g. (1,2) denotes extracting unigrams and bigrams.
- `token_pattern` : a string to be used within a regular expression to extract all tokens.
  Note that data is already tokenised so you could opt for a simple white space tokenisation.
- `stop_words` : a list of stop words
- `vocab` : a given vocabulary. It should be used to extract specific features.
- `char_ngrams` : boolean. If true the function extracts character n-grams

and returns:

- `x`: a list of all extracted features.

See the examples below to see how this function should work.

```python
In [8]:
1  def extract_ngrams(x_raw, ngram_range, token_pattern,
2                     stop_words, vocab, char_ngrams):
3
4      #Set the smallest value of ngram types
5      min_ = ngram_range[0]
6
7      #Set the biggest value of ngram types
8      max_ = ngram_range[-1]
9
10     #Initialise output values
11     output_ngram = []
12     output_char_gram =[]
13
14     #Produce Character ngrams or regular ngrams
15     if char_ngrams == False:
16
17         #Go through every type of ngram(i.e. unigram, bigram)
18         for rn in range(min_,max_+1):
19             print(rn)
20
21             # Replace all none alphanumeric characters with spaces
22             x_sub = re.sub(r'[^a-zA-Z0-9\s]', ' ', str(x_raw))
23
24             x_sub.replace("'", " ")
25
26             # Break sentence in the token, remove empty tokens
27             tokens = [token for token in x_sub.split(token_pattern)
```

```
28
29                  # Use the zip function to help generate n-grams
30                  # Concatentate the tokens into ngrams and return
31                  ngrams = zip(*[tokens[i:] for i in range(rn)])
32                  final_ngrams = [" ".join(ngram) for ngram in ngrams]
33
34
35                  #Remove stop words from ngrams
36                  no_stop_ngram = [word for word in final_ngrams if word
37                  #if rn == 3:
38                      #print('This is the stop_words',*no_stop_ngram, sep
39
40                  #filter ngrams in vocabulary
41                  for word_o in no_stop_ngram:
42                      if word_o in vocab:
43                          output_ngram.append(word_o)
44
45          print(output_ngram)
46          return output_ngram
47
48      else:
49
50          #Generate character ngrams
51
52          #Go through every type of ngram(i.e. unigram, bigram)
53          for rn in range(min,max+1):
54
55                  final_char =[]
56                  #b[i:i+n] for i in range(len(b)-n+1)
57
58                  # Replace all none alphanumeric characters with spaces
59                  x_sub = re.sub(r"[^a-zA-Z0-9\s]", "", str(x_raw))
60
61                  x_sub.replace("'", "")
62                  x_sub.replace(" ","")
63                  # tokens = [token for token in x_sub.split(" ") if toke
64
65
66                  # Use the zip function to help generate character n-gra
67                  # Concatentate the tokens into ngrams and return
68                  char_grams = zip(*[x_sub[i:] for i in range(rn)])
69
70                   #Split words by character, not by whitespace
71                  final_char = ["".join(char_gram) for char_gram in char_
72
73
74                  #Remove stopwords
75                  output_char_gram = [word for word in final_char if word
76
77          print(output_char_gram)
78          return output_char_gram
```

Note that it is OK to represent n-grams using lists instead of tuples: e.g. `['great',
['great', 'movie']]`

For extracting character n-grams the function should work as follows:

In [163]:    1  `##### Will Keep this cell commneted, for preview purposes #####`

```
 2
 3   ##### To check the real running code go to cell above #####
 4
 5   # def extract_ngrams(x_raw="movie",
 6   #                     ngram_range=(2,4),
 7   #                     stop_words=[],
 8   #                     char_ngrams=True):
 9
10   #       min = ngram_range[0]
11
12   #       max = ngram_range[-1]
13
14   #       output_char_gram, no_stop_char =[]
15
16
17   #       for rn in range(min,max+1):
18
19   #               #b[i:i+n] for i in range(len(b)-n+1)
20
21   #               # Replace all none alphanumeric characters with space
22   #               x_sub = re.sub(r'[^a-zA-Z0-9\s]', ' ', str(x_raw))
23
24   #               # Break sentence in the token, remove empty tokens
25   #               # tokens = [token for token in x_sub.split(token_patte
26
27   #               # Use the zip function to help generate character n-g
28   #               # Concatentate the tokens into ngrams and return
29   #               char_grams = zip(*[x_sub[i:] for i in range(rn)])
30   #               final_char = ["".join(char_gram) for char_gram in cha
31
32
33   #               #Remove stopwords
34   #               no_stop_char = [word for word in final_char if word r
35
36
37   #               #search in vocab
38   #               output_char_gram = [word_o for word_o in no_stop_char
39
40
41
42   #       return output_char_gram
```

## Create a vocabulary

The `get_vocab` function will be used to (1) create a vocabulary of ngrams; (2) count the document frequencies of ngrams; (3) their raw frequency. It takes as input:

- `X_raw` : a list of strings each corresponding to the raw text of a document
- `ngram_range` : a tuple of two integers denoting the type of ngrams you want to extract, e.g. (1,2) denotes extracting unigrams and bigrams.
- `token_pattern` : a string to be used within a regular expression to extract all tokens. Note that data is already tokenised so you could opt for a simple white space tokenisation.
- `stop_words` : a list of stop words
- `min_df` : keep ngrams with a minimum document frequency.
- `keep_topN` : keep top-N more frequent ngrams.

and returns:

- `vocab` : a set of the n-grams that will be used as features.
- `df` : a Counter (or dict) that contains ngrams as keys and their corresponding document frequency as values.
- `ngram_counts` : counts of each ngram in vocab

