
Sound of Data

Riccardo Cervero 794126
Marco Ferrario 795203
Pranav Kasela 846965
Federico Moiraghi 799735

Universit degli Studi di Milano Bicocca

Anno Accademico 2018/19

Obiiettivo del progetto è analizzare la discussione mediatica sul social network Twitter studiando in tempo reale l'andamento della discussione nell'ambito musicale, presente e passato, per verificare la presenza di variazioni all'interno del ciclo giornaliero e la presenza di comunità digitali e le loro interconnessioni. Il presente lavoro può essere di aiuto nella costruzione di sistemi di raccomandazione interni a Twitter (quali pagine consigliate) e nella programmazione di BOT, a scopo pubblicitario, che intendano infilarsi in una discussione mediatica già presente (ovvero argomenti di tendenza) per inserire i propri messaggi.

Indice

I	Introduzione	2
II	Costruzione dello <i>knowledge graph</i>	3
III	Analisi dei tweet	5
IV	Visualizzazioni dei dati	9
V	Risultati e conclusioni	11

Parte I

Introduzione

I traguardi posti dal progetto *Sound of Data*¹, hanno innanzitutto richiesto l'uso di una *knowledge graph* adeguata alla materia: scaricato un *dump* di musicbrainz.org, versione *semantic web* di last.fm, disponibile pubblicamente in formato *.tsv* (*Tabular Separated Values*), si è provveduto ad importarlo all'interno di Apache Hadoop² per effettuare una rapida pulizia preliminare. Fatto ciò, i dati sono stati importati in Neo4J³, grazie al pratico programma *neo4j-import*, in modo da costruire la *knowledge graph* richiesto per l'analisi dei *tweet*, raccolti ed analizzati in tempo reale grazie ad Apache Kafka⁴: grazie ad un efficiente plug-in che permette l'uso di Neo4j come *Consumer*, il singolo *tweet* è archiviato nativamente in uno schema a grafo. Durante il processo di *streaming*, questi subiscono varie operazioni di *preparation* e vengono filtrati da un rudimentale strumento di *instance matching*.

Dopo la complessa fase di gestione e raccolta dei dati, è stato possibile condurre un'analisi finale sulle dinamiche della discussione musicale instauratesi fra i singoli utenti all'interno di comunità, astratte ed estremamente mutevoli, e sulle relazioni, più o meno intense, esistenti fra di esse. La progressiva composizione che i nodi associati ai *tweet* assumono all'interno del grafo è studiata ad intervalli più o meno regolari, corrispondenti a quattro fasce orarie relative al mattino, pomeriggio, sera, notte. Si è preferito adoperare tale ciclo giornaliero, anziché un più ampio ciclo settimanale, o addirittura mensile, per due ragioni precise. Innanzitutto, questo garantisce una più facile interpretazione: risulta di estremo interesse osservare i trend ricorrenti nei vari momenti della giornata, piuttosto che riferiti all'arco settimanale. Inoltre, il ciclo giornaliero

ha il vantaggio di non risentire della distorsione provocata dagli eventi, previsti o meno. Infatti, essendo gli eventi in grado di assorbire l'interesse della gran parte dell'opinione pubblica per interi giorni e settimane, questi penalizzano l'analisi e concentrano tutti gli utenti verso un numero limitato di argomenti musicali. Due esempi riguardanti tale problematica, e le sue conseguenze negative, possono essere citati per quanto concerne il genere *drone*, per cui durante il periodo di raccolta, si sono succeduti due eventi particolarmente rilevanti: il primo riguarda l'abbattimento di un drone americano in territorio Iraniano, causando nuove tensioni fra i due Paesi⁵, mentre il secondo consiste in uno spettacolo di droni luminosi a Torino⁶ per la festa patronale.

Raccogliendo l'attenzione di centinaia di utenti, i due eventi e la polisemia della parola hanno fatto sì che la comunità di ascoltanti il genere musicale *drone* fosse sproporzionatamente più grande rispetto alle altre nei giorni degli eventi; ma tale distorsione è diluita dalla divisione in fasce orarie e dalla moltitudine dei giorni in cui è stata effettuata la raccolta dati.

¹<https://github.com/pkasela/Sound-of-Data>

²<https://hadoop.apache.org>

³<https://neo4j.com>

⁴<https://kafka.apache.org>

⁵http://www.ansa.it/sito/notizie/mondo/mediooriente/2019/06/20/iran-pasdaran-abbattuto-drone-usa_b8922ea0-e61c-436b-8f66-60f3be32ed21.html

⁶<https://www.guidatorino.com/eventi-torino/san-giovanni-torino-droni-2019/>

Parte II

Costruzione dello *knowledge graph*

1 Data Cleaning con Python e Pig

I dati sono scaricati dal sito di musicbrainz attraverso lo script python `get_data.py`⁷ che cerca sul sito l'ultima versione del database e la scarica col comando `wget` dei sistemi Unix. Decompressa quindi la cartella compressa con il dump del database e selezionati solamente i file necessari allo scopo, sono eseguite delle operazioni `cat-sed` per effettuare una rapida pulizia preliminare (la gestione dei valori nulli \N). È salvata inoltre una lista di generi musicali, già presente in una pagina HTML del sito, tramite uno *scraper*, che ne analizza l'elenco per importarla in Python.

