

# PDT - Neo4j

[https://github.com/MennoCoevoorn/PDT\\_6.git](https://github.com/MennoCoevoorn/PDT_6.git)

Marek Štrba

December 2021

## Obsah

- 1 Zoberte si nami vytvorený dataset z vašich tweetov a importujte ho cez neo4j-admin 4
- 2 Vypíšte 5 Accountov s najvyšším množstvom followerov. Nezoraďujte Accounty podľa poľa followers\_count, hodnota je prevzatá z Twitteru a nezodpovedá našim vzťahom v datasete. Zaujímajú nás followujúce Accounty v našom datasete cez vzťah FOLLOWS. 4
- 3 Nájdite najkratšie cesty medzi Katy Perry (screen\_name: 'katyperry') a Kim Kardashian (screen\_name: 'KimKardashian') cez vzťah FOLLOWS. Všetky cesty, kde Donald Trump followuje niekoho, kto followuje niekoho, kto..., kto followuje Katy Perry. 5
- 4 Vyhľadajte neúspešné tweety influencerov. Vyhľadajte 10 najmenej retweetovaných tweetov od Accountov, ktoré sú na prvých 10 miestach v celkovom počte retweetov. 6
- 5 Vytvorte volaním iba jednej query nový Account s Vaším menom, ktorý bude followovať Donalda Trumpa (screen\_name: "realDonaldTrump") a v tom istom volaní vytvorte tweet, ktorý bude retweetom Donaldovho najretweetovanejšieho tweetu. 7
- 6 Odporučte používateľovi (screen\_name: "777stl") followovanie ďalších Accountov, na základe followovania rovnakých Accountov: Vyhľadajte 10 Accountov, ktoré followujú najviac rovnakých Accountov ako náš používateľ, ale náš používateľ ich ešte nefollowuje. 8
- 7 Odporučte používateľovi (screen\_name: "DaynerWilson") followovanie ďalších Accountov na základe zhody v retweetovaní rovnakých tweetov: Vyhľadajte 10 accountov, ktoré retweetli najviac tých istých tweetov, ako náš používateľ. Ak tweet ktorý retweetujeme, je už tiež retweetom, rátajte za zhodu aj retweetovanie jeho parent tweetu – retweetovanie teda zohľadňujte rekurzívne. 9
- 8 Vyhľadajte 5 tweetov ostatných Accountov, ktoré do hĺbky 5 followujú account, ktorý napísal tweet (id: "1289380305728503808"), ktoré síce nie sú retweetom vybraného tweetu, ale napriek tomu majú čo najviac rovnakých slov v poli content zhodných s vybraným tweetom (stačí rozdeliť content na slová cez split(tweet.content, " "). Account, ktorý followuje Account, ktorý followuje nami vybraný Account rozumieme hĺbkou 2. Odporúčam pozrieť si procedúry v knižnici APOC pracujúce s collections, ale nie je to podmienkou na zvládnutie úlohy. 10

- 9 **BONUS:** Nájdite najkratšie cesty medzi Katy Perry (katyperry) a Donaldom Trumpom cez vzťah RETWEETS (a tým pádom aj POST). Všetky cesty, kde Katy Perry retweetla post Accountu, ktorý retweetol post Accountu, ktorý..., ktorý retweetol post Donalda Trumpa. 11

## 1 Zoberte si nami vytvorený dataset z vašich tweetov a importujte ho cez neo4j-admin

S importom som bojoval asi 3,5 hodiny (som ten spolužiak, ktorý plakal - nie naozaj ale nervy som mal). Nakoniec som s pomocou Martina prišiel na to, že problém bol v tom, že dump bol verzia 4.2.1 a ja som to chcel naimportovať do DMBS verzie 4.4.0. Za trest mu nedonesiem pivo z New Yorku.

```
D:\FIIT\Apps\N4J\relate-data\dbmss\dbms-014a5628-ded8-42e8-9ded-959eccebbd60>D:\FIIT\Apps\N4J\relate-data\dbmss\dbms-014a5628-ded8-42e8-9ded-959eccebbd60\bin\neo4j-admin load --from=D:\FIIT\ING\1_Semester\PDT\PDT_6\tweets2021.dump --database=tweets --force --verbose
```

