

Università di Roma Tor Vergata
Corso di Laurea magistrale in Ingegneria Informatica
Dipartimento di Ingegneria Informazione



Analisi della polarizzazione di Endorsement
Graph, attraverso sentiment analysis

Relatore:

Giuseppe F. Italiano

Correlatore:

Nikos Parotsidis

Candidato:

Alessandro Valenti

matricola 0228709

Anno Accademico 2016-2017

Sommario

La diffusione delle informazioni e di opinioni sin dai tempi antichi ha generato conflitti di ogni genere. Tali problematiche sono sempre più evidenti all'interno delle reti sociali che si vengono a creare mettendo in contatto individui con pensieri ed idee differenti tra loro. I conflitti vengono generati in base al tipo di argomento e quanto tale è "caldo" per gli utenti in questione. La polarizzazione è un utilissimo strumento per lo studio e l'analisi delle opinioni in differenti aree di ricerca all'interno di una rete sociale. Generalmente, la polarizzazione può essere applicata all'interno di contesti politici, sociali e culturali permettendo di comprendere al meglio quali siano gli schieramenti delle persone riguardo tali argomenti. La possiamo definire come: *Divisione in due gruppi fortemente contrastanti per una serie di opinioni o credenze*. Questo processo di analisi può assumere diversi significati a seconda dello scenario studiato. Come ad esempio: *Polarizzazione Politica* (divergenza di opinione su estremi ideologici) o *Polarizzazione Sociale* (differenza di opinione all'interno delle società che possono nascere da disuguaglianze sociali ed economiche).

Un problema che sta affliggendo i social media è la formazione degli *Echo-chambers*, cioè quelle comunità che condividono le stesse opinioni rafforzando il proprio punto di vista senza cercare di verificare la presenza o meno di altre opinioni. La polarizzazione è uno strumento che può essere utilizzato per identificare questi gruppi di utenti, ma può assumere anche altri compiti come lo studio delle opinioni e la loro diffusione all'interno della rete. Questa tesi è stata sviluppata per cercare di studiare il comportamento delle informazioni espresse nei social media, in questo caso si è scelto *Twitter*. I dati, ovvero i *Tweet*, sono stati raccolti attraverso una ricerca per *hashtag* e classificati in due gruppi di pensiero distinti rispetto all'argomento analizzato. La raccolta dei dati è stata effettuata all'interno di una finestra temporale per poter studiare la polarizzazione e la sua evoluzione nel tempo. Il primo problema che si è cercato di risolvere è stato quello di classificare le notizie attraverso una comprensione testuale. Tale operazione è stata resa possibile attraverso la

tecnica della *Sentiment Analysis*, che consente di classificare i messaggi degli utenti attraverso un'analisi del sentimento. Una volta suddivise le opinioni all'interno della rete è stato costruito un grafo contenente i tweet ed i relativi retweet, questo viene definito *endorsement graph*. Partendo dal grafo sono stati applicati due algoritmi per il calcolo della polarizzazione definiti come segue:

- *Algoritmo basato sul grado del grafo*: Una volta generato il grafo vengono catalogati i nodi in due categorie:
 - *Elite*: l'utente che ha *tweettato* un'opinione.
 - *Listener*: l'utente che ha *retweettato* il tweet di uno o più nodi *Elite*.

Partendo da queste categorie viene calcolata la polarizzazione sfruttando il grado di ogni nodo, assegnando un valore numerico, in base al gruppo di appartenenza (*Elite*) e poi calcolare la polarizzazione attraverso la formula espressa nel paper [13].

- *Algoritmo basato sulla topologia*: La differenza principale con il precedente algoritmo, consiste nell'utilizzare la tecnica dei *Random Walk*. La polarizzazione dipende fortemente dalla topologia dell'*endorsement graph* e dalla probabilità sugli archi. Questa è definita probabilità di *retweet*, ed è pari al rapporto tra il numero di retweet fatti da un utente su un certo tweet sul numero totale di retweet che l'utente ha effettuato su quel topic.

Una volta calcolata la polarizzazione sono stati utilizzate tecniche di *Forecasting* come: *Double exponential smoothing*, *linear regression* e *moving average*. Queste tecniche consentono di poter predire il comportamento della rete nel futuro e quindi potrebbe essere utile in ottica di prevenzione di *Echo-chamber*. Tutte queste funzionalità sono state realizzate mediante librerie *Python*, mentre la raccolta dei dati è stata effettuata all'interno di una istanza *EC2*. I topic utilizzati per sperimentare il comportamento della polarizzazione sono stati due cioè: *elezioni regionali in Sicilia* ed il *Biotech*. Il motivo di tale scelta è dettata dalla curiosità di testare queste funzionalità per due argomenti che ricoprivano due contesti differenti tra loro; anche perché molto dibattuti. I risultati ottenuti per entrambi i topic sono stati molto soddisfacenti perché la polarizzazione ha assunto un comportamento in linea con la realtà. Per il primo il trend della diffusione delle opinioni ha rispecchiato quanto osservato nei risultati delle elezioni, infatti attraverso la raccolta dei tweet è stato possibile notare come nel mese di Novembre (mese in cui ci sono state le elezioni), prima del giorno delle elezioni

gli utenti che si schieravano verso la coalizione del centro-destra fosse simile a quelle del Movimento 5 Stelle. Subito dopo il giorno delle elezioni in conformità con i risultati ottenuti si è potuto constatare come gli utenti aderenti al centro-destra superassero quelli del Movimento 5 Stelle. Per sperimentare tale topic attraverso la polarizzazione è stato necessario considerare solo due forze politiche, in questo caso quelle vincenti. Per quanto riguarda il biotestamento si è potuto notare come il retaggio culturale e religioso di un paese avessero una forte influenza anche nei social-media, infatti la polarizzazione della coalizione degli utenti contrari era di gran lunga superiore a quella dei favorevoli. Lo studio di questo topic si può definire *hot*, perché ancora di attualità e lo si è potuto riscontrare raccogliendo i dati nel mese successivo al 14 Dicembre (giorno in cui è stata emanata la legge sul Biotestamento) notando come il numero dei tweet aumentava. Da questo progetto è possibile effettuare nuove ricerche per risolvere il problema delle comunità fortemente polarizzate attraverso una analisi basata sulla *controversia* cioè cercare di far comunicare le comunità isolate in qualche modo. Altri spunti potrebbero essere quello di raffinare la classificazione dei tweet attraverso uno studio delle immagini, media e link pubblicati all'interno dei tweet, perché la sentiment analysis per quanto possa essere raffinata non consente di poter percepire l'ironia all'interno del testo, ironia che può essere espressa attraverso elementi multimediali.

Indice

Sommario	i
1 Introduzione	1
1.1 Struttura Tesi	9
2 Stato dell'arte e Background	10
2.1 Stato dell'arte	10
2.1.1 Rete Sociale	11
2.1.2 Le reti	11
2.1.3 Twitter	12
2.1.4 Polarizzazione	14
2.1.5 Sentiment Analysis	17
2.1.6 Forecasting	18
3 Progetto logico della soluzione del problema	20
3.1 Raccolta Dati	20
3.2 Sentiment Analysis	22
3.3 Endorsement Graph	23
3.4 Polarizzazione	24
3.5 Predizione	25

4	Implementazione e realizzazione del sistema	28
4.1	Raccolta dati	28
4.1.1	Twitter Api	29
4.1.2	Get Old Tweet	31
4.2	Sentiment Analysis	34
4.3	Endorsement Graph	38
4.4	Polarizzazione	41
4.4.1	Polarizzazione basata sul grado	42
4.4.2	Polarizzazione basata sulla topologia	46
4.5	Predizione	50
4.5.1	Double Exponential Smoothing	50
4.5.2	Linear Regression	53
4.5.3	Moving Average	54
5	Test sperimentali e valutazione	56
5.1	Elezioni Regionali Sicilia 2017	58
5.2	Biotestamento	71
6	Sviluppi futuri e conclusioni	84
	Bibliografia	89

Capitolo 1

Introduzione

La diffusione delle informazioni e di opinioni sin dai tempi antichi ha generato conflitti di ogni genere, per contrapposizioni sociali, culturali, religiosi ed economici. Tali problematiche sono sempre più evidenti all'interno delle reti sociali che si vengono a creare mettendo in contatto individui con pensieri ed idee differenti tra loro. I conflitti vengono generati in base al tipo di argomento e quanto tale è "caldo" per gli utenti in questione. La polarizzazione è un utilissimo strumento per lo studio e l'analisi delle opinioni in differenti aree di ricerca all'interno di una rete sociale. Generalmente, la polarizzazione può essere applicata all'interno di contesti politici, sociali e culturali permettendo di comprendere al meglio quali siano gli schieramenti delle persone riguardo tali argomenti. Una generica definizione della polarizzazione è la seguente:

Divisione in due gruppi fortemente contrastanti per una serie di opinioni o credenze.

Questo processo di analisi può assumere diversi significati a seconda dello scenario studiato.

- *Polarizzazione Politica*: divergenza di opinione su estremi ideologici.

- *Polarizzazione Sociale*: differenza di opinione all'interno delle società che possono nascere da disuguaglianze sociali ed economiche.

La polarizzazione può comportare diversi cambiamenti sullo scenario in questione, in quanto mette in luce come la formazione di due grandi gruppi non consenta una diffusione democratica delle opinioni. A tal proposito è interessante notare come la divisione in queste due grandi partizioni generi alcune problematiche quali:

- La frammentazione della rete stessa.
- L'isolamento delle opinioni.

In conclusione potremmo definire la polarizzazione come un processo sociale per cui gli utenti che vi partecipano vengono divisi in due grandi sottogruppi aventi visioni, punti di vista ed opinioni differenti del problema in questione, con alcuni individui che rimangono neutrali tra i due grandi gruppi.

La formazione di due comunità isolate che non comunicano tra loro, comporta un problema di isolamento delle opinioni cioè un utente che appartiene a quel gruppo difficilmente potrà ricevere informazioni o aderire alle idee del gruppo antagonista. Otteniamo la formazione degli *Echo-Chambers*, definita come:

Una situazione in cui le informazioni, le idee e le credenze vengono rinforzate e amplificate perché espresse all'interno dello stesso ambiente, rimanendo isolato.

Un altro problema che può generare una forte polarizzazione delle opinioni e delle informazioni sono i *Filter Bubble* ovvero:

Uno stato di un isolamento intellettuale che può essere ottenuto a partire da risultati di ricerche su siti che registrano la storia del comportamento dell'utente. Questi siti sono in grado di utilizzare informazioni sull'utente per scegliere selettivamente tra tutte le risposte quelle che l'utente vorrà visualizzare. L'effetto è di isolare l'utente da informazioni che sono in contrasto con il suo punto di vista, isolandolo nella sua bolla culturale o ideologica.

Come precedentemente anticipato la polarizzazione è uno strumento che può essere facilmente utilizzato per individuare tutte queste problematiche all'interno dei moderni Social Network come *Facebook*, *Twitter* e molti altri. Questo perché gli utenti si sentono sempre più liberi di poter esprimere le proprie opinioni all'interno di queste piattaforme riguardo problematiche sociali, culturali, politiche ed economiche. Non è sempre possibile poter uscire dalle *Filter Bubble*, perché gli stessi social network tendono a indirizzare l'utente a visualizzare informazioni che potrebbero interessargli senza farli confrontare con opinioni divergenti. Alla luce di questo grande problema il calcolo di una polarizzazione può consentire agli amministratori dei social network di individuare i topic più polarizzati garantendo una diffusione democratica delle informazioni, facendo comunicare gli utenti con opinioni divergenti.

L'obiettivo della mia tesi consiste nell'utilizzare la polarizzazione per poter individuare gli argomenti fortemente polarizzati e comprendere come tali informazioni si diffondano all'interno della rete sociale. Lo sviluppo di questo strumento è stato effettuato attraverso due algoritmi, presentati nei seguenti paper:

- ***Measuring Political Polarization: Twitter shows the two sides of Venezuela***([13]: Studia la diffusione delle informazioni all'interno di un *endorsement graph* collezionando i dati relativi alle elezioni

in Venezuela all'interno del *social network Twitter*. Viene effettuato uno studio della polarizzazione all'interno di un contesto politico attraverso la diffusione delle opinioni sui candidati politici durante le ultime elezioni presidenziali, l'*endorsement graph* viene costruito partendo da un utente che pubblica nella rete un *Tweet* esprimendo la propria opinione, formando un nuovo nodo, mentre eventuali follower di quell'utente che *retwettano* tale notizia sono nuovi nodi all'interno del grafo con archi uscenti verso il nodo che hanno *retwettato*. In questo modo viene generato un grafo basato sul *retweet*. Una volta generato il grafo vengono catalogati i nodi in due categorie:

- *Elite*: l'utente che ha *tweettato* un'opinione.
- *Listener*: l'utente che ha *retwettato* il tweet di uno o più nodi *Elite*.

Partendo da queste categorie viene calcolata la polarizzazione sfruttando il grado di ogni nodo, tale operazione viene eseguito iterativamente fino ad ottenere una stabilizzazione della polarizzazione.

- ***Reducing Controversy by Connecting Opposing Views***[15]: Identifica la polarizzazione sfruttando la topologia del grafo. Il grafo viene generato utilizzando la medesima tecnica del precedente paper, così come il social network di riferimento *Twitter*. La differenza principale con la soluzione proposta in precedenza, consiste nell'utilizzare la tecnica dei *Random Walk*. La polarizzazione adottando questo approccio dipende fortemente dalla topologia dell'*endorsement graph*.

Prima di poter effettuare il calcolo vero e proprio della polarizzazione occorre effettuare una prima scrematura, da intendersi come una prima classificazione delle opinioni in due gruppi contrastanti, nel dettaglio attraverso

la *Sentiment Analysis*. Questa particolare tecnica consente di partizionare il grafo in due gruppi che per semplicità chiameremo **Rossi** e **Blu**, nel dettaglio viene analizzato il testo contenuto in un tweet o in un post (a seconda del social network adottato) catalogandolo per un gruppo piuttosto che un altro a seconda del contenuto e all'affinità col topic in questione. Per meglio comprendere cosa viene effettuato presentiamo la definizione di *Sentiment Analysis*:

L'Analisi del sentiment o Sentiment analysis (ma anche opinion mining) é la maniera a cui ci si riferisce all'uso dell'elaborazione del linguaggio naturale, analisi testuale e linguistica computazionale per identificare ed estrarre informazioni soggettive da diverse fonti.

In conclusione l'analisi semantica consente di poter catalogare le informazione in base alla loro vicinanza alle opinioni di un gruppo piuttosto che ad un altro, ed eventualmente scartare quelle informazioni che non sono di alcun interesse per il calcolo della polarizzazione. Tale operazione è possibile soltanto se la macchina é stata precedentemente istruita sul topic in questione, infatti si definisce *training set* l'insieme delle informazioni di riferimento che consentono alla macchina di poter distinguere le opinioni a seconda del loro contenuto.

Dopo aver effettuato questa separazione o catalogazione delle informazioni é possibile identificare quali utenti siano più o meno vicini ai due poli di un'opinione. Ricapitolando partendo da *post* o *topic* viene effettuata la *Sentiment analysis*, viene costruito l'*endorsement Graph* ed infine calcolata la **polarizzazione**. Diamo ora qualche informazione in più sulla polarizzazione, a livello matematico la polarizzazione può assumere valori compresi tra $[-1, 1]$ definendo in questo modo due poli opposti. I nodi che avranno una

polarizzazione pari a 0 sono da ritenersi nodi neutrali ovvero che non sono soggetti ad una forte polarizzazione, ma sono l'emblema della democrazia, in quanto ricevono informazioni da entrambi i gruppi.

Per concludere è stato effettuato anche uno studio per poter consentire la predizione del valore della polarizzazione in un periodo futuro. In questo modo gli amministratori dei social network possono effettuare degli accorgimenti alla rete consentendo una democratica diffusione delle opinioni, senza creare *Echo-Chambers* e *Filter Bubble*. La predizione è stata realizzata attraverso tecniche di *Forecasting* molto utilizzate in contesti economici, in quanto consentono, attraverso delle serie numeriche, di poter predirne il valore nell'istante temporale successivo. Sfruttando queste particolarità è stato possibile effettuare una predizione, nel dettaglio le tecniche utilizzate sono tre:

- *Double exponential smoothing*
- *Linear regression*
- *Moving average*

Terminiamo questa sezione presentando i casi di studio utilizzati. Per lo sviluppo della mia tesi ho deciso di analizzare la polarizzazione all'interno di due contesti differenti, attraverso questi due Topic:

- **Elezioni Regionali in Sicilia nel 2017:** analisi in un contesto politico.
- **Biotestamento:** analisi in un contesto sociale.

I dati relativi a questi due topic sono stati raccolti attraverso il social network *Twitter*, dal 01/09/2017 al 20/12/2017, per il primo topic, mentre per il secondo il periodo di raccolta è compreso tra il 01/09/2017 al 31/01/2018. Il

periodo indicato é stato scelto per analizzare l'evoluzione della polarizzazione nel tempo comprensiva della conclusione di questi due topic. Il 05/11/2017 si sono svolte le elezioni regionali in Sicilia ed il 14/12/2017 il parlamento italiano ha approvato la legge sul biotestamento. I dati sono stati raccolti effettuando una ricerca attraverso i 5 *hashtags* più utilizzati per entrambi i topic. Questi *hashtag* non esprimono nessuna opinione o parere presi singolarmente ma sono delle parole chiavi necessarie per catalogare il contesto del *tweet*. Per classificare i dati raccolti e poi valutarne la polarizzazione sono state adottate le tecniche precedentemente illustrate.

Per quanto riguarda le elezioni regionali in Sicilia si è deciso di raccogliere i tweet relativi alle due grandi fazioni che hanno dominato la scena politica siciliana:

- Il *Movimento 5 Stelle*.
- *Forza Italia* (la coalizione del centro destra).

Innanzitutto è stata una scelta dettata dai risultati conseguiti durante le suddette elezioni e dal fatto che in Italia non sono presenti soltanto due fazioni politiche, come in molti altri paesi del mondo, quindi sarebbe risultato impossibile definire un valore polarizzato se avessimo considerato più di due fazioni politiche. A tal proposito sono stati scartati i dati relativi a quei candidati appartenenti ad altri partiti politici e coalizione rispetto a quelli sopra elencati, utilizzando la *Sentiment Analysis*. I risultati ottenuti da questo topic hanno un comportamento interessante, cioè il cambiamento nel tempo della polarizzazione, seguendo il trend riscontrato durante i sondaggi effettuati mensilmente. Nel dettaglio si può facilmente assistere ad un cambiamento di trend col passare del tempo. Si comincia con una polarizzazione di $\simeq 79\%$ a favore del *Movimento 5 Stelle* per concludere alla fine del suddet-

to periodo con un capovolgimento di fronte con il $\simeq 59\%$ della polarizzazione a favore della coalizione di *Forza Italia*, mostrando un cambiamento radicale nel tempo conforme con quanto accaduto nei sondaggi.

Per quanto riguarda il *Biotestamento* si è deciso di raccogliere i tweet relativi alla legge approvata il 14 Dicembre 2017 dal parlamento italiano, per analizzare la polarizzazioni in un contesto sociale. La polarizzazione riguardava l'adesione o meno a questa legge, riscontrando una fortissima polarizzazione verso i contrari all'attuazione di tale legge. Questi risultati sono conformi al contesto sociale e religioso presente in Italia, confermando quanto spiegato in precedenza e cioè quanto un retaggio culturale o religioso possa influenzare il pensiero e le opinioni di una persona. Questa considerazione non è applicabile soltanto nella vita di tutti i giorni ma anche all'interno di un social network e ciò viene dimostrato dai risultati ottenuti all'interno di questo topic.

In conclusione questi due Topic hanno contribuito a confermare quanto precedentemente spiegato all'interno di questo capitolo e cioè che la polarizzazione è un potentissimo strumento che consente di poter individuare gli *Echo-chambers* presenti nella rete. Eventuali sviluppi futuri possono riguardare l'eliminazione degli *Echo-chambers* abbassando il livello di controversia tra i due gruppi ottenuti attraverso la polarizzazione, attraverso un congiungimento tra quei gruppi di nodi che condividono sempre le stesse opinioni.