I file sono successivamente importati e processati all'interno del *file system* HDFS. Nel dettaglio, questa fase viene eseguita attraverso l'utilizzo di Apache Pig, piattaforma atta all'elaborazione di grosse quantità di dati. La decisione di utilizzare Pig anziché SQL dipende principalmente dalle migliori prestazioni del primo nella fase di caricamento, e dalla maggiore versatilità della suite Hadoop per la fase di elaborazione. L'*engine* adoperato in combinazione con Pig non è MapReduce ma Tez: costruito sopra YARN, migliora drasticamente le prestazioni mantenendo la stessa scalabilità evitando alcune fasi di scrittura di risultati parziali sull'HDFS (vedere Figura 1 per maggiori dettagli).

Più precisamente, la velocità di esecuzione è ottenuta grazie all'invio diretto dei dati da un processo all'altro, evitando la scrittura all'interno del file system, eccezion fatta per i *checkpoints*. Inoltre, la definizione del *job* sfrutta i *Directed Acyclic Graph*, dove i verti-

ci rappresentano gli step del processo e gli archi la connessione tra i vertici *Producer* e *Consumer* (come mostrato in Figura 2).

Conclusa la fase di pulizia, i dati sono trasferiti all'esterno di HDFS, sul filesystem reale, unendo i risultati di Pig mediante il comando `cat` di *Unix*. A questo punto i dati sono pronti per essere importati in Neo4j.

2 Import dei dati in Neo4J

Infine, il risultato dell'elaborazione svolta da Pig viene caricato all'interno del *DBMS* *Neo4j*, usando il comando `neo4j-import`, indicizzando i *gid* delle varie entità, in modo da velocizzare la ricerca da svolgere successivamente.

Il risultato è un grafo dal peso approssimativo di 5Gb contenente i vari artisti, i relativi album e le relative case discografiche collegate in modo opportuno. Si è preferito ridurre la capillarità del database non inserendo le singole canzoni contenute negli album, sia perché una simile precisione sarebbe stata difficile da gestire in fase di analisi, sia a causa della mancanza nel dump del database di tale informazione (comunque presente nella versione online), che avrebbe richiesto un tempo superiore alla settimana (dovuto perlopiù alla latenza supportata dal sito) per poter essere recuperata.

