



INSTITUTO  
UNIVERSITÁRIO  
DE LISBOA

---

**Anti-Russia or Anti-Ukraine: how do Twitter users feel about the ongoing conflict between August 2022 and February 2023? A Sentiment Analysis Approach**

Márcia Pico Rodrigues

Master in Data Science

Supervisors:

Professor Ph.D. Diana E. Aldea Mendes, Associate Professor,  
Iscte - University Institute of Lisbon

Ph.D. Cátia Miriam Costa, Integrated Researcher,  
Iscte - University Institute of Lisbon

October, 2023



---

Department of Quantitative Methods for Management  
and Economics

Department of Information Science and Technology

**Anti-Russia or Anti-Ukraine: how do Twitter users feel  
about the ongoing conflict between August 2022 and  
February 2023? A Sentiment Analysis Approach**

Márcia Pico Rodrigues

Master in Data Science

Supervisor:  
Professor Ph.D. Diana E. Aldea Mendes, Associate  
Professor,  
Iscte - University Institute of Lisbon

Ph.D. Cátia Miriam Costa, Integrated Researcher,  
Iscte - University Institute of Lisbon

October, 2023



## **Acknowledgments**

I would like to express my heartfelt gratitude to Professor Diana Aldea Mendes and Professor Cátia Miriam Costa for their valuable guidance and constructive feedback throughout the past year of my research journey.

To my family, with a special mention to my parents and my sister, Matilde, for always teaching me that with integrity, respect, humility, and hard work we can achieve anything we set our minds to. Above all else, they represent a constant reminder that love and family stand as the most treasured values of all.

To my partner, Pedro, whose unconditional support and encouragement have been a constant presence in my life since the day we first met.



## Resumo

A investigação apresentada nesta tese teve como objetivo analisar a evolução do sentimento dos utilizadores do Twitter face ao conflito Ucrânia-Rússia entre agosto de 2022 e fevereiro de 2023. Para melhor compreender esta evolução de sentimento e da opinião pública, pesquisámos literatura relativa às relações entre a Ucrânia e a Rússia desde 1991, o ano da dissolução da União Soviética. Utilizando uma combinação de técnicas de análise descritiva, Análise de Sentimento, Topic Modelling e algoritmos de Machine Learning, como Regressão Logística, Árvore de Decisão, Naïve Bayes, AdaBoost e XGBoost, analisámos a evolução dos sentimentos Anti-Ucrânia e Anti-Rússia expressos pelos utilizadores do Twitter durante o segundo semestre do conflito. Concluímos que, dentro dos nossos conjuntos de dados, existe uma maior prevalência de tweets que expressam sentimentos Anti-Ucrânia em comparação com sentimentos Anti-Rússia. O modelo XGBoost apresentou as melhores métricas de performance, com uma taxa de accuracy de 90% para o dataset com dados de agosto e setembro de 2022 e uma taxa de accuracy de 93% para o dataset com dados de fevereiro de 2023.

**Palavras-chave:** Análise de sentimento; Topic Modelling; Machine Learning; Twitter; Anti-Ucrânia; Anti-Rússia





## Abstract

The research presented in this thesis aimed to investigate the shifting sentiment among Twitter users regarding the Ukraine-Russia conflict between August 2022 and February 2023. To comprehend this sentiment variation and public opinion, we travelled back to 1991, the year of the Soviet Union's dissolution, and reviewed literature to gain deeper insights into the Ukraine-Russia relationship. Employing a combination of descriptive analysis techniques, Sentiment Analysis, Topic Modelling, and Machine Learning algorithms such as Logistic Regression, Decision Tree, Naïve Bayes, AdaBoost, and XGBoost, we examined the evolving Anti-Ukraine and Anti-Russia sentiments expressed by Twitter users during the second semester of the conflict. Our findings revealed that, within our datasets, there was a higher prevalence of tweets expressing Anti-Ukraine sentiments than those expressing Anti-Russia sentiments. Notably, the XGBoost model exhibited the most promising performance metrics, achieving an accuracy rate of 90% for the dataset with data from August and September 2022 and 93% accuracy for the dataset with data from February 2023.

**Keywords:** Sentiment Analysis; Topic Modelling; Machine Learning; Twitter; Anti-Ukraine; Anti-Russia



## List of Contents

|  |      |
|--|------|
| Acknowledgments .....  | i    |
| Resumo .....   | iii  |
| Abstract.....  | v    |
| List of Contents .....   | vii  |
| List of Figures.....   | xi   |
| List of Tables .....   | xiii |
| Glossary of Acronyms .....   | xv   |
| Chapter 1 .....  | 1    |
| Introduction .....   | 1    |
| Chapter 2 .....  | 3    |
| Context .....  | 3    |
| 2.1. The years after the Soviet Union’s dissolution in 1991 .....        | 3    |
| 2.2. The War in 2014 .....   | 5    |
| 2.3. The Ukraine-Russian relationship after the Minsk-II agreement ..... | 6    |
| 2.4. Volodymyr Zelenskyy election.....                                   | 7    |
| 2.5. The 2022 Conflict .....   | 7    |
| 2.6. Conclusion .....  | 11   |
| Chapter 3 .....  | 13   |
| Literature Review .....  | 13   |
| 3.1. The Network Society .....   | 13   |
| 3.2. Social Media .....  | 13   |
| 3.3. Sentiment Analysis .....  | 14   |
| 3.4. Related Work .....  | 15   |
| Chapter 4 .....  | 25   |
| Methodology .....  | 25   |
| 4.1. Descriptive Analysis .....  | 25   |
| 4.2. Data and Text Pre-processing .....                                  | 28   |
| Chapter 5 .....  | 31   |
| Results and Discussion.....  | 31   |
| 5.1. Sentiment Analysis using VADER.....                                 | 31   |
| 5.2. WordCloud.....  | 32   |

|                            |    |
|----------------------------|----|
| 5.3. Topic Modelling ..... | 38 |
| 5.4. Models .....          | 42 |
| 5.5. Conclusion .....      | 48 |
| Chapter 6 .....            | 49 |
| Conclusions .....          | 49 |
| References .....           | 51 |
| Appendix Index .....       | 59 |
| Appendix .....             | 60 |





## List of Figures

|   |    |
|---|----|
| Figure 4.1: Methodology Flowchart.....  | 25 |
| Figure 4.2: Number of tweets per day in August and September 2022 .....       | 26 |
| Figure 4.3: Number of tweets per day in February 2023 .....                   | 26 |
| Figure 4.4 - Number of tweets per location in August and September 2022 ..... | 27 |
| Figure 4.5 - Number of tweets per location in February 2023: .....            | 27 |
| Figure 5.1: Number of tweets per sentiment (2022 dataset) .....               | 32 |
| Figure 5.2: Number of tweets per sentiment (2023 dataset) .....               | 32 |





## List of Tables

|   |    |
|---|----|
| Table 3.1: List of languages .....  | 15 |
| Table 3.2: List of Pre-processing methods .....   | 16 |
| Table 3.3: List of Feature Extraction methods .....   | 17 |
| Table 3.4: List of ML algorithms used .....   | 18 |
| Table 3.5: List of DL algorithms used.....  | 19 |
| Table 3.6: Best results from the literature .....   | 19 |
| Table 4.1: Example of a tweet before and after data pre-processing .....  | 28 |
| Table 4.2: Example of a tweet before and after tokenization.....  | 28 |
| Table 4.3: Example of a tweet before and after stemming.....  | 29 |
| Table 4.4: Example of a tweet before and after lemmatization.....   | 29 |
| Table 5.1: Most frequent words for each sentiment on both datasets .....  | 33 |
| Table 5.2: Most frequent words and their co-occurrences for the ‘Anti-Ukraine’<br>sentiment (2022 dataset)..... | 33 |
| Table 5.3: Most frequent words and their co-occurrences for the ‘Anti-Ukraine’<br>sentiment (2023 dataset)..... | 34 |
| Table 5.4: Most frequent words and their co-occurrences for the ‘Anti-Russia’ sentiment<br>(2022 dataset) ..... | 35 |
| Table 5.5: Most frequent words and their co-occurrences for the ‘Anti-Russia’ sentiment<br>(2023 dataset) ..... | 36 |
| Table 5.6: Days with more tweets and the most frequent words.....   | 37 |
| Table 5.7: Highest-scoring words for each topic on the 2022 dataset.....  | 39 |
| Table 5.8: Highest-scoring words for each topic on the 2023 dataset.....  | 40 |
| Table 5.9: Performance Metrics for Decision Tree .....  | 43 |
| Table 5.10: Performance Metrics for Logistic Regression.....  | 44 |
| Table 5.11: Cross-validation values for the 2022 dataset.....   | 45 |
| Table 5.12: Cross-validation values for the 2023 dataset.....   | 45 |
| Table 5.13: Performance Metrics for Naïve Bayes .....   | 45 |
| Table 5.14: Performance Metrics for AdaBoost.....   | 46 |
| Table 5.15: Performance Metrics for XGBoost.....  | 47 |



## Glossary of Acronyms

|                 |   |
|-----------------|---|
| <b>ADABoost</b> | Adaptive Boosting                                       |
| <b>BERT</b>     | Bidirectional Encoder Representations from Transformers |
| <b>BoW</b>      | Bag of Words  |
| <b>CIS</b>      | Commonwealth of Independent States                      |
| <b>CNN</b>      | Convolutional Neural Networks                           |
| <b>CV</b>       | Cross-validation  |
| <b>DL</b>       | Deep Learning   |
| <b>DT</b>       | Decision Tree   |
| <b>ETC</b>      | Ensemble Tree Classifier                                |
| <b>EU</b>       | European Union  |
| <b>GRU</b>      | Gated Recurrent Networks                                |
| <b>KNN</b>      | K-Nearest Neighbors                                     |
| <b>LDA</b>      | Latent Dirichlet Allocation                             |
| <b>LR</b>       | Logistic Regression                                     |
| <b>LSTM</b>     | Long Short-Term Memory                                  |
| <b>ML</b>       | Machine Learning  |
| <b>MLP</b>      | Multilayer Perceptron                                   |
| <b>NATO</b>     | The North Atlantic Treaty Organization                  |
| <b>NB</b>       | Naïve Bayes   |
| <b>NLP</b>      | Natural Language Processing                             |
| <b>RF</b>       | Random Forest   |
| <b>RNN</b>      | Recurrent Neural Network                                |
| <b>SGD</b>      | Stochastic gradient descent                             |
| <b>SVM</b>      | Support Vector Machine                                  |

|                |  |
|----------------|--|
| <b>SWIFT</b>   | Society for Worldwide Interbank Financial Telecommunications |
| <b>TF-IDF</b>  | Term Frequency-Inverse Document Frequency                    |
| <b>UK</b>      | United Kingdom   |
| <b>US</b>      | United States  |
| <b>VADER</b>   | Valence Aware Dictionary and sEntiment Reasoner              |
| <b>XGBoost</b> | Extreme Gradient Boosting                                    |



# Chapter 1

## Introduction

February 24<sup>th</sup>, 2022, will always be drawn in history as the day when a full-scale invasion of Ukraine by Russia marked the beginning of a conflict that continues to impact Eastern Europe. While this war is primarily between these two countries, its repercussions are felt on a global scale. The ongoing conflict has led to multiple consequences, including a surge in Ukrainian refugees, significant human casualties and injuries, an economic crisis, and food shortages, to name a few. Across the world, individuals have turned to social media platforms, particularly Twitter, not only to stay informed about developments in Ukrainian territory but also to share their sentiments and opinions concerning this conflict. As a result, Twitter became a valuable resource for conducting Sentiment Analysis, enabling us to examine how people's feelings towards the war evolve across different locations and over time.

Sentiment Analysis, also known as Opinion Mining, is a NLP technique that intends to “analyse people’s sentiments or opinions toward entities such as topics, events, individuals, issues, services, products, organizations, and their attributes” (Yue et al., 2019). Bobichev et al. (2017) state that “sentiment analysis aims to determine the attitude of speaker or writer with respect to some topic or the overall contextual polarity of a document”. This field of NLP is closely associated with ML, DL, Data Mining, and Text Mining techniques. Its practical applications are diverse, ranging from gaining insights into business and marketing strategies to assessing public opinion and sentiment on various matters.

The primary objective of this thesis is to conduct Sentiment Analysis on two distinct datasets consisting of tweets. The first dataset comprises tweets collected in August and September 2022, while the second dataset contains tweets from February 2023. The purpose of this analysis is to discern how Twitter users' sentiment regarding the conflict has evolved over a six-month period.

To accomplish our research goal, we initiated our study with descriptive data analysis. This phase allowed us to gain a comprehensive understanding of our datasets, their composition and characteristics. Following this analysis, we proceeded with the pre-processing of our unstructured data, preparing it for use as inputs in our chosen classification models. During the pre-processing phase, we employed a series of essential data transformation techniques, including tokenization, stemming, lemmatization, lowercasing, removal of stop words, punctuation, and special characters. For sentiment labelling, we employed VADER (Valence Aware Dictionary and sEntiment Reasoner). After careful examination and consideration, we categorized our data into two primary sentiments: "Anti-Ukraine" and "Anti-Russia". To gain further insights into our data and identify prevalent topics and themes, we utilized two additional techniques: the WordCloud library and Latent Dirichlet Allocation (LDA), a Topic Modelling approach. These tools allowed us to analyse visually and quantitatively the most frequently occurring words associated with each sentiment category. Lastly, we performed a set of

classification models to classify tweets into "Anti-Russia" and "Anti-Ukraine" sentiments. The chosen ML models for this classification task were LR, DT, NB, AdaBoost, and XGBoost. Then, we evaluated and analysed the performance results of these models to draw conclusions regarding their effectiveness in accurately classifying tweets based on sentiment.

This thesis is structured as follows: Chapter 2 describes the decisive events leading up to the 2022 Ukraine-Russia conflict, providing crucial background information. In Chapter 3, we review existing literature relevant to our research. Chapter 4 outlines the steps involved in our descriptive analysis and the data pre-processing. Chapter 5 details the techniques and tools employed in our analysis, including VADER, WordCloud, LDA, and our chosen classification models. Chapter 6 presents the key findings and conclusions drawn from our research. Chapter 7 provides the list of references cited throughout the thesis and is followed by the Appendix, which serves as the repository for supplementary materials that enhance the understanding of the research.

## Chapter 2

### Context

#### 2.1. The years after the Soviet Union's dissolution in 1991

With the dissolution of the Soviet Union in December 1991, Ukraine became an independent country. After its independence, the country tried to establish its own Ukrainian identity by the “renunciation of all links with Russia” and by moving its orientation westward (Odey & Bassey, 2022). Ukraine, as a sovereign nation, began to detach from Russia and tried to strengthen “its partnership relations with NATO and with individual members of this organization to provide external assurance of national sovereignty” since the former Soviet Union was still “considered a severe security threat to Ukrainian independence” (Odey & Bassey, 2022). Even when most Ukrainian analysts viewed the integration of their country into the EU as “highly appealing”, the country could not isolate itself entirely from Russia.

