

# The Data

The study was conducted using Italian X's tweets from November 2022 to April 2024. Due to recent changes in X's policies regarding APIs and data collection, this part of the research turned out to be difficult and time-consuming. To overcome the data collection limitation and the high fees charged by Elon Musk's platform, we used Nitter, a discontinued free and open-source alternative viewer for Twitter. However, because of the new changes in X's data policies, which make it impossible to access the social network without a registered account, many instances of Nitter have been either closed or limited. Since all the APIs or the pre-built packages that allowed communication with Nitter relied on limited instances, we had to scrape the data manually.

First of all, we found an instance that was still fully working and then, after using Selenium and studying the structure of the web page, we managed to automatically click the "load more" button present in the page. Although the Nitter instance was not limited by any rate, it only made available data up to November 2022. Despite this little inconvenience, we still managed to collect 62,624 tweets, which was enough to conduct our research. The tweets collected are all the ones made between November 2022 and April 2024 containing the word "femminicidio" (femicide) either in the tweet's text, hashtags, or username. Our scraping program collected the username, the text, the hashtags, the dates, the number of likes, comments, and retweets for each tweet. Luckily, this data set of more than 60,000 tweets allowed us to split the data exactly as we wanted (more details will be provided in the next section), creating two balanced corpuses for the "before" and "after" periods.

In the end, since our research aims to analyze the effect that the event that presumably polarized public opinion, we decided to remove from the dataset all the tweets and comments coming from mass media, newspapers, radios, tv channels and local information services. This step was conducted using a mixture of machine classification and human annotation. More in detail, ChatGPT was used to point out the users who could potentially belong to one of the aforementioned institutions, based on some criteria (for example the presence of the word "Tv" in the username). Later, to avoid mistakes and the presence of "false positives" we manually annotated only the users that were labeled as "media", therefore there is a possibility that some users belonging to organizations are still present in the final dataset. However, we conducted one last check based on a simple heuristic: the users with high number of tweets might probably be institutions or media channels. Checking the accounts with the highest posting stats we then removed the ones that turned out not to be private users. With this step we concluded our data collection and moved on with the preprocessing.

## Methodology

### Exploratory Data Analysis

As the first step of our research, an initial exploratory data analysis is essential for understanding the material collected. Simple visualization techniques, as well as summary statistics, may give us useful information that can be later used.

Due to the nature of mass media, we believed that the distribution of the number of tweets would be greatly linked to the notoriousness of the femicide cases that are covered by news. Plotting the number of tweets by their date reveals that their volumes do indeed spike in the period right after the coverage of gender-related violence cases (Figure 1). This information provides us with a useful starting point for choosing the right polarizing event that may shift the public perception of the phenomenon.

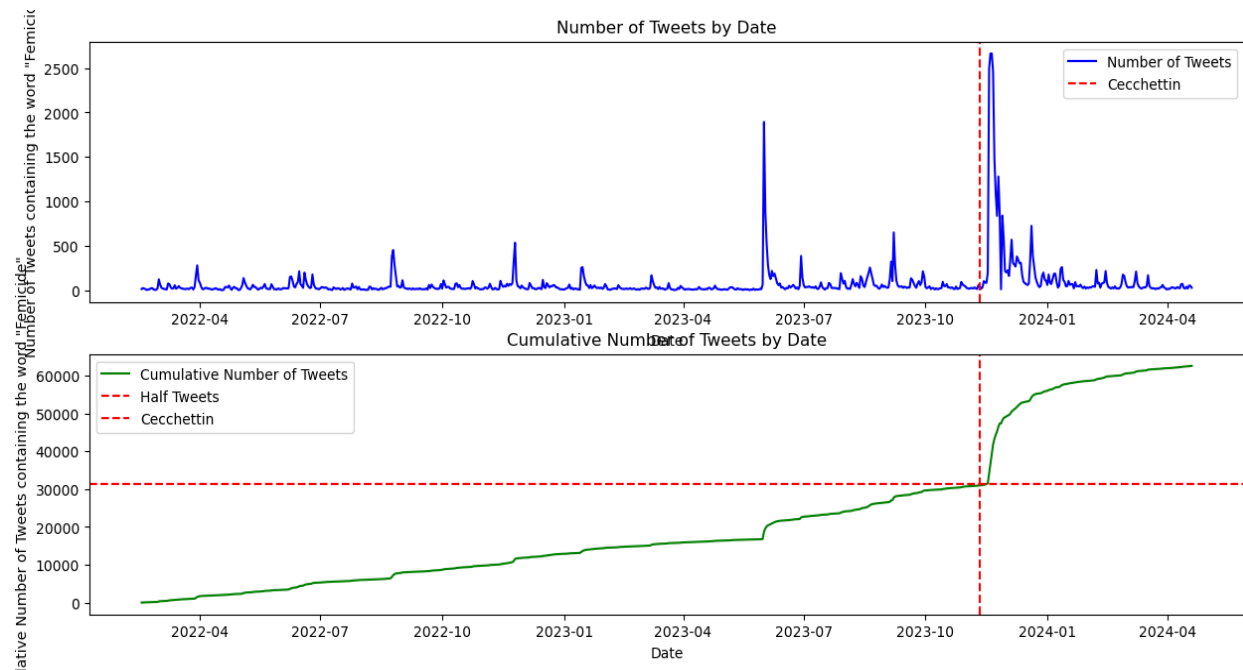


Figure 1: Number of tweets containing the word 'femicide' sparks right after the news coverage of famous gender-related violence cases

