

Skeleton of final report

COMS3

Group 9

Tracking Interconnected Twitter Links
Using Graph Database Neo4j

Lindiwe, Clifford, Thomas

2017

Contents

1	Introduction	1
1.1	Overview	1
1.2	Neo4j	1
1.3	requirements	1
1.4	Running the code	1
2	Description	1
2.1	Design	2
2.2	Back-end	6
2.3	Front-end	6
2.4	application	6
2.5	looking forward	7
3	Members	7

1 Introduction

1.1 Overview

This purpose of this project is to graph Twitter data so that hidden trends and patterns may be revealed.

1.2 Neo4j

Neo4j is the leading graphic database management system, that distinguishes itself from other systems through its ease of use and its speed. Its was initially launched in 2010. Its fundamental design is to store data as nodes, edges and attributes. Nodes are connected by edges, both can have any number of attributes.

1.3 requirements

To run the code as desired the following programs or packages are required Neo4j v3.2.4

Python v2.7.13

PyCharm 2017.2.3

TextBlob v0.13.0

Tweepy v3.6.0

py2neo v3.1.2

[CipherQueries](#)(*GitHub*)

[Twitterneofinal](#)(*GitHub*)

1.4 Running the code

The specifics of running the code as desired can be found on [GitHub in the ReadMe](#) document

2 Description

We import public data from twitter and store it in a Neo4j database. This database is graphically represented through the Neo4j Browser, which shows nodes and connections, in different colours and sizes to highlight and differentiate various things. Filters can be applied to customise or limit what the graph shows. For example only tweets with a certain phrase or hashtag, "#neo4j", can be shown. This can give an idea about the age, location or gender of people tweeting about Neo4j.

The JSON for "Wonderfulvideo2" isn't shown, as it is very similar to that of "FightinFourteen", which is contained in the second set of bracket. Details such as friends, followers and time aren't shown on the graph, but hovering over an element, or selecting it, will reveal these and other properties.

2.1 Design

Figure 1 shows a demo of the final look.

Green nodes are people on twitter, labelled by their twitter handle, they share an edge labelled "created by" to red nodes, labelled "T", with anything they tweet. Should anyone retweet this tweet, then a "RT" node is created in purple, sharing "created by" edges with the original tweet and the retweter. A blue node is created in a similar fashion, when someone replies to or mentions someone else on Twitter. The purple node in the center is the topic we chose to filter the tweets through. We thought "Trump" should always be trending, although in different volumes, based on the time of day. The central topic shares an edge withh all shows tweet, labelled "generated by".

Every time the python script is run it increases the database size, the figure below was created after running the code just once. Twitter limits what we can access, running the script again an hour later brings in a similar quantity of data.

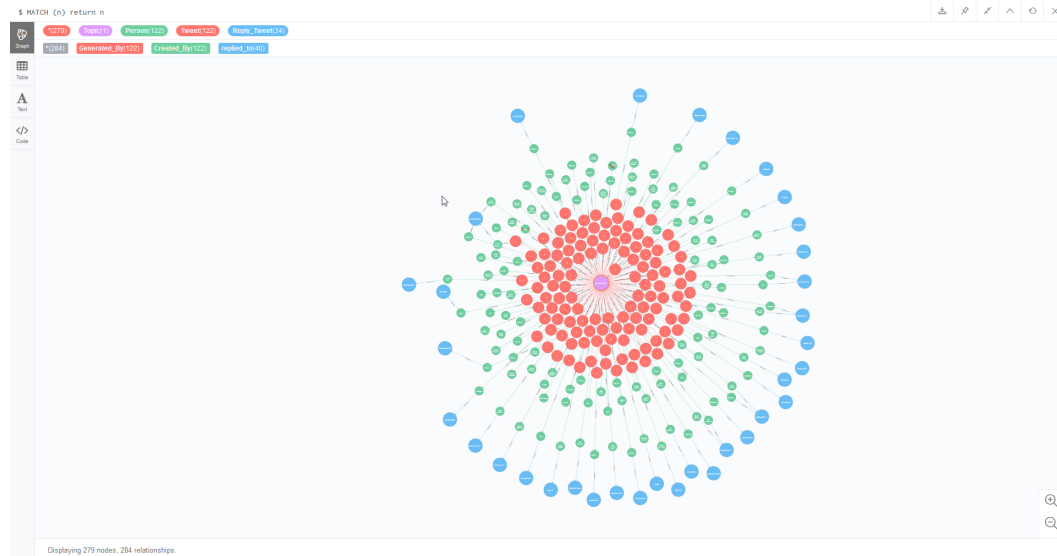


Figure 1: Neo4j Browser look

Figure 2 shows the data process, from Twitter to the final front-end interface.

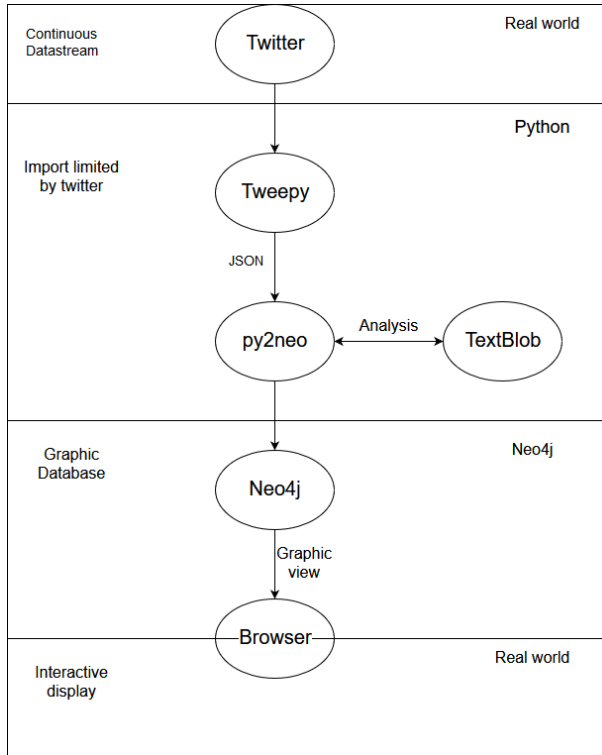


Figure 2: Data process

The Tweepy API collects data from the continuous flow of information that is Twitter, and converts it to a JSON string. This string is then processed to neo4j data and analysed by TextBlob. The data is then passed to Neo4j, where it is stored and displayed as a graph. The graph can be manipulated and filtered through the Neo4j Browser

Figure 3 shows the basic idea behind the set up of the database.

Database overview

How Neo4j stores data

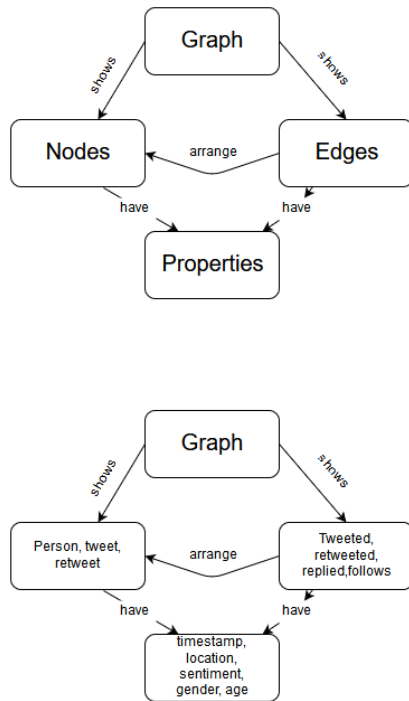


Figure 3: Neo4j database set up

Our nodes are people, tweets, retweets, replies and the central topic, they have various properties, such as age, gender, sex (where appropriate), time created, location. The edges "generated by", "created by" and "replied to" link the nodes. Edges can have properties too, properties not only provide additional information, that can be seen by clicking or hovering over a node/edge, they also can be used as filters to search/limit the data.

To emphasize how Neo4j simplifies the data being worked with to a graph Figure 4 shows the JSON string, the raw information, and the equivalent graph.

Neo4j

JSON raw string to graphic Neo4j view

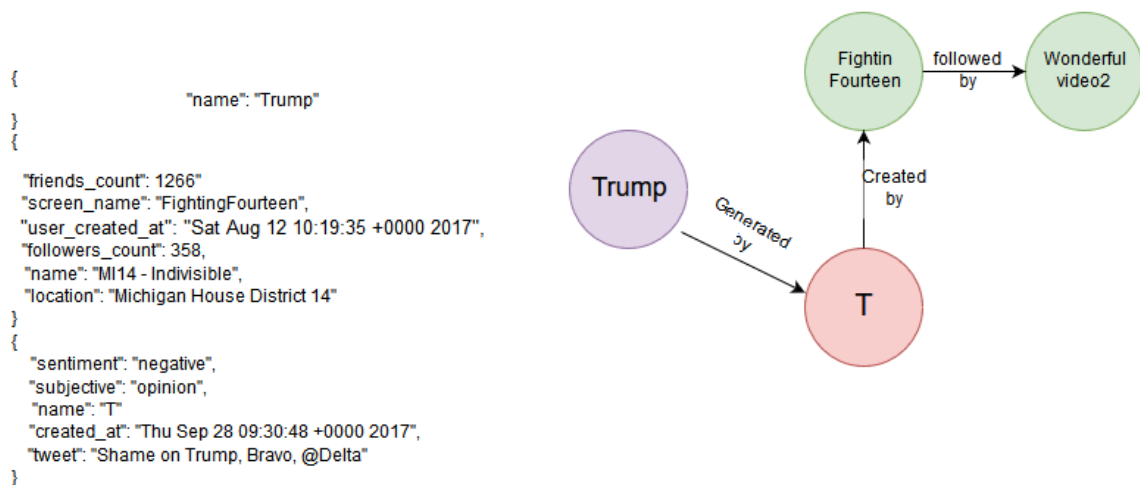


Figure 4: Raw data vs graph

The graph is a better presentation, that offers a better overview and logical structure of interconnections. This is mainly because it is easy to look at a graph showing the connection, than trying to imagine it for oneself.

Additionally the graph shows a lot less information at the first glance, data not relevant to the connection network aren't shown in the overview, but can easily be accessed by selecting a node/edge. The viewer can choose what information he wants, and doesn't need to sift through lines and lines of text.

We discussed including the location of where the tweet was made, however this is problematic for several reasons

- Not all users share their location with twitter
- Location isn't always accurate
- Different providers share locations differently (3 locations from our testing include, " #GameofThrones", "Basin, WY, United States", "Kenya")

Using a text analysis tool (TextBlob) we can look at the content of a tweet to see whether it's an opinion, report or advertising. We can also look at the sentiment of the content to determine whether a tweet is positive or negative.

2.2 Back-end

The back-end of the project imports data from Twitter to python as a JSON stream. The JSON is then analysed and converted to neo4j data. In the process Nodes and edges are created and properties are assigned. external analysis such as sentiment are done before the conversion to neo4j, and passed along as properties.' We don't have a server to host the Neo4j data, which means that when a different machine runs the code they will collect separate data from twitter, based on what is trending at the time. This can result in differences based on when and how often Twitter data is collected, however by using a topic like "trump" we hope to always have a fair amount of hits.

2.3 Front-end

The front end interface is the Neo4j browser. Unfortunately we don't have a server to host this, so system settings can't be saved, such as the graphic template for colour, size and the like. This also means that changes to the data being viewed are done as Ciper queries. The browser itself allows for movement of nodes, zooming and inspection of elements as well as style changes, all of which can be done from the tool bar around the view panel.

2.4 application

A tool like this has many applications, certain tweet patters could signal to an account being bot operated, it can be used to find out what a certain age, race, gender, location are tweeting about a certain topic. The first application that pops to mind is targeted advertising, but many more exist. The main advantage a graphic data base offers over a conventional one is that we humans can understand them better, leading to a better overview. This allows otherwise hidden pttterns to emerge. Once we know of a few hidden patterns, we can start writing algorithms to find and analys more.

2.5 looking forward

The first improvement we would like to make would be to host the Neo4j database on a server, so that visual designs and templates can be stored for a uniform look. The next would be a script, either through a website or a program to automate the initial set up in Python, so that the user just needs to select his preferences or enter a topic into a search field and the script does the rest and displays the results in the browser, this also eliminates having to manually enter Ciper queries.

3 Members

who did what