

PROGETTO RETI GEOGRAFICHE 2021/2022



Twitter Sentiment
Analysis during italian
TV show: "Amici 2021"

Docente:

*Delfina Malandrino
Rocco Zaccagnino*

Componenti gruppo:

*Antonio Zizzari
Vincenzo Marrazzo
Simone Giglio*



Sommario

- 1. Abstract 3
- 2. Introduzione 4
- 3. Modalità di Analisi 5
- 4. Progettazione e Implementazione..... 6
 - Twitter Capture6
 - Sentiment Analysis9
 - Lemmatizzatore 10
 - Gestione dei dati 11
- 5. Analisi in R dei dati Raccolti 12
 - Indici di Sintesi e Grafici dei Dati 14
 - Regressione lineare 15
 - Regressione lineare semplice 15
- 6. Conclusioni..... 21
- 7. Riferimenti 21

1. Abstract

Ogni giorno all'interno di Twitter avvengono moltissimi tweets, ed essi giocano un ruolo fondamentale all'interno degli show televisivi, poiché il voto del pubblico ha un'importanza molto rilevante.

Essendo che le opinioni dei telespettatori su internet sono in qualche modo un indice di gradimento essenziale, valutare se un'opinione è positiva o negativa su un concorrente di "Amici 2021" è l'obiettivo del nostro progetto.

Per una panoramica generale sull'indice di gradimento dei concorrenti se si andasse a identificare manualmente i tweets positivi o negativi risulterebbe un lavoro laborioso se non impossibile. È per questo motivo che attraverso tecniche per l'analisi della linguistica, frutto anche di un lavoro di psicologia, di lemmatizzazione e di machine-learning, si va a sviluppare un classificatore che determini se un tweet su un concorrente è positivo, tendente positivo, neutro, tendente negativo o negativo, cercando di ottenere un'accuratezza quanto prossima al 100%.

2. Introduzione

Con la crescita della popolarità del Web, i post e le opinioni degli utenti hanno un impatto considerevole sull'immagine dello show televisivo, e ovviamente sono frutto di un potenziale guadagno monetario.

Le opinioni derivanti dai tweets tramite la scelta accurata di hashtags mirati al tipo di analisi che si vuole effettuare, sono fondamentali per l'individuazione dei gusti degli utenti nei confronti dei concorrenti del programma. Una scelta sbagliata degli hashtags da ricercare potrebbe portare ad un'analisi non coerente oppure con risultati scarni.

Lo scopo del progetto è comprendere attraverso dei classificatori se un commento è positivo, tendente positivo, neutro, tendente negativo o negativo, utilizzando le conoscenze acquisite al corso di "Reti Geografiche". Il motivo per il quale è stato scelto questo argomento è dovuto alla scarsa quantità di elaborati relativi a tale ambito.

3. Twitter vs Instagram

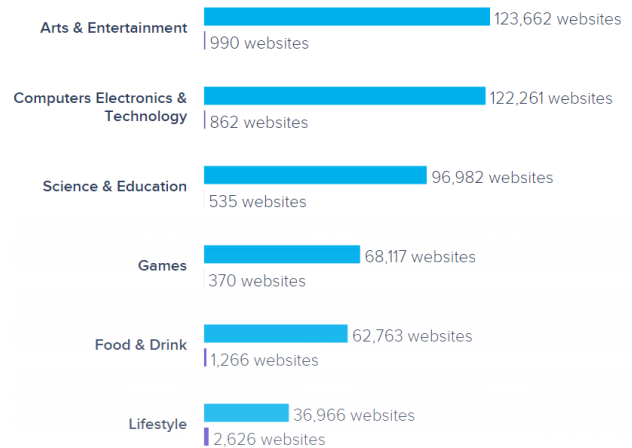




Twitter ha una migliore copertura dell'utilizzo in più categorie di siti web. Tra cui Arte e intrattenimento, Elettronica e tecnologia per computer, Scienza e istruzione, Giochi e altre 20 categorie.



L'API di **Instagram** non ha un vantaggio su Twitter in nessuna categoria di siti Web.



È stato scelto, quindi, Twitter come base per la cattura dei commenti che riguardano il mondo di Amici2021 perché secondo i dati illustrati precedentemente Twitter viene utilizzato maggiormente in diversi categorie di siti.

4. Modalità di Analisi

In questa fase lo studio relativo a quali tecnologie, a quale social media, a come raccogliere i dati e l'analisi statistica di questi ultimi, è significativo ai fini di un'analisi corretta che possa dare dei risultati ottimali per il nostro scopo.

Il social media preso in considerazione è stato Twitter; la scelta di tale social è dipesa dal fatto che Twitter mette a disposizione delle API gratuite e semplici da usare, a differenza degli altri social media.

Il linguaggio di programmazione scelto per la cattura, analisi e valutazione è stato Python, in quanto risulta vantaggioso per la semplice implementazione e supporto alla moltitudine di librerie e modelli per il machine-learning, in particolare le librerie Tweepy, Vader e Spacy (progetti open source) ci hanno permesso di elaborare in modo semplificato ed efficace i dati catturati.

Per quanto riguarda la gestione dei dati raccolti, si è optato per il formato “.csv”, che ci ha permesso di strutturare ed immagazzinare dati, per poi importare tali file di questo formato sull’ambiente di sviluppo RStudio.

Tramite RStudio, si fanno analisi e statistiche dei dati immagazzinati per evidenziare eventuali situazioni non facilmente intuibili con la semplice visualizzazione dei dati, ma attraverso strumenti grafici fatti ad hoc per il tipo di problema analizzato.

5. Progettazione e Implementazione

In questa fase analizzeremo come e con quali strumenti è stato implementato il nostro problema.