Among the femicide cases of the past years, the murder of Giulia Cecchettin in November 2023 is undoubtedly the most notorious one. Countless debates and protests sparked both online and in the streets of Italy, with people arguing on which are the systemic mechanisms that may favor the occurrence of such gender-related violence episodes. The public response for this case appeared to be much stronger than for previous murders, and it covered a great number of aspects of this phenomenon, including gender responsibilities, the role of news media in shaping the public perception, the institutions' inactions. An analysis of how language has evolved after this case could therefore provide useful insights on the actual public perception of femicides.

Despite the significant volume of tweets after the polarizing event, imbalances in data may represent a significant threat to the robustness of our results. After splitting the data into two different corpora containing tweets prior and post the polarizing event, respectively, it appeared they both contained a similar number of around 30,000 tweets each.

A second threat to our research was imposed by the representativeness of our sample. The number of individual users may give us a hint on how diverse is the user-base we are analyzing. For the

62,624 tweets we scraped, there are more than 23,000 different users, the great majority of which have only produced one tweet. A few of them have done between 2 and 4 tweets on the topic, and even less are those that have made 5 or more tweets. We believe these numbers suggest that the sample we collected is representative enough of the population.

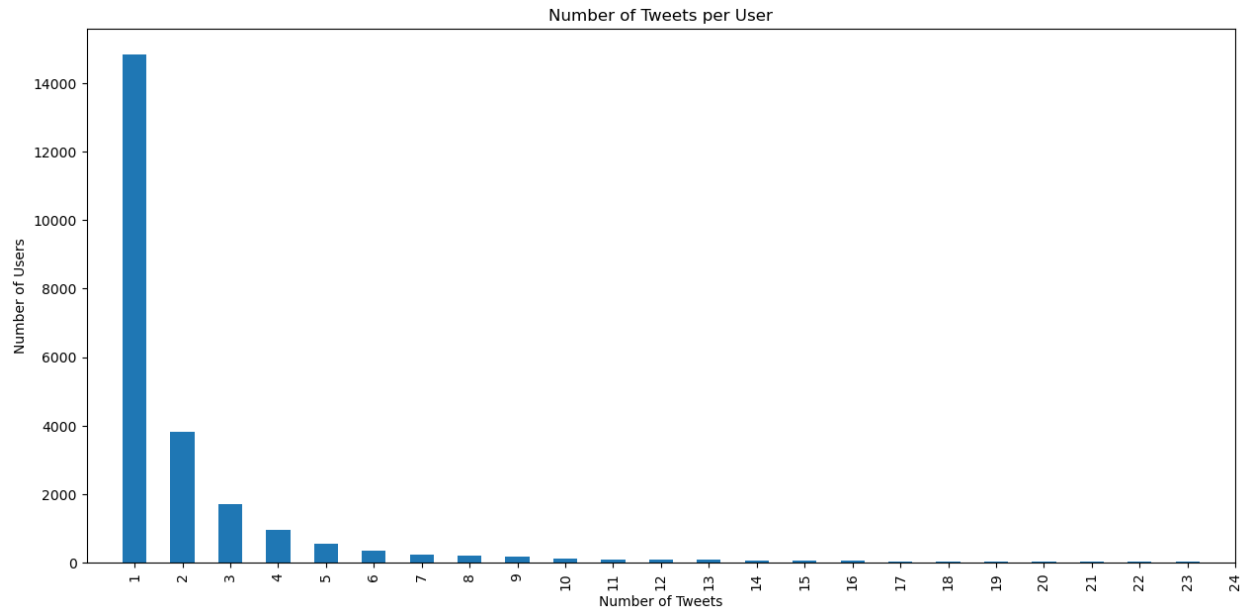


Figure 2: Great majority of users have done only 1 tweet containing the word 'femicide'. Outliers are present, but their number is non-significant

## Data Cleaning & Preprocessing

Due to the nature of tweets, many of them contain undesired characters, as well as different spelling techniques and punctuations that may not be identified properly by the tokenizers we will use in the Word2Vec models. For this reason, we defined different functions that aim at removing special characters, such as “#”, “@”, as well as formulas like “https”, which do not contain relevant information. Furthermore, we removed all words whose length is less than 3, in an effort to get rid of misspellings, acronyms and other uninformative terms. Lastly, we used the ‘stopwords’ function from the NLTK package.