Hint: it should make use of the `extract_ngrams` function

```python
In [162]:
1  def get_vocab(X_raw, ngram_range, token_pattern,
2                  min_df, keep_topN, stop_words, char_ngrams):
3
4      #Set the smallest value of ngram types
5      min_ = ngram_range[0]
6
7      #Set the biggest value of ngram types
8      max_ = ngram_range[-1]
9
10     ngrams = np.array([0])
11
12     #Go through every type of ngram(i.e. unigram, bigram)
13     for rn in range(min_, max_+1):
14         #     n = ngram-range
15         print("Filter rn.... ", rn ,"\n")
16
17
18         special_char=[",",":"," ",";",".","?","'"]
19
20         # Replace all none alphanumeric characters with spaces
21         s = re.sub(r'[^a-zA-Z0-9\s]', ' ', str(X_raw))
22
23         # Break sentence in the token, remove empty tokens
24         tokens = [token for token in s.split(" ") if token != ""]
25
26         # Use the zip function to help generate n-grams
27         # Concatentate the tokens into ngrams and return
28         n_grams = zip(*[tokens[i:] for i in range(rn)])
29         ngrams = np.append(ngrams, [" ".join(ngram) for ngram in n
30
31     print("Filter vocab....","\n")
32
33
34     #Remove stop words and special charcters from the list of ngra
35     filtered_vocab = [w for w in ngrams if w not in stop_words and
36
37
38
39     print("Start extract_ngram....","\n")
40
41     #Initialise and pass the filtered vocab as a set
42     original_vocab = set(filtered_vocab)
43
44     #Extract ngrams from the vocabulary
45     ngram = extract_ngrams(X_raw, ngram_range, token_pattern, stop
46
47     #print(ngram)
48
49     #Count all the ngrams
50     df_count = Counter()
51     df_count.update(ngram)
```

```python
    def Compute_DF(ngrams):

        DF = {}

        print("Started DFs small.. \n")
        for i in range(len(ngrams)):

            for w in ngrams[i]:
                try:
                    DF[w].add(i)
                except:
                    DF[w] = {i}


        for i in DF:
            DF[i] = len(DF[i])
        print(DF)
        return DF


    def find_doc_freq(word,DF):

        #Method to get a specific ngram's Document frequency

        c = 0

        try:
            c = DF[word]
        except:
            pass
        return c



    DF = Compute_DF(ngram)

    vocab_init = set(ngram)

#     found_gram = np.array([0])
    #Filter ngrams through vocabulary PROBLEM
#     found_gram = [w for w in ngram if w in vocab]


    N = len(ngram)
#     print(N)

    #Initialise a new Counter to pass into ngrams with a count hig
    df_final = Counter()


    df ={}


    for i in range(N):

        tokens = found_gram[i]
        counter = Counter(tokens) #Replace with count vector
```

```
112            words_count = len(tokens)
113
114            #df_final.update(np.unique(tokens))
115
116            for token in np.unique(tokens):
117                tf = counter[token]/words_count
118                df_word = find_doc_freq(token,DF)
119                if df_word >= min_df:
120                    df.update({token: df_word}) #was df
121                    df_final.update(token)
122
123        vocab = set()
124        ngram_counts = []
125
126        #Go through the the top n most common ngrams,
127        # and extract their raw frequency and word
128
129        for word, count in df_final.most_common(keep_topN):
130            vocab.add(word)
131            ngram_counts.append(count)#Count is raw frequency
132
133
134        print(vocab)
135        print(df)
136        print(ngram_counts)
137 #      print(type(top_ngrams))
138
139 #      count = Counter()
140 #      count.update(top_ngrams)
141
142
143 #      ngram_counts = count.values()
144 #      print(types(ngram_counts))
145
```

Now you should use `get_vocab` to create your vocabulary and get document and raw frequencies of n-grams:

```
In [164]:   1 test_vocab, test_df, test_count = get_vocab(test_text, ngram_range=
            2                    min_df=2, keep_topN=500,
            3                    stop_words = stop_words,char_ngrams = False)
            4
            5
            6
            7 # print('TEST VOCAB: ', test_vocab, '\n')
```

```
Filter rn....  1

Filter rn....  2

Filter rn....  3

Filter vocab....

Start extract_ngram....

1
2
3
```

```
IOPub data rate exceeded.
The notebook server will temporarily stop sending output
to the client in order to avoid crashing it.
To change this limit, set the config variable
`--NotebookApp.iopub data rate limit`.
```

Then, you need to create 2 dictionaries: (1) vocabulary id -> word; and (2) word -> vocabulary id so you can use them for reference:

```python
In [24]:   1  def create_2dict(df_dict):
           2      id2word = {}
           3      word2id = {}
           4      dic_id = 0
           5
           6      for word in test_df.keys():
           7
           8              #(1) vocabulary id -> word
           9              id2word.update({dic_id : word})
          10
          11              # (2) word -> vocabulary id
          12              word2id.update({word: dic_id})
          13
          14              dic_id += 1
          15
          16      print('Dictionary [ID : WORD] : ',id2word, "\n")
          17      print('Dictionary [WORD : ID] : ',word2id, "\n")
          18
          19      return id2word , word2id
```

Now you should be able to extract n-grams for each text in the training, development and test sets:

```python
In [165]:   1  #TEST
            2  # test_vocab, test_df, test_count = get_vocab(test_text, ngram_rang
            3  #                          min_df=2, keep_topN=500,
            4  #                          stop_words = stop_words)
            5
            6  #train
            7  train_vocab, train_df, train_count = get_vocab(train_text, ngram_ra
            8                          min_df=10, keep_topN=100,
            9                          stop_words=stop_words,char_ngrams=False)
           10  #Dev
           11  dev_vocab, dev_df, dev_count = get_vocab(dev_text, ngram_range=(1,3
           12                          min_df=10, keep_topN=100,
           13                          stop_words=stop_words,char_ngrams=False)
           14
```

Filter rn....  1

IOPub data.rate2exceeded.
The notebook server will temporarily stop sending output
to the client in order to avoid crashing it.
To change this limit, set the config variable
`--NotebookApp.iopub_data_rate_limit`.

```
In [30]:   1  print("Test Dictionary >>> ", "\n")
           2  id_w_test, w_id_test = create_2dict(test_df)
           3
           4  print("Train Dictionary >>> ", "\n")
           5  id_w_train, w_id_train =create_2dict(train_df)
           6
           7  print("DEV Dictionary >>> ", "\n")
```

Test Dictionary >>>

Dictionary [ID : WORD] :  {0: 'know', 1: 'but', 2: 'got', 3: 'around
', 4: 'last', 5: 'one', 6: 'about', 7: 'final', 8: 'scene', 9: 'out
', 10: 'enough', 11: 'watch', 12: 'such', 13: 'good', 14: 'behind',
15: 'show', 16: 'most', 17: 'well', 18: 'gets', 19: 'school', 20: 's
', 21: 'plays', 22: 'him', 23: 'plot', 24: 'help', 25: 'very', 26: '
finds', 27: 'himself', 28: 'love', 29: 'fun', 30: 'begins', 31: 'go
', 32: 'while', 33: 'goes', 34: 'like', 35: 'too', 36: 'young', 37:
't', 38: 'way', 39: 'two', 40: 'people', 41: 'really', 42: 'up', 43:
'year', 44: 'completely', 45: 'if', 46: 'me', 47: 'maybe', 48: 'best
', 49: 'picture', 50: 'instead', 51: 'film', 52: 'reason', 53: 'so',
54: 'point', 55: 'next', 56: 'turn', 57: 'gives', 58: 'movie', 59: '
how', 60: 'having', 61: 'bit', 62: 'bad', 63: 'not', 64: 'being', 6
5: 'big', 66: 'woman', 67: 'performance', 68: 'had', 69: 'just', 70:
'because', 71: 'both', 72: 'into', 73: 'great', 74: 'character', 75:
'long', 76: 'when', 77: 'course', 78: 'actually', 79: 'even', 80: 'm
an', 81: 'every', 82: 'something', 83: 'life', 84: 'first', 85: 'fri
end', 86: 'where', 87: 'couple', 88: 'getting', 89: 'see', 90: 'anyt

