

Do kedy?

Keďže všetci robíte radi s predstihom, deadline je 31.12.2021 o 23:59.

Čo odovzdať?

Dotazy v separátnom textovom editore + dotaz a screenshot výsledku v pdf.

Hodnotenie

Úlohy spolu za 10 bodov rozdelené nasledovne:

- 1,2,3,4,5 a 6 sú za jeden bod,
- 7, 8 za 2 body.
- Bonus za 2 extra body.

Zadanie – Neo4j

V odpovediach nepoužívajte na zoradovanie existujúce polia `followers_count`, `retweet_count`, etc. Tieto hodnoty sú prevzaté z Twitteru a nezodpovedajú našim vzťahom v datasete. Pri počítaní početnosti followerov, retweetov a podobne vždy počítajte zodpovedajúce vzťahy medzi uzlami.

1. Zoberte si nami vytvorený dataset z vašich tweetov:

https://drive.google.com/file/d/1kufTAFP5DH8hiCT-BRLHqAicoulJ_0Qi/view?usp=sharing

a importujte ho cez neo4j-admin

2. Vypíšte 5 Accountov s najvyšším množstvom followerov. Nezoraďujte Accounty podľa poľa `followers_count`. Zaujímajú nás followujúce Accounty v našom datasete cez vzťah `FOLLOWS`.

3. Nájdite najkratšie cesty medzi Katy Perry {`screen_name`: 'katyperry'} a Kim Kardashian {`screen_name`: 'KimKardashian'} cez vzťah `FOLLOWS`. Všetky cesty, kde Kim Kardashian followuje niekoho, kto followuje niekoho, kto..., kto followuje Katy Perry.

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

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.

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.

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.

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 follow uje nami vybraný Account rozumieme hĺbkou 2. Odporúčam pozrieť si procedúry v knižnici APOC pracujúce s collections, ale nie je to podmienkou na zvládnutie úlohy.

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.

(Account) -> [POSTS] -> (Tweet) – [RETWEETS] -> (Tweet) <- [POSTS] – (Account) – [POSTS] -> (Tweet) – [RETWEETS] -> (Tweet) <- [POSTS] - (Account) -> ...

Github: <https://github.com/jaruij/neo4j-PDT>

1. Neo4j sa mi podarilo rozchodiť v dockeri, dump som musel definovať cez docker-compose.



2. Query je zapísané aj v separátnom textovom súbore, výsledok vyšiel nasledovný:

```
1 MATCH (acc1:Account)-[fol:FOLLOWS]-(acc2:Account) RETURN acc1.id as id, acc1.screen_name as name, acc1.followers_count as actual_follower_count, count(fol) as neo4j_follower_count
2 ORDER BY count(fol) DESC LIMIT 5
```

	id	name	actual_follower_count	neo4j_follower_count
1	"813286"	"BarackObama"	122782797	12725
2	"21447363"	"katyperry"	108521347	11461
3	"25873877"	"realDonaldTrump"	86091475	6720
4	"25366536"	"KimKardashian"	66740871	6680
5	"18839785"	"narendramodi"	62440782	6425

Started streaming 5 records after 2 ms and completed after 2694 ms.

Vypísal som si pre zaujímavosť aj pôvodný follower_count (zoraďoval som však podľa countu follow vzťahov z neo4j db).

```
1 MATCH (katyP:Account),(kimK:Account),
2 p = shortestPath((katyP)-[:FOLLOWS*]->(kimK))
3 WHERE katyP.screen_name = "katyperry" AND kimK.screen_name = "KimKardashian"
```

4. Úlohu som spravil cez subquery prostredníctvom CALL funkcie. Najprv si vyberieme 10 najviac retweetovaných účtov, následne nad touto množinou hľadáme 10 najmenej retweetovaných tweetov. Selectoval som zaujímavé hodnoty ako názov účtu, tweet_id, tweet content, počet retweetov účtu a počet retweetov samotného tweetu.

1 CALL

2 MATCH (acc:Account)-[p:POSTS]->(t:Tweet)-[ret:RETWEETS]-()

3 RETURN acc, count(ret) as acc_retweet_count ORDER BY acc_retweet_count DESC LIMIT 10

4

5 MATCH (acc)-[p:POSTS]->(t:Tweet)-[ret:RETWEETS]-()

6 RETURN acc.screen_name as name, acc_retweet_count, t.id as tweet_id, t.content as content, count(ret) as tweet_retweet_count ORDER BY tweet_retweet_count ASC

