Pen and Paper Task 1

Subtask A

- $sun \rightarrow 1, 2, 3, 4, 7$
- nice \rightarrow 1, 5, 6
- water \rightarrow 5, 6, 8, 9
- is \rightarrow 6
- beer $\rightarrow 10$

Subtask B

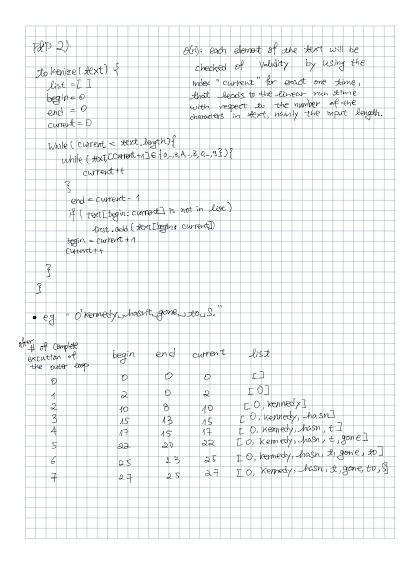
- $sun \to 1, 2, 3, 4, 7$
- nice \rightarrow 1, 5, 6
- water \rightarrow 5, 6, 8, 9
- is \rightarrow 6
- beer $\rightarrow 10$

Example query: nice AND is

- 1. Comparisons without skip pointers: 1 & 6, 5 & 6, 6 & 6 \Rightarrow 3 comparisons
- 2. Comparisons with skip pointers: 1 & 6, 6 & 6 \Rightarrow 2 comparisons

Without skip pointers we must compare the terms step by step, although 5 in **nice** is still smaller than 6 in **is**. With skip pointers we can skip the 5 in **nice** and directly go to the 6 in **nice**.

Pen and Paper Task 2



Alberto Saponaro - saponaroalberto 97@gmail.com Walter Väth - walter.vaeth@gmail.com Chong Shen - st143575@stud.uni-stuttgart.de Xin Pang - st145113@stud.uni-stuttgart.de

IRTM Wi 20/21 Assigment 1

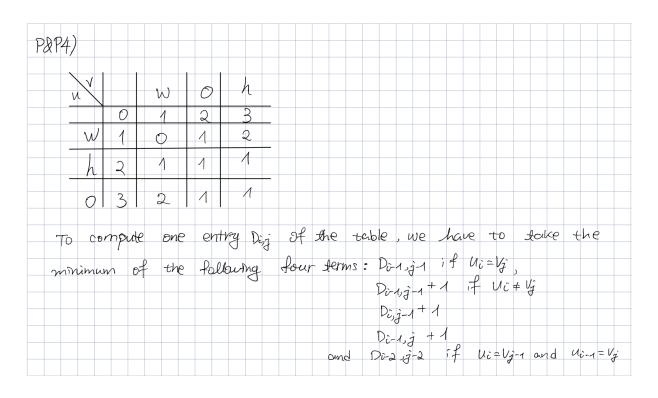
Pen and Paper Task 3

```
Query: Gates /2 Microsoft
```

```
(Gates, 4): [{1:[3], 2:[6], 3:[2, 17], 4:[1]}]
(Microsoft, 4): [{1:[1], 2:[1, 21], 3:[3], 5:[16, 22, 51]}]
cross product = { 1:[(3,1)], 2:[(6,1), (6,21)], 3:[(2,3), (17,3)], 4:[(1,16), (1,22), (1,51)] }
```

From all tuples in the cross product, the tuples (3,1) and (2,3) fulfill the query's condition $abs(tuple[1]-tuple[0]) \le 2$. So the answer is: document 1, document 3.