During the first years of independence, Ukraine's reliance on Russian oil, gas, industrial production, and economy prevented it from severing ties with Russia since the Western democracies did not demonstrate availability in providing aid to Ukraine's economic reform. Besides, Ukraine and Russia kept political and military connections, due to the past both countries share, as well as “their geographic position (...) [and] trade connections” (Ratten, 2022).

For Russians, especially political and cultural elites, it was difficult to accept the independence of the former “Little Russia” and that Ukraine was no longer in Russia's sphere of influence (D’Anieri, 2023), although they had dissimilar political systems. The capital of Ukraine, Kyiv, was seen as the centre of the Russian Foundation for hundred years “in which much of Ukraine was part of the Russian Empire and Soviet Union” (D’Anieri, 2023). Odey and Bassey (2022) state that “after independence from Soviet rule in December 1991, Russia’s imperialistic revisionist policies also worried Ukrainians, as Russian politicians (...) have rejected or shown insufficient respect for Ukrainian independence”. The Russian aim to keep Ukraine on its territory is supported by several geopolitical reasons, for example, the natural resources of oil and gas, which represent a strategic advantage “for Russia being a non-European Union member to have control of Ukraine”, since the EU citizens represent “a large percentage of consumers” (Ratten, 2022). A democratized Ukraine was seen as an “anti-Russian weapon” (D’Anieri, 2023) because it meant that Ukrainians, as happened with several former Soviet republics, would more likely turn towards democracy and western Europe and its organizations.

After three years of independence, both countries were able to maintain peace and “word out their differences” and “Russia was forced to accept the reality that Ukraine was developing into a fully sovereign state” (Odey & Bassey, 2022), despite their disparities persisted. In 1997 both countries signed the ‘Friendship, Cooperation, and Partnership’ treaty with the intent to solve their issues regarding the Black Sea Fleet, a fleet with significant historical and political value for Russia. However, this good relationship



between both countries didn't last long. The Commonwealth of Independent States (CIS), a regional intergovernmental organization created in 1991, was transformed into "a single Eurasian Economic Space" (Odey & Bassey, 2022). This organization created "provisions to prohibit members of the CIS Economic Union from becoming members of other countries' economic and custom unions" (Odey & Bassey, 2022). While insisting on its sovereignty, Ukraine did not become a full member of CIS and "firmly closed its doors towards Russia and took a pro-Western orientation" (Odey & Bassey, 2022). In 1992, the transfer of Crimea to Ukraine, made official in 1954, was considered illegal by the Russian Supreme Soviet, later abolished in 1993. Hereupon, the Ukrainian citizens started to feel threatened by Russian influence in their country and, in order to defend themselves, they started an internal "policy of 'de-Russification'" and the nationalist parties gained more predominance (Odey & Bassey, 2022).

Ukraine's EU membership was a "state goal" (D'Anieri, 2023) during the presidency of Leonid Kuchma, from 1994 to 2005, and the presidency of Viktor Yushchenko until 2010. In 2005, Yushchenko promised the Council of Europe that his government would be "reorganized to add a real, rather rhetorical, dimension and content to the process of integration into the European Union" (D'Anieri, 2023). During this year, the EU was determined to improve its relationship with Ukraine, promoting "stability, security and well-being (...) shared values, joint ownership and differentiation" (D'Anieri, 2023), however, there was no mention of membership. Since there was no hope for membership, Ukraine didn't feel the need to "generate the sacrifices needed for far-reaching reform" and, on the other hand, the EU "could not make a commitment to a country that needed so much reform" (D'Anieri, 2023). Along with this, there were also concerns regarding Russia's response. Furthermore, Ukraine was more determined to complete its "accession to the World Trade Organization (WTO)" (D'Anieri, 2023) than to the EU. The relationship between Ukraine and the EU only progressed in 2009 with the establishment of the Eastern Partnership, an agreement between the member states and other Eastern European countries, such as Ukraine, Moldova, and Georgia.

Regarding membership of Ukraine in NATO, President Viktor Yushchenko "reinstated the goal of NATO" and prepared the country for possible membership. However, D'Anieri (2023) argues that the president "was not only ahead of most NATO members. He was also ahead of the Ukrainian electorate". The Ukrainians were not supportive of NATO's membership, nor was the country ready, as Yanukóvytch stated in 2006 when visiting Brussels. Despite "there was no immediate danger of Ukraine joining (...), the Russians could not accept being surrounded" by the organization in all its borders (Noonan, 2023). On the northern border, Russia is surrounded by the Baltics, Canada, and the US, westward there is Poland, eastward there is the US and with the possibility of Ukraine joining NATO, Russia would be surrounded on its southern border as well. The Russian opposition to NATO's enlargement is more related to geopolitical interests than the threat of a growing militarization, since "Russia's most important demand of NATO is that Ukraine and the Baltic States be recognized as a zone of Russia's vital interests" (Odey & Bassey, 2022).

## 2.2. The War in 2014

Viktor Yanukóvytch won the presidential election in Ukraine in 2010. Ukraine's society was characterized as a "lopsided economy dominated by oligarchs, heavy reliance on Russia" (McMahon, 2014). After being elected, Yanukóvytch "sought to permanently eliminate competition for power" (D'Anieri, 2023) against citizens' and elites' will. During this period, the president took "control of the country's Constitutional Court", "invalidate[d] the crucial "pact" limiting presidential power" and "used other illegal means to forge a majority in the parliament" (D'Anieri, 2023). There was also a lot of corruption and strong divergences between East and West Ukraine since the eastern zone is more industrialized and has several Russian populations and the rest of the territory is more western oriented.

Yanukóvytch favored keeping solid ties with Russia while continuing its conversations with the EU on "a trade association agreement" (McMahon, 2014). Russia pressured the Ukrainian president to resign its agreement with the EU. After dropping off from that agreement in 2013, in favor of Ukraine's relationship with Russia, Ukraine's parliament "voted unanimously to remove [Yanukóvytch] from office" (O'Connell 2017). By the end of January 2014, there was an outbreak of protests throughout the country, especially in Kyiv. These protests intensified and there were several deaths. Faced with this situation, Yanukóvytch tried to flee Ukraine and "requested Russian military assistance in restoring order" (O'Connell, 2017).

In February, Russia helped Yanukóvytch flee and moved its armed forces to Crimea, the port city of Sevastopol, and the regional capital of Simferopol. The peninsula of Crimea is characterized by its "commercial and military significance" (D'Anieri, 2023). Besides its historical relevance (the city of Sevastopol was labelled "Hero City" during World War II), this region is known for its coal mining industry, which would benefit Russia's economy, and its access to the Black Sea, extremely important for Russia's military and economic strategy, because it represents Russia's sole entry point to a temperate sea region. This territory was transferred to Ukraine in 1954 by the Soviet leader Nikita Khrushchev, however, after the dissolution of the Soviet Union, "Russians had always viewed Crimea as the most humiliating loss of all the territories left outside of Russia" (D'Anieri, 2023).

In March, Crimea's parliament requested a referendum for the peninsula's future and the citizens "opted overwhelmingly for union with Russia" (McMahon, 2014). Nevertheless, O'Connell (2017) argues that previously the 'Berkut', Ukrainian Special Police, and pro-Russian militias "took over the Crimean parliament building, holding it until MPs voted in favour of seceding from Ukraine". Then, Russia "passed a resolution nullifying Ukrainian laws in Crimea and putting in force Russian legislation" (McMahon, 2014).

On the 18<sup>th</sup> of March, Vladimir Putin announced the annexation of Crimea. The conflict between both states had officially started. Russia marked "presence in the Donbas region of eastern Ukraine" in support of pro-Russian separatists' militias (O'Connell, 2017). The EU and the United States (US) supported Ukraine with economic, energy and

military aid and sanctions on Russia. As the conflict continued, Ukraine registered human losses and feared the "spread beyond the Donbas to central Ukraine" (O'Connell, 2017). In September, Ukraine reached a ceasefire with Russia.

After the ceasefire, there was an agreement that "lasted about a week before fighting resumed", known as Minsk I (O'Connell, 2017). On January 2015, despite the Minsk I agreement, the fighting continued. In February, Ukraine and Russia agreed on Minsk II, intending to "provide comprehensive roadmap to peace in eastern Ukraine" (O'Connell, 2017). O'Connell (2017) argues that the agreement was a "complicated set of interrelated stages that have been interpreted differently by the parties" and was seen as "heavily favouring Russia".

### **2.3. The Ukraine-Russian relationship after the Minsk-II agreement**

The signing of the Minsk agreements marked the intent of ending the Donbas conflict between Ukrainian armed forces and the Pro-Russian separatist groups. D'Anieri (2023) claimed that Ukraine "signed the Minsk agreements at gunpoint", however, the conflict in the eastern region of Ukraine decreased in the intensity of hostility. Disregarding the signing of these agreements, both countries still faced a lot of disagreements, especially, "about the order in which various steps would be carried out" since no one wanted to be left feeling "exploited" (D'Anieri, 2023). The seizure of Crimea, which was one of the triumphs for Russia, "did not satisfy its ambitions regarding Ukraine" (D'Anieri, 2023). D'Anieri (2023) claimed that "there is a big difference between a peace agreement that leads to lasting peace and one like the Minsk agreements that leads to a new war". In an interview to *Die Zeit* newspaper, the former German chancellor argued that "the 2014 Minsk agreement was an attempt to give Ukraine time" to rearm the country (Schwarz, 2022). The signing of the "Minsk disagreements", as some authors call it, failed because both parties wanted to succeed on their terms and neither of them wanted to compromise.

In 2015, Ukraine created a "new set of laws on "decommunization"" (D'Anieri, 2023) to distance itself from Russia. These laws were related to the ban of communist and national socialist propaganda, which meant the "banning of the Communist Party of Ukraine" and was seen "as impeding freedom of speech" (D'Anieri, 2023), and the reform on education, which meant that all Ukrainian students needed to improve their skills in the Ukrainian language, and they were limited to using the Russian language exclusively within their household.

Ukraine's government tried to increase and strengthen its relationship with the EU. The EU-Ukraine Association Agreement, previously resigned by Yanukóvytch, came into full effect in September 2017. This pact establishes the relationship between Ukraine and the rest of the EU's member states. The content of this agreement focused on "a program of technical and financial assistance focused on decentralization and anti-corruption/rule of law support, as well as harmonization to bring Ukraine's trade rules in line with those of the EU" (D'Anieri, 2023). The result was that Russia's share of Ukrainian exports significantly decreased, while the EUs almost doubled it.

The cooperation with NATO also expanded. The support for Ukraine's membership increased throughout the country since 2014. In 2016, Ukraine and NATO settled on a new "Comprehensive Assistance Package" intending to support Ukraine's Armed Forces and facilitate the "essential reforms, in particular in the security and defence sector" (Comprehensive Assistance Package for Ukraine). In 2017 Ukraine's parliament announced that NATO membership was "an objective of foreign policy" and this objective was expressed in Ukraine's constitution by 2019 (D'Anieri, 2023). In 2020 Ukraine, alongside Bosnia and Herzegovina, Finland, and Sweden, became a candidate member of NATO.

#### **2.4. Volodymyr Zelenskyy election**

After Volodymyr Zelenskyy won the 2019 presidential election, he claimed that Ukraine will "continue in the direction of the Minsk [peace] talks and head towards concluding a ceasefire", ensuring the continuation of the best relationship possible "with the Russia-backed separatists" (D'Anieri, 2023). He also assured strengthening the ties between Ukraine, NATO, and the EU. However, Vladimir Putin refused to give Zelenskyy "the traditional congratulatory statement" and announced that all Ukrainian citizens who lived in the Donbas areas would receive Russian passports (D'Anieri, 2023).

The two presidents finally met at the end of 2019 and they both agreed on "several potentially important measures, including a complete ceasefire by the end of the year, a plan to clear land mines, identification of new areas in which troops would disengage, and additional prisoner exchanges" (D'Anieri, 2023). Nevertheless, these measures were never employed, and the tension persisted during the following years. In 2021, Zelenskyy was advised by the National Security Council to apply sanctions to Russian media. D'Anieri (2023) stated that Zelenskyy's spokesperson affirmed that these sanctions are "in order to protect national security". This event caused Russia to move the first troops to the border with Ukraine.

#### **2.5. The 2022 Conflict**

The war between Ukraine and Russia started, unofficially, in February 2021, when Russia started "large-scale exercises" with thousands of troops on the Ukrainian border. During this month, Russia televised a meeting with Putin and his National Security Council to discuss the situation in Donbas. The goal of this meeting "was to consider appeals by the leader of the Donetsk and Luhansk republics to recognize their sovereignty" and abandon the Minsk agreements (D'Anieri, 2023). Putin also claimed that "Ukraine was led by far-right nationalists and neo-Nazis who were waging genocide against Russians" (D'Anieri, 2023) and that the country was preparing itself for an attack against Russia by cooperating with NATO. These claims became the focus of Russian propaganda.

In August 2021, Vladimir Putin published an article on the Kremlin website, titled "On the Historical Unity of Russians and Ukrainians". Analysts around the world interpreted this article as a "Declaration of War" (Peter Dickinson), "call to arms" (Anne Appelbaum) and "final ultimatum to Ukraine" (Mikhail Rostovsky) (D'Anieri, 2023). On

the other hand, Volodymyr Zelenskyy “poked fun at Putin, saying that he was “envious that the president of such a great power can permit himself to spend so much time [writing] such a volume of detailed work”” (D’Anieri, 2023). By December, Vladimir Putin claimed that the threat on the Ukrainian border has risen, and the US Intelligence alerted that “Russia was planning to invade Ukraine in early 2022” (D’Anieri, 2023). Putin's motivations have been the same since Ukraine’s independence, adding to the threat of NATO and/or EU membership.

The leaders of Germany, France, and the UK made efforts to engage in dialogue with Putin, aiming to reassure that Ukraine would not receive assistance from NATO and that, if an invasion occurred, the Western countries would respond with economic sanctions. These sanctions "would be immediate and far-reaching potentially including banning Russia from the SWIFT payments system" (D’Anieri, 2023). Discreetly, Ukraine was preparing itself for another war without causing any panic among its citizens. The Ukrainian president "rejected claims that an invasion was imminent and criticized Western governments for claiming it was, while simultaneously desperately seeking support for them" (D’Anieri, 2023).

In the final months of 2021 and at the beginning of 2022, Russian increased military exercises near Ukraine’s borders and NATO “accused Russia of planning an invasion” (Koroutchev, 2023). The recognition of the self-proclaimed separatist states in Donbas (Donetsk and Luhansk) and the sending of troops into those territories marked the triggering of the tensions between the two countries.

Prior to the outbreak of the war, President Vladimir Putin contacted the US and NATO "with a written document outlining concrete security guarantees that they hoped to use as the basis of negotiations" (Noonan, 2023). The author adds that these guarantees were "ignored by the Biden administration" (Noonan, 2023).