7 LIMIT 10

	name	acc_retweet_count	tweet_id	content
1	"ewarren"	1883	"1287140698131750913"	"Last night, the federal evictions moratorium expired, and rent is due next week—the same week coronavirus unemployment benefits are set to end. This is a completely preventable crisis. Congress must act immediately to extend these critical protections. https://t.co/9gALMcR3K "
2	"CarlosLoret"	1486	"1289243214428303360"	"Este estudio en el @washingtonpost dice que para reabrir las universidades habria que hacer una prueba de Covid cada dos dias a los alumnos https://t.co/69T4wIguXG "
3	"replouegohmert"	1613	"1288554276231548929"	"Lots of #FakeNews going around about this https://t.co/OQHYZJZvh "
4	"ewarren"	1883	"1289255470293757955"	"We know what we need to do to contain the virus and save lives and our economy—but Republicans refuse to invest enough in widespread testing and contact tracing. Trump and his Republican buddies don't have what it takes to get us out of this crisis. https://t.co/1CDH428mu "
5	"CarlosLoret"	1486	"1288962040401399815"	"Para morirse de envidia: cómo Francia está viviendo su nueva normalidad. Este articulista sale a restaurantes, va a conciertos y a centros comerciales. El truco: pruebas y rastreo de contactos. https://t.co/XcVG6T7E0P "

Started streaming 10 records after 79 ms and completed after 935 ms.

Výsledky:

"name"	"acc_retweet_count"	"tweet_id"	"content"	"tweet_retweet_count"
"ewarren"	1883	"1287140698131750913"	"Last night, the federal evictions moratorium expired, and rent is due next week—the same week coronavirus unemployment benefits are set to end. This is a completely preventable crisis. Congress must act immediately to extend these critical protections. https://t.co/9gALMcR3K "	1
"CarlosLoret"	1486	"1289243214428303360"	"Este estudio en el @washingtonpost dice que para reabrir las universidades habria que hacer una prueba de Covid cada dos dias a los alumnos https://t.co/69T4wIguXG "	2
"replouegohmert"	1613	"1288554276231548929"	"Lots of #FakeNews going around about this https://t.co/OQHYZJZvh "	3
"ewarren"	1883	"1289255470293757955"	"We know what we need to do to contain the virus and save lives and our economy—but Republicans refuse to invest enough in widespread testing and contact tracing. Trump and his Republican buddies don't have what it takes to get us out of this crisis. https://t.co/1CDH428mu "	3
"CarlosLoret"	1486	"1288962040401399815"	"Para morirse de envidia: cómo Francia está viviendo su nueva normalidad. Este articulista sale a restaurantes, va a conciertos y a centros comerciales. El truco: pruebas y rastreo de contactos. https://t.co/XcVG6T7E0P "	6

"ewarren"	1883	"1289287273880547328"	"We need to make sure schools have all the resources they need to determine whether and how to safely reopen. Anything less is recklessly endangering lives for political gain. https://t.co/XWqQ7yRf4 "	10
"CarlosLoret"	1486	"1289354415372025858"	"688 fallecimientos documentados en 24 horas, ya son 46 mil 688 decesos por #Covid en México. https://t.co/gH6a8hYU2 "	12
"CarlosLoret"	1486	"1289417499654541312"	"Murio Paco Valverde, un gran luchador por la naturaleza, valiente defensor de la Vaguita Marina. Hubo una enorme solidaridad para tratar de salvarlo. Gracias a todos los que estuvieron pendientes. Descanse en Paz el buen pescador. Abrazo entrañable para Alan y toda su familia. https://t.co/tazvILqgo "	50
"maddievelasco"	3143	"1288978547898322945"	"2. You put employees at risk for getting sick. Yes we wear a mask, but we are there to serve you and have families and friends we are afraid to be around now because we don't know how long ago we came in contact with someone or if we are infected until it's too late."	54
"maddievelasco"	3143	"1288978061501722629"	"1. You put yourself at unnecessary risk of contracting covid-19. We sanitize as often as we can, but you still take your mask off to eat around strangers and you don't know where they have been or who they have been in contact with."	66

5. V tejto úlohe bolo potrebné najprv nájsť najviac retweetovaný tweet účtu Donalda Trumpa prostredníctvom subquery (CALL rovnako ako predošlá úloha). Následne pomocou CREATE môžeme vytvoriť nový účet, ktorému cez vzťah FOLLOWS pridáme follow na Trumpov účet a cez vzťah POSTS vytvoríme nový tweet s ľubovoľným obsahom.

```

1 CALL {
2   MATCH (acc1:Account)-[p:POSTS]→(t:Tweet)←[ret:RETWEETS]-()
3   WHERE acc1.screen_name = 'realDonaldTrump'
4   RETURN acc1, t, count(ret) as neo4j_retweet_count
5   ORDER BY count(ret) DESC
6   LIMIT 1
7 }
8 CREATE(acc2:Account {screen_name: 'jaruji', name: 'Juraj'})-[f:FOLLOWS]→(acc1), (acc2)-[p:POSTS]→(myT:Tweet {content: "Hello TRUMP"})←[r:RETWEETS]-()
9 RETURN acc1.screen_name, neo4j_retweet_count, t.id, t.content, acc2.screen_name, myT.content