## Vectorise documents

Next, write a function `vectoriser` to obtain Bag-of-ngram representations for a list of documents. The function should take as input:

- `X_ngram` : a list of texts (documents), where each text is represented as list of n-grams in the `vocab`
- `vocab` : a set of n-grams to be used for representing the documents

and return:

- `X_vec` : an array with dimensionality Nx|vocab| where N is the number of documents and |vocab| is the size of the vocabulary. Each element of the array should represent the frequency of a given n-gram in a document.

```
In [152]:   1  def vectorise(X_ngram, vocab):
            2
            3
            4      def count_vectorize(tokens):
```

```python
        ''' This function takes list of words in a sentence as inp
        and returns a vector of size of filtered_vocab.It puts 0 i
        word is not present in tokens and count of token if presen

        vector = np.array([0])
        for w in np.array(filtered_vocab):
            vector = np.append(vector, tokens.count(w))


        return vector

    def Compute_DF(ngrams):

        DF = {}

        print("Started DFs small.. \n")
        for i in range(len(ngrams)):

            for w in ngrams[i]:
                try:
                    DF[w].add(i)
                except:
                    DF[w] = {i}


        for i in DF:
            DF[i] = len(DF[i])
        print(DF)
        return DF


    def find_doc_freq(word,DF):

        #Method to get a specific ngram's Document frequency

        c = 0

        try:
            c = DF[word]
        except:
            pass
        return c



    def compute_tf_IDF(ngram):

        #Calculate the Document Frequency

        DF = Compute_DF(ngram)




        doc = 0
        token_counter = 0

        print("Started TF.IDF...\n")

        #Calculate TF.IDF
```

```python
 65
 66            # N=Total number of documents in the dataset
 67
 68
 69        found_gram = np.array([0])
 70
 71        #Filter ngrams through vocabulary PROBLEM
 72        found_gram = [w for w in ngram if w in filtered_vocab]
 73
 74
 75
 76        print(found_gram)
 77
 78        N = len(found_gram)
 79        print(N)
 80
 81        vocab_size = len(filtered_vocab)
 82
 83
 84        dim_row = N
 85        dim_columns = vocab_size
 86
 87
 88        tf_idf = [[0 for j in range(dim_columns)] for i in range(d
 89        print(tf_idf)
 90
 91        for i in range(N):
 92
 93            tokens = found_gram[i]
 94            counter = Counter(tokens)#Replace with count vector
 95            words_count = len(tokens)
 96
 97            token_counter =0
 98            for token in np.unique(tokens):
 99
100
101                tf = counter[token]/words_count
102                df = find_doc_freq(token,DF)
103                idf = np.log(N/(df+1)) #numerator is added 1 to av
104
105                # df=total number of documents in which nth word o
106
107                tf_idf[i][token_counter] = tf*idf
108                token_counter +=1
109 #                 tf_idf[doc, token] = tf*idf
110
111 #            doc += 1
112 #            token_counter += 1
113
114        print(np.array(tf_idf))
115        return np.array(tf_idf)
116
117
118
119
120
121
122
123    #list of special characters.You can use regular expressions to
124    special_char=[",",":"," ",";",".","?","'"]
```

```
125
126
127        #split the sentences into tokens
128        x_sub = re.sub(r"[^a-zA-Z0-9\s]", " ", str(X_ngram))
129
130        tokens1 = [token for token in x_sub.split(" ") if token != ""]
131
132
133        #filter the vocabulary list
134        filtered_vocab = [w for w in vocab if w not in stop_words and
135
136        #print(filtered_vocab)
137
138        print("Count Vector...\n")
139        vector1=count_vectorize(tokens1)
140
141        print("Start compute_tf_IDF...\n")
142        TF_IDF_vector=compute_tf_IDF(tokens1)
143
144
```

Finally, use `vectorise` to obtain document vectors for each document in the train, development and test set. You should extract both count and tf.idf vectors respectively:

## Count vectors

In [153]:
```
1  #COPY COUNT VECTORIZER HERE
2
3
4
5
6
7  #UNIGRAMS, SET_UNIGRAMS
8  print("Vectorise test text....","\n")
9  test_count, test_vect = vectorise(test_text, test_vocab)
10
11 print("Vectorise train text....","\n")
12 train_count, train_vect = vectorise(train_text, train_vocab)
13
14 print("Vectorise dev text....","\n")
```

Vectorise test text....

```
1  print('Array Shape = ',np.shape(test_vect) )   # test_vect.shape[0]
2  print('Array Shape = ',np.shape(train_vect) )
3
4  # print(np.shape(train_vect[0:19724, 0:]))
5
6  train_vect_sliced = train_vect[0:19724, 0:]#reshape
7  print('Array Shape = ',np.shape(train_vect_sliced) )
8
```

```
Array Shape =   (72016, 500)
Array Shape =   (141944, 100)
Array Shape =   (19724, 100)
Array Shape =   (19724, 100)
```

**TF.IDF vectors**

First compute `idfs` an array containing inverted document frequencies (Note: its elements should correspond to your `vocab` )

In [ ]:

```
1  ######COMPUTE TF.IFD using Term Frequency, Document Frequency and i
2
3  #Copy the TF.IDF method
4
5  #      tf_idf = {}
6  #      for i in range(N):
7      #       tokens = processed_text[i]
8      #       counter = Counter(tokens)
9      #       words_count = len(tokens)
10
11     #       for token in np.unique(tokens):
12  #               tf = count vector
13     #            df = doc_freq(token)
14
15     # ------> idf = np.log(N/(df+1)) <------
16
17
18
19
20  #  Formula can be one of these two:
21  #
22  #  IDF = 1+log(N/dN)
23  #
24  #  idf = log(N/(dN+1))
25
26  # Where
27
28  # N=Total number of documents in the dataset
```

Then transform your count vectors to tf.idf vectors:

In [ ]:

```
1  #          tf = counter[token]/words_count
2  # Replace with the count vector
```

In [ ]:

```
1  #          tf_idf[doc, token] = tf*idf
```

# Binary Logistic Regression

After obtaining vector representations of the data, now you are ready to implement Binary Logistic Regression for classifying sentiment.

