

Prof. Dr. Agnès Voisard, Nicolas Lehmann

Datenbank Systeme, SoeSe 2017

Project 1

TutorIn: Hoffman Christian

Tutorium Tutorium 3

Ingrid Tchilibou, Emil, Boyan Hristov

11. Mai 2017

1 Aufgabe: Projektdokumentation

Link zum Projekt: https://github.com/gancia-kiss/dbs_projekt

Das Team: Emil Milanov, Boyan Hristov, Ingrid Tchilibou
Wir sind alle Bachelor Informatik Studenten im 4. Semester.

Alle drei von uns haben gründlich das ER Modell besprochen und mit Aufmerksamkeit eine Entscheidung getroffen. Ingrid hat zusätzlich die Daten gründlich analysiert, Emil und Ingrid haben das relationales Modell erstellt und Boyan hat die Datenbank mithilfe von Postgres erstellt und eine Möglichkeit gefunden, eine Dumpdatei zu erstellen, als Beweis, dass wir eine Datenbank fertig haben. Zusätzlich haben Emil und Ingrid die Dokumentation verfasst.

2 Aufgabe: Explorative Datenanalyse

Beschreibung von der Datensatz

Wir haben eine Tabelle mit 10 Attributen von verschiedenen Tweets aus der America Election von 2016. Jede Tupel beschreibt genau :

Von wem war der Tweet gepostet? Dafür gibt es den Attribut *Handle*

Was war die Nachricht von diesem Tweet? Dafür ist der Attribut *Text*

War die Authentizität von dem Tweet verletzt (Das heißt ob der Tweet, der gepostet wird, von diesem Person geschrieben)? dafür wird der Attribut Original Author erstellt Außerdem hat jeder Tweet einen *Source-url* als Attribut (D.h. Wohin wurde der Tweet erstellt)

Wie viel wurde diesem Tweet weitergeleitet? *retweet-count*

Wie viel mal wurde diesem Tweet "gelikt" *favorite-count* ?

Andererseits gibt auch Tweets die ein *is-quote* Attribut haben, um zu sagen, dass 2 Tweets von verschiedenen Autoren zusammengefasst werden *is-quote status*

Zusätzlich gibt es auch manche Tweets, die 'truncated' worden sind. Das heißt, dass sie ein zitiertes Tweet oder irgendwelche Media enthalten, deswegen überschritten sie das Twitter Zeichenlimit und werden abgeschnitten.

3 Aufgabe: ER-Modellierung

Alle drei von uns haben dieses Modell gründlich besprochen und eine Entscheidung getroffen. Wir glauben, dass dieses Modell das beste für unsere Vorstellung des Projekts ist.

Probleme und Lösungen

- Um 'Welche Hashtags werden am meistens verwendet' zu beantworten, muss man für jeden Hashtag überprüfen, in welchen Tweets er vorkommt und dann diese aufzählen. Deswegen haben wir uns entschieden, in dem Hashtag Entitätstyp ein Attribut TotalCount zu speichern. So werden wir auch weniger Server-Anfragen bearbeiten müssen.
- Der Datenumfang ist über mehreren Monaten. Um die Fragen 'Wann traten insgesamt am meisten Hashtags auf?' und 'Wie hat sich die Häufigkeit der Verwendung eines speziellen Hashtags im Laufe der Zeit entwickelt?' zu beantworten müssen wir jeden Tweet zu einem zeitlichen Bereich zuordnen.

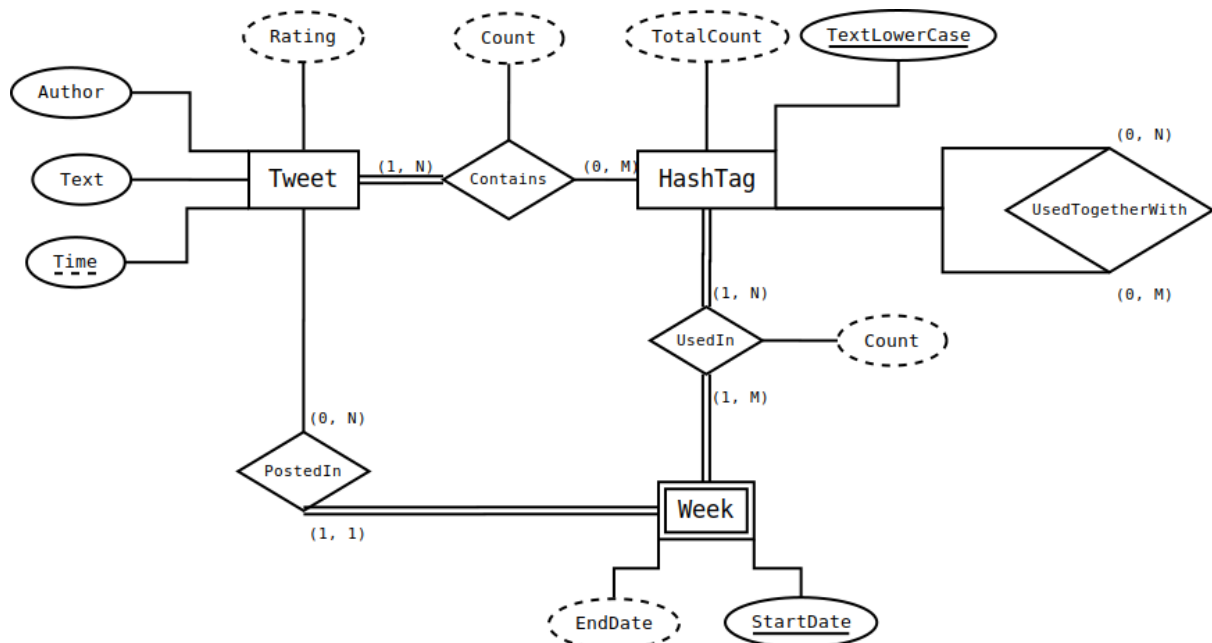
Wir haben uns deswegen entschieden, einen Entitätstyp 'Week' zu haben, wo wir an einer Woche benutzten Hashtags speichern werden. (Natürlich auch einen Attribut "Count" für höhere effizienz verwenden)

- Unserer Meinung nach, wäre eine Metrik für wichtige Tweets die Summe von Retweets und Favou-rites. So haben wir einen Attribut 'Rating', der nach der folgende Schema berechnet wird;

$$\text{retweets} * \text{weight}_1 + \text{favourites} * \text{weight}_2 = \text{rating}$$

wobei weight_1 und weight_2 später genau bestimmt werden.

- Es wird ziemlich aufwändig herauszufinden, welche Paar von Hashtags am häufigsten gemeinsam auftritt, weil wir über alle Tweets iterieren müssen. Deswegen haben wir eine N:M rekursive Relation auf 'Hashtag' erstellt.
- Eigentlich haben wir zu fast jedem Entitätstyp einen Attribut 'count' hinzugefügt, weil die Berechnung von Anzahl viel einfacher zu Parse-Zeit wäre, als jedes mal eine Anfrage zu machen, und dann alle verschiedene Entitäten aufzuzählen.



4 Aufgabe: Relationales Modell

Nachdem wir die ER-Relation hatten, war es einfach das relationale Modell zu erstellen. Es kann sein, dass wir nicht alle Attributen brauchen, oder nicht genug Attributen haben, da aber wir noch keine Daten im Datenbank haben, wäre es einfach, die Schema anzupassen.

```
Tweet(ID :: int, Handle :: char(20), Text :: char(200), Time :: Date,  
      Rating :: int, Count :: int, replyTo :: char(20),  
      isRetweet :: bool, isQuote :: bool, isTruncated :: bool )
```

```
Hashtag(ID:: int, Text :: char(30), TotalCount :: int)  
Week(ID:: int, StartDate :: Date, EndDate :: Date)
```

```
HashtagTweet(Tweet.ID, Hashtag.ID)  
WeeklyHashtag(Week.ID, Hashtag.ID, Count :: int)  
WeeklyTweers(Week.ID , Hashtag.ID)  
HashtagPairs(HT1.ID, HT2.ID)
```

5 Aufgabe: Datenbank erstellen