```

acc1.screen_name	neo4j_retweet_count	t.id	t.content	acc2.screen_name	myT.content
"realDonaldTrump"	487	"1289260338135699456"	"Great job by Jim Jordan, and also some very good statements by Tony Fauci. Big progress being made! https://t.co/8Oeca9H3yq "	"jaruji"	"Hello TRUMP"

Added 2 labels, created 3 nodes, set 3 properties, created 3 relationships, started streaming 1 records after 60 ms and completed after 235 ms.

6. Cez CALL som si vybral účet používateľa 777stl a účty ktoré followuje. Následne som hľadal iné účty ktoré followujú tieto účty a počítal prekryv, resp. početnosť zhodných účtov ktoré followujú s prvotne vybraným účtom. Bolo potrebné ošetriť aj situáciu, aby sme nenašli rovnaký účet (cez $acc1 \neq acc2$).

```

1 CALL {
2   MATCH (acc1:Account)-[fo1:FOLLOWS]→(accounts:Account)
3   WHERE acc1.screen_name = '777stl'
4   RETURN acc1, accounts
5 }
6 MATCH (acc2:Account)-[fo2:FOLLOWS]→(accounts)
7 WHERE acc1 <> acc2 AND NOT (acc1)-[:FOLLOWS]→(acc2)
8 RETURN acc2.screen_name as screen_name, count(fo2) as identical_follow_count
9 ORDER BY count(fo2) DESC
10 LIMIT 10

```

screen_name	identical_follow_count
"aa49677901"	2
"TwitteraActivat1"	1
"Steelersgirl690"	1
"aacoolatore94"	1
"jessiear_wurjant"	1
"moesnap78"	1
"ACrowdedPlanet"	1
"heruyaheru"	1
"han011014"	1
"Anonymous5426"	1

MAX COLUMN WIDTH:

7. Pomocou CALL som si vybral všetky tweety ktoré retweetuje používateľ DaynerWilson, následne som hľadal zhodu prostredníctvom tejto predom vybranej množiny tweetov. Výsledkov mi našlo iba 7.

```

1 CALL{
2   MATCH (acc1:Account)-[:POSTS]->(:Tweet)-[:ret1:RETWEETS*]->(t:Tweet)
3   WHERE acc1.screen_name = 'DaynerWilson'
4   RETURN acc1, t
5 }
6 MATCH (acc2:Account)-[:POSTS]->(:Tweet)-[:ret2:RETWEETS*]->(:t)
7 WHERE acc1 <> acc2
8 RETURN acc2.screen_name, count(ret2) as identical_retweets_count
9 ORDER BY count(ret2) DESC
10 LIMIT 10
11
12

```

acc2.screen_name	identical_retweets_count
"Aguilezreyes9"	3
"GusRodr05589797"	3
"GoretliLiza"	3
"BetyRod50219672"	2
"OCCARSOLIO"	1
"elizagomalcsla"	1
"viralvideovlogs"	1

MAX COLUMN WIDTH:

8. V prvom rade bolo potrebné spojiť APOC, keďže treba mať v plug-ins .jar file inak tie funkcie nezbehnú. Nainštaloval som to cez docker CLI príkazom `wget` na najnovšiu APOC verziu a následne reštartom neo4j. Pri úlohe som postupoval nasledovne: najprv som si vybral tweet, všetky účty ktoré followujú účet ktorý tweetol tweet do hĺbky 5 (neviem či mám dávať 0-4 alebo 1-5 hĺbku podľa toho čo je napísané v zadani - dával som 1-5 ale skúšal som aj 0-4 a výsledky boli trochu odlišné, 1-5 našlo o jednu 6 slovnú zhodu viac) a všetky tweety týchto účtov, ktoré nie sú retweetom nášho pôvodného tweetu. Nasledovne som cez APOC funkciu `intesection` vyhľadával počet zhodných slov medzi dvoma listami (získanými cez `split` s parametrom `" "`) a podľa počtu zhôd som následne našiel najpodobnejšie tweety. Najlepší match bol 6 slov.

```

1 MATCH (acc1:Account)-[:POSTS]->(t1:Tweet),
2 (acc2:Account)-[:FOLLOWS*1..5]->(acc1),
3 (acc2)-[:POSTS]->(t2:Tweet)
4 WHERE t1.id = '1289380305728503808' AND NOT (t2)-[:RETWEETS*]->(t1)
5 RETURN t1.content as tweet1, t2.content as tweet2, apoc.coll.intersection(split(t1.content, " "), split(t2.content, " ")) as matching_words,
6 size(apoc.coll.intersection(split(t1.content, " "), split(t2.content, " "))) as matching_words_count
7 ORDER BY matching_words_count DESC
8 LIMIT 5