First, you need to implement the `sigmoid` function. It takes as input:

- `z` : a real number or an array of real numbers

and returns:

- `sig` : the sigmoid of `z`

```
In [34]:   1  def sigmoid(z):
           2
           3      sig = 1 / (1 + np.exp(-z))
           4      return sig
           5
           6  # z = np.dot(X, theta)
           7  # h = sigmoid(z)
```

Then, implement the `predict_proba` function to obtain prediction probabilities. It takes as input:

- `X` : an array of inputs, i.e. documents represented by bag-of-ngram vectors $(N \times |vocab|)$
- `weights` : a 1-D array of the model's weights $(1, |vocab|)$

and returns:

- `preds_proba` : the prediction probabilities of X given the weights

```
In [35]:   1  def predict_proba(X, weights):
           2
           3      preds_proba = sigmoid(np.dot(X, weights))
           4
           5
```

Then, implement the `predict_class` function to obtain the most probable class for each vector in an array of input vectors. It takes as input:

- `X` : an array of documents represented by bag-of-ngram vectors $(N \times |vocab|)$
- `weights` : a 1-D array of the model's weights $(1, |vocab|)$

and returns:

- `preds_class` : the predicted class for each x in X given the weights

```
In [36]:   1  def predict_class(X, weights):
           2
           3
```

```
 4
 5        """
 6 #        Predict the class between 0 and 1 using learned logistic
 7 #        Using threshold value 0.5 to convert probability value to
 8
 9 #        I/P
10 #        ----------
11 #        X : 2D array where each row represents a docuemnt  and ea
12 #
13 #        weights : 1D array of weights. Dimension (1 x |vocab|)
14
15 #        O/P
16 #        -------
17 #        Class type based on threshold
18 #        """
19
20     p = preds_proba(X,weights) >= 0.5
21
22     preds_class = p.astype(int)
23
24 #   if y_pred_tr>=0.5: #LABELS
25
26 #        predictions.append(1)
27 #   else:
28 #        predictions.append(0)
29
30     return preds_class
```

To learn the weights from data, we need to minimise the binary cross-entropy loss. Implement `binary_loss` that takes as input:

- `X` : input vectors
- `Y` : labels
- `weights` : model weights
- `alpha` : regularisation strength

and return:

- `l` : the loss score

```
In [37]:  1 def binary_loss(X, Y, weights, alpha=0.00001):
          2
          3     """
          4 #        Compute cost for logistic regression.
          5
          6 #        I/P
          7 #        ----------
          8 #        X : 2D array where each row represents a document  and ea
          9 #
         10 #        y : 1D array of labels/target value for each traing examp
         11
         12 #        weights : 1D array of fitting parameters or weights. Dime
         13
         14 #        alpha: regularisation strengths to be added when calculat
         15
         16 #        O/P
         17 #        -------
```

```
18  #          J : The cost of using theta as the parameter for linear r
19  #          """
20
21
22      m = len(X)
23      yhat = sigmoid(np.dot(X, weights) + alpha)
24
25      predict = Y * np.log(yhat) + (1 - Y) * np.log(1 - yhat)
26
27      l = -sum(predict) / m
28
29
30
31
32      return l
33
34
35
```

Now, you can implement Stochastic Gradient Descent to learn the weights of your sentiment classifier. The `SGD` function takes as input:

- `X_tr` : array of training data (vectors)
- `Y_tr` : labels of `X_tr`
- `X_dev` : array of development (i.e. validation) data (vectors)
- `Y_dev` : labels of `X_dev`
- `lr` : learning rate
- `alpha` : regularisation strength
- `epochs` : number of full passes over the training data
- `tolerance` : stop training if the difference between the current and previous validation loss is smaller than a threshold
- `print_progress` : flag for printing the training progress (train/validation loss)

and returns:

- `weights` : the weights learned
- `training_loss_history` : an array with the average losses of the whole training set after each epoch
- `validation_loss_history` : an array with the average losses of the whole development set after each epoch