1.1 Struttura Tesi

All'interno di questa tesi, verranno illustrati e sviluppati:

- dettagli implementativi
- dettagli teorici
- dettagli sperimentali
- sviluppi futuri

definiti all'interno dell'introduzione.

La parte teorica relativa a tutti gli argomenti precedentemente illustrati verranno trattati all'interno del capitolo 3 *Progetto logico della soluzione del problema*. Nel dettaglio la parte implementativa dei due algoritmi, della raccolta dati, della sentiment analysis e della predizione verrà trattata all'interno del capitolo 4 *Implementazione e realizzazione del sistema*. Gli esperimenti e le motivazioni che hanno mosso alla definizione dei due topic scelti per l'analisi della polarizzazione verranno trattate all'interno del capitolo 5 *Test sperimentali e valutazione*. Infine gli sviluppi futuri e le conclusioni verranno trattati all'interno del capitolo 6 *Sviluppi futuri e conclusioni*.

Capitolo 2

Stato dell'arte e Background

2.1 Stato dell'arte

All'interno delle reti sociali sta sempre più prendendo piede il problema della polarizzazione delle opinioni. Nel linguaggio comune il confronto tra individui ha sempre generato una forte controversia oppure un punto d'incontro. I social network hanno permesso all'utente di poter diffondere attraverso post, messaggi o espressioni audio video le proprie opinioni e/o pensieri all'interno di una comunità sociale. A tal proposito per favorire la diffusione delle diverse correnti di pensiero i social network stanno sempre più sviluppando algoritmi per permettere di identificare le comunità isolate, che condividono un unico punto di vista di un determinato argomento. La polarizzazione è un algoritmo matematico che applicato all'interno delle reti sociali consente di capire quanto un utente, che accede per la prima volta all'interno di una rete sociale, venga influenzato dagli altri utenti e quanto una news o un giudizio si propaga all'interno di una rete sociale. Prima di poter illustrare questo algoritmo, con le relative problematiche, verrà illustrata una definizione di rete sociale.

2.1.1 Rete Sociale

Una rete sociale consiste in un qualsiasi gruppo di individui connessi tra loro da diversi legami. Per gli esseri umani i legami vanno dalla conoscenza casuale, ai rapporti di lavoro, ai vincoli familiari. Le reti sociali sono spesso usate come base di studi interculturali in sociologia, in antropologia, in etologia.

L'analisi delle reti sociali, ovvero la mappatura e la misurazione delle reti sociali, può essere condotta con un formalismo matematico usando la teoria dei grafi. In generale, il corpus teorico ed i modelli usati per lo studio delle reti sociali sono compresi nella cosiddetta *social network analysis*.

La ricerca condotta nell'ambito di diversi approcci disciplinari ha evidenziato come le reti sociali operino a più livelli e svolgano un ruolo cruciale nel determinare le modalità di risoluzione di problemi e i sistemi di gestione delle organizzazioni, nonché le possibilità dei singoli individui di raggiungere i propri obiettivi.

2.1.2 Le reti

La diffusione del web e del termine social network ha creato negli ultimi anni alcune ambiguità di significato. La rete sociale è infatti storicamente, in primo luogo, una rete fisica.

Rete sociale è, ad esempio, una comunità di lavoratori, che si incontra nei relativi circoli dopolavoristici e che costituisce una delle associazioni di promozione sociale. Esempi di reti sociali sono inoltre le comunità di sportivi, attivi o sostenitori di eventi, le comunità unite da problematiche strettamente lavorative e di tutela sindacale del diritto nel lavoro, le confraternite e in generale le comunità basate sulla pratica comune di una religione e il ritrovo in chiese, templi, moschee, sinagoghe e altri luoghi di culto.

La rete sociale ha una sua versione all'interno di **Internet** definita anche come *Social media*, questa è una delle forme più evolute di comunicazione in rete. La rete delle relazioni sociali che ciascuno di noi tessesse ogni giorno, in maniera più o meno casuale, nei vari ambiti della nostra vita, si può così "materializzare", organizzare in una "mappa" consultabile, e arricchire di nuovi contatti. I principali social network sono: *Facebook*, *MySpace*, *Instagram*, *Twitter*, *Google+*, *LinkedIn*, *Ask.fm*, *Pinterest*, *Formspring*, *Bebo*, *Friendster*, *Hi5*, *Ning*, *Tagged*, *Meetup*, *Tumblr*.

2.1.3 Twitter

Twitter è il social network di riferimento utilizzato per la realizzazione di questa tesi di laurea, ne diamo una breve descrizione, è una rete sociale, creata il 21 marzo 2006 dalla *Obvious Corporation* di San Francisco, che fornisce agli utenti, attraverso l'omonima piattaforma, una pagina personale aggiornabile tramite messaggi di testo con lunghezza massima di 140 caratteri; nel 2017 l'azienda ha aumentato la lunghezza dei tweet a 280 caratteri per alcuni paesi. Gli aggiornamenti di stato possono essere effettuati tramite il sito stesso, via SMS, con programmi di messaggistica istantanea, posta elettronica, oppure tramite varie applicazioni basate sulle API di Twitter.

Il nome "Twitter" deriva dal verbo inglese *to tweet* che significa "cinguiettare". *Tweet* è anche il termine tecnico degli aggiornamenti del servizio. I *Tweet* che contengono esattamente 140 caratteri vengono chiamati *Twoosh*. Gli aggiornamenti sono mostrati nella pagina di profilo dell'utente e comunicati agli utenti che si sono registrati per riceverli. È anche possibile limitare la visibilità dei propri messaggi oppure renderli visibili a chiunque.

Questo social network è molto utilizzato da personaggi famosi, sportivi, ma anche capi di stato, esponenti politici ed economici, perfino il Papa,



Figura 2.1: Logo di Twitter

consentendo a tutti gli utenti iscritti alla piattaforma di poter visualizzare le idee ed opinioni espresse da questi personaggi. Proprio per questo motivo *Twitter* consente una diffusione su larga scala di idee ed opinioni da parte di moltissimi utenti, gli stessi giornalisti utilizzano questa piattaforma per diffondere *scoop* o notizie di vario genere, proprio per la grande visibilità che fornisce questo social network. Vengono generati moltissimi dati ed è proprio per questo motivo che ho deciso di utilizzarlo come strumento per la raccolta dati e cercare di calcolare la polarizzazione di alcuni topic selezionati, permettendomi di visualizzare come le informazioni espresse si propagano all'interno della rete e come queste influenzano il pensiero degli utenti.

Tweet

L'elemento alla base di questo social network è il *Tweet*, che non è altro che un messaggio contenente le opinioni, le idee, i pensieri dell'utente. Tutte queste cose possono essere arricchite utilizzando elementi multimediali come video, immagini oppure link a pagine web. Il tweet una volta pubblicato potrà essere accessibile per tutti gli utenti delle rete che hanno una relazione di "amicizia" con l'utente in questione, potendo anche effettuare delle operazioni aggiuntive quale commentare la notizia, apprezzarla oppure *retwettarla* cioè condividere queste informazioni a tutti i propri collegamenti. I Tweet possono essere anche etichettati utilizzando due strumenti messi a disposizione da questo

social network:

- *Mentions*: sono dei riferimento ad utenti, cioè come assegnare un mittente ad un email, affermando a tutti i membri all'interno della rete, che quel tweet è destinato esclusivamente a quella persona. Può essere utile per condividere più informazioni con gli amici delle rete più stretti, oppure per formulare degli attacchi, riflessioni o considerazioni anche con utenti famosi. Per fare questa operazione è necessario inserire la `@` ed il nome dell'utente in questione per poter collegare al tweet in questione alla persona desiderata.
- *Hashtag*: parole o combinazioni di parole concatenate precedute dal simbolo cancelletto (`#`). Etichettando un messaggio con un hashtag si crea un collegamento ipertestuale a tutti i messaggi recenti che citano lo stesso hashtag. Molto vengono utilizzati durante *Show Tv*, *comizi*, *pubblicità*, *eventi politici* per poter consentire alle persone che vogliano esprimere la loro opinioni sugli argomenti espressi durante questi eventi.

2.1.4 Polarizzazione

Definiamo in maniera concettuale il significato, le funzionalità della polarizzazione. Citando il paper *A Measure of Polarization on Social Media Networks Based on Community Boundaries* di *Pedro H. Calais Guerra, Wagner Meira Jr. Claire Cardie, Robert Kleinberg* ([17]) la polarizzazione viene definita come un processo sociale in cui un gruppo viene diviso: in due sottogruppi aventi un conflitto o una visione differente del problema in questione, e da alcuni individui che rimangono neutrali. Comprendere la polarizzazione e quantificarla è una sfida a lungo termine per i ricercatori all'interno di

diverse aree, proprio perché risulta essere un potente strumento per l'analisi delle opinioni.

Nel paper *Measuring Political Polarization: Twitter shows the two sides of Venezuela* di A. J. Morales, J. Borondo, J. C. Losada e R. M. Benito ([13]), la polarizzazione assume una connotazione molto simile, vista in chiave politica, infatti all'interno di questo paper tale analisi viene effettuata all'interno di un contesto politico. La definizione che viene data della polarizzazione è: un fenomeno sociale che si verifica quando gli individui contrappongono le proprie credenze ed opinioni in una posizione conflittuale tra loro, mentre alcuni mantengono una posizione neutrale. Per citare *Jonh Turner*¹ : "Come le molecole polarizzate, i membri di un gruppo allineano il proprio pensiero nella direzione in cui erano originalmente diretti".

Adottare questo strumento all'interno di un contesto politico fornisce una chiara visione della preferenza o meno verso i candidati politici, verificando il trend del candidato. Dal punto di vista matematico possiamo definire la polarizzazione tra due gruppi come la distanza che intercorre tra due numeri. Definendo in questo modo quanto due idee siano diverse e contrastanti utilizzando la distanza tra due valori numerici di riferimento. Le problematiche evidenziate dalla polarizzazione vengono trattate all'interno del paper : *Reducing Controversy by Connecting Opposing Views* di Kiran Garimella ,Gianmarco De Francisci Morales, Aristides Gionis e Michael Mathioudakis [15]. Quando una popolazione si divide in due gruppi con visioni opposte quello che spesso si può verificare è la creazione degli *Echo-chambers* cioè: una situazione dove persone che la pensano allo stesso modo rafforzano le proprie convinzioni a vicenda, senza esporsi verso idee contrastanti.

¹Noto psicologo sociale britannico, con altri colleghi ha sviluppato la teoria dell'Auto-categorizzazione, che afferma tra l'altro che il sé non è un aspetto fondamentale della cognizione, ma è il risultato di processi cognitivi ed interazioni tra la persona e il contesto sociale.

In conclusione la polarizzazione può aiutare ad individuare e a prevenire la formazioni degli *Echo-chambers*.

Echo-Chambers

Brevemente illustriamo un gravissimo problema che abbiamo precedentemente riscontrato all'interno della polarizzazione. Gli *Echo-chambers* è una descrizione metaforica di una situazione in cui le informazioni, le idee, o credenze sono rinforzate dalla comunicazione e dalla ripetizione all'interno di un ambiente definito. Praticamente viene diffuso lo stesso messaggio, come un "eco" che riproduce lo stesso suono più e più volte all'interno di un ambiente[4]. Nei social media il fenomeno degli echo-chambers è un problema piuttosto comune e ciò è dovuto dalla forte presenza di comunità di utenti che condividono sempre le stesse informazioni con un forte tasso di pubblicazione, quindi un nuovo utente che appoggia una di quelle idee entra all'interno di quel circolo vizioso. Ovviamente gli amministratori di questi social network sviluppano algoritmi che possano prevenire la formazione di queste comunità. Le conseguenze di questo problema sono molto evidenti ovvero:

- *Impossibilità della crescita culturale*
- *Impossibilità della crescita sociale*
- *Rendere la comunicazione poco democratica*

La polarizzazione come precedentemente accennato è uno strumento che può aiutare la prevenzione per la formazione di queste comunità.

2.1.5 Sentiment Analysis

La Sentiment Analysis è lo strumento utilizzato per poter catalogare le informazioni degli utenti attraverso un'analisi del testo. Non è possibile citare la sentiment analysis senza menzionare l' "opinion mining" ovvero una branca del machine learning che studia il comportamento delle opinioni nella rete. Il termine sentiment viene utilizzato in riferimento all'analisi automatica effettuata per valutare il testo ed esprimere un giudizio predittivo sul contenuto [14]. Il compito base di questo processo è quello di classificare la polarità di un testo e quindi il suo contenuto e catalogarlo come *positivo*, *negativo* o *neutrale*. Gli approcci più comunemente individuati rispetto alla sentiment analysis possono suddividersi in tre macro-categorie:

- *rilevamento delle keyword*: questo metodo consente di classificare il testo tramite categorie emotive facilmente riconoscibili, individuate in base alla presenza di parole emotive non ambigue, come felice, triste, e annoiato.
- *affinità lessicale*: questo metodo non rileva solo le keyword emotive, ma assegna anche a parole arbitrarie *un'affinità probabile* a emozioni particolari. Rispetto alla prima metodologia vista, l'affinità lessicale consente di affinare la selezione e l'attribuzione della polarità.
- *metodi statistici*: questi metodi si basano su elementi di apprendimento automatico. Per misurare l'opinione nel contesto e trovare la caratteristica che è stata giudicata, sono usate le relazioni grammaticali delle parole utilizzate. Le relazioni di dipendenza grammaticale sono ottenute attraverso la scansione approfondita del testo. Il processo di apprendimento da parte della macchina (anche detto machine learning) non è immediato, devono infatti essere costruiti dei modelli che asso-

ciano a diverse tipologie di commenti una polarità e se necessario ai fini dell'analisi anche un topic.

Quando parliamo della sentiment analysis non possiamo non esporre una delle grandi limitazioni offerte da questo strumento, cioè l'incapacità della macchina di poter comprendere concetti emotivi complessi come l'ironia. Ad esempio un commento ad un ritardo nel volo:

"Il mio volo è in ritardo. Splendido!"

verranno interpretati e classificati dalla macchina come post dalla polarità positiva mentre invece dovrebbero essere assegnate delle negatività. Quindi in un contesto lessicale pieno di ironia l'attendibilità sarà inferiore rispetto ad un documento con informazioni oggettive. Per migliorare l'accuratezza della predizione è necessario l'intervento dell'uomo, aumentando il volume dei dati di riferimento per l'analisi, banalmente milioni di post, possono alleviare le preoccupazioni sull'affidabilità a livello granulare, ossia di un singolo post. In conclusione la sentiment analysis è uno strumento che può essere molto utile all'interno dei social media in quanto può consentire di catalogare le opinioni degli utenti attraverso l'analisi dei contenuti pubblicati all'interno della rete. Questa potenzialità è ciò che mi ha spinto nell'adottare tale tecnica all'interno della mia tesi, garantendo una classificazione dei diversi *Tweet* collezionati.

2.1.6 Forecasting

Il Forecasting è un processo per fare predizioni del futuro basandosi sui dati passati e presenti mediante l'analisi delle tendenze. Questo viene utilizzato molto all'interno delle aziende per poter analizzare i guadagni ed i futuri ricavi. [5] Spesso viene associato ai processi di budgeting e pianificazione, si serve di dati passati e presenti, analisi dei trend e informazioni esecutive per

prevedere la situazione futura di ogni indicatore. Il forecasting prevede tre diverse tecniche:

- *Qualitativa*
- *Analisi e proiezione delle serie temporali*
- *Modelli causali*

La previsione finanziaria non è una scienza esatta e l'incertezza ne è un tipico aspetto. Per limitare gli errori e migliorare l'esattezza delle previsioni è necessario utilizzare precisi dati storici e in tempo reale. Inoltre, la possibilità di combinare previsioni rolling, previsioni a lungo raggio, simulazione di scenari possibili e stress test supporta i risultati delle previsioni e rende le previsioni più agili e reattive ai cambiamenti economici e di business. Ai fini della mia tesi ho deciso di adottare uno studio basato sulle analisi e proiezioni delle serie temporali, perché permette di estrarre statistiche significative ed altre caratteristiche dei dati. Predice i valori futuri attraverso l'osservazione dei dati precedentemente raccolti.

Capitolo 3

Progetto logico della soluzione del problema

In questo capitolo verrà presentato il flusso logico della tesi con la soluzione proposta per la suddivisione dei gruppi partendo dai dati raccolti, il calcolo della polarizzazione ed infine la sua predizione nel tempo.

3.1 Raccolta Dati

La prima parte del flusso logico della mia tesi si basa sulla raccolta dei dati. Questa operazione è stata effettuata utilizzando il social network **Twitter**. Nel dettaglio sono stati raccolti tutti i tweet relativi a due topic, usati per effettuare le analisi, le motivazioni della scelta di questi due argomenti verranno illustrate più avanti, cioè:

- **La elezioni regionali in Sicilia**
- **Biotestamento**

Come precedentemente illustrato la scelta di questi due topic è dovuta al fatto che sono in primis due argomenti molto recenti e di attualità all'interno del nostro paese, in secundis perché riferiti a due contesti differenti tra loro ovvero quello politico e quello sociale. Prima di effettuare la raccolta di tweet per ognuno dei due topic, è stato effettuato uno studio sugli hashtag, cioè la ricerca veniva effettuata per una serie di hashtag per cui sono stati catalogati i 5 più utilizzati dagli utenti per esprimere le loro opinioni sull'argomento in questione. Lo studio di questi hashtag è stato improntato ricercando quelli che non esprimessero un giudizio, bensì che aiutassero l'utente a connotare i loro pensieri sul topic avendo una connotazione generica. La ricerca è stata fatta in maniera del tutto equilibrata soprattutto per quanto concerne le elezioni regionali in Sicilia in quanto è facile cadere in preda di hashtag utilizzati dalle fazioni politiche per attirare gli elettori, ne sono un esempio:

- *#diventeràbellissima*: utilizzato dal centro destra come motto all'interno dei social media per pubblicizzare il proprio piano politico.
- *#impresentabili*: utilizzato dal Movimento 5 Stelle per denunciare i candidati degli altri partiti politici.

Per evitare quindi di raccogliere dati già fortemente polarizzati, si è deciso di adottare una strategia più neutrale cercando 5 hashtags generici che rendessero l'idea del topic in questione. Twitter ha una politica di protezione per i dati, che sono accessibile a qualsiasi utente che abbia effettuato l'abilitazione allo sviluppo attraverso le Api messe a disposizioni, impedendo di effettuare più di 100 richieste ogni 15 minuti al server, impedendo un uso improprio e maligno con i dati pubblicati dagli utenti. Per richieste basta considerare la raccolta di un singolo Tweet. Per ottimizzare i tempi di raccolta si è deciso di utilizzare un'istanza *EC2*, che eseguisse uno script *Python* per la raccolta dei

dati in questione, rimanendo attivo anche durante le ore notturne. I dati in questione venivano salvati all'interno di file binari, in modo da ottimizzare lo spazio che avrebbero occupato sulla macchina. Il motivo che mi ha spinto ad effettuare tale operazione è dettata dai costi che ha l'istanza EC2 per poter mantenere i dati fisici al suo interno, perché i consumi economici non sono generati soltanto dall'utilizzo delle risorse fisiche della stessa, ma anche dalla quantità di dati presente al suo interno.



Figura 3.1: Servizi AWS



Figura 3.2: Python

3.2 Sentiment Analysis

La sentiment Analysis è stata utilizzata per catalogare i *Tweet* delle persone in riferimento ai due topic d'interesse. Utilizzando questo strumento l'operazione effettuata è stata quella di suddividere in 3 categorie il contenuto pubblicato nella rete. In precedenza sono state illustrate le numerose tecniche per poter utilizzare la sentiment analysis, per questo motivo la tecnica che più si avvicinasse alle mie esigenze è quella basata su metodi statistici, cioè vengono applicati principi del machine learning per poter identificare la vicinanza o meno del testo con il *training set* di riferimento.[8] Come descritto in precedenza maggiore è la quantità dei dati che contiene il training set e maggiore ne risulterà la precisione nella predizione del contenuto desiderato. All'interno del flusso di esecuzione della mia tesi la sentiment analysis viene collocata all'interno della raccolta dati, cioè nel momento in cui viene

collezionato un tweet questo viene immediatamente analizzato e collocato nel gruppo di riferimento. Successivamente se il tweet in questione presenta dei *retweet*, ovvero una condivisione del contenuto con altri utenti, questi saranno a loro volta analizzati attraverso la sentiment analysis perché c'è la possibilità di aggiungere commenti al retweet.