Twitter Capture

Il primo problema che ci siamo posti è stato: “come catturare i tweets?”. Tra i vari tool analizzati, l’uso delle API di Twitter è stato di particolare importanza, infatti l’utilizzo della libreria *tweepy* ci ha permesso di catturare tutti i post su Twitter attraverso delle query specifiche.

Queste API hanno dei limiti imposti da Twitter che sono i seguenti:

Elevated		View detailed features
Apps	3 environments per project	
Tweets	2M Tweets per month / Project	
Cost	free	

Infatti il limite risulta essere di 2mln di catture al mese, e le catture sono limitate anche temporalmente, infatti il limite massimo è di 1 settimana antecedente al giorno corrente.

Come primo approccio abbiamo creato le variabili alle quali sono state associate le credenziali delle API:

```
1 # Keys e Tokens per accedere a Twitter
2 consumer_key = 
3 consumer_secret = 
4 
5 access_token = 
6 access_token_secret = 
7
```

In un secondo momento, abbiamo scritto la query che ci ha permesso di catturare i tweets. I criteri di ricerca hanno seguito delle regole ben precise, che sono:

- **Hashtags:** Abbiamo creato uno snippet di codice che prende in input i concorrenti aggiornati in modo dinamico e li mette all'interno della query di *tweepy*, garantendo una maggiore precisione e accuratezza sulla cattura dei tweets.

```
20 # Inserimento Hashtag da cercare
21 inputTerm = "#Amici21 AND ("
22 for i in range(len(concorrenti_amici)):
23     if (i < (len(concorrenti_amici)-1)):
24         inputTerm += "#" + concorrenti_amici[i] + " OR "
25     elif (i == (len(concorrenti_amici)-1)):
26         inputTerm += "#" + concorrenti_amici[i] + ")"
```

- **Retweets:** Abbiamo impostato il parametro `-filter:retweets` per filtrare dalla ricerca i retweet, per evitare di avere tweets ridondanti;
- **Giorno limite cattura tweet:** Abbiamo usato il parametro `until` che ci ha permesso di settare il limite temporale per la cattura di un tweet. Nella fase di cattura è importante specificare che la query viene eseguita sui tweet antecedenti alla data in input fino ad una settimana, limiti imposti dalle API;
- **Numero max tweets:** Uno dei limiti che ci impone l'API di Twitter è il numero massimo di tuple che possiamo ottenere da una query. Il limite è fissato a 100 per cattura, anche se si andasse ad aumentare tale limite, la cattura massima rimarrebbe invariata.

Dopo queste analisi preliminari si procede con l'effettiva cattura dei tweets mediante il seguente codice Python:

```
34 for tweet in tweepy.Cursor(api.search_tweets,q=searchTerm, lang="it",
35                             until = date_tweepy,
36                             result_type="mixed",
37                             tweet_mode="extended").items(100):
38     rt_count = tweet.retweet_count
39     p2 = ["", tweet.full_text]
40     hast = ""
41     a = tweet.full_text.count("#")
42     if (a >= 0):
43         for c in range(a):
44             p1 = p2[1].split("#", 1)
45             try:
46                 if (p1[1].find(" ") < 0):
47                     hast += "#" + p1[1]
48                     hast = hast.replace("\n", "")
49                     break
50                 else:
51                     p2 = p1[1].split(" ", 1)
52                     hast += "#" + p2[0]
53                     hast = hast.replace("\n", "")
54             except IndexError:
55                 hast += "#" + p1[0]
56     tweet.full_text = cleanTweet(tweet.full_text)
57     hast = cleanHashtags(hast)
59     tt = {
60         "Id": tweet.id,
61         "Screen_name": tweet.user.name,
62         "Created_at": tweet.created_at,
63         "Retweet": str(rt_count),
64         "Text": tweet.full_text,
65         "Hashtags": hast,
66         "Sentiment": 0,
67         "Compound": 0
68     }
69     tt = Tweet(tweet.id, tweet.user.name, tweet.created_at, str(rt_count), tweet.full_text, hast, 0, 0)
70     if (str(tt.getCreated_at()).find(date_to_search) != -1):
71         tweets.append(tt)
72         tweetsTranslated.append(translatorTweet(tt))
```

In questo snippet di codice è stato possibile ricavare la tupla del tweet, con i criteri sopra citati. Il testo del tweet prima di esser salvato viene elaborato tramite un processo di filtraggio. In questo processo vengono rimosse dal tweet, attraverso i metodi `cleanHashtags()` e `cleanTweet()`, le componenti non adatte come i simboli, i caratteri speciali, stringhe http, ecc..., in modo da ottenere un testo pulito.

Infine, si crea un dizionario che raccoglie tutti gli elementi principali dei tweets e successivamente li si aggiunge in una lista.

Dato che la libreria per i sentiment funziona solamente con il dizionario inglese, per prevenire tale situazione, si traduce in anticipo il tweet, con la funzione `translatorTweet()`.

Una volta ottenuta la nostra lista di tweets tradotti, procediamo con il Sentiment Analysis.

Sentiment Analysis

Una volta eseguita la traduzione del messaggio del tweet, nello stesso snippet, viene avviata la procedura di sentiment, tramite la funzione `analyzer.polarity_scores()`, si fa uso della libreria Vader per ottenere un parametro in output che corrisponde al sentimento del testo passato; successivamente il valore ottenuto in output viene passato al metodo `sentiment()` e tale funzione permette di fornire come output i punteggi relativi a 5 classi di sentiment:

1. Negativo;
2. Tendente Negativo;
3. Neutro;
4. Tendente Positivo;
5. Positivo.

Queste classi vengono restituite e salvate all'interno di un parametro appartenente al dizionario che si stava in precedenza strutturando.