In [174]:
```
1   def SGD(X_tr, Y_tr, X_dev, Y_dev, lr,
2           alpha, epochs,
3           tolerance, print_progress):
4
5   #          X = # data points with some features which we want to tr
6   #          y = # labels of all datapoints
7   #          # Initialize the weights and bias i.e. 'm' and 'c'
8   #          m = np.zeros_like(X[0]) # array with shape equal to no.
9   #          c = 0#regularisation
10  #          LR = 0.0001  # The learning Rate
11  #          epochs = 50 # no. of iterations for optimization
12
13
14  #     w=np.zeros(shape=(1,train_data.shape[1]-1))
```

```
15
16  #     C = f_integ(np.array([1]))
17  #     print "C", C
18      m_tr = np.zeros_like(X_tr)
19
20      m_dev = np.zeros_like(X_dev)
21
22      alpha_tr = alpha
23      alpha_dev = alpha
24
25
26      training_loss_history = np.array([0])
27      validation_loss_history = np.array([0])
28
29      training_loss_prev = np.array([0])
30      validation_loss_prev = np.array([0])
31
32      training_loss_current = np.array([0])
33      validation_loss_current = np.array([0])
34
35
36      # for every epoch
37      for epoch in range(1,epochs+1):
38
39          ####TRAINNING####
40          # for every data point(X_train,y_train)
41          for i in range(len(X_tr)):
42
43              #compute gradient w.r.t 'm'
44              form_train = np.dot(X_tr[i], m_tr.T) + alpha_tr
45
46              gr_wrt_m_tr = X_tr[i]*(Y_tr[i] - sigmoid(form_train))
47
48              #compute gradient w.r.t 'c'
49              gr_wrt_c_tr = Y_tr[i] - sigmoid(form_train)          #up
50
51              m_tr = m_tr - lr * gr_wrt_m_tr
52
53              alpha_tr = alpha_tr - lr * gr_wrt_c_tr# At the end of
54
55
56
57
58          if training_loss_prev == np.array([0]):
59
60              training_loss_prev = binary_loss(X_tr,Y_tr,m_tr,alpha_
61              training_loss_history = np.append(training_loss_histor
62
63          else:
64
65              training_loss_current = binary_loss(X_tr,Y_tr,m_tr,alp
66
67              if (training_loss_current - training_loss_prev) >= tol
68
69                  training_loss_history = np.append(training_loss_hi
70                  training_loss_prev = training_loss_current
71
72
73          #          if i % 10000 == 0:
74          if print_progress == True:
```

```python
 75                print("Loss after %d steps is: %.10f " % (epoch,traini
 76
 77            ####Development####
 78            # for every data point(X_train,y_train)
 79            for j in range(len(X_dev)):
 80
 81                #compute gradient w.r.t 'm'
 82                form_train = np.dot(X_dev[j], m_dev.T) + alpha_dev
 83
 84  #              In [1]: import numpy
 85
 86  #              In [2]: numpy.dot(numpy.ones([97, 2]), numpy.ones([2
 87  #              Out[2]: (97, 1)
 88
 89                gr_wrt_m_dev = X_dev[j]*(Y_dev[j] - sigmoid(form_train
 90
 91                #compute gradient w.r.t 'c'
 92                gr_wrt_c_dev = Y_tr[j] - sigmoid(form_train)        #u
 93
 94                m_dev = m_dev - lr * gr_wrt_m_dev
 95
 96                alpha_dev = alpha_dev - lr * gr_wrt_c_dev# At the end
 97
 98
 99
100
101            if validation_loss_prev == np.array([0]):
102
103                validation_loss_prev = binary_loss(X_dev,Y_dev,m_dev,a
104                validation_loss_history = np.append(validation_loss_hi
105
106            else:
107
108                validation_loss_current = binary_loss(X_dev,Y_dev,m_de
109
110                if (validation_loss_current - validation_loss_prev) >=
111
112                    validation_loss_history = np.append(validation_los
113                    validation_loss_prev = validation_loss_current
114
115
116
117  #          validation_loss_history = np.append(validation_loss_hist
118
119            if print_progress == True:
120                print("Loss after %d steps is: %.10f " % (epoch,valida
121
122
123  #      weights
124
125  #      binary_loss(X_tr,Y_tr,m_tr,alpha_tr)
126
127  #      binary_loss(X_dev,Y_dev,m_dev,alpha_dev)
128
129      if print_progress == True:
130          print("Final loss after %d steps is: %.10f " % (epoch,trai
131          print("Loss after %d steps is: %.10f " % (epoch,validation
132          print("Final weights for trainning: ", m_tr,"\n")
133          print("Final weights for development: ", m_dev,"\n")
134
```

```python
135     weigths = np.array([0])
136     weigths = np.append(weigths, m_tr)
137     weigths = np.append(weigths, m_dev)
138
139     # So by using those optimum values of 'm' and 'c' we can perfo
140     #############MAYBE CALL predict class ###########
141 #     for i in range(len(X_tr)):
142 #         z_tr = np.dot(X_tr[i], m) + alpha
143 #         y_pred_tr = sigmoid(z_tr)
144
145 #         if y_pred_tr>=0.5: #LABELS
146 #             predictions.append(1)
147 #         else:
148 #             predictions.append(0)
149
150 #     for i in range(len(X_dev)):
151
152 #         z_dev = np.dot(X_dev[i], m) + alpha
153 #         y_pred_dev = sigmoid(z_dev)
154
155 #         if y_pred_dev>=0.5:#LABELS
156 #             predictions.append(1)
157 #         else:
158 #             predictions.append(0)
159
160
161
162
163
164     # Make a prediction with coefficients
165 #     def predict(row, coefficients):
166 #         yhat = coefficients[0]
167 #         for i in range(len(row)-1):
168 #             yhat += coefficients[i + 1] * row[i]
169 #         return 1.0 / (1.0 + exp(-yhat))
170
171
172 #     # Estimate logistic regression coefficients using stochastic
173 #     def coefficients_sgd(train, l_rate, n_epoch):
174 #         coef = [0.0 for i in range(len(train[0]))]
175 #         for epoch in range(n_epoch):
176 #             for row in train:
177 #                 yhat = predict(row, coef)
178 #                 error = row[-1] - yhat
179 #                 coef[0] = coef[0] + l_rate * error * yhat * (1.0
180 #                 for i in range(len(row)-1):
181 #                     coef[i + 1] = coef[i + 1] + l_rate * error *
182 #         return coef
183
184 #     # Linear Regression Algorithm With Stochastic Gradient Desce
185 #     def logistic_regression(train, test, l_rate, n_epoch):
186 #         predictions = list()
187 #         coef = coefficients_sgd(train, l_rate, n_epoch)
188 #         for row in test:
189 #             yhat = predict(row, coef)
190 #             yhat = round(yhat)
191 #             predictions.append(yhat)
192 #         return(predictions)
193
194 #     def MyCustomSGD(train_data,learning_rate,n_iter,k,divideby):
```

```python
# 		# Initially we will keep our W and B as 0 as per the Tra
# 		w=np.zeros(shape=(1,train_data.shape[1]-1))
# 		b=0

# 		cur_iter=1
# 		while(cur_iter<=n_iter):

# 			# We will create a small training data set of size K
# 			temp=train_data.sample(k)

# 			# We create our X and Y from the above temp dataset
# 			y=np.array(temp['price'])
# 			x=np.array(temp.drop('price',axis=1))

# 			# We keep our initial gradients as 0
# 			w_gradient=np.zeros(shape=(1,train_data.shape[1]-1))
# 			b_gradient=0

# 			for i in range(k): # Calculating gradients for point
# 				prediction=np.dot(w,x[i])+b
# 				w_gradient=w_gradient+(-2)*x[i]*(y[i]-(predictio
# 				b_gradient=b_gradient+(-2)*(y[i]-(prediction))

# 			#Updating the weights(W) and Bias(b) with the above
# 			w=w-learning_rate*(w_gradient/k)
# 			b=b-learning_rate*(b_gradient/k)

# 			# Incrementing the iteration value
# 			cur_iter=cur_iter+1

# 			#Dividing the learning rate by the specified value
# 			learning_rate=learning_rate/divideby

# 		return w,b #Returning the weights and Bias
#   ############################################################
#     class LogisticRegressionCustom():

#     def __init__(self, l_rate=1e-5, n_iterations=50000):
#         self.l_rate = l_rate
#         self.n_iterations = n_iterations

#     def initial_weights(self, X):
#         self.weights = np.zeros(X.shape[1])

#     def sigmoid(self, s):
#         return 1/(1+np.exp(-s))

  #        m = len(X)
#     yhat = sigmoid(np.dot(X, weights) + alpha)

#     predict = Y * np.log(yhat) + (1 - Y) * np.log(1 - yhat)

#     l = -sum(predict) / m
#     return l

#     def binary_cross_entropy(self, X, y):
#         return -(1/len(y))*(y*np.log(self.sigmoid(np.dot(X,self.

#     def gradient(self, X, y):
```

```
255  #          return np.dot(X.T, (y-self.sigmoid(np.dot(X,self.weights
256  #
257  #      def fit(self, X, y):
258  #          self.initial_weights(X)
259  #          for i in range(self.n_iterations):
260  #              self.weights = self.weights+self.l_rate*self.gradien
261  #              if i % 10000 == 0:
262  #                  print("Loss after %d steps is: %.10f " % (i,self
263  #          print("Final loss after %d steps is: %.10f " % (i,self.b
264  #          print("Final weights: ", self.weights)
265  #          return self
266  #
267  #      def predict(self, X):
268  #          y_predict = []
269  #          for t in X:
270  #              y_predict.append(1) if self.sigmoid(np.dot(self.weig
271  #          return y_predict
272  #
273  #      def predict_proba(self, X):
274  #          y_predict = []
275  #          for t in X:
276  #              y_predict.append(self.sigmoid(np.dot(self.weights,t)
277  #          return y_predict
278  #
279  ###############################################################
280  #
281  #
282  #
283  #          def sigmoid(z):
284  #            sig = 1/(1+np.exp(-z))
285  #             return sig# Performing Gradient Descent Optimization
286  #
287  #
288  #
289  #
290  #
291
292      return weights, training_loss_history, validation_loss_history
```