3.3 Endorsement Graph

L'Endorsement Graph è il grafo di riferimento su cui verranno effettuate le successive operazioni. Concettualmente è un grafo diretto basato sui tweet e sui retweet fatti dagli utenti della rete. Gli elementi costituenti del grafo sono:

- *Tweet*: sono i nodi che formano il grafo.
- *Retweet*: sono assimilabili ad archi e noti, cioè essendo delle estensioni alle pubblicazioni fatte dai nodi, allora all'interno del grafo verranno generati nuovi nodi. Questi sono gli utenti che hanno pubblicato tali informazioni, ed i nuovi archi che collegano i tweet con i retweet rappresentano il collegamento tra i due nodi. Una menzione particolare su questi archi che sono orientati dal Retweet verso il Tweet.

L'endorsement graph viene popolato da tutti i dati raccolti attraverso la ricerca per hashtag, rendendo in questo modo il grafo molto popolato. Per dare una parvenza grafica della polarizzazione i nodi appartenenti ai due gruppi distinti sono stati colorati in due colori di riferimento:

- **Rossi**
- **Blu**

3.4 Polarizzazione

La polarizzazione viene calcolata subito dopo la creazione del grafo in questione. La realizzazione di questo strumento è stata fatta implementando gli algoritmi proposti all'interno di questi due paper:

- *Measuring Political Polarization: Twitter shows the two sides of Venezuela*
- *Reducing Controversy by Connecting Opposing Views*

Per facilitare la dichiarazione all'interno dell'elaborato chiameremo il primo algoritmo come *polarizzazione basata sul grado*, mentre il secondo lo chiameremo *polarizzazione basata sulla topologia*. La polarizzazione basata sul grado analizza il grafo effettuando due operazioni in più passi:

1. Separare tutti i nodi in due categorie:
 - **Elite**: sono quei nodi che pubblicano notizie, in questo caso coloro che postano dei Tweet sul topic esprimendo una loro opinione.
 - **Listener**: sono quei nodi che condividono le notizie altrui attraverso una operazione di retweet.
2. Calcolare la polarizzazione utilizzando il grado in ingresso del nodo, incentrando il calcolo sui nodi *Elite*. Per grado del nodo si intende il numero di archi in entrata verso il nodo di riferimento.
3. Ripetere il passo precedente fino a quando non si stabilizzano i valori della polarizzazione

Per quanto riguarda la polarizzazione basata sulla topologia, come suggerisce il nome scelto, dipende molto dalla forma che assume l'endorsement

graph durante la sua creazione, è soggetta quindi a variazioni dettate dai collegamenti che si vengono a creare tra i diversi nodi presenti. I passi effettuati per la realizzazione di questo algoritmo sono i seguenti:

1. Per ogni nodo verificare se è possibile raggiungere i nodi di grado massimo, appartenenti ai due gruppi distinti che chiameremo *Rossi* e *Blue*.
2. Effettuare dei *Random Walk* ¹, utilizzando i pesi sugli archi per calcolare la polarizzazione, per raggiungere quei nodi di grado massimo appartenenti alle due categorie sopra citate.

Queste sono le operazioni effettuate dall'algoritmo basato sulla topologia. Una volta applicati i due algoritmi i risultati vengono salvati per poi essere rappresentati graficamente, mostrando a livello visivo il grafo con le colorazioni in conformità con i risultati raggiunti attraverso la polarizzazione. Questa parte è il focus della mia tesi in quanto utilizzando questi algoritmi possiamo visualizzare con i nostri occhi quanto queste informazioni siano polarizzate, quanto le persone (i nodi sono utenti della rete) hanno dibattuto sul quel topic. I risultati ed i dettagli implementativi verranno mostrati in seguito, come anticipazione posso affermare che i topic hanno mostrato attraverso la polarizzazione lo stesso trend mostrato nel mondo reale.

3.5 Predizione

L'ultima fase della tesi consiste nella predizione della polarizzazione in un periodo successivo rispetto a quello considerato durante la fase di raccolta dati.

¹Si definisce Random Walk la formalizzazione dell'idea di prendere passi successivi in direzioni casuali. Matematicamente parlando, è il processo stocastico più semplice, il processo markoviano, la cui rappresentazione matematica più nota è costituita dal processo di Wiener.

Questa esigenza nasce per poter consentire una prevenzione sul problema della formazione degli *Echo-chambers*, questi molto spesso vengono a formarsi con argomenti fortemente polarizzati, quindi poter prevenire per tempo la formazione di queste comunità isolate può comportare un enorme vantaggio per gli amministratori dei social media. Per la realizzazione di questa parte della tesi si è deciso di adottare una predizione basata sull'analisi dei dati raccolti, nel dettaglio si è deciso di utilizzare algoritmi di *Forecasting* basati sull'analisi delle serie temporali.[5] Nel dettaglio ho effettuato un'analisi basata su 3 tecniche differenti tra loro:

- *Double exponential smoothing*: è una tecnica che sfrutta la serie temporale dei dati raccolti, assegnando una crescita esponenziale sui dati nel tempo, tenendo conto anche del trend di crescita o di decrescita nel tempo. In questo modo si cerca di predire il contenuto nel primo istante successivo.
- *Linear regression*: questa tecnica formalizza e risolve il problema di una relazione funzionale tra variabili misurate sulla base di dati campionari estratti da un'ipotetica popolazione infinita. In statistica l'analisi della regressione è associata alla risoluzione del modello lineare.
- *Moving average*: è una tecnica che estende la media aritmetica dei risultati, andando a mediare i risultati all'interno di una finestra mobile, prelevando gli ultimi risultati ottenuti in precedenza mediandoli e definendo come valore futuro questo risultato.

I risultati della predizione verranno trattati in seguito così come i dettagli implementativi. La predizione all'interno di questo contesto può risultare utilissimo per l'analisi futura dei topic nel tempo, perché ovviamente la for-

mazione degli *Echo-chambers* può essere arginata facendo dei controlli e/o previsioni sul valore della polarizzazione nel futuro.

Capitolo 4

Implementazione e realizzazione del sistema

All'interno di questo capitolo verranno illustrate tutte le tecniche implementative per realizzare quanto è stato enunciato nel capitolo precedente. Nel dettaglio verranno spigate ed argomentate le implementazioni effettuate su tutte le fasi della tesi, i riferimenti teorici che hanno permesso la realizzazione di questa tesi. Di seguito una breve descrizione grafica del flusso con le operazioni effettuate durante la redazione del lavoro di tesi.(vedi fig:4.1)

4.1 Raccolta dati

La raccolta dati è la fase iniziale della tesi di laurea. Come precedentemente argomentato il social network di riferimento utilizzato è *Twitter*.Prima di spiegare quanto fatto occorre citare le funzionalità della *Twitter Api*. Successivamente verranno illustrate le tecniche utilizzate per la raccolta dei *Tweet* del passato.

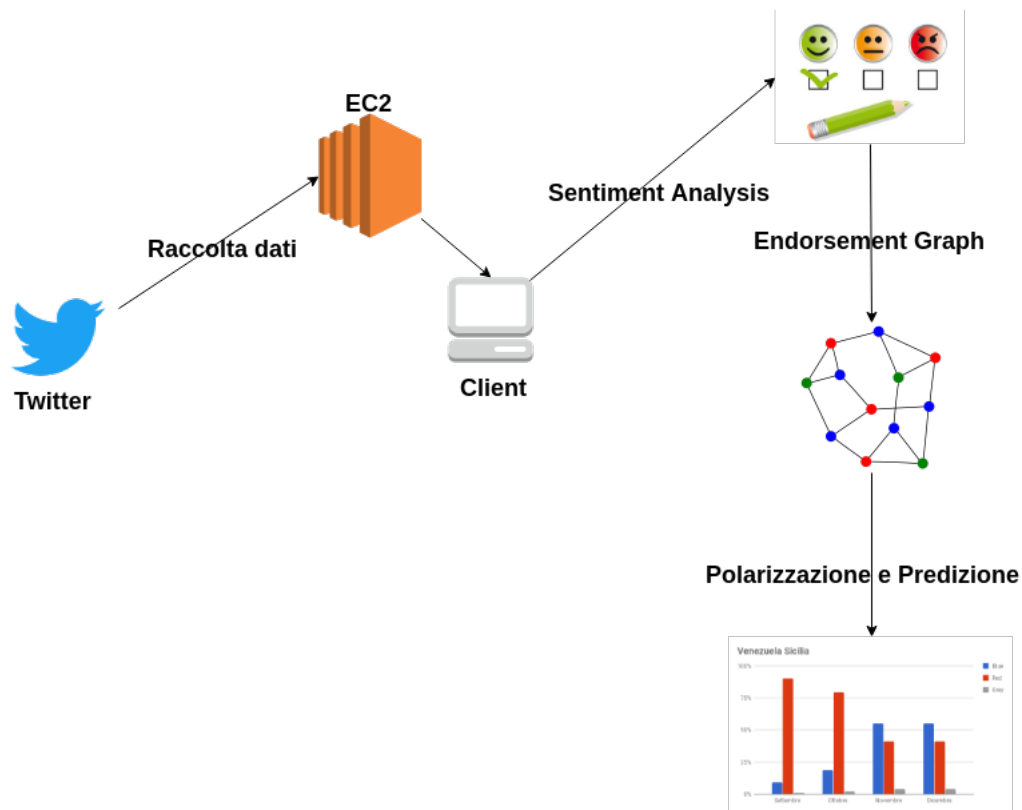


Figura 4.1: Flusso

4.1.1 Twitter Api

Sono delle api messe a disposizione dal social network Twitter, per poterle utilizzare occorre necessariamente effettuare l'iscrizione al reparto sviluppatori di Twitter, mettendo a disposizione il proprio account personale, quindi come prima cosa occorre predisporre di un proprio account. La raccolta dei dati può essere effettuata soltanto dopo aver abilitato l'account. Per raccogliere i dati pubblicati dagli utenti all'interno della rete occorre creare una *Twitter Apps*, la quale rilascerà delle credenziali, che dovranno essere utilizzate per poter effettuare le richieste al server attraverso le Twitter Api. Viene inizializzato un processo di streaming tra il client ed il server che consenta la

ricezione dei dati, diamo una definizione delle credenziali, che sono divise in:

- *Token*: identificano il token di accesso ai servizi messi a disposizione da Twitter, a sua volta è composto da:
 - *Access token*
 - *Access secret*
- *Consumer*: è un codice di accesso ai servizi di consumo, ovvero consente lo streaming dei dati dal server, si suddivide in due chiavi:
 - *Consumer key*
 - *Consumer secret*

Queste chiavi di accesso possono essere generate più e più volte, una volta creata un'applicazione per lo sviluppo di Twitter. Le Api sono disponibili in diversi linguaggi di programmazione tra cui il *Python*, utilizzato per lo sviluppo della tesi. Per richiamare tutte le api attraverso il codice è necessario sempre autenticarsi, caricando le credenziali di accesso fornite attraverso i permessi da sviluppatore precedentemente illustrate. Lo streaming per la raccolta dei dati è soggetto a stringenti regole per evitare un uso improprio delle informazioni diffuse dagli utenti all'interno della rete. Infatti c'è un numero massimo di richieste che un utente, con le credenziali da sviluppatore, può effettuare, cioè 100 ogni 15 minuti. Allo scadere di tale tempo potrà ricominciare, altrimenti nel caso in cui un utente non tenesse conto di questa limitazione ed effettuasse una nuova richiesta durante il tempo di pausa le sue credenziali verrebbero bloccate per circa un'ora non potendo più interagire con le api di Twitter. La risposta da parte del server sarà un file *JSON* contenente tutti i metadati dell'utente in questione, come il testo del Tweet, il suo id, lo username dell'utente che ha pubblicato il tweet, i mentions gli

hashtag ed il numero di retweet con la lista degli utenti che hanno retwettato la notizia in questione.

4.1.2 Get Old Tweet

Questa sezione illustrerà in che modo sono stati raccolti tutti i tweet del passato. Avendo spiegato nella sezione precedente le problematiche relative al numero di richieste da effettuare nell'arco temporale, è stata seguito un approccio differente che consentiva di limitare il numero di richieste alla volta. L'approccio in questione consiste nell'effettuare una *GET* sulla pagina di ricerca di Twitter, specificando l'argomento ricercato, ed il periodo di validità della ricerca. Twitter suddivide gli account abilitati allo sviluppo in 3 categorie:

- *Standard*: account gratuito soggetto a limitazione temporali e sui contenuti, non è possibile ricercare tweet precedenti a 7 giorni. I dati richiesti non sono completi.
- *Premium*: account a pagamento con sole limitazioni temporali, non è possibile ricercare tweet precedenti a 30 giorni. Non ha nessun problema di completezza sui contenuti.
- *Enterprise*: account a pagamento senza alcuna limitazione temporale e sui contenuti.

Avendo utilizzato un account *Standard* sono stati ricercati dei metodi per poter risolvere le problematiche precedentemente descritte, cercando di ottenere più informazioni possibili.

Tutte queste problematiche sono state risolte la libreria *Get Old Tweet*, il quale crea un indirizzo http con i parametri di ricerca richiesti, nel nostro caso:

- *La query*: la parola chiave, l'hashtag da ricercare all'interno dei Tweet.
- *Data inizio*: la data in cui inizio a raccogliere i dati.
- *Data Fine*: la data in termino la ricerca e la raccolta dei dati.

Una volta definiti questi parametri all'interno della *url*, viene restituita la pagina html contenente tutti i tweet presenti nel periodo indicato. Il risultato verrà convertito in un formato *json*, per estrarre le informazioni basterà parsare la pagina ottenendo i seguenti dati:

- Lo *username* della persona che ha postato il tweet.
- Il numero di *retweet* al tweet.
- Il *testo* del tweet pubblicato.
- La lista degli *hashtag* pubblicati dall'utente all'interno del tweet
- La lista dei *mentions* pubblicati dall'utente all'interno del tweet.
- La *data* di pubblicazione del tweet.

Per far comprendere meglio queste informazioni nell'immagine sovrastante viene presentato un tweet di esempio con le informazioni precedentemente elencate, per far comprendere al meglio gli elementi in questione.

Tutti questi dati sono stati successivamente elaborati ed analizzati. Particolare attenzione è stata posta verso due di essi ovvero: il testo ed il numero di retweet. Il primo per effettuare la sentiment analysis e poter definire una prima partizione dei dati, che andranno a popolare il grafo. Il secondo perché costituisce la base per i nuovi nodi del grafo e quindi la diffusione del tweet con altri utenti. Per identificare la lista degli utenti che hanno effettuato il retweet sul tweet in questione è stato necessario usare una libreria chiamata



Figura 4.2: esempio Tweet

Tweepy. Questa libreria sfrutta le *Twitter Api*, effettuando delle chiamate rest sul server di Twitter per ricevere le informazioni richieste. Nel nostro caso è stato utilizzato il metodo *retweet* che ricevendo come argomento l'id relativo al tweet pubblicato permette di ricevere la lista degli username che hanno pubblicato quell'argomento all'interno della propria rete. Il problemi presentati in precedente nell'utilizzare queste chiamate al server sono sempre validi, infatti se vengono superate le 100 richieste viene lanciata un'eccezione: *Tweepy Error* che consente di mettere in pausa l'applicazione attraverso una sleep per 15 minuti. Una volta terminato il tempo di attesa la richiesta viene ripresa dal punto richiesto, continuando a collezionare gli username richiesti. Tutti i dati raccolti durante le operazioni di streaming sono stati salvati all'interno di un file binario per poter ottimizzare lo spazio fisico.

EC2 La raccolta dati utilizzando le implementazioni precedentemente citate, è stata eseguita all'interno di una istanza **EC2**. Il motivo di tale scelta è dettata dai tempi di attesa via via sempre più lunghi per la raccolta degli username degli utenti che hanno retweettato i tweet raccolti, per via delle

politiche stringenti dettate da **Twitter**. Utilizzando un'istanza a pagamento, si è dovuto ottimizzare il salvataggio fisico dei dati, perché oltre ad essere un servizio a consumo di risorse di calcolo, vengono pagate anche i dati salvati all'interno dello storage fisico della macchina. Attraverso la libreria *pickle*¹ si è potuto risolvere questo problema, in quanto ottimizza il salvataggio dei dati attraverso una conversione in binario degli stessi.

4.2 Sentiment Analysis

La sentiment analysis è l'approccio utilizzato per la divisione dei diversi tweet, raccolti per un determinato topic, all'interno di una determinata finestra temporale, analizzando il "sentimento" espresso nel contenuto del testo pubblicato dagli utenti. Questa fase viene effettuata al momento della raccolta dei dati, cioè una volta collezionati l'insieme dei dati vengono sottoposti alla funzione implementata e salvati in due categorie distinte mediante la libreria *pickle*, precedentemente illustrata.

L'implementazione di questo metodo è stata effettuata seguendo un approccio basato su *metodi statistici*: questi metodi si basano su elementi di apprendimento automatico. Per misurare l'opinione nel contesto e trovare la caratteristica che è stata giudicata, sono usate le relazioni grammaticali delle parole utilizzate. Le relazioni di dipendenza grammaticale sono ottenute attraverso la scansione approfondita del testo. Il processo di apprendimento da parte della macchina (anche detto machine learning) non è immediato, devono infatti essere costruiti dei modelli che associano a diverse tipologie di commenti una polarità e se necessario ai fini dell'analisi anche un topic.[[11]

Si può riassumere quanto sviluppato attraverso la seguente immagine:

¹libreria per il salvataggio dei dati fornita da *Python*

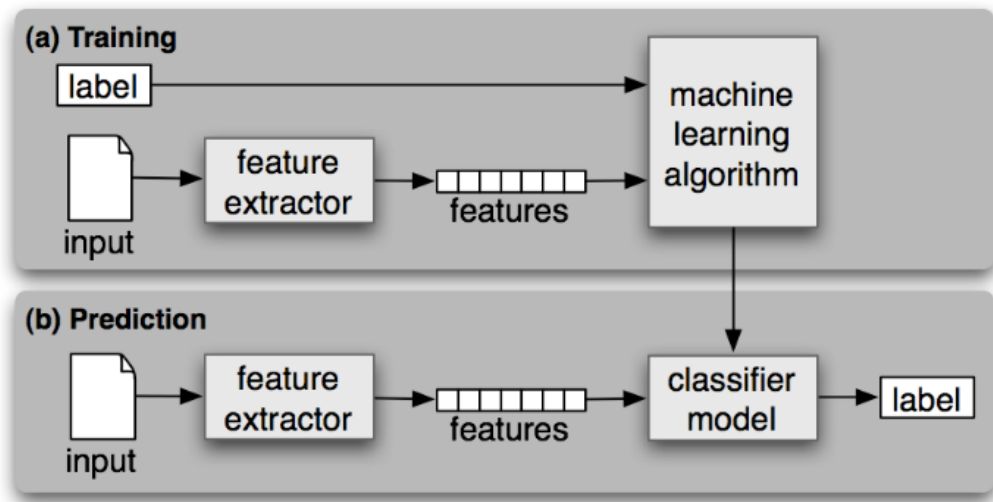


Figura 4.3: schema Sentiment Analysis

Da quanto si evince dall'immagine mostrata risulta evidente quanto un meccanismo di *machine learning* necessiti di un classificatore molto preciso. La precisione può essere ottenuta attraverso dei *training set*, ovvero dei file che istruiscono la macchina permettendo una classificazione ottimale dei dati. Ci sono diversi classificatori, quello utilizzato per l'implementazione della Sentiment Analysis è descritto di seguito.

- **classificatore Naive Bayes**[16]: Un classificatore bayesiano è un classificatore basato sull'applicazione del teorema di Bayes. Richiede la conoscenza delle probabilità a priori e condizionali relative al problema, quantità che in generale non sono note ma sono tipicamente stimabili. Se è possibile ottenere delle stime affidabili delle probabilità coinvolte nel teorema, il classificatore bayesiano risulta generalmente affidabile e potenzialmente compatto. Per costruzione, il classificatore bayesiano minimizza il rischio di classificazione.

