IRTM Wi 20/21 Assigment 4

Task 1

Vocabulary = { happy, new, year, holiday, term, starts, work, celebrations } |Vocabulary| = 8

Prior:

$$P(c_1) = \frac{3}{5} = 0.6$$

 $P(c_2) = \frac{2}{5} = 0.4$

Posterior: (with Add-One-Smoothing)

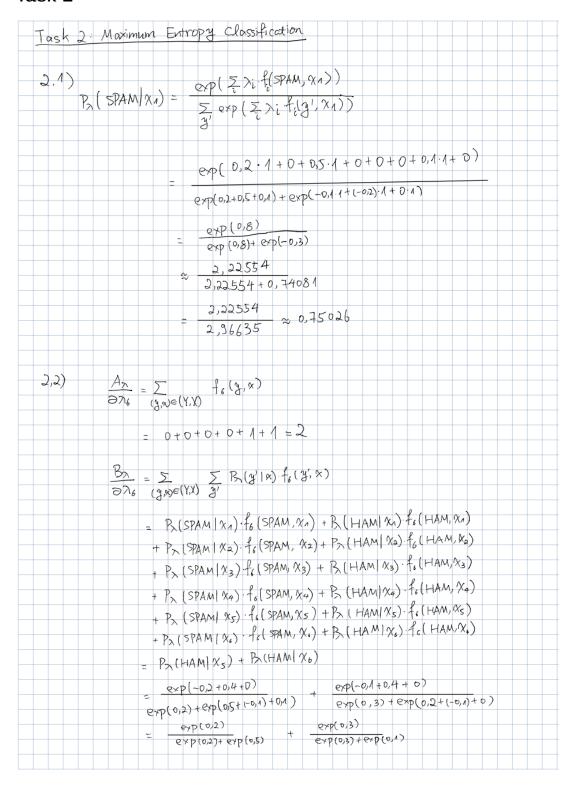
$$\begin{array}{ll} P(happy|c_1) = \frac{2+1}{7+8} = \frac{3}{15} & P(happy|c_2) = \frac{0+1}{4+8} = \frac{1}{12} \\ P(new|c_1) = \frac{2+1}{7+8} = \frac{3}{15} & P(new|c_2) = \frac{0+1}{4+8} = \frac{1}{12} \\ P(year|c_1) = \frac{2+1}{7+8} = \frac{3}{15} & P(year|c_2) = \frac{0+1}{4+8} = \frac{1}{12} \\ P(holiday|c_1) = \frac{1+1}{7+8} = \frac{2}{15} & P(holiday|c_2) = \frac{0+1}{4+8} = \frac{1}{12} \\ P(term|c_1) = \frac{0+1}{7+8} = \frac{1}{15} & P(term|c_2) = \frac{1+1}{4+8} = \frac{2}{12} \\ P(starts|c_1) = \frac{0+1}{7+8} = \frac{1}{15} & P(starts|c_2) = \frac{2+1}{4+8} = \frac{1}{12} \\ P(work|c_1) = \frac{0+1}{7+8} = \frac{1}{15} & P(work|c_2) = \frac{1+1}{4+8} = \frac{2}{12} \\ P(celebrations|c_1) = \frac{0+1}{7+8} = \frac{1}{15} & P(celebrations|c_2) = \frac{0+1}{4+8} = \frac{1}{12} \end{array}$$

$$P(c_1|d) = \frac{3}{5} * \frac{3}{15} * \frac{3}{15} * \frac{3}{15} * \frac{1}{15} = \frac{81}{253125} = c_{map}$$

$$P(c_2|d) = \frac{1}{12} * \frac{1}{12} * \frac{1}{12} * \frac{1}{12} = \frac{2}{103680}$$

The model assigned the class c_1 to the document.

Task 2



	1,22140 1,34986 1,34986 1,34986 1,34986
	1,22140 1,34986 1,45503
	≈ 0,65311 + 0,92 → 72 = 1,58083
2A2 -	$\frac{\partial B_n}{\partial \lambda_6} = 2 - 1,58083 = 0,41917$

Task 3

```
• k=1

P(black|x) = 0

P(white|x) = \frac{1}{1} = 1

• k=3

P(black|x) = \frac{1}{3}

P(white|x) = \frac{2}{3}

• k=5

P(black|x) = \frac{2}{5}

P(white|x) = \frac{3}{5} = 1
```

In case of k=5 the kNN classifier hast the lowest classification confidence.