## Train and Evaluate Logistic Regression with Count vectors

First train the model using SGD:

```
In [175]:   1  print(type(train_vect))
            2  print(np.shape(train_label))
            3  print(type(train_count))
            4
            5  #BOW-count
            6
            7  weights, training_loss_history, validation_loss_history = SGD(train
            8                                                          alpha
            9                                                          toler
           10
           11  # (X_tr, Y_tr, X_dev, Y_dev, lr=0.1,
           12  #       alpha=0.00001, epochs=5,
           13  #       tolerance=0.0001, print_progress=True):
           14
           15  # print("Vectorise test text....","\n")
           16  # test_vect = vectorise(test_text, test_vocab)
```

```
17  # print("Vectorise train text....","\n")
18  # train_vect = vectorise(train_text, train_vocab)
19  # print("Vectorise dev text....","\n")
20  # dev_vect = vectorise(dev_text, dev_vocab)
21
22
23  # print(train_label,"\n")
24
25
26  # ###########################################
27
28
29  # #put the testing raw texts into Python lists
30  # test_text = list(test[test_column_names[0]])
31
32  # #print the text for verification
33  # #print(test_text,"\n")
34
35  # #put the testing labels into a NumPy arrays
36  # test_label = test[test_column_names[1]].values
37
38  # #print the label for verification
39  # print(test_label,"\n")
40
41
42  # ###########################################
43
44
45  # #put the development raw texts into Python lists
46  # dev_text = list(dev[dev_column_names[0]])
47
48  # #print the text for verification
49  # #print(dev_text,"\n")
50
51  # #put the development labels into a NumPy arrays
52  # dev_label = dev[dev_column_names[1]].values
53
54  # #print the label for verification
55  # # print(dev_label,"\n")
```

```
<class 'numpy.ndarray'>
(1399,)
<class 'list'>

<ipython-input-34-6527dd331435>:3: RuntimeWarning: overflow encounte
red in exp
  sig = 1 / (1 + np.exp(-z))
<ipython-input-37-8d7bf092e40c>:8: RuntimeWarning: divide by zero en
countered in log
  predict = Y * np.log(yhat) + (1 - Y) * np.log(1 - yhat)

--------------------------------------------------------------------
-------
ValueError                                Traceback (most recent cal
l last)
<ipython-input-175-e176f2c34a9e> in <module>
      5 #BOW-count
      6
----> 7 weights, training_loss_history, validation_loss_history = SG
D(train_count, train_label,dev_count, dev_label, lr=0.1,
      8
```

```
                alpha=0.00001, epochs=5,
                        9
                tolerance=0.0001, print_progress=True)

                <ipython-input-174-58f0580d6eb2> in SGD(X_tr, Y_tr, X_dev, Y_dev, l
                r, alpha, epochs, tolerance, print_progress)
                        58              if training_loss_prev == np.array([0]):
                        59
                ---> 60                 training_loss_prev = binary_loss(X_tr,Y_tr,m_tr,
                alpha_tr)
                        61                 training_loss_history = np.append(training_loss_
                history, training_loss_prev)
                        62

                <ipython-input-37-8d7bf092e40c> in binary_loss(X, Y, weights, alpha)
                        6       yhat = sigmoid(np.dot(X, weights) + alpha)
                        7
                ----> 8       predict = Y * np.log(yhat) + (1 - Y) * np.log(1 - yhat)
                        9
                        10      l = -sum(predict) / m

                ValueError: operands could not be broadcast together with shapes (13
                99,) (100,)
```

In [ ]:

Now plot the training and validation history per epoch for the best hyperparameter combination. Does your model underfit, overfit or is it about right? Explain why.

In [ ]:
```python
 1  # #plot
 2
 3  # from sklearn.metrics import roc_curve, roc_auc_score
 4  # fpr, tpr, _ = roc_curve(y_test,  y_prob)
 5  # auc = roc_auc_score(y_test, y_prob)
 6
 7  # plt.figure(figsize=(10,8))
 8  # plt.plot(fpr,tpr,label="data, auc="+str(round(auc,4)))
 9
10  # plt.xlabel("False Positive Rate")
11  # plt.ylabel("True Positive Rate")
12
13  # plt.title("ROC Curve for Model from Sratch")
14  # plt.legend(loc=4)
15  # plt.show()
16
17  ##################### MAIN ####################
18  # training_loss_history
19  # validation_loss_history
20
21  plt.figure(figsize=(25,6))
22
23  plt.title('Cost Function Slope')
24  plt.plot(training_loss_history, label='Training Loss History')
25  plt.plot(validation_loss_history, label='Validation Loss History')
26  plt.legend(prop={'size': 16})
27  plt.xlabel('Number of Iterations')
28  plt.ylabel('Error Values')
29  plt.show()
```

```
30
31  ##################### MAIN #####################
32
33  # plt.figure(figsize=(10,8))
34  # plt.title('Cost Function Slope')
35  # plt.plot(cost)
36  # plt.xlabel('Number of Iterations')
37  # plt.ylabel('Error Values')
```

Explain here...

```
In [ ]:   1  #Underfit??
          2
          3  #Overfit??
          4
```