Nel gergo della classificazione di testi o *Text Categorization*, con il ter-

mine classificatore bayesiano ci si riferisce convenzionalmente al classificatore Naive Bayes, ossia un classificatore bayesiano semplificato con un modello di probabilità sottostante, che assume l'indipendenza delle *feature*, ovvero che la presenza o l'assenza di una particolare *feature* in un documento testuale non è correlata alle altre.

L'esperienza dimostra che il metodo funziona in molti problemi pratici, come per esempio il filtraggio antispam adattivo. Un vantaggio del classificatore Naive Bayes è che richiede solo un training set di esigue dimensioni per stimare i parametri necessari per la classificazione.

Tornando allo sviluppo, la prima operazione fatta è stata quella di raccogliere dati per classificare i tweet secondo un determinato topic. Questa operazione è stata effettuata manualmente anche per permettere alla macchina di poter riconoscere l'ironia, cosa che è possibile solo per la mente umana. Per questo motivo la precisione dei classificatori non sarà mai assoluta, conseguentemente le macchine necessitano dell'intervento dell'uomo per risolvere tali problematiche. Il training set in questione è stato redatto attraverso un file csv (vedi tabella: 4.1) costruito con due colonne:

- *label*: ovvero l'identificatore che verrà utilizzato nell'analisi del classificatore per analizzare il contenuto di tweet ed associarlo ad un gruppo piuttosto che ad un altro.
- *Tweet*: il testo del tweet pubblicato, contenente anche caratteri speciali, elementi multimediali, link etc.

Una volta definito un training set occorre definire un *vettore di Feature*:

- Rappresenta l'elemento cardine di un classificatore, maggiore è la sua efficienza e maggiore sarà la sua precisione. Il vettore di feature viene

Training set	
Label	Tweet
5stelle	Chi ama la sua terra non può che votare il #M5S #Regionali #Sicilia pic.twitter.com/gv3CQWhCeI
ForzaItalia	Con elezioni @matteosalvinimi e @Musumeci_Staff regionali 2017 in #Sicilia !!#elezioniregionali2017 #andiamoagovernare #forzalega
Altri	Tutti gli schieramenti alle #elezioniregionali in Sicilia.L'articolo di Pierangelo Bonanno

Tabella 4.1: Esempio di training set

utilizzato per costruire un modello che consenta, attraverso un training set, di poter fare una predizione sui dati che la macchina non ha mai analizzato prima. Dovendo analizzare dei dati basati su Twitter si è deciso di adottare degli schemi basati sull'assenza o meno di alcune parole che appaiono nei tweet come feature. Attraverso il training set, composto da tweet suddivisi in 3 gruppi (positivi, negativi e neutrali)², suddividiamo ciascun tweet in parole e aggiungeremo ogni parola la vettore di feature. Tale approccio viene definito *unigrams*. Prima di andare ad inserire tali parole all'interno di questo vettore, andremo ad effettuare delle operazioni di filtraggio scartando quelle parole che non sono necessarie per comprendere il sentimento. Nel nostro caso le congiunzioni, articoli, preposizioni semplici ed articolate, ma anche altri caratteri perché generalmente i tweet contengono caratteri speciali come # , @, *link*, vengono inseriti all'interno della lista di parole di *Stop*, cioè una lista popolata da tutte quelle parole che potrebbero far saturare il vettore di feature e che non esprimono un sentimento. Il vettore di feature viene popolato a partire da tutte le parole che non sono presenti all'interno della lista di *Stop* e che appartengono a tutti

²Questa classificazione può essere modificata in base alle esigenze dell'utente, modificando i label nel file csv del *training set*

i tweet del training set.

In conclusione una volta definito il classificatore, il vettore di Feature ed il training set non resta che illustrare il calcolo della sentiment Analysis. Per l'implementazione del classificatore *Naive Bayes* si è utilizzata la libreria *Python NLTK*³, la quale una volta istruito il classificatore attraverso il training set, precedentemente popolato in base alle proprie esigenze (come nell'esempio 4.1) utilizza il vettore di feature per ricercare la vicinanza del tweet al training set, restituendo il label corrispettivo (inserito dall'utente nella definizione del csv sopra citato). In questo modo sarà possibile sottoporre ogni testo dei diversi tweet raccolti al sistema il quale restituirà il responso in breve tempo. Per garantire una maggiore accuratezza nelle predizioni sarà necessario arricchire il training set con molteplici tweet. All'interno della sezione dedicata agli esperimenti verranno illustrati alcuni esempi dei label e dei tweet utilizzati per il calcolo dell'analisi illustrata all'interno di questa sezione.

4.3 Endorsement Graph

In questa sezione verranno illustrate le tecniche e gli strumenti utilizzati per la realizzazione dell'*endorsement graph*, cioè un grafo diretto che consente la diffusione e la pubblicazione delle notizie. Un grafo è un insieme di elementi chiamati nodi connessi tra di loro attraverso degli archi, essendo il nostro un grafo diretto significa che gli archi in questione hanno dei versi di orientamento da un nodo ad un altro. Definiamo gli elementi costituenti:

- *Nodo*: per ogni nodo è stato associato un utente che abbia pubblicato una notizia sul topic analizzato.

³Natural Language Toolkit.

- *Arco*: per ogni arco è stata associata una relazione tra il tweet pubblicato da un nodo ed il retweet effettuato da un altro utente. Questi archi sono possono essere pesati, cioè avere una probabilità che una qualsiasi coppia di nodi sia adiacente. Nel nostro caso tale valore viene definito dalla *probabilità di retweet*.

La creazione di questo grafo diretto dipende strettamente dai risultati ottenuti durante la raccolta dati. Nel dettaglio sono stati collezionati tutti i tweet che contenessero i 5 *hashtag* più utilizzati per l'espressione del topic scelto, una volta raccolta tutti i dati in questione questi sono stati fusi ed utilizzati per la definizione del grafo. I dati sono stati salvati in tre gruppi differenti grazie ai risultati ottenuti durante la *Sentiment Analysis*, questa classificazione risulta necessaria per la colorazione ed il calcolo della polarizzazione sul grafo. Tornando alla creazione è stata applicata una limitazione e cioè che i nodi isolati non venissero considerati all'interno del grafo. La motivazione che ha spinto nell'adottare tale politica è dettata dall'esigenza di ricercare quanto un'idea espressa attraverso un tweet venga diffusa in un grafo e quanto questa sia polarizzata, quindi un nodo isolato che non viene ripubblicato dagli utenti risulta inutile ai fini dell'analisi in questione. In precedenza nella definizione di Arco si è definita la probabilità di retweet:

Si definisce tale il rapporto che intercorre tra il numero dei retweet che l'utente ha effettuato su un nodo⁴ con il numero di retweet totali effettuati dall'utente stesso.

$$P(\text{retweet})_{ij} = \frac{\#retweet_{ij}}{\#retweet_j} \quad (4.1)$$

con i,j = nodi adiacenti appartenenti al grafo.

⁴Per nodo ci si riferisce ad un tweet pubblicato da un utente

Avendo raccolto dati da *hashtag* differenti è possibile che alcuni utenti potessero retwettare lo stesso utente che avesse pubblicato nuovi tweet, per mantenere questa relazione si è incrementato il numero totale di retweet effettuati complessivamente, così come il numero di volte che l'utente ha retwettato le opinioni di quello specifico utente.

Prima di popolare il grafo sono state effettuate tutte le precedenti operazioni. La realizzazione del grafo è stata ultimata attraverso la libreria **Networkx**⁵, scritta in *Python*.^[9] Il grafo in questione può essere modificato graficamente, nel nostro caso per meglio far comprendere la polarizzazione si è deciso di adottare i colori: *Rosso* e *Blu*; per rappresentare le due partizioni raccolte attraverso la sentiment analysis, successivamente calcolarne la polarizzazione (Vedi fig:4.4).

⁵NetworkX è una libreria Python per la creazione, manipolazione e studio delle strutture, dinamiche, e delle funzioni di una rete complessa.

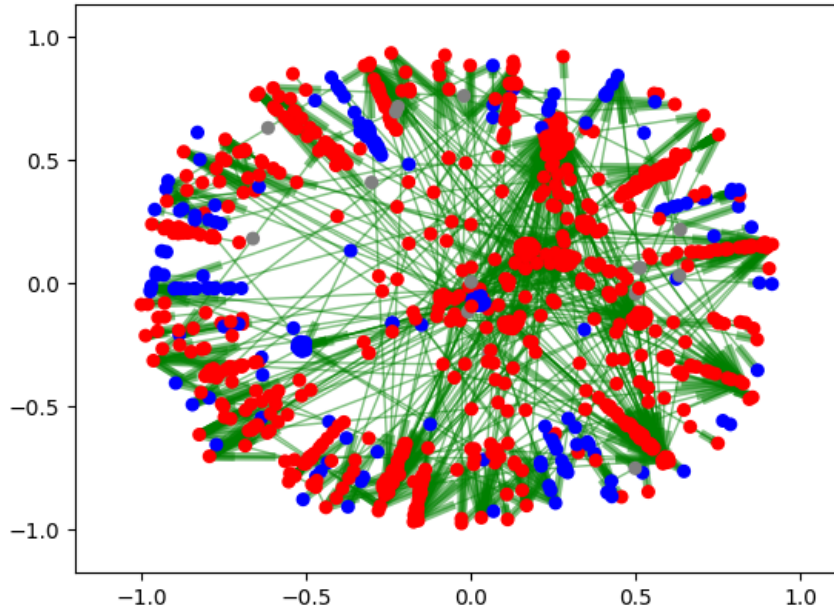


Figura 4.4: esempio di Endorsement Graph

4.4 Polarizzazione

La polarizzazione viene definita come: un processo sociale in cui un gruppo viene diviso in due sottogruppi aventi un conflitto o una visione differente del problema in questione, e da alcuni individui che rimangono neutrali. Dal punto di vista matematico la polarizzazione è una assegnazione di un valore compreso all'interno di un range:

$$[-1, 1] \quad (4.2)$$

I valori in questione identificano la vicinanza ad un gruppo piuttosto che ad un altro, in conclusione occorre assegnare alle due visioni contrastanti, che si vogliono analizzare, due valori di riferimento ovvero gli estremi. Di

seguito verranno illustrate le due implementazioni realizzate per il calcolo della polarizzazione all'interno di un endorsement graph.

4.4.1 Polarizzazione basata sul grado

All'interno di questa sezione verrà illustrata l'implementazione della polarizzazione attraverso l'algoritmo illustrato all'interno del paper: *Measuring Political Polarization: Twitter shows the two sides of Venezuela*. [13]

Come precedentemente spiegato all'interno di questa tesi, il primo passo consiste nell'individuare due tipi di nodi:

- *Elite*: sono i nodi che hanno pubblicato una notizia all'interno del grafo, sul topic selezionato.
- *Listener*: sono i nodi che seguono le informazioni pubblicate dai nodi *Elite*.

Questa suddivisione è indipendente dai risultati ottenuti dalla sentiment analysis, cioè la prima partizione effettuata analizzando i contenuti, in questo modo si analizza la diffusione delle opinioni. Illustriamo ora i passi algoritmici implementati per la realizzazione di questa polarizzazione:

1. Per prima cosa individuiamo tutti i nodi *Elite* e quelli *Listener*, per poter assegnare il primo valore della polarizzazione.

$$-1 \leq X_s \leq 1 \quad (4.3)$$

La formula precedentemente illustrata indica il range dei valori che la polarizzazione può assumere all'interno del grafo.

2. Il passo successivo, definisce le condizioni iniziali dell'algoritmo, consiste nell'assegnare ai nodi Elite e Listener rispettivamente:

$$\begin{cases} X_e = \pm 1 \\ X_l = 0 \end{cases}$$

Il valore dei noti elite dipende dall'appartenenza o meno ad i gruppi ottenuti attraverso la classificazione effettuata dalla sentiment analysis, per esempio Rossi = +1, Blu = -1.

3. I nodi elite propagheranno le loro opinioni verso i nodi listener, tale operazione verrà effettuata iterativamente fino al verificarsi di alcune condizioni, cioè diffonderà le proprie notizie ai propri vicini. Calcoliamo la polarizzazione dell'opinione per ogni listener appartenente al grafo: La polarizzazione all'istante temporale t , di un dato listener i , è data dalla seguente espressione:

$$X_i(t) = \frac{\sum_j A_{ij} X_j(t-1)}{k_i^{in}} \quad (4.4)$$

Dove A_{ij} definisce gli elementi della matrice di adiacenza del grafo, il cui valore è pari a 1 se esiste un collegamento da j a i , e k_i^{in} corrisponde al proprio *indegree*⁶. Tale formula è stata modificata poiché è stata modificata la topologia del grafo, cioè all'interno del paper gli archi congiungeva i nodi elite verso i nodi listener, mentre nella tesi i collegamenti sono stati creati nel verso opposto. L'endorsement graph collega i nodi in base alla relazione di retweet, quindi si è deciso di creare gli archi dai retweet verso il tweet. Alla luce di questa considerazione

⁶Indica il numero di archi entrati nel nodo selezionato

la formula in questione viene modificata in questo modo:

$$X_i(t) = \frac{\sum_j A_{ij} X_j(t-1)}{k_i^{out}} \quad (4.5)$$

Dove k_i^{out} corrisponde all'*outdegree*⁷ del nodo.

4. La formula illustrata nel passo precedente deve essere eseguita fino a quando avremo una stabilizzazione della polarizzazione.

Brevemente ora verrà illustrato un esempio dell'applicazione dell'algoritmo del grafo. Come prima cosa sono stati impostati i pesi ai nodi *elite*, nel nostro caso rappresentati dai nodi : A, B, C (vedi fig:5.26). Successivamente è stata applicata la formula precedentemente illustrata (vedi 4.5), quindi utilizzando i valori della polarizzazione nel passo iniziale dell'algoritmo insieme alla seguente matrice di adiacenza del grafo:

		Verso			
		A	B	C	D
Da	A	0	1	0	0
	B	0	0	0	0
	C	0	1	0	0
	D	1	1	1	0

Tabella 4.2: Matrice di adiacenza

Come possibile notare dall'esempio assistiamo ad un cambiamento di polarizzazione poiché mentre in un primo momento i nodi C e D che erano di un colore Blu (risultato ottenuto attraverso la sentiment analysis) possiamo notare come il valore cambi drasticamente in funzione del grado del nodo e agli archi, conseguentemente cambia il colore del nodo stesso. In questo

⁷Indica il numero di archi uscenti dal nodo selezionato

modo risulta evidente come un grafo cambi drasticamente la propria polarizzazione. L'algoritmo come si evince dall'esempio termina nel momento in cui l'algoritmo converge stabilizzando i valori della polarizzazione nei nodi (vedi fig:5.31).

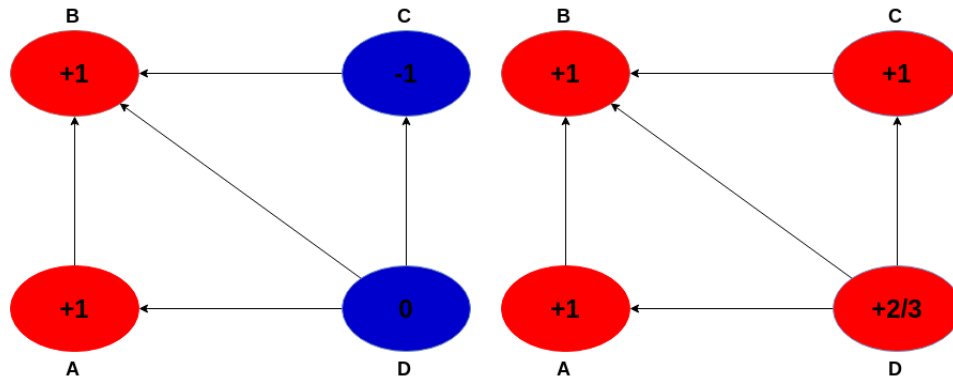


Figura 4.5: Passo 0

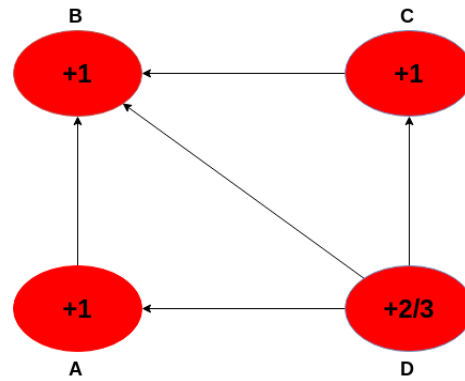


Figura 4.6: Passo 1

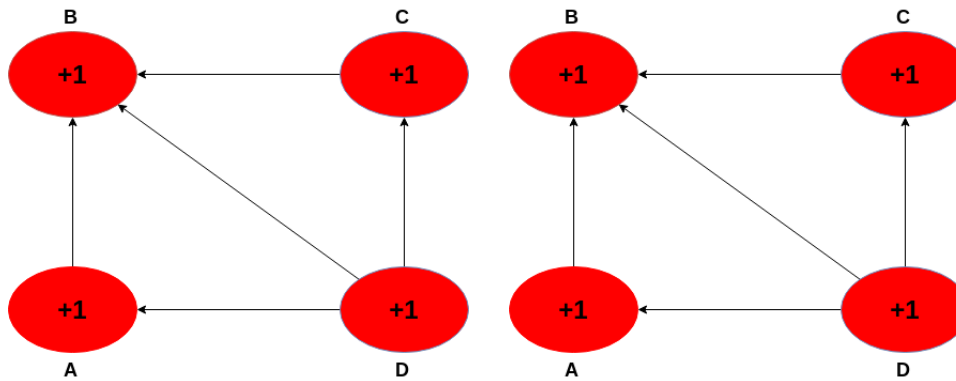


Figura 4.7: Passo 2

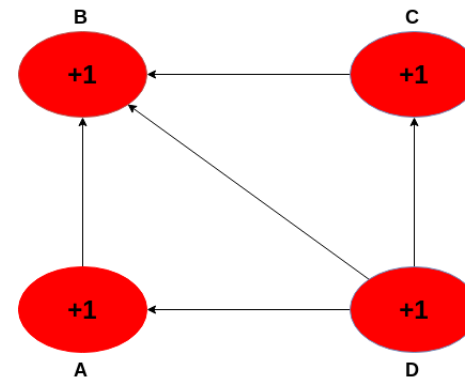


Figura 4.8: Passo 3

Figura 4.9: Esempio algoritmo basato sul grado

Dal punto di vista implementativo l'algoritmo trattato è stato implementato mediante le funzionalità contenute all'interno della libreria *NetworkX*, la quale garantisce ottime prestazioni nel caricamento dei dati relativa alla matrice di adiacenza, senza far saturare la memoria volatile del sistema.

4.4.2 Polarizzazione basata sulla topologia

La polarizzazione basata sulla topologia del grafo è un algoritmo definito all'interno di questo paper: *Reducing Controversy by Connecting Opposing Views*. L'algoritmo in questione sfrutta la topologia del grafo per poter valutare la polarizzazione di ogni singolo nodo a fronte dell'opinione espressa all'interno del tweet pubblicato. Per prima cosa una volta costruito l'*endorsement graph* un ruolo fondamentale lo svolge la scelta delle due partizioni che compongono il grafo. Nel paper in questione le partizioni venivano effettuate attraverso una catalogazione degli hashtags, cioè degli elementi testuali che esprimessero in poche parole, o una loro concatenazione, un parere su una determinata notizia.[15] Per fare un esempio se il topic dell'analisi fosse stata una partita di calcio ed un utente avesse pubblicato un tweet con un hashtag come *#ForzaBlu*, questo sarebbe stato associato come un parere positivo verso la squadra blu. Però se nel testo ci fosse stata una frase negativa insieme all'hashtag precedente, sarebbe stato un errore perché non più a favore della squadra blu, bensì per il suo avversario. Per risolvere questa problematica la suddivisione è stata effettuata attraverso la *Sentiment Analysis* (vedi sezione n:4.2, per ulteriori spiegazioni). Tornando alla definizione della polarizzazione anche in questo caso tale strumento può assumere un valore contenuto nel seguente range:

$$[-1, 1] \tag{4.6}$$

La polarizzazione applicando quanto illustrato in questo paper risulta essere:

Il tempo atteso per un random walk per raggiungere un nodo di grado massimo appartenente ad un insieme X ed Y partendo

rispettivamente da un nodo u . Otteniamo così:

$$\begin{cases} \rho^X(u) \in [0, 1] \\ \rho^Y(u) \in [0, 1] \end{cases}$$

Rispettivamente i risultati dei random walk per raggiungere i nodi di grado massimo di X ed Y⁸. Infine definiamo la polarizzazione come:

$$P_u = \rho^X(u) - \rho^Y(u) \in [-1, 1] \quad (4.7)$$

In caso di nodi che non hanno uscite, cioè quei nodi che non hanno alcun arco uscente che gli permetta di comunicare con altri, questi assumeranno il comportamento di nodi di *dangling*, cioè quei nodi che non hanno possibilità di comunicazione. Per dar loro un valore della polarizzazione gli verrà assegnato il valore massimo a seconda della partizione a cui sono associati. Un discorso analogo è valido anche per i nodi di grado massimo, questi al passo iniziale hanno una polarizzazione pari al valore massimo del loro insieme di appartenenza, però se hanno degli archi in uscita, che li metta in contatto con nodi di opinione opposta, allora tale valore verrà aggiornato attraverso l'algoritmo precedentemente illustrato.