Programming Task

```
import nltk
  import pandas as pd
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
 from collections import Counter
 from operator import itemgetter
  class SVM_Classifier:
      path = "04-assignment/code/"
10
      def __init__(self):
11
          self.corpus = self.generate_corpus("games-train.csv")
          self.term_weights_dictionary = self.calc_term_weights_dictionary()
13
14
15
16
      def generate_corpus(self, filename: str) ->pd.DataFrame:
17
            "" Take data as an imput and generate the corpus for classification.
18
              The data-structure is a pandas DataFrame with two columns.
19
              The labels are converted in binary values (1 = 'gut', -1 = 'schlecht')
20
              and the text is normalized. Some stopsword will be also removed.
21
22
          Args:
23
              filename (str): name of the dataset
24
2.5
26
              pd.DataFrame: Pandas DataFrame with two columns(label, text)
27
28
          # import data
          dataframe = pd.read_csv(self.path + filename, sep="\t", header=None)
31
32
          # stops that will remain in the corpus
```

```
relevant_stops = set(['nicht', 'nichts', 'kein', 'kein', 'keine', 'keinem', '
34
               keinen', 'keiner', 'keines'])
          # stops to remove
36
           stops = set(stopwords.words('german')) - relevant_stops
37
38
39
40
          # process data to create corpus
41
          label_col , text_col = (dataframe[1].tolist(), dataframe[3].tolist())
42
          binary_label_col, normalized_text_col = ([], [])
43
44
           for index in range(len(label_col)-1):
45
46
               label, text = (label_col[index], text_col[index])
47
48
               # covert the labels into binary values
49
               \# 'gut' = 1, 'schlecht' = -1
50
               if label == 'gut':
51
                   binary_label_col.append(1)
52
               else:
                   binary_label_col.append(-1)
54
55
56
               # text normalization
57
               normalized_text = word_tokenize(str(text).lower())
58
59
60
               # remove stops
               for token in normalized_text:
61
                   if token in stops: normalized_text.remove(token)
62
63
               normalized_text_col.append(word_tokenize(str(text).lower()))
64
65
66
          # create the corpus's data-structure
67
          corpus = {'label': binary_label_col , 'text': normalized_text_col}
68
69
           return pd. DataFrame (corpus)
73
      def calc_term_weights_dictionary(self) ->dict:
           """ Calculate the weight of each term and store them in a dictionary.
74
          the weight is calulated with term-frequency * label.
75
76
           Returns:
               dict: store the weight as a value of a term
78
79
80
          term_weights_dict = {}
81
           for row in self.corpus.iterrows():
82
83
               label, text = row[-1]
84
               term_freq_dict = Counter(text)
85
86
               for term in term_freq_dict:
87
88
```

```
if term not in term_weights_dict:
89
                        term_weights_dict.update({ term: label * term_freq_dict[term] })
90
91
                         # sum the term fequencies
92
                         # if the class is 'schlecht' the term's frequency will be
93
                             negative
                         term_weights_dict[term] += label * term_freq_dict[term]
94
95
           return term_weights_dict
96
97
98
99
100
102
103
   if __name__ == "__main__":
104
       classifier = SVM_Classifier()
105
       dictionary = classifier.term_weights_dictionary
106
107
108
       sorted_term_weight_list = sorted(dictionary.items(), key=itemgetter(1))
109
       print("-"*50, "\n Top 100 result for class 'gut': \n", "-"*50)
111
       for term, weight in sorted_term_weight_list[::-1][:100]:
           print(" ", weight, "\t", term)
113
114
       print("\n" * 3)
115
116
       print("-"*50, "\n Top 100 result for class 'schlecht': \n", "-"*50)
117
       for term, weight in sorted_term_weight_list[:100]:
118
            print(" ", weight, "\t", term)
119
```