Pen and Paper Task 4



Programming Task 1

```
import csv, re, nltk
  def index(filename: str):
      index = \{\}
      dictionary = {}
      postings_lists = []
      tokenizer = nltk.RegexpTokenizer(r"\w+")
      with open(filename, 'r') as file:
           reader = csv.reader(file, delimiter = '\t')
11
           postings = []
13
          #iterate through each row of the table
14
           for row in reader:
15
               (doc_id, url, pub_date, title, news_text) = row
16
               #tokenize and normalize news text
18
19
               #this procedure will remove symbols like !?() etc.
               #the set data structure will remove all duplicates
21
               news_text_norm = set(tokenizer.tokenize(news_text.lower()))
               #generate postings
23
               #iterate through each term
24
               for term in news_text_norm:
25
                   postings.append((term, doc_id))
26
          #sort postings
28
           postings = sorted (postings [1:], key = lambda tup: tup [0])
29
30
31
32
           post_id = 0
33
          post\_size = 0
34
          #iterate through postings
           for posting in postings:
               term, doc_id = posting
36
               if term not in dictionary:
38
                   #upate the dictionary with the new term
39
                   #initialize the postings size
40
                   #save the postings id,
41
                   #witch is the position of the postings list
42
                   #into the postings lists
43
                   dictionary.update({term: [post_size+1, post_id]})
44
45
                   #initialize a new postings list
46
                   postings_lists.append([doc_id])
47
48
                   #update postings id
49
                   post_id +=1
50
51
                   #update size of posting
52
                   dictionary [term][0] += 1
```

```
54
                    #update postings list
                    postings_lists[-1].append(doc_id)
56
57
       return dictionary, postings_lists
58
59
   class Search:
60
61
       def __init__(self, filename: str, index: tuple):
62
            self.filename = filename
63
            self.index = index
64
65
       def getPostingList(self, postID):
66
            postings_lists = self.index[1]
67
68
           #set the index of the first postings list
69
           idx = 0
70
           #iterate through postings lists
71
            for postings_list in postings_lists:
72
                if postID == idx:
73
                    return postings_list
74
                    break
75
                else:
76
                    #update index
                    idx += 1
78
79
       def query(self, term: str):
80
81
            dictionary, postings_lists = self.index
            postID = dictionary[term][1]
82
            postings_list = self.getPostingList(postID)
83
84
           #retrive text
85
           with open(filename, 'r') as file:
86
                reader = csv.reader(file, delimiter = '\t')
87
                postings = []
88
89
90
                #iterate through each row of the table
91
                for row in reader:
                    (docID, url, pub_date, title, news_text) = row
92
                    if docID in postings_list:
93
                         return(docID, news_text)
94
95
            pass
96
97
       def query(self, term1: str, term2: str = ''):
98
            dictionary, postings_lists = self.index
99
            out_list = []
100
102
           #CASE 1: only one term
            if term2 == '':
103
104
                postID = dictionary [term1][1]
                postings_list = self.getPostingList(postID)
106
                #retrive text
                with open(filename, 'r') as file:
108
                    reader = csv.reader(file, delimiter = '\t')
109
```

```
postings = []
                    #iterate through each row of the table
                    for row in reader:
113
                        (docID, url, pub_date, title, news_text) = row
114
                         if docID in postings_list:
115
                             out_list.append((docID, news_text))
116
           #CASE 2: two terms
           else:
118
119
                intersection_list = []
                term1_postID = dictionary[term1][1]
                term2_postID = dictionary[term2][1]
124
                term1_postings_list = self.getPostingList(term1_postID)
125
                term2_postings_list = self.getPostingList(term2_postID)
126
127
                #intersection algorithm
128
                for term1_docID in term1_postings_list:
129
                    for term2_docID in term2_postings_list:
130
                        if term1_docID == term2_docID :
                             intersection_list.append(term1_docID)
                #retrive text
134
                with open(filename, 'r') as file:
135
                    reader = csv.reader(file, delimiter = '\t')
136
                    postings = []
138
                    #iterate through each row of the table
139
140
                    for row in reader:
                        (docID, url, pub_date, title, news_text) = row
141
                        if docID in intersection_list:
142
                             out_list.append((docID, news_text))
143
           return out_list
144
145
146
      __name__ == "__main__":
147
       filename = 'assignment1/code/postillon.csv'
148
       index = index (filename=filename)
149
150
       search = Search(filename=filename, index=index)
       #queries
       print('weiß AND maß')
       for item in search.query('weiß', 'maß'):
154
155
           print(item)
156
       print('weiß AND masse')
158
       for item in search.query('weiß', 'masse'):
159
           print(item, '\n')
160
       print('weiss AND maße')
161
       for item in search.query('weiss', 'maße'):
162
           print(item, '\n')
163
164
       print('weiss AND masse')
165
```

IRTM Wi 20/21 Assigment 1 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

```
for item in search.query('weiss', 'masse'):
print(item, '\n')
```

code/script.py

The result shows the docID with the news text. With the use of an iterator -in Python simply a for loop- I iterate through the first and second postings list to find if the two terms are found in both documents.