All'interno della tesi sono state affrontate diverse problematiche per la realizzazione di questo algoritmo, per prima cosa è stata inserita la *probabilità di retweet* per calcolare la polarizzazione. In questo modo è stata assegnata una forte relazione tra i nodi che retwettano notizie e quelli che le pubblicano. Nel calcolo del random walk il numero di passi dipende dalla probabilità che relaziona i due nodi, rendendo il calcolo più verosimile alla realtà. Per

⁸Per X e Y sono da intendere come le due partizioni ottenute attraverso la sentiment analysis.

completezza si definisce *probabilità di retweet*:

Si definisce tale il rapporto che intercorre tra il numero dei retweet che l'utente ha effettuato su un nodo⁹ con il numero di retweet totali effettuati dall'utente stesso.

$$P(\text{retweet})_{ij} = \frac{\#retweet_{ij}}{\#retweet_j} \quad (4.8)$$

con i,j = nodi adiacenti appartenenti al grafo.

Detto ciò illustriamo un esempio dell'applicazione di questo algoritmo, per completezza consideriamo lo stesso grafo utilizzato in precedenza (vedi fig:5.31):

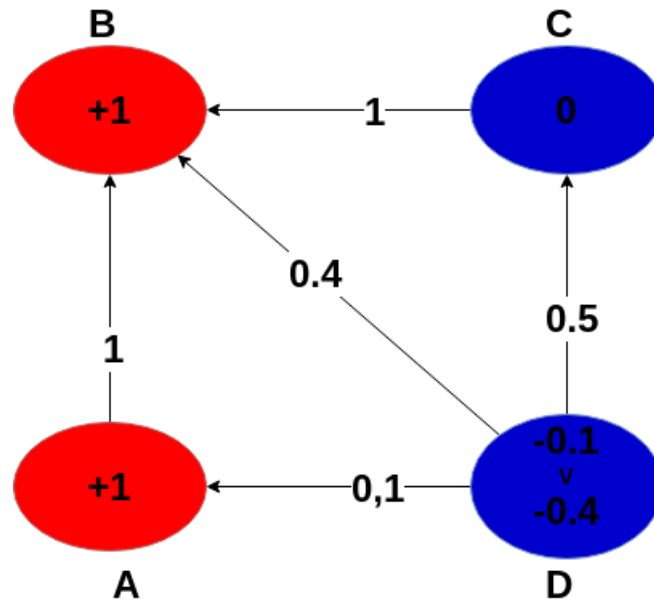


Figura 4.10: esempio di polarizzazione basata sulla topologia

Come possiamo notare i nodi di grado massimo della partizione *Rossa* e *Blu* sono rispettivamente i nodi: B e C, perché quelli aventi *indegree* massima

⁹Per nodo ci si riferisce ad un tweet pubblicato da un utente

per le rispettive partizioni. Le colorazioni sono il risultato della sentiment analysis. Analizzando i risultati presentati nell'immagine 4.10, possiamo notare come la probabilità di retweet giochi un ruolo cruciale nell'analisi in questione. Possiamo notare un comportamento particolare nel nodo B il quale assume una polarizzazione pari a zero. Essendo un nodo di grado massimo appartenente alla partizione Blu, ha pubblicato notizie a favore di quel gruppo, però ha anche condiviso notizie a favore del gruppo opposto modificando la sua natura originale. In conclusione applicando la formula 4.7, otteniamo un annullamento della sua polarizzazione. Un altro comportamento particolare è osservabile nel nodo D, questi avrà due possibili valori a seconda del percorso scelto durante il random walk. Se consideriamo il percorso per raggiungere il nodo di grado massimo appartenente alla partizione Blu è facile notare che tale valore è pari a:

$$\rho^Y(D) = (D, C) = 0.5$$

Per quanto riguarda la partizione Rossa il discorso non è semplice infatti:

$$\begin{cases} \rho^X(D) = (D, B) = 0,4 \\ \rho^X(D) = (D, A) \times (A, B) = 0,1 \times 1 = 0,1 \end{cases}$$

In conclusione la sua polarizzazione può assumere questi due valori:

$$\begin{cases} P_B = \rho^X(D) - \rho^Y(D) = 0,4 - 0,5 = -0.1 \\ P_B = \rho^X(D) - \rho^Y(D) = 0,1 - 0,5 = -0.4 \end{cases}$$

Per concludere dal punto di vista implementativo essendo molto oneroso dal punto di vista della memoria volatile, l'esecuzione di tutti questi random walk per ogni nodo, sono state ottimizzate le operazioni. Piuttosto che utilizzare la matrice di adiacenza si è utilizzato un metodo definito nella libreria *NetworkX*

che restituisce i primi vicini dell'algoritmo. Dovendo mantenere in memoria una lista di nodi visitati durante l'esecuzione, al raggiungimento di un nodo di grado massimo viene effettuata una operazione di *free* che consente al processo di deallocare le risorse occupate da tale lista.

4.5 Predizione

La predizione è una tecnica che consente di analizzare il comportamento di serie di dati nel tempo, permettendo di poter predirne il comportamento nell'istante temporale successivo a quello attualmente in vigore. Per la realizzazione degli algoritmi che hanno consentito di poter realizzare una tale operazione sono state adottate delle tecniche molto conosciute nell'ambito del *Forecasting*. Nello sviluppo di questa tesi queste tecniche sono state implementate per consentire una predizione della polarizzazione di ogni nodo nel mese successivo a quello calcolato. Per mantenere i dati e averli sempre a disposizione, questi sono stati salvati all'interno di file *csv* per poter consentire all'utente di poterli analizzare e consultare nel tempo.

4.5.1 Double Exponential Smoothing

Il Double Exponential Smoothing è una tecnica che sfrutta la serie temporale dei dati raccolti, assegnando una crescita esponenziale sui dati nel tempo, tenendo conto anche del trend di crescita o di decrescita nel tempo. In questo modo si cerca di predire il contenuto nel primo istante successivo.[3] Questa tecnica è uno sviluppo della semplice Exponential Smoothing che faceva la previsione del singolo punto senza tenere conto del trend che veniva a generarsi tra i diversi punti. La polarizzazione cambiava in base ai nuovi collegamenti che venivano a generarsi volta per volta nel tempo, per cui avere

una cognizione sul trend della serie numerica risulta molto importante. Per meglio comprendere il funzionamento di questa tecnica, verrà illustrata la *Exponential Smoothing*:

$$\hat{y}_x = \alpha y_x + (1 - \alpha) \hat{y}_{x-1} \quad (4.9)$$

Che definisce la relazione che intercorre tra l'istante attuale con quello precedente. Possiamo notare la presenza della costante α che viene definita come *fattore di smoothing*. Fondamentalmente viene effettuata una *Moving average* (vedi sezione n:4.5.3), utilizzando i fattori: α ed $(1 - \alpha)$. Analizzando la formula si può notare la presenza di una operazione ricorsiva assumendo il comportamento di un esponenziale. Il fattore di smoothing assegna un peso all'osservazione più recente rispetto all'ultimo valore precedentemente osservato. In conclusione più α è alto e maggiore e più velocemente il metodo dimentica il passato¹⁰. Alla luce di quanto illustrato con la Exponential Smoothing, è intuibile che la presenza di un trend di analisi diventa fondamentale per un'analisi a lungo termine. La polarizzazione è uno strumento a lungo termine che può rivelare il futuro andamento delle opinioni degli utenti nella rete. Definiamo la Double Exponential Smoothing come:

$$\begin{cases} l_x = \alpha y_x + (1 - \alpha)(l_{x-1} + b_{x-1}) & \text{livello} \\ b_x = \beta(l_x - l_{x-1}) + (1 - \beta)b_{x-1} & \text{trend} \\ \hat{y}(x + 1) = l_x + b_x & \text{forecast} \end{cases}$$

Il *livello* espresso nella prima equazione del sistema è simile alla equazione 4.9, con la differenza che l'osservazione precedente tiene conto anche dell'andamento (trend) dell'istante precedente. Il *trend* indica la pendenza del-

¹⁰Indica un tasso di decadenza della memoria

la funzione, dal punto di vista matematico e geometrico è assimilabile al *coefficiente angolare* di una funzione.

$$m = \frac{\Delta_y}{\Delta_x} \quad (4.10)$$

Dove Δ_y e Δ_x sono rispettivamente le differenze delle ordinate e delle ascisse tra due punti. Però dal punto di vista del Double Exponential Smoothing, possiamo assumere come la differenza tra le ascisse sia pari ad uno essendo una serie temporale, la differenza tra l'istante attuale e quello precedente è un 1. in conclusione avremo che

$$m = \frac{\Delta_y}{1} = y_x - y_{x-1} \quad (4.11)$$

Per comodità verrà indicata con b_x . Un altro fattore che risulta evidente nel sistema precedentemente illustrato è la presenza di β , questo viene chiamato *fattore di trend*. Dal punto di vista funzionale si comporta allo stesso modo di α soltanto che si riferisce al trend e non al livello. Per la realizzazione di tale algoritmo è stata implementata una funzione *Python* che effettuasse tale a calcolo a partire dalla serie di dati calcolati durante le operazioni della polarizzazione. La serie di dati in questione risultava contenere al massimo 4 elementi, ovvero i mesi scelti per il periodo dell'analisi. Per calcolare il valore successivo laddove mancassero osservazioni precedenti sono state effettuate delle assunzioni:

- In caso di una singola osservazione, per predire il futuro, non è possibile applicare tale algoritmo per cui si è scelto di restituire come valore il valore della precedente osservazione.
- Il numero minimo di osservazioni per effettuare il calcolo con questo

algoritmo deve essere pari a 2.

In seguito verranno illustrati i risultati ottenuti mediante questo approccio.

4.5.2 Linear Regression

La Regressione Lineare formalizza e risolve il problema di una relazione funzionale tra variabili misurate sulla base di dati campionari estratti da un'ipotetica popolazione infinita. In statistica l'analisi della regressione è associata alla risoluzione del modello lineare.[6] La regressione consiste nel costruire un modello attraverso cui prevedere i valori di una variabile dipendente o risposta (quantitativa) a partire dai valori di una o più variabili indipendenti o esplicative. La relazione tra due o più variabili può essere effettuata attraverso modelli matematici, nel nostro caso quello di una *retta*. La definizione di regressione lineare è:

$$Y_i = \beta_0 + \beta_1 X_i + e_i \quad (4.12)$$

- i indica le diverse osservazioni, $i = 1, \dots, n$, con n pari al numero di osservazioni
- Y_i è la variabile dipendente.
- X_i è la variabile indipendente.
- $\beta_0 + \beta_1 X$ è la retta di regressione.
- β_0 è l'intercetta, corrisponde al valore medio di Y quando X è pari a zero.
- β_1 è il coefficiente angolare, indica come varia Y in corrispondenza di una variazione unitaria di X .
- e_i è l'errore statistico.

Applicando questa formula viene generata una retta che sia più vicina possibile ai punti analizzati. Nel caso della predizione partendo dalla retta, si calcola il valore nell'istante temporale successivo individuato sulla retta. La medesima operazione è stata effettuata per la serie di dati della polarizzazione. L'implementazione di questa tecnica è stata effettuata mediante la libreria *Python numpy*, la quale consentiva attraverso una lista numerica di calcolare la regressione lineare. Una volta calcolata la retta e quindi la relativa *quota* e *coefficiente angolare*, si è applicata la classica formula geometrica:

$$y = mx + q \quad (4.13)$$

sostituendo il valore di x con l'istante temporale che si voleva calcolare, ottenendo la predizione desiderata.

4.5.3 Moving Average

La moving average è una tecnica che estende la media aritmetica dei risultati, andando a mediare i risultati all'interno di una finestra mobile. Concettualmente vengono prelevati i risultati più recenti all'interno della serie temporale, successivamente questi verranno mediati. Il risultato sarà il valore della predizione dell'istante temporale richiesto.[7] Viene anche definita *finestra mobile*, proprio perché calcola la media dei valori all'interno di una finestra limitata rispetto alla serie di dati contenuti. La definizione matematica è la seguente: Data una serie storica y_t , con $t = 1, \dots, T$, contenente i valori osservati di una variabile Y dal tempo 1 al tempo T , siano:

- m_1 il numero dei periodi precedenti a t ;
- m_2 il numero dei periodi successivi a t ;

- θ_i il peso da attribuire all' *i-esimo* valore osservato;

$$Media_t = \frac{1}{k} \sum_{i=-m_1}^{m_2} \theta_i y_{t+i} \quad (4.14)$$

In conclusione attraverso questa tecnica di forecasting è possibile fare una predizione del valore successivo, attraverso una media tra i valori antecedenti all'evento richiesto. In questo modo viene definita una forte correlazione tra i valori vicini. La moving average è molto utile nel caso si abbiano pochi dati a disposizione. Dal punto di vista implementativo, utilizzando i dati ottenuto attraverso la polarizzazione, è stato necessario fare delle assunzione:

- Non è possibile utilizzare questa tecnica con serie numeriche aventi un solo elemento, infatti per sopperire a tale problematica si suppone che la previsione coincida con l'ultimo valore acquisito.
- La lunghezza massima della finestra temporale, è stata impostata pari a 2 in questo modo si cerca di tener conto del trend ottenuto.

In conclusione l'implementazione di tale funzionalità è stata implementata attraverso il *Python*.

Capitolo 5

Test sperimentali e valutazione

All'interno di questo capitolo verranno illustrati i test sperimentali con le relative valutazioni. Tali operazioni sono stati svolti per due topic distinti:

- *Elezioni Regionali Sicilia 2017*
- *Biotestamento*

Prima di illustrare i risultati ottenuti, verranno presentate le motivazioni della scelta di questi due topic differenti tra loro. Sono stati scelti perché identificano differenti contesti nel dettaglio:

- Politico
- Sociale

Questa differenza è molto importante perché dimostra come la polarizzazione possa essere utilizzata in contesti completamente differenti tra loro. Infatti questo strumento consente di avere una visione generale sulla diffusione di opinioni e l'eventuale formazione di *Echo-Chambers*. Le motivazioni che mi hanno spinto a scegliere questi due contesti così differenti tra loro è stato dettato dalla necessità di cercare degli argomenti, che consentissero di avere

una forte divergenza di opinione. Nel contesto politico le divergenze di opinioni vengono generate già dagli esponenti dei diversi partiti politici, quindi utilizzando questi dibattiti è possibile definire due gruppi distinti di opinione. Un gruppo che aderisce ad un'idea espressa da un partito piuttosto che da un'altra. In conclusione studiare ed analizzare un topic in un contesto del genere risulta essere molto interessante. Per quanto riguarda il contesto sociale, questo è molto interessante, perché gli argomenti che appartengono a questa categoria hanno presa su un campione maggiore di utenti, perché molto spesso collegate a problematiche religiose, politiche, etiche ed economiche. La scelta di questo contesto consente di poter verificare quanto un topic, sia appetibile e polarizzato. Nel dettaglio è interessante notare come utilizzare uno strumento come la polarizzazione possa essere usato non solo per misurare quanto un argomento divida due gruppi di utenti, ma anche per fare indagini di mercato, statistiche da poter utilizzare in futuro.

5.1 Elezioni Regionali Sicilia 2017

Le elezioni regionali in Sicilia sono state svolte il 5 Novembre del 2017, è stato scelto come topic di studio per calcolare la polarizzazione all'interno di un contesto politico. Le suddette elezioni hanno mostrato un cambiamento di trend, poiché vinte dalla coalizione del centro-destra a discapito della coalizione del centro-sinistra che era il partito in carica. I risultati di questa elezione hanno mostrato che gli antagonisti principali della coalizione del centro-destra fosse il Movimento 5 stelle. I risultati ottenuti attraverso queste elezioni rispecchiano ampiamente il panorama politico dello stato italiano, poiché da molti anni non è più presente un'unica linea guida, ma ci sono molteplici visioni espresse, che possono generare molteplici conflitti e quindi una forte polarizzazione delle opinioni. Nel dettaglio è interessante notare come in Italia la presenza dei dibattiti innescati da parte dei diversi esponenti della classe politica abbia ripercussioni sugli utenti appartenenti ad una rete sociale. In questo periodo storico vige un certa sfiducia da parte degli elettori verso la classe politica, infatti nelle elezioni regionali in Sicilia la percentuale degli astenuti è stata di $\simeq 54\%$ a discapito del $\simeq 53\%$ ottenuto nelle precedenti elezioni nel 2012. Queste informazioni rendono poco accurate le predizioni fatti dai diversi enti per calcolare gli *exit pool*. [12] Dal punto di vista della polarizzazione le differenti correnti politiche generano conflitti, divergenze o anche convergenze di opinioni tra gli ascoltatori; quindi fare uno studio attraverso questi dati raccolti può mostrare quanto la controversia attiri gli elettori. La classe politica sta sempre più utilizzando i social media per poter esprimere le loro idee, anche per catturare l'attenzione di diversi elettori, per rendere accessibili a più utenti le proprie visioni politiche. Proprio per questo motivo si è deciso di utilizzare le elezioni regionali in Sicilia come caso di studio all'interno del social media di *Twitter*, sempre

più utilizzato dalla classe politica.

I test sono stati effettuati dopo una raccolta di dati attraverso una ricerca per *hashtag* (vedi sezione 4.1). La lista degli hashtag utilizzati è la seguente:

- **#elezionisicilia2017**
- **#EleSicilia**
- **#elezionisicilia**
- **#regionalisicilia2017**
- **#regionalisicilia**

Per ognuno di questi hashtag sono stati collezionati i tweet pubblicati dai diversi utenti del social media durante il periodo dal 01/09/2017 al 31/12/2017. Il motivo della scelta di un periodo del genere è dettata dalla necessità di analizzare la polarizzazione del grafo nel tempo. Il periodo utilizzato per la raccolta dei dati è stato scelto per analizzare il comportamento della polarizzazione in un periodo di tempo precedente ed antecedente all'evento in questione. In questo modo è possibile effettuare considerazioni sul comportamento temporale del topic. La scelta degli hashtag non è puramente casuale, infatti per garantire una imparzialità nella raccolta dei tweet sono state fatte delle ricerche all'interno del social media per identificare gli hashtag utilizzati da giornalisti, redazioni giornalistiche per poter parlare in maniera neutrale delle elezioni regionali in Sicilia. Una volta identificati attraverso una ricerca manuale sulla rete è stato utilizzato un framework online per ricercare le loro eventuali correlazioni con altri:

- **hashtagify** (vedi fig:5.1):consente di ricercare degli hashtag attraverso una barra di ricerca. Una volta terminata, verranno illustrati all'interno di una interfaccia grafica alcuni risultati quali: gli hashtag correlati,

la popolarità nel tempo, i maggiori *influencer*. Attraverso questo servizio è stato possibile identificare altri hashtag generici utilizzati per arricchire la lista degli hashtag sopra elencati.