D’Anieri (2023) asserts that the intentions of Russia are based on "three basic things". The first is to "regain control over most, if not all, of the territory of the former Soviet Union". Secondly, Russia wants to have a say over European security issues. Lastly, the country wants to legitimize its government and "trade freely" with others without them trying to democratize Russia. However, the main cause for the beginning was the eastward expansion of NATO, despite the organization’s “assurances explicitly given in 1991 that it would not expand east the German-Polish border, as well as explicit warnings by senior US diplomats that all Russian political forces (...) would read that expansion as threatening, NATO has continued to expand towards Russia” (Noonan, 2023). From the economic point of view, NATO had no benefit in extending its sphere of influence, “post-Soviet Union Russia was hardly in a position to economically dominate its ex-satellites” (Noonan, 2023). Noonan (2023) argues that the only reason that could have been for NATO to augment its influence were “*raison d’etat*”, i.e., “Cold War triumphalism, American hubris, European inertia, and paranoid fears that post-Soviet Russia would revert to some sort of Tsarist imperium”. Noonan (2023) claims that the war between Russia and Ukraine “must be considered in the context of the post-Cold War history of Europe” and that “while these political reasons are readily comprehensive from within a realist perspective”, the invasion of Ukraine is still “politically irrational”. From the Russian point of view, the decision of the country to invade Ukraine is rational, in the

sense that “the Russian leadership had a set of goals and calculated that war was more likely to achieve them than peace and would do so at an acceptable cost” (D’Anieri, 2023). Furthermore, after the outbreak of the war in 2022, many believe that Putin was “unhinged” (D’Anieri, 2023) and “delusional” (Noonan, 2023) when taking the decision to invade Ukraine.

On February 24, 2022, Russia began its "special military operation" intending to "demilitarize and denazify Ukraine” (D’Anieri, 2023). At the beginning of the invasion, Putin appealed to the cooperation of Ukrainian citizens and the military to end the war as soon as possible. He also threatened the West saying that if a country tried to threaten Russia back, he "will respond immediately, and the consequences will be such as you have never seen" (D’Anieri, 2023). The Western countries concluded that Russia was prepared to utilize its nuclear weapons if any nations attempted to intervene. Putin’s goal “was to depose the Zelenskyy government and put a puppet regime in place” (D’Anieri, 2023).

The Russian military operation in southern Ukraine was deemed successful, "by early March, Russian forces had captured the entire coast between Crimea and Russia" (D’Anieri, 2023). Ukrainian army had to move out of Crimea, towards Berdiansk, Mariupol, Zaporizhzhya, Dnipro, Kherson, Mykolayiv, and Odesa. Putin was hoping that this invasion was going to be easy; however, Ukrainian forces kept their persistence and determination and "misunderstanding the level of Ukrainian patriotism and commitment was likely Russia's biggest mistake" (D’Anieri, 2023). D’Anieri (2023) argues that since 2014 "a whole generation of soldiers and officers had gained combat experience and the army had refined its tactics". Furthermore, the citizens maintained strong support and cooperation with the Ukrainian military.

The West responded with coordinated sanctions that “are causing serious problems in the Russian economy” (Noonan, 2023) while avoiding war with Russia. The West’s goal was “to inflict such a heavy blow on Russia’s economy and on its oligarchs that pressure would build within Russia to end the war” (D’Anieri, 2023). The Western countries also helped Ukraine with military equipment since the beginning of the war and continued to do so in the subsequent months.

One month after the outbreak of the conflict, the President of Ukraine stated that Ukraine's membership in NATO "was unlikely" (D’Anieri, 2023) and showed openness to discussing the future of the sovereignty of Donetsk and Luhansk.

Peace talks started in late March and took place in Turkey. Among Russian demands there was "a commitment to neutrality", "the surrender of Crimea and Donbas", and "Ukraine downsizing its military to 50,000 troops", which meant leaving Ukraine almost “defenceless” (D’Anieri, 2023). Ukraine showed interest in maintaining the negotiations, although it rejected these demands. Countries from outside proposed “a return to the status quo of February 24, ceding to Russia both Crimea and the portions of Donbas occupied since 2015” (D’Anieri, 2023). However, neither Ukraine nor Russia was open to accepting that. Ukraine was eager to keep fighting for its territory and Russia was ready to conquer more territory and was not open to "surrendering all the territory conquered in 2022" (D’Anieri, 2023).

The ongoing war “has caused a humanitarian catastrophe” (Wadhvani et al., 2023). Millions of Ukrainian citizens fled to other European countries, “resulting in a humanitarian and political crisis” (Ratten, 2022). After a week the conflict started, according to IOM Response 2022, “more than 600.000 people fled Ukraine” (Koroutchev, 2023). Since the welcoming countries were part of the EU, this organization felt the need to create a program that gave all the Ukrainian citizens that left the country on February 24 the right “to live in the EU, (...) to work, access to social security, medicine, and education” (Koroutchev, 2023). There are also reports of war crimes and mass graves in Russian-controlled territory. The International Criminal Court statute defines war crimes committed by Russia as “crimes against humanity” (D’Anieri, 2023).

The war also strengthened “the power of Ukrainian nationalists” and made the idea of a “Russian-leading government unthinkable” (Noonan, 2023). President Volodymyr Zelenskyy has been seen as a “global spokesperson for promoting peace in Ukraine” (Ratten, 2022) and along with the “Spirit of Ukraine” was chosen as the 2022 Person of the Year by Time Magazine. On the other hand, the “weaknesses” of Russia were exposed, as well as its economy, which is more vulnerable than ever (Noonan, 2023). The author adds that “Russia is less secure politically and economically” than ever (Noonan, 2023).

Regarding the EU and the US, these regions “are likely going to face a recession (...) [and] they are at the very least no more secure than they were before the war” (Noonan, 2023). From Ukraine’s perspective, the country “is suffering most of all of these politically irrational choices” (Noonan, 2023). In June 2023, Ukraine was announced as a candidate for the EU after the submission of the application for membership by Volodymyr Zelenskyy after the beginning of the invasion. If Ukraine ever becomes an EU member, this will significantly help “transform the prospects of reform in post-war Ukraine” (D’Anieri, 2023). Thinking about the future, the main causes of concern are the economy, since the economic help provided by the West will be insufficient for a post-war country, the duration of the war, while “fighting could continue in one form or another for years”, and the idea of EU or NATO membership is unlikely (Noonan, 2023).

D’Anieri (2023) proposes three categories of countries based on their perceptions of the ongoing war. The first group consisted of the “global community of democracies” (D’Anieri, 2023), where the author includes Europe, North America, Japan, Australia, Finland, and Sweden, who were not part of NATO and applied for membership after the outbreak of the war. The second group comprises Russia, China and “a few others who clearly sided with Russia” (D’Anieri, 2023), for example, Hungary. The third and last group includes the countries of the southern hemisphere, such as India, “which sought to align itself neither with Russia's war nor with those who were fighting it” (D’Anieri, 2023).

## 2.6. Conclusion

Ratten (2022) claims that the conflict between Russia and Ukraine leaves us with several concerns for the future: "How long the conflict will persist, (...) how Russia will cope with sanctions, and (...) the impact of the crisis on the global economy".

Besides the massive number of losses from both sides of the conflict, two of the greatest consequences are food shortage and emigration. Since Ukraine and Russia were "two of the world's largest food experts" (D'Anieri, 2023), the war is threatening the food supply in various countries, especially in Africa. Regarding emigration, Ukraine is suffering from a "major population loss and a major rearrangement of the population" (D'Anieri, 2023), because the more time citizens remain in different countries, "the less likely they are to return", and for the Ukrainian population that stayed in the country, they will fight major "loss and trauma" (D'Anieri, 2023). Furthermore, Ukrainians will have to go through "a massive economy reconstruction effort", while counting on the West's support to tackle a hostile Russia (D'Anieri, 2023), which has now been pictured "as a threat to global peace and security" (Ratten, 2022).

D'Anieri (2023) claims that the Russian invasions in 2014 and 2022 "share the same set of underlying causes". Among these causes are "Russia's desire to regain control of Ukraine, its conception of itself as a great power, the security dilemma in Europe, and the impact of democratization" (D'Anieri, 2023). According to the author, these factors are as valid for Russia in 2022 as they were in 2014. However, the invasion of 2014 had precedents that are easier for analysts to understand, unlike what has been happening since 2022. D'Anieri (2023) argues that "the timing of the 2022 attack is harder to understand" because it "reflects (...) a mix of groupthink, misperception, and the personal traits of Vladimir Putin".

The author also argues that the ongoing war started "eight years earlier, in 2014, when Russia seized Crimea and attacked Donetsk and Luhansk Oblasts" and is embedded in "deep disagreements about the post-Cold War world should look like" (D'Anieri, 2023).





## Chapter 3

### Literature Review

#### 3.1. The Network Society

The emergence of technology allowed the disclosure of a new form of organizational structure where computers, network technologies, telecommunications, and nanotechnologies, for example, became a part of our society. The social, cultural, economic, and institutional transformation did not exist without this technology.

Castells (2022) defined this new social structure as the Network Society. The Network Society is "the social structure of our age, the Information Age, (...) it is a global social structure, and so it refers to all societies, albeit with extreme cultural and institutional diversity" (Castells, 2022).

Technology allowed the expansion of the network society with its Internet-based networks. Although in the digital era, the public sphere, which is "an essential component of socio-political organization" (Castells, 2008), can have various forms, the Network Society organizes itself around the media communication networks, i.e., social media, and the Internet. Wankhade et al. (2022) argued that "the growing popularity of the Internet has lifted the web to the rank of the principal source of universal information".

#### 3.2. Social Media

Öztürk and Ayvaz (2018) argue that social media "creates virtual bonds between users, in which people express opinions and develop relationships through posts, comments, messages and likes (...) allows people to share their thoughts, feelings, opinions with other people instantly and easily". Hence, social media is an essential part of people's daily routine in the Network Society. However, it has its advantages and consequences or challenges, just like Castells (2022) called attention to – "the social media became the global repository of inhuman values whose spread around the planet came to haunt us". Some of its positive aspects are quality education and collaboration between students, information sharing and updates, social movements organization and the creation of online communities. Nonetheless, social media has become the ideal place for the proliferation of cyber-harassment, cyberbullying, disinformation, fake news and hate speech, especially on social media sites where users can maintain their anonymity.

In times of conflict, social media is essential for reporting events in real-time, providing live updates and being "used for political purposes as people tend to share political opinions" (Öztürk & Ayvaz, 2018). According to Peter Wallensteen, the concept of conflict has "three main elements [...] actors, incompatibility, and action" and "contemporary conflicts are characterized by the active involvement of non-state actors" (Šerstka, 2021). These non-state actors can be governmental or non-governmental organizations, but also "terrorist organizations, insurgents, rioters, private armies,

individuals, and so forth" (Šerstka, 2021). The rise of the Internet introduced a new actor to recent conflicts, which is the user of social media. During these events, it is also expected that social media becomes a vehicle for political controversy, the spread of fake news and disinformation, propaganda campaigns, and the organization of social movements, potentially affecting the lives of thousands of citizens. Some research confirms that the role mass media plays in armed conflicts "complicates this situation", because, usually, "conflict parties and their supporters have tried to use this development [modern technology] to their advantage and turned the internet into a propaganda battleground" (Hauter, 2021). Twitter has been used for these purposes during military crises and conflicts such as the Syrian civil war, the war in Afghanistan, the Arab Spring, the annexation of Crimea by the Russian Federation back in 2014 and the invasion of Ukraine in 2022.

According to Smart et al. (2022), social media is "playing a large role in the 2022 Russian invasion of Ukraine" with "people around the world [...] turning to various social media platforms to access the latest news and express their opinion and sentiment on the Russo-Ukrainian War" (Garcia & Cunanan-Yabut, 2022). Citizens resort to these platforms to share their opinions on the conflict and, occasionally, "these expressions of diversified opinion carry support for either of the parties" (Zhu et al., 2022). Along with sharing opinions and live updates, social media platforms have also been "playing a larger role in hate crimes" (MacAvaney et al., 2019), and the current conflict is no exception. Besides, social media platforms, as well as other forms of communication in past decades, "become an integral part of warfare affecting the ways in which the conflict is perceived by the general public, [and] what decisions are made by policy-makers (...)" (Makhortykh & Sydorova, 2017).

### **3.3. Sentiment Analysis**

The emergence of social networks "has generated a slew of fields devoted to analysing these networks and their contents in order to extract necessary information". (Wankhade et al., 2022). Sentiment Analysis is one of these fields whose main aim is to "analyse the reviews and examine the scores of sentiments" (Hussein, 2018).

Sentiment Analysis uses NLP "to extract, convert and interpret opinion from a text" (Drus & Khalid, 2019). This approach is also called opinion mining and "should be treated as a branch of machine learning, data mining, natural language processing, and computational linguistics, which borrows elements from sociology and psychology" (Yue et al., 2019). Sentiment Analysis is instrumental in studying public opinion on various topics, such as politics or product reviews. However, the main obstacle is the "ambiguity of word polarity" (Wankhade et al., 2022).

There are four levels when investigating Sentiment Analysis: Document Level, Sentence Level, Phrase Level, and Aspect Level. The first one is "performed on a whole document, and single polarity is given to the whole document" (Wankhade et al., 2022) and it's used for chapter classification, for example. On the Sentence Level, "each sentence is analysed and (...) [found] with corresponding polarity" (Wankhade et al., 2022). This level is used when there is a "wide range" of sentiment within a document.

Finally, each sentence "may contain multiple aspects or single aspects" and "multiple aspects", that is why it is pertinent to analyse the sentiments on Phrase and Aspect Levels (Wankhade et al., 2022).

On Sentiment Analysis researchers have three options: follow a Lexicon Based Approach, a ML Approach, or an Ensemble Method, which is a combination of Lexicon Based and ML approaches. The Lexicon Based Approach doesn't require any training data and "the aggregation of scores of each token is performed, i.e., positive, negative, neutral scores are summed separately" (Wankhade et al., 2022). This technique aims to assign the overall polarity to the text being analysed, "based on the highest value of individual scores" (Wankhade et al., 2022). In the ML Approach, the researcher uses algorithms to label sentiments. ML algorithms are "used to predict a class label for a given instance of an unknown class" (Wankhade et al., 2022). The trained algorithms are usually classification models, "where sentiment detection is framed as a binary (i.e., positive or negative)" (Gonçalves et al., 2013), so the neutral category is often discarded. Yue et al. argue that "the neutral category is ignored under the assumption that neutral objects lie near the boundary of the binary classifier", however, there are some researchers that suggest that "in every polarity classification task, three categories must be identified" (Yue et al., 2019).

### 3.4. Related Work

A sentiment "describes an opinion or attitude expressed by an individual, the opinion holder, about an entity, the target" (Kiilu et al., 2018). Sentiment Analysis investigates the public's sentiments and opinions through "computational treatment of subjectivity in text" (Hutto & Gilbert, 2014). The literature reviewed in this article ranges from ML to more recent DL approaches, "to automatically detect hate speech in text documents" (Mutanga et al., 2022).

Regarding the languages of the datasets used in the reviewed literature, they range from English, Spanish, Italian, German, Indonesian, Sinhala, Arabic, Bangla, Russian and Ukrainian (see Table 3.1). In the reviewed literature, the analysis was conducted on tweets, Facebook comments, and news.

