



Amirkabir University of Technology
(Tehran Polytechnic)

Report of the Project about Spam detection in Persian

Professor:

Dr. M.Akbari (akbari.ma@aut.ac.ir)

Teaching Assistants:

A.Malekzade (malekzadeh@ieee.org) , Y.Ommi(yassi.ommi@gmail.com)

Student:

Mohammadreza Ardestani (ardestani.zr@gmail.com)

5, Dec, 2020

0) **Introduction**

Phase 1) **Preprocessing the data and preparing it**

Phase 2) **Best feature selection with Chi-square method**

Phase 3) **Finding Cosine-similarity and tf-idf**

Phase 4) **Implementing KNN from scratch**

Phase 5) **Implementing Naïve Bayes Method**

Phase 6) **Calculating Precision, Recall, F1, Confusion Matrix**

6.1) Finding the measures for KKN, KKN with K-Best-Features ,& Naive Bayes

6.2) Comparing KKN & KKN with k-Best

6.3) Comparing KKN ,& Naive Bayes

(I've implemented all of **Bonus questions**)

0) Introduction

This report will be concise but thorough. If you need more detail, please contact me.

How to run:

- I. At first you should have been installed all required libraries and Jupyter notebook.
- II. Create a folder & paste 'Ass2_spamDetectionInPersianLang.ipynb' and 'stopwords.txt' and 'emails.folder'.
- III. At the address bar write "cmd" then click "Enter".
- IV. At the Cmd console type "jupyter notebook" and then "Enter".
- V. Open Jupyter source code and run each cell in order
- VI. Notice some cells take a few second to run completely

Pre analysis of Input data and deciding which part are useful for model:

Before using Test set I have removed all words that are not in trainset :

```
in cosine similarity function)

In [13]: def lenght(listOfDics):
# this funcn get a list of dictionaries and calculate total number of words in all dics
s = 0
for dic in listOfDics:
    s += len(dic)
return s

def unknownRemover(testDic ,DFA):
# we go through all dictionaries in testdic and remover unknown words
# notice testDic is a list of dics and DFA is a dic
for dic in testDic:
    CFR = [] #CandidateForRemoving
    for k,v in dic.items():
        if k not in DFA:
            CFR.append(k)
    for k in CFR:
        del dic[k]

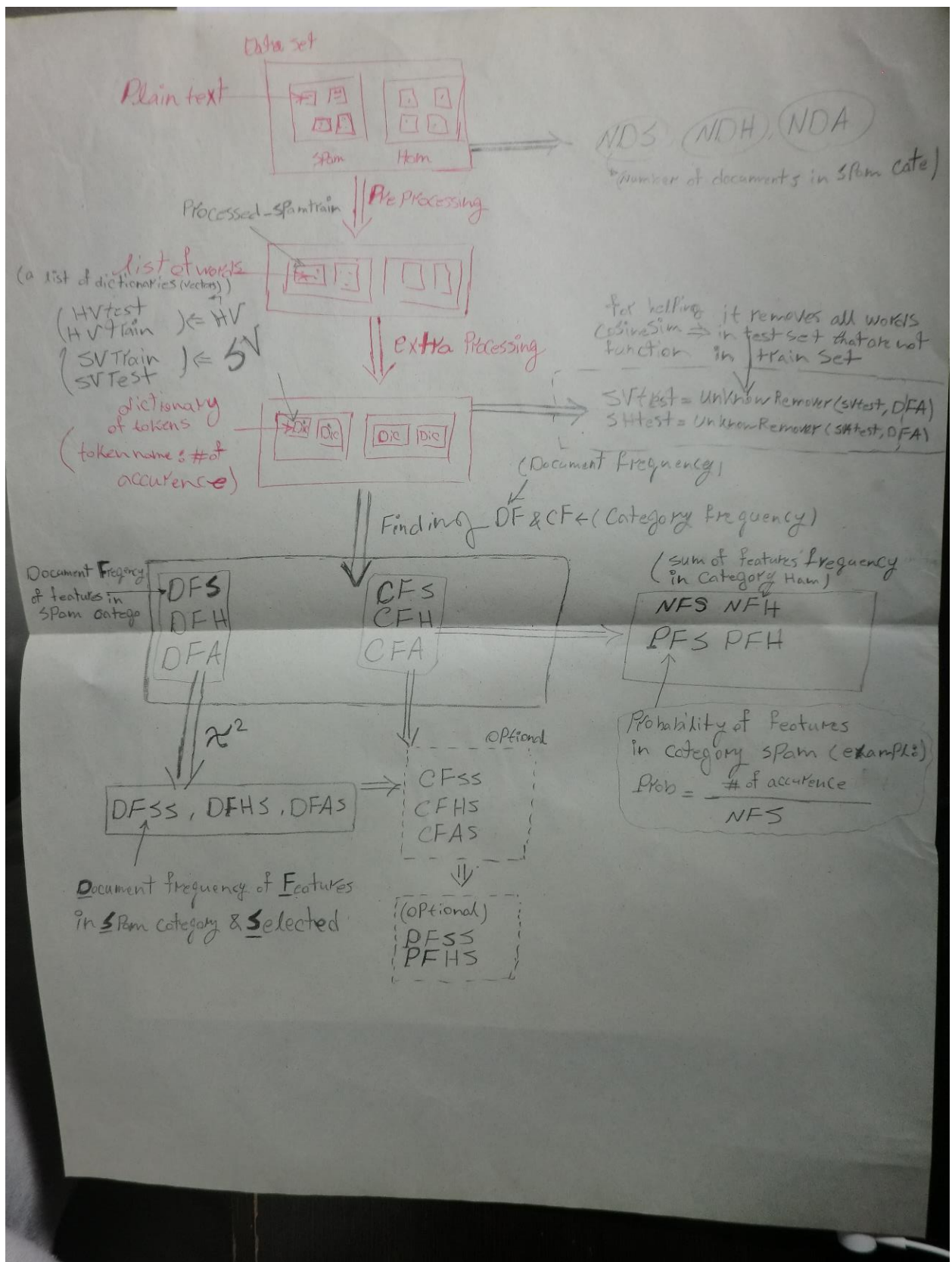
In [14]: print("total words in SVTest and SHTest before removing unknowns",lenght(SVTest),lenght(HVTest))
unknownRemover(SVTest,DFA)
unknownRemover(HVTest,DFA)
print("total words in SVTest and SHTest after removing unknowns",lenght(SVTest),lenght(HVTest))

total words in SVTest and SHTest before removing unknowns 26033 16735
total words in SVTest and SHTest after removing unknowns 23654 14783
```