Figura 5.1: Hashtagify

Ricerca la neutralità attraverso gli hashtag è molto importante, soprattutto in un contesto politico, perché molti partiti o rappresentanti politici utilizzano questi strumenti per diffondere le loro opinioni, quindi utilizzare quelle parole chiave come elementi per la raccolta dati non avrebbe permesso di avere uno studio imparziale sul topic. Per fare un esempio basta considerare che il partito del candidato *Musumeci* si chiamava: *#Diventerà bellissima - Per la Sicilia*. Un altro esempio può essere l'hashtag lanciato dal Movimento 5 Stelle: *#SceglieteIlFuturo*; utilizzato per incoraggiare gli elettori a votare per il loro candidato *Cancellieri*. In conclusione per evitare argomenti fortemente polarizzati in un unico senso, i test e quindi la raccolta dati sono stati effettuati attraverso una lista di parole chiave che descrivessero il contenuto senza esprimere una opinione di partenza.

Una volta che i dati sono stati raccolti attraverso la ricerca per hashtag, è stata effettuata una prima operazione di catalogo attraverso la *Sentiment Analysis* (vedi sezione:4.2). Attraverso questa operazione verranno classificati i tweet in 3 gruppi distinti, di cui soltanto per 2 di loro verrà calcolata la polarizzazione. Nel dettaglio i gruppi in questione sono stati classificati attraverso 3 *label*, utilizzati dalla sentiment analysis per catalogare il testo, analizzando il loro contenuto:

- **ForzaItalia** identifica tutti i tweet e retweet relativi alla coalizione del centro-destra.
- **5Stelle** identifica tutti i tweet e retweet relativi al Movimento 5 Stelle.
- **Altri** identifica tutti i tweet e retweet relativi alle restanti coalizioni, informazioni giornalistiche oppure altri tweet che non avevano alcun interesse con l'analisi in questione.

Si è scelto di calcolare la polarizzazione solo per le due forze politiche che hanno dominato la scena politica siciliana: la coalizione del centro destra ed il Movimento 5 Stelle. Il centro-destra rappresenta un'insieme di partiti, si è deciso di raggrupparli all'interno di un'unica coalizione, mentre l'altra forza politica in gioco, ovvero il Movimento 5 Stelle, rappresenta un unico partito politico. La sentiment analysis utilizza un training set per catalogare le informazioni in base al sentimento che esprimono i tweet pubblicati dagli utenti. Per questo motivo è stata fatta una operazione di raccolta e catalogo di diversi tweet di esempio per consentire alla macchina di poter comprendere a quale gruppo associare il tweet che avrebbe analizzato. Qui di seguito mostriamo una tabella contenente 3 tweet di esempio per ogni label per meglio rendere l'idea.

Training set	
Label	Tweet
5stelle	Chi ama la sua terra non può che votare il #M5S #Regionali #Sicilia
5stelle	A disposizione del #MoVimento5Stelle per cambiare la nostra fantastica #Sicilia! Contattatemi #regionaliSicilia #m5s
5stelle	#Regionalisicilia2017 con la vittoria del #M5S i partiti capiranno tutto il potere che ha un popolo unito!
ForzaItalia	Con elezioni @matteosalvinimi e @Musumeci_Staff regionali 2017 in #Sicilia !!#elezioniregionali2017 #andiamoagovernare #forzalega
ForzaItalia	Regionali Sicilia Berlusconi a Palermo: "No a M5s chi li vota non ragiona"
ForzaItalia	'Siamo qui a #Palermo per vincere con @Musumecistaff , #M5s non hanno arte né parte #regionalisicilia
Altri	Tutti gli schieramenti alle #elezioniregionali in Sicilia.L'articolo di Pierangelo Bonanno
Altri	#Londra ricevo questa card per gli #emigrati : è valida solo per i #treni in #Italia . E le migliaia di #Siciliani all' #estero ? #elesicilia
Altri	#cambiamoabitudini Basta! Adesso si agisce.Il #5novembre #io voto #siciliafutura #micaripresidente #elesicilia

Tabella 5.1: Esempi di tweet del Training Set

	Settembre	Ottobre	5 Novembre	Novembre	Dicembre
Utilizzati	350	968	2145	5183	5270
Totali	464	1158	2631	6779	6876

Tabella 5.2: Numero di nodi del grafo

Una volta partizionato i dati, l'endorsement graph relativo alle elezioni siciliane è stato popolato. Questo grafo è cresciuto nel tempo, cioè via via che venivano raccolti i tweet ed i relativi retweet durante la raccolta dati effettuata mese per mese. I tweet isolati cioè quelle informazioni, con espressione di idee, che non sono state ripubblicate dagli altri utenti della rete sono stati scartati, ma presi in riesame nel caso in cui nei mesi successivi ci fosse stato qualche utente che avesse retweettato quelle notizie. Per meglio rendere l'idea in questione illustriamo una tabella (vedi tabelle: 5.2) contenente i dati raccolti nei diversi mesi, e quelli effettivamente utilizzati per il calcolo della polarizzazione.

Per quanto riguarda la colorazione del grafo i colori di riferimento delle due partizioni sono: **Blu** per la coalizione del centro-destra; **Rosso** per il Movimento 5 Stelle. Essendo la polarizzazione un'assegnazione di un valore all'interno di un range tra $[-1, +1]$, i valori di riferimento per i nodi *Blu* sono definiti nel range $(0, +1]$ mentre per i nodi *Rossi* in $[-1, 0)$. Per meglio contraddistinguere quegli utenti che hanno mantenuto uno schieramento neutrale, quindi quei nodi del grafo avente polarizzazione esattamente pari a 0, si è deciso di assegnare loro il colore **Grigio**. Mostriamo i grafi ottenuti durante la raccolta dei dati durante il periodo di tempo sopra indicato. Nel dettaglio la colorazione è stata modificata in base al valore ottenuto durante

la polarizzazione, quindi con i nodi colorati di *rosso*, *blu* e grigio. I risultati mostrati (vedi fig:5.7) illustrano la popolazione del grafo attraverso il calcolo della polarizzazione attraverso l'algoritmo basato sulla topologia. I risultati sono cumulativi e mostrano il cambiamento del grafo nel tempo. Ogni grafo mostra risultati cumulativi; cioè ogni mese contiene anche i dati raccolti nel periodo ad esso precedente. Le immagini successive (vedi fig:5.32), mostrano le medesime considerazioni, soltanto che è stato applicato un algoritmo differente per il calcolo della polarizzazione.

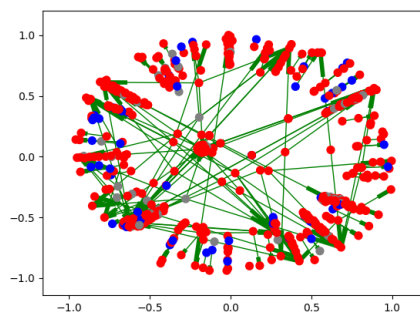


Figura 5.2: Settembre

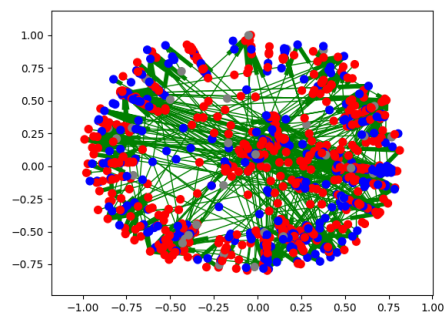


Figura 5.3: Ottobre

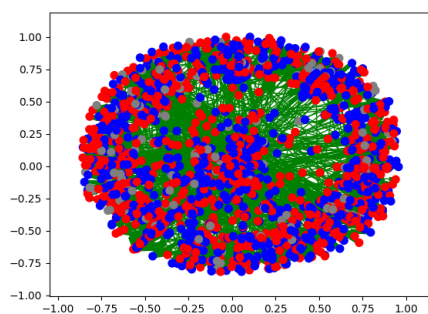


Figura 5.4: 5 Novembre

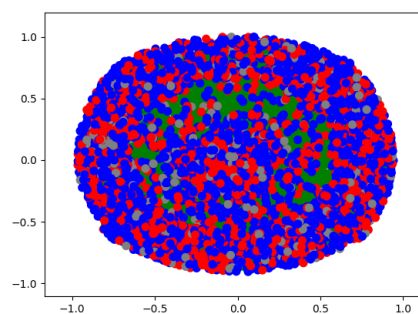


Figura 5.5: Novembre

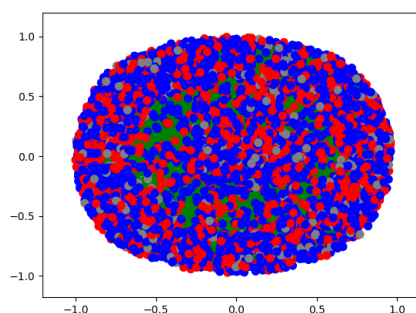


Figura 5.6: Dicembre

Figura 5.7: Popolazione del grafo attraverso l'algorithmo basato sulla topologia

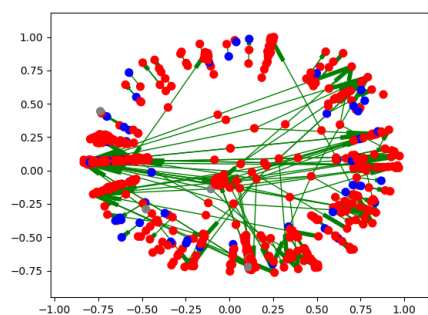


Figura 5.8: Settembre

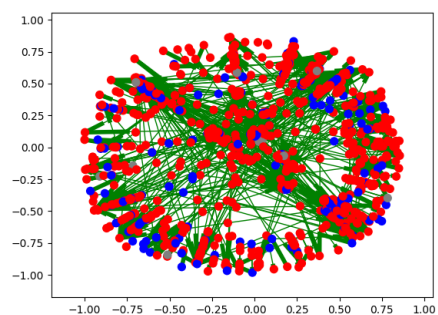


Figura 5.9: Ottobre

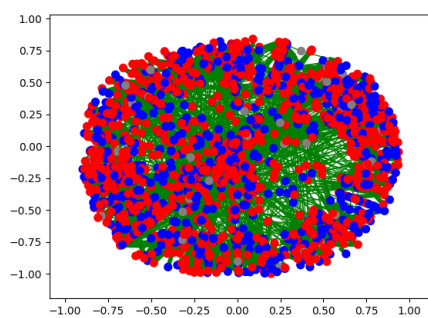


Figura 5.10: 5 Novembre

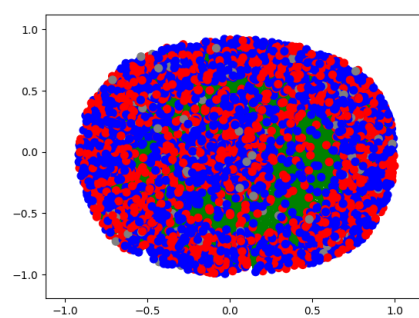


Figura 5.11: Novembre

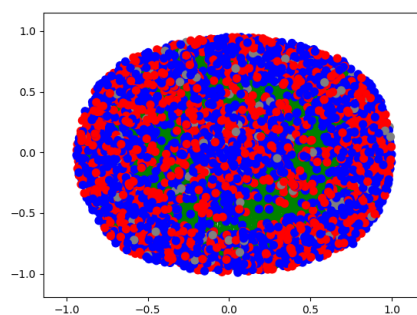


Figura 5.12: Dicembre

Figura 5.13: Popolazione del grafo attraverso l'algoritmo basato sul grado del nodo

Il calcolo della polarizzazione è stato effettuato con i due algoritmi presentati in precedenza (vedi sezione:4.4), per meglio comprendere i risultati ottenuti sono stati fatti dei grafici che mostrano l'andamento della polarizzazione nel tempo. Dopo aver analizzato i grafici si è dimostrato come la *polarizzazione basata sulla topologia* fosse più accurata di quella *basata sul grado del nodo*. La motivazione per cui si verifica questa situazione è dettata dalla *probabilità di retweet*, cioè il valore che è stato assegnato ad ogni arco del grafo in base alle preferenze degli utenti; cioè dal numero di volte che un utente ripubblicava le notizie di un altro utente. Per riscontrare in maniera visiva il comportamento di questa informazione basterà confrontare i due diagrammi mostrati.

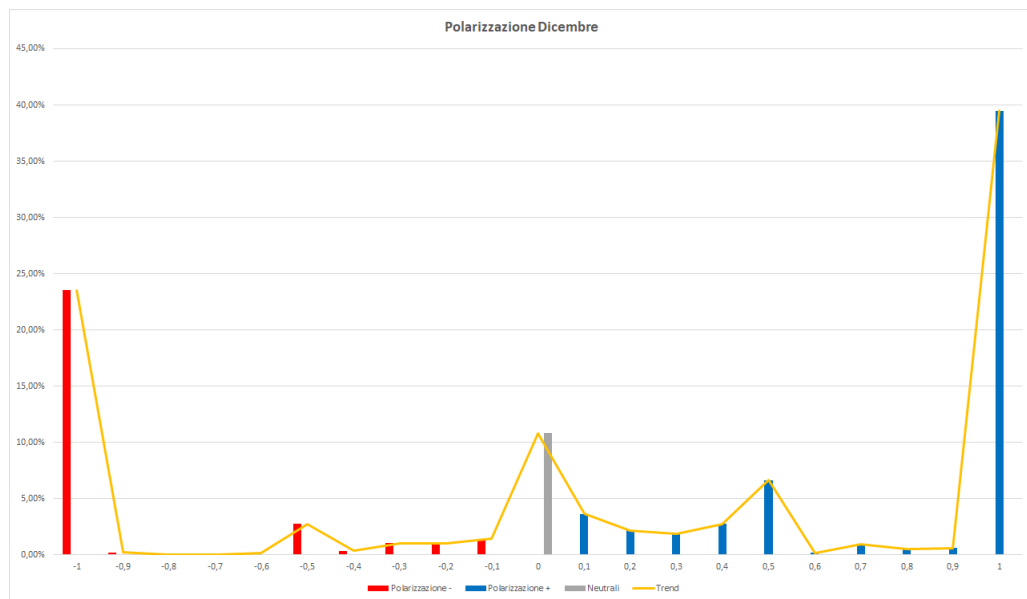


Figura 5.14: Risultati Dicembre algoritmo basato sulla topologia

Il risultati più interessanti sono stati ottenuti durante la raccolta dati nel mese di Novembre, il giorno 5 Novembre ci sono state le elezioni in Sicilia. Nel grafico (vedi fig:5.16) sono stati confrontati i risultati ottenuti da Settembre fino a Novembre e quelli ottenuti da Settembre fino al 5 Novembre.

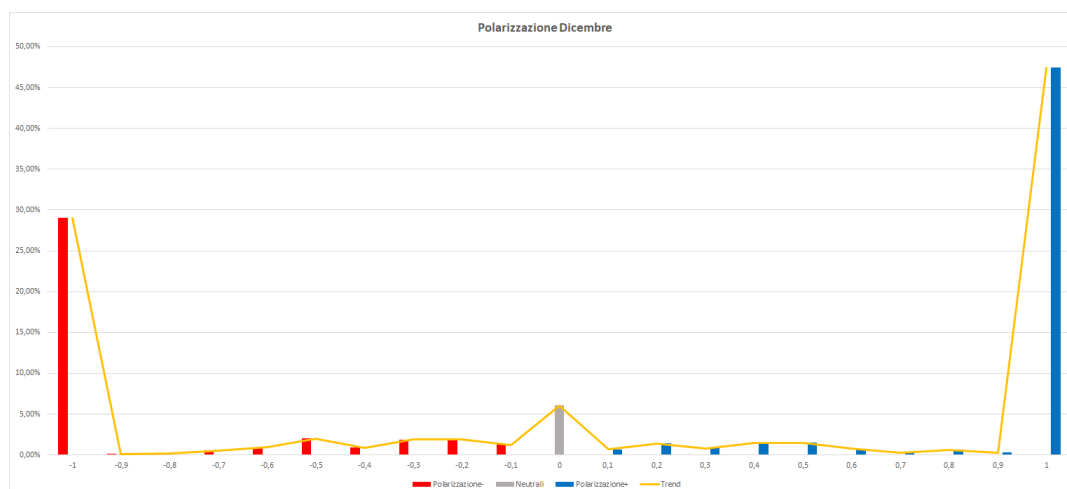


Figura 5.15: Risultati Dicembre algoritmo basato sul grado del nodo

Risulta evidente come i dati raccolti prima delle elezioni, mostrino che le opinioni a favore del Movimento 5 Stelle fossero molto simile a quelle espresse dalla coalizione del centro destra. Questo risultato è conforme agli exit-pool espressi poco prima delle elezioni fatti dalla *Rai*, i quali predicavano una lotta serrata tra il candidato *Cancellieri*¹ e il candidato *Musumeci*².^[10] I risultati delle elezioni mostrano un capovolgimento di fronte cioè la vittoria della coalizione del centro-destra a discapito di quella dei 5 Stelle. Un'altra importante considerazione che possiamo fare alla luce dei test effettuati all'interno di questo topic è i cambiamento di trend nel tempo. Avendo raccolto dati per un periodo molto lungo si è potuto assistere all'evoluzione del trend della polarizzazione ottenendo un capovolgimento di fronte; infatti mentre nei mesi di Settembre ed Ottobre la rete era maggiormente influenzata dalle opinioni del Movimento 5 Stelle, a Novembre abbiamo ottenuto un considerevole capovolgimento in conformità con quanto avvenuto nella realtà. A fini didattici si è deciso di analizzare il comportamento della rete anche nel mese

¹Il candidato presidente della regione Sicilia per il Movimento 5 Stelle

²Il candidato presidente della regione Sicilia per la coalizione del centro destra

di Dicembre per verificare se la rete cambiasse il proprio comportamento e quindi la propria polarizzazione; come indicato nella tabella 5.2, i nodi raccolti tra il mese di Novembre e di Dicembre era di 87 nuovi tweet e/o retweet, notando che il comportamento della polarizzazione non era cambiato e che ormai il topic in questione non era più considerato appetibile.

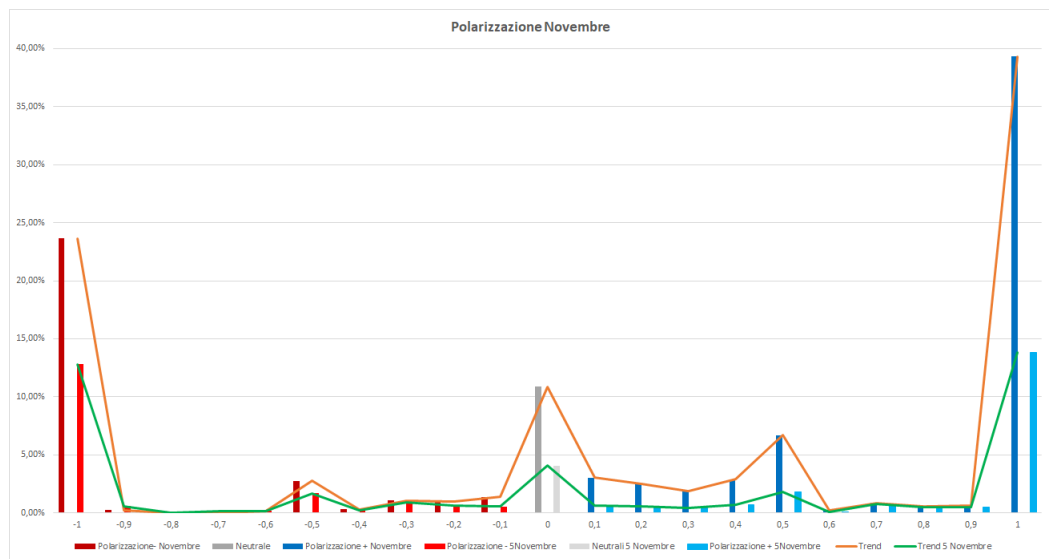


Figura 5.16: Risultati Novembre

Durante lo sviluppo della tesi, analizzando i dati ottenuti, si è deciso di realizzare un principio di predizione della polarizzazione in mesi successivi utilizzando i dati calcolati nei periodi precedenti. Il motivo di tale scelta è dettata dalla necessità di poter intuire il comportamento della rete, quindi capire come saranno partizionate le opinioni degli utente, nel futuro. Le tecniche utilizzate per la la predizione sono state illustrate all'interno della sezione: 4.5. Durante la fase di test di queste tecniche di predizione, sono state fatte alcune considerazioni:

- In caso di una singola osservazione, per predire il futuro, non è possibile applicare tale algoritmo per cui si è scelto di restituire come valore il

valore della precedente osservazione.

