

UNIVERSITÀ DEGLI STUDI DI VERONA

DIPARTIMENTO DI INFORMATICA

Corso di Laurea in Informatica

Tesi di Laurea

**Un approccio temporale per
determinare la proprietà di
notizie diffuse su Twitter**

Relatore

prof. ssa Elisa Quintarelli

Laureando

Nicola AGRESTI
VR407685

ANNO ACCADEMICO 2019 - 2020

Indice

Dataset	3
Descrizione	3
Progettazione Logica	5
Personalizzazione	6
Database	9
Creazione	9
Upload Dati	10
Grafo	12
Struttura	12
Creazione	14
Salvataggio	23
Caricamento	23
Algoritmo Path-Consistency	24
Gestione	30
Sitografia	32

Dataset

Descrizione

Il dataset di partenza scelto (si veda [1]) contiene i Tweet degli utenti che hanno usato i seguenti hashtags: #coronavirus, #coronavirusoutbreak, #coronavirusPandemic, #covid19, #covid_19, #epitwitter, #ihavecorona, #Stay-HomeStaySafe, #TestTraceIsolate, a partire dal giorno 30/04/2020 alle ore 00:00:00 fino al giorno 01/05/2020 ore 01:59:00. Esso è composto da un unico file in formato CSV che contiene 22 colonne e 355387 righe. In dettaglio le 22 colonne sono:

1. status_id: un id univoco che identifica un tweet
2. user_id: un id univoco che identifica un utente
3. created_at: un timestamp di quando è stato creato un tweet
4. screen_name: il nome utente collegato allo user_id
5. text: il testo del tweet
6. source: la piattaforma usata per creare il tweet
7. reply_to_status_id: l'id del tweet a cui il tweet è in risposta
8. reply_to_user_id: l'id dell'utente a cui il tweet è in risposta
9. reply_to_screen_name: il nome utente a cui il tweet è in risposta
10. is_quote: un booleano per descrivere se il tweet è di tipo quotes
11. is_retweet: un booleano per descrivere se il tweet è di tipo retweet
12. favorites_count: il numero di mi piace ricevuti
13. retweet_count: il numero di retweet ricevuti
14. country_code: il codice dello stato da dove è stato creato il tweet
15. place_full_name: il nome della città da dove è stato creato il tweet
16. place_type: una descrizione del tipo di città da dove è stato creato il tweet

17. `followers_count`: il numero dei followers dell'account che ha creato il tweet
18. `friends_count`: il numero di amici dell'account che ha creato il tweet
19. `account_lang`: la lingua usata dall'account che ha creato il tweet
20. `verified`: un booleano per descrivere se è un account verificato
21. `lang`: la lingua usata nel testo del tweet

Alcune di queste informazioni non sono utili per quanto riguarda l'obiettivo di questa tesi, di conseguenza sono state eseguite delle scelte progettuali per capire quali attributi sono da mantenere nel dataset.

Progettazione Logica

Prima di andare a personalizzare il dataset è stato ideato un opportuno schema ER che potesse descrivere al meglio il problema, lo schema ER risultante è: Mentre lo schema relazionale risulta:

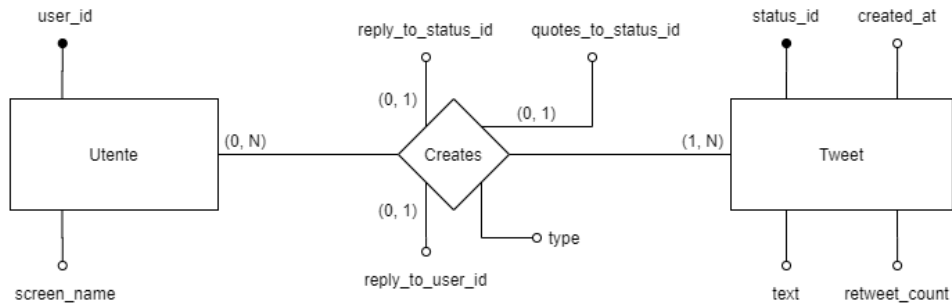


Figura 1. Schema ER

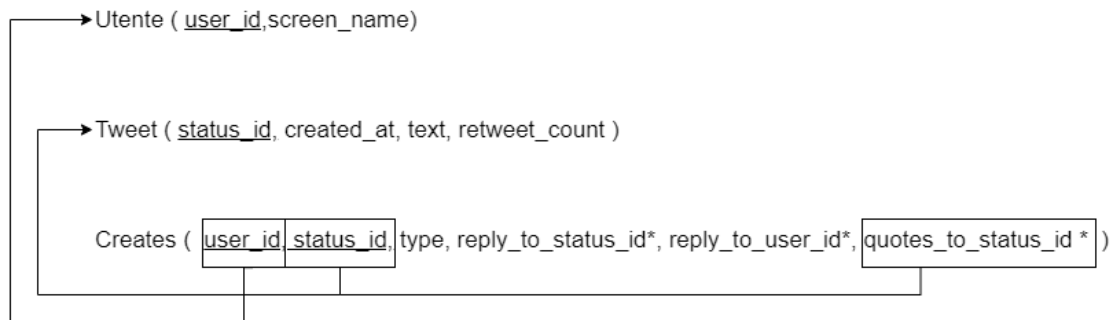


Figura 2. Schema Relazionale

Entrambe le figure sono state realizzate con ER.drawio [3].

Personalizzazione

In base alla progettazione logica ho creato 3 file csv, per prima cosa ho aperto con Excel il dataset originale e quindi ho copiato l'intera colonna `user_id` e l'intera colonna `screen_name` in un nuovo file denominato "Utente.csv", poi ho controllato tramite un' apposita funzione di Excel che non ci fossero duplicati nella chiave primaria "user_id".

	A	B
1	user_id	screen_name
2	3171712086	RahulGandhi
3	93794912	seoulmania
4	705694814612938752	NCDGov
5	102071743	GovWhitmer
6	16989178	JamesOKeefeIII
7	531041640	KKMPutrajaya
8	1030099481533014017	MazzoleniJulio
9	147994804	rishibagree
10	21059255	tedlieu

Figura 3. Esempio preso dal file utente.csv

Ho poi proseguito copiando le colonne `status_id`, `created_at`, `text` e `retweet_count` in un altro file di nome "Tweet.csv", ho controllato nuovamente la presenza di eventuali duplicati nella chiave primaria `status_id`.

