



UNIVERSITA' DEGLI STUDI DI SALERNO

DIPARTIMENTO DI SCIENZE AZIENDALI - MANAGEMENT &  
INNOVATION SYSTEMS

CORSO DI LAUREA MAGISTRALE IN DATA SCIENCE & GESTIONE  
DELL'INNOVAZIONE

---

## Project Work PxDS

---

A.A. 2024-2025  
Docente: F. Orciuoli

**Gruppo N.1**  
**Candidati:**

DENISE BRANCACCIO - MAT. 0222800163  
LUCIA BRANDO - MAT.0222800162  
BRUNO MARIA DI MAIO - MAT.0222800149

# Indice

<b>1</b>	<b>Introduzione</b>	<b>2</b>
1.1	Project Work . . . . .	2
1.2	Definizioni utilizzate . . . . .	2
1.3	Organizzazione del gruppo . . . . .	5
1.4	Tecnologie utilizzate . . . . .	5
<b>2</b>	<b>Descrizione del codice</b>	<b>7</b>
2.1	Modelli . . . . .	8
2.1.1	Sentiment Analysis . . . . .	9
2.1.2	Emotion Analysis . . . . .	9
2.1.3	Propaganda Classifier . . . . .	9
2.1.4	Toxicity Classifier . . . . .	9
2.1.5	Word Embedding . . . . .	9
<b>3</b>	<b>Sequence Diagram</b>	<b>10</b>
3.1	Funzione def speech_info(politician,speech) . . . . .	10
3.2	Funzione def detect_propaganda_type(text, classifier=propaganda_classifier, chunk_size=512, overlap=128) . . . . .	11
3.3	Funzione def offsets(speech,data) . . . . .	12
<b>4</b>	<b>Funzioni</b>	<b>13</b>
4.1	Verify Politician . . . . .	13
4.2	Get Politician Info . . . . .	13
4.3	Speech Info . . . . .	16
4.4	Speech Data . . . . .	17
4.5	Classify Text . . . . .	17
4.6	Split Sentences . . . . .	17
4.7	Propaganda Criteria . . . . .	18
4.8	Detect Propaganda Type . . . . .	19
4.9	Contains Toxicity Criteria . . . . .	20
4.10	Classify Toxicity . . . . .	20
4.11	Offsets . . . . .	22
4.12	Formality . . . . .	23
4.13	Readablity . . . . .	23
4.14	Analyze Narrative Structure . . . . .	23
4.15	TF-IDF . . . . .	25
4.16	Word Embeddings . . . . .	25
4.17	Visualize Embeddings . . . . .	27
4.18	Analyze Persuasion . . . . .	29
<b>5</b>	<b>Analisi delle Prestazioni</b>	<b>31</b>
<b>6</b>	<b>Conclusioni</b>	<b>32</b>
<b>7</b>	<b>Collegamenti esterni</b>	<b>32</b>

# 1 Introduzione

Questa documentazione descrive il lavoro svolto per l'esame di Programming for Data Science sviluppato dal nostro gruppo di lavoro, volto all'analisi automatizzata di discorsi politici. Dopo una panoramica della traccia assegnata, vengono presentate le tecnologie utilizzate, l'organizzazione del lavoro e le ipotesi applicative del sistema.

In seguito, nei successivi capitoli, vengono illustrate le funzionalità del codice, i diagrammi di sequenza per la soluzione dei requisiti, gli esempi pratici di utilizzo e l'analisi delle prestazioni. Infine, nell'ultima sezione, vengono riportati i risultati raggiunti e i link di riferimento per la repository Github del progetto.

## 1.1 Project Work

La traccia del progetto ha richiesto lo sviluppo di una data pipeline per l'arricchimento di un dataset chiamato `speech-a.tsv`, composto da tre tipi di contenuto: nomi di politici, i loro discorsi e un valore binario ad essi associato.

L'obiettivo finale è creare un nuovo dataset che supporti l'analisi e la detection di forme di propaganda all'interno dei discorsi. L'arricchimento prevede l'integrazione di informazioni sugli autori e i loro discorsi, l'aggiunta di feature linguistiche e contenutistiche come leggibilità, sentiment, parole chiave, e narrative sottostanti, oltre alla possibilità di individuare offset e tecniche specifiche di propaganda.

## 1.2 Definizioni utilizzate

In una prima fase, il team si è occupato dello studio e della definizione di due concetti fondamentali per lo svolgimento del lavoro: 'propaganda' e 'storytelling'. Per comprendere appieno cosa sia la propaganda, ci siamo riferiti alla seguente definizione fornita dal Dizionario di Politica di N. Bobbio, N. Matteucci e G. Pasquino (Torino, Utet Libreria, 2004, p. 775):

”La propaganda può essere definita come la diffusione deliberata e sistematica di messaggi indirizzati a un determinato uditorio, con l'obiettivo di creare un'immagine positiva o negativa di determinati fenomeni (persone, movimenti, eventi, istituzioni, ecc.) e di stimolare determinati comportamenti. La propaganda è quindi uno sforzo consapevole e sistematico volto a influenzare le opinioni e le azioni di un pubblico o di un'intera società. In questo contesto, il termine 'propaganda' viene originariamente utilizzato dalla Chiesa cattolica per indicare attività di proselitismo, senza connotazioni esplicitamente negative.”

Successivamente, e' stato ampliato lo studio sulla persuasione e le tecniche di propaganda, approfondendo la lettura di testi specifici.

Lo studio condotto sulla propaganda si basa principalmente su due lavori: **'The Science of Persuasion' di Robert B. Cialdini** e **'A Survey on Computational Propaganda Detection' di G. Da San Martino, S. Cresci, A. Barron-Cede, S. Yu, R. Di Pietro e P. Nakov**.

Il primo paper esplicita come la persuasione sia la base della propaganda e si fonda su sei principi chiave:

- **Validazione sociale:**

Le persone tendono a seguire cio' che vedono fare dagli altri.

Se molte persone fanno qualcosa, e' facile pensare che sia la cosa giusta o accettabile. Tuttavia, se questo principio viene usato male, ad esempio in una campagna che promuove comportamenti negativi, puo' portare a un aumento dell'accettazione di tali comportamenti, proprio perche' altre persone li praticano.

- **Affinita':**

Le persone sono pi inclini a dire di si a chi gli piace. Quando ce' una connessione personale o un'affinita', e' pi facile che una persona accetti una richiesta. Complimenti sinceri e interazioni positive possono rafforzare questo legame, aumentando la predisposizione ad ascoltare e a rispondere favorevolmente.

- **Autorita':**

Le persone tendono a seguire i consigli di chi percepiscono come un'autorita' o un esperto. L'autorita' puo' essere riconosciuta da titoli, dall'abbigliamento o dal comportamento di una persona, e questo influisce molto sulle decisioni e le azioni degli altri.

- **Scarsita':**

Quando qualcosa e' raro o difficile da ottenere, sembra piu' prezioso. La percezione di scarsita' aumenta il desiderio di quella cosa. Questo principio e' spesso usato nel marketing per creare un senso di urgenza e spingere i consumatori a competere per un'opportunita' limitata.

- **Conoscenza:**

La conoscenza e' potere. Sapere come funzionano i principi di persuasione puo' dare un vantaggio nelle interazioni sociali e nelle decisioni quotidiane. Chi sa come usare queste tecniche in modo efficace puo' ottenere risultati migliori, mentre chi non le conosce puo' trovarsi in difficolta'.

Il secondo paper, *A Survey on Computational Propaganda Detection*, ha offerto al team una panoramica esaustiva sulle tecniche di rilevamento della propaganda computazionale, attraverso l'importanza di un approccio integrato che combini **l'analisi del linguaggio naturale (NLP)** e **l'analisi delle reti (Network Analysis)**.

