# ▾ Download Data

```
1  !wget http://www2.aueb.gr/users/ion/data/lingspam_public.tar.gz
2  !tar -zxf lingspam_public.tar.gz

--2021-11-01 19:27:09--  http://www2.aueb.gr/users/ion/data/lingspam_public.tar.gz
Resolving www2.aueb.gr (www2.aueb.gr)... 195.251.255.138
Connecting to www2.aueb.gr (www2.aueb.gr)|195.251.255.138|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 11564714 (11M) [application/x-gzip]
Saving to: 'lingspam_public.tar.gz'

lingspam_public.tar 100%[===================>]  11.03M  1.73MB/s    in 7.0s

2021-11-01 19:27:17 (1.59 MB/s) - 'lingspam_public.tar.gz' saved [11564714/11564714]
```

# ▾ Feature selection using information gain(IG) matrix

```
1  import os
2  import re
3  import math
4  import numpy as np
5
6  from collections import Counter
7
8  # Copy from readme.txt
9  # Each one of the 10 subdirectories contains both spam and legitimate
10  # messages, one message in each file. Files whose names have the form
11  # spmsg*.txt are spam messages. All other files are legitimate messages.
12
13  path = "lingspam_public/lemm_stop"
14  word_filter = re.compile(r'^[a-z]+-?[a-z]+[0-9]*$')
15
16  spam_words = []
17  ham_words = []
18  train_X = []
19  train_Y = []
20  spam_num = 0
21  ham_num = 0
22
23  # Build IG and training dataset
24  for dir in [os.path.join(path, 'part' + str(i)) for i in range(1, 10)]:
25      mails = [(os.path.join(dir, file_name), 'spmsg' in file_name) for file_name in os.li
26
27      for file_path, spam in mails:
```

```
28             with open(file_path) as f:
29                 content = f.readlines()[2]
30                 words = content.split()
31                 filtered_words = list(set([w for w in words if (len(w) > 1 and re.match(word
32                 # Add words to training set
33                 train_X.append(filtered_words)
34                 # Label spam as 1, ham as 0
35                 if spam:
36                     train_Y.append(1)
37                     spam_words += filtered_words
38                     spam_num += 1
39                 else:
40                     train_Y.append(0)
41                     ham_words += filtered_words
42                     ham_num += 1
43
44
45    train_all_counter = Counter(spam_words + ham_words)
46    train_spam_counter = Counter(spam_words)
47    train_ham_counter = Counter(ham_words)
48
49    print(f'Loading {spam_num + ham_num} emails, {spam_num} of them are spam, \
50    and {ham_num} of them are ham.')
51    print(f'Within the total of {len(train_all_counter)} words, \
52    getting {len(train_spam_counter)} words occur in spam, \
53    and {len(train_ham_counter)} words occur in ham email.')
```

```
Loading 2602 emails, 432 of them are spam, and 2170 of them are ham.
Within the total of 48583 words, getting 8657 words occur in spam, and 45091 words occu
```

◄ ▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓ ►

```
1 print(train_X[0])
```

```
['letter', 'establish', 'convert', 'participate', 'start', 'customer', 'simple', 'inter
```