	A	B	C	D
1	status_id	created_at	text	retweet_count
2	1255701270901387264	30/04/2020 05:31	A conversation with Dr Raghuram Rajan, former RBI Governor, on dealing with the #Co	10146
3	1255852107673907202	30/04/2020 15:30	i,~iſCEi ê+*e -eſ" è'ſ€eſ€Eâ+*eſ'iCEi€ ê*€ê'CEi iſµê<`è«dŸ' —#ſ*e -_e*~e'ê'e_j")i...iž	9521
4	1255993387603329026	01/05/2020 00:52	204 new cases of #COVID19 reported;80-Kano45-Lagos12-Gombe9-Bauchi9-Sokoto7-Bc	8852
5	1255689432738349059	30/04/2020 04:44	At 9 AM today, you can watch my conversation with Dr Raghuram Rajan, former RBI Go	8507
6	1255978307641905154	30/04/2020 23:52	lâ€™ve said it before, and lâ€™ll say it again â€” Michigan is an extraordinary place to l	6961
7	1255976688594419718	30/04/2020 23:45	UNBELIEVABLE! Our team just caught @Twitter red handed stealing retweets off our né	5504
8	1255778934668398593	30/04/2020 10:40	Terkini 30 April #COVID19 ðŸ†ðŸ†#Kes baharu 5âŒž7âŒž€ (25 kes import)Kematian	4629
9	1255682896020541440	30/04/2020 04:18	Informe #COVID19: hoy procesamos 563 muestras, 10 dieron positivo: 5 del exterior, 1	4400
10	1255857575783919617	30/04/2020 15:52	Didi Stays in a Gated Community at 41st floor.Now tell me which beggar/ Unknown per	4239

Figura 4. Esempio preso dal file tweet.csv

Per quanto riguarda il file `creates.csv` ho copiato le colonne `user_id`, `status_id`, `reply_to_status_id`, `reply_to_user_id` e successivamente ho aggiunto una colonna nuova di nome "type" che conterrà una delle seguenti tipologie di tweet:

- TW: un tweet originale creato da un utente

- RE: il tweet creato è in risposta ad un altro tweet
- QT: il tweet creato è un retweet con un commento
- REQT: il tweet creato è un retweet con un commento in risposta ad un altro tweet
- RT: un retweet di un tweet di uno dei 4 tipi precedenti

	user_id	status_id	type	reply_to_status_id	reply_to_user_id
28	921525775169458176	1255968042787524609	TW		
29	22703339	1255922986726707200	TW		
30	109263459	1255909313316601858	RE	1255905597976969216	109263459
31	16910746	1255860128231895044	TW		
32	11347122	1255949158751571968	TW		
33	705694814612938752	1255997250544812035	RE	1255993387603329026	705694814612938752
34	77372998	1255852155967332352	TW		
35	12272322	1255930519298224134	TW		
36	42407972	1255981948872413184	TW		
37	1154180084078501889	1255947666866843648	QT		
38	928677206	1255894000827011072	TW		

Figura 5. Esempio preso dal file creates.csv

La colonna type nel file creates l'ho popolata nel seguente modo:

- Assegno RE ad ogni riga in cui è presente un valore in reply_to_status_id o reply_to_user_id
- Assegno QT ad ogni riga in cui is_quote è true
- Assegno REQT ad ogni riga in cui è presente un valore in reply_to_status_id o reply_to_user_id e is_quote è true
- Assegno TW ai restanti tweet

Come si può notare il tipo RT non è presente, questo perché nel dataset è salvato solamente il numero dei retweet e non chi e quando li ha fatti. Per ovviare a questa mancanza ho provveduto a generare N tuple da inserire in creates tante quante la somma di tutti i retweet_count della tabella tweet. Per fare ciò ho creato un file creates.csv dove per convenienza ho aggiunto la colonna dei retweet_count. Poi grazie a un convertitore online, si veda [2], ho trasformato il file creates_rtcount.csv in creates_rtcount.json, quindi

	user_id	status_id	type	reply_to_status_id	reply_to_user_id	retweet_count
28	921525775169458176	1255968042787524609	TW			1947
29	22703339	1255922986726707200	TW			1691
30	109263459	1255909313316601858	RE	1255905597976969216	109263459	1617
31	16910746	1255860128231895044	TW			1616
32	11347122	1255949158751571968	TW			1598
33	705694814612938752	1255997250544812035	RE	1255993387603329026	705694814612938752	1591
34	77372998	1255852155967332352	TW			1563
35	12272322	1255930519298224134	TW			1443
36	42407972	1255981948872413184	TW			1429
37	1154180084078501889	1255947666866843648	QT			1428
38	928677206	1255894000827011072	TW			1428

Figura 6. Esempio preso dal file creates_rtcount.csv

ho proceduto a programmare in python questo "generatore" di retweet, la porzione di codice più importante è:

```

1 for row in reader:
2     if row['retweet_count'] > 0:
3         for tupla in range(row['retweet_count']):
4             stringa = {"user_id": user[posUtente],
5                        "status_id": row['status_id'],
6                        "type": "RT",
7                        "reply_to_status_id": None,
8                        "reply_to_user_id": None}
9             retweet.append(stringa)
10            if posUtente != len(user) - 1:
11                posUtente = posUtente + 1
12            else:
13                posUtente = 0

```

Ogni tupla che rappresenta un retweet di un tweet avrà un user_id casuale e lo status_id uguale a quello del tweet originale, le ho inserite in una lista che poi ho salvato in un file JSON. Sempre grazie a [2] ho proseguito a convertire gli altri file CSV in modo da prepararmi a popolare un eventuale database.

Database

Creazione

Grazie al lavoro fatto in precedenza, vedi Progettazione Logica, ho creato un Database locale sul mio PC tramite PostgreSQL 12 [4], successivamente ho creato le tre tabelle previste dalla Progettazione Logica con una semplice query:

```
1 CREATE TABLE utente(  
2     user_id VARCHAR PRIMARY KEY CHECK (user_id <> ''),  
3     screen_name VARCHAR NOT NULL CHECK (screen_name <> '')  
4 );  
5 CREATE TABLE tweet(  
6     status_id VARCHAR PRIMARY KEY CHECK (status_id <> ''),  
7     text VARCHAR NOT NULL CHECK (text <> ''),  
8     created_at TIMESTAMP NOT NULL CHECK (created_at <=  
9         CURRENT_TIMESTAMP),  
10    retweet_count INTEGER CHECK (retweet_count >= 0)  
11 );  
12 CREATE DOMAIN TYPETWEET VARCHAR(4) CHECK(  
13     VALUE IN ('TW', 'RE', 'RT', 'QT', 'REQT')  
14 );  
15 CREATE TABLE creates(  
16     user_id VARCHAR REFERENCES utente,  
17     status_id VARCHAR REFERENCES tweet,  
18     type TYPETWEET NOT NULL,  
19     reply_to_status_id VARCHAR DEFAULT NULL,  
20     reply_to_user_id VARCHAR DEFAULT NULL,  
21     quotes_to_status_id VARCHAR DEFAULT NULL REFERENCES  
22         tweet,  
23     PRIMARY KEY (user_id, status_id)  
24 );
```