**Table 3.1: List of languages**

| <i>Language</i> | <i>References</i>   |
|-----------------|---|
| English         | MacAvaney et al., 2019; Davidson et al., 2017; Prasetyo & Samudra, 2022; Turki & Roy, 2022; Sultan et al., 2023; Corazza et al., 2020; Hasan et al., 2022; Cahyana et al., 2022; Mutanga et al., 2022; Marchellim & Ruldeviyani, 2021; Ruwandika & Weeransinghe, 2018; Hajibabae et al., 2022; Kiilu et al., 2018; Maruf et al., 2022 |
| Spanish         | Amores et al., 2021   |
| Italian         | Del Vigna et al., 2017; Corazza et al., 2020  |
| German          | Corazza et al., 2020  |
| Indonesian      | Putri et al., 2020  |
| Arabic          | Al-Hassan & Al-Dossari, 2022  |
| Sinhala         | Ruwandika & Weeransinghe, 2018  |

|                       |                                      |
|-----------------------|--------------------------------------|
| Afaan Oromo           | Defersha & Tune, 2021                |
| Bangla                | Hasan et al., 2023                   |
| Russian and Ukrainian | Bobichev et al., 2017                |
| Not mentioned         | Balahur, 2013; Wadhvani et al., 2023 |

When it comes to data retrieved from social media, Data Pre-Processing is crucial because "users of Social Media platforms employ a special "slang" (i.e. informal language, with special expressions, such as "lol", "omg"), emoticons, and often emphasize words by repeating some of their letters" (Kiilu et al., 2018). The goal of this process is "to obtain clean data in order to improve the accuracy of the detection process" (Prasetyo & Samudra, 2022). On Twitter, "the language employed (...) has specific characteristics, such as the markup of tweets that were reposted by other users with "RT", the markup of topics using "#" (hash sign) and of the users using the "@" sign" (Kiilu et al., 2018).

The methods employed vary between case folding, tokenization or user and topic labelling and are all presented in the following table (Table 3.2):

**Table 3.2: List of Pre-processing methods**

| <i>Pre-processing</i>         | <i>References</i>   |
|-------------------------------|---|
| Case Folding/Lowercase        | MacAvaney et al., 2019; Putri et al., 2020; Turki & Roy, 2022; Mutanga et al., 2022; Davidson et al., 2017; Prasetyo & Samudra, 2022; Marchellim & Ruldeviyani, 2021; Defersha & Tune, 2021; Hajibabae et al., 2022; Wadhvani et al., 2023; Balahur, 2013               |
| Tokenization                  | MacAvaney et al., 2019; Turki & Roy, 2022; Mutanga et al., 2022; Prasetyo & Samudra, 2022; Cahyana et al., 2022; Ruwandika & Weeransinghe, 2018; Marchellim & Ruldeviyani, 2021; Amores et al., 2021; Defersha & Tune, 2021; Hajibabae et al., 2022; Maruf et al., 2022 |
| Punctuation removal           | MacAvaney et al., 2019; Turki & Roy, 2022; Sultan et al., 2023; Al-Hassan & Al-Dossari, 2022; Ruwandika & Weeransinghe, 2018; Defersha & Tune, 2021; Wadhvani et al., 2023; Maruf et al., 2022  |
| Normalisation                 | Putri et al., 2020; Balahur, 2013   |
| Normalisation of hashtags     | Mutanga et al., 2022; Corazza et al., 2020  |
| Removal of hashtags           | Turki & Roy, 2022; Hajibabae et al., 2022; Maruf et al., 2022   |
| Removal of special characters | Mutanga et al., 2022; Prasetyo & Samudra, 2022; Cahyana et al., 2022; Sultan et al., 2023; Corazza et al., 2020; Marchellim & Ruldeviyani, 2021; Defersha & Tune, 2021; Hajibabae et al., 2022; Kiilu et al., 2018; Maruf et al., 2022                                  |
| Removal of short words        | Mutanga et al., 2022; Maruf et al., 2022  |
| Removal of stop words         | Putri et al., 2020; Turki & Roy, 2022; Prasetyo & Samudra, 2022; Sultan et al., 2023; Ruwandika & Weeransinghe, 2018; Marchellim & Ruldeviyani, 2021; Defersha & Tune, 2021; Wadhvani et al., 2023  |
| Filtering                     | Putri et al., 2020  |
| Stemming                      | Putri et al., 2020; Davidson et al., 2017; Prasetyo & Samudra, 2022; Sultan et al., 2023; Marchellim & Ruldeviyani, 2021; Wadhvani et al., 2023; Maruf et al., 2022   |
| Spelling check                | Cahyana et al., 2022; Defersha & Tune, 2021   |

|   |   |
|---|---|
| Remove of non-English text              | Cahyana et al., 2022  |
| Lemmatization                           | Cahyana et al., 2022; Ruwandika & Weeransinghe, 2018; Wadhvani et al., 2023; Maruf et al., 2022 |
| Normalisation of Arabic text            | Sultan et al., 2023   |
| Removal of repeated characters          | Sultan et al., 2023; Hajibabae et al., 2022   |
| Description of emojis                   | Corazza et al., 2020  |
| Expanding abbreviations                 | Hajibabae et al., 2022  |
| Removal of quotes                       | Kiilu et al., 2018  |
| Removal of the word RT from tweets      | Kiilu et al., 2018  |
| Removal of HTML and links               | Wadhvani et al., 2023   |
| Removal of numbers                      | Maruf et al., 2022  |
| Removal of Twitter handles (@user)      | Maruf et al., 2022  |
| Segmentation                            | Maruf et al., 2022  |
| Repeated punctuation sign normalization | Balahur, 2013   |
| Emoticon replacement                    | Balahur, 2013   |
| Slang replacement                       | Balahur, 2013   |
| Affect word matching                    | Balahur, 2013   |
| Modifier word matching                  | Balahur, 2013   |
| User and topic labelling                | Balahur, 2013   |

Feature Extraction is "a key task in sentiment classification as it involves the extraction of valuable information from the text data" (Wankhade et al., 2022). The techniques applied impact the performance of the algorithms. The methods applied to the reviewed literature range from TF-IDF, Part-of-speech, Count Vectorizer, Word Cloud Representation, Word2Vec Embeddings, BoW, Bag of Features, Word Embeddings, Emoji Embeddings, Ngrams, Emotion lexica, Glove (global vector), Transformer-based Embedding, FastText, Correlation-based Feature Subset Selection and Information Gain (see Table 3.3).

**Table 3.3: List of Feature Extraction methods**

| <i>Feature Extraction</i> | <i>References</i>  |
|---------------------------|--|
| TF-IDF                    | MacAvaney et al., 2019; Prasetyo & Samudra, 2022; Cahyana et al., 2022; Sultan et al., 2023; Ruwandika & Weeransinghe, 2018; Davidson et al., 2017; Hasan et al., 2022; Marchellim & Ruldeviyani, 2021; Defersha & Tune, 2021; Hajibabae et al., 2022; Wadhvani et al., 2023; Maruf et al., 2022 |
| Part-of-speech            | Davidson et al., 2017; Del Vigna et al., 2017; Kiilu et al., 2018  |
| Count Vectorizer          | Turki & Roy, 2022; Maruf et al., 2022  |
| Word Cloud Representation | Turki & Roy, 2022  |
| Word2Vec Embeddings       | Sultan et al., 2023; Hajibabae et al., 2022  |
| BoW                       | Sultan et al., 2023; Ruwandika & Weeransinghe, 2018; Mutanga et al., 2022; Amores et al., 2021; Wadhvani et al., 2023; Maruf et al., 2022; Bobichev et al., 2017   |
| Bag of Features           | Ruwandika & Weeransinghe, 2018   |

|  |   |
|--|---|
| Word Embeddings                            | Corazza et al., 2020; Hasan et al., 2022; Mutanga et al., 2022; Amores et al., 2021                                       |
| Emoji Embeddings                           | Corazza et al., 2020  |
| N-grams                                    | Corazza et al., 2020; Putri et al., 2020; Defersha & Tune, 2021; Wadhvani et al., 2023; Maruf et al., 2022; Balahur, 2013 |
| Emotion lexica                             | Corazza et al., 2020  |
| Glove                                      | Hasan et al., 2022  |
| Transformer-based Embedding                | Hasan et al., 2022  |
| FastText                                   | Hajibabae et al., 2022  |
| TF-IDF Vectorizer                          | Maruf et al., 2022  |
| Inverse Document Frequency                 | Maruf et al., 2022  |
| Correlation-based Feature Subset Selection | Bobichev et al., 2017   |
| Information Gain                           | Bobichev et al., 2017   |

Finally, regarding the algorithms applied by these authors, we will present two tables with the set of ML (Table 3.4) and DL (Table 3.5) models performed:

**Table 3.4: List of ML algorithms used**

| <i>Algorithm</i>   | <i>Count</i> |
|--|--------------|
| Random Forest (RF)   | 10           |
| Naïve Bayes (NB)   | 10           |
| Logistic Regression (LR)                                     | 9            |
| Decision Tree (DT)   | 9            |
| AdaBoost Classifier  | 4            |
| K-Nearest Neighbor (KNN)                                     | 4            |
| Bagging  | 2            |
| Multi-view SVM   | 1            |
| K-Means Clustering   | 1            |
| AdaBoost-DT  | 1            |
| Bagging-SVM  | 1            |
| LR+SVM+DT  | 1            |
| Natural Language Processing-SVM (NLP-SVM)                    | 1            |
| Support Vector Classifier (SVC)                              | 1            |
| Linear-SVC   | 1            |
| Gradient Boosting  | 1            |
| XGBoost  | 1            |
| Gaussian Naïve Bayes (GNB)                                   | 1            |
| AdaBoost   | 1            |
| Ensemble tree classifier (ETC)                               | 1            |
| Stochastic gradient descent (SGD)                            | 1            |
| Linear Regression  | 1            |
| Ensemble method (bagging and boosting)                       | 1            |
| Discriminative Multinomial Naïve Bayes classifier (DMNBtext) | 1            |
| Multinomial NB   | 1            |

**Table 3.5: List of DL algorithms used**

| <i>Algorithm</i>   | <i>Count</i> |
|--|--------------|
| Support Vector Machine (SVM)                                   | 15           |
| Long Short-Term Memory (LSTM)                                  | 5            |
| Convolutional Neural Networks (CNN)                            | 5            |
| Multilayer Perceptron (MLP)                                    | 4            |
| Bidirectional LSTM   | 4            |
| LSTM+CNN   | 1            |
| Gated Recurrent Networks (GRU)                                 | 1            |
| GRU+CNN  | 1            |
| Bidirectional GRU  | 1            |
| BiLSTM+CNN+MLP   | 1            |
| MLP+SVM+XGB  | 1            |
| Transformer-CNN+MLP  | 1            |
| Recurrent Neural Network (RNN)                                 | 1            |
| Neural Ensemble  | 1            |
| Bidirectional Encoder Representations from Transformers (BERT) | 1            |
| multilingual BERT (mBERT)                                      | 1            |
| Distil-mBERT   | 1            |
| XLM-RoBERTa  | 1            |
| XLM-RoBERTa-large  | 1            |
| BanglaBERT   | 1            |
| SVM SMO  | 1            |

In summary, within the ML approaches used, RF (count=10), NB (count=10), LR (count=9) and DT (count=9) stand out. Regarding DL models, the most common are SVM (count=15), LSTM (count=5), CNN (count=4), MLP (count=3) and Bidirectional LSTM (count=2). Besides singular algorithms, there are a few hybrid models, for example, the ensemble DL model with LSTM with a layer of CNN.

Considering that the majority of the algorithms are classification models, which means that the goal of the authors is to distinguish two or more types of categories, the metrics used to evaluate the performance were Precision, “the number of true classifications (...) divided by the total number of elements labelled as that class” (Hutto & Gilbert, 2014), Recall, “the number of true classifications divided by the total number of elements that are known to belong to the class” (Hutto & Gilbert, 2014), Accuracy, that represents the overall performance of the model, and F1 Score, “the harmonic mean of precision and recall (...) represents the overall accuracy” (Hutto & Gilbert, 2014).

The following table summarizes the best results obtained by the authors in the reviewed literature:

**Table 3.6: Best results from the literature**

| <i>Reference</i>         | <i>Best model</i> | <i>Accuracy</i> | <i>Precision</i> | <i>Recall</i> | <i>F1 score</i> |
|--------------------------|-------------------|-----------------|------------------|---------------|-----------------|
| Prasetyo & Samudra, 2022 | KNN               | 67.86%          |                  |               |                 |
| Cahyana et al., 2022     | KNN               | 59.68%          |                  |               |                 |



|                                |                       |        |        |        |        |
|--------------------------------|-----------------------|--------|--------|--------|--------|
| Davidson et al., 2017          | LR                    |        | 91%    | 90%    | 90%    |
| Putri et al., 2020             | MLP                   | 83.40% |        |        | 75.90% |
| Turki & Roy, 2022              | RF                    | 95%    |        |        | 95%    |
| Ruwandika & Weeransinghe, 2018 | NB                    | 73.90% | 75%    | 73.90% | 71.90% |
| Marchellim & Ruldeviyani, 2021 | RF                    | 76.70% | 82.70% | 64%    | 73.30% |
| Corazza et al., 2020           | LSTM                  |        |        |        | 82.30% |
| Del Vigna et al., 2017         | SVM                   | 80.60% |        |        |        |
| Mutanga et al., 2022           | LR + SVM+ DT          | 94.21% |        |        |        |
| MacAvaney et al., 2019         | Neural Ensemble (CNN) | 92.17% |        |        | 91.18% |
| Sultan et al., 2023            | CNN                   | 90.20% | 91.60% | 90.40% | 89.90% |
| Al-Hassan & Al-Dossari, 2022   | LSTM + CNN            |        | 72%    | 75%    | 73%    |
| Hasan et al., 2022             | CNN + MLP             | 96.23% | 95.23% |        | 94.19% |
| Amores et al., 2021            | RNN                   | 81%    |        |        |        |
| Defersha & Tune, 2021          | Linear-SVC            |        | 66%    | 66%    | 64%    |
| Hajibabaei et al., 2022        | SVM                   | 95%    | 95%    | 95%    | 95%    |
| Kiilu et al., 2018             | NB                    | 70%    |        |        |        |
| Wadhvani et al., 2023          | ETC                   | 84%    |        |        |        |
| Maruf et al., 2022             | LR                    | 88%    |        |        |        |
| Hasan et al., 2023             | BanglaBERT            | 86%    |        |        | 82%    |
| Balahur, 2013                  | SVM SMO               | 85.07% |        |        |        |
| Bobichev et al., 2017          | NB                    |        |        |        | 84.50% |

Four of the papers reviewed performed only one model – two of them performed KNN, one of them performed Näive Bayes (NB) and the last one performed Support Vector Machine Sequential Minimal Optimization (SVM SMO). KNN algorithm “works by assigning new unclassified examples to the class that contains the majority of its k-nearest neighbours” (Prasetyo & Samudra, 2022). The goal of the article by Prasetyo and Samudra was "to create a system that can detect potential violations of content on Twitter, particularly content containing hate speech” (Prasetyo & Samudra, 2022). The results of

their study showed that, with a  $K$  value equal to 10, the highest accuracy obtained was 67.86% for a KNN method that uses Euclidean as a distance metric. On the other hand, the article by Cahyana et al. (2022) regarding Semi-supervised Text Annotation for Hate Speech Detection had a worse result with an accuracy of 59.68%. In the conclusion of their paper, the authors note that, in future work, there must be an improvement in the vectorization methods, as they only used TF-IDF.