And also we find out whether there is any empty email or not. And there is not any empty email.

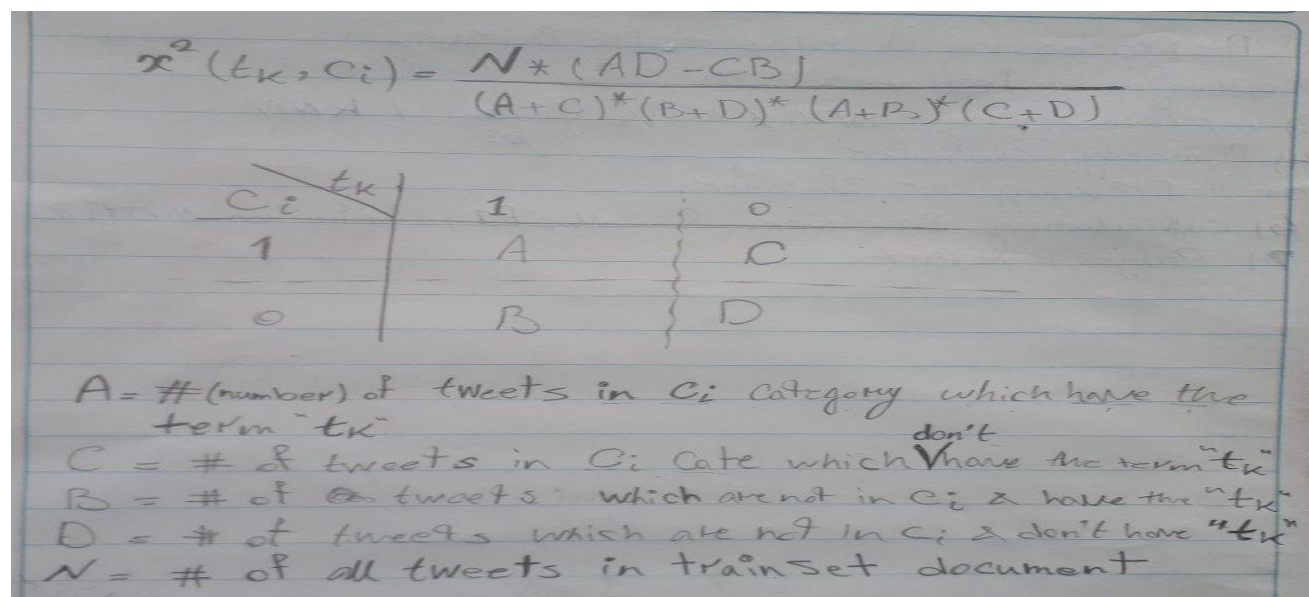
Phase 1) Preprocessing the data and preparing it