Upload Dati

Durante alcune prove per popolare le tabelle del database ho riscontrato problemi con il testo dei tweet, questo perché erano presenti degli apici e delle emoticon non codificate. Per quanto riguarda gli apici, SQL li interpreta solo se sono doppi, mentre le emoticon dovevano essere tradotte nella loro codifica esadecimale. Di conseguenza ho creato in python (si veda [5] e [6]) un programma che va a sostituire un apice con un doppio apice e le emoticon con la loro codifica, in generale va a correggere un intero file JSON dato in input. Il codice è:

```
1 import json
2
3 parser = json.JSONDecoder()
4 parsed = []
5 with open("tweet.json", encoding='utf8') as f:
6     data = f.read()
7     data = data.replace("'", '"'+"'")
8
9 head = 0
10 for row in data:
11     head = (data.find('{', head) + 1 or data.find('[', head)
12             + 1) - 1
13     try:
14         struct, head = parser.raw_decode(data, head)
15         parsed.append(struct)
16     except (ValueError, json.JSONDecodeError):
17         break
18
19 with open("tweet_output.json", "w") as outfile:
20     json.dump(parsed, outfile, separators = (',', ': '),
21             indent = 1)
22
23 print("\nFile JSON formattato correttamente!")
```

Fatto ciò ho tutti gli ingredienti per andare a popolare il database, ho creato (si veda [7]) un uploader in python che va a creare una connessione e fare un INSERT per ogni oggetto del file JSON dato in input. Il codice è:

```
1 import json
2 import psycopg2
3 from myAppConfig import myHost, myDatabase, myUser, myPsw
4
5 connessione = psycopg2.connect (host = myHost, database =
    myDatabase, user = myUser, password = myPsw)
6
7 with connessione:
8     with connessione.cursor() as cursore:
9
10         with open("createsRT.json", encoding='utf8') as f:
11             data = json.load(f)
12
13             keys = []
14             for row in data:
15                 for key in row.keys():
16                     if key not in keys:
17                         keys.append(key)
18
19             for row in data:
20                 query = "INSERT INTO creates (user_id,
                    status_id, type, reply_to_status_id,
                    reply_to_user_id)
                    VALUES({0});".format(",".join(map(lambda key:
                    "'{0}'".format(row[key]) if key in row else
                    "NULL", keys)))
21                 cursore.execute(query)
22
23             print("Caricamento completato, tuple inserite!")
24
25 connessione.close()
26 print("Connessione chiusa con successo!")
```

Modificando il file in input con quello desiderato e la stringa query ho eseguito in modo corretto l'upload di tutti i file JSON che ho creato in precedenza.

Grafo

Struttura

Prima di poter sviluppare un programma che generi un determinato grafo c'è la necessità di capire come andarlo a strutturare. Dopo alcuni colloqui abbiamo scelto di implementare un grafo che avesse i seguenti nodi

- DB, il nodo radice del grafo
- Tweet, un nodo che rappresenta un'istanza di un tweet
- Retweet, un nodo che rappresenta un'istanza di un retweet

Ogni nodo può essere connesso ad un altro tramite un arco, in particolare nel grafo troveremo questi archi:

- $DB \rightarrow Tweet$
- $DB \rightarrow Retweet$
- $Tweet \rightarrow Retweet$
- $Tweet \rightarrow Tweet$ successivo (stesso creatore)
- $Retweet \rightarrow Tweet$ successivo (se presente)
- $Tweet$ precedente $\rightarrow Retweet$ (se presente)

Ogni arco è descritto da un'etichetta del tipo $\langle Start, End, Author, Weight \rangle$ dove $Weight$ rappresenta un valore che descrive quanto un autore ha collaborato nella creazione del tweet o retweet, in particolare può assumere questi valori:

- 1 se l'arco è entrante in un Tweet
- 0.1 se l'arco è entrante in un Retweet
- 0.5 se l'arco è entrante in un Tweet di tipo quotes

Start ed End rappresentano un intervallo temporale in minuti per descrivere la creazione del nodo a cui punta l'arco.

Author rappresenta il proprietario del Tweet o Retweet.

I nodi con degli intervalli di tempo indeterminati vengono descritti con un'etichetta qualitativa con il valore 0 in Start e il valore ∞ in End. Per distinguere meglio i nodi e archi ho scelto di assegnare alcuni colori:

- Per il nodo radice DB
- Per il nodo che rappresenta un tweet di tipo TW, il suo arco collegato a DB e l'arco collegato al tweet successivo
- Per il nodo che rappresenta un tweet di tipo RE, il suo arco collegato a DB e l'arco collegato al tweet successivo
- Per il nodo che rappresenta un tweet di tipo QT, il suo arco collegato a DB e l'arco collegato al tweet successivo
- Per il nodo che rappresenta un tweet di tipo REQT, il suo arco collegato a DB e l'arco collegato al tweet successivo
- Per il nodo che rappresenta un retweet e il suo arco collegato a un nodo tweet

Creazione

Per la creazione del grafo ho scelto di usare nuovamente Python, in particolare la libreria NetworkX [8] mi è venuta in aiuto grazie alle possibilità di creazione e personalizzazione di nodi, archi ed etichette. Ho scelto di inizializzare un grafo di tipo DiGraph, questo perché dà la possibilità di dare una direzione agli archi e fornisce alcuni metodi aggiuntivi che ho usato nel mio programma. Ad ogni avvio la prima cosa che fa il programma è il creare la connessione con il database e successivamente inizializzare il cursore che mi permetterà di eseguire le query necessarie a prelevare determinate informazioni. La prima query che andrà a fare è la seguente:

```
1 DROP VIEW IF EXISTS user_id_quotes;
2 CREATE VIEW user_id_quotes AS (
3     SELECT U.screen_name, U.user_id, T.status_id
4     FROM utente U
5         JOIN creates C ON U.user_id = C.user_id
6         JOIN tweet T ON C.status_id = T.status_id
7     WHERE C.type = 'TW' AND
8           T.status_id IN (SELECT Ci.quotes_to_status_id
9                           FROM creates Ci
10                          WHERE Ci.quotes_to_status_id
11                             IS NOT NULL)
12 );
13 SELECT T.status_id, T.created_at, U.screen_name, U.user_id,
14        C.type,
15        C.quotes_to_status_id, V.screen_name
16 FROM tweet T
17     JOIN creates C ON T.status_id = C.status_id
18     JOIN utente U ON C.user_id = U.user_id
19     LEFT JOIN user_id_quotes V ON C.quotes_to_status_id =
20        V.status_id
21 WHERE C.type <> 'RT' AND
22        T.status_id IN ('1255877065758388224',
23                       '1255779231776206849',
24                       '1255751113053265920',
25                       '1255930357691453442',
26                       '1255744610523021312')
27 ORDER BY U.user_id, T.created_at;
```