In the study of Kiilu et al. (2018), the authors only performed Naïve Bayes on their data to detect the hate sentiment on tweets. The best result was achieved with Bigram feature with a value of accuracy of 70%. The authors point out the success of Bigram feature because it "incorporates some contextual information which is important for sentiment classification and also generally contain a large number of noisy features and sparse matrix of terms" (Kiilu et al., 2018). Regarding Balahur (2013), the article focuses on Sentiment Analysis of tweets. The work of this author stands out from most of the reviewed articles due to the pre-processing techniques performed. Among the techniques applied are Repeated Punctuation Sign Normalization, which consists of the replacement of "multiple consecutive punctuation signs" by the labels "multistep", "multiexclamation" and "multiquestion" (Balahur, 2013), Emoticon Replacement, where emojis are replaced by the labels "positive", "negative" and "neutral" (the last ones are deleted), and Slang Replacement, where the author "employed the list of slang from a specialized site" (Balahur, 2013). Further on, the author matched tokens with several sentiment lexicons and replaced the word with the correspondent sentiment (Affect Word Matching). Then Modifier Word Matching was employed, where the author used "a list of expressions that negate, intensify or diminish the intensity of the sentiment expressed to detect such words in the tweets" and replaced them with "negator", "diminisher" and "intensifier" (Balahur, 2013). Finally, Balahur (2013) resorted to User and Topic Labelling to replace "@" with "PERSON" and "#" with "TOPIC". In its analysis, Balahur (2013) chose to perform only SVM SMO. This acronym stands for Support Vector Machine Sequential Minimization Optimization (SMO). Sequential Minimization Optimization is "a fast implementation of SVM" that has "better scaling properties than other SVM implementations" (Bobichev et al., 2017). With this method, the author achieved 85.07% of accuracy.

In the reviewed literature, a few articles perform and compare results among ML algorithms. That is the case of Turki & Roy (2022) and Defersha & Tune (2021). Turki and Roy (2022) used Count Vectorizer and Word Cloud Representation for the extraction of the features. The authors performed three models: RF, AdaBoost and Bagging (Bootstrap aggregating), which is an ensemble learning technique that aims to improve the accuracy of ML models. Their best result was achieved with RF with an accuracy of 95%. Defersha and Tune (2021) compared results between Linear-SVC, NB, RF, LR, SVC and DT. The best performance metrics were obtained with Linear-SVC, with precision and recall values of 66% and an F1 score of 64%.

Corazza et al. (2020), MacAvaney et al. (2019) and Hasan et al. (2023) employed merely DL models. Corazza et al. (2020) performed a multilingual evaluation of hate speech detection. The authors used LSTM, Bidirectional LSTM (BiLSTM), Bidirectional Gated Recurrent Unit (BiGRU) and CNN to detect hate sentiment on speech in English,

Italian and German. The English dataset obtained the best result with LSTM (F1 score = 82.30%). MacAvaney et al. (2019)'s study focused on the proposal of a multi-view SVM approach. Among the models used, the best results were obtained with BERT, Neural Ensemble, "which combines the decisions of ten convolutional neural networks with different weight initializations" (MacAvaney et al., 2019), and the proposed classifier. Although the multi-view SVM model had a great performance, the authors conclude that neural networks are "better suited for Twitter data" (MacAvaney et al., 2019). Thus, the best result was achieved with the Neural Ensemble in the HatabaseTwitter dataset with an accuracy of 92.17%. In the article of Hasan et al. (2023), the authors "created a Bangla annotated sentiment analysis dataset based on the context of Ukraine-Russia War" (Hasan et al., 2023). The authors' goal was to analyse how the war was perceived by Bangladeshi people and their sentiments towards the event. This analysis was conducted with the following models: multilingual BERT, Distil-mBERT, XLM-RoBERTa, XLM-RoBERTa-large, BanglaBERT and BiLSTM. The last algorithm was their baseline model. With an accuracy of 86% and an F1 score of 82%, BanglaBERT was the model with the best performance.

The following articles compared results between ML and DL approaches. Regarding the work of Davidson et al. (2017), whose goal was to conduct a multi-class classifier to distinguish between hate speech, offensive speech and neither, they performed the following algorithms: LR, NB, DT, RF and SVM. Their best result was obtained with LR (with L2 regularization). Putri et al. (2020) aimed to compare different classification algorithms to detect hate speech in Indonesian tweets. To achieve their goal, the authors performed NB, DT, MLP, SVM and AdaBoost Classifier. In this paper, Putri et al. (2020) introduced SMOTE, which aims to "balance the data" and "improve classification model" (Putri et al., 2020) and performed the mentioned models with and without the process of SMOTE. The results showed that broadly speaking, the accuracy of the models was better using SMOTE, except for the SVM algorithm. With an accuracy of 83.40%, the best algorithm was MLP.

Regarding the work of Ruwandika & Weeransinghe (2018), the authors executed the following models: SVM, NB, LR, DT and K-Means Clustering. The best-fitted model was the NB model with an accuracy of 73.90%. The NB model labels "probabilities directly with the assumption that the features do not interact with one another" (MacAvaney et al., 2019). Marchellim & Ruldeviyani (2021) performed a Sentiment Analysis of hate speech on Indonesian tweets towards the Indonesian president. The models compared were RF and SVM, and the best results were achieved with RF with an accuracy value of 76.70%. The goal of the article by Maruf et al. (2022) was to identify "emotion, racism, and sentiment analysis" (Maruf et al., 2022). The authors wanted to detect emotions like anger, sadness or love, detect racism and analyse "how people feel about the Ukraine-Russia conflict" (Maruf et al., 2022). The supervised ML and DL models performed were DT, Linear Regression, NB, LR, RF, SVM and an ensemble method. Ensemble methods consist of "a procedure that involves integrating models or classifiers to address a specific issue" (Maruf et al., 2022). On ensemble learning, the authors employed various classifiers and applied bagging and boosting. The best result was obtained with LR, with 88% accuracy.

Bobichev et al. (2017) performed Sentiment Analysis on Ukrainian and Russian news. The authors compared results between ML and DL models such as NB, SVM SMO, Discriminative Multinomial Naive Bayes classifier (DMNBtext) and Multinomial NB. With an F1 score of 84.50%, the model with the best performance was NB. Wadhvani et al. (2023) used tweets to analyse the sentiments of users towards the Ukrainian-Russian War. The authors compared the results within the following algorithms: RF, DT, LR, SVM, XGBoost, Gaussian Naïve base (GNB), AdaBoost, KNN, ETC and SGD. ETC is “an ensemble learning model and operates like RF” and it “fits/trains various randomised DTs to categorise data” (Wadhvani et al., 2023). With an accuracy value of 84%, this was the best-performing model. The work of Sultan et al. (2023) focuses on cyberbullying-related hate sentiment detection and the models performed were DT, RF, NB, LR, KNN, SVM, LSTM, BiLSTM and CNN. Among these, the best-suited approach was a CNN with an accuracy of 90.2%. Al-Hassan & Al-Dossari (2022) focused on detecting hate speech in Arabic tweets and with five distinct categories of hate – General Hate, Religious, Racial, Sexism and None. The models performed were SVM, LSTM, LSTM with a layer of CNN, GRU and GRU and a layer of CNN. The authors concluded that the best model was LSTM with a layer of CNN.

Concerning the paper of Hasan et al. (2022), the authors performed SVM, MLP, RF, CNN, and the ensemble models BiLSTM+CNN+MLP, MLP+SVM+XGBoost, and CNN+MLP. The best result was achieved with CNN+MLP and with the Transformer-based Embedding for the feature extraction process. Regarding the work of Amores et al. (2021), whose goal was to construct an automatic hate speech detector on Twitter in Spanish, the algorithms performed were NB, LR, Stochastic Gradient Descent LR, SVM and RNN. The last one had the best result with 81% of accuracy. The article of Del Vigna et al. (2017), which aimed to detect hate sentiment on Facebook comments on public pages, performed three experiences: the first had three classes of hate (Strong hate, Weak hate and No hate), the second had two classes of hate (Hate and No hate) and the latter had two classes of hate and the annotators were 70% in agreement when classifying the comments. The best-fitting model was SVM (Accuracy=80.60%), which outperformed LSTM.

Mutanga et al. (2022) performed AdaBoost, AdaBoost-DT, Bagging, Bagging-SVM, CNN, DT, LR, LSTM, RF, SVM and Voting Ensemble learning model, “that harnesses the capabilities of LR, DT, and SVM to address overfitting, accommodate new data, and allow for interpretability of a hate speech detection system” (Mutanga et al., 2022). The best-fitting model was the ensemble method (LR+SVM+DT) with an accuracy of 94.21%. Lastly, we present the results of the article by Hajibabae et al. (2022). The authors performed the following algorithms: NB, DT, LR, RF, AdaBoost, SVM, Gradient Boosting and MLP. The best result was achieved with SVM, with TF-IDF as the selected feature extraction method (Accuracy=95%).

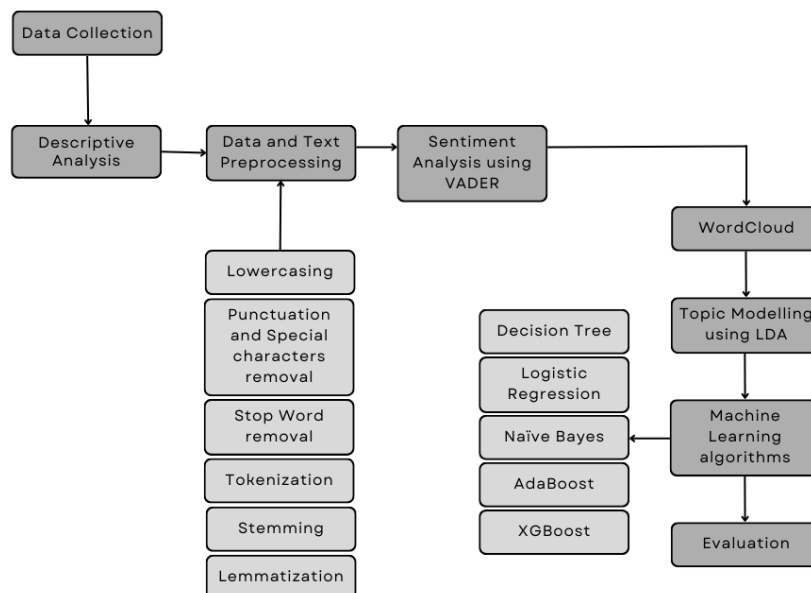


## Chapter 4

### Methodology

In this chapter, we discuss the selection of the dataset and the steps employed for data and text pre-processing. The following flowchart (see Figure 4.1) represents the steps taken during the current phase, which will be clarified later. This visual representation also contains the steps relevant to the next chapter, ‘Results and Discussion’, thereby providing an integrated overview of the sequential steps spanning both phases.

**Figure 4.1: Methodology Flowchart**



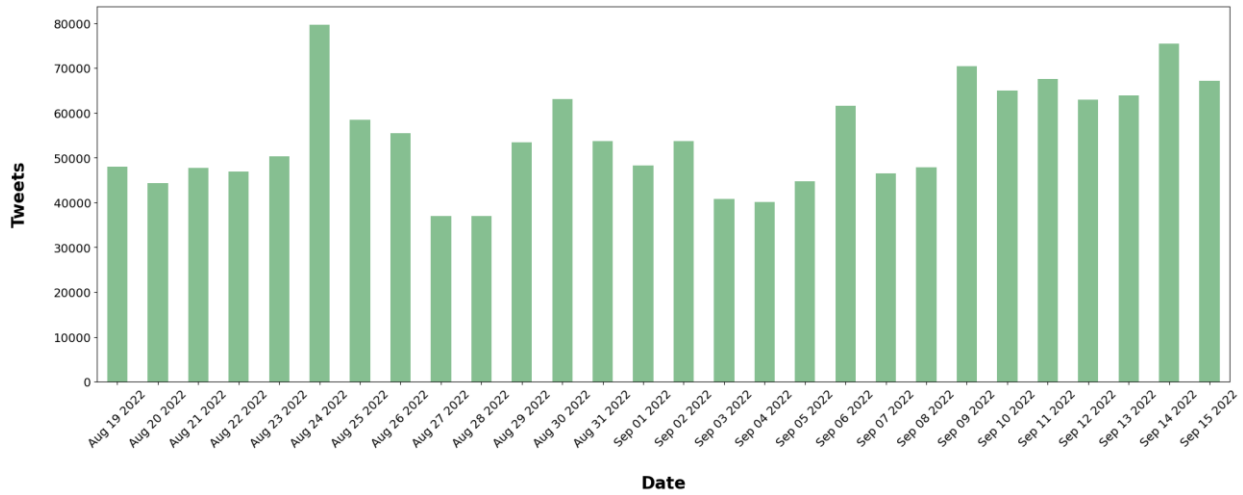
#### 4.1. Descriptive Analysis

To classify tweets between Anti-Russia and Anti-Ukraine sentiments, we resorted to Kaggle and found daily datasets of tweets regarding the ongoing war between Ukraine and Russia. We focused on tweets from August and September of 2022 and February 2023 because these time periods allowed us to analyse sentiment changes over six months. The dataset from 2022 had 1,531,034 records from August 19<sup>th</sup> to September 15<sup>th</sup>, and the dataset from 2023 had 2,364,590 tweets from February 1<sup>st</sup> to February 28<sup>th</sup>.

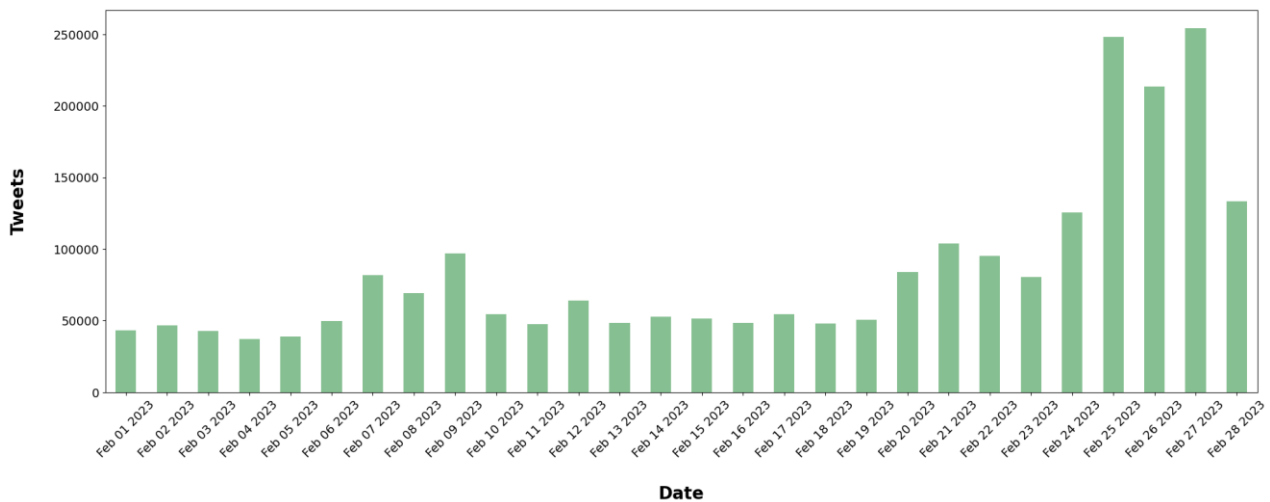
Before pre-processing the datasets, we performed some descriptive analysis. We started with the August and September dataset, whose day with the most tweets was August 24<sup>th</sup> (see Figure 4.2). This day is a significant day for Ukraine as it commemorates 31 years of independence from the Soviet Union and marks the six-month milestone since the Russian invasion. The dataset had tweets in 64 unique languages, with 48.58% in English. In February 2023, the highest Twitter activity coincided with the days after the

first anniversary of the invasion, which was February 24<sup>th</sup> (see Figure 4.3). This dataset also included tweets in 64 unique languages, with 52.62% in English.