We have a long road in Preprocessing and further processing. To sum up this long story look at the following graph of preprocessing:



Phase 2) Opting the best features with CHI-SQUARE method

Following terminology describes CHI-Square method and we how to calculate it.



After we calculate Chi-square for all words, we need to pick up some of them that have more value (or more correlation).

In Final_Bag of features we select those words which have more than 0.1 Chi-square value.

In [18]: `print(important_spam_words)`

```
[('پست', 23.81), ('آموز', 23.619), ('استفاده', 18.468), ('ارسال', 10.115), ('سا', 8.952), ('کلیک', 8.322), ('د', 8.034), ('باشد', 7.99), ('افزار', 7.909), ('مو', 7.192), ('زیبا', 7.137), ('مجموعه', 6.734), ('ایمیل', 6.4), ('رو', 6.201), ('صور', 5.433), ('افزا', 5.416), ('دارا', 5.253), ('فر', 4.955), ('آموزش', 4.899), ('دکمه', 4.71), ('فرو', 4.833), ('د', 4.794), ('دستگاه', 4.789), ('منزل', 4.751), ('خده', 4.724), ('مشاهده', 4.71), ('ساع', 4.663), ('طراح', 4.65), ('العاده', 4.63), ('تجایید', 4.621), ('فوق', 4.611), ('توانید', 4.61), ('ق', 4.575), ('شامل', 4.521), ('کالا', 4.515), ('زب', 4.484), ('هزینه', 4.44), ('همراه', 4.419), ('خود', 4.39), ('بگیرید', 4.363), ('کسب', 4.363), ('رایگ', 4.318), ('تیا', 4.303), ('س', 4.3), ('سفار', 4.24), ('هدیه', 4.224), ('بالا', 4.183), ('شماره', 4.167), ('فروشگاه', 4.143), ('مناسب', 4.09), ('رنگ', 4.07), ('کودک', 4.066), ('ایس', 4.066), ('میباشد', 4.052), ('ایس', 4.028), ('اینترنت', 4.028), ('نظیر', 4.0), ('یسته', 3.997), ('پوس', 3.971), ('وب', 3.968), ('شوید', 3.961), ('فیل', 3.958), ('جذاب', 3.952), ('فروشنده', 3.899), ('عدد', 3.883), ('بپردازید', 3.883), ('کاربر', 3.877), ('دیجیتال', 3.876), ('یکج', 3.861), ('عینک', 3.846), ('الکترونیک', 3.8), ('کلیه', 3.839), ('اندا', 3.817), ('اورجینال', 3.795), ('بشتیان', 3.795), ('درب', 3.79), ('لاغر', 3.774), ('کیف', 3.773), ('بجوز', 3.766), ('شیک', 3.752), ('درآمد', 3.75), ('وزن', 3.742), ('کاه', 3.74), ('دریاف', 3.7), ('ویندوز', 3.731), ('پاور', 3.724), ('بازار', 3.716), ('سایز', 3.711), ('انواع', 3.71), ('مدل', 3.7), ('ق', 3.697), ('باز', 3.695), ('دند', 3.683), ('انلین', 3.676), ('استاندارد', 3.676), ('دیدن', 3.666), ('دهید', 3.666), ('بیشرفته', 3.663), ('طبیع', 3.66), ('مگنا', 3.65), ('یکس', 3.65), ('فارس', 3.646), ('بازاریاب', 3.64), ('موسیق', 3.636), ('بازیکن', 3.623), ('حرفه', 3.622), ('قد', 3.622), ('تناسل', 3.617), ('اتو', 3.617), ('ام', 3.617), ('روش', 3.615), ('میتوانید', 3.615), ('تعمیر', 3.61), ('بذر', 3.61), ('تیلیفات', 3.606), ('انزال', 3.604), ('وجه', 3.598), ('تضمین', 3.597), ('امور', 3.584), ('پنل', 3.584), ('تغیث', 3.58), ('تناسب', 3.5), ('ناخن', 3.578), ('بالنس', 3.578), ('دینا', 3.571), ('سانت', 3.571), ('تبل', 3.571), ('تصاویر', 3.568), ('پروژه', 3.568), ('صدا', 3.562), ('اینستایبلر', 3.559), ('حسابدار', 3.559), ('تیشز', 3.559), ('پدن', 3.555), ('ک', 3.554), ('فایل', 3.554), ('معتبر', 3.554), ('دانلود', 3.551), ('مصرف', 3.548), ('ادایم', 3.546), ('چرب', 3.546), ('هدفون', 3.546), ('گرافیک', 3.546), ('یوس', 3.545), ('خسنگ', 3.54), ('تب', 3.54), ('مزه', 3.54), ('ا', 3.537), ('مخصوص', 3.537), ('مزارتو', 3.534), ('گزینه', 3.527), ('چین', 3.521), ('اتوماتیک', 3.521), ('ام', 3.518), ('اسپر', 3.515), ('صاف', 3.515), ('متونددین', 3.515), ('شامپو', 3.515), ('صدف', 3.515), ('شارژر', 3.509), ('دکوراسیون', 3.509), ('کمیان', 3.503), ('قوشاپ', 3.503), ('عطر', 3.503), ('گن', 3.497), ('رژنا', 3.4), ('97), ('استون', 3.497), ('گلد', 3.497), ('ادارید', 3.495), ('منگا', 3.494), ('ظلا', 3.491), ('ماساز', 3.49), ('معد', 3.478), ('جاق', 3.484), ('سایت', 3.484), ('مانه', 3.484), ('تفیس', 3.484), ('انگ', 3.484), ('مزار', 3.478), ('9), ('گروه', 3.479), ('فوط', 3.478), ('دوبله', 3.478), ('دستبند', 3.478), ('صا', 3.478), ('لواز', 3.478), ('لیزر', 3.478), ('اکان', 3.472), ('سرایک', 3.472), ('یاف', 3.472), ('پد', 3.472), ('زل', 3.472), ('ارشیو', 3.472), ('چادو', 3.467), ('کانادا', 3.467), ('مکمل', 3.466), ('اشانتیون', 3.466), ('مایع', 3.466), ('امروزه', 3.466), ('3.466), ('گردنبند', 3.466), ('رادیو', 3.466)]
```