◄ ▓ ►

```
1 # Build test data set
2 test_X = []
3 test_Y = []
4
5 test_path = 'lingspam_public/lemm_stop/part10'
6 test_mail_path = [(os.path.join(test_path, file_name), 'spmsg' in file_name) for file_name
7 for file_path, spam in test_mail_path:
8     with open(file_path) as f:
9         content = f.readlines()[2]
10        words = content.split()
11        filtered_words = list(set([w for w in words if (len(w) > 1 and re.match(word_filte
12        # Add words to test set
13        test_X.append(filtered_words)
14        # Label spam as 1, ham as 0
```

```
15          if spam:
16              test_Y.append(1)
17          else:
18              test_Y.append(0)
```

```
1 # Helper parameter and function for IG
2 def Log(i):
3     return np.log(i) if i != 0 else 0
4
5
6 total_num = spam_num + ham_num
7 P_SPAM = spam_num / total_num
8 P_HAM = 1 - P_SPAM
9 HC = -P_SPAM * Log(P_SPAM) -P_HAM * Log(P_HAM)
```

```
1 # Generating IG for all words
2
3 IG = {}
4
5 for word in train_all_counter:
6     all_occur = train_all_counter[word]
7     spam_occur = train_spam_counter[word]
8     ham_occur = train_ham_counter[word]
9
10    # X=1, C = ham
11    P1h = (ham_occur / ham_num) * P_HAM
12    # X=0, C = ham
13    P0h = (1 - (ham_occur / ham_num)) * P_HAM
14    # X=1, C = spam
15    P1s = (spam_occur / spam_num) * P_SPAM
16    # X=0, C = spam
17    P0s = (1 - (spam_occur / spam_num)) * P_SPAM
18    # P(X=1)
19    P1 = all_occur / total_num
20    # P(X=0)
21    P0 = 1 - P1
22
23    # Note: H(C|X) = - sum(P(X,C) log(C|X))
24    # But it is easier to do addition operation here :)
25    HCX = P1h*Log(P1h/P1) + P0h*Log(P0h/P0) + P1s*Log(P1s/P1) + P0s*Log(P0s/P0)
26    IG[word] = HC + HCX
27
28 top_words = [k for k, v in sorted(IG.items(), key=lambda item: -item[1])]
29 top_10 = top_words[:10]
30 top_100 = top_words[:100]
31 top_1000 = top_words[:1000]
```

```
1 print(top_10)
```

```
['language', 'remove', 'free', 'linguistic', 'university', 'money', 'click', 'market',
```

# ▾ Classifiers Implementation

# ▾ Generate Feature Matrix

```
1   def get_feature(N, dataset, term_frequency=False):
2       """Generate feature matrix based on top N words."""
3       # print(term_frequency)
4       top_n = {v:i for i,v in enumerate(top_words[:N])}
5       feature_matrix = np.zeros((len(dataset), N))
6       for j, content in enumerate(dataset):
7           for word in content:
8               if word in top_n:
9                   if term_frequency:
10                      feature_matrix[j, top_n[word]] += 1
11                  else:
12                      feature_matrix[j, top_n[word]] = 1
13
14      return feature_matrix
```

# ▾ Bernoulli NB Classifier

```
1 class BNB:
2     """Naive Bayes classifier for multivariate Bernoulli models."""
3     def __init__(self, alpha=1.0):
4         self.alpha = alpha  # Actually, no need in this homework :)
5         self.spam_num = 0
6         self.ham_num = 0
7         self.P_X_spam = []
8         self.P_X_ham = []
9
10
11    def fit(self, train_x, train_y):
12        """Fit Naive Bayes classifier according to train_x, train_y."""
13        self.spam_num = sum(train_y)
14        self.ham_num = len(train_y) - sum(train_y)
15        rows, feature_num = train_x.shape
16
17        self.P_X_spam = np.zeros(feature_num)
18        self.P_X_ham = np.zeros(feature_num)
19        for i in range(feature_num):
```

```
20                x_spam = 0
21                x_ham = 0
22                for j in range(rows):
23                    if train_x[j,i]:
24                        if train_y[j]:
25                            x_spam += 1
26                        else:
27                            x_ham += 1
28            self.P_X_spam[i] = (1 + x_spam) / (self.spam_num + 2)
29            self.P_X_ham[i] = (1 + x_ham) / (self.ham_num + 2)
30
31
32    def _row_predict(self, one_sample):
33        """Predict one record of the data set."""
34        prob = 1
35        for i, feature in enumerate(one_sample):
36            p0 = self.P_X_ham[i]
37            p1 = self.P_X_spam[i]
38            prob *= (p1 / p0) if feature == 1 else (1 - p1) / (1 - p0)
39        return prob * self.spam_num / self.ham_num
40
41
42    def predict(self, test_x):
43        """Perform classification on an array of test vectors test_x."""
44        pred = np.zeros((test_x.shape[0]))
45        for i in range(test_x.shape[0]):
46            pred[i] = int(self._row_predict(test_x[i]) > 1)
47        return pred
48
49
50    def score(self, test_x, test_y):
51        """Return the mean accuracy on the given test data and labels."""
52        y_pred = self.predict(test_x)
53        count = 0
54        for i in range(len(test_y)):
55            count += 1 if y_pred[i] == test_y[i] else 0
56        return count / len(test_y)
57
```