**Evaluation**

Compute accuracy, precision, recall and F1-scores:

```
In [170]:   1  X_te_count = train_count
            2
            3  w_count = weights
            4
            5  preds_te_count = predict_class(X_te_count, w_count)
            6
            7  # train_count, weights
            8
            9
           10  Y_te = dev_count
           11
           12  print('Accuracy:', accuracy_score(Y_te,preds_te_count))
           13  print('Precision:', precision_score(Y_te,preds_te_count))
           14  print('Recall:', recall_score(Y_te,preds_te_count))
```

```
---------------------------------------------------------------------
-------
NameError                                 Traceback (most recent cal
l last)
<ipython-input-170-f6a9ad36b6d1> in <module>
      1 X_te_count = train_count
      2
----> 3 w_count = weights
      4
      5 preds_te_count = predict_class(X_te_count, w_count)

NameError: name 'weights' is not defined
```

Finally, print the top-10 words for the negative and positive class respectively.

```
In [171]:   1  # id_w_test, w_id_test = create_2dict(test_df)
            2  # print("Train Dictionary >>> ", "\n")
            3
            4  # id_w_train, w_id_train =create_2dict(train_df)
            5  # print("DEV Dictionary >>> ", "\n")
            6
```

```
 7  # id_w_dev, w_id_dev =create_2dict(dev_df)
 8
 9  top_neg = w_count.argsort()[:10]
10  for i in top_neg:
11  #     print(id2word[i])
```

---------------------------------------------------------------
-------
NameError                                Traceback (most recent cal
l last)
<ipython-input-171-e06a79cbdd75> in <module>
      7 # id_w_dev, w_id_dev =create_2dict(dev_df)
      8
----> 9 top_neg = w_count.argsort()[:10]
     10 for i in top_neg:
     11 #     print(id2word[i])

NameError: name 'w_count' is not defined


In [172]:
```
1  top_pos = w_count.argsort()[::-1][:10]
2  for i in top_pos:
3  #     print(id2word[i])
```

---------------------------------------------------------------
-------
NameError                                Traceback (most recent cal
l last)
<ipython-input-172-b8de8e434f3d> in <module>
----> 1 top_pos = w_count.argsort()[::-1][:10]
      2 for i in top_pos:
      3 #     print(id2word[i])
      4     print(id_w_train)

NameError: name 'w_count' is not defined


If we were to apply the classifier we've learned into a different domain such laptop reviews or restaurant reviews, do you think these features would generalise well? Can you propose what features the classifier could pick up as important in the new domain?


Provide your answer here...

Sentiment Analysis


**Discuss how did you choose model hyperparameters (e.g. learning rate and regularisation strength)? What is the relation between training epochs and learning rate? How the regularisation strength affects performance?**


Enter your answer here...

(e.g. learning rate and regularisation strength)


# Train and Evaluate Logistic Regression with TF.IDF

## vectors

Follow the same steps as above (i.e. evaluating count n-gram representations).

**Now repeat the training and evaluation process for BOW-tfidf, BOCN-count, BOCN-tfidf, BOW+BOCN including hyperparameter tuning for each model...**

## BOW-tfidf:

```
In [ ]:  1  ###########BOW-tfidf###########
         2
         3  # #TEST
         4  # test_vocab, test_df, test_count = get_vocab(test_text, ngram_rang
         5  #                           min_df=2, keep_topN=500,
         6  #                           stop_words = stop_words,char_ngrams=False)
         7
         8  # #train
         9  # train_vocab, train_df, train_count = get_vocab(train_text, ngram_
        10  #                           min_df=10, keep_topN=100,
        11  #                           stop_words=stop_words,char_ngrams=False)
        12  # #Dev
        13  # dev_vocab, dev_df, dev_count = get_vocab(dev_text, ngram_range=(1
        14  #                           min_df=10, keep_topN=100,
        15  #                           stop_words=stop_words,char_ngrams=False)
        16
        17  # #Test Vectorisation
        18  # print("Vectorise test text....","\n")
        19  # test_count, test_vect = vectorise(test_text, test_vocab)
        20
        21  # #Train Vectorisation
        22  # print("Vectorise train text....","\n")
        23  # train_count, train_vect = vectorise(train_text, train_vocab)
        24
        25  # #Dev Vectorisation
        26  # print("Vectorise dev text....","\n")
        27  # dev_count, dev_vect = vectorise(dev_text, dev_vocab)
        28
        29
        30  weights_tfidf, training_loss_history_tfidf, validation_loss_history
        31                                              alpha
        32                                              toler
        33
        34  X_te_count = train_vect
        35
        36  w_count = weights_tfidf
        37
        38  preds_te_count = predict_class(X_te_count, w_count)
        39
        40  # train_count, weights
        41
        42
        43  Y_te = dev_vect
        44
        45  print('Accuracy:', accuracy_score(Y_te,preds_te_count))
        46  print('Precision:', precision_score(Y_te,preds_te_count))
```

```python
47  print('Recall:', recall_score(Y_te,preds_te_count))
48  print('F1-Score:', f1_score(Y_te,preds_te_count))
49
50  # training_loss_history
51  # validation_loss_history
52
53  plt.figure(figsize=(25,6))
54
55  plt.title('Cost Function Slope')
56  plt.plot(training_loss_history_tfidf, label='Training Loss History'
57  plt.plot(validation_loss_history_tfidf, label='Validation Loss Hist
58  plt.legend(prop={'size': 16})
59  plt.xlabel('Number of Iterations')
60  plt.ylabel('Error Values')
61  plt.show()
62
63  # print("Test Dictionary >>> ", "\n")
64  # id_w_test, w_id_test = create_2dict(test_df)
65
66  # print("Train Dictionary >>> ", "\n")
67  # id_w_train, w_id_train =create_2dict(train_df)
68
69  # print("DEV Dictionary >>> ", "\n")
70  # id_w_dev, w_id_dev =create_2dict(dev_df)
71
72
73
74  top_neg = w_count.argsort()[:10]
75  for i in top_neg:
76  #     print(id2word[i])
77      print(id_w_train[i])
78
79  top_pos = w_count.argsort()[::-1][:10]
80  for i in top_pos:
81  #     print(id2word[i])
82      print(id_w_train[i])
```