⁷<https://github.com/pkasela/>

Sound-of-Data/blob/master/musicbrainz_data/Data_Cleaning/get_data.py

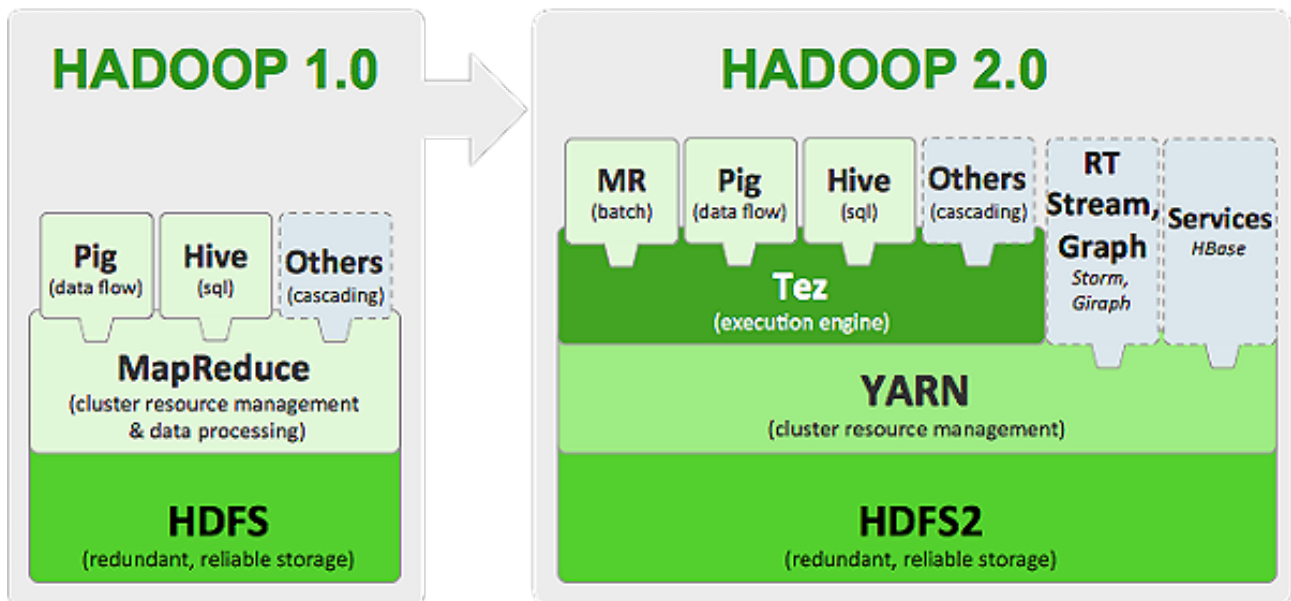


Figura 1: Differenza architetturale di HADOOP 1.0 con MapReduce rispetto alla versione 2.0 con Tez.

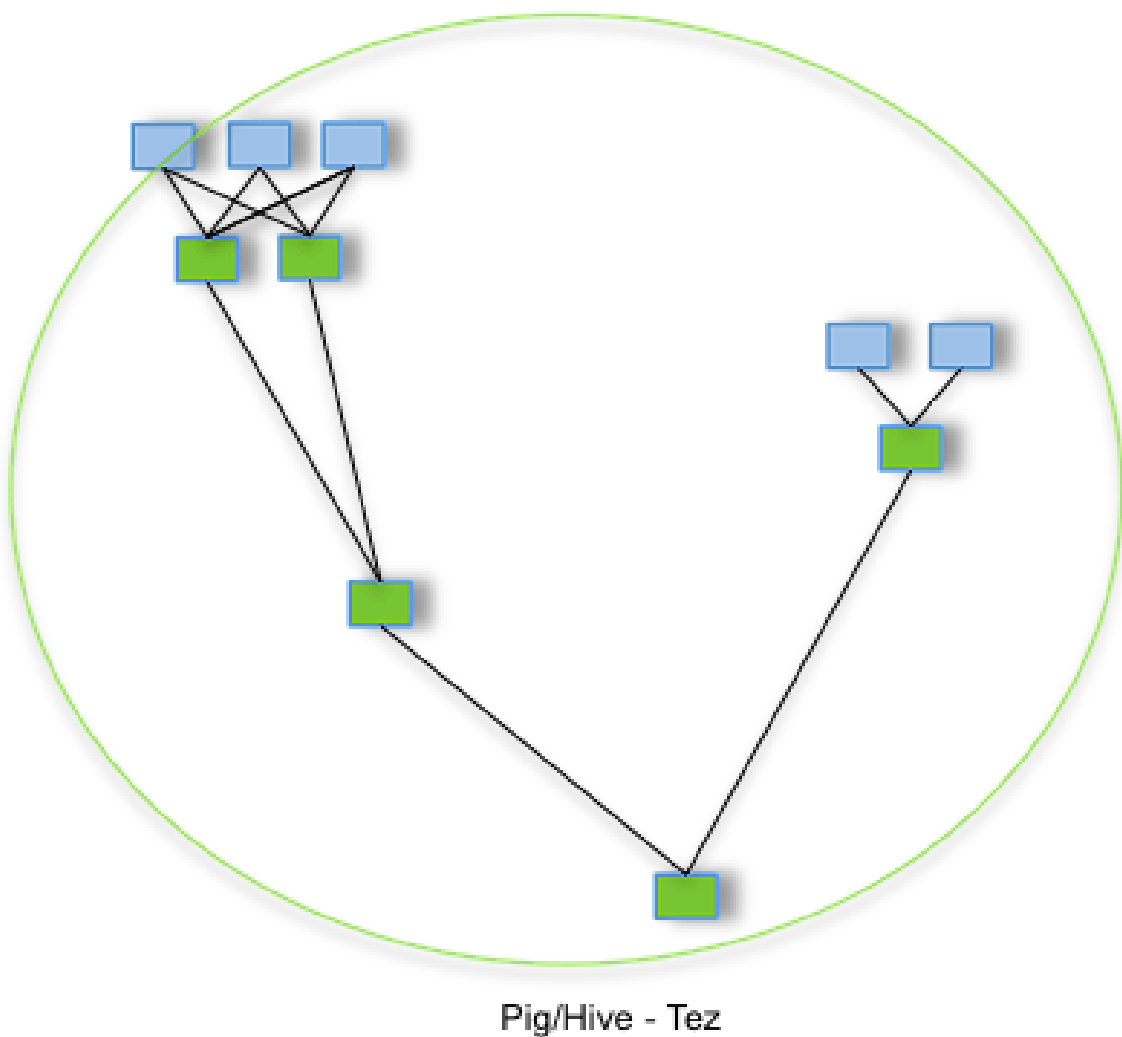


Figura 2: Funzionamento di Tez.

