**Homework 1**

To make the file handling easier, every program receives two parameters via command line, the input and output file. In this way, it makes faster and easier to define which file is going to be used and where the output would be stored.

I decided to use python as it has many libraries that makes reading and writing in files, decoding, encoding and making requests easier.

The first step was to parse the data was to encode the data to a supported encoding, the parserCsvToJson.py program encoded the text of the tweets to latin-1 encoding, stored the class as a boolean 'isBot' and dumped the result into a file with json format. It also checked for errors and stored them in a different file. Three tweets were deleted from the data because they were empty strings.

```
carlos@carlos-Lenovo-Y50-70:~/Projects/BotsIdentification$ python3 parserCsvToJson.py Example/train.csv Example/train.json
```

```
train.csv                                              {} train.json  ×
  1   Tweet,Class                                       1   [
  2   RT @CHlLDHOODRUINER: basically how my life is     2       {
      going http://t.co/gIHzuq                          3           "id": 0,
  3   @bripriscilla No. I'm your sisters sister. Not    4           "isBot": false,
      yours lol,0                                        5           "text": "RT @CHlLDHOODRUINER: basically how my life is goin
  4   Cute! The dog looks just like my dog! Watch:      6       },
      Seal Makes A New Friend h                          7       {
  5   RT @WLYeung: @fryan Checked vice.cn. Its link     8           "id": 1,
      to Youku is blocked. Dis                           9           "isBot": false,
  6   i should get my copy of smash bros by the 10th.  10           "text": "@bripriscilla No. I'm your sisters sister. Not you
      :0,0                                              11       },
  7   @jordanbks so glad you enjoyed the play!!,0       12       {
  8   @aenertia that would be massively expensive and  13           "id": 2,
      probably environmental                            14           "isBot": false,
  9   @normanisbeyonce me,0                             15           "text": "Cute! The dog looks just like my dog! Watch: Seal
 10   @markyeg @mushion22 #SizeQueen,0                  16       },
 11   https://t.co/QJMGz1ZDOZGender equality.,0         17       {
                                                        18           "id": 3,
                                                        19           "isBot": false,
                                                        20           "text": "RT @WLYeung: @fryan Checked vice.cn. Its link to Y
                                                        21       },
                                                        22       {
```

Once the data was parsed to a json format and encoded to latin-1, the meaning.py program made a request for each tweet to the meaning cloud API. It waited .2 seconds between each request and stored the errors so I could check what was failing. The program printed the count of tweets processed each 50 tweets and write them on a file in order to record every advance.

```
carlos@carlos-Lenovo-Y50-70:~/Projects/BotsIdentification$ python3 meaning.py Data/trainws.json Data/new.json
Total:   0
Total:   50
Total:   100
Total:   150
Total:   200
Total:   250
Total:   300
Total:   350
Total:   400
Total:   450
Total:   500
Total:   550
Total:   600
Total:   650
```

It also translated the Polarity value to an int value, as it is a qualitative ordinal value,  the scale goes from -2 to 2 where the negative values represent negative polarity just as the positive values represent positive polarities and the 0 is neutral, confidence was stores as provided by the API. The other three values where stored as boolean values: isAgreement, isSubjective and isIronic.

```
{} trainws.json ✕
 1   [
 2       {
 3           "id": 0,
 4           "isBot": true,
 5           "text": "i remember in 7th grade we used to fight the 8th graders at recess . dem were the days man lol",
 6           "isAgreement": true,
 7           "isSubjective": true,
 8           "isIronic": false,
 9           "confidence": 100,
10           "polarity": 1
11       },
12       {
13           "id": 1,
14           "isBot": true,
15           "text": "You cannot perform in a manner inconsistent with the way you see yourself. - Zig Ziglar",
16           "isAgreement": true,
17           "isSubjective": false,
18           "isIronic": false,
19           "confidence": 92,
20           "polarity": 1
21       },
```

The emotion.py program uses indico python library to request the 5 emotion values, it makes a call to the api for each tweet, also writes them in a file every 100 tweets and again at the end.

```
carlos@carlos-Lenovo-Y50-70:~/Projects/BotsIdentification$ python emotion.py Data/trainws.json Data/trainwsae.json
```

The values of the emotion API where stored as provided.

```
{} trainwsae.json ✕
 1   [
 2       {
 3           "polarity": 1,
 4           "confidence": 100,
 5           "text": "i remember in 7th grade we used to fight the 8th graders at recess . dem were the days man lol",
 6           "sadness": 0.7035837173461914,
 7           "isIronic": false,
 8           "isBot": true,
 9           "isSubjective": true,
10           "isAgreement": true,
11           "anger": 0.1284189075231552,
12           "surprise": 0.008168450556695461,
13           "joy": 0.04591129720211029,
14           "id": 0,
15           "fear": 0.11391761898994446
16       },
17       {
18           "polarity": 1,
19           "confidence": 92,
20           "text": "You cannot perform in a manner inconsistent with the way you see yourself. - Zig Ziglar",
21           "sadness": 0.19053326547145844,
22           "isIronic": false,
```

Finally, the parserJsonToCsv.py program re-formats the data from json to csv.



 The trainwsae.json and trainwsae.csv contain the tweets with the new 10 features.