Important ham words

Phase3) Finding Cosine-similarity and tf-idf

Simply we use following formulas to find out Tfidf value or CosineSimilarity

$$Sim(A, B) = \frac{AB}{|A||B|} = \frac{x_{1A}x_{1B} + x_{2A}x_{2B} + \dots}{\sqrt{x_{1A}^2 + x_{2A}^2 + \dots} + \sqrt{x_{1B}^2 + x_{2B}^2 + \dots}} \quad (1)$$

$$Score_{tf-idf}(\hat{e}, e_i) = \sum_{i=1}^T tf(w_i, \hat{e})idf(w_i) \quad (2)$$

در عبارت ۲ مقادیر $tf(w_i, \hat{e})$ و $idf(w_i)$ به شکل زیر محاسبه می‌گردند:

$$tf(w_i, \hat{e}) = \log(\text{count}(w_i, \hat{e}) + 1) \quad (3)$$

$$idf(w_i) = \log\left(\frac{n}{df(w_i)}\right) \quad (4)$$

```
def cosine(dic10 , dic20 , selected):
    # dic1 is query and dic2 is already exist in train set and Selected is a dic of important words
    # if selected is a null dictionary, it means we should consider all words in our trainSet
    s = 0 # score
    dic1 = dic10 # prevent damaging the Train/test set documents by removing words (optimize later)
    dic2 = dic20 # prevent damaging the Train/test set documents by removing words (optimize later)
    if(selected=={}):
        # just we check size of dics
        if(len(dic1)==0 or len(dic2)==0 ):# preventing divide by 0
            #print("null doc!")
            #print("dic1",dic1)
            #print("dic2",dic2)
            return 0
    else:
        unknownRemover([dic1],selected)
        unknownRemover([dic2],selected)
        if(len(dic1)==0 or len(dic2)==0 ): # preventing divide by 0
            return 0
    norm1 = normCal(dic1) # |A|
    norm2 = normCal(dic2) # |A|
    if(len(dic1)==0 or len(dic2)==0 ): # Remove this later(i.e we shouldn't have hull dic1 or dic2 as input)
        return 0
    dotProduct = dot(dic1,dic2) # A.B
    s = (dotProduct/(norm1*norm2)) # Cosine of A&B
    return s
```

```

def tfidf(dic10 , dic20 ,selected):
    # dic1 is query and dic2 is already exist in train set and Selected is a dic of important words
    # if selected is null, it means we should consider all words in our trainSet
    s = 0 # score
    dic1 = dic10 # prevent damaging the Train/test set documents by removing words (optimize later)
    dic2 = dic20 # prevent damaging the Train/test set documents by removing words (optimize later)
    if(selected == {}):
        pass
    else:
        unknownRemover([dic1],selected)
        unknownRemover([dic2],selected)

    for k,v in dic1.items():
        if k in dic2: #Check if Key Exists in Dic2
            s += math.log(dic2[k]+1) * math.log(NDA/DFA[k]) # tf*idf

    return s

# optional:
def distance(dic1,dic2): # Euclidean distance
    pass