Parte III

Analisi dei tweet

3 Elaborazione mediante Apache Kafka

Nella fase di gestione e analisi in tempo reale, i *tweet* sono filtrati ed elaborati mediante una sequenza di operazioni, attraverso una *pipeline* costituita dalla concatenazione di due coppie di processi Producer e Consumer. L'architettura di elaborazione in *streaming* è stata così configurata per garantire una maggiore scalabilità ed una migliore performance, riuscendo a reggere maggiori volumi di *tweet* in ingresso. Questi vengono estratti grazie all'API *Tweepy*⁸ per Python, utilizzando come parole chiave di ricerca una selezione di 400 generi musicali fra i 419 ottenuti dallo scraping della pagina '<https://musicbrainz.org/genres>', e ricevendo soltanto i testi pubblicati in lingua italiana.

Quindi sono inviati ad una prima Topic come messaggi di un *Producer* di Kafka, inizializzato mediante la libreria *Kafka-python*⁹. All'interno di un secondo file¹⁰, associato ad un nuovo processo Producer, dunque ad una seconda Topic connessa in *streaming* con la precedente, il JSON grezzo che costituisce ciascun *tweet* in ingresso viene processato in vari step dalla funzione `tweet_preparations`.

La prima fase coincide con l'estrazione dello *screen name* dell'utente, del contenuto postato e dell'ora e data precisa di pubblicazione. In particolare, il testo viene modificato in modo da rendere leggibili i caratteri speciali e i cosiddetti *escape characters*.

La seconda fase sfrutta l'API di *Botometer*¹¹

⁸<https://tweepy.readthedocs.io/en/latest/>

⁹<https://kafka-python.readthedocs.io/en/master/>

¹⁰https://github.com/pkasela/Sound-of-Data/blob/master/NL/Kafka_Producer_2.py

¹¹<https://rapidapi.com/OSoMe/api/botometer/details>

per calcolare la probabilità - definita *score* - che un profilo sia in realtà gestito da un BOT, osservandone il comportamento passato. Poiché i profili gestiti in modo automatico alterano, spesso fortemente, la discussione mediatica aumentando arbitrariamente il flusso di determinati contenuti, si è deciso di non archiviare i *tweet* appartenenti a tali profili. Nel dettaglio, se lo *score* associato all'utente di ogni *tweet* supera una soglia (fissata arbitrariamente al 90%), lo *screen name* viene memorizzato all'interno di una *blacklist*, altrimenti smistato in una *whitelist*, in modo che il programma possa riconoscere più velocemente ogni profilo ed evitare inutili ricalcoli e così da ridurre il numero di richieste fatte all'API, gratuita fino ad un certo volume e poi a pagamento. Per aumentare ulteriormente l'efficienza di questo processo e garantire la scalabilità, si è scelto di adoperare il *data store Riak*, di tipo *NoSQL Key-Value*, attraverso la propria libreria in Python¹²: è così archiviato lo *score* di ogni utente e la richiesta all'API è effettuata solamente se non è trovata la chiave nel database.

La terza fase verifica che il testo sia effettivamente legato al tema musicale grazie ad un rudimentale *istance matcher* che tenta di riconoscere le parole riguardanti il contesto musicale, aggiungendo al dizionario prodotto dalle fasi precedenti tre chiavi corrispondenti rispettivamente alla lista di artisti, album o canzoni estratti dal testo analizzato. Ciò è permesso dall'uso dell'API¹³ fornita dal sito *musicbrainz.org* e basata su Apache Lucene¹⁴. Per maggiori dettagli si rimanda alla sezione dedicata.

Ogniquale volta si verificano le due condizioni note, cioè che l'utente non abbia probabilità elevata di essere un BOT e che il testo pubblicato si riferisca effettivamente al mondo della musica, ciascun risultato della funzione `tweet_preparations` è ricodificato e passato come messaggio al secondo *Producer* di *Apa-*

¹²<https://github.com/basho/riak-python-client>

¹³<https://python-musicbrainzngs.readthedocs.io/en/v0.6/api/>

¹⁴<https://lucene.apache.org/>

che *Kafka* - menzionato in precedenza - che effettuerà l'invio del dato finale alla *topic* nominata *KafkaTopic*. Grazie all'utilizzo di questa piattaforma, il flusso di *feed* è gestito in tempo reale, assicurando bassissima latenza e grande scalabilità.

4 Consuming dei Tweets in Neo4j

Un principale vantaggio di *Apache Kafka* consiste nella capacità di connettersi efficientemente a sistemi esterni, garantendo in tal modo la continuità dello *stream* di dati da una fonte - la class *Listener* adoperata in Python per "ascoltare" i tweet grazie all'API Tweepy - ad una "*landing zone*", che in questo caso corrisponde al DBMS *Neo4j*. Si è scelto di effettuare lo *storage* dei tweets in un *GraphDB*, e in particolare *Neo4j*, tra le varie motivazioni, per le sue caratteristiche di estrema scalabilità, efficienza nella gestione dei dati, elevata capacità di adattamento al fenomeno di studio e facile interpretabilità. L'esportazione diretta dei tweet da *Apache Kafka* all'interno del *data store* avviene grazie a un *plug-in* di *Neo4j*, "*Neo4j Streams Consumer*". Questo si configura come un'applicazione che effettua un'ingestione diretta e automatizzata all'interno di *Neo4j*, permettendo la lettura dei dati presenti nella *topic* identicamente a qualsiasi applicazione *Consumer* di *Kafka*.

Neo4j Streams Consumer permette all'utente di specificare arbitrariamente le relazioni, entità e proprietà in cui i *payloads* di *Apache Kafka*, corrispondenti al JSON di ciascun tweet importato, dovranno essere organizzati per costruire progressivamente il grafo. La struttura di quest'ultimo viene dichiarata attraverso un *Cypher template*, ovvero un insieme di *queries* semantiche di *Cypher*, all'interno di un file di configurazione, che nel caso presentato è `docker-compose.yml`¹⁵, poiché il progetto viene eseguito mediante *Docker De-*

sktop, in quanto quest'ultimo presenta il vantaggio di riuscire a far comunicare automaticamente e in maniera trasparente i containers coinvolti. Una volta specificato il *Cypher template*, installato il *plugin* all'interno del *Docker*, e successivamente montati i container e i relativi collegamenti, la fase di mera esecuzione del progetto avviene mediante la funzione denominata `streams.consume`, la quale crea immediatamente la struttura del grafo come indicato.

Nello specifico, all'arrivo di un *tweet* verranno dunque aggiunti i nodi "*Tweet*", "*User*" ed un terzo che riporterà ora e data di pubblicazione. Il primo conterrà le *property keys* del testo "*text*", e dello "*screen_name*". Il secondo, invece, conterrà solamente la proprietà "*name*", che riporta lo *screen name*. Tra primi i due nuovi nodi verrà generata la relazione `BELONGS_TO`, ad indicare il profilo di appartenenza di ciascun post, mentre "*User*" sarà connesso all'ora/data di pubblicazione da un *edge* etichettato come "*Twitted*".

In più il programma effettua automaticamente un'operazione di *merge* della coppia appena generata coi nodi già presenti nel grafo: quanto è individuato l'id di un'entità presente nel testo, il software collega il nodo del tweet all'entità con una relazione `TALKS_ABOUT`. L'entità presente su *musicbrainz* è collegata al proprio genere musicale, così da svolgere da collegamento tra i nodi riguardanti genere e tweet, oggetto della presente analisi.

Dunque il grafo conterrà una serie di nodi utente, raggruppati in comunità attorno ad artisti o generi a seconda dell'oggetto dei propri *tweet*. Grazie a questa semplice struttura, l'analisi può avvenire già in maniera grafica e trasparente, e l'utente è in grado di distinguere la grandezza di queste comunità, intendere quali legami esistano fra di esse, e quanto questi siano intensi. Per ottenere queste informazioni, è infatti sufficiente scaricare un *dump* del grafo a intervalli regolari e confrontare i vari risultati. Il risultato finale, dall'interfaccia web di *Neo4j*, apparirà come nella Figura 3; invece la composizione dei cluster, cioè delle

¹⁵<https://github.com/pkasela/Sound-of-Data/blob/master/Neo4j\%20\%26\%20kafka/docker-compose\%20configuration.yml>

comunità ben separate, verrà mostrata come nella Figura 4.

5 Riconoscimento delle istanze nel testo

Data la quantità e qualità dei dati, si è deciso di costruire un strumento di *instance matching* creato *ad hoc* per i tweet: le API disponibili basate su *deep learning* sono spesso allenate su altri tipi di testo (articoli di giornale o opere letterarie) e spesso non offrono una velocità adeguata per l'analisi in tempo reale. Il problema della qualità del testo è fondamentale: serve uno strumento che riesca ad analizzare il testo e a trovare le entità corrispondenti senza bisogno di contestualizzazione. Per tale scopo, filtrati già in partenza i tweet grazie a parole chiave a tema musicale, si costruisce un modello che identifichi le istanze restituendone gli *id* grazie all'API di musicbrainz disponibile per Python.

5.1 Identificazione delle entità

Le entità sono riconosciute non mediante *machine learning* ma grazie a semplici stratagemmi linguistici e l'aiuto di un correttore ortografico. Prima di tutto il testo del tweet è ripulito da eventuali abbreviazioni gergali (che vengono sostituite dalla propria forma estesa), link e caratteri speciali. Subito dopo sono ricercate le parole nel testo che non risultano essere italiane: analizzando tweet in lingua italiana, il cui tema è sicuramente musicale, si presume che una stringa in lingua diversa abbia una certa importanza; ciò è fatto sia analizzandone l'ortografia sia verificando con un'espressione regolare che le parole rispecchino la regola di costruzione delle sillabe in italiano¹⁶. In tal caso, o la parola è un sostantivo della quinta classe (parole straniere entrate nell'uso comune) o un termine

in lingua non italiana e quindi degno di attenzione. Infine sono analizzate le parole la cui iniziale è maiuscola, col supporto di un albero sintattico (costruito grazie ad un programma estero, TreeTagger¹⁷, interfacciato con Python grazie ad un wrapper¹⁸ fornito dagli stessi sviluppatori) per verificare che siano nomi propri.

Grazie all'uso della punteggiatura e delle preposizioni, si tenta di stabilire se la presunta entità è un artista o un'opera, verificando poi con certezza confrontando con lo *knowledge graph*.

Tramite le API di musicbrainz, attraverso una serie di funzioni *search* è possibile consultare il database sfruttando un server di ricerca costruito utilizzando la tecnologia *Lucene*. Analizzate le entità e tentato di stabilire cosa rappresentino, si interroga la base di conoscenza con le funzioni *search_artists* e *search_recordings* per le entità riconosciute rispettivamente come artisti e opere, a cui si aggiunge *search_release_groups* (funzione relativa alla ricerca degli album), per quelle di tipologia incerta. Per verificare la corrispondenza tra entità e risultati di ricerca, si è scelta la politica di accettare il risultato la cui similarità coseno è maggiore, a patto che superi una soglia arbitraria di 0.95. Inoltre, per aumentare la precisione, si sono incrociati i risultati parziali delle ricerche per testare se nello stesso testo compaia l'opera accompagnata dal proprio autore.

5.2 Prestazioni del modello e margini di miglioramento

La problematica principale emersa nella costruzione del modello è stata la corretta attribuzione dell'entità nei casi di omonimia: questi casi infatti sono molto frequenti tra canzoni ed album relativi ad artisti differenti

¹⁶Una sillaba in italiano è composta da un massimo di tre consonanti iniziali, vocale ed eventuale consonante finale.

¹⁷<https://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

¹⁸<https://perso.limsi.fr/poital/doku.php?id=dev:treetaggerwrapper>

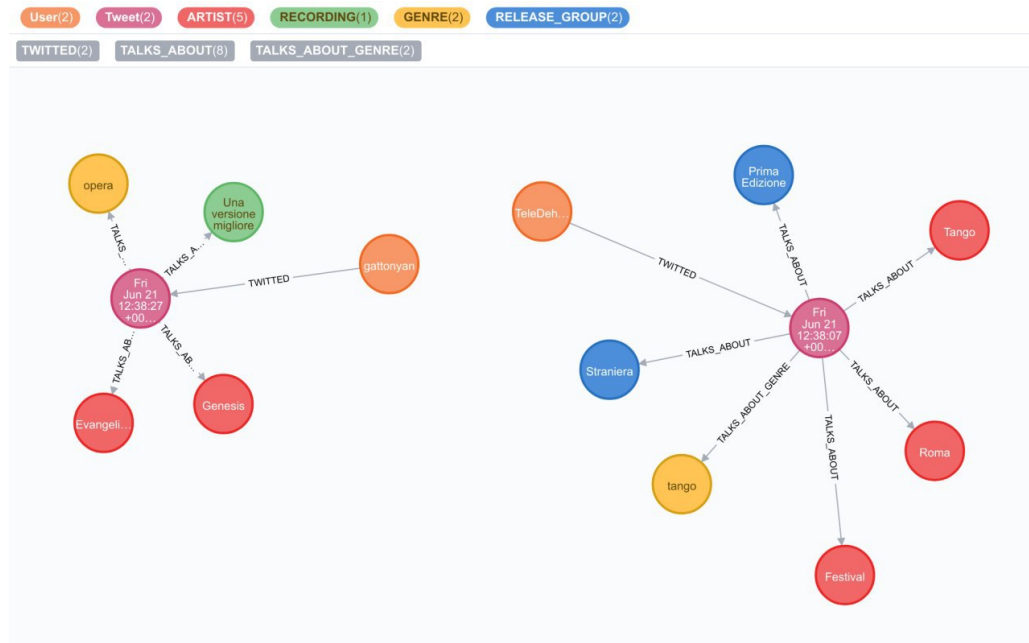


Figura 3: Configurazione della struttura del grafo.

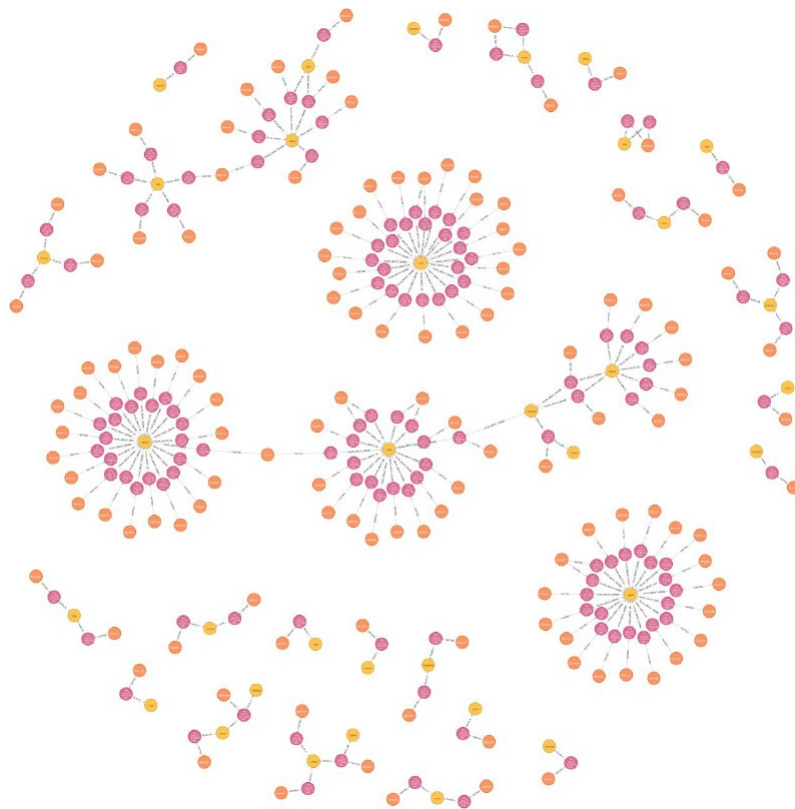


Figura 4: Composizione parziale delle comunità all'interno del grafo.

ed è quindi possibile che il risultato ottenuto non sia quello desiderato.

Inoltre, nonostante il modello offra una buona *performance* nella selezione dei testi riferiti al mondo musicale e nel riconoscimento delle

istanze nel testo, continua a presentare alcune imperfezioni. Un primo difetto deriva dalla difficoltà nella gestione dei casi polisemici: alcuni nomi di genere (in particolare *drone*, *dance*, *club* e *opera*) sono polisemici,

distorcendo i valori dell'analisi. Il programma utilizzato è troppo rudimentale per distinguere il contesto e quindi mostra numerosi falsi positivi: il termine *club* in particolare viene spesso citato all'interno di testi riguardanti qualsiasi tipo di associazione; la parola *drone*, come già accennato nell'introduzione, indica, nella maggior parte dei casi, l'aeromobile; l'espressione "dance" è frequentemente legata alla musica e *opera* possiede troppi significati per poterli elencare.

Un'altra problematica, non legata al modello ma a tweepy è sorta dopo aver rilevato alcuni *tweet* in lingua straniera (perlopiù spagnola) che citavano il termine *trap* senza riferirsi al genere musicale.

Pertanto, il modello raggiunge buoni risultati nell'ambito di tale progetto, rendendo comunque necessari ulteriori tentativi di miglioramento della sua capacità di riconoscere il contesto all'interno del quale le parole vengono citate.

Parte IV

Visualizzazioni dei dati

5.3 Distribuzione dei principali generi musicali nelle quattro fasce orarie

La prima visualizzazione consiste nella tipologia *barchart*, realizzata attraverso il software *Tableau*. In ascissa, il grafico raccoglie i principali quarantasette generi musicali estratti dall'analisi delle pubblicazioni avvenute tra il 21 e il 28 giugno, ordinati in senso crescente secondo la frequenza assoluta totale. L'ordinata, invece, è segmentata nelle quattro sezioni dedicate rispettivamente al periodo notturno, mattutino, pomeridiano e serale. In ogni sezione, la dimensione della "barra" corrisponderà pertanto al numero di volte in cui il dato genere è stato menzionato durante la specifica fascia oraria. Pertanto, sommando verticalmente i valori, si otterrà la dimensione globale associata ad ogni genere. Tuttavia, sulle misure parziali, relative quindi al singolo intervallo temporale, è stato necessario effettuare un filtraggio: dato che si è deciso di mantenere la scala naturale dei valori, le barre associate alle frequenze minori di 5 in valore assoluto sono state omesse, per evitare che l'ascissa risultasse carica di un eccessivo numero di generi, peggiorando la fruizione, e poichè, in ogni caso, esse avrebbero assunto dimensioni eccessivamente ridotte e non sarebbero comunque state sufficientemente visibili.