The differences between the two queries is that in the first one we don't need an intersect algorithm. In this case we can use the postings list of the given term directly. In the second one with two terms we have to calculate the intersection of the two postings lists and the return the text of the documents that have bouth terms in it.

Terminal OUTPUT:

```
weiß AND maß
('691', ' +++ Würde am liebsten im Erdboden versinken: Maulwurf hat das GraACben
    verlernt +++ +++ Aus Fluss: Schwimmerin weiß, woher sie Entzündungskrankeit hat
    +++ +++ Es ist ein Mädchen!: Fliege bekommt Nachwuchs +++ +++ Nur an der Nase
    herumgeführt: Masochist unzufrieden mit Domina-Dienstleistung +++ +++ Maß Capone
    : Schneider von Gangsterboss erhält Frischkäse zum Dank +++ +++ Hat Eindruck
    hinterlassen: Jürgen W. Möllemann +++ +++ Stehen auf der Gersteliste: Gerd Reide
    und Johann Hafer zu Roggen-Roll-Party eingeladen +++ Jetzt bestellen! Die
    besten Newsticker als Buch (nur 9,9 A): Der Postillon: +++ Newsticker +++ ')
weiß AND masse
('864', '
            Berlin (dpo) - Die Polizei spricht von einem gefährlichen neuen Drogen-
    Trend: Immer mehr Berliner begeben sich seit gestern mit Bananenresten in den
    Atemwegen in ärztliche Behandlung. Intensive Recherchen des Postillon im
    Berliner Drogenmilieu haben ergeben, dass der Ursprung der mysteriösen Fälle
    offenbar bei einer 140-Kilo-Lieferung eines neuen Mode-Kokains aus Kolumbien
    liegt, die seit Dienstagvormittag im Umlauf ist. Wie der Stoff in die Szene
    geraten ist, bleibt unklar.
                                  Nach Angaben von Dealer Henrik K.* (*Name von der
    Redaktion geändert), der mehrere Kilogramm des neuartigen Kokains kaufte, hat
    die Lieferung gravierende Mängel. "Nicht einmal gerade Lines lassen sich mit dem
    Zeug legen", klagt der 32-Jährige.
                                           Ist unzufrieden mit dem neuen Stoff:
    Henrik K. Einige seiner Kunden, die spekulierten, dass es sich bei der neuen,
    gelb-gebogenen Darreichungsform um Verpackung handeln könnte und nach
    stundenlangem Kampf zum Inneren vordrangen, berichten von einer matschigen Masse
    , die zwar einigermaßen weiß, aber kein bisschen feinkörnig sei und nur die Nase verstopfe. Lediglich ein Abnehmer, der das neue Kokain im LSD-Rausch aß, habe
    von einem wohligen Völlegefühl berichtet und wolle unbedingt mehr von dem "
    Wunderzeug" kaufen. Henrik K. schöpft daraus jedoch wenig Trost: Seit gestern
    trocknet er 3 Kilogramm der Masse in seinem Backofen, um sie anschließend zu
    zermahlen und doch noch irgendwie nutzbar zu machen. "Wenn das nicht klappt,
    kann ich nur noch schauen, ob es vielleicht intravenös funktioniert. Sonst weiß
    ich auch nicht weiter." Zwar scheint eine Überdosis mit dem neuen Stoff nicht mö
    glich, acht Konsumenten werden aber wohl bleibende Schäden davontragen. Sie
    leierten ihre Nasenlöcher beim Versuch, das Kokain ungeschält zu schnupfen,
    hoffnungslos aus.
             Osnabrück (dpo) - Müssen Fans bald auf eine liebgewonnene Tradition
    verzichten? Nachdem Schiedsrichter Martin Petersen beim Pokalspiel zwischen dem
    VfL Osnabrück und RB Leipzig von einem Feuerzeug getroffen wurde, erwägt der DFB
```