## BOCN-count:

```python
In [ ]:   1  ############ BOCN-count ############
          2
          3  # #TEST
          4  # test_vocab, test_df, test_count = get_vocab(test_text, ngram_rang
          5  #                     min_df=2, keep_topN=500,
          6  #                     stop_words = stop_words,char_ngrams=True)
          7
          8  #train
          9  train_vocab_BOCN, train_df_BOCN, train_count_BOCN = get_vocab(trair
         10                      min_df=10, keep_topN=100,
         11                      stop_words=stop_words,char_ngrams=True)
         12  #Dev
         13  dev_vocab_BOCN, dev_df_BOCN, dev_count_BOCN = get_vocab(dev_text, r
         14                      min_df=10, keep_topN=100,
         15                      stop_words=stop_words,char_ngrams=True)
         16
         17  # #Test Vectorisation
         18  # print("Vectorise test text....","\n")
         19  # test_count, test_vect = vectorise(test_text, test_vocab)
         20
```

```python
21  #Train Vectorisation
22  print("Vectorise train text....","\n")
23  train_count_BOCN, train_vect_BOCN = vectorise(train_text, train_voc
24
25  #Dev Vectorisation
26  print("Vectorise dev text....","\n")
27  dev_count_BOCN, dev_vect_BOCN = vectorise(dev_text, dev_vocab_BOCN)
28
29
30  weights_BOCN, training_loss_history_BOCN, validation_loss_history_E
31
32
33
34  X_te_count = train_count_BOCN
35
36  w_count = weights_BOCN
37
38  preds_te_count = predict_class(X_te_count, w_count)
39
40  # train_count, weights
41
42
43  Y_te = dev_count_BOCN
44
45  print('Accuracy:', accuracy_score(Y_te,preds_te_count))
46  print('Precision:', precision_score(Y_te,preds_te_count))
47  print('Recall:', recall_score(Y_te,preds_te_count))
48  print('F1-Score:', f1_score(Y_te,preds_te_count))
49
50  # training_loss_history
51  # validation_loss_history
52
53  # plt.figure(figsize=(25,6))
54
55  plt.title('Cost Function Slope')
56  plt.plot(training_loss_history_BOCN, label='Training Loss History')
57  plt.plot(validation_loss_history_BOCN, label='Validation Loss Histc
58  plt.legend(prop={'size': 16})
59  plt.xlabel('Number of Iterations')
60  plt.ylabel('Error Values')
61  plt.show()
62
63  # print("Test Dictionary >>> ", "\n")
64  # id_w_test, w_id_test = create_2dict(test_df)
65
66  print("Train Dictionary >>> ", "\n")
67  id_w_train_BOCN, w_id_train_BOCN =create_2dict(train_df_BOCN)
68
69  print("DEV Dictionary >>> ", "\n")
70  id_w_dev_BOCN, w_id_dev_BOCN =create_2dict(dev_df_BOCN)
71
72
73
74  top_neg = w_count.argsort()[:10]
75  for i in top_neg:
76  #     print(id2word[i])
77      print(id_w_train_BOCN[i])
78
79  top_pos = w_count.argsort()[::-1][:10]
80  for i in top_pos:
```

```
81  #      print(id2word[i])
82       print(id_w_train_BOCN[i])
```

## BOCN-tfidf:

```python
In [ ]:  1  ############ BOCN-tfidf ############
         2
         3  # #TEST
         4  # test_vocab, test_df, test_count = get_vocab(test_text, ngram_rang
         5  #                     min_df=2, keep_topN=500,
         6  #                     stop_words = stop_words,char_ngrams=True)
         7
         8  #train
         9  # train_vocab, train_df, train_count = get_vocab(train_text, ngram_
        10  #                     min_df=10, keep_topN=100,
        11  #                     stop_words=stop_words,char_ngrams=True)
        12  # #Dev
        13  # dev_vocab, dev_df, dev_count = get_vocab(dev_text, ngram_range=(1
        14  #                     min_df=10, keep_topN=100,
        15  #                     stop_words=stop_words,char_ngrams=True)
        16
        17  # #Test Vectorisation
        18  # print("Vectorise test text....","\n")
        19  # test_count, test_vect = vectorise(test_text, test_vocab)
        20
        21  #Train Vectorisation
        22  # print("Vectorise train text....","\n")
        23  # train_count, train_vect = vectorise(train_text, train_vocab)
        24
        25  # #Dev Vectorisation
        26  # print("Vectorise dev text....","\n")
        27  # dev_count, dev_vect = vectorise(dev_text, dev_vocab)
        28
        29
        30  weights_BOCN_tfidf, training_loss_history_BOCN_tfidf, validation_lo
        31
        32
        33
        34  X_te_count = train_vect_BOCN_tfidf
        35
        36  w_count = weights_BOCN_tfidf
        37
        38  preds_te_count = predict_class(X_te_count, w_count)
        39
        40  # train_count, weights
        41
        42
        43  Y_te = dev_vect_BOCN
        44
        45  print('Accuracy:', accuracy_score(Y_te,preds_te_count))
        46  print('Precision:', precision_score(Y_te,preds_te_count))
        47  print('Recall:', recall_score(Y_te,preds_te_count))
        48  print('F1-Score:', f1_score(Y_te,preds_te_count))
        49
        50  # training_loss_history
        51  # validation_loss_history
        52
        53  plt.figure(figsize=(25,6))
        54
```

```
55 plt.title('Cost Function Slope')
56 plt.plot(training_loss_history_BOCN_tfidf, label='Training Loss His
57 plt.plot(validation_loss_history_BOCN_tfidf, label='Validation Loss
58 plt.legend(prop={'size': 16})
59 plt.xlabel('Number of Iterations')
60 plt.ylabel('Error Values')
61 plt.show()
62
63 # print("Test Dictionary >>> ", "\n")
64 # id_w_test, w_id_test = create_2dict(test_df)
65
66 # print("Train Dictionary >>> ", "\n")
67 # id_w_train, w_id_train =create_2dict(train_df)
68
69 # print("DEV Dictionary >>> ", "\n")
70 # id_w_dev, w_id_dev =create_2dict(dev_df)
71
72
73
74 top_neg = w_count.argsort()[:10]
75 for i in top_neg:
76 #     print(id2word[i])
77     print(id_w_train_BOCN_tfidf[i])
78
79 top_pos = w_count.argsort()[::-1][:10]
80 for i in top_pos:
81 #     print(id2word[i])
82     print(id_w_train_BOCN_tfidf[i])
```

## BOW+BOCN:

In [ ]:

## Full Results

Add here your results:

| LR | Precision | Recall | F1-Score |
|---|---|---|---|
| BOW-count | | | |
| BOW-tfidf | | | |
| BOCN-count | | | |
| BOCN-tfidf | | | |
| BOW+BOCN | | | |

Please discuss why your best performing model is better than the rest.

In [ ]:

In [ ]: