

Prof. Dr. Agnès Voisard, Nicolas Lehmann

# Datenbank Systeme, SoeSe 2017

## Project 2.Iteration

TutorIn: Hoffman Christian

Tutorium 3, Gruppe 22

Ingrid Tchilibou, Emil Milanov, Boyan Hristov

2. Juni 2017

---

## Allgemein

Link zum Projekt: [https://github.com/gancia-kiss/dbs\\_projekt](https://github.com/gancia-kiss/dbs_projekt)

## 1.Aufgabe: Datenbankschema erstellen

Link zum .sql Datei:

[https://github.com/gancia-kiss/dbs\\_projekt/blob/master/DatabaseDump.sql](https://github.com/gancia-kiss/dbs_projekt/blob/master/DatabaseDump.sql)

Die letzte Iteration hatten wir ein paar 'Count' Attributen addiert, wir haben aber uns später entschieden, dass wir sie nicht brauchen. Diese Iteration, aber wir haben bemerkt, dass wir eigentlich diese Attribute brauchen.

Zusätzlich haben wir auch ein Count Attribut, damit wir aufzählen können, welche Paar von Hashtags am häufigsten vorkommt.

Wir haben uns entschieden selber IDs von Tweets zu erstellen und zwar mit Python UUID. Das Ergebnis kann größere Integer sein, deswegen haben wir als TWEET.ID 'bigint' als Datentyp benutzt.

Bei uns ist der Inhalt der Hashtags die primäre Schlüssel. Wir haben auch als Erinnerung der Attribut 'textlowercase' genannt, um zu wissen, nur kleine Buchstaben zu verwenden.

```
pg_dump election > DatabaseDump.sql      # Dump Datei erstellen
```

```
psql election < DatabaseDump.sql          # Datenbank vom .sql Datei importieren
```

## 2.Aufgabe: Datenbereinigung

Link zum Datenreinigungsprogramm:

[https://github.com/gancia-kiss/dbs\\_projekt/blob/master/programs/cleaner.py](https://github.com/gancia-kiss/dbs_projekt/blob/master/programs/cleaner.py)

Nach tiefere Datenanalyse haben wir zwei Hauptprobleme identifiziert:

1. Manche Tweets sind abgeschnitten und ein Teil von denen sind 'Truncated'. Wir wussten aber nicht, wie wir schnell die Tweets rekonstruieren können. Deswegen haben wir die gezählt und haben festgestellt, dass nur 2% der Tweets abgeschnitten sind. Es wäre für uns am leichtesten, alle solchen Tweets einfach zu löschen.

Wir haben bemerkt, dass alle abgeschnittenen Tweets auf '...' dann einen Link endeten. Es ist aber herausgekommen, dass '...' in solchen Tweets nur einen Unicode Zeichen war. Deswegen müssten wir im Python Program überprüfen, ob der Zeichen `u'\u2026'` im Körper des Tweets vorkommt.

2. Die vorgegebene Einstellungen der Libre Office Calculator hat versucht mit Codierung 'UTF-8' geöffnet. Dann gab es aber Probleme mit manchen Symbolen. Apostrophe, manche Sonderzeichen und vermutlich auch Emojis wurden nicht angezeigt. Nach kurze Analyse haben wir festgestellt, dass wenn wir die Tabelle mit Codierung 'Windows 1252' öffnen, werden Apostrophen und wichtige Sonderzeichen normal dargestellt, und Emojis und nicht für den Datensatz relevante Zeichen wurden gelöscht.

Die Python Program war sehr einfach. Wir haben die Standardbibliothek 'csv' benutzt um ein *csv\_writer* und *csv\_reader* zu erzeugen. Wir haben dann die .csv Datei gelesen, und Zeile für Zeile haben wir geprüft ob der Tweet Truncated ist, oder abgeschnitten ist. Wenn nicht, haben wir den Tweet in einer anderen Datei *test.csv* gespeichert.

## 3.Aufgabe: Datenimport

- Was & Wie?