Figure 1: Jeden z asi 11 importov, ktoré som robil. Finálny vyzeral rovnako, len nemal verbose a išiel priamo do neo4j databázy, nevytváral som dodatočne ďalšiu

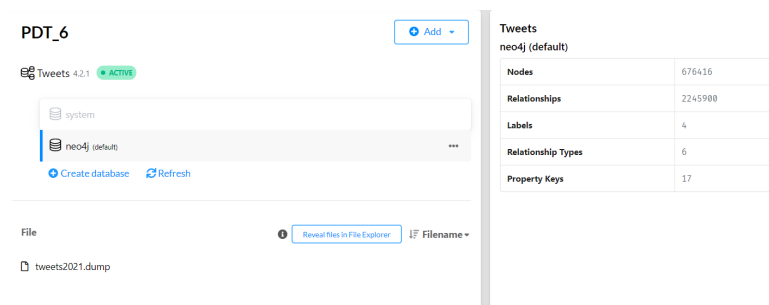


Figure 2: Ukážka úspešne naimportovaného dumpu

## 2 Vypíšte 5 Accountov s najvyšším množstvom followerov. Nezoraďujte Accounty podľa poľa followers\_count, hodnota je prevzatá z Twitteru a nezodpovedá našim vzťahom v datasete. Zaujímajú nás followujúce Accounty v našom datasete cez vzťah FOLLOWS.

Ako prvé som si matchol všetky accounts, následne som ich cez with spolu s výpočtom počtu vzťahov hocijaký node follows account a preniesol k order by followerCnt desc a to som limitom orezal na top 5. Výsledok aj query sú v priloženom screenshote.

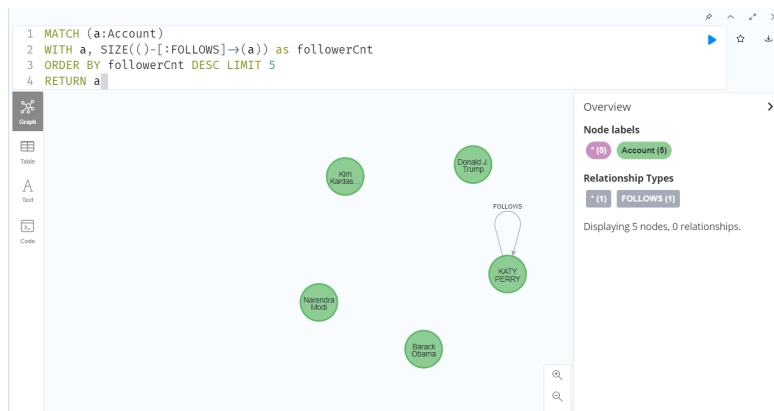


Figure 3: Ukážka query a výsledku

### 3 Nájdite najkratšie cesty medzi Katy Perry (screen\_name: 'katyperry') a Kim Kardashian (screen\_name: 'KimKardashian') cez vzťah FOLLOWS. Všetky cesty, kde Donald Trump followuje niekoho, kto followuje niekoho, kto..., kto followuje Katy Perry.

Dostať sa od Katy Perry ku Kim Kardashian bolo priamočiare, stačilo použiť funkciu `shortestPath` naviazanú na vzťah `FOLLOWS`, výsledok a query priložené v screenshote:

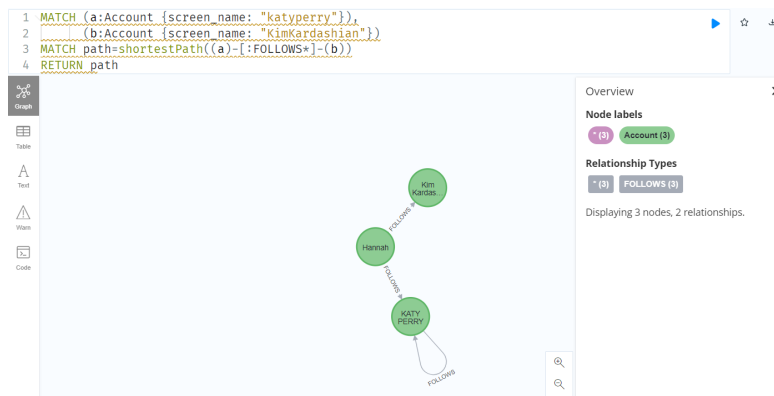


Figure 4: Katy -> Kim speedrun ukážka

Na to aby som zistil prepojenia hĺbky 4 medzi Donaldom a Katy som použil vzťah `FOLLOWS` s obmedzením na hĺbku 4 (dúfam, že som správne pochopil,

že to tak bolo myslené). Síce vo výsledkoch môžeme vidieť aj path Donald - BBC - Katy, ktorý má hĺbku 2 no ten je tam kvôli ceste Donald - Comport - Ashdax - BBC - Katy, čo je hĺbka 4.

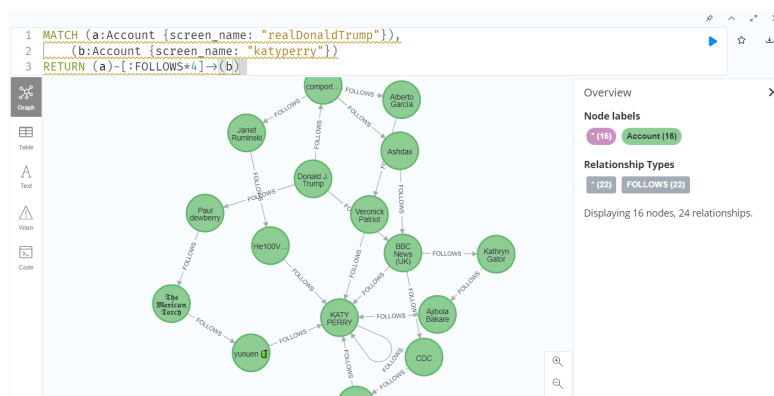


Figure 5: Ukážka Donaldových ciest ku Katy

#### 4 Vyhľadajte neúspešné tweety influencerov. Vyhľadajte 10 najmenej retweetovaných tweetov od Accountov, ktoré sú na prvých 10 miestach v celkovom počte retweetov.

Prvým krokom bolo zistiť, ktoré účty majú najviac retweetov. Preto som si matchol accounts cez post na tweets, to som si s with preniesol spolu so súčtom retweet relationships na dané tweety a zoradil som to podľa tohoto počtu s limitom 10, takto som dostal top infulencerov. Nasledne som si znova matchol tweety na týchto influencerov cez POSTS relationship. Zrátal im znova RETWEETS relationships a zoradil teraz opačne, od najmenšieho a limitoval znovu na 10. Tieto tweety som následne vrátil ako najneúspešnejšie tweety influencerov.

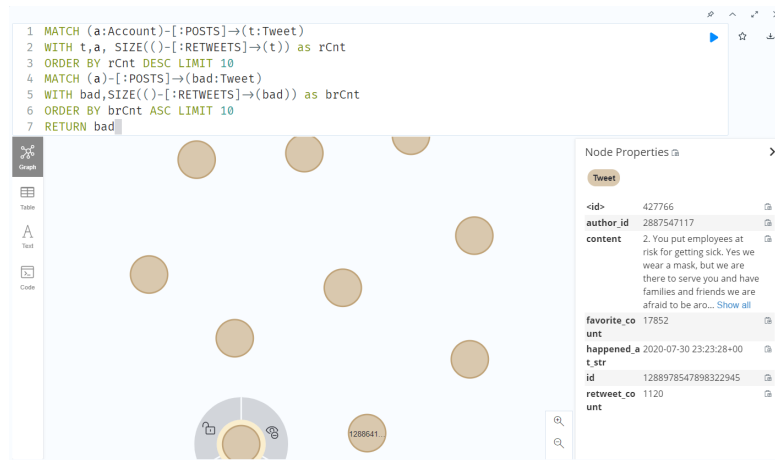


Figure 6: Query a výsledok najslabších tweetov

## 5 Vytvorte volaním iba jednej query nový Account s Vaším menom, ktorý bude followovať Donalda Trumpa (screen\_name:"realDonaldTrump") a v tom istom volaní vytvorte tweet, ktorý bude retweetom Donaldovho najretweetovanejšieho tweetu.

Prvým krokom bolo získať najretweetovanejší Donaldov tweet. Na to som si matchol všetky Donaldové tweety (vzťah POSTS), zrátal im RETWEETS relationship, zoradil, limitoval na 1 a tak som dostal jeho najretweetovanejší tweet. Tento tweet som si pomocou with preniesol do zvyšku query. Tam som si matchol Donalda, vytvoril som si svoj účet a svoj tweet. Potom som pomocou premenných vytvoril relationships, kde ja followujem Donalda, môj tweet retweetuje jeho tweet a ja postujem svoj tweet.

```

1 MATCH (a:Account {screen_name: 'realDonaldTrump'})-[:POSTS]->(t:Tweet)
2 WITH t, SIZE()-[:RETWEETS]->(t) as rCnt
3 ORDER BY rCnt DESC LIMIT 1
4 WITH t
5 MATCH (donald:Account {screen_name: 'realDonaldTrump'})
6 CREATE (me:Account {screen_name: 'MennoCoeHoorn', name: 'Marek Štrba', id: '18111998'})
7 CREATE (mytweet:Tweet {content: t.content, author_id: me.id})
8 CREATE (me)-[:FOLLOWS]->(donald)
9 CREATE (mytweet)-[:RETWEETS]->(t)
10 CREATE (me)-[:rp:POSTS]->(mytweet)

```

Added 2 labels, created 2 nodes, set 5 properties, created 3 relationships, completed after 466 ms.

Figure 7: Query vytvarajúca môj účet, tweet a potrebné relationships

Následne som vykonal kontrolné dopyty, na to či followujem správny účet a či som retweetol správny tweet:

```

1 MATCH (a:Account {screen_name: 'MennoCoeHoorn'})-[:FOLLOWS]->(d:Account)
2 RETURN d

```

```

{
  "identity": 28626,
  "label": {
    "Account": {
      "properties": {
        "followers_count": 58,
        "screen_name": "realDonaldTrump",
        "statuses_count": 1880,
        "followers_count": 8091476,
        "name": "Donald J. Trump",
        "description": "45th President of the United States of America",
        "id": "20873877"
      }
    }
  }
}

```

Started streaming 1 records after 1 ms and completed after 107 ms.

```

1 MATCH (t:Tweet {author_id: '18111998'})-[:RETWEETS]->(dt:Tweet)
2 RETURN dt

```

```

{
  "identity": 18808,
  "label": {
    "Tweet": {
      "properties": {
        "retweeted_at": "2020-07-31 18:00:12+00",
        "retweeted_count": 76776,
        "id": "1280482325999454",
        "author_id": "20873877",
        "content": "Never join the 1st border, and also some very good statements by Tony Fawc. Big progress being made! https://t.co/8W0u0u0u0u",
        "retweet_count": 18808
      }
    }
  }
}

```

Started streaming 1 records after 1 ms and completed after 122 ms.

## 6 Odporučte používateľovi (screen\_name:"777stl") followovanie ďalších Accountov, na základe followovania rovnakých Accountov: Vyhľadajte 10 Accountov, ktoré followujú najviac rovnakých Accountov ako náš používateľ, ale náš používateľ ich ešte nefollowuje.

Ako prvé som si matchol accounty, ktoré followuje stl, nazval som ich followers (mohlo byť aj lepšie meno ale tak čo už). Následne som matchol fofollowers (ľudí, ktorí followujú followerov, ale ako podmienku som pridal, že nie sú stl a ani jeden z followerov). Fofollowerov spolu s vypočítaným počtom vzťahov fofollower FOLLOWS follower som si pomocou with preniesol ďalej, tam som ich zoradil podľa countu zostupne a limitoval na 10. Výsledných fofollowerov som vrátil ako účty, ktoré by stl mal followovať.



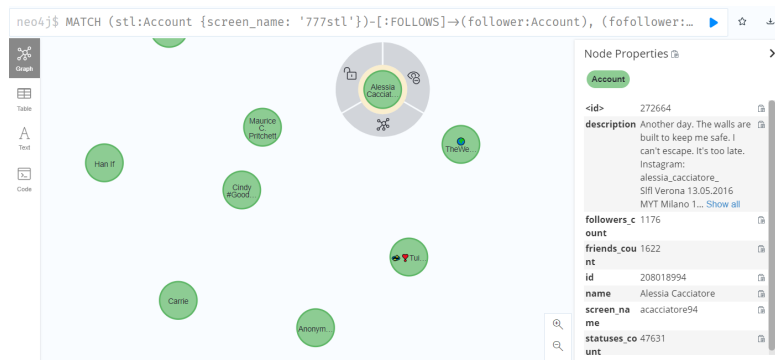


Figure 8: Query vracajúca recommendation pre používateľa 777stl

## 7 Odporučte používateľovi (screen\_name:"DaynerWilson") followovanie ďalších Accountov na základe zhody v retweetovaní rovnakých tweetov: Vyhľadajte 10 accountov, ktoré retweetli najviac tých istých tweetov, ako náš používateľ. Ak tweet ktorý retweetujeme, je už tiež retweetom, rátajte za zhodu aj retweetovanie jeho parent tweetu – retweetovanie teda zohľadňujte rekurzívne.

Keďže asi Dayner nemá záujem followovať sám seba a hľadá nových ľudí na followovanie, tak som do podmienky dal, že account nemôže byť ani Dayner a ani followovaný Daynerom. Následne som cez rekurzívny vzťah RETWEETS našiel všetky tweety, ktoré retweetujú rovnaké tweety ako Dayner. Cez vzťah POSTS som na dané tweety naviazal accounty. Potom som cez size a sum zisťoval, ktoré accounty majú najviac tweetov, ktoré retweetujú rovnaké tweety ako Dayner. Podľa tohoto súčtu som potom zoradil výsledok a limitoval ho na 10, čo ale nebolo potrebné, keďže mi vyšlo iba 7 accountov. Query a výsledok v priloženom screenshote:

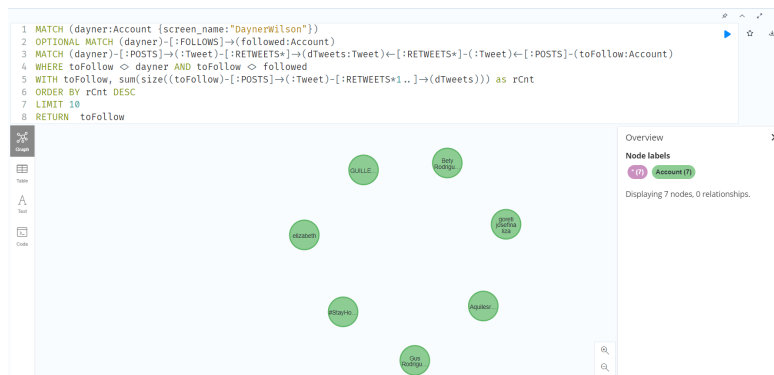


Figure 9: Noví kamaráti pre Daynera

- 8 **Vyhľadajte 5 tweetov ostatných Accountov, ktoré do hĺbky 5 followujú account, ktorý napísal tweet (id: "1289380305728503808"), ktoré síce nie sú retweetom vybraného tweetu, ale napriek tomu majú čo najviac rovnakých slov v poli content zhodných s vybraným tweetom (stačí rozdeliť content na slová cez `split(tweet.content, " ")`). Account, ktorý followuje Account, ktorý followuje nami vybraný Account rozumieme hĺbkou 2. Odporúčam pozrieť si procedúry v knižnici APOC pracujúce s collections, ale nie je to podmienkou na zvládnutie úlohy.**