```

Phase 4) Implementing KNN from scratch

```

In [29]: def KNN(dic1, flag ,k, selected):
    # input=(queryDocument,function is used for finding similarity ,number of neighbors, selected features)
    # type of inputs (dictionary, num , num , dictionary)
    # we have access to our train set, So we don't pass train set
    # selected feature or in other word, most important words
    # if flag is 0 we use Cosine func
    # if flag is 1 we use tfidf

    neighbors = [(0,0)]*k # k nearest neighbors and their similarity value
    # ( name of category of neighbor , similarity value)

    if(flag==0):
        # we use cosine similarity function

        # passing all dics in SVTrain
        for dic2 in SVTrain:
            neighAdder(neighbors,k, 1 , cosine(dic1 , dic2 , selected) )
        # passing all dics in HVTrain
        for dic2 in HVTrain:
            neighAdder(neighbors,k, 0 , cosine(dic1 , dic2 , selected) )

    if(flag==1):
        # we use tfidf similarity function

        # passing all dics in SVTrain
        for dic2 in SVTrain:
            neighAdder(neighbors,k, 1 , tfidf(dic1 , dic2 , selected) )
        # passing all dics in HVTrain
        for dic2 in HVTrain:
            neighAdder(neighbors,k, 0 , tfidf(dic1 , dic2 , selected) )

    # at the end we return 0 (Ham) or 1 (Spam)
    return bestNeighbor(neighbors)

```

Simply we get a query and we find its similarity comparing to all documents in the train set and then we report the nearest class.

Phase 5) Implementing Naïve Bayes from scratch

Best Category for document_j = $\text{Max}_{i \in \text{Categories}} \{ P(c_i | d_j) \}$

Bayes rule: $P(A|B) = \frac{P(A \cap B)}{P(B)}$ } $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$

$P(B|A) = \frac{P(B \cap A)}{P(A)}$ } = $\frac{P(A \cap B)}{P(A)}$

So:

$$P(c_i | d) = \frac{P(c_i \cap d)}{P(d)} = \frac{P(d | c_i) \cdot P(c_i)}{P(d)} \quad ** \quad \frac{P(w_1, w_2, \dots, w_n | c_i)}{P(d)}$$

$$\alpha \approx P(w_1, w_2, \dots, w_n | c_i) \cdot \frac{|c_i|}{|N|} = \frac{|c_i|}{|N|} \cdot (P(w_1 | c_i) \dots P(w_n | c_i))$$

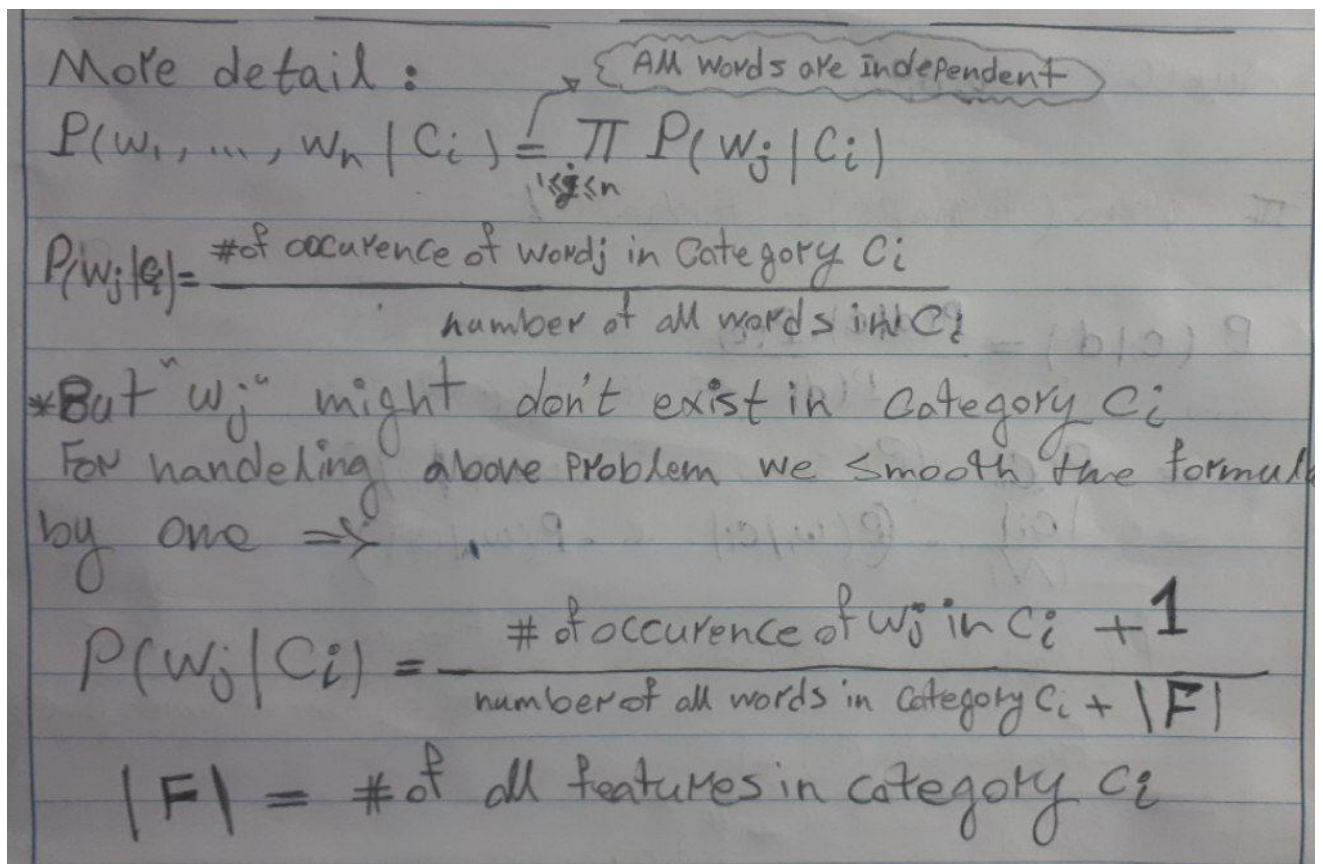
** { without damaging accuracy, we assume all words in a document are Independent.

$P(c_i) = \frac{|\text{all tweets in category } c_i|}{|\text{all tweets}|}$

** : } since $P(d)$ is a fixed & nonchanging value & it doesn't depend on any specific category, so we can omit $P(d)$ from formula

Finally:

Best Category for d_j = $\text{Max}_{i \in \text{cates}} \{ (P(c_i) \cdot P(w_1, w_2, \dots, w_n | c_i)) \}$



Based on the knowledge that we have gotten from last phase, now we can design a classifier. In order to do that, at first we find the value of all words in our bag of words given corresponding category ($P_{w_j \in C_i}(w_j | C_i)$).

We store the result in Prob_P, Prob_N, Prob dictionaries.

Prob_P is dictionary for selected features in category positive and its probability given class "positive". And same description for Prob_N that is a dictionary for class "negative" and Prob is related to class "neutral".

Classifier function gives a string and then it calculate its probability regarding to each category that we have and then return the most likely category.

```
## Naive Bayse Classifier :
def Classifier(string):
    CleanList = TextCleaner(string) #clean list of words
    resultDic = {}
    resultDic["positive"] = Probability(CleanList, "positive")
    resultDic["negative"] = Probability(CleanList, "negative")
    resultDic["neutral"] = Probability(CleanList, "neutral")
    resultDic = sorted(resultDic.items(), key=lambda x: x[1], reverse=True) #resultDic is no longer a dic (It's a list)
    return resultDic[0][0]
```

→ sorting the result and then returning the cate's name for the greatest probability