I codici in Python che permettono tale analisi, sono:

```

62 def calculate_and_set_compound_score_to_tweet(tweet):
63     analyzer = SentimentIntensityAnalyzer()
64     sentiment = analyzer.polarity_scores(tweet.getText())
65     compound = sentiment["compound"]
66     tweet.setCompound(compound)
67
68 def calculate_and_set_sentiment_score_to_tweet(tweet):
69     sentimentAnalysis = Sentiment(tweet.getCompound())
70     tweet.setSentiment(sentimentAnalysis[0])
71
72 def print_emotion_score(tweet):
73     sentimentAnalysis = Sentiment(tweet.getCompound())
74     print("\033[1m" + "Risultato: " + "\033[0m" + sentimentAnalysis[1])
75
76 def Sentiment(compound):
77     emotion = ""
78     if(compound <= -0.5501 and compound >= -1):
79         sentimentAnalysis = 1
80         emotion = "Negativo"
81     if(compound <= -0.2001 and compound >= -0.5500):
82         sentimentAnalysis = 2
83         emotion = "Tendente negativo"
84     if(compound >= -0.2000 and compound <= 0.2000):
85         sentimentAnalysis = 3
86         emotion = "Neutro"
87     if(compound >= 0.2001 and compound <= 0.5500):
88         sentimentAnalysis = 4
89         emotion = "Tendente positivo"
90     if(compound >= 0.5501 and compound <= 1):
91         sentimentAnalysis = 5
92         emotion = "Positivo"
93
94     return sentimentAnalysis, emotion
95

```

Lemmatizzatore

Adesso si passa alla 2ª parte del progetto, che include l'analisi grammaticale del tweet passato. Attraverso la libreria Spacy, è possibile avere un'analisi grammaticale del testo passato abbastanza accurata dovuto dal modello di machine-learning che ha una accuratezza seguente:

NAMED ENTITY RECOGNITION SYSTEM	ONTONOTES	CONLL '03
spaCy RoBERTa (2020)	89.8	91.6
Stanza (StanfordNLP) ¹	88.8	92.1
Flair ²	89.7	93.1

Poiché da essa è possibile ricavare il soggetto o i soggetti del testo in questione, è stato possibile associare uno score ai concorrenti del programma televisivo.

Il modello, durante l'analisi, è stato allenato, per renderlo ad hoc riguardo il problema analizzato (sfruttando la tecnologia CUDA di NVIDIA).

Lo snippet di codice che ha permesso tale implementazione è il seguente:

```
71 # SECONDA PARTE DEL CODICE
72 print("Analisi più dettagliata dei Tweet")
73 nlp = spacy.load("it_core_news_sm")
74
75 print(concorrenti_amici)
76 for tweet in tweets:
77     amici_objects_clone = copy.deepcopy(concorrenti_amici_objects)
78     doc = nlp(tweet.getText())
79     print("TWEET: ")
80     print("Testo tweet: " + tweet.getText())
81     print("Sentiment: " + str(tweet.getSentiment()))
82
83     ents = []
84     for entity in doc.ents:
85         for concorrente in concorrenti_amici_objects:
86             # non mettiamo il break, a fine 'if', poichè all'interno di un soggetto potrebbero esserci più
87             # concorrenti e quindi vogliamo
88             # prima estrapolare tutti i concorrenti dal soggetto e poi andare avanti col prossimo soggetto
89             if ((concorrente.getName().upper() in entity.text.upper()) and (concorrente in amici_objects_clone)):
90                 ents.append(concorrente.getName())
91                 print("text: " + entity.text, ", concorrente_name: " + concorrente.getName())
92                 concorrente.addScore(tweet.getSentiment())
93                 amici_objects_clone.pop(amici_objects_clone.index(concorrente))
94                 #print("FINE " + "text: " + entity.text, "type/label: " + entity.label_)
95
96     tweet.setEnts(ents)
97     print("Ents: " + str(tweet.getEnts()))
98     print()
```

Gestione dei dati

I dati dei concorrenti e dei tweets sono stati immagazzinati all'interno di un file formato ".csv". Esso ci servirà più avanti, per effettuare l'elaborazione dei dati attraverso il linguaggio R.

Lo snippet di codice è il seguente:

```
4 # Nomi dei file
5 name_file_tweets = "./Dataset_tweets/Amici(15-12-2021).csv"
6 name_file_concorrenti_score = "./Dataset_tweets/Amici_Score_Concorrenti(15-12-2021).csv"
7
8 # 1 parte
9 print(f"\nWriting in file '{name_file_tweets}' ...")
10 # Creazione nostro file
11 csv_file = open(name_file_tweets, 'w', encoding='utf-8', newline='')
12 writer = csv.writer(csv_file)
13 head1 = ["Id", "Screen_name", "Created_at", "Retweet", "Text", "Hashtags", "Sentiment", "Compound", "Ents"]
14 writer.writerow(head1)
15
16 # Aggiunta dei tweets nel file CSV
17 for t in tweets:
18     lista_ents = converto_in_list(t.getEnts())
19     writer.writerow([str(t.getId()), str(t.getScreen_name()), str(t.getCreated_at()), str(t.getRetweet()), str(t.getText()),
20                     str(t.getHashtags()), str(t.getSentiment()), str(t.getCompound()), lista_ents])
21
22 csv_file.close()
```

I file che creiamo sono:

- ❖ i dati relativi agli score dei singoli concorrenti;
- ❖ i tweet con le relative informazioni, suddivisi per giorni del mese

6. Analisi in R dei dati Raccolti

Prima di procedere con l'analisi vera e propria i dati sono stati elaborati ulteriormente.

L'elaborazione dei dati ha incluso vari step quali:

- Si è raccolto tutti i file .csv delle catture giornaliere e sono stati uniti in un unico file contenente tutti i tweets (1 mese di cattura) denominato "File_finale.csv" attraverso il seguente snippet di codice:

```
1 import csv
2 import os
3
4 array_stampa=[]
5
6 # Cattura delle righe dei csv dei tweets 2021
7 folder_path_tweets_2021 = os.getcwd() + "\\tweets (2021)"
8 for filename in os.listdir(folder_path_tweets_2021):
9     if (filename=="desktop.ini"):
10         continue
11     print(f"File name: {filename}")
12     with open(folder_path_tweets_2021 + "\\" + filename, 'r', encoding='utf-8', newline='') as csv_file:
13         csv_reader = csv.reader(csv_file, delimiter=",", quotechar='"')
14         next(csv_reader)
15         for row in csv_reader:
16             array_stampa.append(row)
17     csv_file.close()
18
19 # Cattura delle righe dei csv dei tweets 2022
20 folder_path_tweets_2021 = os.getcwd() + "\\tweets (2022)"
21 for filename in os.listdir(folder_path_tweets_2021):
22     if (filename=="desktop.ini"):
23         continue
24     print(f"File name: {filename}")
25     with open(folder_path_tweets_2021 + "\\" + filename, 'r', encoding='utf-8', newline='') as csv_file:
26         csv_reader = csv.reader(csv_file, delimiter=",", quotechar='"')
27         next(csv_reader)
28         for row in csv_reader:
29             array_stampa.append(row)
30     csv_file.close()
31
32 # Scrittura di tutti i tweets nel file 'File_finale.csv'
33 csv_file = open("File_finale.csv", 'w', encoding='utf-8', newline='')
34 writer = csv.writer(csv_file)
35 head = ["Id", "Screen_name", "Created_at", "Retweet", "Text", "Hashtags", "Sentiment", "Compound", "Ents"]
36 writer.writerow(head)
37 writer.writerows(array_stampa)
38 csv_file.close()
```

Questa unione ci ha permesso di ottenere gli scores complessivi di 1 mese di cattura di tutti i concorrenti lanciando nuovamente l'algoritmo di sentiment analysis e di lemmatizzazione su di esso.

- Creazione di un metodo di ranking per la valutazione degli sfidanti di Amici volto a classificarli in base a come ha reagito il pubblico di Twitter nei loro confronti.

```
56 # Metodi d'appoggio
57 def calculate_rank(self):
58     rank = 0
59     rank += (self.getScoreNegativo() * -1)
60     rank += (self.getScoreTendentePositivo() * -0.5)
61     rank += (self.getScoreNeutro() * 0)
62     rank += (self.getScoreTendentePositivo() * +0.5)
63     rank += (self.getScorePositivo() * +1)
64
65     return rank
66
```

Questo valore è ponderato, facciamo un esempio sul concorrente Alex con i seguenti valori (ottenuti tramite l'applicazione di VaderSentiment sui tweets):

	Negativo	Tend. Negativo	Neutro	Tend. Positivo	Positivo
Alex	10	2	35	23	20

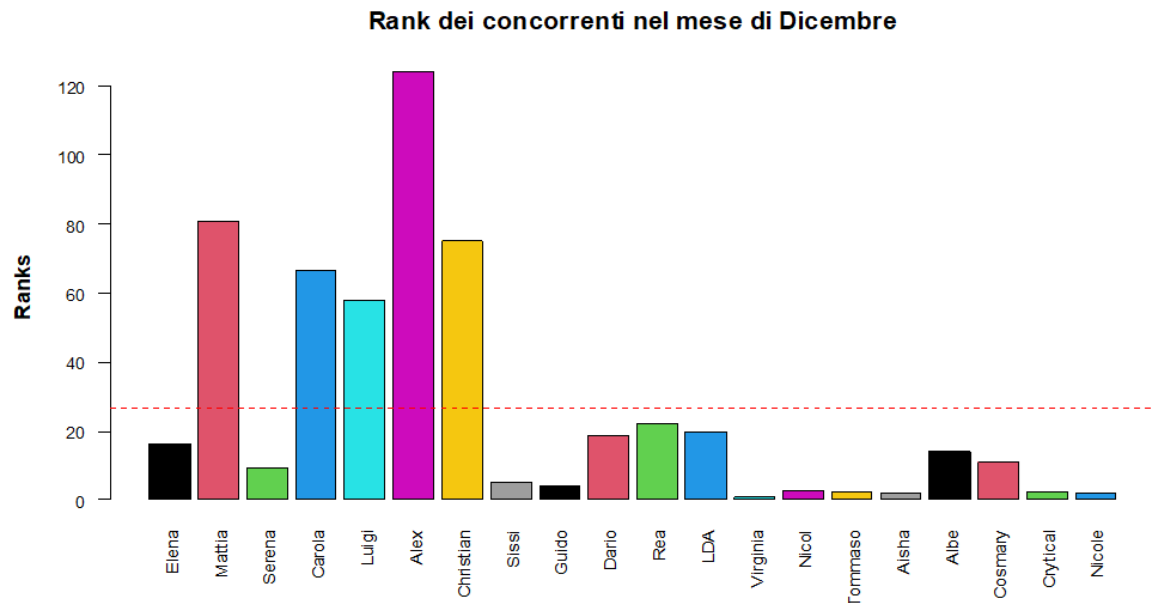
Il calcolo che viene eseguito è il seguente:

$$\begin{aligned} RANK = & (\text{Negativo} * (-1)) + (\text{Tend. Negativo} * (-0,5)) \\ & + (\text{Neutro} * (0)) + (\text{Tend. Positivo} * (+0,5)) \\ & + (\text{Positivo} * (+1)) = 20,5 \end{aligned}$$

- E la creazione di un file .csv contenente tutti i punteggi di ranking che verrà analizzata in dettaglio con R usando la retta di regressione con gli scarti complementari.

Indici di Sintesi e Grafici dei Dati

Prima di mostrare uno studio approfondito, è meglio avere una panoramica preliminare dei dati inerenti i punteggi derivati dall'algoritmo di ranking per concorrente che andremo ad elaborare:



Il codice in R che ci ha permesso di ottenere il seguente grafico è:

```
1 createBarplot<-function(ranks, main){
2
3   x<-barplot(ranks, ylab="Ranks", main=main,
4             col=1:40,
5             las=2,
6             names.arg=`amici_rank_concorrenti(TOTALE)`[,1:1],
7             cex.axis=0.80,
8             cex.names=0.80,
9             font.lab=2)
10  abline(h=mean(ranks), col="red", lty=2, xpd=FALSE)
11 }
12
13 createBarplot(`amici_rank_concorrenti(TOTALE)`[,2:2], "Rank dei concorrenti
14               nel mese di Dicembre")
```

Come si può notare vi è una linea rossa tratteggiata che taglia il grafico, essa rappresenta la media, e ci permette di capire fin da subito i partecipanti allo show televisivo che hanno un punteggio ben superiore alla media. Infatti notiamo fin da subito che i concorrenti: Alex, Mattia, Christian, Carola e Luigi sono quelli con i punteggi(Rank) più alti. Successivamente, andremo a studiare i concorrenti che hanno ottenuto un punteggio superiore alla media.

Gli indici di sintesi dei nostri dati sono i seguenti:

Indice	Valore
Media	26.875
Mediana	12.5
Varianza	1201.181
Deviazione standard	34.65806

Tramite essi descriviamo:

- La **media** rappresentante la media aritmetica dei dati;
- La **mediana** che ci aiuta nella statistica descrittiva a capire se i dati sono simmetrici (quando media e mediana coincidono) o meno
- La **varianza** per valutare se le osservazioni analizzate sono dipendenti da uno o più fattori al fine di decidere quali sono i dati più rilevanti
- La **deviazione standard** è l'indice statistico che misura la dispersione delle singole osservazioni di un fenomeno intorno alla media.

Regressione lineare

In statistica si parla di regressione lineare indicando un approccio che modella le relazioni tra una variabile detta dipendente e una o più variabili dette indipendenti.

Il caso in cui la variabile indipendente sia una sola è detto regressione lineare semplice, mentre quando ci sono più variabili indipendenti abbiamo a che fare con la regressione lineare multipla.

Regressione lineare semplice

Nella regressione lineare le relazioni sono modellate usando funzioni di predizione lineare i cui parametri del modello sono stimati dai dati. Questi modelli sono detti modelli lineari.

Il modello di regressione lineare semplice è esprimibile attraverso l'equazione di una retta che riesce meglio di qualunque altra ad interpolare la nuvola di punti.

Data dunque l'equazione

$$Y = \alpha + \beta X$$

Dove

1. α è l'intercetta, cioè indica l'ordinata con la quale la retta di regressione si interseca con l'asse delle ordinate
2. β è il coefficiente angolare, cioè indica la “pendenza” della retta.
 - i. Se positivo allora la retta di regressione è crescente
 - ii. Se negativo allora la retta di regressione è discendente
 - iii. Se nullo allora la retta di regressione è orizzontale

Questa retta viene ottenuta con il metodo dei minimi quadrati. Per calcolare i coefficienti di regressione è necessario considerare la **somma Q dei quadrati degli errori**

$$Q = \sum_{i=1}^n [y_i - (\alpha + \beta x_i)]^2$$

e minimizzarla. Quindi la variabile dipendente y viene “spiegata” attraverso una **relazione lineare della variabile indipendente x** (cioè: $\alpha + \beta x$).

Il problema della regressione si traduce nella determinazione di α e β in modo da esprimere al meglio la relazione funzionale tra y e x .

Derivando Q rispetto α e β e ponendo le derivate parziali ottenute a 0, il metodo dei minimi quadrati conduce a:

$$\beta = \frac{S_y}{S_x} r_{xy}, \quad a = \bar{y} - \beta \bar{x}$$

Notiamo come il coefficiente di correlazione influenza fortemente β e se quest'ultimo è 0, allora la retta è **orizzontale**.

Procediamo dunque con l'analisi delle rette di regressione e residui dei concorrenti di Amici che hanno riscosso un maggiore successo secondo le nostre analisi.

Si prendono in considerazione 2 vettori inerenti ai punteggi(rank) del mese di Dicembre: Il primo contiene i primi

15 giorni del mese e il secondo gli ultimi 15. Ed in base a ciò andremo ad analizzare 2 casi:

1. Se i dati risultano molto dispersi attorno alla retta di regressione, analizzare i fenomeni su tale concorrente risulta essere poco affidabile, poiché frutto di casualità.
2. Mostrare come all'aumentare del tempo, il giudizio del pubblico sul concorrente aumenti positivamente o negativamente.

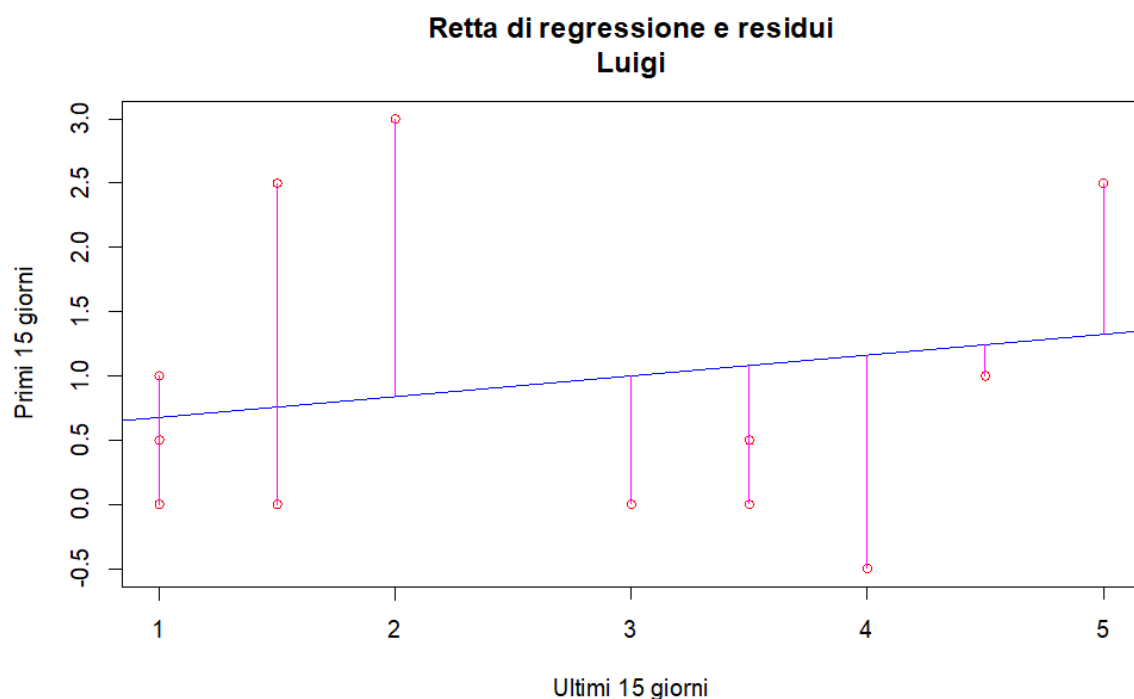
Dato il seguente codice:

```

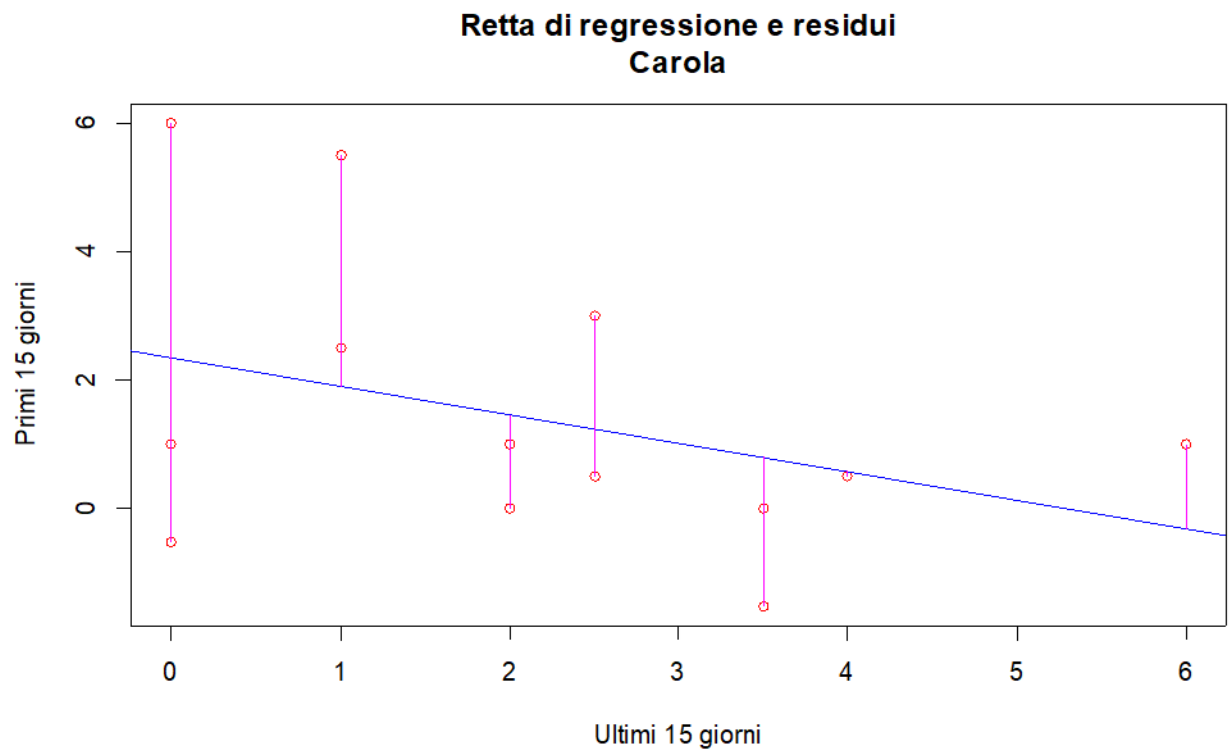
1 ▾ function_retta_di_regressione <- function(indice_concorrente, main) {
2   matrix <- as.matrix(score_days)
3   concorrente<-indice_concorrente
4   start_column<-3
5   end_column <- 32
6   first_15_rank <- c()
7   last_15_rank <- c()
8 ▾   for (start_column in start_column:end_column) {
9 ▾     if (start_column<=17) {
10       first_15_rank = c(
11         first_15_rank,
12         as.double(matrix[[concorrente:concorrente,start_column:start_column]])
13       )
14 ▾     } else {
15       last_15_rank = c(
16         last_15_rank,
17         as.double(matrix[[concorrente:concorrente,start_column:start_column]])
18       )
19 ▾     }
20     start_column = start_column + 1
21 ▾   }
22
23   x1 <- first_15_rank
24   x2 <- last_15_rank
25
26   plot(x1,x2,main=main,
27     xlab="Ultimi 15 giorni",ylab="Primi 15 giorni",
28     col="red")
29   abline(lm(x2~x1),col="blue")
30   stime<-fitted(lm(x2~x1))
31   segments(x1,stime,x1,x2,col="magenta")
32 ▾ }
33
34 function_retta_di_regressione(5, "Retta di regressione e residui\nLuigi")
35 function_retta_di_regressione(4, "Retta di regressione e residui\nCarola")
36 function_retta_di_regressione(7, "Retta di regressione e residui\nChristian")
37 function_retta_di_regressione(2, "Retta di regressione e residui\nMattia")
38 function_retta_di_regressione(6, "Retta di regressione e residui\nAlex")

```

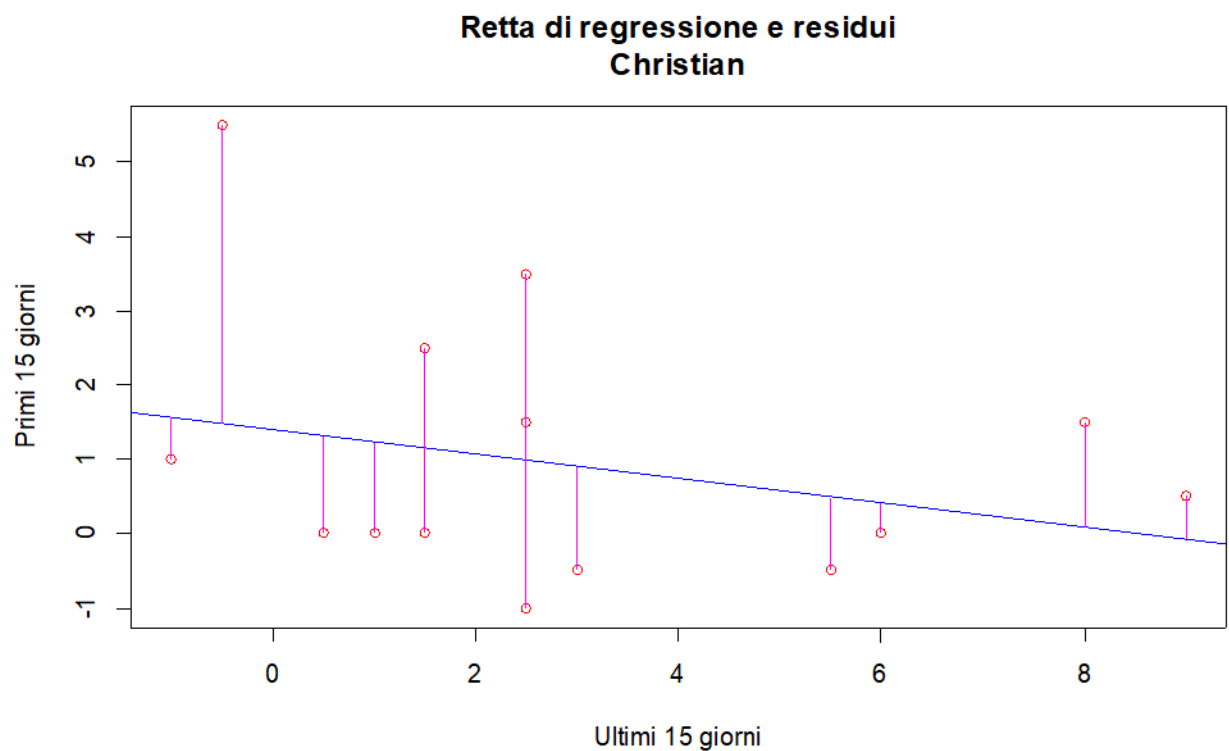
Otteniamo le seguenti stampe:



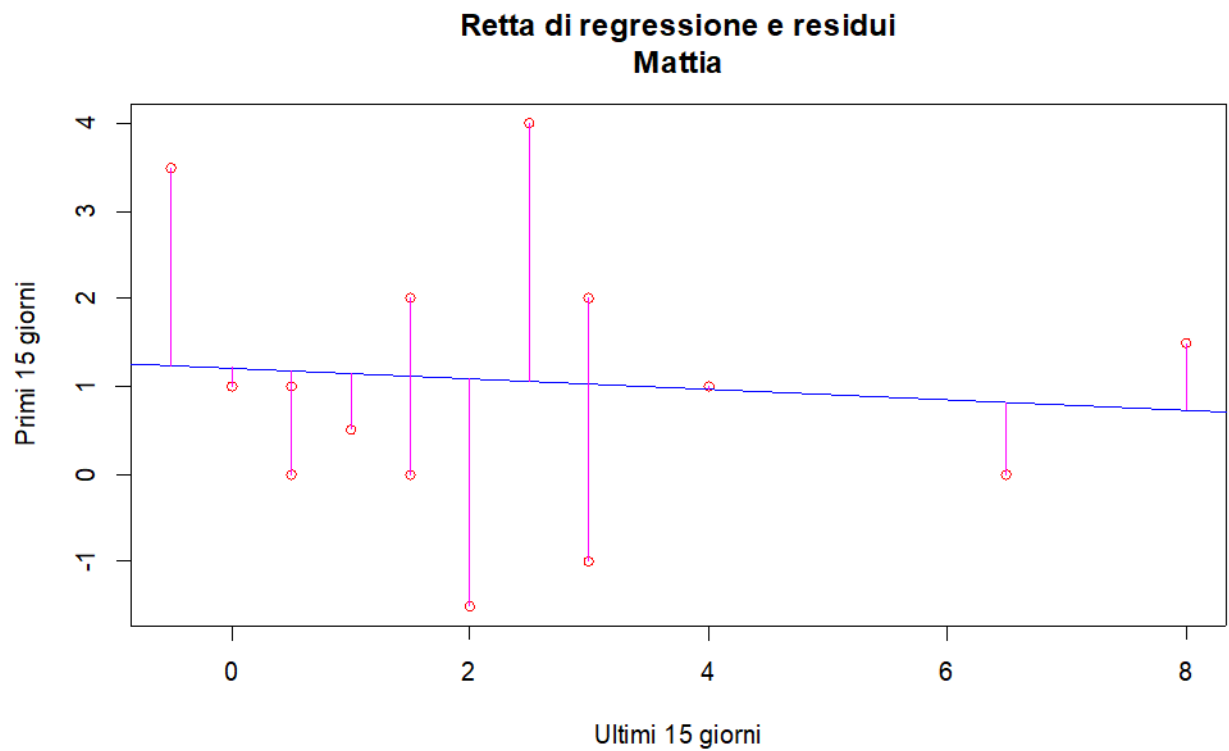
Possiamo affermare che su Twitter Luigi negli ultimi giorni del mese è stato valutato più positivamente.