- Il numero minimo di osservazioni per effettuare il calcolo con questo algoritmo deve essere pari a 2.

Una volta ottenuti i risultati si è calcolato l'errore medio e la relativa varianza in base ai risultati conseguiti nel calcolo della polarizzazione attraverso i due algoritmi implementati. Di seguito mostriamo i risultati conseguiti attraverso un diagramma (vedi fig:5.36)

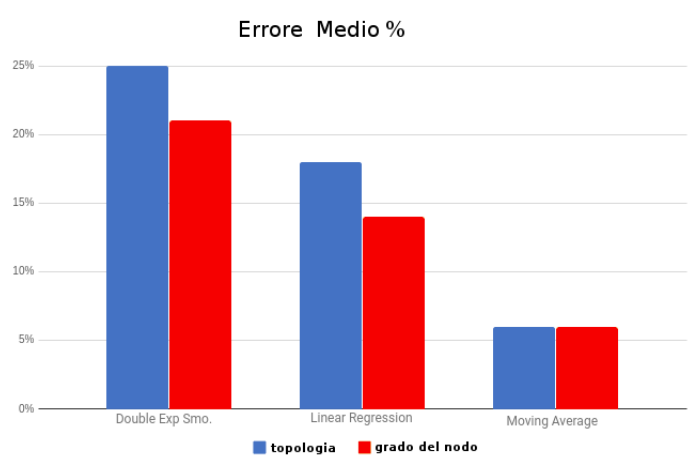


Figura 5.17: Errore medio predizione

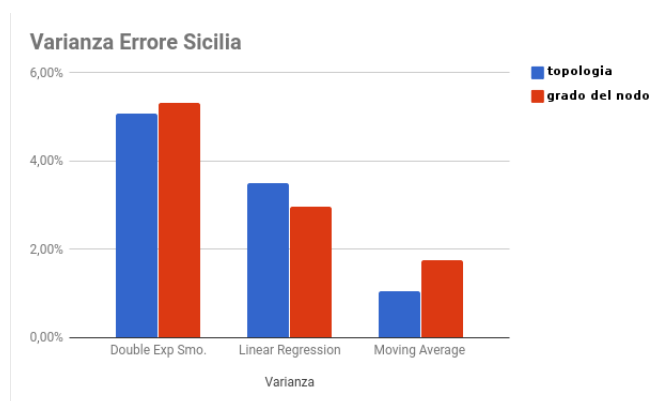


Figura 5.18: Varianza tasso di errore predizione

5.2 Biotestamento

Il biotestamento è una legge approvata dal governo italiano il 14 Dicembre 2017. Questa legge consente a qualsiasi individuo che abbia compiuto i 18 anni di poter depositare presso sede comunali e notarili, un testamento che regola la sanità dell'individuo. Nel dettaglio all'interno di questo testamento l'individuo può inserire le proprie volontà, per regolamentare la propria sanità. Queste volontà sono definite *DAT*, nel dettaglio *disposizione anticipate di trattamento*, sono disposizioni che possono essere impugnate dal fiduciario dichiarato dall'intestatario del biotestamento che dovranno essere eseguite in mancanza di volontà del paziente. Questa legge non consente alcuna forma di *eutanasia*, però consente di poter richiedere la sospensione idrica e nutrizionale da non confondere con una morte assistita perché potrebbero essere prescrizioni assegnate da un medico e che il paziente potrebbe rifiutare. Le volontà redatte all'interno del documento possono anche essere revocate dai pazienti, i quali essendo in grado di intendere e di volere possono anche accettare le cure offerte dal medico indipendentemente da quanto dichiarato nel documento. Per quanto riguarda i medici, questi possono anche rifiutarsi di adempiere all'impegno scritto nel biotestamento dal malato, però la struttura ospedaliera dovrà garantire che un altro medico si prenda l'onere di applicare tali dettami.[1]

Come descritto in questo breve paragrafo il topic scelto accomuna molteplici problematiche che essere fonte di dibattito all'interno dello stato italiano. Possiamo assimilarlo come un problema di polarizzazione delle opinioni all'interno di un contesto etico e sociale. Affrontando l'analisi di questo topic, ci si era posti l'obiettivo di analizzare l'opinione degli utenti all'interno della rete di *Twitter* in un periodo molto lungo; sono stati collezionati tweet e retweet dal 01/09/2017 al 31/01/2018. La scelta di un tale topic è nata dalla curiosità

di testare le potenzialità della polarizzazione in ottica sociale, vista anche la risonanza che aveva riscontrato nella popolazione attraverso i mass media. Dal punto di vista scientifico è interessante anche analizzare il comportamento degli utenti italiani, vista anche la connotazione religiosa che assume lo stato italiano essendo uno paese prevalentemente *Cattolico*. Ovviamente un dibattito acceso su questo argomento pone in gioco differenti correnti di pensiero che spaziano dalla politica, alla religione fino all'etica. Lo studio che è stato affrontato su questo topic verte nel calcolare la polarizzazione semplicemente tra i favorevoli e contrari alla legge presentata.

I test sono stati effettuati dopo una raccolta di dati effettuati una ricerca per *hashtag* (vedi sezione 4.1). La lista degli hashtag utilizzati è la seguente:

- **#biotestamento**
- **#eutanasia**
- **#finevita**
- **#suicidioassistito**
- **#testamentobiologico**

Per ognuno di questi hashtag sono stati collezionati i tweet pubblicati dai diversi utenti del social media durante il periodo dal 01/09/2017 al 31/01/2018. Il motivo della scelta di un periodo del genere è dettata dalla necessità di analizzare la polarizzazione del grafo nel tempo. Il periodo utilizzato per la raccolta dei dati è stato scelto per analizzare il comportamento della polarizzazione in un periodo di tempo precedente ed antecedente all'evento in questione. In questo modo è possibile effettuare considerazioni sul comportamento temporale del topic. La scelta degli hashtag non è puramente casuale, infatti per garantire una imparzialità nella raccolta dei tweet sono state fatte

delle ricerche all'interno del social media per identificare gli hashtag utilizzati da giornalisti e testate giornalistiche per poter parlare in maniera neutrale del biotestamento. Anche per questo topic si è utilizzato un framework che consentisse di poter ricercare quegli hashtag che non esprimessero un giudizio favorevole o contrario sull'argomento. Il servizio in questione è *hashtagify* (vedi fig:5.1) precedentemente illustrato all'interno della sezione: 5.1. Questo argomento essendo di grande interesse per gli utenti della rete necessita di un approccio molto neutrale nella raccolta dei tweet, con cui esprimono le proprie opinioni, per non drogare i risultati della polarizzazione. In questo topic sono stati lanciati da molti utenti o da associazioni religiose hashtag di denuncia in cui gli utenti si appellavano in maniera negativa riguardo l'attuazione del biotestamento; per esempio: *#NoDAT*, *#NoBiotestamento*. Vere e proprie campagne per denunciare la legge.

Per cercare di rimanere più neutrali e cercare di catalogare le informazioni attraverso le opinioni scritte dagli utenti è stata utilizzata la *sentiment analysis*, che consente di catalogare le informazioni attraverso una analisi sui contenuti pubblicati. Quindi una volta raccolta tutti i dati per gli hashtag, i tweet sono stati classificati in due gruppi di distinti. Nel dettaglio i gruppi in questione sono stati classificati attraverso 2 *label*, utilizzati dalla sentiment analysis per catalogare il testo, analizzando il loro contenuto:

- **Blue:** identifica tutti quei tweet e/o retweet che esprimono un dissenso sulla legge emanata il 14 Dicembre sul biotestamento.
- **Red:** identifica tutti quei tweet e/o retweet che esprimono un consenso favorevole sulla legge emanata il 14 Dicembre sul biotestamento.

La sentiment analysis utilizza un training set per catalogare le informazioni in base al sentimento che esprimono i tweet pubblicati dagli utenti. Per

questo motivo è stata fatta una operazione di raccolta e catalogo di diversi tweet di esempio per consentire alla macchina di poter comprendere a quale gruppo associare il tweet che avrebbe analizzato. Qui di seguito mostriamo una tabella (vedi 5.3) contenente 3 tweet di esempio per ogni label per meglio rendere l'idea.

Training set	
Label	Tweet
Blue	#biotestamento #testamentobiologico #eutanasia battete le mani in nome del progresso e morirete felici a contenti Dove stiamo andando?
Blue	L'aspetto forse più aberrante è il divieto di obiezione di coscienza. Se il paziente chiede di morire il medico dovrà obbedire, ossia sarà costretto a compiere un assassinio pena il licenziamento e, forse, una condanna penale. Altro che Nazismo! #eutanasia #suicidioassistito
Blue	Una volta sancito che vi è il diritto di togliersi la vita, la posizione di chi pretendere di limitarlo ai malati terminali diventa, come minimo, arbitraria. #DAT #eutanasia #finevita
Red	ti dico solo che la sedazione profonda è accettata anche dai medici obiettori e da molti cattolici tradizionalisti. Non è eutanasia o sospensione cure e viene fatta in accordo fra paziente e medico. Poi se tu vuoi soffrire, libera, ma mi pare un proposito masochista. #biotestamento'
Red	Sul #biotestamento o, meglio, sull'autodeterminazione del malato. Una legge assolutamente positiva, che ha assai poco a che fare con l' #eutanasia o #finevita

Red	#eutanasia Tutti parlano di diritto alla vita, ma pochi si preoccupano del diritto a morire in modo dignitoso. Eppure non è difficile immaginare cosa sia peggiore della morte immedesimandosi con partecipazione, e non superficialmente, in certe situazioni.....
-----	---

Tabella 5.3: Esempi di tweet del Training Set del biotestamento

Una volta partizionato i dati, l'endorsement graph relativo al biotestamento è stato popolato. La cardinalità del grafo aumentava nel tempo, cioè via via che venivano raccolti i tweet ed i relativi retweet attraverso la raccolta dati effettuata mese per mese. I tweet isolati cioè quelle informazioni, con espressione di idee, che non sono state ripubblicate dagli altri utenti della rete sono stati scartati, ma presi in riesame nel caso in cui nei mesi successivi ci fosse stato qualche utente che avesse retweettato quelle notizie. Per meglio rendere l'idea in questione illustriamo una tabella (vedi tabelle: 5.4) contenente i dati raccolti nei diversi mesi, e quelli effettivamente utilizzati per il calcolo della polarizzazione.

	Settembre	Ottobre	Novembre	14 Dicembre	Dicembre	Gennaio
Utilizzati	723	1877	4145	5103	6050	7994
Totali	934	2346	5514	6749	7630	9916

Tabella 5.4: Numero di nodi del grafo

La colorazione del grafo in riferimento delle due partizioni sono: **Blu** per gli utenti contrari al biotestamento; **Rosso** per favorevoli. Essendo la pola-

rizzazione un'assegnazione di un valore all'interno di un range tra $[-1, +1]$, i valori di riferimento per i nodi *Blu* sono definiti nel range $(0, +1]$ mentre per i nodi *Rossi* in $[-1, 0)$. Per meglio contraddistinguere quegli utenti che hanno mantenuto uno schieramento neutrale, quindi quei nodi del grafo avente polarizzazione esattamente pari a 0, si è deciso di assegnare loro il colore **Grigio**. Mostriamo i grafi ottenuti durante la raccolta dei dati durante il periodo di tempo sopra indicato. Nel dettaglio la colorazione è stata modificata in base al valore ottenuto durante la polarizzazione, quindi con i nodi colorati di *rosso*, *blu* e *grigio*. I risultati mostrati (vedi fig:5.25) illustrano la popolazione del grafo attraverso il calcolo della polarizzazione attraverso l'algoritmo basato sulla topologia. I risultati sono cumulativi e mostrano il cambiamento del grafo nel tempo. Ogni grafo mostra risultati cumulativi; cioè ogni mese contiene anche i dati raccolti nel periodo ad esso precedente. Le immagini successive (vedi fig:??), mostrano le medesime considerazioni, soltanto che è stato applicato un algoritmo differente per il calcolo della polarizzazione.

Il calcolo della polarizzazione è stato effettuato con i due algoritmi presentati in precedenza (vedi sezione:4.4), per meglio comprendere i risultati ottenuti sono stati fatti dei grafici che mostrano l'andamento della polarizzazione nel tempo. Dopo aver analizzato i grafici si è dimostrato come la *polarizzazione basata sulla topologia* fosse più accurata di quella *basata sul grado del nodo*. La motivazione per cui si verifica questa situazione è dettata dalla *probabilità di retweet*, cioè il valore che è stato assegnato ad ogni arco del grafo in base alle preferenze degli utenti; cioè dal numero di volte che un utente ripubblicava le notizie di un altro utente. Per riscontrare in maniera visiva il comportamento di questa informazione basterà confrontare i due diagrammi mostrati. (vedi fig:5.33 e 5.34)

Il risultati più interessanti sono stati ottenuti durante la raccolta dati nel

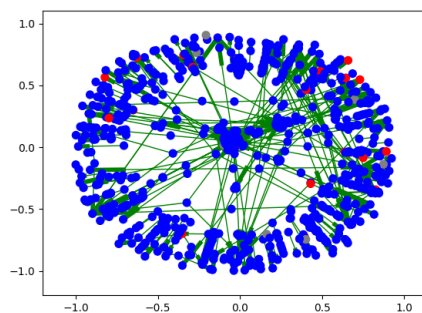


Figura 5.19: Settembre

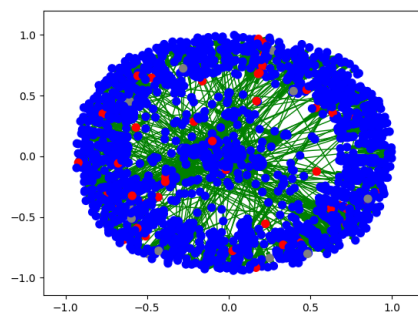


Figura 5.20: Ottobre

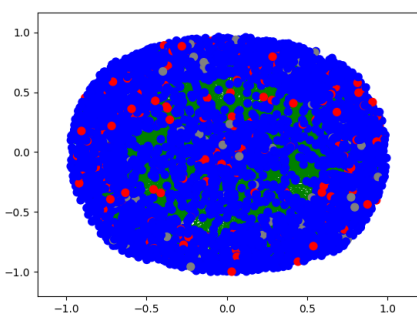


Figura 5.21: Novembre

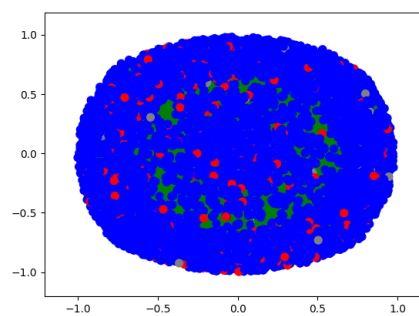


Figura 5.22: 14 Dicembre

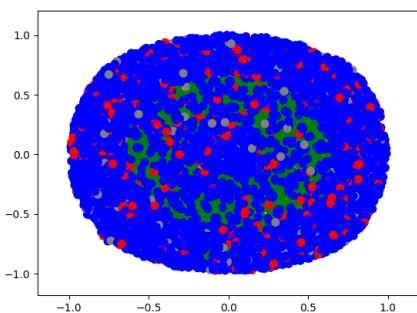


Figura 5.23: Dicembre

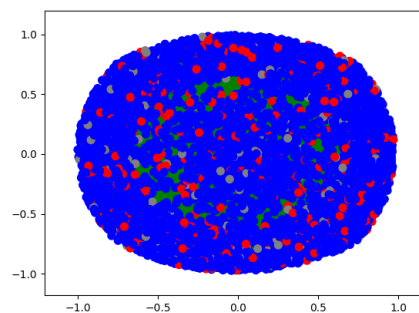


Figura 5.24: Gennaio

Figura 5.25: Popolazione del grafo attraverso l'algoritmo basato sulla topologia

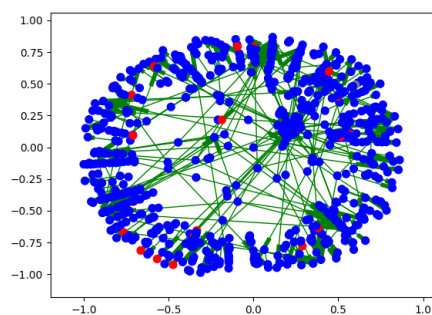


Figura 5.26: Settembre

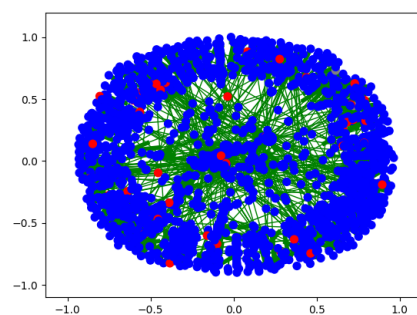


Figura 5.27: Ottobre

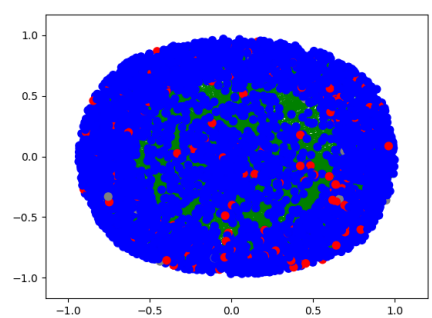


Figura 5.28: Novembre

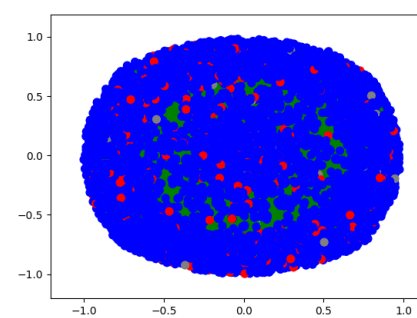


Figura 5.29: 14Dicembre

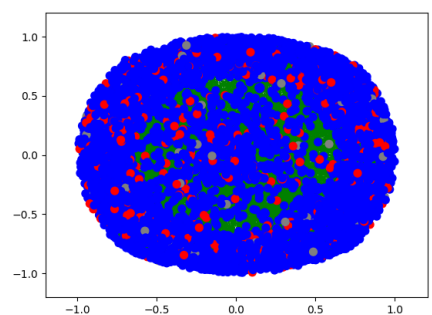


Figura 5.30: Dicembre

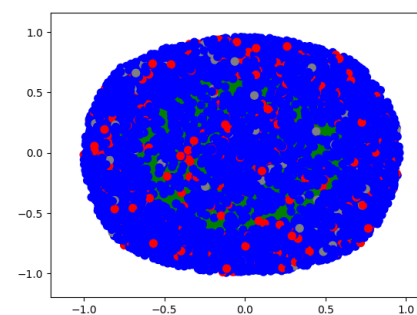


Figura 5.31: Gennaio

Figura 5.32: Popolazione del grafo attraverso l'algoritmo basato sul grado del nodo