Since our research aims to analyze the effect of chosen polarizing event on the public perception of this phenomenon, we decided to remove from the dataset all the tweets and comments coming from mass media, newspapers, radios, tv channels and local information services. This step was conducted using a mixture of machine classification and human annotation. More in detail, ChatGPT was used to point out the users who could potentially belong to one of the aforementioned institutions, based on some criteria (for example the presence of the word “Tv” in the username). Later, to avoid mistakes and the presence of “false positives” we manually annotated only the users that were labeled as “media”, meaning it is possible that some users belonging to organizations are still present in the final dataset. However, we conducted one last check based on a simple heuristic: the users with high number of tweets might probably be institutions or media channels. Checking the accounts with the highest posting stats we then

removed the ones that turned out not to be private users. With this step we concluded our data preprocessing and moved on with training the models.

## Word Embeddings Models

The main part of our research consists of two different Word2Vec models, trained on the corpuses we previously defined. Word embeddings created through those models may allow us to investigate how the semantic meaning of specific words has changed over time. More specifically, we want to understand if the polarizing event we identified has caused a shift in the public perception on the phenomenon of femicides.

To fully utilize the embedding models, the first step involves extracting the relevant word vectors from the matrix generated by Word2Vec. We did this by creating a function that returns a numpy array for each word embedding, as well as two dictionaries that associate each word with its respective matrix row and vice versa. Following this, we developed another function to identify the 'k' nearest neighbors of a given word by computing the dot product between the embeddings of the query word and those of other words. By combining these two functions we can thus analyze what are the words that are semantically most similar to a given query, for a specific Word2Vec model.

For the purpose of this study, we intend to compare the semantic evolution of words over time, which involves a comparison of word embeddings across the two models. However, the embeddings for a specific words are created separately by the Word2Vec we trained, so the geometrical positions represented by their vectors are not directly comparable.

To overcome this obstacle, we refer to the work of *Hamilton et al. 2016*, who used the orthogonal Procrustes problem. Specifically, Procrustes analysis is a set of methods that aim at comparing the shapes of two objects after their optimal superimposition. Their positioning and shape undergo adjustments through translation, rotation, and uniform scaling to maximize their similarity while minimizing the Procrustes distance between the objects. In the orthogonal Procrustes problem, given two different matrices A and B, we find the orthogonal matrix mapping their shapes by applying singular value decomposition (SVD). Consequently, we coded three distinct functions, each serving a specific purpose. Firstly, one function learns the optimal rotation matrix through SVD. Secondly, another function intersects the vocabularies of the two matrices and reindexes them accordingly. Finally, a third function utilizes the previous two functions to yield the newly aligned matrices. Thanks to this procedure, we are thus able to make significant comparisons of words' semantic similarities across models.

The public's perception of gender-related violence is a complex and multifaceted phenomenon. To truly understand it, we devised different sets of words that are related to specific hot topics, which may capture the different dimensions of femicides. Our intention is to understand how the terms contained in such lists differ across the two models, in an effort to infer possible cultural shifts after the polarizing event. This is not done not only by comparing their similarity directly, but by also analyzing which are their nearest neighbors for each of the two models, respectively. The sets of words considered are the following:

- Victim related words: 'vittima', 'assassinata', 'innocente'. They may reveal which are the common associations with the victim of femicides, and how the public perceives and describes them.
- Murderer related words: 'omicida', 'assassino', 'mostro', 'colpevole', 'delinquente', 'criminale', 'pazzo'. As with the previous case, such words are indicative of the general associations that the killers are subject to. Moreover, we included terms like 'mostro' and 'pazzo', as it is often argued that newspapers tend to downplay the frequency of such extreme cases of violence. This is done through the portrayal of the killer as a 'mad man', as an anomaly that is not representative of the norm.
- Positive emotional words: 'passione', 'felici', 'amore', 'rispetto'. Their similarities may point to the presence of sensationalist framing techniques which focus on the tragedy of the events
- Negative emotional words: 'gelosia', 'tradimento', 'rifiuto', 'vendetta', 'rabbia'. Like in the previous example, they may indicate sensationalist narratives. However, they also implicitly justify or downplay the gravity of femicide by portraying it as a consequence of strong emotions
- Crude words: 'sangue', 'cadavere', 'strangolata', 'violenza', 'tortura'. Victims are often dehumanized by referring to their body parts or to other macabre details. The evolution of their similarity may reveal how such techniques have evolved over time
- Modal words: 'dovere', 'necessario', 'bisogno', 'occorre', 'urgente', 'fondamentale'. They may reveal what actions or changes are deemed as necessary by users and if such needs shift their focus over time
- Patriarchy words: 'uomo', 'patriarcato', 'oppressione', 'discriminazione', 'sessismo'. The center of debates surrounding the Cecchettin murder was the concept of patriarchy, in all its aspects, and how it shapes current society. Significant changes over time are therefore expected. Û

For each of these words, we use the previously defined functions to understand if their semantic meaning changed before and after the polarizing event. Furthermore, we will run within model comparisons to see what their nearest neighbors are, for each of the two models. After that, we overlap their 300 most similar words per model to see if there are any intersections.