Prvým krokom bolo nájsť accounty, ktoré cez vzťah FOLLOWS followujú autora nášho tweetu. Následne som zo všetkých tweetov, ktoré sú na dané accounty naviazané cez vzťah POSTS odfiltroval všetky, ktoré sú vo vzťahu RETWEETS s pôvodným tweetom. Potom som cez funkciu split rozdelil aj pôvodný tweet aj vyfiltrované tweety na listy slov, (splitter " "). Potom som pomocou `apoc.col.intersection` a `size()` vyrátal `splitScore` - intersection vracia pole, ktoré je prienikom dvoch polí a `size` vracia dĺžku poľa, čiže `splitScore` bol počet rovnakých slov. Následne som už len zoradil tweety zostupne podľa `splitScore` a limitoval výsledok na 5. Query a výsledok priložené v screenshots nižšie:

```

1 MATCH (a:Account)-[:POSTS]->(t:Tweet {id: "1289388305728503808"})
2 MATCH (followers:Account)-[:FOLLOWS*1..5]->(a)
3 MATCH (followers)-[:POSTS]->(ftweets:Tweet) WHERE NOT (ftweets)-[:RETWEETS]->(t)
4 WITH split(t.content, ' ') as tSplit, split(ftweets.content, ' ') as ftSplit, ftweets
5 WITH tSplit, ftSplit, ftweets, SIZE(apoc.coll.intersection(tSplit, ftSplit)) as splitScore
6 ORDER BY splitScore DESC
7 LIMIT 5
8 RETURN ftweets.id, splitScore

```

ftweets.id	splitScore
"1289382408753401865"	6
"128938944824862211"	6
"1289437406047993850"	5
"1289433063796148737"	5
"1289389590262699538"	5

Started streaming 5 records after 101 ms and completed after 200 ms.

Figure 10: Tweety s najlepším splitScore

## 9 BONUS: Nájdite najkratšie cesty medzi Katy Perry (katyperry) a Donaldom Trumpom cez vzťah RETWEETS (a tým pádom aj POST). Všetky cesty, kde Katy Perry retweetla post Accountu, ktorý retweetol post Accountu, ktorý..., ktorý retweetol post Donalda Trumpa.

Cesty som našiel pomocou `apoc.path.expandConfig`. Do `relationshipFiltera` som dal `POSTS` a `RETWEETS` (museli ostať obojstranné). Ako `terminatingNode` som dal Donalda, lebo k nemu sme sa chceli dostať. Postupne som nastavoval `maxLevel` na číslo, ktoré dalo aspoň 2 výsledky rôznej dĺžky aby som si vedel overiť funkčnosť zvyšku query. V produkcii by toto bolo byť veľké číslo alebo by sa ten atribút omitol a tým pádom by tam nebol limit. Toto mi vrátilo zoznam všetkých ciest do danej dĺžky. Cez min som zistil dĺžku najkratšej cesty a následne som znovu použil `apoc.path.expandConfig` s tým, že `maxLevel` bolo predtým získané minimum.

```

MATCH (katy:Account {screen_name: 'katyperry'}), (donald:Account {screen_name: 'realDonaldTrump'})
CALL apoc.path.expandConfig(katy, {
  relationshipFilter: "POSTS|RETWEETS",
  minLevel:1,
  maxLevel:15,
  terminatorNodes: [donald]
})
YIELD path
WITH katy,donald,path,length(path) as hops
WITH katy,donald,min(hops) as minHops
CALL apoc.path.expandConfig(katy, {
  relationshipFilter: "POSTS|RETWEETS",
  minLevel:1,
  maxLevel:minHops,
  terminatorNodes: [donald]
})
YIELD path
RETURN path

```

Figure 11: Celá volaná query

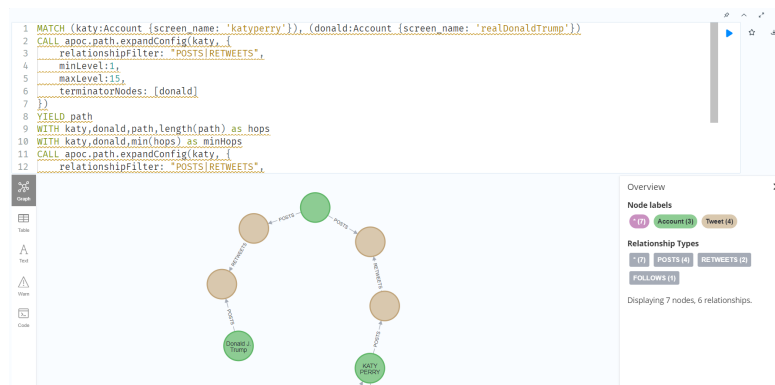


Figure 12: Výsledok volania query, najkratšia bola len 1 cesta a mala 6 hopov (relationshipov)