# Assignment № 1

Bishnu Sarker, Khulna Univesity of Engineering  Technology 26/10/2016

## Naive Bayes Classifier

Naive Bayes classifier is one of the simplest and popular probabilistic machine learning techniques. Naive Bayes classifier works based on Bayesian Theorem.

$$\text{According to Bayes theorem, } P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

To use this theorem in machine learning, let's define a classification problem: Given, a data point d, and a set of classes $C = c_1, c_2, c_3, .., c_n$, the likelihood of d being belongs to the any of the classes is P(C|d).

To expand this using Bayes theorem, we have:

$$P(C|d) = \frac{P(d|C)P(C)}{P(d)}$$

P(d) is same for all of the classes. Therefore, we can omit the this and we have following equation:

$$P(C|d) = P(d|C)P(C)$$

To decide on appropriate class, we have to take class for which P(C|d) is maximum. So the problem can be reduced to the follwoing form:

$$argmax_C P(C|d)$$
$$= argmax_C P(d|C)P(C)$$

Lets split the document in it's constituent tokens such as $d = w_1 w_2 w_3 w_4 w_5 .. w_n$

$$= argmax_C P(d = w_1 w_2 w_3 w_4 w_5 .. w_n | C)P(C)$$

Assuming terms are independent, we have:

$$= argmax_C P(C)P(w_1|C)P(w_2|C)P(w_3|C)..P(w_n|C)$$
$$= argmax_C P(C) \prod_{w_i \in V} P(w_i|C) \text{ , V is the Vocabulary}$$
$$= argmax_C \log P(C) + \sum_{w_i \in V} \log P(w_i|C) \text{ , V is the Vocabulary}$$

We can compute the class probability using following formula:

$$P(C) = \frac{N_c}{N}$$
$$P(w_i|C) = \frac{count(w_i,c)+1}{V + \sum_{w_i \in V} count(w_i,C)}$$

Here is the a python code for working with Naive Bayes:

## Implementation

Listing 1: Sample Python code – Naive Bayes for Product Review Analysis.

```python
'''
Author: Bishnu Sarker, dept. of Computer Sciecne and Engineering, Khulna ↩
    University of Engineering & Technology
'''
import random
import math

def write_object(object,filename):
    '''
    To write a python list object in external file
    :param object: Object prefereably a List object to be written in a ↩
        file
    :param filename: Filename in which data will be written
    :return: NA
    '''
    with open(filename,'w') as W:
        for item in object:
            W.write(str(item)+'\n')

def read_review(filename):
    '''
    Reading raw review and prepraring dataset seperating pos and neg data
    :param filename: Dataset file name
    :return:tuple, positive dataset and negative dataset
    '''
    pos_dataset=[]
    neg_dataset=[]
    with open(filename) as R:
        for line in R:
            #print(line)
            line=line.split()
            if line[0]=="Neg":
                neg_dataset.append(list(line[1:]))
            else:
                pos_dataset.append(list(line[1:]))
    write_object(pos_dataset,"positve.txt")
    write_object(neg_dataset,"negative.txt")
    return (pos_dataset,neg_dataset)

def split_test_train(pos, neg):
    '''
    Splits the datasets into train and test set
```

```python
41      :param pos: positive dataset
42      :param neg:  negative dataset
43      :return:  postive training, negative training, positive testing and ↵
            negative training
44      '''
45      pl=len(pos)
46      nl=len(neg)
47
48      pSample_train=random.sample(pos,int(pl*0.80))
49      nSample_train=random.sample(neg,int(nl*0.80))
50      pSample_test=random.sample(pos,int(pl*0.20))
51      nSample_test=random.sample(neg,int(nl*0.20))
52
53      return (pSample_train,nSample_train,pSample_test,nSample_test)
54
55  def build_Vocab(pTrain,nTrain):
56      '''
57      Building a vocabulary  V from taining dataset
58      :param pTrain: positive training dataset
59      :param nTrain: negative training dataset
60      :return: Vocabulary of the form {"Token":(p, n)} p for positve ↵
            frequency and n for negative frequency
61      '''
62      Vocabulary={}
63      for sample in pTrain:
64          for word in sample:
65
66              word=word.lower().strip("'!.:)(?-")
67              if word == "":
68                  continue
69              if Vocabulary.__contains__(word):
70                  freq,T=Vocabulary[word]
71                  freq=freq+1
72                  Vocabulary[word]=(freq,T)
73              else:
74                  Vocabulary.__setitem__(word,(1,0))
75      for sample in nTrain:
76
77          for word in sample:
78              word = word.lower().strip("'!.:)(?-[]+")
79              if word=="":
80                  continue
81              if Vocabulary.__contains__(word):
82                  T,freq=Vocabulary[word]
83                  freq=freq+1
84                  Vocabulary[word]=(T,freq)
85              else:
```

```python
86                    Vocabulary.__setitem__(word,(0,1))
87        print_dic(Vocabulary)
88        return Vocabulary
89
90    def Maxmimum_likelihood_Estimation(pTrain,nTrain):
91        '''
92        Compute the required probabilities using maximum likelihood estimation←
                .
93        :param pTrain:
94        :param nTrain:
95        :return:
96        '''
97
98        pN=len(pTrain)
99        nN=len(nTrain)
100       N=pN+nN
101       pCount=0
102       nCount=0
103       vocab=build_Vocab(pTrain,nTrain)
104       V=len(vocab)
105       for words in vocab:
106           pCount=pCount+vocab[words][0]
107           nCount=nCount+vocab[words][1]
108       return {"vocab":vocab,"pN":pN,"nN":nN,"N":N,"V":V,"pC":pCount,"nC":←
              nCount}
109
110
111   def NB_classifier(test_sample,Vocab,pN,nN,N,V,pCount,nCount):
112       '''
113       Naive bayes classifier
114       :param test_sample:
115       :param Vocab:
116       :param pN:
117       :param nN:
118       :param N:
119       :param V:
120       :param pCount:
121       :param nCount:
122       :return:
123       '''
124       #tokens=test_sample.strip().split()
125       N=float(N)
126       prob_pos=pN/N
127       prob_neg=nN/N
128       pos_prob=0.0
129       neg_prob=0.0
130       for token in test_sample:
```

```python
131              token=token.lower().strip("'!.:)(?-[]+")
132              #computing p(token|pos)
133              if token=="":
134                  continue
135              if Vocab.__contains__(token):
136                  pos_prob=pos_prob+math.log10(Vocab[token][0]+1.0/(pCount+V))
137                  neg_prob=neg_prob+math.log10(Vocab[token][1]+1.0/(nCount+V))
138              else:
139                  pos_prob = pos_prob + math.log10( 1.0 / (pCount + V))
140                  neg_prob = neg_prob + math.log10( 1.0 / (nCount + V))
141
142          pP=math.log10(prob_pos)+pos_prob
143          nP=math.log10(prob_neg)+neg_prob
144          return (pP,nP)
145
146  def testing(pTest,nTest,Param):
147      '''
148      Find the probabale class for a test sample
149      :param pTest:
150      :param nTest:
151      :param Param:
152      :return:
153      '''
154      with open("Rsult.txt",'w') as W:
155          for item in pTest:
156              pP,nP=NB_classifier(item,Param["vocab"],Param["pN"],Param["nN"↩
                      ],Param["N"],Param["V"],Param["pC"],Param["nC"])
157              print("Prediction-->\tpos:"+str(pP)+" \tNeg:"+str(nP)+" \↩
                      tActual-->Positive")
158              W.write("Prediction--> pos:"+str(pP)+" \tNeg:"+str(nP)+" \↩
                      tActual-->Positive\t"+str(item)+'\n')
159          for item in nTest:
160              pP, nP = NB_classifier(item,Param["vocab"],Param["pN"],Param["↩
                      nN"],Param["N"],Param["V"],Param["pC"],Param["nC"])
161              print("Prediction--> pos:" + str(pP) + " Neg:" + str(nP)+" ↩
                      Actual-->Negative")
162              W.write("Prediction--> pos:" + str(pP) + " \tNeg:" + str(nP) +↩
                      " \tActual-->Negative\t"+str(item)+"\n")
163  def print_dic(D):
164      with open("Vocab.txt",'w') as W:
165          for d in D.keys():
166              W.write(d+"-->"+" Pos:"+str(D[d][0])+" Neg:"+str(D[d][1])+'\n'↩
                      )
167
168  if __name__ == '__main__':
169      pos,neg=read_review("AppleReview.txt")
170      pos_train,neg_train,pos_test,neg_test=split_test_train(pos,neg)
```

```
171     Params=Maxmimum_likelihood_Estimation(pos_train,neg_train)
172     print (Params)
173     testing(pos_test,neg_test,Params)
```

Code 1... shows a typical implementation of naive Bayes following the formula described above.

## Tasks to be done

- Select A data set from following sources:

    - Twitter Data: http://www.sananalytics.com/lab/twitter-sentiment/

    - Movie Review Data https://www.cs.cornell.edu/people/pabo/movie-review-data/

    - Kaggle Competition: https://www.kaggle.com/c/sentiment-analysis-on-movie-reviews

    - https://inclass.kaggle.com/c/si650winter11/data

- Experiemnt with given datasets and document the classification error to evaluate the model performance. The model performance is computed using confusion Matrix.

$$Error = \frac{No\ of\ Missclassified\ examples}{Total\ No\ of\ Examples\ in\ the\ test\ set}$$

## Deliverable

A short report (max 3 page docx and pdf) and Code. Things that should be stated in your report:

- General Introduction

- Description of the datasets you have tried

- Experiment setup and Result Analysis

- Conclusion and discussion

- References