IRTM Wi 20/21 Assigment 1

> nun ernsthaft, das Werfen von Feuerzeugen in Fußballstadien unter Verbot zu stellen. Man befürchtet, dass von den Projektilen Verletzungsgefahr ausgeht. Fußball ist ein Erlebnis für die ganze Familie. Deshalb spielen wir mit dem Gedanken, Feuerzeugwürfe - so schön sie auch anzuschauen sind - zu untersagen", erklärte DFB-Präsident Wolfgang Niersbach. "Derzeit werben wir für die Zustimmung der Vereine zu dieser sicherlich umstrittenen Maßnahme." nur noch zum Entzünden von Pyrotechnik verwendet werden: Feuerzeug Fan-Initiativen zeigen sich empört: "Fußball lebt von Emotionen", erklärt etwa RenÄC Finkel vom Osnabrück-Fanclub Lila-Weiß Westerkappeln. "Da gehört das Werfen von Feuerzeugen aus einer anonymen Masse einfach dazu. Was kommt als Nächstes? Ein Verbot, den Platz zu stürmen?" Tatsächlich hat der Feuerzeugwurf Tradition. Er geht zurück bis Mitte des 19. Jahrhunderts, als es üblich war, Feuersteine und Zunderdosen auf das Spielfeld zu werfen. Ziel der Würfe war es schon damals, gegnerische Spieler oder den Schiedsrichter möglichst schwer zu verletzen und einen Spielabbruch herbeizuführen. Während der Weimarer Republik geriet der Brauch nahezu in Vergessenheit. Zu wenig Wucht entwickelten die damals üblichen Streichhölzer. Erst mit der massenhaften Einführung des Feuerzeugs erlebte der Z ündmittelwurf eine Renaissance. Sollte das Feuerzeugwurfverbot tatsächlich kommen, müssten zahlreiche Fanclubs und Ultra-Gruppierungen ihre Feuerzeugwurf-Bataillone auflösen oder auf andere Wurfgegenstände umlernen lassen. Auch die beliebten Wurffeuerzeug-Stände in den Stadien würden dann der Vergangenheit angehören.

('4933', ' Köln, Stuttgart, München (dpo) - Ein meteorologisches Rätsel sorgt derzeit in weiten Teilen Deutschlands für große Verunsicherung. Der Grund: In vielen Regionen fiel im Laufe des Tages eine bislang noch unbekannte nasskalte, helle Masse vom Himmel. Vielerorts blieb die Substanz liegen und färbte die Landschaft leuchtend weiß ein. Laut Behörden besteht keine Gefahr für die Bevö lkerung. "Es war unheimlich", berichtet Hausfrau Renate P. aus Stuttgart dem Postillon. "Plötzlich hab ich was Kaltes auf der Nase gespürt. Ich schau nach oben und seh, dass da so weißes Zeug runterrieselt. Ganz leise. Sah aus wie bei mir daheim in der Tiefkühltruhe." Frau P. flüchtete sich daraufhin gemeinsam mit anderen Passanten in einen Laden, wo sie mehrere Stunden lang ausharrte. Maaammmiiiii! Ich hab Angst!" Inzwischen haben Meteorologen eine Erklärung für das Phänomen parat. Demnach soll es sich nicht wie zunächst vermutet um Vulkanasche oder radioaktiven Fallout handeln. "Vielmehr haben wir es bei der weißen Substanz mit einer speziellen Form von Niederschlag zu tun, wie man sie sonst aus dem skandinavischen Raum oder aus alten Erzählungen kennt", erklärt Susanne Brauer vom Deutschen Wetterdienst. "Früher hat diese Substanz in der Winterzeit oft weite Teile Europas bedeckt." Zwar halten die Behörden den Niederschlag bislang für ungiftig, doch zumindest auf Autofahrer scheint die unbekannte Masse einen bizarren Effekt zu haben. Viele von ihnen reagieren panisch und verwirrt und scheinen innerhalb kürzester Zeit selbst die grundlegendsten Regeln des Straßenverkehrs vollständig zu vergessen. Vorerst sollen deshalb die kommunalen Winterdienste sämtliche Straßen mit extrem langsam fahrenden Räum- und Streufahrzeugen blockieren, bis sich die Lage wieder

weiss AND maße weiss AND masse

beruhigt hat.