Perciò, considerando la sola variabile della frequenza assoluta, dal punto di vista dell'"overview" il fruitore è in grado di individuare la distribuzione dei vari generi, e intuire immediatamente quali siano i più presenti. Ad esempio, è chiaro come "dance" sia il genere comparso più volte durante il pomeriggio, o ancora "rave" il più presente durante la mattina. Inoltre, è possibile notare come la notte rappresenti ovviamente il periodo meno

produttivo per quanto riguarda il numero di pubblicazioni, e, viceversa, come la mattina accolga un flusso di tweet a tema musicale nettamente maggiore. Infine, la componente dell'interattività permette di ottenere i dati di ogni singolo genere, approfondendo le sue caratteristiche e valutando la sua disposizione nell'arco della giornata.

Una seconda variabile è stata aggiunta al grafico: la misura di appartenenza di ogni genere musicale ad ogni fascia oraria, ovvero la tendenza percentuale di ogni genere a comparire in una particolare fascia oraria, rispetto al totale dei tweet che lo riguardano. Quanto maggiore è la concentrazione delle citazioni ad esso riferite in uno dei quattro periodi, tanto superiore sarà il suo grado di appartenenza al dato intervallo temporale. L'indice, che accoglie valori da zero a uno, è rappresentato attraverso il dettaglio della saturazione del colore blu in ogni barra. A questo punto, il fruitore non è soltanto in grado di percepire visivamente la presenza di un particolare genere, ma anche il suo diverso legame con i quattro momenti della giornata. Con l'obiettivo di evidenziare i valori più interessanti sotto questo punto di vista, le prime tre maggiori percentuali di appartenenza, per ogni intervallo temporale, vengono specificate sopra le relative barre. In più, il massimo e il minimo vengono descritti testualmente. È così possibile osservare come la musica classica rappresenti il genere musicale più specifico, poichè citata nel 77% dei casi durante la sera. Una sezione speciale del grafico, sul lato destro, viene dedicata alle quattro tipologie musicali più presenti. Innanzitutto, questo *zoom* permette al fruitore di cogliere in maniera migliore la differenza e la composizione interni dei quattro maggiori elementi, potendo effettuare un miglior confronto visivo. In secondo luogo, la rappresentazione dell'estremo inferiore e superiore dell'intervallo di confidenza calcolato al 95% permette di capire se e quali differenze fra le frequenze assolute siano significative, in modo da attribuire una validità statistica ai parametri estratti dalla distribuzione campionaria. Nel dettaglio, il numero

di volte in cui la categoria "*dance*" è comparsa all'interno delle dinamiche della discussione su Twitter appare significativamente maggiore rispetto a quello relativo a "*rave*", ma non rispetto a "*jazz*" e "*rock*".

5.4 Frequenza assoluta dei generi nella discussione a tema musicale

Un approfondimento grafico sulla presenza delle principali categorie all'interno delle dinamiche della discussione virtuale è offerto mediante realizzazione di un *wordcloud* mediante linguaggio di programmazione *Python*. Grazie alla semplicità di questa visualizzazione, l'osservatore può percepire visivamente la differenza fra le varie grandezze, ordinate a comporre la forma di una nota musicale.

5.5 Critiche alla visualizzazione e margini di miglioramento

Le critiche poste dal campione di fruitori durante l'osservazione e la valutazione del grafico *barchart* hanno riguardato maggiormente la mancanza di un adeguato livello di intuitività della visualizzazione, dovuta principalmente alla presenza di molti elementi grafici in uno spazio ridotto. Nonostante ciò, questa problematica non pare aver mai penalizzato gravemente la fruizione del grafico e la comprensione finale dell'oggetto descritto dall'analisi. Infatti, durante l'esecuzione di un'attività interattiva - che consisteva nel chiedere di selezionare alcuni generi in base a precise caratteristiche -, gli utenti coinvolti sono riusciti in ogni caso a soddisfare correttamente le richieste.

Per quanto riguarda il *wordcloud*, non sono state poste critiche, e la visualizzazione è stata in generale apprezzata per la sua semplicità.

5.6 Risultati dei questionari psicometrici

Parte V

Risultati e conclusioni

Nel breve periodo in cui la macchina virtuale è rimasta online a raccogliere i *tweets*, sono stati raccolti