# ▾ Multinominal NB Classifier

```
1 class MNB:
2     """Naive Bayes classifier for multinomial models."""
3     def __init__(self, alpha=1.0):
4         self.alpha = alpha
5         self.spam_num = 0
6         self.ham_num = 0
7         self.P_X_spam = []
8         self.P_X_ham = []
```

```python
 9
10
11     def fit(self, train_x, train_y):
12         """Fit Naive Bayes classifier according to train_x, train_y."""
13         self.spam_num = sum(train_y)
14         self.ham_num = len(train_y) - sum(train_y)
15         rows, feature_num = train_x.shape
16
17         self.P_X_spam = np.zeros(feature_num)
18         self.P_X_ham = np.zeros(feature_num)
19         for i in range(feature_num):
20             x_spam = 0
21             x_ham = 0
22             for j in range(rows):
23                 if train_x[j,i]:
24                     if train_y[j]:
25                         x_spam += 1
26                     else:
27                         x_ham += 1
28             self.P_X_spam[i] = (1 + x_spam) / (self.spam_num + 2)
29             self.P_X_ham[i] = (1 + x_ham) / (self.ham_num + 2)
30
31
32     def _row_predict(self, one_sample):
33         """Predict one record of the data set."""
34         prob = 1
35         for i, feature in enumerate(one_sample):
36             p0 = self.P_X_ham[i]
37             p1 = self.P_X_spam[i]
38             prob += (Log(p1) - Log(p0)) if feature == 1 else 0
39         return prob + Log(self.spam_num) - Log(self.ham_num)
40
41
42     def predict(self, test_x):
43         """Perform classification on an array of test vectors test_x."""
44         pred = np.zeros((test_x.shape[0]))
45         for i in range(test_x.shape[0]):
46             pred[i] = int(self._row_predict(test_x[i]) > 0)
47         return pred
48
49
50     def score(self, test_x, test_y):
51         """Return the mean accuracy on the given test data and labels."""
52         y_pred = self.predict(test_x)
53         count = 0
54         for i in range(len(test_y)):
55             count += 1 if y_pred[i] == test_y[i] else 0
56         return count / len(test_y)
```