Wir haben hier versucht zu beschreiben, wie das Programm für die Datenimportierung funktioniert. Den Code kann man ganz am Ende der Dokumentation finden. Wir haben versucht eine gute Modularisierung zu schaffen und gut lesbaren Code zu schreiben, deswegen sollen die hier nicht erklärte Funktionalitäten gut im Code verfolgbar sein.

- Grobe Struktur des Programms

Unser server schickt on StartUp ein Request für alle Hashtags zu dem DBS. Falls keine zurückgeliefert werden, müssen wir die Daten importieren.

Alles fängt bei TableParser.py an. Die schon bereinigte Datei data-cleaned.csv wird gelesen und es werden für jede Zeile in der Tabelle dynamisch die dazugehörige Tabellen als Dictionaries erstellt. Diese Module sorgt dafür, dass die dazugehörige Funktionen in Utils.py und Extractor.py aufgerufen werden und fasst dann alle Ergebnisse in Wörterbücher, die von dem DBController.py weiter bearbeitet werden. Hier wird definiert was in jede Tabelle reinkommt, es werden aber keine große Bearbeitungen gemacht.

Für die Bearbeitung der Daten ist Extractor.py verantwortlich. Diese Module macht alle Operationen, die Daten aus der Tabelle zu 'derived attributes' konvertieren. Hier werden die Hashtags aus dem Text genommen, das Datum in Woche umgewandelt und das Rating für ein Tweet kalkuliert.

Als wir schon erkennen, gibt es viele Abhängigkeiten zwischen das Programm, die Tabellen in den Datenbank und die CSV Dateien. Deswegen haben wir die Klasse Contract.py hergestellt, wo alle Strings definiert werden. Diese Klasse haben wir als eine Art von Struct benutzt, da es in Python keine Structs gibt.

Es gibt einige Umformungen, die in der Informatik sehr bekannt sind und ganz abstrakt sind, wie z.B. die Erzeugung von allen Paaren ohne Wiederholungen (2-Sets). Solche Methoden haben wir in der Module Utils.py definiert.

Am wichtigsten ist sicherlich die Module DBController.py. Hier wird die Verbindung mit dem DBS hergestellt und die ganze Kommunikation mit DBS implementiert. Es werden hier auch alle Seiteneffekten und Datenlogiken behandelt, wie z.B. inkrementierung von Count Attributen.

- Feststellung von in einem Tweet vorhandene Hashtags

Worttrennung ist eine schwierige Operation, als wir in Softwaretechnik gelernt haben. Deswegen sind wir davon ausgegangen, was wir bereits sicherlich kennen - nach dem Symbol '#' sind alle Buchstaben und Ziffern Teil eines Hashtags, alles anderes beendet das Hashtag. Deswegen haben wir, so lange noch '#' Symbol weiter in dem String zu finden ist, den String von '#' bis nach einem Symbol das nicht in 0-9, a-z und A-Z zu finden ist als Hashtag genommen und weiter iterativ nach '#' gesucht. Mehr kann man in Extractor::extractHashtags() erfahren

- Woche bestimmen

Wir haben die bereits in den Bibliotheken 'time' und 'datetime' Funktionalität benutzt, um erstmal ein Datenobjekt aus dem String zu erzeugen und danach davon den 'isocalender' zu erzeugen, was ein Trippel vom Jahr, Wochennummer und Tag ist. Dann haben wir den Tag entsprechend zu 1 (für 'Week::StartDate') oder zu 7 gesetzt und davon wieder ein Datum Objekt erzeugt.

- Rating bestimmen

Wie schon in der 1. Projektiteration erleutert, haben wir die Retweets und Favourites benutzt, um ein Float für den Rating zu erzeugen. Da alle Einträge Retweets und Favourites definiert haben, hatten wir hier keine Schwierigkeiten.

- Count setzen

Nachdem wir jede Zeile geparsed haben in TableParse.py, haben wir für den neuen Eintrag in Datenbank immer Count auf 1 gesetzt. Wir haben aber die Ausführung von dem SQL Expression in einem try-except (try-catch in Python) gewrapped. Falls eine Fehler erzeugt wird, haben wir überprüft, ob das konkrete Eintrag ein 'Count' Attribut hat und in dem Fall ein UPDATE SQL Expression ausgeführt wo 'Count' hochgesetzt wird. An der Stelle haben wir auch überprüft, ob die erzeugte Fehlermeldung nach dem Einfügung einer Woche erzeugt wurde und in dem Fall nichts gemacht (auch keine logs), da eine Woche nicht doppelt vorkommen soll und kein 'Count' hat.

Wir haben uns entschieden, keine SQL Expressions für Überprüfung auszuführen, da wir so uns eine SQL Ausführung sparen können, falls den entsprechenden Hashtag / Woche / Relation noch nicht erzeugt wurde. Sonst können wir trotzdem ein weiteres SQL Expression ausführen um Count hochzusetzen. SQL kann schon solche Situation gut behandeln, deswegen können wir uns die Überprüfungen sparen.

- Muster Abfragen

Wir haben schon eine Abfrage für die relevantesten Tweets implementiert, damit wir zeigen, dass die Daten erfolgreich implementiert wurden. Als schon für die 1. Projektiteration erwähnt, kann man alle da als Beispiel gegebene Abfragen gut mit unserem Datenmodell implementieren.

Wann man in dem Browser zu <http://127.0.0.1:5234/topTweets> navigiert, kann man schon die 10 relevantesten Tweets sehen (noch nicht schön visualisiert, aber irgendwie).

## 4. Aufgabe: Webserver

Link zum Webserver:

[https://github.com/gancia-kiss/dbs\\_projekt/tree/master/programs/server](https://github.com/gancia-kiss/dbs_projekt/tree/master/programs/server)

Wir haben uns entschieden selbe einen Webserver zu entwickeln und zwar wollten wir die 'Flask' Framework von Python zu verwenden. Die Installation war sehr einfach, und es ist sehr einfach ein kurzes Programm zu entwickeln.

```
1 from flask import Flask
3 app = Flask(__name__)
4 port = 5234
5 host = '127.0.0.1'
7 @app.route('/')
8 def index():
9     return 'Hello world'
11 if __name__ == '__main__':
12     app.run(debug=True, host=host, port=port)
```

Das Programm erzeugt einen einfachen Webserver auf Port 5234. Wenn man die Seite *localhost:5234* öffnet, dann sieht man unformatiert 'Hello world'.

Es wäre sehr einfach das zu erweitern. Wir können .html Dateien im selben ordner legen, und dann neue Wege mit *@app.route()* definieren. Vielleicht werden wir die Programme für Hashtag-analyse in den Server einbauen, damit wir dynamisch den Inhalt anpassen können.

## CODE

- server.py

```
1 from flask import Flask
2 from cleaner import cleanData
3 from DBController import DBController
4 from TableParser import TableParser
6 app = Flask(__name__)
7 port = 5234
8 host = '127.0.0.1'
10 @app.route('/')
11 def index():
12     return 'Hello world'
14 @app.route('/topTweets')
15 def getTopTweets():
16     topTweets = dbController.getTopTweets()
17     return 'Top tweets: <br><br>' + '<br><br>'.join(str(e) for e in topTweets)
19 if __name__ == '__main__':
21     dbController = DBController()
```

```

22 filled = dbController.checkFilled()

24 if not filled:
25     cleanData()
26     TableParser.parseTables()
27 else:
28     print('Data already imported :')

30 app.run(debug=True, host=host, port=port)

```

- TableParser.py

```

1 import csv
2 from DBController import DBController
3 from Extractor import Extractor
4 from Utils import Utils
5 from Contract import Contract

7 class TableParser:

9     @staticmethod
10    def getTweetFromEntry(entry):

12        return {

14            Contract.ADD_TO_TABLE_KEY: Contract.TABLE_TWEET,
15            Contract.ID_COLUMN: Utils.getRandom8ByteInt(),
16            Contract.AUTHOR_COLUMN: entry[Contract.HANDLE_ENTRY],
17            Contract.TEXT_COLUMN: entry[Contract.TEXT_ENTRY],
18            Contract.TIME_COLUMN: Extractor.extractTime(entry),
19            Contract.RATING_COLUMN: Extractor.calculateRatingForTweet(entry)
20        }

22    @staticmethod
23    def getHashTagsFromHTTexts(hashtagTexts):
24        hashtags = []

26        for ht in hashtagTexts:
27            hashtags.append(
28                {
29                    Contract.ADD_TO_TABLE_KEY: Contract.TABLE_HASHTAG,
30                    Contract.TEXT_LOWER_CASE_COLUMN: ht,
31                    Contract.COUNT_COLUMN: 1
32                }
33            )

35        return hashtags

37    @staticmethod
38    def getWeekFromEntry(entry):
39        return {
40            Contract.ADD_TO_TABLE_KEY: Contract.TABLE_WEEK,
41            Contract.START_DATE_COLUMN: Extractor.getWeekStart(entry),
42            Contract.END_DATE_COLUMN: Extractor.getWeekEnd(entry)
43        }

45    @staticmethod
46    def getUsedIns(week, hashtags):

48        usedIns = []
49        startDate = week[Contract.START_DATE_COLUMN]

51        for ht in hashtags:
52            usedIns.append(
53                {
54                    Contract.ADD_TO_TABLE_KEY: Contract.TABLE_USED_IN,

```

```

55         Contract.HASHTAG_TEXT_COLUMN: ht,
56         Contract.WEEK_START_DATE_COLUMN: startDate,
57         Contract.COUNT_COLUMN: 1
58     }
59 )

61     return usedIns

64 @staticmethod
65 def getPostedIn(week, tweet):
66     return {
67         Contract.ADD_TO_TABLE_KEY: Contract.TABLE_POSTED_IN,
68         Contract.TWEET_ID_COLUMN: tweet[Contract.ID_COLUMN],
69         Contract.WEEK_START_DATE_COLUMN: week[Contract.START_DATE_COLUMN]
70     }

72 @staticmethod
73 def getContains(tweet, hashtags):

75     contains = []
76     tweetid = tweet[Contract.ID_COLUMN]

78     for ht in hashtags:
79         contains.append(
80             {
81                 Contract.ADD_TO_TABLE_KEY: Contract.TABLE_CONTAINS,
82                 Contract.HASHTAG_TEXT_COLUMN: ht,
83                 Contract.TWEET_ID_COLUMN: tweetid
84             }
85         )

87     return contains

89 @staticmethod
90 def getUsedTogetherWithFromHTTexts(hashtagTexts):

92     pairs = Utils.getUniquePairs(hashtagTexts)
93     usedTogetherWiths = []

95     for pair in pairs:
96         usedTogetherWiths.append(
97             {
98                 Contract.ADD_TO_TABLE_KEY: Contract.TABLE_USED_TOGETHER_WITH,
99                 Contract.PRIMARY_HASHTAG_COLUMN: pair[0],
100                 Contract.TOGETHER_WITH_HASHTAG_COLUMN: pair[1],
101                 Contract.COUNT_COLUMN: 1
102             }
103         )

105     return usedTogetherWiths

107 @staticmethod
108 def parseRow(row, dbController):

110     hashtagTexts = Extractor.extractHashtags(row)

112     tweet = TableParser.getTweetFromEntry(row)
113     week = TableParser.getWeekFromEntry(row)
114     hashtags = TableParser.getHashTagsFromHTTexts(hashtagTexts)

116     usedIns = TableParser.getUsedIns(week, hashtagTexts)
117     postedIn = TableParser.getPostedIn(week, tweet)
118     contains = TableParser.getContains(tweet, hashtagTexts)
119     usedTogetherWiths = TableParser.getUsedTogetherWithFromHTTexts(hashtagTexts)

121     dbController.addMultiple(tweet, week, hashtags, usedIns, postedIn, contains,
        usedTogetherWiths)

```

```

123 @staticmethod
124 def parseTables():
125     dbController = DBController()

127     filepath = Contract.CSV_CLEAN
128     csvfile = open(filepath, 'r', encoding='cp1252')
129     csv_reader = csv.DictReader(csvfile, delimiter=';', quotechar='"')

131     for idx, row in enumerate(csv_reader):
132         TableParser.parseRow(row, dbController)
133         print("Parsed {} rows...".format(idx))

135     dbController.close()

```

- Contract.py

```

1 class Contract:

3     CSV_INITIAL = 'american-election-tweets.csv'
4     CSV_CLEAN = 'data-cleaned.csv'

6     TABLE_WEEK = 'week'
7     TABLE_TWEET = 'tweet'
8     TABLE_HASHTAG = 'hashtag'
9     TABLE_USED_IN = 'usedin'
10    TABLE_POSTED_IN = 'postedin'
11    TABLE_USED_TOGETHER_WITH = 'usedtogetherwith'
12    TABLE_CONTAINS = 'contains'

14    ADD_TO_TABLE_KEY = 'addToTable'

16    COUNT_COLUMN = 'count'
17    ID_COLUMN = 'id'
18    TWEET_ID_COLUMN = 'tweetid'
19    AUTHOR_COLUMN = 'author'
20    TEXT_COLUMN = 'text'
21    START_DATE_COLUMN = 'startdate'
22    END_DATE_COLUMN = 'enddate'
23    WEEK_START_DATE_COLUMN = 'weekstartdate'
24    TIME_COLUMN = 'time'
25    RATING_COLUMN = 'rating'
26    PRIMARY_HASHTAG_COLUMN = 'primaryhashtag'
27    TOGETHER_WITH_HASHTAG_COLUMN = 'togetherwithhashtag'
28    TEXT_LOWER_CASE_COLUMN = 'textlowercase'
29    HASHTAG_TEXT_COLUMN = 'hashtagtext'

31    TEXT_ENTRY = TEXT_COLUMN
32    HANDLE_ENTRY = 'handle'

34    IGNORE_DUPLICATES_IN_TABLES = [TABLE_WEEK]

```

- Extractor.py

```

1 import time
2 from datetime import datetime
3 import re

5 class Extractor:

7     RATING_WEIGHT = {
8         "retweetCount": 0.8,
9         "favouriteCount": 1.2
10    }

```

```

12 TIME_FORMAT = '%Y-%m-%dT%H:%M:%S'

14 @staticmethod
15 def calculateRatingForTweet(entry):
16     retweetCount = float(entry['retweet_count'])
17     favouriteCount = float(entry['retweet_count'])

19     return retweetCount*Extractor.RATING_WEIGHT['retweetCount'] + retweetCount*Extractor
        .RATING_WEIGHT['favouriteCount']

21 @staticmethod
22 def extractTime(entry):
23     time = entry['time']

25     return datetime.strptime(time, Extractor.TIME_FORMAT)

27 @staticmethod
28 def getNthDayOfWeek(entry, n):
29     time = Extractor.extractTime(entry)
30     isoCalender = time.date().isocalendar()

32     return datetime.strptime('{0} {1} {2}'.format(isoCalender[0], isoCalender[1], n), '%
        G %V %u').date()

34 @staticmethod
35 def getWeekStart(entry):
36     return Extractor.getNthDayOfWeek(entry, 1)

38 @staticmethod
39 def getWeekEnd(entry):
40     return Extractor.getNthDayOfWeek(entry, 7)

42 @staticmethod
43 def extractHashtags(entry):
44     text = entry['text']
45     searchedUntil = 0
46     hashtags = []

48     try:
49         while (True):

51             crntHashtagStart = text.index('#', searchedUntil) + 1

53             try:
54                 crntHashtagEnd = re.search('[^a-zA-Z0-9]', text[(crntHashtagStart):]).start() +
                    crntHashtagStart
55             except AttributeError as atrErr:
56                 crntHashtagEnd = len(text) - 1

58             crntHashtag = text[crntHashtagStart:crntHashtagEnd].lower()
59             searchedUntil = crntHashtagEnd

61             if not crntHashtag in hashtags:
62                 hashtags.append(crntHashtag)

64         except ValueError as err:
65             # no more hashtags found, we are ready
66             pass

68     return hashtags