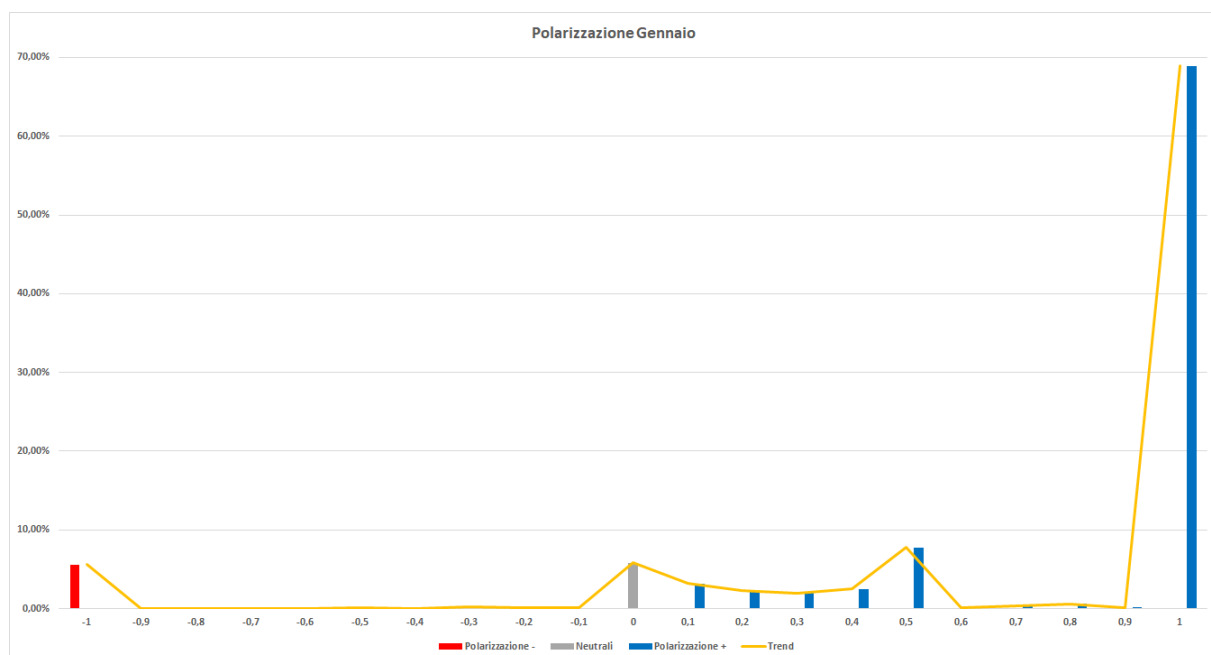


Figura 5.33: Risultati Gennaio algoritmo basato sulla topologia

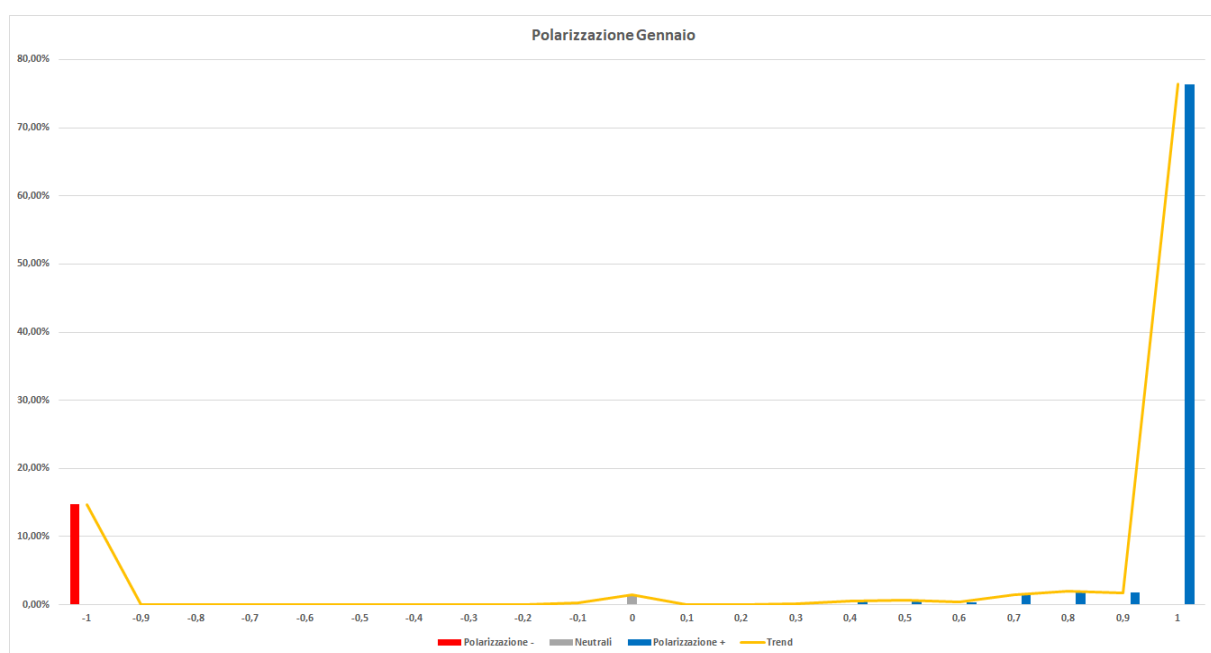


Figura 5.34: Risultati Gennaio algoritmo basato sul grado del nodo

mese di Dicembre, il giorno 14 Novembre la legge sul biotestamento è stata approvata dallo stato italiano. Nel grafico (vedi fig:5.35) sono stati confrontati i risultati ottenuti da Settembre fino a Dicembre e quelli ottenuti da Settembre fino al 14 Dicembre. Risulta evidente come i tweet collezionati siano fortemente polarizzati verso un senso unico e cioè verso i contrari all'attuazione della legge. Come riscontrabile nel diagramma è possibile notare come dopo la legge ci sia stata una lieve crescita da parte degli utenti favorevoli. Tali risultati possono descrivere come la sfera religiosa possa interferire nel calcolo della polarizzazione; basti pensare che all'interno dei tweet raccolti fossero presenti molti esponenti del mondo cattolico. A tal proposito la Chiesa Cattolica attraverso il *CEI*, ha invitato tutti i medici a fare obiezione di coscienza, nel caso si verificassero casi di biotestamento. Tali affermazioni hanno un gran seguito nella comunità religiosa italiana, essendo lo stato italiano prevalentemente cattolico; per cui possono in qualche modo drogare i risultati della polarizzazione.[2] Essendo la polarizzazione uno studio parziale questo ha permesso di comprendere, studiando il mese di Dicembre, quale fosse all'interno della rete il pensiero degli utenti. Tutto ciò dimostra come anche una religione possa influenzare le masse e quindi anche i social media. La dimostrazione di una forte polarizzazione genera la formazione di un gruppo molto chiuso di utenti, un *Echo-chambers*. Si è deciso di studiare il comportamento della polarizzazione durante un periodo successivo a Dicembre per verificare se questo argomento fosse di pubblico interesse anche successivamente l'applicazione della legge; come possibile notare dalla tabella: 5.4, è facile notare come da Dicembre a gennaio il numero di utenti che hanno continuato a pubblicare tweet e a dibattere su quel topic sono molti infatti notiamo una crescita dei dati pari al $\simeq 24\%$ in un solo mese.

Durante lo sviluppo della tesi, analizzando i dati ottenuti, si è deciso di

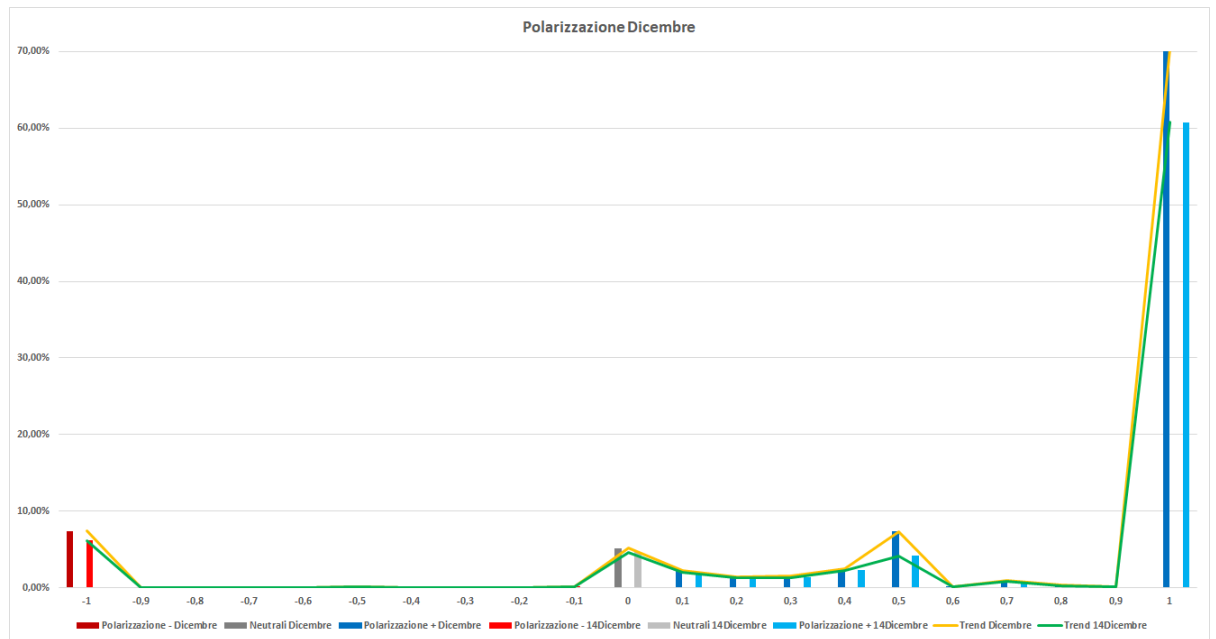


Figura 5.35: Risultati Dicembre algoritmo basato sulla topologia

realizzare un principio di predizione della polarizzazione in mesi successivi utilizzando i dati calcolati nei periodi precedenti. Il motivo di tale scelta è dettata dalla necessità di poter intuire il comportamento della rete, quindi capire il partizionamento delle opinioni degli utenti, nel futuro. Le tecniche utilizzate per la la predizione sono state illustrate all'interno della sezione: 4.5. La necessità di fare uno studio del genere è dettata dalla continua crescita degli utenti che pubblicano tweet sul biotestamento. Quindi è stato fatta un'analisi per poter garantire in futuro una possibile predizione delle polarizzazione nei mesi successivi. Durante la fase di test di queste tecniche di predizione, sono state fatte alcune considerazioni:

- In caso di una singola osservazione, per predire il futuro, non è possibile applicare tale algoritmo per cui si è scelto di restituire come valore il valore della precedente osservazione.
- Il numero minimo di osservazioni per effettuare il calcolo con questo

algoritmo deve essere pari a 2.

I test in questione per il topic analizzato sono stati fatti utilizzando i dati raccolti; cioè sono stati utilizzati i risultati della polarizzazione nei mesi di Settembre Ottobre Novembre e Dicembre per poter predire Gennaio; avendo così modo di verificare la bontà dello studio effettuato. Una volta ottenuti i risultati si è calcolato l'errore medio e la relativa varianza in base ai risultati conseguiti nel calcolo della polarizzazione attraverso i due algoritmi implementati. Di seguito mostriamo i risultati conseguiti attraverso un diagramma (vedi fig:5.37).

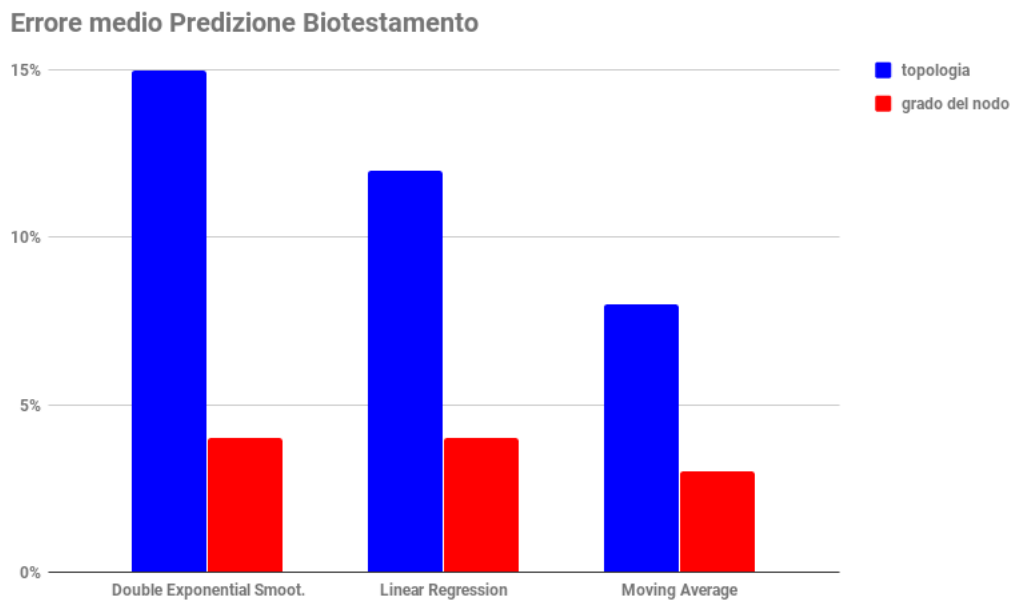


Figura 5.36: Errore medio predizione

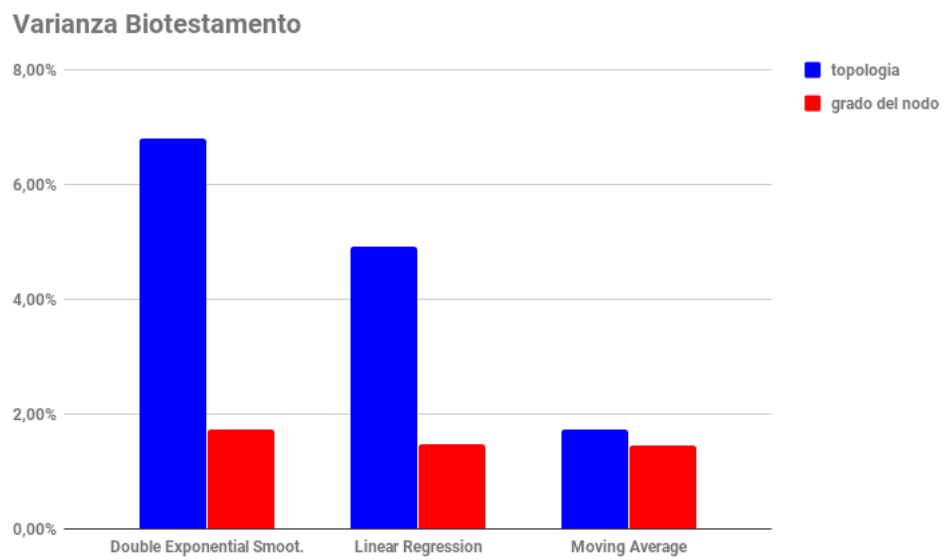


Figura 5.37: Varianza tasso di errore predizione

Capitolo 6

Sviluppi futuri e conclusioni

La polarizzazione è uno strumento che può essere utilizzato all'interno dei social media per poter analizzare alcuni comportamenti della rete a seconda del topic analizzato. Nel dettaglio attraverso i test è stata verificata la presenza di numerose comunità di utenti aventi una forte polarizzazione verso un'unica visione. Attraverso questo strumento è possibile identificare questo tipo di comportamento perché in caso di reiterazione nel tempo può portare alla formazione di *Echo-Chambers*. Una problematica che sta sempre più coinvolgendo gli amministratori dei social media. Ci sono diversi meccanismi per poter impedire la formazione di questi gruppi isolati. Uno di questi è quello di cercare di abbassare il livello di controversia del grafo nel tempo. Il modo per farlo è quello di individuare i due gruppi fortemente polarizzati e cercare di farli comunicare tra di loro attraverso un congiungimento, cioè un arco. Una soluzione sarebbe quella di suggerire ad utente che si trova all'interno di una rete fortemente polarizzata, di aggiungere, condividere o anche semplicemente mostrare dei post aventi una polarizzazione opposta a quella in cui ci si trova. Ovviamente non è un approccio vincente al 100% perché molto spesso dipende dalla volontà dell'utente di incuriosirsi delle

opinioni avverse. Però attraverso la polarizzazione si è potuto notare come molto spesso fossero presenti nodi grigi che risultavano imparziali verso i topic studiati, oppure come ci fossero anche nodi che non erano completamente polarizzati.

Un possibile impiego della polarizzazione potrebbe essere quello delle indagini di mercato all'interno dei social media, basterebbe che un'azienda lanciasse un *hashtag* (nel caso di *Twitter*), e studiasse che impatto abbia sulla rete. Le informazioni ottenute potrebbero facilmente essere utilizzate per intuire l'interesse del mercato di un nuovo prodotto.

Come dimostrato attraverso lo studio ed i test effettuati attraverso le *elezioni regionali in Sicilia*, è possibile scoprire la presa che ha un partito sulla rete piuttosto che un altro, in questo modo sarebbe fattibile fare delle previsioni sul risultato. La previsione della polarizzazione potrebbe essere utilizzata in tutti i precedenti casi illustrati all'interno di questa sezione proprio perché consente a chiunque di poter in parte conoscere quanto un certo argomento sarà oggetto di dibattito, e quindi capirne la presa sugli utenti.

In conclusione la polarizzazione attraverso un processo di analisi in profondità come la *Sentiment Analysis* consente di comprendere al meglio la contrapposizione di idee in una rete sociale indipendentemente dal contesto di studio, sia che sia politico, sociale, etico ed economico. Raffinare lo studio sui social media mediante algoritmi che non analizzino solo il testo dei messaggi pubblicati, ma che analizzino:

- immagini
- media vari
- link a siti internet

consentirebbe uno studio più preciso sulla polarizzazione e magari riuscirebbe anche a risolvere il problema dell'ironia che una tecnica come la sentiment analysis non sempre è in grado di classificare nella maniera più corretta.

Ringraziamenti

- Un ringraziamento speciale ai miei genitori, due persone speciali, a cui voglio dedicare tutto il mio lavoro, mi avete permesso di studiare, mi avete supportato e sopportato in tutti questi anni. Mi avete spronato ogni giorno ed insegnato che nella vita le difficoltà vanno sempre affrontate a testa alta. Cercherò di affrontare la vita ed il mondo del lavoro nella stessa maniera con cui ho affrontato questo periodo utilizzando tutti i vostri insegnamenti. Non potrò mai ringraziarvi abbastanza, spero di avervi reso orgogliosi di me per il traguardo che ho appena raggiunto.
- A mio fratello per avermi insegnato che bisogna vivere la vita al massimo senza perdersi d'animo alle prime difficoltà.
- Ai miei nonni per tutte le piacevoli telefonate, chiacchierate e momenti di leggerezza trascorsi durante questo periodo.
- A Stefano e Paolo, con cui ho condiviso tutti gli anni della magistratura, per le ore trascorse insieme ai dibattiti e le prese in giro mentre facevamo i progetti.
- A tutti i miei amici dell'università presenti e passati per le battute gli scherzi e le gioie che abbiamo trascorso insieme.

- A tutti i miei amici per il calore e l'affetto dimostrato.
- Al prof. Italiano per il supporto e l'aiuto nella realizzazione di questo lavoro di tesi

Bibliografia

- [1] Biotestamento. http://www.repubblica.it/cronaca/2018/01/31/news/bioestamento_da_oggi_i_desideri_dei_malati_sono_legge-187697062/.
- [2] Biotestamento. <http://www.famigliacristiana.it/articolo/legge-sulle-dat-un-occasione-persa.aspx>.
- [3] Double exponential smoothing. <http://www.itl.nist.gov/div898/handbook/pmc/section4/pmc433.html>.
- [4] Echo-chambers. [https://en.wikipedia.org/wiki/Echo_chamber_\(media\)](https://en.wikipedia.org/wiki/Echo_chamber_(media)).
- [5] Forecasting. <https://en.wikipedia.org/wiki/Forecasting>.
- [6] Linear regression. http://www.ecostat.unical.it/Costanzo/Didattica/Probabilit%C3%A0%20ed%20Inferenza%20Statistica/Prob_4.pdf.
- [7] Moving average. <https://www.thebalance.com/simple-exponential-and-weighted-moving-averages-1031196>.
- [8] Naive bayes. https://en.wikipedia.org/wiki/Naive_Bayes_classifier.

- [9] Networkx.
<https://networkx.github.io/documentation/stable/index.html>.
- [10] Risultati elezioni regionali sicilia 2017. <http://www.repubblica.it/static/speciale/2017/elezioni/regionali/sicilia.html>.
- [11] Sentiment analysis.
<https://www.celi.it/soluzioni/opinion-mining/>.
- [12] Sondaggi elezioni regionali sicilia 2017. <http://www.sondaggipoliticoelettorali.it/GestioneSondaggio.aspx>.
- [13] J. C. Losada A. J. Morales, J. Borondo and R. M. Benito. Measuring political polarization: Twitter shows the two sides of venezuela. May 2015.
- [14] Lillian Lee Bo Pang. Opinion mining and sentiment analysis.
- [15] Aristides Gionis e Michael Mathioudakis Kiran Garimella, Gianmarco De Francisci Morales. Reducing controversy by connecting opposing views.
- [16] Kevin P. Murphy. Naive bayes classifiers.
- [17] Robert Kleinberg Pedro H. Calais Guerra, Wagner Meira Jr. Claire Cardie. A measure of polarization on social media networksbased on community boundaries. 2015.