Successivamente avrò come risultato una tabella contenente lo status_id, created_at, screen_name, user_id e type di ogni tweet specificato nella clausola WHERE, inoltre se il tweet è di tipo QT avrò il quotes_to_status_id e lo screen_name di chi ha creato il tweet a cui fa riferimento il quotes, negli altri tipi di tweet queste due colonne sono popolate da NULL. La fase successiva è la creazione della prima parte di grafo, il codice è il seguente:

```

1 tweetIdOriginale = cursore.fetchall()
2 status_id_tweet = [tupla[0] for tupla in tweetIdOriginale]
3 created_at_tweet = [tupla[1] for tupla in tweetIdOriginale]
4 screen_name_tweet = [tupla[2] for tupla in tweetIdOriginale]
5 user_id_tweet = [tupla[3] for tupla in tweetIdOriginale]
6 type_tweet = [tupla[4] for tupla in tweetIdOriginale]
7 original_tw_quoted = [tupla[5] for tupla in tweetIdOriginale]
8 author_tw_quoted = [tupla[6] for tupla in tweetIdOriginale]
9
10 startDB = datetime.strptime('2020-04-30 00:00:00', '%Y-%m-%d
    %H:%M:%S')
11
12 for x in range(len(tweetIdOriginale)):
13     node_created_at =
14         datetime.strptime(str(created_at_tweet[x]),
15                             '%Y-%m-%d %H:%M:%S')
16     interval = node_created_at - startDB
17     minuteInterval = int(interval.total_seconds()/60)
18     attr_dict = {}
19     time_list = []
20     if type_tweet[x] != 'QT':
21         attr_dict["Start"] = minuteInterval
22         attr_dict["End"] = minuteInterval
23         attr_dict["Authors"] = [(screen_name_tweet[x],
24                                 1)]
25         stringa = tuple(x for x in attr_dict.values())
26         time_list.append(stringa)
27         graph.add_edge('DB', status_id_tweet[x],
28                         attr_dict = attr_dict, time = stringa,
29                         time_list = time_list)
30     elif type_tweet[x] == 'QT':
31         attr_dict["Start"] = minuteInterval
32         attr_dict["End"] = minuteInterval
33         attr_dict["Authors"] = [(screen_name_tweet[x],
34                                 0.5)]
35         stringa = tuple(x for x in attr_dict.values())
36         time_list.append(stringa)

```

```

31     graph.add_edge('DB', status_id_tweet[x],
32                     attr_dict = attr_dict, time = stringa,
33                     time_list = time_list)
34     attr_dict["Start"] = 0
35     attr_dict["End"] = math.inf
36     attr_dict["Authors"] = [(screen_name_tweet[x],
37                               0.5)]
38     stringa = tuple(x for x in attr_dict.values())
39     time_list = []
40     time_list.append(stringa)
41     graph.add_edge(original_tw_quoted[x],
42                     status_id_tweet[x], attr_dict = attr_dict,
43                     time = stringa, time_list = time_list)
44 if x < len(tweetIdOriginale) - 1:
45     if user_id_tweet[x] == user_id_tweet[x+1]:
46         attr_dict = {}
47         attr_dict["Start"] = 0
48         attr_dict["End"] = math.inf
49         attr_dict["Authors"] =
50             [(screen_name_tweet[x], 1)]
51         stringa = tuple(x for x in
52                         attr_dict.values())
53         time_list = []
54         time_list.append(stringa)
55         graph.add_edge(status_id_tweet[x],
56                         status_id_tweet[x+1], attr_dict =
57                         attr_dict, time = stringa, time_list =
58                         time_list)

```

Questa porzione di codice salva ogni colonna della tabella ottenuta tramite la query in una lista, successivamente procede a scorrere ogni elemento, controlla se è un quote o un altro tipo di tweet e di conseguenza crea l'etichetta adatta in base alle scelte fatte in precedenza, vedi Struttura, lo collega al nodo radice DB e se sono creati dallo stesso utente connette i tweet tra di loro.

Dopo questa fase il programma procede con la seconda query:

```
1 DROP VIEW IF EXISTS User_Id_Originale;
2 CREATE VIEW User_Id_Originale AS (
3     SELECT U.user_id, T.status_id, U.screen_name
4     FROM tweet T
5         JOIN creates C ON T.status_id = C.status_id
6         JOIN utente U ON C.user_id = U.user_id
7     WHERE C.type <> 'RT' AND
8         T.status_id IN ('1255877065758388224',
9                        '1255779231776206849',
10                       '1255751113053265920',
11                       '1255930357691453442',
12                       '1255744610523021312')
13     ORDER BY U.user_id
14 );
15 SELECT U.user_id, T.status_id, T.created_at, U.screen_name,
16        V.user_id, V.screen_name
17 FROM utente U
18     JOIN creates C ON U.user_id = C.user_id
19     JOIN tweet T ON C.status_id = T.status_id
20     JOIN User_Id_Originale V ON T.status_id = V.status_id
21 WHERE C.type = 'RT' AND
22        T.status_id IN ('1255877065758388224',
23                       '1255779231776206849',
24                       '1255751113053265920',
25                       '1255930357691453442',
26                       '1255744610523021312')
27 ORDER BY V.user_id, T.created_at;
```

La quale restituirà una tabella contenente user_id, status_id, created_at, screen_name di chi ha creato il retweet e user_id, screen_name del creatore del tweet che è stato retweetato, il tutto in base agli status_id inseriti nella clausola WHERE.

Per creare la seconda parte del grafo ho scritto il seguente codice:

```
1 utenti_tweet_RT = cursore.fetchall()
2 user_id = [tupla[0] for tupla in utenti_tweet_RT]
3 status_id = [tupla[1] for tupla in utenti_tweet_RT]
4 status_id_no_duplicates =
    list(OrderedDict.fromkeys(status_id))
5 created_at = [tupla[2] for tupla in utenti_tweet_RT]
6 screen_name = [tupla[3] for tupla in utenti_tweet_RT]
7 creator_user_id = [tupla[4] for tupla in utenti_tweet_RT]
8 creator_screen_name = [tupla[5] for tupla in utenti_tweet_RT]
9
10 user_tweet_duplicates = defaultdict(list)
11 for i, j in zip(creator_screen_name, status_id):
12     user_tweet_duplicates[i].append(j)
13 user_tweet = {a:list(OrderedSet(b)) for a, b in
    user_tweet_duplicates.items()}
14
15 for i in range(len(utenti_tweet_RT)):
16     rt_id = status_id[i]+'-'+str(i)
17     retweet_id.append(rt_id)
18     attr_dict = {}
19     attr_dict["Start"] = 0
20     attr_dict["End"] = math.inf
21     attr_dict["Authors"] = [(screen_name[i], 0.1)]
22     stringa = tuple(x for x in attr_dict.values())
23     time_list = []
24     time_list.append(stringa)
25     graph.add_edge(status_id[i], retweet_id[i],
        attr_dict = attr_dict, time = stringa, time_list
        = time_list)
26     graph.add_edge('DB', retweet_id[i], attr_dict =
        attr_dict, time = stringa, time_list = time_list)
27
28 posRT = 0
29 rt_id_precedenti = []
30 for user, tweet in user_tweet.items():
31     posTW = 0
32     for tw in tweet:
33         nodes = []
34         nodes = nx.dfs_preorder_nodes(graph, tw,
            depth_limit = 1)
35         for n in nodes:
36             if n == 'DB':
37                 continue
38             elif n in status_id_no_duplicates:
39                 continue
```

```

40         elif n in rt_id_precedenti:
41             continue
42         else:
43             for x in range(0, posRT):
44                 if retweet_id[x] not in
45                     rt_id_precedenti:
46                         rt_id_precedenti.append(retweet_id[x])
47             posRT = posRT + 1
48             if len(tweet) > 1 and posTW < len(tweet)
49                 - 1:
50                 attr_dict = {}
51                 attr_dict["Start"] = 0
52                 attr_dict["End"] = math.inf
53                 attr_dict["Authors"] =
54                     [(creator_screen_name[posTW+1],
55                      1)]
56                 stringa = tuple(x for x in
57                     attr_dict.values())
58                 time_list = []
59                 time_list.append(stringa)
60                 graph.add_edge(n, tweet[posTW+1],
61                     attr_dict = attr_dict, time =
62                     stringa, time_list = time_list)
63             if len(tweet) > 1 and posTW > 0:
64                 attr_dict = {}
65                 attr_dict["Start"] = 0
66                 attr_dict["End"] = math.inf
67                 attr_dict["Authors"] =
68                     [(screen_name[posTW], 0.1)]
69                 stringa = tuple(x for x in
70                     attr_dict.values())
71                 time_list = []
72                 time_list.append(stringa)
73                 graph.add_edge(tweet[posTW-1], n,
74                     attr_dict = attr_dict, time =
75                     stringa, time_list = time_list)
76             posTW = posTW + 1

```

Come in precedenza salvo ogni colonna della tabella ottenuta in liste, in aggiunta ho creato la lista `status_id_no_duplicates` che corrisponde alla lista `status_id` ma senza duplicati, in pratica conterrà tutti gli id univoci dei tweet originali che sono stati retweetati. Il programma procede con la creazione di un dizionario che avrà come chiavi tutti gli `user_id` dei creatori dei tweet originali e come valori ogni `status_id`, senza duplicati (si veda [9], [10], [11]), che corrisponde ad un tweet creato dall'utente. Vengono poi

creati gli archi che dal tweet originale si collegano al nodo di ogni retweet e poi da ogni retweet un arco verso DB, il tutto con le etichette decise nella Struttura. Poi il programma inizia a scorrere ogni utente ed analizzare ogni tweet creato dallo stesso, in particolare va ad estrarre i nodi successivi al tweet fino alla profondità massima di 1 arco (vedi [12]), salta i nodi che non servono per creare nuovi archi ed infine se l'utente ha creato più di 1 tweet verranno connessi i retweet precedenti/successivi secondo quanto stabilito nella Struttura.

Dopo queste fasi la struttura del grafo è completa, la fase seguente è la colorazione dei nodi e archi in base alla loro tipologia e la loro stampa sotto forma di immagine, per fare ciò ho creato delle liste per poter capire di quale tipologia fosse il nodo e arco, ho usato queste query:

```
1 SELECT DISTINCT C.user_id
2 FROM creates C
3 WHERE C.type = 'RT';
4
5 SELECT C.status_id
6 FROM creates C
7 WHERE C.type = 'TW';
8
9 SELECT C.status_id
10 FROM creates C
11 WHERE C.type = 'QT';
12
13 SELECT C.status_id
14 FROM creates C
15 WHERE C.type = 'RE';
16
17 SELECT C.status_id
18 FROM creates C
19 WHERE C.type = 'REQT';
```

Le loro liste verranno usate in questa funzione creata appositamente:

```
1 def mapping_colors_nodes_edges(G):
2     edge_colors.clear()
3     node_colors.clear()
4     for node in G:
5         if node in check_tweetTW:
6             node_colors.append('#269AED')
7         elif node in check_tweetQT:
8             node_colors.append('#EB8D00')
9         elif node in check_tweetRE:
10            node_colors.append('#FFE100')
11        elif node in check_tweetREQT:
12            node_colors.append('#AA00FF')
13        elif node == 'DB':
14            node_colors.append('red')
15        else:
16            node_colors.append('#0DA342')
17
18    for edge in G.edges:
19        if any(ele in edge for ele in retweet_id) == True:
20            edge_colors.append('#0DA342')
21        elif any(ele in edge for ele in check_tweetQT) ==
22            True:
23            edge_colors.append('#EB8D00')
24        elif any(ele in edge for ele in check_tweetRE) ==
25            True:
26            edge_colors.append('#FFE100')
27        elif any(ele in edge for ele in check_tweetREQT) ==
28            True:
29            edge_colors.append('#AA00FF')
30        elif any(ele in edge for ele in check_tweetTW) ==
31            True:
32            edge_colors.append('#269AED')
```

In dettaglio questa funzione va a salvare in due liste, una per i nodi e una per gli archi, il codice del colore corrispondente alla posizione del nodo/arco.

Questa verrà usata in un'altra funzione per la stampa del grafo:

```

1 def print_graph (sub_G):
2     posSub_G = nx.shell_layout(sub_G)
3     mapping_colors_nodes_edges(sub_G)
4     nx.draw_networkx(sub_G, pos = posSub_G, with_labels =
5         True, node_color = node_colors, edge_color =
6         edge_colors, font_size = 9)
7     nx.draw_networkx_edge_labels(sub_G, pos = posSub_G,
8         edge_labels = nx.get_edge_attributes(sub_G, 'time'),
9         font_color = 'red', font_size = 7)
10    plt.plot(1)
11    plt.show(block = False)

```

La funzione di libreria `nx.shell_layout(sub_G)`, si veda [13], mappa la posizione dei nodi in cerchi concentrici. In generale verrà creata l'immagine completa del grafo e la si potrà visionare grazie alla libreria `matplotlib.pyplot` (si veda [14]).

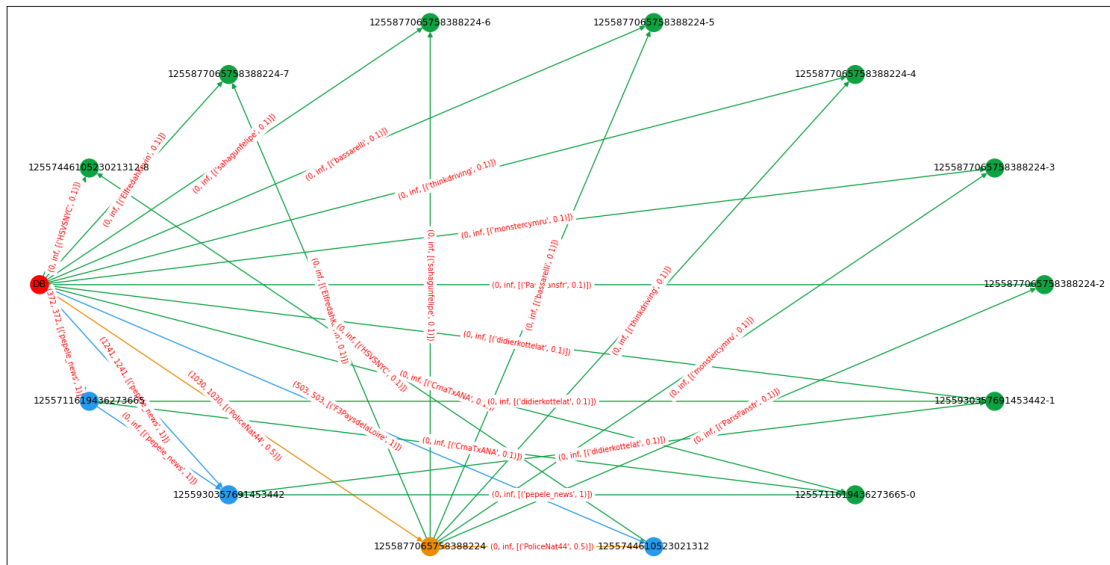


Figura 7. Grafo di esempio

Salvataggio

Il salvataggio del grafo può avvenire in due modi, tramite il plot del grafo avremo la possibilità di salvare un file PNG del grafo, oppure tramite una apposita funzione che ho creato si potrà salvare un file che conterrà la lista completa dei nodi e archi del grafo, la funzione è:

```
1 def save_file_edgelist(G):
2     if not nx.is_empty(G):
3         nome = input("Inserire un nome per il salvataggio
4             del grafo: ")
5         nome = nome + '.edgelist'
6         f = open(nome, 'wb')
7         nx.write_edgelist(G, f)
8         print("Sotto-grafo salvato correttamente.")
9         f.close()
10    else:
11        print("Non c'e' nessun sotto-grafo da salvare.")
```

La funzione `is_empty(G)` (si veda [16]) controlla che il grafo non sia vuoto, quindi si procede con la creazione del file tramite `write_edgelist(G, f)` che riceve in input il grafo da salvare e il file (si veda [15]).

Caricamento

Un grafo salvato può essere caricato dal programma tramite la seguente funzione:

```
1 def read_file_edgelist():
2     nome = input("Inserire il nome del grafo da caricare: ")
3     nome = nome + '.edgelist'
4     try:
5         f = open(nome, 'rb')
6         G = nx.read_edgelist(f, create_using = nx.DiGraph())
7         f.close()
8         print("Sotto-grafo letto correttamente.")
9         return G
10    except:
11        print("Il nome inserito non esiste!")
```

La funzione `read_edgelist` (si veda [17]) riceve in input il file.edgelist e l'indicazione della tipologia di grafo, creerà e copierà il grafo creato nella variabile `G` la quale verrà inserita nel return della funzione principale.

Algoritmo Path-Consistency

Prima di dare la definizione dell'algoritmo è necessario dare la definizione di altri tre concetti che ci serviranno per capirlo meglio.

Definizione 1 (Inversione) *Data una tupla $T_k = \langle Start, End \rangle$ il suo inverso è $T_k^{-1} = \langle -End, -Start \rangle$*

Definizione 2 (Composizione \circ) *Date due tuple $T_{k1} = \langle Start_1, End_1 \rangle$ e $T_{k2} = \langle Start_2, End_2 \rangle$ la composizione è $T_{k12} = \langle Start_1 + Start_2, End_1 + End_2 \rangle$, in pratica la somma di due tuple.*

Definizione 3 (Congiunzione \otimes) *Date due tuple $T_{k1} = \langle Start_1, End_1 \rangle$ e $T_{k2} = \langle Start_2, End_2 \rangle$ la congiunzione è $T_{k12} = \langle MAX(Start_1, Start_2), MIN(End_1, End_2) \rangle$, in pratica crea una tupla con il valore Start massimo e il valore End minimo tra le due tuple.*

Ora possiamo dare la definizione dell'algoritmo

Definizione 4 (Path-Consistency) *Date tre tuple $R_{ij} = \langle Start, End \rangle$, $R_{ik} = \langle Start, End \rangle$ e $R_{kj} = \langle Start, End \rangle$ l'algoritmo esegue quanto segue $R'_{ij} = R_{ij} \otimes (R_{ik} \circ R_{kj})$*

Può accadere che in alcuni casi di studio serva applicare l'inversione a una tupla per ottenere un risultato coerente con il grafo.

Definizione 5 (Intersezione Intervalli \cap) *Dati due intervalli $T_1 = \langle Start_1, End_1 \rangle$ e $T_2 = \langle Start_2, End_2 \rangle$ l'intersezione $T_{12} = T_1 \cap T_2$ è data dal valore minimo di Start $\min(Start_1, Start_2)$ e quello massimo di End $\max(End_1, End_2)$, nel caso End sia maggiore o uguale di Start il risultato è la sottrazione $End - Start$ altrimenti 0.*

Definizione 6 (Unione Intervalli \cup) *Dati due intervalli $T_1 = \langle Start_1, End_1 \rangle$ e $T_2 = \langle Start_2, End_2 \rangle$ l'unione $T_{12} = T_1 \cup T_2$ è data dal valore massimo di Start $\max(Start_1, Start_2)$ e quello minimo di End $\min(End_1, End_2)$, nel caso End sia maggiore o uguale di Start il risultato è la sottrazione $End - Start$ altrimenti 0.*

Definizione 7 (Similarità) *Dati due intervalli $T_1 = \langle Start_1, End_1 \rangle$ e $T_2 = \langle Start_2, End_2 \rangle$ la similarità è data da $\text{sim}(T_1, T_2) = \frac{T_1 \cap T_2}{T_1 \cup T_2}$, nel caso l'unione risulti infinita il risultato è 0.1 (poca similarità), nel caso sia diversa da 0 il risultato è dato dal rapporto tra l'intersezione e l'unione, altrimenti risulta 0.*

Al grafo viene applicato l'algoritmo con il seguente codice:

```
1 def path_consistency(G):
2     changed = True
3     count = 0
4     while(changed == True):
5         count = count + 1
6
7         G_edgelist = open('G_file.txt', 'wb')
8         nx.write_edgelist(G, G_edgelist)
9         G_edgelist.close()
10
11        G_new = copy.deepcopy(path_consistency_step(G))
12
13        G_new_edgelist = open('G_new_file.txt', 'wb')
14        nx.write_edgelist(G_new, G_new_edgelist)
15        G_new_edgelist.close()
16
17        G_edgelist = open('G_file.txt')
18        G_new_edgelist = open('G_new_file.txt')
19
20        if filecmp.cmp('G_file.txt', 'G_new_file.txt',
21            shallow=False) == True:
22            G_edgelist.close()
23            G_new_edgelist.close()
24            print("La "+str(count)+"      esecuzione non ha
25                apportato modifiche.")
26            changed = False
27        else:
28            G_edgelist.close()
29            G_new_edgelist.close()
30            G = copy.deepcopy(G_new)
31            print(str(count) + "      esecuzione
                dell'algoritmo.")
32
33    return G_new
```