code/script.py

OUTPUT:

```
Top 100 result for class 'gut':
     24534
                       spiel
     22102
                       !
     19037
                       ist
     18307
                       es
     15812
                       cool
     14978
                       das
     13549
     12106
                      macht
     10989
     10854
                      super
13
     10760
                      ich
     10104
15
                      und
     10019
16
                       )
     9793
             geil
17
```

```
8686
               aber
      8514
               gut
19
               sehr
      8346
20
      7933
               einfach
      7425
22
               man
23
      7323
               spaß
      7083
24
      6499
               ein
25
      6411
               nur
26
      6316
               die
27
      4627
               echt
28
      4427
29
      4038
               finde
30
      3973
31
               so
      3931
               zu
32
      3847
               noch
33
      3601
34
               viel
      3489
               nachbarn
35
      3476
               auch
36
      3024
               süchtig
37
      2906
               kann
38
      2805
               suche
39
      2802
               für
40
      2692
               richtig
41
      2586
               der
42
      2569
               mich
43
      2531
               toll
      2478
45
               wenn
      2464
46
      2463
               game
47
      2439
               mit
48
      2390
               gutes
49
      2365
               dieses
50
      2252
               voll
51
52
      2128
               adden
53
      2127
               spiele
      2058
               liebe
54
      2051
               cooles
55
      2023
               top
56
      1967
               hammer
57
      1934
               tolles
58
      1900
               bin
59
      1765
60
      1717
               mir
61
      1686
               neue
62
      1679
               hat
63
      1668
               geiles
64
65
      1647
               beste
66
      1637
               sind
67
      1623
               sonst
      1616
68
               immer
      1602
69
               . . .
      1541
               weiter
70
      1531
               klasse
71
      1471
               muss
72
73
      1451
               alles
```

```
langweilig
       1431
       1428
                bitte
 75
       1420
                freunde
 76
 77
       1411
                manchmal
 78
       1394
                wie
                zeitvertreib
 79
       1385
       1350
 80
                app
       1337
                also
 81
       1329
                spass
 82
       1298
                ganz
 83
       1271
                gibt
 84
       1250
                lol
 85
       1240
                is
 86
       1227
                5
 87
       1202
                sterne
 88
       1198
                mega
 89
       1189
               wäre
       1176
 91
                was
       1153
                eine
 92
       1131
                auf
 93
       1114
                spielen
 94
       1087
                läuft
 95
       1062
                könnte
 96
       1057
                empfehlen
 97
       1057
                dass
 98
       1040
                ihr
       1000
                wirklich
       991
                weil
101
       987
                jeden
102
       985
                nie
103
104
106
108
109
      Top 100 result for class 'schlecht':
110
       -2016
111
                          nicht
       -1837
                         mehr
112
       -1183
113
                          seit
       -987
114
                update
       -953
                beheben
115
       -902
               komme
116
       -721
                rein
117
       -711
118
       -711
                ins
119
       -678
                weg
120
121
       -641
                event
122
       -632
                war
123
       -619
                stürzt
       -588
124
                geht
       -564
               dem
125
       -542
                nichts
126
       -534
                ständig
127
       -498
                soll
128
       -473
                keinen
129
```

```
-453
                heute
130
       -419
                letzten
131
       -408
                verbindung
       -407
133
                gar
       -367
                neu
134
135
       -359
               komm
       -355
136
                ab
       -353
137
                server
       -350
                deinstalliert
138
       -343
                anfangen
139
       -321
                tagen
140
       -318
                fehler
141
       -308
                starten
142
       -303
                vorne
143
       -303
                scheiße
144
       -302
                scheiß
145
       -299
146
                nix
       -279
                gelöscht
147
       -269
                startet
148
       -257
                behebt
149
       -255
                lässt
150
       -246
                werde
       -226
                öffnen
       -221
153
       -218
                wieder
154
       -218
                steht
155
       -217
156
                stern\\
       -210
157
               möglich
       -208
                obwohl
158
       -202
                nun
159
       -198
                mist
160
       -191
                dann
161
       -190
                support
162
       -189
163
       -176
                gestern
164
165
       -170
                spielstand
166
       -168
                null
167
       -165
       -157
                schlüssel
       -154
                lädt
       -150
170
                kotzen
       -145
       -143
                wurde
       -137
173
       -137
                gekauft
174
       -136
                anmelden
175
       -136
                zynga
176
177
       -134
                behoben
178
       -133
                enttäuscht
179
       -132
                sauer
       -130
180
                bekomme
       -129
                wollte
181
       -127
                raus
182
       -125
                investiert
183
       -121
                zurück
184
       -121
                passiert
185
```

Alberto Saponaro - saponaroalberto97@gmail.com Walter Väth - walter.vaeth@gmail.com Chong Shen - st143575@stud.uni-stuttgart.de Xin Pang - st145113@stud.uni-stuttgart.de

```
-119
               bildschirm
186
       -118
               beendet
187
       -115
               sofort
188
               deinstallieren
       -115
189
190
       -114
               abzocke
       -113
               belohnung
191
       -113
               scheisse
192
               installiert
       -112
193
       -112
               dauernd
194
       -111
               beim
195
       -108
               laden
196
       -106
               langsam
197
       -106
               installieren
198
       -105
               ärgerlich
199
       -104
               andauernd
200
       -103
               wochen
201
       -99
               seid
202
       -98
               trotz
203
       -97
               sekunden
204
       -96
               morgen
205
       -95
               zeigt
206
       -94
               angezeigt
207
       -93
               plötzlich
208
       -91
               bricht
209
       -89
               mü l l
210
```