# ▾ Classifiers Comparison

```python
1  from sklearn.naive_bayes import MultinomialNB, BernoulliNB
2  from sklearn.metrics import classification_report
3
4  top_N = [10, 100, 1000]
5  result_accuracy = []
6
7  for N in top_N:
8      Xprint = get_feature(N, train_X)
9      # X label with term frequency
10     X_tf = get_feature(N, train_X, term_frequency=True)
11     Y = np.array(train_Y)
12
13     X_test = get_feature(N, test_X)
14     # X label with term frequency
15     X_test_tf = get_feature(N, test_X, term_frequency=True)
16     Y_test = np.array(test_Y)
17
18     model1 = BNB()
19     model2 = MNB()
20     model3 = MNB()
21
22     # train models
23     model1.fit(X, Y)
24     model2.fit(X, Y)
25     model3.fit(X_tf, Y)
26
27     # predict
28     y_1 = model1.predict(X_test)
29     y_2 = model2.predict(X_test)
30     y_3 = model3.predict(X_test_tf)
31
32     print(f'--------------- Using top {N} features --------------')
33     print( '--------------- Bernoulli Naive Bayes ---------------')
34     print(classification_report(Y_test,y_1, target_names = ["ham", "spam"]))
35     print('--- Multinominal Naive Bayes with binary features ---')
36     print(classification_report(Y_test,y_2, target_names = ["ham", "spam"]))
37     print('---- Multinominal Naive Bayes with term frequency ---')
38     print(classification_report(Y_test,y_3, target_names = ["ham", "spam"]))
39
40     N_accuray = [model1.score(X_test,Y_test),
41                  model2.score(X_test,Y_test),
42                  model3.score(X_test_tf,Y_test)]
43     result_accuracy.append(N_accuray)
```

```
--------------- Using top 10 features --------------
--------------- Bernoulli Naive Bayes ---------------
              precision    recall  f1-score   support

         ham       0.96      0.98      0.97       242
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| spam         | 0.89      | 0.82   | 0.85     | 49      |
| accuracy     |           |        | 0.95     | 291     |
| macro avg    | 0.93      | 0.90   | 0.91     | 291     |
| weighted avg | 0.95      | 0.95   | 0.95     | 291     |

--- Multinominal Naive Bayes with binary features ---

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ham          | 0.98      | 0.95   | 0.97     | 242     |
| spam         | 0.80      | 0.92   | 0.86     | 49      |
| accuracy     |           |        | 0.95     | 291     |
| macro avg    | 0.89      | 0.94   | 0.91     | 291     |
| weighted avg | 0.95      | 0.95   | 0.95     | 291     |

---- Multinominal Naive Bayes with term frequency ---

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ham          | 0.98      | 0.95   | 0.97     | 242     |
| spam         | 0.80      | 0.92   | 0.86     | 49      |
| accuracy     |           |        | 0.95     | 291     |
| macro avg    | 0.89      | 0.94   | 0.91     | 291     |
| weighted avg | 0.95      | 0.95   | 0.95     | 291     |

-------------- Using top 100 features --------------
-------------- Bernoulli Naive Bayes --------------

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ham          | 0.94      | 1.00   | 0.97     | 242     |
| spam         | 1.00      | 0.67   | 0.80     | 49      |
| accuracy     |           |        | 0.95     | 291     |
| macro avg    | 0.97      | 0.84   | 0.89     | 291     |
| weighted avg | 0.95      | 0.95   | 0.94     | 291     |

--- Multinominal Naive Bayes with binary features ---

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ham          | 1.00      | 0.98   | 0.99     | 242     |
| spam         | 0.89      | 1.00   | 0.94     | 49      |
| accuracy     |           |        | 0.98     | 291     |
| macro avg    | 0.95      | 0.99   | 0.96     | 291     |
| weighted avg | 0.98      | 0.98   | 0.98     | 291     |

---- Multinominal Naive Bayes with term frequency ---

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ham          | 1.00      | 0.98   | 0.99     | 242     |
| spam         | 0.89      | 1.00   | 0.94     | 49      |
| accuracy     |           |        | 0.98     | 291     |

▼ Spam precision and spam recall

| Params | Spam Precision | Spam Recall |
|---|---|---|
| Top 10 Words, Bernoulli Naive Bayes | 0.89 | 0.82 |
| Top 10 Words, Multinominal Naive Bayes with binary features | 0.80 | 0.92 |
| Top 10 Words, Multinominal Naive Bayes with term frequency | 0.80 | 0.92 |
| Top 100 Words, Bernoulli Naive Bayes | 1.00 | 0.67 |
| Top 100 Words, Multinominal Naive Bayes with binary features | 0.89 | 1.00 |
| Top 100 Words, Multinominal Naive Bayes with term frequency | 0.89 | 1.00 |
| Top 1000 Words, Bernoulli Naive Bayes | 1.00 | 0.61 |
| Top 1000 Words, Multinominal Naive Bayes with binary features | 1.00 | 1.00 |
| Top 1000 Words, Multinominal Naive Bayes with term frequency | 1.00 | 1.00 |

## Compare with sklearn

```
1 result_accuracy_sk = []
2
3 for N in top_N:
4     model1_sk = BernoulliNB()
5     model2_sk = MultinomialNB()
6     model3_sk = MultinomialNB()
7
8     X = get_feature(N, train_X)
9     X_tf = get_feature(N, train_X, term_frequency=True)
10    Y = np.array(train_Y)
11
12    X_test = get_feature(N, test_X)
13    X_test_tf = get_feature(N, test_X, term_frequency=True)
14    Y_test = np.array(test_Y)
15
16    # train models
17    model1_sk.fit(X, Y)
18    model2_sk.fit(X, Y)
19    model3_sk.fit(X_tf, Y)
20
21    # predict
22    N_accuray = [model1_sk.score(X_test,Y_test),
23                 model2_sk.score(X_test,Y_test),
24                 model3_sk.score(X_test_tf,Y_test)]
25    result_accuracy_sk.append(N_accuray)
26
27
28 print(result_accuracy)
29 print(result_accuracy_sk)
```

```
[[0.9518900343642611, 0.9484536082474226, 0.9484536082474226], [0.9450171821305842, 0.9
[[0.9518900343642611, 0.9518900343642611, 0.9518900343642611], [0.9450171821305842, 0.9
```

◀         ▭                                                                                        ▶

# ▾ SVM based spam filter