```

- Utils.py

```

1 import uuid

```



```

3 class Utils:
4
5     @staticmethod
6     def getUniquePairs(arr):
7
8         if len(arr) < 2:
9             return []
10
11         pairs = []
12
13         for i, itemA in enumerate(arr):
14
15             for j, itemB in enumerate(arr[(i+1):]):
16
17                 crntPair = [itemA, itemB]
18                 crntPair.sort()
19                 pairs.append(crntPair)
20
21         return pairs
22
23     @staticmethod
24     def getRandom8ByteInt():
25         return uuid.uuid4().int & (1 << (8*8 - 1))-1
26
27     @staticmethod
28     def getValuesFromDict(entriesDict):
29         sqlValuesArr = []
30         for key in entriesDict:
31             escapedString = str(entriesDict[key]).replace("'", '\\\'').replace('\\', '\\\\')
32             sqlValuesArr.append('\\'{0}\\''.format(escapedString))
33         return ', '.join(sqlValuesArr)

```

- DBController.py

```

1 import psycopg2
2 import json
3 from Contract import Contract
4 from Utils import Utils
5
6
7 class DBController:
8
9     dbUser = "testuser"
10     dbName = "Election"
11     host = "localhost"
12     dbUserPassword = "testpass"
13
14     connectionStringTemplate = "dbname='{0}' user='{1}' host='{2}' password='{3}'"
15     connectionString = connectionStringTemplate.format(dbName, dbUser, host,
16                                                         dbUserPassword)
17
18     def __init__(self):
19         print("Connected to DB!")
20         self.connection = psycopg2.connect(DBController.connectionString)
21         self.cursor = self.connection.cursor()
22
23     def addToDB(self, table, entriesDict):
24
25         sqlInsertExpr = "INSERT INTO {0} ({1}) VALUES ({2})"
26         sqlInsertCommand = sqlInsertExpr.format(table, ",".join(list(entriesDict.keys())),
27                                                         Utils.getValuesFromDict(entriesDict))
28
29         try:
30             self.cursor.execute(sqlInsertCommand)
31         except Exception as err:

```

```

31     if Contract.COUNT_COLUMN in entriesDict:
32         entriesDictCopy = entriesDict.copy()
33         del entriesDictCopy[Contract.COUNT_COLUMN]

35         sqlUpdateExpr = "UPDATE {0} SET {1}={1}+1 WHERE ({2})"
36         sqlUpdateCommand = sqlUpdateExpr.format(table, Contract.COUNT_COLUMN, DBController
37             .getWhereConditionsForUpdate(entriesDictCopy))
38         self.connection.commit()
39         self.cursor.execute(sqlUpdateCommand)

40     elif table in Contract.IGNORE_DUPLICATES_IN_TABLES:
41         pass
42     else:
43         print("Something went wrong while adding new DB entry to table {0}: {1}".format(
44             table, err))

45     self.connection.commit()

47 def close(self):
48     self.connection.close
49     print("DB connection closed!")

51 def handleAddTable(self, data):
52     table = data[Contract.ADD_TO_TABLE_KEY]
53     del data[Contract.ADD_TO_TABLE_KEY]
54     self.addToDB(table, data)
55     # try:
56     #     self.addToDB(table, data)
57     # except Exception as err:
58     #     print("Something went wrong: " + str(err))

61 def addMultiple(self, *tableLists):

63     for tableList in tableLists:

65         if isinstance(tableList, list):

67             for table in tableList:
68                 self.handleAddTable(table)
69         else:
70             self.handleAddTable(tableList)

72 def checkFilled(self):

74     self.cursor.execute("SELECT * FROM {0}".format(Contract.TABLE_WEEK))
75     weeks = self.cursor.fetchall()

77     return len(weeks) > 0

79 def getTopTweets(self):
80     self.cursor.execute("SELECT * FROM {0} ORDER BY {1} LIMIT 10".format(Contract.
81         TABLE_TWEET, Contract.RATING_COLUMN))
82     return self.cursor.fetchall()

83 @staticmethod
84 def getWhereConditionsForUpdate(columnsDict):

86     conditions = []

88     for key in columnsDict:
89         conditions.append(key + '=' + "{0}".format(columnsDict[key]))

91     return ' AND '.join(conditions)

```