Le tecniche di rilevamento della propaganda si basano principalmente sull'analisi testuale, sfruttando strumenti avanzati di NLP per identificare e classificare le tecniche propagandistiche presenti nei testi.

Tra le principali attività spiccano la classificazione binaria, che determina se una frase contiene elementi propagandistici, e la classificazione multi-etichetta, che identifica frammenti testuali specifici utilizzando determinate tecniche e categorizzandoli.

Modelli contestuali come **BERT** e **RoBERTa**, utilizzati all'interno del codice del progetto, si sono dimostrati particolarmente efficaci nel migliorare la precisione e l'affidabilità della classificazione.

Un'ulteriore studio sulla propaganda e la narrativa è stato effettuato tramite il paper **Reports of personal experiences and stories in argumentation: Datasets and Analysis** di N. Faalk e G. Lapesa, Institute for Natural Language Processing, University of Stuttgart.

Il paper analizza la narrativa come un elemento essenziale dell'argomentazione e della discussione pubblica, ponendo particolare attenzione ai racconti personali e alle storie condivise dalle persone. Queste narrazioni, spesso basate su esperienze individuali e arricchite da elementi emotivi, si rivelano strumenti potenti per catturare l'attenzione, favorire l'empatia e rendere gli argomenti più accessibili e relazionabili.

La ricerca ha sottolineato come le narrazioni personali possano migliorare la qualità del dibattito, aumentando la persuasività degli argomenti e facilitando la comprensione reciproca tra i partecipanti.

In contesti complessi, come le discussioni su questioni sociali, le sole informazioni fattuali potrebbero non essere sufficienti a comunicare l'urgenza o l'importanza di una tematica. Le storie, in questi casi, diventano un mezzo indispensabile per costruire un dialogo più profondo e consapevole.

Queste conoscenze hanno guidato il team nello sviluppo del codice, assicurando che le soluzioni implementate tenessero conto delle dinamiche narrative e delle loro caratteristiche, per un'analisi più accurata e sensibile del contenuto testuale.

### 1.3 Organizzazione del gruppo

Il progetto si è svolto nell'arco di 20 giorni, con un impegno medio di circa 5 ore al giorno, seguendo una struttura collaborativa che ha ottimizzato l'assegnazione dei compiti.

Bruno Maria Di Maio ha iniziato immediatamente a lavorare sui primi due punti di scraping, concentrandosi sull'aggiunta di informazioni sull'autore e dei metadati relativi ai discorsi. Lucia Brando, cooperando con Denise Brancaccio, si è occupata dei punti successivi, con un'attenzione particolare alla definizione e selezione delle feature da includere nel dataset.

Successivamente, Bruno Maria Di Maio ha integrato e ottimizzato il codice, gestendo la creazione degli offset e implementando modelli di linguaggio per il trattamento avanzato del testo.

Infine Denise Brancaccio e Lucia Brando, hanno curato la redazione della documentazione finale, sintetizzando le scelte metodologiche e le fasi operative.

Questa suddivisione ha garantito un'organizzazione efficace e ha permesso di completare il progetto nei tempi previsti, soddisfacendo i requisiti richiesti dalla traccia.

### 1.4 Tecnologie utilizzate

Il codice è stato sviluppato utilizzando **Python 3.12.7** come versione di riferimento. Per lo sviluppo è stato utilizzato l'ambiente **VSCode**, configurato con i requisiti specifici del progetto.

È stata inoltre integrata un'intelligenza artificiale basata sui modelli di OpenAI per svolgere alcune sezioni del progetto. L'accesso ai servizi di OpenAI avviene tramite una chiave API dedicata, memorizzata in un file locale per garantire la sicurezza, caricata nel seguente modo:

```
os.environ["OPENAI_API_KEY"] = open("src/api_key/API-GLHF.txt", "r").read()
```

Le librerie utilizzate sono le seguenti:

- **os**: serve per interagire con il sistema operativo, permettendo di gestire file, directory, variabili d'ambiente e processi.
- **pandas**: e' utilizzata per analisi e manipolazione dei dati, offrendo strutture come DataFrame e Series per lavorare con dati tabellari.
- **openai**: e' una libreria per interagire con le API di OpenAI, usata per implementare modelli di intelligenza artificiale come ChatGPT e DALL-E.
- **concurrent.features**: serve per gestire l'esecuzione parallela di attivita', facilitando l'utilizzo di thread o processi per operazioni asincrone.
- **itertools**: offre strumenti per lavorare con iteratori, come combinazioni, permutazioni e generatori infiniti, utili per ottimizzare algoritmi.
- **transformers**: e' una libreria di Hugging Face per utilizzare modelli di machine learning avanzati (es. BERT, GPT) in compiti di NLP e altro.
- **time**: fornisce funzioni per lavorare con il tempo, come misurare la durata, formattare date o gestire ritardi temporali.
- **threading**: consente di creare e gestire thread in Python, utili per eseguire piu' operazioni simultaneamente in un'applicazione.
- **torch**: libreria principale di PyTorch, usata per il machine learning, con strumenti per creare e addestrare reti neurali.
- **nltk**: libreria per l'elaborazione del linguaggio naturale che include strumenti per analisi testuali, tokenizzazione, stemming, analisi grammaticale e creazione di modelli linguistici.
- **sklearn**: libreria per il machine learning che offre strumenti per classificazione, regressione, clustering e preprocessing dei dati, con un'interfaccia coerente e modulare.
- **wikipedia**: libreria Python che consente di interagire con l'API di Wikipedia per estrarre informazioni, riassunti, contenuti delle pagine e link correlati.
- **requests**: libreria per effettuare richieste HTTP in modo semplice e intuitivo. Supporta GET, POST, autenticazione, gestione di sessioni e manipolazione di header.
- **bs4 (Beautiful Soup)**: libreria per il parsing di documenti HTML e XML, progettata per estrarre dati strutturati da pagine web. Utile per il web scraping.
- **collections**: modulo della libreria standard Python che offre tipi di dati specializzati come Counter, defaultdict, deque e OrderedDict per la gestione efficiente di collezioni di dati.
- **re**: modulo della libreria standard Python per la gestione delle espressioni regolari, utilizzato per ricerca, sostituzione e analisi avanzata di stringhe.
- **textstat**: libreria per l'analisi della leggibilita' della complessita' dei testi. Fornisce metriche come l'indice di Flesch-Kincaid, SMOG, Gunning Fog e altri punteggi di leggibilita'.
- **textblob**: libreria per l'elaborazione del linguaggio naturale che include strumenti per analisi del sentiment, traduzione, correzione grammaticale e analisi testuali semplificate.
- **spacy**: libreria avanzata per l'elaborazione del linguaggio naturale, ottimizzata per grandi dataset e applicazioni di deep learning. Include strumenti per tokenizzazione, analisi sintattica, riconoscimento di entita' e embedding.
- **typing**: modulo della libreria standard Python per fornire annotazioni di tipo statico, inclusi tipi complessi come List, Tuple, Dict, Union e altro, per una programmazione piu' leggibile e sicura.

## 2 Descrizione del codice

Il codice sviluppato prende in input un file TSV e rinomina le colonne in surname, code e speech per uniformare i dati.

Successivamente, sfruttando il modello Llama, raccoglie informazioni aggiuntive come fullname, birthday, birthplace, death day, death place, political party, speech date, speech location, speech occasion, topic, cognitive bias e produce il riassunto dello speech.

Inoltre, genera analisi avanzate come la narrativa di sottofondo, il contesto, la struttura retorica e le keywords.

Parallelamente, una serie di modelli pre-addestrati viene utilizzata per analisi specifiche: sentiment analysis, emotion analysis, propaganda detection (incluso il tipo di propaganda e la localizzazione degli offset), toxicity classification, persuasion analysis, TF-IDF, narrative structure analysis, readability, formality e word embedding. Infine, i risultati vengono visualizzati graficamente per facilitare l'interpretazione delle relazioni e delle caratteristiche del testo.



## 2.1 Modelli

Nel progetto sono stati utilizzati modelli pre-addestrati disponibili sulla piattaforma Hugging Face che hanno permesso al team di effettuare analisi complesse come la classificazione del testo, il riconoscimento di entità e l'analisi delle narrazioni propagandistiche.

I modelli sono stati caricati attraverso la seguente funzione:

```
1 def load_model_and_tokenizer(model_name, device=None):
2     tokenizer = AutoTokenizer.from_pretrained(model_name)
3     model = AutoModelForSequenceClassification.from_pretrained(model_name)
4
5     if device and device != "cpu":
6         model = model.to(device)
7
8     return tokenizer, model
```

Per evitare deadblock e' stato disattivato il parallelismo tra i tokenizers.

```
1 os.environ["TOKENIZERS_PARALLELISM"] = "false"
```

Per gestire il limite di 512 caratteri e' stata implementata una funzione che suddivide automaticamente i testi pi lunghi in segmenti compatibili con il modello, assicurando che ogni frammento sia elaborato correttamente.

Successivamente, i risultati parziali vengono aggregati per fornire un'analisi completa e coerente del testo originale, garantendo cosi la gestione efficace di documenti di grandi dimensioni.

```
1 # Funzione per creare chunk con finestra scorrevole
2 def split_text_sliding_window(text, chunk_size=512, overlap=128):
3     words = text.split()
4     chunks = []
5     start = 0
6
7     while start < len(words):
8         end = start + chunk_size
9         chunk = ' '.join(words[start:end])
10        chunks.append(chunk)
11        start = end - overlap
12
13    return chunks
```

### 2.1.1 Sentiment Analysis

```
1 sentiment_tokenizer, sentiment_model =  
2     load_model_and_tokenizer("nlptown/bert-base-  
3     multilingual-uncased-sentiment")  
sentiment_labels = ["very negative", "negative", "neutral", "positive", "very  
positive"]
```

### 2.1.2 Emotion Analysis

```
1 emotion_tokenizer, emotion_model =  
2     load_model_and_tokenizer("bhadresh-savani/distilbert-  
3     base-uncased-emotion")  
emotion_labels = ["sadness", "joy", "love", "anger", "fear", "surprise"]
```

### 2.1.3 Propaganda Classifier

```
1 propaganda_classifier = pipeline("text-classification",  
2     model="IDA-SERICS/PropagandaDetection", device=0 if device != "cpu" else  
3     -1)
```

### 2.1.4 Toxicity Classifier

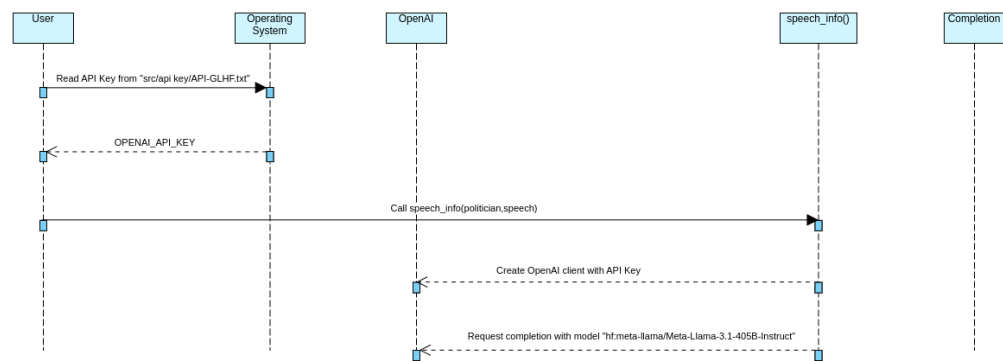
```
1 toxicity_classifier = pipeline("text-classification",  
2     model="citizenlab/distilbert-base-multilingual-cased-toxicity", device=0  
3     if device != "cpu" else -1)
```

### 2.1.5 Word Embedding

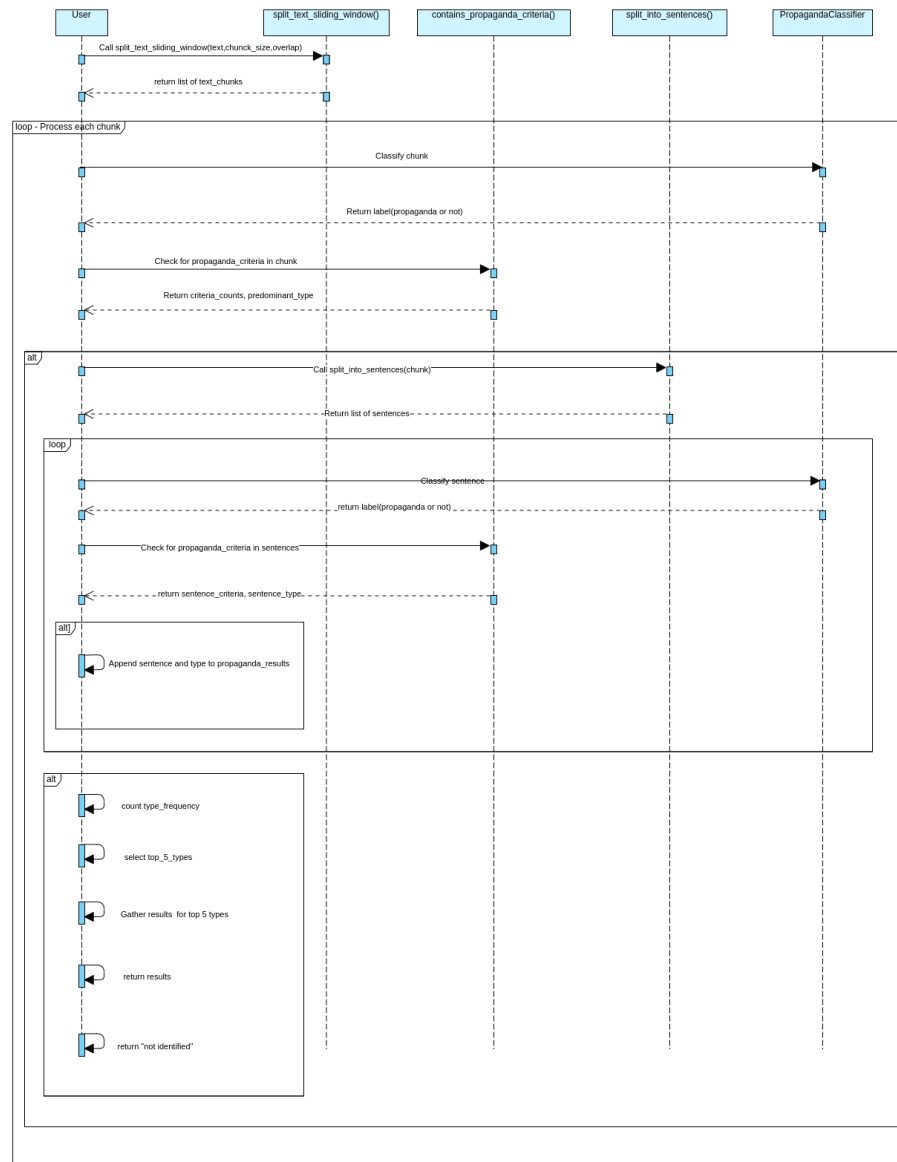
```
1 tokenizer = AutoTokenizer.from_pretrained("sentence-transformers/  
2     all-MiniLM-L6-v2")  
3 model = AutoModel.from_pretrained("sentence-transformers/  
4     all-MiniLM-L6-v2")
```

### 3 Sequence Diagram

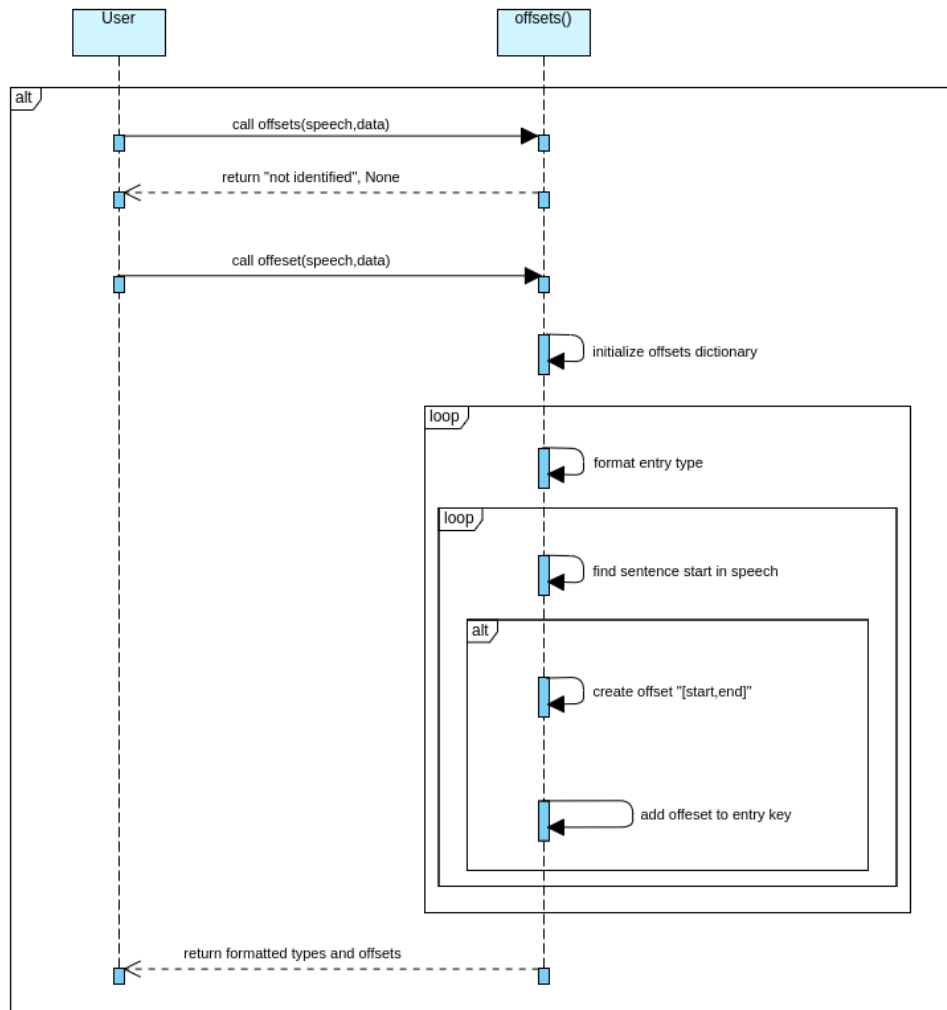
#### 3.1 Funzione def speech\_info(politician,speech)



### 3.2 Funzione def detect\_propaganda\_type(text, classifier=propaganda\_classifier, chunk\_size=512, overlap=128)



### 3.3 Funzione def offsets(speech,data)



## 4 Funzioni

### 4.1 Verify Politician

La funzione verifica se un dato titolo della pagina web corrisponde a una figura politica. Utilizza il riassunto della prima frase della pagina di Wikipedia e controlla se contiene parole chiave come "who", "politician" o "statesman". Se si verifica un errore di disambiguazione o la pagina non esiste, restituisce False.

```
1 def verify_politician(titolo):
2
3     try:
4         # Prende la prima frase del riassunto della pagina
5         summary = wikipedia.summary(titolo, sentences=1)
6         # Filtra per le parole chiave contenute nella prima frase
7         return any(keyword in summary.lower() for keyword in ["who",
8             "politician", "statesman"])
9
10    except (wikipedia.exceptions.DisambiguationError,
11        wikipedia.exceptions.PageError):
12        return False
```

### 4.2 Get Politician Info

La funzione effettua una ricerca, verifica se il cognome e' associato a un politico (tramite verify politician), e scarica l'HTML della pagina per analizzarlo con BeautifulSoup. Successivamente, in base al parametro richiesto, estrae e restituisce informazioni come nome, data e luogo di nascita, data e luogo di morte, o partito politico.

Questi dati vengono salvati nel dizionario per eventuali richieste future.

```
1 # Ottiene le informazioni dal politico
2 def get_info(cognome, parametro):
3     # Verifica se il parametro esiste
4     if cognome in politici and parametro in politici[cognome]:
5         return politici[cognome][parametro]
6
7     # Verifica se il cognome esiste e crea il dizionario
8     elif cognome not in politici:
9         politici[cognome] = {}
10
11     # Verifica se il soup non esiste e lo deve creare
12     if "soup" not in politici[cognome]:
13         # Cerca su wikipedia il cognome dato
14         results = wikipedia.search(cognome)
15
16         # Trova l'indice del risultato del politico tramite la funzione
17         # verify_politician
18         index = next((i for i, result in enumerate(results) if
19             verify_politician(result)), None)
20
21     if index is not None:
22         # Prende l'URL della pagina wikipedia
23         url = wikipedia.page(results[index]).url
24
25         # Effettua una richiesta GET alla pagina
26         response = requests.get(url)
27
28         # Verifica se la richiesta ha avuto successo
29         if response.status_code == 200:
30             # Crea l'oggetto BeautifulSoup per analizzare l'HTML e lo
31             # aggiunge dizionario
32             politici[cognome]["soup"] = BeautifulSoup(response.text,
33                 features="html.parser")
34         else:
```

```

31         print("Errore nella richiesta della pagina:",
32               response.status_code)
33
34     if "soup" in politici[cognome]:
35         soup = politici[cognome]["soup"]
36         # Trova l'infobox contenente le rows in cui sono presenti diversi dati
37         infobox = soup.find("table", {"class": "infobox"}).find_all("tr")
38
39     # Cerca il nome
40     if parametro == "name":
41         # Estrai il nome completo
42         name = soup.find("div", {"class": "nickname"}).text
43
44         if name:
45             # Aggiunge il nome al dizionario
46             politici[cognome]["name"] = name
47             # Restituisce il nome
48             return name
49         return "N/A"
50
51     # Cerca la data di nascita
52     elif parametro == "birthday":
53         # Estrae la data di nascita
54         birth_date = soup.find("span", {"class": "bday"}).text
55
56         if birth_date:
57             # Aggiunge la data di nascita al dizionario
58             politici[cognome]["birthday"] = birth_date
59             # Restituisce la data di nascita
60             return birth_date
61         return "N/A"
62
63     # Cerca il luogo di nascita
64     elif parametro == "birthplace":
65         # Trova il luogo di nascita in caso questo sia un link
66         try:
67             # Estrae il luogo di nascita
68             birth_place = next(
69                 # Prende il link del luogo di nascita e lo unisce allo Stato
70                 # di nascita
71                 ("%s" % (row.find("a").text,
72                             row.find("a").find_next_sibling(string=True).text)
73                 for row in infobox if "Born" in row.get_text()),
74                 None
75             )
76
77         # Trova il luogo di nascita in caso questo sia una scritta
78         # semplice
79         except AttributeError:
80             # Estrae il luogo di nascita
81             birth_place = next(
82                 # Prende la stringa del luogo di nascita
83                 ("%s" % row.find("br").find_next_sibling(string=True).text
84                 for row in infobox if "Born" in row.get_text()),
85                 None
86             )
87
88         if birth_place:
89             # Aggiunge il luogo di nascita al dizionario
90             politici[cognome]["birthplace"] = birth_place
91             # Restituisce il luogo di nascita
92             return birth_place
93         return "N/A"
94
95     # Cerca la data di morte
96     elif parametro == "deathday":
97         # Estrae la data di morte
98         death_day = next(
99             # Prende lo span della data di morte
100             (row.find("span").text
101             for row in infobox if "Died" in row.get_text()),

```

```

97         None
98     )
99
100     if death_day:
101         # Riformatta la data per togliere le parentesi ed avere il
102             formato aaaa-mm-gg
103         death_day = death_day[1:-1]
104         # Aggiunge la data di morte al dizionario
105         politici[cognome]["deathday"] = death_day
106         # Restituisce la data di morte
107         return death_day
108     # Aggiunge None al dizionario nel caso in cui la persona non sia
109     morta
110     politici[cognome]["deathday"] = None
111     return None
112
113 # Cerca il luogo di morte
114 elif parametro == "deathplace":
115     # Trova il luogo di morte in caso questo sia un link
116     try:
117         # Estrae il luogo di morte
118         death_place = next(
119             # Prende il link del luogo di morte e lo unisce allo Stato di
120             morte
121             ("%s" % (row.find("a").text,
122                 row.find("a").find_next_sibling(string=True).text)
123             for row in infobox if "Died" in row.get_text()),
124             None
125         )
126     # Trova il luogo di nascita in caso questo sia una scritta
127     semplice
128     except AttributeError:
129         death_place = next(
130             # Prende la stringa del luogo di morte
131             ("%s" % row.find("br").find_next_sibling(string=True).text
132             for row in infobox if "Died" in row.get_text()),
133             None
134         )
135     if death_place:
136         # Aggiunge il luogo di morte al dizionario
137         politici[cognome]["deathplace"] = death_place
138         # Restituisce il luogo di morte
139         return death_place
140     # Aggiunge None al dizionario nel caso in cui la persona non sia
141     morta
142     politici[cognome]["deathplace"] = None
143     return None
144
145 # Cerca il partito politico
146 elif parametro == "party":
147     # Estrae il partito politico
148     party = next(
149         # Prende il link del partito politico
150         ("%s" % row.find("a").text
151         for row in infobox if "Political party" in row.get_text()),
152         None
153     )
154     if party:
155         # Aggiunge il partito politico al dizionario
156         politici[cognome]["party"] = party
157         # Restituisce il partito politico
158         return party
159     return "N/A"
160
161 else:
162     print("Parametro selezionato non valido.")

```



### 4.3 Speech Info

La funzione utilizza llama per ottenere varie informazioni sul discorso.

Oltre alla data e al luogo in cui si e' tenuto il discorso, attraverso il prompt vengono riportati anche l'occasione in cui si tenuto, il topic di cui tratta, i cognitive bias, una summary e un'analisi dello storytelling.

```
1 def speech_info(politician, speech):
2     try:
3         client = openai.OpenAI(
4             api_key=os.environ.get("OPENAI_API_KEY"),
5             base_url="https://glhf.chat/api/openai/v1"
6         )
7         completion = client.chat.completions.create(
8             model="hf:meta-llama/Meta-Llama-3.1-405B-Instruct"
9             messages=[Prompt]
10            temperature=0
11
12        return speech, completion.choices[0].message.content
13
14    except (openai.RateLimitError, Exception):
15        raise
```

## 4.4 Speech Data

Questa funzione e' progettata per estrarre i dati associati a un determinato discorso (speech) da un dizionario (results dict) e formattarli in un array con esattamente 8 valori, gestendo eventuali errori.

```
1 def get_speech_data(speech, results_dict):
2     try:
3         parts = results_dict[speech].split('$')
4         parts += ["N/A"] * (8 - len(parts))
5         return parts[:8]
6     except (KeyError, Exception):
7         return ["N/A", "N/A", "N/A", "N/A", "N/A", "N/A", "N/A", "N/A"]
```

## 4.5 Classify Text

La funzione classifica un testo in base a un modello pre-addestrato.

```
1 # Funzione unificata per classificare il testo
2 def classify_text(text, model, tokenizer, labels):
3     inputs = tokenizer(text, return_tensors="pt", truncation=True)
4     outputs = model(**inputs)
5     probabilities = torch.nn.functional.softmax(outputs.logits, dim=1)
6     prediction = torch.argmax(probabilities, dim=1).item()
7     return labels[prediction].capitalize()
```

## 4.6 Split Sentences

La funzione suddivide un testo in singole frasi utilizzando una semplice espressione regolare.

```
1 # Dividi il testo in frasi usando una semplice regex
2 def split_into_sentences(text):
3     sentence_endings = re.compile(r'(?<=[.!?])\s+')
4     return sentence_endings.split(text)
```

## 4.7 Propaganda Criteria

La funzione analizza un testo per rilevare la presenza di elementi riconducibili a tecniche di propaganda.

Confronta il contenuto del testo con un insieme di parole o frasi chiave definite per ciascun criterio, come toni accusatori, distorsione dell'opposizione o appelli alla paura.

Per ogni criterio, conta quante volte le frasi chiave sono presenti nel testo.

Se vengono rilevati uno o più criteri, restituisce un dizionario con il conteggio delle occorrenze e il criterio predominante (quello con più corrispondenze).

Se non viene rilevato nulla, restituisce 'None'.

```
1 # Funzione per rilevare i criteri di propaganda
2 def contains_propaganda_criteria(text):
3     criteria = {
4         "accusatory_tone": ["fault of"],
5         "opposition_distortion": ["enemies", "traitor", "dishonest",
6         "corrupt"],
7         "slogan_repetition": ["you have to believe", "the truth is", "don't
8         forget", "do not forget"],
9         "fear_appeals": ["danger", "threat", "crisis", "chaos"],
10        "flagwaving": ["our nation", "homeland", "defend values"],
11        "black_white_fallacy": ["either with us or against us", "only
12        choice"],
13        "clichés": ["golden times", "like once upon a time", "old times"]
14    }
15    # Conta quante frasi/parole chiave sono presenti per ciascun criterio
16    type_counts = {key: sum(phrase in text.lower() for phrase in phrases) for
17    key, phrases in criteria.items()}
18    # Filtra i criteri che hanno almeno una corrispondenza
19    detected_types = {key: count for key, count in type_counts.items() if
20    count > 0}
21    # Se non ci sono criteri rilevati, ritorna None
22    if not detected_types:
23        return None, None
24    predominant_type = max(detected_types, key=detected_types.get)
25    return detected_types, predominant_type
```

## 4.8 Detect Propaganda Type

La funzione suddivide il testo in chunk sovrapposti per garantire che tutte le parti siano valutate senza perdere contesto.

Ogni chunk viene analizzato sia tramite un modello di classificazione di propaganda che con una funzione che cerca criteri definiti di propaganda nel testo.

Se un chunk e' classificato come propaganda o contiene criteri rilevanti, viene ulteriormente diviso in frasi per un'analisi pi' dettagliata.

Per ciascuna frase classificata come propaganda, vengono salvati il tipo di propaganda identificato e la frase stessa.

Infine, se vengono rilevati risultati, la funzione restituisce i cinque tipi di propaganda piu' comuni, ciascuno accompagnato dalle frasi associate.

Se non viene rilevato nulla, restituisce "Not identified".

```
1 def detect_propaganda_type(text, classifier=propaganda_classifier,
2   chunk_size=512, overlap=128):
3     # Divide il testo in chunk
4     text_chunks = split_text_sliding_window(text, chunk_size, overlap)
5     propaganda_results = []
6
7     for chunk in text_chunks:
8         # Classifica il chunk
9         is_propaganda_by_model = classifier(chunk, truncation=True,
10          max_length=chunk_size)[0]['label'] == "propaganda"
11         criteria_counts, predominant_type =
12          contains_propaganda_criteria(chunk)
13
14         # Analizza le frasi del chunk
15         if is_propaganda_by_model or criteria_counts:
16             sentences = split_into_sentences(chunk) # Dividi il chunk in
17             frasi
18             for sentence in sentences:
19                 # Classifica la singola frase
20                 sentence_is_propaganda = classifier(sentence,
21                  truncation=True)[0]['label'] == "propaganda"
22                 sentence_criteria, sentence_type =
23                  contains_propaganda_criteria(sentence)
24
25                 if sentence_is_propaganda or sentence_criteria:
26
27                     propaganda_results.append({
28                         "type": sentence_type,
29                         "sentence": sentence
30                     })
31
32     # Se esistono risultati di propaganda
33     if propaganda_results:
34         # Conta la frequenza di ogni tipo di propaganda
35         types = [result["type"] for result in propaganda_results]
36         type_counts = Counter(types)
37         top_5_types = type_counts.most_common(5)
38         # Estrai i risultati per le 5 tipologie piu comuni
39         results = []
40         for propaganda_type, _ in top_5_types:
41             sentences_of_type = [
42                 {"sentence": result["sentence"]}
43                 for result in propaganda_results if result["type"] ==
44                  propaganda_type
45             ]
46             results.append({
47                 "type": propaganda_type,
48                 "sentences": sentences_of_type
49             })
50
51     return results
52
53     return "Not identified"
```

## 4.9 Contains Toxicity Criteria

La funzione analizza un testo per identificare criteri di tossicità basati su un insieme di categorie predefinite.

Ogni categoria, come insulti, minacce o linguaggio violento, viene associata a una lista di parole o frasi chiave indicative di quella forma di tossicità.

Confronta il testo con le chiavi, calcola il numero di corrispondenze per ciascuna categoria, e restituisce un dizionario con i criteri rilevati e il loro conteggio.

Inoltre, identifica quale categoria è predominante, ovvero quella con il maggior numero di corrispondenze.

Se nessun criterio viene rilevato, restituisce None.

```
1 # Funzione per rilevare i criteri di tossicità
2 def contains_toxicity_criteria(text):
3     criteria = {
4         "insults": ["idiot", "stupid", "moron", "fool", "scum", "loser",
5         "crooked", "degenerate", "ignoramus"],
6         "threats": ["destroy", "kill", "eliminate", "eradicate", "you'll
7         regret", "annihilate", "purge", "defeat utterly"],
8         "violent_language": ["blood", "war", "fight", "smash", "knock them
9         out", "battle", "no quarter"],
10        "divisive_statements": ["us vs them", "enemy of the people", "deep
11        state", "traitors", "fifth column"],
12        "personal_attacks": ["you always fail", "your fault", "worthless",
13        "incompetent", "phony", "betrayal"],
14        "exaggerations": ["always lie", "never truthful", "the worst",
15        "greatest challenge", "unprecedented threat"],
16        "fear_appeals": ["danger", "threat", "crisis", "chaos", "they're
17        coming for you", "crime wave"],
18        "flagwaving": ["our nation", "homeland", "defend values", "make
19        america great again", "fatherland", "motherland", "maga"]
20    }
21    # Conta quante frasi/parole chiave sono presenti per ciascun criterio
22    type_counts = {key: sum(phrase in text.lower() for phrase in phrases) for
23    key, phrases in criteria.items()}
24    # Filtra i criteri che hanno almeno una corrispondenza
25    detected_types = {key: count for key, count in type_counts.items() if
26    count > 0}
27    # Se non ci sono criteri rilevati, ritorna None
28    if not detected_types:
29        return None, None
30    # Identifica il tipo di tossicità
31    predominant_type = max(detected_types, key=detected_types.get)
32    return detected_types, predominant_type
```

## 4.10 Classify Toxicity

La funzione classifica un testo identificando eventuali criteri di tossicità.

Utilizza la funzione contains toxicity criteria per analizzare il testo e determinare quali categorie di tossicità sono presenti, insieme al conteggio delle loro occorrenze.

Se vengono rilevati criteri, questi vengono formattati come una lista leggibile in cui ogni riga mostra il nome del criterio (con spazi al posto degli underscore) e il numero di volte che è stato riscontrato.

```
1 # Funzione per classificare la tossicità
2 def classify_toxicity(text):
3     try:
4         # Rilevamento criteri
5         detected_criteria, _ = contains_toxicity_criteria(text)
6         # Formattazione della lista leggibile
```

```
7         # Formattare i criteri rilevati uno per riga
8         if detected_criteria:
9             return '\n'.join(["%s: %s" % (k, v) for k, v in
10                                detected_criteria.items()]).replace("_", " ").title()
11         else:
12             return ""
13     except Exception:
14         raise
```

## 4.11 Offsets

La funzione calcola gli offset delle frasi identificate come propaganda all'interno di un discorso, fornendo la posizione iniziale e finale di ogni frase rilevata per ciascun tipo di propaganda attraverso il formato [inizio:fine].

Se il parametro data contiene il valore "Not identified", la funzione restituisce direttamente questo valore con None, indicando che non sono state rilevate frasi di propaganda.

```
1 def offsets(speech, data):
2
3     # Verifica che la propaganda sia stata correttamente identificata
4     if data == "Not identified":
5         return "Not identified", None
6
7     # Crea un dizionario per definire gli offset
8     offsets = {
9         # Crea una chiave per tipo e la formatta correttamente
10        entry["type"].replace("_", " ").capitalize(): [
11            # Crea un offset per ogni frase
12            "[%s:%s]" % (start, start + len(sentence["sentence"]))
13            for sentence in entry["sentences"]
14            # Trova le frasi e le misura
15            if (start := speech.find(sentence["sentence"])) != -1
16        ]
17    }
18    for entry in data
19
20    # Restituisce il tipo di propaganda e gli offset per ogni tipo di
21    propaganda
22    return "\n".join(offsets.keys()), "\n".join("%s: %s" % (key, ",
23        ".join(element)) for key, element in offsets.items())
```

## 4.12 Formality

La funzione determina se un discorso e' formale o informale in base al suo livello di soggettivita' utilizzando il modulo TextBlob per analizzare il sentiment del testo.

Se la lunghezza del discorso e' inferiore o uguale a 150 caratteri, la funzione restituisce "N/A", indicando che la valutazione non e' applicabile.

Se il discorso e' piu' lungo, la funzione analizza la soggettivita' del testo: una soggettivita' inferiore a 0,5 suggerisce che il testo e' formale, mentre una soggettivita' piu' alta indica che il discorso e' informale.

```
1 def is_formal(speech):
2     if len(speech) <= 150:
3         return "N/A"
4
5     blob = TextBlob(speech)
6     # Bassa soggettivita indica formalita
7     if blob.sentiment.subjectivity < 0.5:
8         return "Formal"
9     else:
10        return "Informal"
```

## 4.13 Readability

La funzione calcola diverse metriche di leggibilita' di un discorso utilizzando il modulo textstat, che fornisce diversi indici per valutare quanto sia facile o difficile leggere un testo.

Le metriche prese in considerazione sono: *flesch reading ease*, *flesch-kincaid grade*, *gunning fog*, *SMOG index*.

```
1 # Funzione per calcolare metriche di leggibilita
2 def calcola_legibilita(speech):
3     return {
4         "flesch_reading_ease": textstat.flesch_reading_ease(speech), #
5             facilita di lettura basandosi sulla lunghezza delle frasi e sul
6             numero di sillabe per parola
7         "flesch_kincaid_grade": textstat.flesch_kincaid_grade(speech), #
8             anni di istruzione necessari per comprendere il testo
9         "gunning_fog": textstat.gunning_fog(speech), # complessita del speech
10            in base alla lunghezza delle frasi e alla percentuale di parole
11            complesse
12         "smog_index": textstat.smog_index(speech), # come quello precedente
13         "text_standard": textstat.text_standard(speech) # sintetizzazione dei
14            risultati in un livello scolastico approssimativo
15     }
```

## 4.14 Analyze Narrative Structure

La funzione analizza la struttura narrativa di un discorso, concentrandosi su vari aspetti legati alla lunghezza e alla variabilita' delle frasi.

In seguito alla tokenizzazione, viene calcolato il numero di parole all'interno di ogni token, generando una lista che rappresenta la lunghezza di ogni frase.

Dopo aver ottenuto le lunghezze delle frasi, la funzione calcola diverse metriche narrative.

La lunghezza media delle frasi viene determinata calcolando la media aritmetica del numero di parole per ciascuna frase.

Viene poi calcolata la deviazione standard, che misura quanto le lunghezze delle frasi variano l'una dall'altra, dando un'idea della coerenza o della varieta' nella struttura del discorso.

Inoltre, la funzione conta il numero totale di frasi nel discorso e, infine, calcola un punteggio di complessita', che e' il rapporto tra la lunghezza media delle frasi e la deviazione standard.



```
1 def analyze_narrative_structure(speech):
2     sentences = sent_tokenize(speech)
3
4     # Analyze sentence length variation
5     sentence_lengths = [len(sentence.split()) for sentence in sentences]
6
7     narrative_metrics = {
8         'Average Sentence Length': np.mean(sentence_lengths),
9         'Sentence Length Variation': np.std(sentence_lengths),
10        'Total Sentences': len(sentences),
11        'Complexity Score': np.mean(sentence_lengths) /
12            (np.std(sentence_lengths) + 1)
13    }
14
15    return "\n".join(["%s: %.3f" % (key, value) if isinstance(value, float)
16        else "%s: %d" % (key, value) for key, value in
17        narrative_metrics.items()])
```

## 4.15 TF-IDF

La funzione estrae le parole chiave piu' rilevanti da un discorso utilizzando il metodo TF-IDF. Imposta un limite sul numero di parole chiave da restituire (top 5) e ignora le stop words. Successivamente, calcola la matrice TF-IDF e seleziona le parole con i punteggi piu' alti. Le parole chiave e i loro punteggi vengono quindi ordinati e restituiti in una stringa formattata con i punteggi arrotondati a tre decimali, mettendo in evidenza le parole piu' significative nel testo.

```
1 def tfidf(speech):
2     # Definisce il numero di parole massimo per le keywords
3     max_ngram_length = 1
4     # Definisce quante keywords dobbiamo trovare
5     top_n = 5
6
7     # Inserisce i parametri di ricerca e il dizionario delle stop words
8     vectorizer = TfidfVectorizer(
9         ngram_range=(1, max_ngram_length),
10        stop_words='english'
11    )
12
13    # Prende la matrice TF-IDF
14    tfidf_matrix = vectorizer.fit_transform([speech])
15
16    # Prende i nomi delle parole
17    feature_names = vectorizer.get_feature_names_out()
18
19    # Prende gli score TF-IDF
20    tfidf_scores = tfidf_matrix.toarray()[0]
21
22    # Crea una lista di tuple (keyword, score)
23    keyword_scores = [
24        (feature_names[idx], score)
25        for idx, score in enumerate(tfidf_scores)
26        if score > 0
27    ]
28
29    # Ordina per score decrescente
30    keyword_scores = sorted(keyword_scores, key=lambda x: x[1], reverse=True)
31
32    # Prende solo le top_n keywords
33    top_keywords = keyword_scores[:top_n]
34
35    # Restituisce le parole chiave con i rispettivi punteggi
36    return "\n".join(["%s: %.3f" % (kw.title(), round(score, 3)) for kw,
37                        score in top_keywords])
```

## 4.16 Word Embeddings

La funzione calcola gli embeddings di un testo suddividendolo in chunk di dimensione specificata, con un sovrapposizione tra i chunk.

Per ogni chunk, il testo viene tokenizzato e passato attraverso il modello per ottenere gli embeddings dei token.

Gli embeddings di ogni token nel chunk vengono mediati (calcolando la media dei vettori) per ottenere un embedding complessivo per il chunk.

Alla fine, gli embeddings medi di tutti i chunk vengono mediati nuovamente per ottenere un embedding finale che rappresenta l'intero testo.

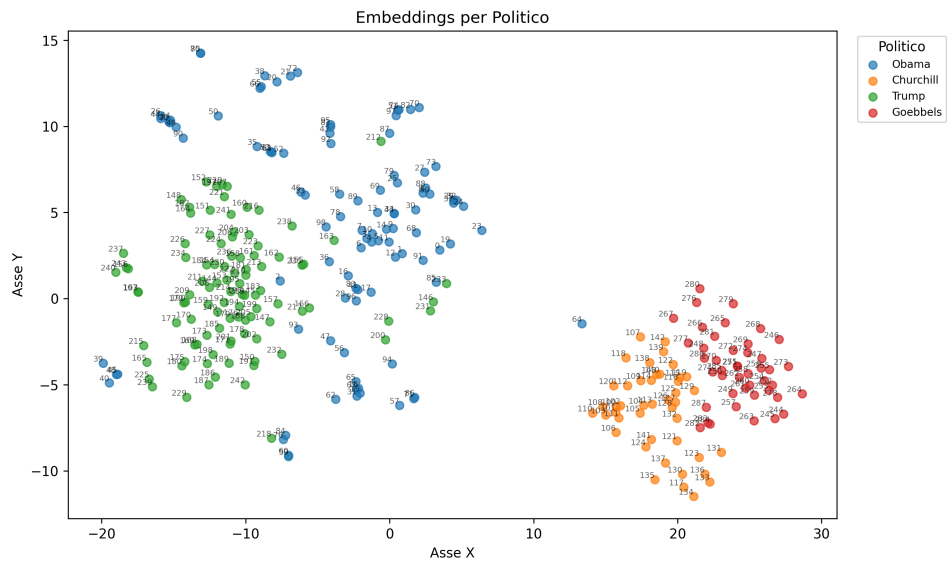
```
1 # Funzione per calcolare gli embeddings sui chunk
2 def calcola_embedding(testo, tokenizer, model, chunk_size=512,
3                        overlap=128):
4     chunks = split_text_sliding_window(testo, chunk_size=chunk_size,
5                                         overlap=overlap)
6     embeddings_chunk = []
7
8     for chunk in chunks:
```

```
7         inputs = tokenizer(chunk, return_tensors="pt", padding=True,
8                               truncation=True, max_length=chunk_size)
9
10        with torch.no_grad():
11            outputs = model(**inputs)
12
13        token_embeddings = outputs.last_hidden_state
14        mean_embedding = token_embeddings.mean(dim=1).squeeze()
15        embeddings_chunk.append(mean_embedding.cpu().numpy())
16
17    if embeddings_chunk:
18        final_embedding = np.mean(embeddings_chunk, axis=0)
19    else:
20        final_embedding = np.zeros(model.config.hidden_size)
21
22    return final_embedding
```

## 4.17 Visualize Embeddings

La funzione utilizza i word embedding per creare uno scatter plot in cui ogni punto rappresenti un discorso e ogni colore un politico. Discorsi semanticamente simili saranno altrettanto vicini sul grafico, discorsi semanticamente differenti verranno rappresentati pi distanti.

```
1 # Visualizza gli embeddings su scatter plot
2 def visualize_embeddings(dataset, plt_show,
3     output_path="word_embedding.png"):
4     tokenizer = AutoTokenizer.from_pretrained("sentence-transformers/
5         all-MiniLM-L6-v2")
6     model = AutoModel.from_pretrained("sentence-transformers/
7         all-MiniLM-L6-v2")
8
9     dataset["Embedding"] = dataset["Speech"].fillna("").apply(lambda x:
10         calcola_embedding(x, tokenizer, model))
11     embeddings = np.stack(dataset["Embedding"].values)
12
13     tsne = TSNE(n_components=2, random_state=42, perplexity=30)
14     embeddings_2d = tsne.fit_transform(embeddings)
15
16     dataset["x"] = embeddings_2d[:, 0]
17     dataset["y"] = embeddings_2d[:, 1]
18
19     politici_unici = dataset["Surname"].unique()
20
21     # Crea una figura di dimensioni (10, 6) e posiziona i punti in base ai
22     # nomi
23     plt.figure(figsize=(10, 6))
24     for politico in politici_unici:
25         subset = dataset[dataset["Surname"] == politico]
26         plt.scatter(subset["x"], subset["y"], label=politico, alpha=0.7)
27
28     # Enumera ogni discorso sui punti
29     for i, (x, y) in enumerate(zip(dataset["x"], dataset["y"])):
30         plt.text(x, y, str(i), fontsize=6, ha='right', va='bottom',
31             color='black', alpha=0.6)
32
33     # Impostazioni del grafico (titolo, label di ascissa e ordinate, legenda
34     # e layout finale)
35     plt.title("Embeddings per Politico", fontsize=12)
36     plt.xlabel("Asse X", fontsize=10)
37     plt.ylabel("Asse Y", fontsize=10)
38     plt.legend(title="Politico", bbox_to_anchor=(1.02, 1), loc="upper left",
39         fontsize=8)
40     plt.tight_layout()
41     # Salva il grafico nella stessa directory del file Python
42     plt.savefig(output_path, dpi=300, bbox_inches='tight')
43     # Chiudi il plot per evitare conflitti in future visualizzazioni
44
45     if plt_show == "s":
46         return plt
47
48     plt.close()
49     return None
```



## 4.18 Analyze Persuasion

La funzione analizza un testo per valutare la presenza di elementi persuasivi.

Usa spaCy per elaborare il testo e calcolare il numero di superlativi, domande retoriche, e parole emotive, di urgenza, autorità e prova sociale.

Ogni metrica è normalizzata in base alla frequenza nel testo e ponderata con un peso specifico che riflette l'importanza persuasiva di ciascun elemento.

Somma gli score normalizzati per ottenere un punteggio totale di persuasione.

I risultati vengono restituiti in un formato leggibile, con ogni metrica e il suo valore separati da una riga.

```
1 def analyze_persuasion(text):
2     # Analisi del testo usando spaCy
3     doc = nlp(text.lower())
4
5     # Calcola metriche persuasive
6     # Per ogni tipo di elemento persuasivo, conta quante volte appare nel
7     # testo
8     metrics = {
9         # conta i superlativi
10        "superlatives": sum(1 for token in doc if token.tag_ in ["JJS",
11        "RBS"]),
12        # conta le domande retoriche, basandosi sul fatto che terminano
13        # con un punto di domanda
14        "rhetorical_questions": sum(1 for sent in doc.sents if
15        sent.text.strip().endswith("?")),
16        # conta le parole emotive che appartengono alla lista predefinita
17        "emotional_words": sum(1 for token in doc if token.text in
18        PERSUASIVE_WORDS["emotional"]),
19        # conta le parole legate all'urgenza dalla lista predefinita
20        "urgency_words": sum(1 for token in doc if token.text in
21        PERSUASIVE_WORDS["urgency"]),
22        # conta le parole legate all'autorità
23        "authority_words": sum(1 for token in doc if token.text in
24        PERSUASIVE_WORDS["authority"]),
25        # conta le parole di prova sociale
26        "social_proof_words": sum(1 for token in doc if token.text in
27        PERSUASIVE_WORDS["social_proof"])
28    }
29
30    # Calcola il numero totale di token
31    total_tokens = len(doc)
32    # Se il testo non contiene token, restituisce un dizionario vuoto
33    if total_tokens == 0:
34        return {key: 0.0 for key in metrics.keys() | {"total_score"}}
35
36    # Definisce il peso per ogni elemento persuasivo, indicando quanto
37    # importante per ciascuna metrica
38    weights = {
39        "superlatives": 1.5, # peso medio
40        "rhetorical_questions": 2.0, # sono più persuasive
41        "emotional_words": 2.5, # sono le più persuasive
42        "urgency_words": 2.0, # sono persuasive
43        "authority_words": 1.8, # peso medio-alto
44        "social_proof_words": 1.7 # peso inferiore
45    }
46
47    # Calcola uno score normalizzato per ogni metrica
48    # Per ogni metrica:
49    # - Divido il conteggio per il numero totale di token (normalizzazione)
50    # - Moltiplico per il peso associato per enfatizzare la sua importanza
51    normalized_scores = {
52        key: (metrics[key] / total_tokens) * weights[key] for key in metrics
53    }
54
55    # Somma tutti gli score normalizzati per ottenere lo score totale di
56    # persuasione
57    total_score = sum(normalized_scores.values())
58    normalized_scores["total_score"] = round(total_score, 3) # arrotondare a
```

```
3 - scelta comune per rappresentare valori numerici in cui la
precisione oltre il millesimo non risulta necessaria
return '\n'.join("%s: %s" % (k.replace('_', ' ').title(), v) for k, v in
normalized_scores.items()) # restituisce gli score
```

## 5 Analisi delle Prestazioni

Il processo di analisi ha richiesto complessivamente *18 minuti e 9,04 secondi*. Questo tempo e' stato suddiviso in diverse fasi:

- scraping dei dati
- estrazione delle informazioni
- analisi approfondita
- creazione di rappresentazioni vettoriali delle parole

La durata complessiva e' influenzata da diversi fattori, tra cui la quantita' di dati da elaborare, la complessita' delle analisi e le risorse computazionali disponibili.

I risultati dell'analisi, inclusi i dati grezzi e le rappresentazioni vettoriali, sono stati salvati in formato strutturato per facilitare ulteriori approfondimenti e sviluppi.

In futuro, si prevede di ottimizzare il processo per ridurre i tempi di esecuzione, ad esempio attraverso l'implementazione di algoritmi piu' efficienti o l'utilizzo di hardware piu' performante.

Inoltre, si valtera' l'opportunita' di espandere l'analisi includendo nuove feature e modelli linguistici piu' avanzati.

```
🔧 Eseguo lo scraping delle informazioni dei politici
✅ Scraping completato
🕒 Tempo di esecuzione dello scraping: 0 ore, 0 minuti, 11.43 secondi

🔍 Ottengo le informazioni sui discorsi
✅ Informazioni ottenute
🕒 Tempo di esecuzione per l'ottenimento delle informazioni: 0 ore, 7 minuti, 27.64 secondi

🧠 Eseguo l'analisi dei discorsi
✅ Analisi completata
🕒 Tempo di esecuzione dell'analisi: 0 ore, 7 minuti, 56.59 secondi

📄 Aggiungo feature text-based
✅ Aggiunta completata
🕒 Tempo di esecuzione dell'analisi: 0 ore, 0 minuti, 5.37 secondi

📊 Visualizzare il grafico del word embedding? (s/n) n
📈 Calcolo il word embedding
✅ Calcolo completato
🕒 Tempo di esecuzione del word embedding: 0 ore, 1 minuti, 14.01 secondi

✅ Grafico del word embedding salvato in src/word_embedding.png
✅ Dataset salvato in src/dataset/speech-b.csv

🕒 Tempo di esecuzione totale: 0 ore, 18 minuti, 9.04 secondi
```



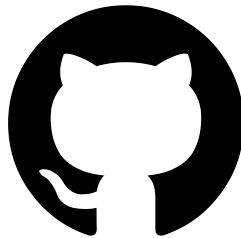
## 6 Conclusioni

Il team e' riuscito a completare con successo il project work assegnato grazie a una solida collaborazione, una gestione efficace del tempo e un'accurata pianificazione delle attivita'.

La continua iterazione, brainstorming e testing dei modelli hanno permesso di superare le difficolta' tecniche iniziali dovute alla novita' della richiesta e rispettare i tempi stabiliti.

Alla fine, il lavoro e' stato completato con il raggiungimento degli obiettivi prefissati in modo soddisfacente.

## 7 Collegamenti esterni



**Codice completo:**

<https://github.com/lbrando/project-work>