```
1  from sklearn.model_selection import GridSearchCV
2  from sklearn.model_selection import train_test_split
3  from sklearn.svm import SVC
4  from sklearn.preprocessing import MinMaxScaler
5  from sklearn.base import clone
6
7
8  feature_sizes = [10, 100, 200]
9  modes = ["bf", "tf"]
10 gamma = [1, 0.1, 0.01, 0.001]
11 C = [0.1, 1, 10, 100]
12 degree = [2,3,10]
13 tuned_parameters = [{'kernel': ['rbf','sigmoid'], 'gamma': gamma,'C': C},
14                     {'kernel': ['linear'], 'C': C},
15                     {'kernel': ['poly'], 'gamma': gamma, "degree" : degree, 'C':C}
16                     ]
17 now_best = {}
18 now_score = 0
19 best_clf = SVC()
20 for N in feature_sizes:
21     for mode in modes:
22         X = get_feature(N, train_X, term_frequency=(mode=='tf'))
23         X_test = get_feature(N, test_X, term_frequency=(mode=='tf'))
24         Y = np.array(train_Y)
25         if mode == 'tf':
26             ss = MinMaxScaler()
27             X = ss.fit_transform(X)
28             X_test = ss.transform(X_test)
29         clf = GridSearchCV(SVC(), tuned_parameters, cv = 5, scoring='recall_macro', verbos
30         clf.fit(X, Y)
31
32         if clf.best_score_ > now_score:
33             now_score = clf.best_score_
34             now_best = {"size":N, "mode":mode, "score":clf.best_score_, "params":clf.best_
35             best_clf = clone(clf.best_estimator_)
36             print (now_best)
37
38 N = now_best["size"]
39 mode = now_best["mode"]
40 X = get_feature(N, train_X, term_frequency=(mode=='tf'))
41 X_test = get_feature(N, test_X, term_frequency=(mode=='tf'))
42 Y = np.array(train_Y)
43 Y_test = np.array(test_Y)
44
45 if mode == "tf":
```

```
46      ss = MinMaxScaler()
47      X = ss.fit_transform(X)
48      X_test = ss.transform(X_test)
49
50 best_clf.fit(X, Y)
51 pred_y = best_clf.predict(X_test)
52 print (now_best)
53 print (classification_report(Y_test, pred_y, target_names = ["ham", "spam"]))
```

```
    {'size': 10, 'mode': 'bf', 'score': 0.9333773098447026, 'params': {'C': 1, 'gamma': 1,
    {'size': 100, 'mode': 'bf', 'score': 0.9611211300362037, 'params': {'C': 1, 'gamma': 0.
    {'size': 200, 'mode': 'bf', 'score': 0.9699466122688307, 'params': {'C': 10, 'gamma': 0
    {'size': 200, 'mode': 'bf', 'score': 0.9699466122688307, 'params': {'C': 10, 'gamma': 0
                  precision    recall  f1-score   support

             ham       0.98      1.00      0.99       242
            spam       0.98      0.88      0.92        49

        accuracy                           0.98       291
       macro avg       0.98      0.94      0.96       291
    weighted avg       0.98      0.98      0.98       291
```

# Methodology

I use `GridSearchCV` to help me find the best parameters for SVM. Also, we've see that higher Information Gain feature matrix can have a better result in prediction. So I compare the score between "BF" and "TF" and the score among "N={10, 100, 200}". The result shows that using **top 200** words, **binary features**, and `{'C': 10, 'gamma': 0.1, 'kernel': 'rbf'}` as the SVM parameter can get the best result.

# ▾ Adversarial Classification

# ▾ Attacker

```
1 class Attacker:
2     def __init__(self):
3         self.Lo = []
4         self.Lo0 = []
5         self.Lo1 = []
6
7
8     def cal_Lo(self, train_x, train_y):
9         rows, feature_num = train_x.shape
```