```

tweet1	tweet2	matching_words	matching_words_count
"\$JadeRhinoes @eturleye @dougluocoy @SenMcSallyAZ Every decision they have made regarding COVID has been economically driven, which has allowed the virus to ravage our communities, and yet, when it came down to actually helping our economy, they abandoned our state's unemployed."	"\$PolitiBunny From what I've read, the teachers will be much more likely to get COVID from other staff than they would from a pupil. \n countries in Europe have schools open w/ no big problems. Day care facilities have been open with no big problems."	["been", "COVID", "have", "they", "the", "to"]	6
"\$JadeRhinoes @eturleye @dougluocoy @SenMcSallyAZ Every decision they have made regarding COVID has been economically driven, which has allowed the virus to ravage our communities, and yet, when it came down to actually helping our economy, they abandoned our state's unemployed."	"\$uacon @SafeReturnUA @UofAlabama This law only applies to people actually infected by COVID, and not to the audience of your heinous email, which referenced "alternative arrangements" for day care and the like. So no. That wasn't your intent, or your email would have explicitly said "if you get sick, etc"	["and", "have", "which", "the", "actually", "to"]	6
"\$JadeRhinoes @eturleye @dougluocoy @SenMcSallyAZ Every decision they have made regarding COVID has been economically driven, which has allowed the virus to ravage our communities, and yet, when it came down to actually helping our economy, they abandoned our state's unemployed."	"RT @natesharoyy: my mom's friend in India called to tell her that her husband insulted the dinner she made & called it bland and tasteless..."	["and", "made", "it", "the", "to"]	5
"\$JadeRhinoes @eturleye @dougluocoy @SenMcSallyAZ Every decision they have made regarding COVID has been economically driven, which has allowed the virus to ravage our communities, and yet, when it came down to actually helping our economy, they abandoned our state's unemployed."	"RT @natesharoyy: my mom's friend in India called to tell her that her husband insulted the dinner she made & called it bland and tasteless..."	["and", "made", "it", "the", "to"]	5
"\$JadeRhinoes @eturleye @dougluocoy @SenMcSallyAZ Every decision they have made regarding COVID has been economically driven, which has allowed the virus to ravage our communities, and yet, when it came down to actually helping our economy, they abandoned our state's unemployed."	"RT @dougmur: the United States is the only country that still has a corona virus problem and Trump wants to focus on tik tok?"	["virus", "and", "have", "the", "to"]	5

BONUS:

Podarilo sa mi nájsť najkratšiu cestu, no nevedel som ako použiť `shortestPath` ak chceme mať komplexnú sekvenciu relationships, preto som vyhľadával medzi množinou tweetov katyP a donaldaT najkraťšiu cestu cez retweets vzťah. Dĺžka cesty bola 1 (katyP retweetla účet ktorý retweetol donalda s obsahom tweetu "Cool Retweet donaldovho Tweetu"). Najkratšia cesta bola iba jedna.

```

1 CALL {
2   MATCH (donaldT:Account)-[:POSTS]->(dtt:Tweet)
3   WHERE donaldT.screen_name = 'realDonaldTrump'
4   return donaldT, dtt
5 }
6 CALL {
7   MATCH (katyP:Account)-[:POSTS]->(t:Tweet)-[:RETWEETS*]->(t:Tweet)-[:POSTS]->(Account)-[:POSTS]->(kpt:Tweet)
8   WHERE katyP.screen_name = 'katyperry'
9   return katyP, kpt
10 }
11 MATCH p = AllShortestPaths((dtt)-[:RETWEETS*]->(kpt))
12 RETURN p

```

Overview

Node labels

(2) Tweet (2)

Relationship Types

(1) RETWEETS (1)

Displaying 2 nodes, 1 relationships.

"p"

```
[{"happened_at_str": "2020-03-28 02:29:57+00", "favorite_count": 122342, "id": "1243726993537073152", "author_id": "25073877", "content": "I love Michigan, one of the reasons we are doing such a GREAT job for them during this horrible Pandemic. Yet your Governor, Gretchen \"Half\" Whitmer is way in over her head, she doesn't have a clue. Likes blaming everyone for her own ineptitude! #MAGA", "retweet_count": 23858}, {}, {"content": "Cool Retweet donaldovho tweetu"}]
```