**Figure 4.2: Number of tweets per day in August and September 2022**

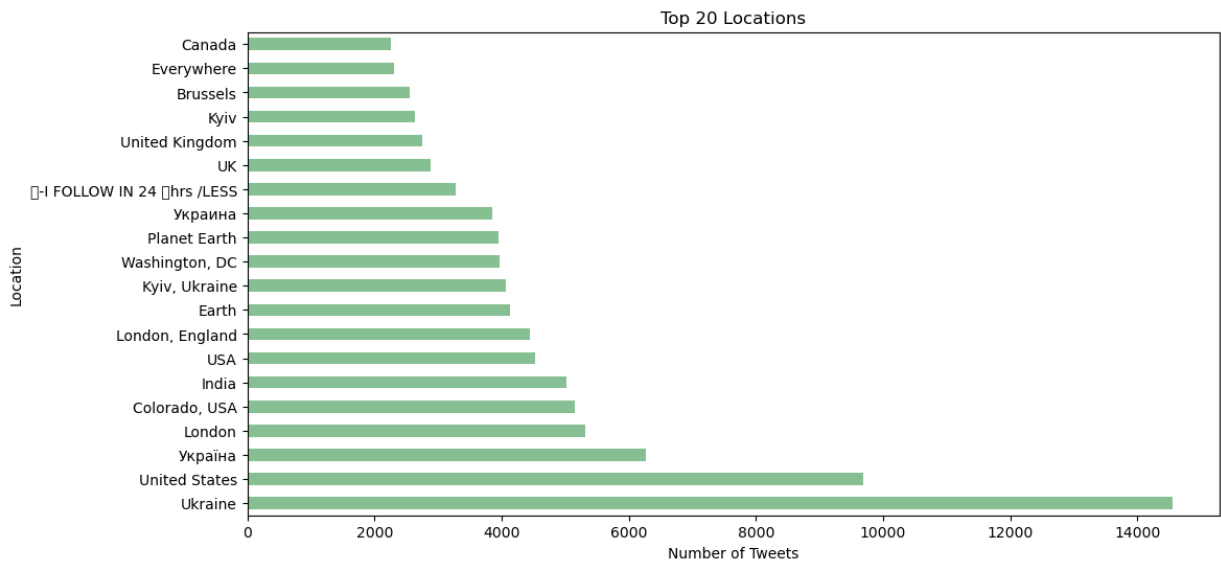


**Figure 4.3: Number of tweets per day in February 2023**

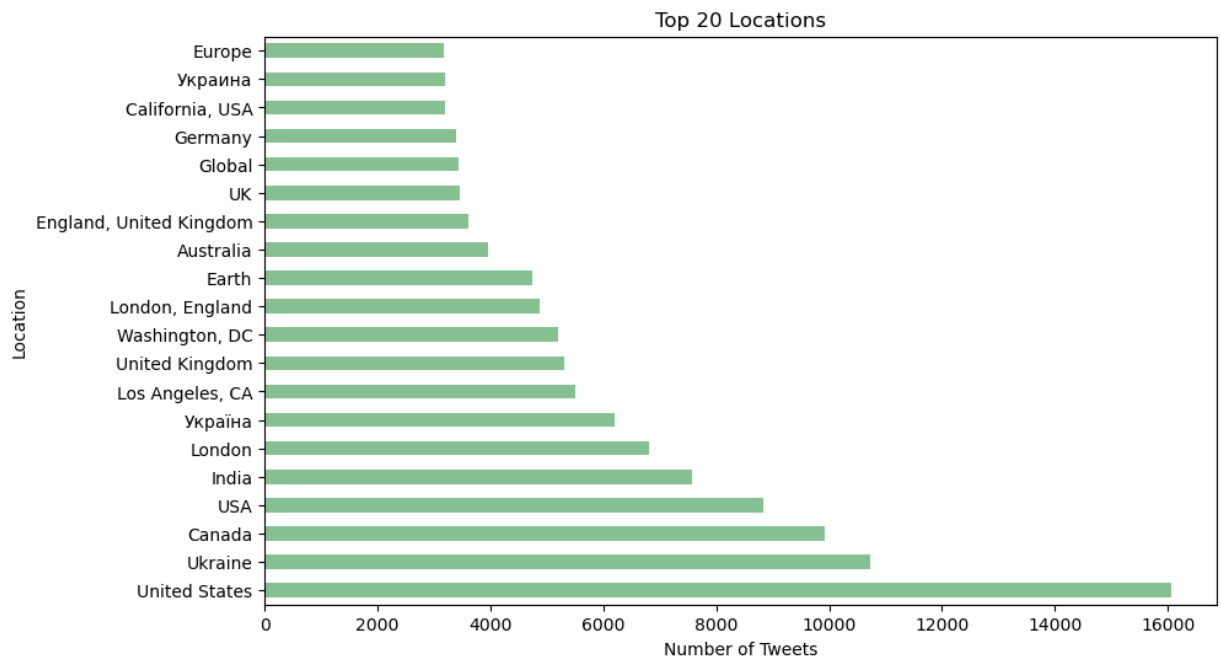


After filtering for tweets in English, we analysed the geographic origins of Twitter users who engaged in discussions regarding the ongoing conflict. Our analysis revealed that, in both datasets, most tweets are from Ukraine and the United States (Figures 4.4 and 4.5). It's interesting to note the geographical shift in Twitter activity between 2022 and 2023, with Ukraine being the primary location for tweets in 2022 and the United States taking the lead in 2023. This shift in location could be related to significant events, such as President Biden's visit to Ukraine in February.

**Figure 4.4 - Number of tweets per location in August and September 2022**



**Figure 4.5 - Number of tweets per location in February 2023:**



The data cleaning process involved removing unnecessary columns and eliminating duplicate records. After cleaning, the dataset from 2022 was left with 743,802 rows and the dataset from 2023 with 1,244,292 rows. The final dataset included the following variables: “userid”, “username”, “following”, “followers”, “tweetid”, “tweetcreatedts”, “retweetcount”, “text”, “hashtags”, “language” and “favorite\_count”. Appendix 3 contains the final list of variables as well as their descriptions.



## 4.2. Data and Text Pre-processing

Before starting our pre-processing, it is important to understand that “preprocessing may affect text mining results” (Hickman et al., 2022). Since our dataset is composed of tweets, i.e., unstructured data, with a lot of noise, such as hashtags or emojis, pre-processing “requires deliberate intervention to make sense of it (...) [as well as] more attention and hence consumes higher processing time” (Nayak et al., 2016). Pre-processing was crucial to clean the data, improve model performance, and make it suitable for input into the models. The pre-processing steps included:

- 1) **Lowercasing:** converting all letters to lowercase to remove differences due to capitalization. This technique “reduces vocabulary size” and “increases statistical power” (Hickman et al., 2022).
- 2) **Punctuation and Special Character removal:** eliminating punctuation, special characters, URLs, HTML, hashtags, mentions, digits, and emojis.
- 3) **Stop Word Removal:** removing common prepositions, conjunctions, pronouns, and articles that do not contribute significantly to the analysis and “make the text look heavier and less important to analysts” and “are not measured as keywords in text mining applications” (Vijayarani et al., 2015). These words frequently “have little discriminative power for the mining process” and “need to be removed” (Aggarwal, 2018).

After these three steps, the tweets were saved into a new column called “cleaned\_text” which was used as input for the next pre-processing steps. Table 4.1 contains an example of the text before and after data preprocessing.

**Table 4.1: Example of a tweet before and after data pre-processing**

|                       |   |
|-----------------------|---|
| <i>Raw tweet</i>      | #Historic visit to Kyiv, Ukraine. " Kyiv has captured a part of my heart " - US President Joe Biden @POTUS<br><br>#JoeBiden in #Kyiv, #Ukraine<br><br>#RussiaIsLosing <a href="https://t.co/blEe158aVV">https://t.co/blEe158aVV</a> |
| <i>“cleaned_text”</i> | historic visit kyiv ukraine kyiv captured part heart us president joe biden potus joebiden kyiv ukraine russiaislosing httpstcoble158avv  |

- 4) **Tokenization:** converting “a character sentence into a sequence of words (or tokens)” (Aggarwal, 2018), by using white spaces as separators. The textual data generated after this step was saved into a new column called “tokens” that can be used as input for further algorithms (see example in Table 4.2).

**Table 4.2: Example of a tweet before and after tokenization**

|                       |  |
|-----------------------|--|
| <i>“cleaned_text”</i> | leftwing academia stand china ukraine stay tuned find                          |
| <i>“tokens”</i>       | ['leftwing', 'academia', 'stand', 'china', 'ukraine', 'stay', 'tuned', 'find'] |

- 5) **Stemming:** extracting the morphological root of a word without changing its semantic interpretation. This method “reduces vocabulary size” and “increases statistical power” (Hickman et al., 2022) by removing suffixes and the number of words “to have accurately matching stems, [and] to save time and memory space” (Vijayarani et al., 2015). Table 4.3 contains an example of a tweet before and after stemming.

**Table 4.3: Example of a tweet before and after stemming**

|                       |   |
|-----------------------|---|
| <i>“cleaned_text”</i> | children treated russia video faint heart typical russian family wonder dont feel sorry soldiers war country devoid human feelings russiaisaterroriststate httpstco8sl3yropvz |
| <i>“stemmed_text”</i> | children treat russia video faint heart typic russian famili wonder dont feel sorri soldier war countri devoid human feel russiaisaterroristst httpstco8sl3yropvz             |

- 6) **Lemmatization:** reducing vocabulary size by converting words to their base form or lemma. Lemmatization differs from Stemming in that “a lemmatizer needs a significant amount of vocabulary and language-specific domain knowledge to carry out its task” (Aggarwal, 2018). Table 4.4 contains an example of a tweet before and after lemmatization.

**Table 4.4: Example of a tweet before and after lemmatization**

|                          |                                |
|--------------------------|--------------------------------|
| <i>“cleaned_text”</i>    | ifenewsagency wants nukes nato |
| <i>“lemmatized_text”</i> | ifenewsagency want nuke nato   |



## Chapter 5

### Results and Discussion

During this chapter, we will discuss the performance of VADER (Hutto & Gilbert, 2014) on Sentiment Analysis and LDA on Topic Modelling (Blei et al., 2003; Jelodar et al., 2019). We also resort to the WordCloud library on Python (<https://pypi.org/project/wordcloud/>) to analyse word frequency and relevance. Further on, we will go over the ML models used for the classification problem.

#### 5.1. Sentiment Analysis using VADER

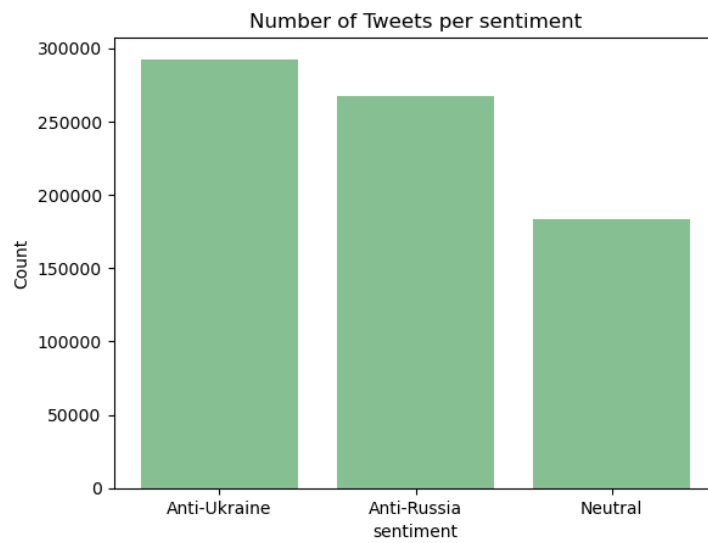
Since the data from Kaggle was collected directly from Twitter, labelling the collected tweets was a crucial step. To achieve this, we employed VADER, a widely used Python library for Sentiment Analysis. VADER uses a lexicon-based method and “performs exceptionally well in the social media domain” (Hutto & Gilbert, 2014).

VADER brings more sensitivity to the social media context and expressions and, as Hutto and Gilbert (2014) highlighted in their article, this sentiment lexicon “is bigger, yet just as simply inspected, understood, quickly applied (without a need for extensive learning/training)”. Furthermore, the lexicon is of high quality and has been validated by human evaluators, making it a valuable resource for Sentiment Analysis tasks.

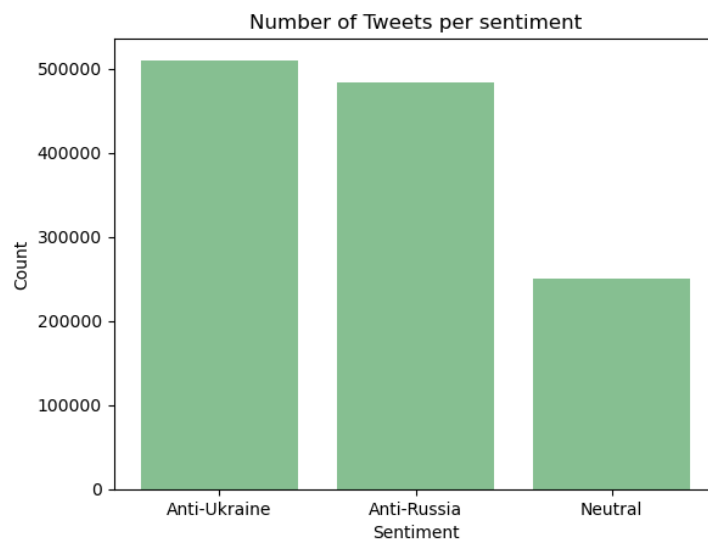
To apply VADER to our dataset, we performed Sentiment Analysis on the “cleaned\_text” column, which had undergone pre-processing steps to ensure data cleanliness and model compatibility. The Sentiment Analysis generated compound scores, ranging from -1 (indicating the most extreme negative sentiment) to 1 (representing the most extreme positive sentiment). Based on these compound scores, we created a new column named “sentiment” to categorize the tweets into ‘Positive’, ‘Negative’, and ‘Neutral’. In order to adapt the sentiment labels to our specific research focus of classifying Anti-Ukraine and Anti-Russia sentiments, we further examined examples of tweets falling into the 'Positive' and 'Negative' sentiment categories. After careful examination, we found that tweets labelled as 'Positive' (compound score  $\geq 0,05$ ) exhibited an 'Anti-Russia' sentiment, while tweets labelled as 'Negative' (compound score  $\leq -0,05$ ) conveyed an 'Anti-Ukraine' sentiment. This approach provided valuable insights into the prevalence of Anti-Ukraine and Anti-Russia sentiments within the collected Twitter data.