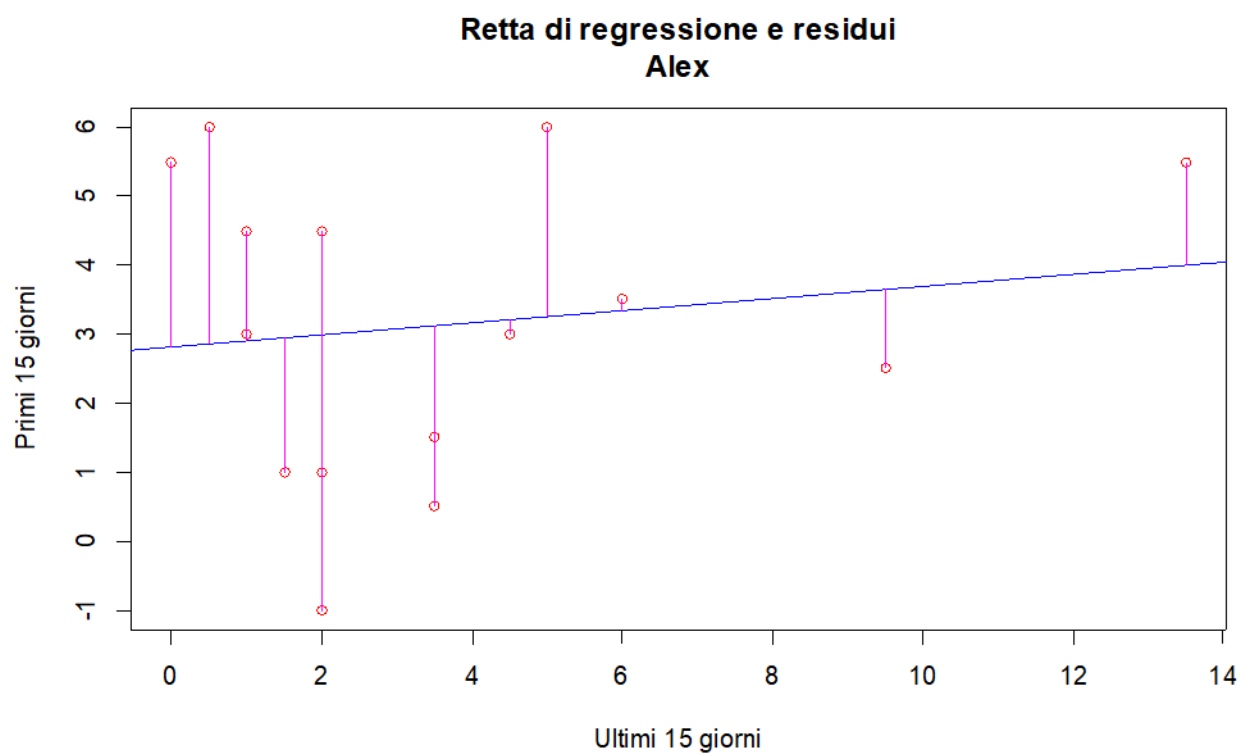
Possiamo affermare che su Twitter Carola negli ultimi giorni del mese è stata valutata negativamente.



Possiamo affermare che su Twitter Christian negli ultimi giorni del mese è stato valutato più negativamente.



Possiamo affermare che su Twitter Mattia negli ultimi giorni del mese è stato valutato leggermente più negativo.



Possiamo affermare che su Twitter Alex negli ultimi giorni del mese è stato valutato più positivamente.

7. Conclusioni

Questo progetto ci ha permesso di estendere le nostre conoscenze sul: Machine Learning, Python, Linguaggio R ed il sentiment analysis. Senza questo progetto non avremmo potuto approfondire le tematiche sopra citate.

Nonostante qualche difficoltà, è stato comunque possibile raccogliere, salvare ed analizzare dati ricavando osservazioni molto interessanti che danno una visione più ampia sui dati raccolti.

Analizzando tutti gli output, forniti dal nostro studio, ne evince come il concorrente di Amici Alex abbia un tasso di popolarità molto più elevato rispetto agli altri partecipanti del talent show, e quindi con quasi certezza finirà per entrare al serale di Amici.

Altri concorrenti che risultano avere punteggi interessanti risultano essere Mattia e Luigi.

8. Riferimenti

- Api Tweeter: <https://developer.twitter.com/en/docs>
- Vader: <https://pypi.org/project/vaderSentiment/>
- Spacy: <https://spacy.io/>
- Cuda: <https://developer.nvidia.com/cuda-downloads>
- RStudio: <https://docs.rstudio.com/>