Dove una porzione significativa della funzione `path_consistency_step` è:

```
1 def path_consistency_step(G):
2     G_new = nx.DiGraph()
3     G_new.add_nodes_from(copy.deepcopy(G.nodes(data =
4         False)))
5
6     for edge in G.edges():
7         triangles = []
8         triangles = get_triangles(G, edge[0], edge[1])
9         for t in triangles:
10             if (G.has_edge(t[0], t[1]) and G.has_edge(t[1],
11                 t[2])):
12
13                 if G_new.has_edge(t[0], t[2]) == True:
14                     Rij_label = {}
15                     Rij_label =
16                         copy.deepcopy(G_new.get_edge_data(t[0],
17                             t[2]))
18                     Rij = (Rij_label['attr_dict']['Start'],
19                         Rij_label['attr_dict']['End'])
20                 else:
21                     Rij_label = {}
22                     Rij_label =
23                         copy.deepcopy(G.get_edge_data(t[0],
24                             t[2]))
25                     Rij = (Rij_label['attr_dict']['Start'],
26                         Rij_label['attr_dict']['End'])
27
28                 Rik_label = {}
29                 Rik_label =
30                     copy.deepcopy(G.get_edge_data(t[0], t[1]))
31                 Rik = (Rik_label['attr_dict']['Start'],
32                     Rik_label['attr_dict']['End'])
33
34                 Rkj_label = {}
35                 Rkj_label =
36                     copy.deepcopy(G.get_edge_data(t[1], t[2]))
37                 Rkj = (Rkj_label['attr_dict']['Start'],
38                     Rkj_label['attr_dict']['End'])
39
40                 autori = []
41                 for autore1 in
42                     Rik_label['attr_dict']['Authors']:
43                     for autore2 in
44                         Rkj_label['attr_dict']['Authors']:
45                         if autore1[0] == autore2[0]:
```

```

32         autori.append((autore1[0],
33                         max(max(autore1[1], autore2[1]))))
34     else:
35         autori.append(max(autore1, autore2))
36
37     Rikj = (Rik[0]+Rkj[0], Rik[1]+Rkj[1])
38
39     if Rij[0] >= Rikj[0]:
40         if Rij[1] >= Rikj[1]:
41             Rij_new = [Rij[0], Rikj[1],
42                       Rij_label['attr_dict']['Authors']]
43         else:
44             Rij_new = [Rij[0], Rij[1],
45                       Rij_label['attr_dict']['Authors']]
46     elif Rij[0] < Rikj[0]:
47         if Rij[1] >= Rikj[1]:
48             Rij_new = [Rikj[0], Rikj[1],
49                       Rij_label['attr_dict']['Authors']]
50         else:
51             Rij_new = [Rikj[0], Rij[1],
52                       Rij_label['attr_dict']['Authors']]
53
54     for index, autore in enumerate(autori):
55         y = list(autore)
56         y[1] = similarity(Rikj[0], Rikj[1],
57                           Rij_new[0], Rij_new[1])
58         autori[index] = y
59
60     for index, autore in enumerate(Rij_new[2]):
61         y = list(autore)
62         y[1] =
63             similarity(Rij_label['attr_dict']['Start'],
64                       Rij_label['attr_dict']['End'],
65                       Rij_new[0], Rij_new[1])
66         Rij_new[2][index] = y
67
68     Rij_new[2] = Rij_new[2]+autori
69
70     for index1, autore1 in enumerate(Rij_new[2]):
71         for index2, autore2 in
72             enumerate(Rij_new[2]):
73             if autore1 == autore2:
74                 continue
75             elif autore1[0] == autore2[0]:

```

```

67         Rij_new[2][index1] =
            [autore1[0], max(autore1[1],
            autore2[1])]
68         Rij_new[2][index2] =
            [autore2[0], max(autore1[1],
            autore2[1])]
69
70
71         Rij_new[2] = [t for t in (set(tuple(i) for i
            in Rij_new[2]))]
72
73         Rij_label['attr_dict']['Start'] = Rij_new[0]
74         Rij_label['attr_dict']['End'] = Rij_new[1]
75         Rij_label['attr_dict']['Authors'] =
            Rij_new[2]
76
77         new_time = tuple(x for x in
            Rij_label['attr_dict'].values())
78
79         G_new.add_edge(t[0], t[2], attr_dict =
            Rij_label['attr_dict'], time = new_time)
80
81     return G_new

```

La funzione triangles è:

```

1 def get_triangles(G, i, j):
2     result = []
3     for nodo in G.nodes():
4         if (G.has_edge(i, nodo) and G.has_edge(nodo, j)) or \
5             (G.has_edge(i, nodo) and G.has_edge(j, nodo)) or \
6             (G.has_edge(nodo, i) and G.has_edge(nodo, j)) or \
7             (G.has_edge(nodo, i) and G.has_edge(j, nodo)):
8
9             triangle = (i, nodo, j)
10            result.append(triangle)
11
12    return result