Regarding the dataset from 2022 (see Figure 5.1), the ‘Anti-Ukraine’ sentiment represents 39.33% of the total (292,516 records), the ‘Anti-Russia’ sentiment represents 35.96% of the tweets (267,478 records) and the ‘Neutral’ sentiment represents 24.71% (183,808 records). Concerning the 2023 dataset (see Figure 5.2), there are 510,100 tweets with ‘Anti-Ukraine’ sentiment (41%), 483,899 tweets with ‘Anti-Russia’ sentiment (38.89%) and 250,293 tweets with ‘Neutral’ sentiment (20.12%). On both datasets, ‘Anti-Ukraine’ is the predominant sentiment, although the difference between the two main sentiments is not very significant.

**Figure 5.1: Number of tweets per sentiment (2022 dataset)**



**Figure 5.2: Number of tweets per sentiment (2023 dataset)**



## 5.2. WordCloud

We employed the WordCloud library for a visual representation of the text data. This method helped us analyse the most frequent words on tweets from both datasets. On the dataset from 2022, the most frequent words were “amp”, "ukraine", "russia", "war", "independence day", “armed force”, “one”, “support ukraine”, “ukrainerussawar” and “slavaukraini” (see Appendix 1). On the dataset from 2023, besides "ukraine", "russia" and “war”, the most frequent words were "canada", "america", "bakhmut", “germany”, “invasion ukraine”, “bidens visit” and “one year” (see Appendix 2).

Furthermore, we analysed the most frequent words for each sentiment – ‘Anti-Ukraine’ and ‘Anti-Russia’. These words are present in Table 5.1 and in the Appendix 4, 5, 6 and 7.

**Table 5.1: Most frequent words for each sentiment on both datasets**

|                     | <i>Anti-Ukraine</i>   | <i>Anti-Russia</i>  |
|---------------------|---|---|
| <i>2022 dataset</i> | “russia”, “amp”, “war ukraine”, “putin”, “biden”, “armed force”, “forces ukraine”, “russiaisaterroriststate”, “country”, “ukrainerussiawar” | “ukraine”, “amp”, “thank”, “one”, “independence day”, “slavaukraini”, “russia”, “support ukraine”, “biden”, “people”, “great”, “today”, “putin” |
| <i>2023 dataset</i> | “ukraine”, “canada”, “germany”, “bakhmut”, “war”, “america”, “russia”, “zelenskys double”, “one year”, “kyiv”, “invasion ukraine”           | “support”, “ukraine”, “dear unhrc”, “online class”, “usa”, “canada”, “justice4tigray”, “one year”, “russia”, “ungeneva intlcrimcourt”           |

The 2022 and 2023 datasets analysis reveals distinctive patterns in sentiment expressions related to the Russo-Ukrainian conflict. In the context of 'Anti-Ukraine' sentiment within the 2022 dataset, numerous references to "war" and armed forces were observed. In contrast, the 2023 dataset frequently mentioned countries that opposed Russia, notably "canada," "germany," and "america". Moreover, the phrase "one year" emerged prominently in the 2023 dataset. Turning to the 'Anti-Russia' sentiment in the 2022 dataset, references to Ukraine's adversaries, specifically “putin" and "russia," were prevalent, alongside allusions to "independence day". Primarily, the 2022 dataset also exhibited a substantial number of supportive expressions, including phrases like "support ukraine" and "slavaukraini", signifying "Glory to Ukraine!". In the 2023 dataset, the term "one year" remained prominent in conjunction with continued indications of support. The variations in frequently mentioned keywords and themes between the two datasets underscore the evolving nature of public sentiment surrounding the Russo-Ukrainian conflict over time.

We analysed the most frequent words associated with each sentiment and further investigated their co-occurrences (see Tables 5.2, 5.3, 5.4 and 5.5).

**Table 5.2: Most frequent words and their co-occurrences for the ‘Anti-Ukraine’ sentiment (2022 dataset)**

| <i>Most frequent words</i> | <i>Co-occurrences</i>  |
|----------------------------|--|
| <i>“russia”</i>            | “ukraine”, “war”, “russian”, “putin”, “ukrainian”, “amp”, “us”, “russiaisaterroriststate”, “nato”, “china” |
| <i>“amp”</i>               | “ukraine”, “russia”, “war”, “russian”, “putin”, “biden”, “us”, “people”, “ukrainian”, “nato”               |
| <i>“war”</i>               | “ukraine”, “russia”, “russian”, “putin”, “amp”, “ukrainian”, “us”, “nato”, “people”, “ukrainewar”          |
| <i>“ukraine”</i>           | “russia”, “war”, “russian”, “ukrainian”, “putin”, “ukrainerussiawar”, “amp”, “ukrainewar”, “us”, “nato”    |

|                                 |   |
|---------------------------------|---|
| <b>“putin”</b>                  | “ukraine”, “russia”, “war”, “russian”, “amp”, “world”, “fuck”, “news”, “ukrainerussiawar”, “china”                  |
| <b>“biden”</b>                  | “trump”, “ukraine”, “amp”, “russia”, “usa”, “us”, “news”, “maga”, “democrats”, “china”                              |
| <b>“armed”</b>                  | “ukraine”, “forces”, “russian”, “ukrainian”, “russia”, “war”, “kherson”, “region”, “ukrainerussiawar”, “ukrainewar” |
| <b>“force”</b>                  | “ukraine”, “russia”, “russian”, “air”, “war”, “world”, “ukrainian”, “military”, “amp”, “putin”                      |
| <b>“russiainterroriststate”</b> | “russia”, “ukraine”, “russian”, “war”, “people”, “terrorist”, “russians”, “state”, “standwithukraine”, “ukrainian”  |
| <b>“country”</b>                | “ukraine”, “russia”, “war”, “russiainterroriststate”, “russian”, “people”, “biden”, “every”, “amp”, “must”          |
| <b>“ukrainerussiawar”</b>       | “ukraine”, “russian”, “russia”, “ukrainewar”, “war”, “ukrainian”, “putin”, “russiaukrainewar”, “kherson”, “nato”    |

**Table 5.3: Most frequent words and their co-occurrences for the ‘Anti-Ukraine’ sentiment (2023 dataset)**

| <i>Most frequent words</i> | <i>Co-occurrences</i>   |
|----------------------------|---|
| <b>“ukraine”</b>           | “war”, “russia”, “russian”, “bakhmut”, “ukrainian”, “germany”, “putin”, “amp”, “kyiv”, “usa”                |
| <b>“canada”</b>            | “ukraine”, “bakhmut”, “germany”, “america”, “russian”, “ukrainian”, “war”, “footage”, “near”, “attack”      |
| <b>“germany”</b>           | “ukraine”, “bakhmut”, “russian”, “canada”, “america”, “kyiv”, “ukrainian”, “war”, “usa”, “footage”          |
| <b>“bakhmut”</b>           | “ukraine”, “germany”, “russian”, “canada”, “america”, “ukrainian”, “war”, “russia”, “wagner”, “forces”      |
| <b>“war”</b>               | “ukraine”, “russia”, “russian”, “amp”, “putin”, “germany”, “year”, “ukrainian”, “us”, “bakhmut”             |
| <b>“america”</b>           | “ukraine”, “bakhmut”, “germany”, “canada”, “russian”, “ukrainian”, “war”, “footage”, “russia”, “attack”     |
| <b>“russia”</b>            | “ukraine”, “war”, “russian”, “putin”, “us”, “amp”, “nato”, “china”, “ukrainian”, “usa”                      |
| <b>“zelenskys”</b>         | “media”, “kyiv”, “visit”, “double”, “polish”, “covered”, “recorded”, “awkward”, “bidens”, “accidentally”    |
| <b>“double”</b>            | “media”, “polish”, “visit”, “kyiv”, “zelenskys”, “recorded”, “covered”, “awkward”, “bidens”, “accidentally” |
| <b>“one”</b>               | “ukraine”, “war”, “year”, “russia”, “russian”, “amp”, “invasion”, “us”, “people”, “russias”                 |
| <b>“year”</b>              | “ukraine”, “war”, “one”, “russia”, “russian”, “ago”, “invasion”, “russias”, “amp”, “last”                   |
| <b>“kyiv”</b>              | “ukraine”, “germany”, “usa”, “russian”, “war”, “canada”, “ukrainian”, “visit”, “bidens”, “media”            |
| <b>“invasion”</b>          | “ukraine”, “russias”, “russian”, “year”, “war”, “russia”, “one”, “feb”, “fullscale”, “standwithukraine”     |

Regarding the ‘Anti-Ukraine’ sentiment, on the 2022 dataset (Table 5.2), the term “amp”, which stands for Ukraine War Amps, a multinational organization that aids the victims of the war since its outbreak, frequently appears alongside words such as "war," "people," "russia," "biden," "putin," "nato," "us", and "ukraine". The terms "ukraine" and "russia" often co-occur in the same context, and the words "usa" and “us” are consistently mentioned in conjunction with both countries. Additionally, references to "russia" are often connected to "china" in the discourse. In 2022, the tweets collected had a lot of mentions of Ukraine's military, such as "force" and "war” which occur with “amp”, “armed” and “military”. Interestingly, the hashtag "russiainterroriststate" is also found in tweets classified as 'Anti-Ukraine,' along with words like "war," "terrorist," and "standwithukraine".

Contrary to what happens in 2022, on the 2023 dataset (Table 5.3), there are more countries as the most common words; besides “ukraine” and “russia”, we have "canada", “germany” and “america”. These countries are often referenced together, possibly due to their joint military support to Ukraine, including the supply of battle tanks, heavy weapons, and troops "to assist Ukrainian soldiers and train them with the skills needed to operate the armoured vehicles” (Bowden, 2023). In the 2023 dataset, the word "bakhmut" stands out as one of the most frequently mentioned terms, appearing alongside "ukraine", "germany", "war", "canada", "america", and "wagner". This association arises from Wagner's capture of part of the Bakhmut territory in February 2023. Furthermore, February 2023 also marks one year since the invasion started and that’s why “war”, “year”, “one” and “invasion” are often mentioned together, generating discussions around the conflict’s ongoing impact. Lastly, the words “zelenskys” and “kyiv” are linked to “media”, “kyiv”, “bidens” and “visit”. Notably, the visit of the US President, Joe Biden, to Kyiv "days before the first anniversary of Russia's full-scale invasion" (Beaumont et al., 2023) has drawn significant attention.

**Table 5.4: Most frequent words and their co-occurrences for the ‘Anti-Russia’ sentiment (2022 dataset)**

| <i>Most frequent words</i> | <i>Co-occurrences</i>  |
|----------------------------|--|
| <b>“ukraine”</b>           | “russia”, “russian”, “ukrainian”, “amp”, “support”, “putin”, “standwithukraine”, “day”, “people”, “us”               |
| <b>“amp”</b>               | “ukraine”, “russia”, “us”, “biden”, “russian”, “putin”, “support”, “people”, “like”, “china”                         |
| <b>“thank”</b>             | “ukraine”, “standwithukraine”, “support”, “slavaukraini”, “much”, “amp”, “independence”, “day”, “us”, “ukrainian”    |
| <b>“one”</b>               | “ukraine”, “russia”, “russian”, “amp”, “like”, “biden”, “ukrainian”, “day”, “putin”, “us”                            |
| <b>“independence”</b>      | “ukraine”, “day”, “happy”, “standwithukraine”, “ukrainian”, “today”, “slavaukraini”, “freedom”, “ukraines”, “people” |



|                       |   |
|-----------------------|---|
| <b>“day”</b>          | “ukraine”, “independence”, “happy”, “standwithukraine”, “ukrainian”, “today”, “slavaukraini”, “russia”, “people”, “amp” |
| <b>“slavaukraini”</b> | “ukraine”, “day”, “ukrainian”, “standwithukraine”, “good”, “happy”, “independence”, “people”, “amp”, “russia”           |
| <b>“russia”</b>       | “ukraine”, “putin”, “russian”, “china”, “amp”, “us”, “world”, “peace”, “usa”, “ukrainian”                               |
| <b>“support”</b>      | “ukraine”, “russia”, “amp”, “people”, “standwithukraine”, “ukrainian”, “russian”, “us”, “today”, “eu”                   |
| <b>“biden”</b>        | “trump”, “amp”, “ukraine”, “usa”, “us”, “president”, “news”, “like”, “russia”, “maga”                                   |
| <b>“people”</b>       | “ukraine”, “russia”, “ukrainian”, “amp”, “support”, “russian”, “world”, “standwithukraine”, “day”, “us”                 |
| <b>“great”</b>        | “ukraine”, “russia”, “amp”, “biden”, “russian”, “putin”, “us”, “slavaukraini”, “standwithukraine”, “people”             |
| <b>“today”</b>        | “ukraine”, “day”, “russia”, “independence”, “ukrainian”, “support”, “people”, “standwithukraine”, “amp”, “russian”      |
| <b>“putin”</b>        | “russia”, “ukraine”, “russian”, “world”, “trump”, “amp”, “china”, “like”, “news”, “biden”                               |

**Table 5.5: Most frequent words and their co-occurrences for the ‘Anti-Russia’ sentiment (2023 dataset)**

| <i>Most frequent words</i> | <i>Co-occurrences</i>  |
|----------------------------|--|
| <b>“support”</b>           | “ukraine”, “standwithukraine”, “amp”, “russia”, “people”, “ukrainian”, “us”, “world”, “war”, “year”                        |
| <b>“ukraine”</b>           | “russia”, “russian”, “amp”, “support”, “standwithukraine”, “ukrainian”, “people”, “peace”, “us”, “like”                    |
| <b>“dear”</b>              | “unhrc”, “tigraygenocide”, “justice4tigray”, “un”, “amp”, “ungeneva”, “justice”, “equal”, “laetitiabader”, “intlcrimcourt” |
| <b>“unhrc”</b>             | “tigraygenocide”, “justice4tigray”, “dear”, “un”, “amp”, “justice”, “ungeneva”, “laetitiabader”, “intlcrimcourt”, “equal”  |
| <b>“online”</b>            | “homework”, “help”, “assignment”, “class”, “exam”, “usa”, “essay”, “canada”, “australia”, “get”                            |
| <b>“class”</b>             | “homework”, “help”, “online”, “assignment”, “exam”, “essay”, “usa”, “canada”, “australia”, “get”                           |
| <b>“usa”</b>               | “canada”, “help”, “homework”, “get”, “uk”, “online”, “ukraine”, “australia”, “assignment”, “class”                         |
| <b>“canada”</b>            | “usa”, “help”, “homework”, “online”, “australia”, “assignment”, “class”, “uk”, “exam”, “essay”                             |
| <b>“justice4tigray”</b>    | “tigraygenocide”, “unhrc”, “amp”, “dear”, “un”, “justice”, “ungeneva”, “laetitiabader”, “intlcrimcourt”, “equal”           |

|                        |  |
|------------------------|--|
| <i>“one”</i>           | “ukraine”, “year”, “russia”, “standwithukraine”, “amp”, “love”, “people”, “two”, “things”, “ukrainian”                 |
| <i>“year”</i>          | “ukraine”, “one”, “russia”, “invasion”, “since”, “standwithukraine”, “support”, “amp”, “people”, “ukrainian”           |
| <i>“russia”</i>        | “ukraine”, “russian”, “putin”, “amp”, “china”, “peace”, “today”, “us”, “people”, “nato”                                |
| <i>“ungeneva”</i>      | “unhrc”, “tigraygenocide”, “dear”, “justice4tigray”, “un”, “equal”, “laetitiabader”, “intlcrimcourt”, “amp”, “justice” |
| <i>“intlcrimcourt”</i> | “unhrc”, “tigraygenocide”, “equal”, “justice4tigray”, “un”, “amp”, “dear”, “ungeneva”, “victims”, “laetitiabader”      |

In the context of 'Anti-Russia' sentiment, the 2022 dataset (Table 5.4) reveals significant co-occurrence of certain keywords. Specifically, mentions of "russia", "putin", "china", "usa", and "amp" are often found together in tweets. The word "trump" is frequently mentioned when users comment about "biden" or "putin". Furthermore, a substantial portion of this dataset's most common words reflects support demonstrations, particularly considering Ukraine's Independence Day commemoration. Phrases and words like "standwithukraine", "slavaukraini", "independence", and "day" are frequently mentioned alongside words like "happy", "support", "people", and "freedom". These associations highlight the strong display of solidarity and backing for Ukraine during this period.