Phase 6) Calculating Precision, Recall, F1, Confusion Matrix

Matrix:

| | correct | not correct |
|--------------|---------|-------------|
| selected | tp | fp |
| not selected | fn | tn |

$$P = \text{Precision} = \frac{tp}{tp+fp}, R = \text{Recall} = \frac{tp}{tp+fn}, F_1 = \frac{2*PR}{P+R}$$

5.1) Naive Bayes

```
# result for classifying Spam testset by using Bayes
SVTestBayes = BayesList(SVTest,{})
# result for classifying Ham testset by using Bayes
HVTestBayes = BayesList(HVTest,{})
# print(SVTestBayes)
# print(HVTestBayes)
temp_ac, temp_pre, temp_recall, temp_f1 = confusion_matrix(SVTestBayes, HVTestBayes)
acc.append(temp_ac)
pre.append(temp_pre)
recall.append(temp_recall)
f1.append(temp_f1)
```

| Confusion Matrix | Predicted: Spam | Predicted: Ham |
|------------------|-----------------|----------------|
| Actual: Spam | 163 | 37 |
| Actual: Ham | 2 | 198 |

5.2) KNN cosine on all words

```
In [35]: # result for classifying Spam testset by using KNN cosine and using All words
SVTestKNNcosine = KNNList(SVTest, 0, 3, {})
# result for classifying Ham testset by using KNN cosine and using All words
HVTestKNNcosine = KNNList(HVTest, 0, 3, {})
# print(SVTestKNNcosine)
# print(HVTestKNNcosine)
temp_ac, temp_pre, temp_recall, temp_f1 = confusion_matrix(SVTestKNNcosine, HVTestKNNcosine)
acc.append(temp_ac)
pre.append(temp_pre)
recall.append(temp_recall)
f1.append(temp_f1)
```

| Confusion Matrix | Predicted: Spam | Predicted: Ham |
|------------------|-----------------|----------------|
| Actual: Spam | 190 | 10 |
| Actual: Ham | 24 | 176 |

5.3) KNN tf-idf on all words

```
In [36]: # result for classifying Spam testset by using KNN tf-idf and using All words
SVTestKNNtfidf = KNNList(SVTest, 1, 3, {})
# result for classifying Ham testset by using KNN tf-idf and using All words
HVTestKNNtfidf = KNNList(HVTest, 1, 3, {})
# print(SVTestKNNtfidf)
# print(HVTestKNNtfidf)
temp_ac, temp_pre, temp_recall, temp_f1 = confusion_matrix(SVTestKNNtfidf, HVTestKNNtfidf)
acc.append(temp_ac)
pre.append(temp_pre)
recall.append(temp_recall)
f1.append(temp_f1)
```

| Confusion Matrix | Predicted: Spam | Predicted: Ham |
|------------------|-----------------|----------------|
| Actual: Spam | 195 | 5 |
| Actual: Ham | 58 | 142 |

5.4) KNN cosine on important words ¶

```
In [37]: # result for classifying Spam testset by using KNN cosine and selected (important) words
SVTestKNNCosineSelected = KNNList(SVTest, 0 , 3 , DFAS)
# result for classifying Ham testset by using KNN cosine and selected (important) words
HVTestKNNCosineSelected = KNNList(HVTest, 0 , 3 , DFA)
# print(SVTestKNNCosineSelected)
# print(HVTestKNNCosineSelected)
temp_ac, temp_pre, temp_recall, temp_f1 = confusion_matrix(SVTestKNNCosineSelected, HVTestKNNCosineSelected)
acc.append(temp_ac)
pre.append(temp_pre)
recall.append(temp_recall)
f1.append(temp_f1)
```