```
10            self.Lo = np.zeros(feature_num)
11            self.Lo0 = np.zeros(feature_num)
12            self.Lo1 = np.zeros(feature_num)
13
14            for i in range(feature_num):
15                x_spam = 0
16                x_ham = 0
17                for j in range(rows):
18                    if train_x[j,i]:
19                        if train_y[j]:
20                            x_spam += 1
21                        else:
22                            x_ham += 1
23                P_X1_ham = x_ham / ham_num
24                P_X1_spam = x_spam / spam_num
25                P_X0_ham = 1 - P_X1_ham
26                P_X0_spam = 1 - P_X1_spam
27                self.Lo1[i] = Log(P_X1_spam/P_X1_ham)
28                self.Lo0[i] = Log(P_X0_spam/P_X0_ham)
29                self.Lo[i] = self.Lo1[i] - self.Lo0[i]
30
31
32    def _row_attack(self, one_sample):
33        sort_Lo = sorted(range(len(self.Lo)), key=lambda i: self.Lo[i])
34
35        cost = 0
36        sigmaLoX = 0
37        for i,v in enumerate(one_sample):
38            sigmaLoX += self.Lo1[i] if v==1 else self.Lo0[i]
39        if sigmaLoX < 0:
40            return 0, one_sample
41
42        for i in sort_Lo:
43            if(self.Lo[i] >= 0):
44                return cost, one_sample
45            if one_sample[i]:
46                continue
47            one_sample[i] = 1
48            sigmaLoX += self.Lo[i]
49            cost += 1
50            if (sigmaLoX < 0):
51                return cost, one_sample
52
53
54    def attack(self, test_x, test_y):
55        costs = []
56        new_test_x = []
57        rows, feature_num = test_x.shape
58        for i in range(rows):
59            if test_y[i]:
60                cost, changed = self._row_attack(test_x[i])
```

```
61                    new_test_x.append(changed)
62                    costs.append(cost)
63                else:
64                    new_test_x.append(test_x[i])
65
66        return np.mean(costs), np.array(new_test_x)
67
```

## ▾ Defender

```
1 class Defender:
2     def __init__(self):
3         self.P_X1_ham = []
4         self.P_X1_spam = []
5         self.P_X0_ham = []
6         self.P_X0_spam = []
7         self.Lo = []
8         self.Lo0 = []
9         self.Lo1 = []
10        self.result = []
11
12    def fit(self, train_x, train_y):
13        rows, feature_num = train_x.shape
14        self.P_X1_ham = np.zeros(feature_num)
15        self.P_X1_spam = np.zeros(feature_num)
16        self.P_X0_ham = np.zeros(feature_num)
17        self.P_X0_spam = np.zeros(feature_num)
18        self.Lo = np.zeros(feature_num)
19        self.Lo0 = np.zeros(feature_num)
20        self.Lo1 = np.zeros(feature_num)
21
22        for i in range(feature_num):
23            x_spam = 0
24            x_ham = 0
25            for j in range(rows):
26                if train_x[j,i]:
27                    if train_y[j]:
28                        x_spam += 1
29                    else:
30                        x_ham += 1
31            self.P_X1_ham[i] = x_ham / ham_num
32            self.P_X1_spam[i] = x_spam / spam_num
33            self.P_X0_ham[i] = 1 - self.P_X1_ham[i]
34            self.P_X0_spam[i] = 1 - self.P_X1_spam[i]
35            self.Lo1[i] = Log(self.P_X1_spam[i] / self.P_X1_ham[i])
36            self.Lo0[i] = Log(self.P_X0_spam[i] / self.P_X0_ham[i])
37            self.Lo[i] = self.Lo1[i] - self.Lo0[i]
38
```

```
39
40     def _row_attack(self, one_sample):
41         sort_Lo = sorted(range(len(self.Lo)), key=lambda i: self.Lo[i])
42
43         cost = 0
44         sigmaLoX = 0
45         for i,v in enumerate(one_sample):
46             sigmaLoX += self.Lo1[i] if v==1 else self.Lo0[i]
47         if sigmaLoX < 0:
48             return 0, one_sample
49
50         for i in sort_Lo:
51             if(self.Lo[i] >= 0):
52                 return cost, one_sample
53             if one_sample[i]:
54                 continue
55             one_sample[i] = 1
56             sigmaLoX += self.Lo[i]
57             cost += 1
58             if (sigmaLoX < 0):
59                 return cost, one_sample
60
61
62     def defence(self, model, attacked_x):
63         for sample in attacked_x:
64             if model._row_predict(sample):
65                 self.result += [1]
66                 continue
67             new_matrix = [[]]
68             for i in sample:
69                 if i:
70                     new_row = [(x + [0]) for x in new_matrix] + [(x + [1]) for x in new_ma
71                 else:
72                     new_matrix_f = [(x + [0]) for x in new_matrix]
73                 new_matrix = new_matrix_f
74             new_matrix = [np.asarray(x) for x in new_matrix]
75             realFroms = []
76             for origin in new_matrix:
77                 _, attacked = self._row_attack(origin)
78                 if np.array_equal(attacked, sample):
79                     if model._row_predict(origin):
80                         realFroms += [origin]
81             if len(realFroms):
82                 P_spam = 0
83                 P_ham = 0
84                 for origin in realFroms:
85                     for i in origin:
86                         P_spam *= self.P_X1_spam[i] if i else self.P_X0_spam[i]
87                         P_ham *= self.P_X1_ham[i] if i else self.P_X0_ham[i]
88                     P_spam += 1
89                     P_ham += 1
```