```

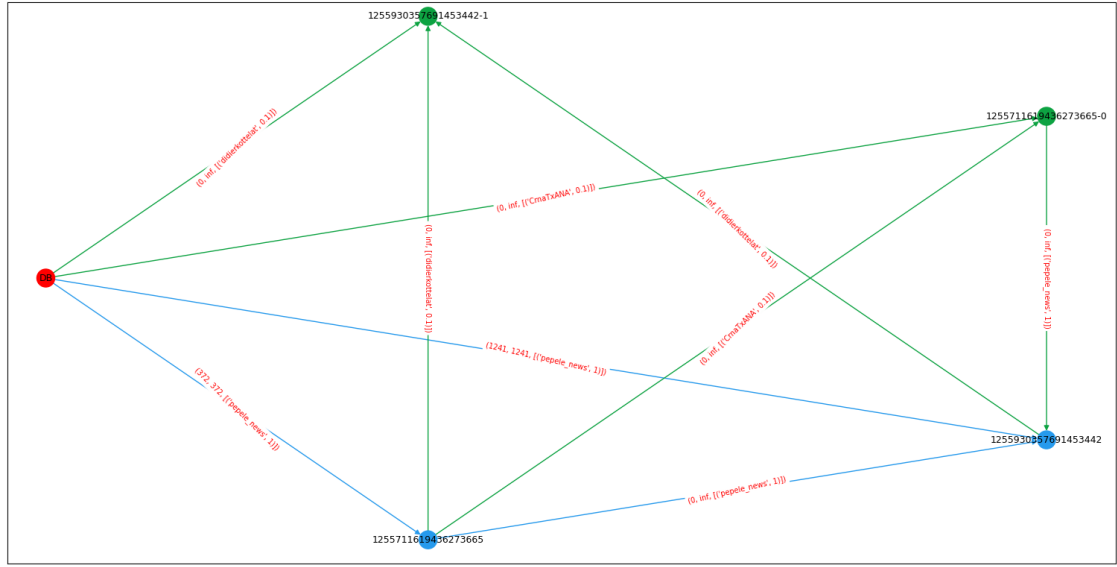


Figura 8. Grafo di esempio senza algoritmi applicati

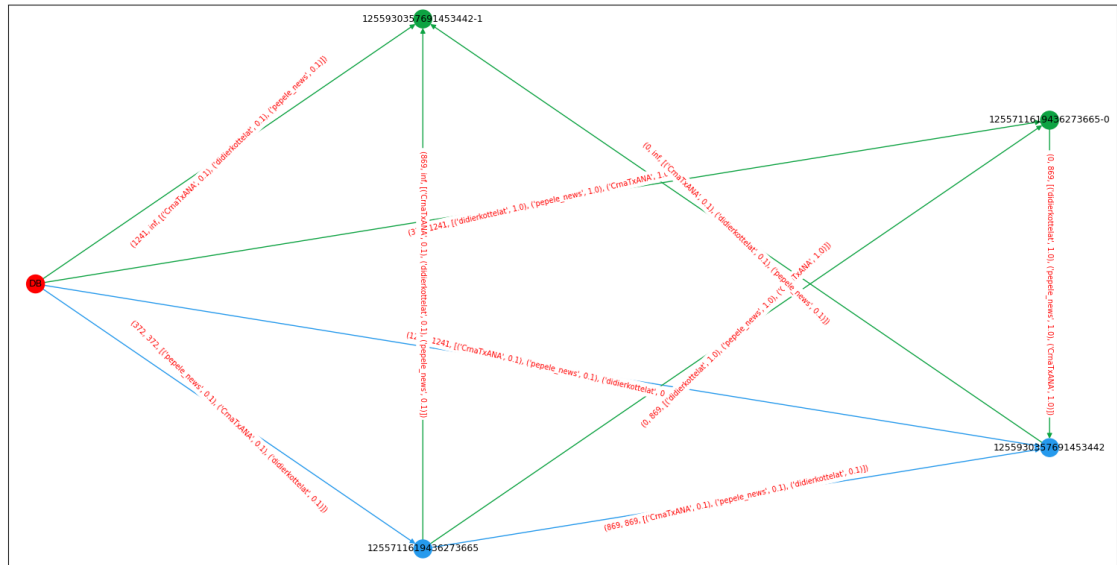


Figura 9. Grafo di esempio con Path-Consistency e Similarità degli intervalli

Gestione

Per fornire una gestione minimale delle funzioni descritte in precedenza ho creato un semplice menù, il codice è il seguente:

```
1 def menu():
2     print('*' * 55)
3     comando = input("Scegliere uno dei seguenti comandi:\n"\
4                       "1) Applica l'algoritmo Path Consistency
5                       al grafo.\n"\
6                       "2) Salvare il grafo creato in un
7                       file.\n"\
8                       "3) Caricare un grafo da un file.\n"\
9                       "4) Visualizzare il grafo
10                      creato/modificato o caricato dal
11                      file.\n"\
12                      "0) Pulire la console.\n"\
13                      "q) exit()\n\n")
14
15     if comando == 'q':
16         print('*' * 55)
17         return comando
18     else:
19         print('*' * 55)
20         return int(comando)
21
22 def loop(G):
23     try:
24         x = menu()
25         while x in {0, 1, 2, 3, 4, 'q'}:
26             if x == 'q':
27                 print("Sto chiudendo il programma...")
28                 break
29             if x == 0:
30                 os.system("cls")
31                 x = menu()
32             if x == 1:
33                 G = path_consistency(G)
34                 print("Algoritmo applicato.")
35                 x = menu()
36             if x == 2:
37                 save_file_edgelist(G)
38                 x = menu()
39             if x == 3:
40                 G = read_file_edgelist()
41                 x = menu()
42             if x == 4:
```

```

39         print_graph(G)
40         x = menu()
41         if x not in {0, 1, 2, 3, 4, 'q'}:
42             raise ValueError
43     except ValueError:
44         print("Hai inserito una scelta del men  non valida,
45             ritenta!\n")
         loop(G)

```

Questo permette di far eseguire un'operazione alla volta tramite l'inserimento di un numero da 0 a 4, nel caso venga inserito il carattere 'q' il programma terminerà.

```

*****
Scegliere uno dei seguenti comandi:
1) Applica l'algoritmo Path Consistency al grafo.
2) Salvare il grafo creato in un file.
3) Caricare un grafo da un file.
4) Visualizzare il grafo creato/modificato o caricato dal file.
0) Pulire la console.
q) exit()

```

Figura 10. Screenshot del menù

Sitografia

- [1] <https://www.kaggle.com/smidth80/coronavirus-covid19-tweets-late-april/data?select=2020-04-30+Coronavirus+Tweets.CSV>
- [2] <http://www.convertcsv.com/csv-to-json.htm>
- [3] <https://app.diagrams.net/>
- [4] <https://www.postgresql.org/download/>
- [5] <https://stackoverflow.com/questions/55555095/how-to-add-commas-in-between-json-objects-present-in-a-txt-file-and-the-n-conver>
- [6] <https://www.geeksforgeeks.org/reading-and-writing-json-to-a-file-in-python/>
- [7] <https://stackoverflow.com/questions/48604563/creating-a-data-structure-from-json-using-python>
- [8] <https://networkx.github.io/documentation/stable/index.html#>
- [9] <https://stackoverflow.com/questions/49268016/creating-a-dictionary-from-2-lists-with-duplicate-keys>
- [10] <https://stackoverflow.com/questions/44395560/how-to-remove-duplicate-values-from-dict>
- [11] <https://stackoverflow.com/questions/1653970/does-python-have-an-ordered-set>
- [12] https://networkx.github.io/documentation/stable/reference/algorithms/generated/networkx.algorithms.traversal.depth_first_search.dfs_preorder_nodes.html#networkx.algorithms.traversal.depth_first_search.dfs_preorder_nodes
- [13] https://networkx.github.io/documentation/stable/reference/generated/networkx.drawing.layout.shell_layout.html?highlight=shell_layout#networkx.drawing.layout.shell_layout
- [14] <https://matplotlib.org/index.html>
- [15] https://networkx.github.io/documentation/stable/reference/readwrite/generated/networkx.readwrite.edgelist.write_edgel

- ist.html?highlight=write_edgelist#networkx.readwrite.edgelist.write_edgelist
- [16] https://networkx.github.io/documentation/stable/reference/generated/networkx.classes.function.is_empty.html?highlight=is_empty#networkx.classes.function.is_empty
- [17] https://networkx.github.io/documentation/stable/reference/readwrite/generated/networkx.readwrite.edgelist.read_edgelist.html?highlight=read_edgelist#networkx.readwrite.edgelist.read_edgelist
- [18] https://networkx.github.io/documentation/stable/reference/algorithms/generated/networkx.algorithms.traversal.breadth_first_search.bfs_edges.html#networkx.algorithms.traversal.breadth_first_search.bfs_edges
- [19] https://networkx.github.io/documentation/stable/reference/classes/generated/networkx.DiGraph.get_edge_data.html?highlight=get_edge_data#networkx.DiGraph.get_edge_data