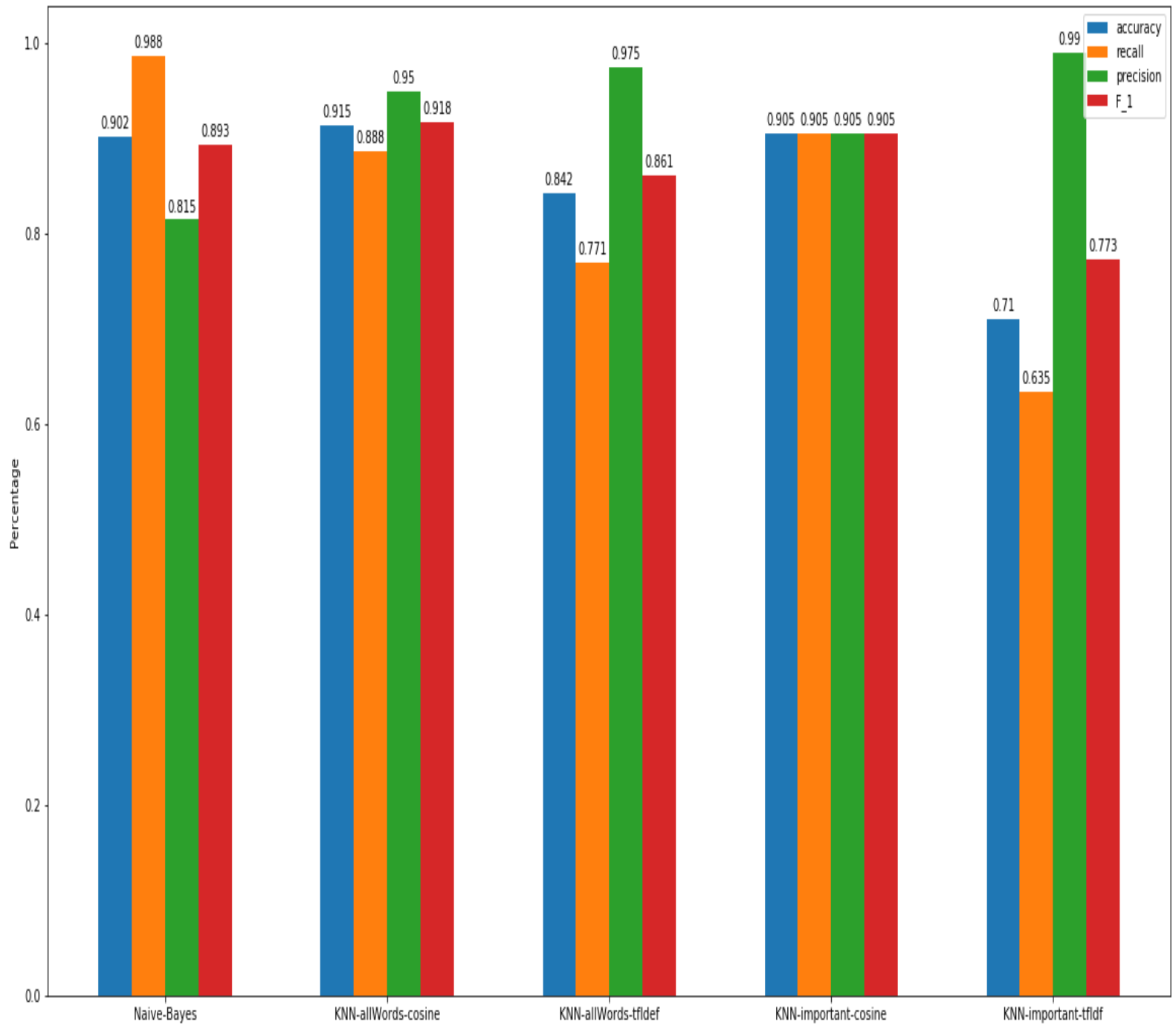
| Confusion Matrix | Predicted: Spam | Predicted: Ham |
|------------------|-----------------|----------------|
| Actual: Spam | 181 | 19 |
| Actual: Ham | 19 | 181 |

5.5) KNN tf-idf on important words

```
[38]: # result for classifying Spam testset by using KNN tf-idf and selected (important) words
SVTestKNNtfIdfSelected = KNNList(SVTest, 1 , 3 , DFAS)
# result for classifying Ham testset by using KNN tf-idf and selected (important) words
HVTestKNNtfIdfSelected = KNNList(HVTest, 1 , 3 , DFA)
# print(SVTestKNNtfIdfSelected)
# print(HVTestKNNtfIdfSelected)
temp_ac, temp_pre, temp_recall, temp_f1 = confusion_matrix(SVTestKNNtfIdfSelected, HVTestKNNtfIdfSelected)
acc.append(temp_ac)
pre.append(temp_pre)
recall.append(temp_recall)
f1.append(temp_f1)
```

| Confusion Matrix | Predicted: Spam | Predicted: Ham |
|------------------|-----------------|----------------|
| Actual: Spam | 198 | 2 |
| Actual: Ham | 114 | 86 |

data and methods' effects on results



This bar chart can summarize everything and presents the power of each method.

We can see that KNN using all words with Cosine similarity measure has the best accuracy in the average .

Thanks for your time.