```
90                    self.result += [0 if P_spam < P_ham else 1]
91            else:
92                    self.result += [0]
93
94        return np.array(self.result)
```

## Adversarial attack analysis

```
1 from sklearn.metrics import confusion_matrix
2
3 N = 10
4 X = get_feature(N, train_X)
5 Y = np.array(train_Y)
6 X_test = get_feature(N, test_X)
7 Y_test = np.array(test_Y)
8
9 # Use original model to predict Y_test
10 botnet = BNB()
11 botnet.fit(X, Y)
12 Y_pred_origin = botnet.predict(X_test)
13 print("Original classification report before attack")
14 print(classification_report(Y_test, Y_pred_origin, target_names = ["ham", "spam"]))
15 print("Confusion Matrix")
16 print(confusion_matrix(Y_test, Y_pred_origin))
17
18 # Assume that Attacker knows
19 # - original model's result: Y_pred_origin
20 # - original test words list: test_X
21 # - oringinal test label: test_Y
22 # Using 'Add-Words' attack
23 hacker = Attacker()
24 hacker.cal_Lo(X, Y)
25 cost, X_test_attacked = hacker.attack(X_test, Y_test)
26 Y_pred_attacked = botnet.predict(X_test_attacked)
27 print("Original classification report after attack")
28 print(classification_report(Y_test, Y_pred_attacked, target_names = ["ham", "spam"]))
29 print("Confusion Matrix")
30 print(confusion_matrix(Y_test, Y_pred_attacked))
31 print(f"Average cost = {cost}")
32
33 # Denfender
34 antispam = Defender()
35 antispam.fit(X, Y)
36 Y_pred_defence = antispam.defence(botnet, X_test_attacked)
37 print("Original classification report after defence")
38 print(classification_report(Y_test, Y_pred_defence, target_names = ["ham", "spam"]))
39 print("Confusion Matrix")
```

```
40 print(confusion_matrix(Y_test, Y_pred_defence))
41 print(f"Average cost = {cost}")
```

Original classification report before attack

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| ham | 0.96 | 0.98 | 0.97 | 242 |
| spam | 0.89 | 0.82 | 0.85 | 49 |
| accuracy |  |  | 0.95 | 291 |
| macro avg | 0.93 | 0.90 | 0.91 | 291 |
| weighted avg | 0.95 | 0.95 | 0.95 | 291 |

Confusion Matrix
```
[[237   5]
 [  9  40]]
```
Original classification report after attack

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| ham | 0.83 | 0.98 | 0.90 | 242 |
| spam | 0.29 | 0.04 | 0.07 | 49 |
| accuracy |  |  | 0.82 | 291 |
| macro avg | 0.56 | 0.51 | 0.49 | 291 |
| weighted avg | 0.74 | 0.82 | 0.76 | 291 |

Confusion Matrix
```
[[237   5]
 [ 47   2]]
```
Average cost = 1.7959183673469388
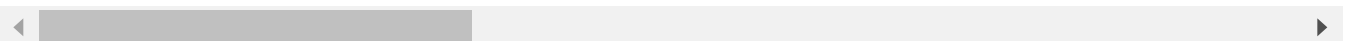Original classification report after defence

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| ham | 0.00 | 0.00 | 0.00 | 242 |
| spam | 0.17 | 1.00 | 0.29 | 49 |
| accuracy |  |  | 0.17 | 291 |
| macro avg | 0.08 | 0.50 | 0.14 | 291 |
| weighted avg | 0.03 | 0.17 | 0.05 | 291 |

Confusion Matrix
```
[[  0 242]
 [  0  49]]
```
Average cost = 1.7959183673469388
/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1272: Undefin
  _warn_prf(average, modifier, msg_start, len(result))

✓  0s      completed at 5:51 PM                                    ●  ✕