On the 2023 dataset (Table 5.5), for the 'Anti-Russia' sentiment, there are also a lot of supportive expressions towards Ukraine. The terms “help”, “canada” and “usa” are frequently mentioned together, likely due to the help Ukraine received from those countries in February. Furthermore, various words in the dataset allude to the marking of the first year since the war started. Words like "support", "year" and "ukraine" co-occur with “standwithukraine”, “people”, “peace”, “one”, “invasion” and “russia”. In addition to the Ukraine-related discussions, the 2023 dataset also contains numerous references to the genocide that occurred in Tigray, Ethiopia.

Also, given the dates with the most tweets from both datasets, we resorted to WordCloud library and analysed the most frequent words present in Table 5.6 and in the Appendix 8, 9, 10, 11, 12, 13, 14 and 15.

**Table 5.6: Days with more tweets and the most frequent words**

| <i>Date</i>              | <i>Most frequent words</i>   |
|--------------------------|--|
| <i>August 24th, 2022</i> | “independence day”, “happy independence”, “ukraine”, “russia”, “biden”, “russiaisaterroriststate”, “day ukraine”, “putin”, “today”, “amp”, “six month” |
| <i>August 30th, 2022</i> | “ukraine”, “russia”, “amp”, “russian”, “slavaukraini”, “biden”, “putin”, “china”, “war”, “kherson”   |

|                            |  |
|----------------------------|--|
| <i>September 6th, 2022</i> | “ukraine”, “russia”, “real madridlive link”, “amp”, “russian”, “putin”, “sevilla”, “ukrainian”, “biden”, “vs juventus”, “celtic v”   |
| <i>September 9th, 2022</i> | “queen elizabeth”, “russia”, “ukraine”, “rest peace”, “russian”, “putin”, “people”, “elizabeth ii”, “slavaukraini”, “amp”  |
| <i>February 7th, 2023</i>  | “ukraine”, “nato”, “usa”, “russian”, “amp”, “russia”, “ukrainian” “america”, “canada”, “germany”, “bakhmut”  |
| <i>February 21st, 2023</i> | “putin”, “russia”, “ukraine”, “germany”, “china”, “amp”, “jimin coming”, “war ukraine”, “russian”, “biden”   |
| <i>February 25th, 2023</i> | “invasion ukraine”, “support ukraine”, “stands ukraine”, “one year”, “ukraine”, “ukrainewillwin”, “war ukraine”, “russia isolated”, “russian propagandist”, “everyone stand” |
| <i>February 27th, 2023</i> | “dear unhr”, “intlcrimcourt”, “ungeneva”, “address justice4tigray”, “assure justice4tigray”, “victims tigraygenocide”, “seek justice4tigray”, “n govt”, “unhr un”            |

On August 24th, 2022, the most common words were related to Ukraine’s Independence Day and six months since the war began. On this day, people also tweeted about “russiaisaterroriststate” probably due to the “rocket attack on a Ukrainian train station (...) killing 22 people” (Gatopoulos & Varenysia, 2022). Besides, “biden” was also mentioned because the president of the US “announced nearly \$3 billion US military aid” (Underwood et al., 2022). On August 30th, 2022, the term “kherson” occurred multiple times because Ukraine decided to launch “a counteroffensive in southern areas” that were occupied by the Russian military, namely Kherson (Underwood et al., 2022). September 6th, 2022, and September 9th, 2022, were marked by some football content and the passing of Elizabeth II, former Queen of the United Kingdom and the Commonwealth.

The words mentioned on February 7th, 2023, are related to Zelenskyy's visit to some European countries and the extended military support given by Canada and Germany to Ukraine. On February 21st, 2023, one of the most common terms was “china”. This is due to the attempt of the US and its allies “to press China not to provide weapons to Russia” (Underwood et al., 2022). On the same day, the words “germany” and “biden” appear together due to the visit of the German Chancellor Olaf Scholz to the US. On February 25th, 2023, the most common words were related to the first year since the invasion of Ukraine. On February 27th, 2023, the words are related to the Tigray genocide in Ethiopia.

### 5.3. Topic Modelling

Topic Modelling methods are extensively utilized in NLP, particularly in the analysis of social media data, due to their ability to understand “the reactions and conversations between people in online communities, as well as extracting useful patterns and

understandable from their interactions in addition to what they share on social media websites (...)" (Jelodar et al., 2019).

Among the various Topic Modelling methods in ML, the LDA model stands out as one of the most widely used and popular approaches. LDA is a “generative probabilistic model of a corpus” (Jelodar et al., 2019) used for classification, summarisation, novelty, and similarity detection. The underlying idea of this model is "that documents are represented as random mixtures over latent topics” (Blei et al., 2003), and each topic is a combination of words. In summary, in LDA, words are represented as word probabilities and a distribution of words represents each topic. Within each topic, the word with the highest probability "usually give[s] a good idea of what the topic is (...)" (Jelodar et al., 2019).

Before implementing the LDA model, we first generated representative samples from both datasets. This step was essential to ensure the efficiency of model training and testing for both the LDA model itself and any subsequent algorithms utilized in the study. The dataset of 2022, which initially consisted of 743,802 records, was reduced to 30%, resulting in 223,141 tweets. Similarly, for the dataset of 2023, we reduced the number of records to 20% of the original records (1,244,292), resulting in 248,858 tweets. In turn, we created a sample of 2,000 records for each dataset and built the LDA model with 10 different topics.

For both datasets, we got a reasonably good coherence score. This score measures the quality of a topic by determining how closely related the significant words within the topic are in terms of meaning. This measure helps distinguish between topics that have meaningful semantic interpretations and topics that are purely for statistical analysis. Coherence scores greater than 0.4 signify a positive result, indicating strong semantic similarity among the high-scoring words within the topic. On the other hand, a coherence score below 0.2 indicates a weak correlation between the high-scoring words and, therefore, a low semantic similarity among them. The coherence score for the 2022 dataset was measured at 0.4965, whereas for the 2023 dataset, it generated a score of 0.4436.

Table 5.7 and Table 5.8 contain the highest-scoring words for each topic on both datasets:

**Table 5.7: Highest-scoring words for each topic on the 2022 dataset**

| <i>Topics</i> | <i>Words (weight)</i>   |
|---------------|---|
| <i>0</i>      | "today" (0.002); "take" (0.002); "biden" (0.002); "man" (0.002); "trump" (0.002); "week" (0.002); "enough" (0.002); "shut" (0.002); "shutdown" (0.002); "blackpink" (0.002) |
| <i>1</i>      | "ukraine" (0.004); "slavaukraini" (0.004); "happy" (0.002); "come" (0.002); "find" (0.002); "thank" (0.002); "get" (0.002); "many" (0.002); "guess" (0.002); "take" (0.002) |
| <i>2</i>      | "day" (0.004); "come" (0.003); "russian" (0.003); "good" (0.003); "nazi" (0.002); "pink" (0.002); "era" (0.002); "bear" (0.002); "shutdown" (0.002); "blackpink" (0.002)    |

|          |  |
|----------|--|
| <b>3</b> | "support" (0.003); "also" (0.002); "ukraine" (0.002); "news" (0.002); "biden" (0.002); "die" (0.002); "mriyareport" (0.002); "true" (0.002); "much" (0.002); "major" (0.001)       |
| <b>4</b> | "help" (0.003); "s" (0.002); "ukraine" (0.002); "new" (0.002); "war" (0.002); "save" (0.002); "military" (0.002); "russian" (0.002); "people" (0.002); "believe" (0.002)           |
| <b>5</b> | "war" (0.004); "say" (0.003); "people" (0.003); "world" (0.002); "amp" (0.002); "think" (0.002); "see" (0.002); "save" (0.002); "go" (0.002); "international" (0.002)              |
| <b>6</b> | "video" (0.002); "dog" (0.002); "way" (0.002); "s" (0.002); "keep" (0.002); "go" (0.002); "refute" (0.002); "fascist" (0.002); "call" (0.002); "need" (0.002)                      |
| <b>7</b> | "lose" (0.002); "military" (0.002); "trump" (0.002); "also" (0.002); "dollar" (0.002); "look" (0.002); "war" (0.002); "seem" (0.001); "say" (0.001); "interest" (0.001)            |
| <b>8</b> | "russian" (0.010); "ukrainian" (0.009); "force" (0.006); "war" (0.005); "destroy" (0.003); "military" (0.003); "today" (0.002); "attack" (0.002); "use" (0.002); "missile" (0.002) |
| <b>9</b> | "ukraine" (0.004); "war" (0.004); "stop" (0.003); "side" (0.003); "country" (0.002); "troop" (0.002); "money" (0.002); "pay" (0.002); "go" (0.002); "kill" (0.002)                 |

On the 2022 dataset, on the first topic, we can conclude that there are words that relate to time, "today", which weighs 0.002 and is considered the most important word within this topic, and "week". This topic is most likely connected to the US because we have mentions of "biden" and "trump". Topic 1 is related to Ukraine's Independence Day since the highest-scoring words are "ukraine", "slavaukraini" and "happy". Topic 2 is also related to Independence Day since we have mentions of "day" and "good"; however, there is also the word "nazi", which expresses an 'Anti-Ukraine' sentiment, as Putin usually uses this term to insult Ukrainians and justify his decision to invade the country. Topic 3 is most likely about the support Ukraine has been receiving, but also about "news" and "mriyareport", a network created to help the Ukrainian victims of the war and that also counts on the help of volunteers to spread the news about the war. Similarly to topic 3, topic 4 is about support for Ukraine, having the word "help" with the highest score. Topic 6 is probably related to insults from Russia to Ukraine – "dog", "fascist" and "call". Topic 8 has "russia" and "ukrainian" as the highest-scoring words and is connected to the violence of the war – "force", "war", "destroy", "military", "attack" and "missile". Finally, on the last topic, there are several action verbs, such as "stop", "pay" and "kill".

**Table 5.8: Highest-scoring words for each topic on the 2023 dataset**

| <i>Topics</i> | <i>Words (weight)</i>  |
|---------------|--|
| <b>0</b>      | "bakhmut" (0.006); "russian" (0.005); "soldier" (0.004); "ukraine" (0.004); "ukrainian" (0.003); "loss" (0.003); "see" (0.003); "slavaukraini" (0.003); "force" (0.002); "today" (0.002) |

|          |  |
|----------|--|
| <b>1</b> | "ukraine" (0.005); "m" (0.002); "httpstco" (0.002); "call" (0.002); "use" (0.002); "russian" (0.002); "domain" (0.001); "make" (0.001); "want" (0.001); "corruption" (0.001)       |
| <b>2</b> | "russian" (0.002); "war" (0.002); "pass" (0.002); "part" (0.002); "deal" (0.002); "missile" (0.002); "time" (0.002); "ukrainian" (0.002); "drone" (0.002); "leader" (0.002)        |
| <b>3</b> | "come" (0.006); "jimin" (0.005); "face" (0.004); "world" (0.002); "year" (0.002); "wake" (0.002); "show" (0.002); "ukraine" (0.002); "people" (0.002); "amp" (0.002)               |
| <b>4</b> | "think" (0.003); "work" (0.002); "ukraine" (0.002); "rescue" (0.002); "kill" (0.002); "believe" (0.002); "power" (0.002); "cause" (0.001); "booktrailer" (0.001); "search" (0.001) |
| <b>5</b> | "body" (0.005); "many" (0.003); "wound" (0.003); "brain" (0.003); "first" (0.002); "vuhledar" (0.002); "country" (0.002); "miss" (0.002); "work" (0.002); "try" (0.002)            |
| <b>6</b> | "war" (0.005); "pay" (0.004); "assignment" (0.004); "peace" (0.004); "homework" (0.004); "help" (0.003); "ukrainewar" (0.003); "get" (0.003); "online" (0.003); "say" (0.002)      |
| <b>7</b> | "russian" (0.003); "sanction" (0.003); "redbubble" (0.003); "get" (0.002); "gift" (0.002); "artist" (0.002); "top" (0.002); "help" (0.002); "sight" (0.002); "joke" (0.001)        |
| <b>8</b> | "russian" (0.007); "war" (0.006); "say" (0.006); "go" (0.003); "military" (0.003); "give" (0.003); "destroy" (0.003); "let" (0.003); "laetitiabader" (0.002); "time" (0.002)       |
| <b>9</b> | "stop" (0.003); "trump" (0.003); "stand" (0.002); "provide" (0.002); "russian" (0.002); "price" (0.002); "support" (0.002); "war" (0.002); "murder" (0.002); "mass" (0.002)        |

The first topic on the 2023 dataset (Table 5.8) has "bakhmut" as the highest-scoring word. Alongside this term there are also "soldier", "loss" and "force", which means that this topic is most likely related to the battle of Bakhmut. Topic 2 is connected to the weapons used in the war – "missile" and "drone". Similarly to what happens on the 2022 dataset, a consistent pattern was observed where the word "russia" had the highest coherence score, and the subsequent words consistently pertained to themes related to war or attacks. This pattern of association indicates a strong semantic relationship between "russia" and the mentioned topics of war and attacks within the dataset. Topic 3 is linked to the marking of one year since the conflict started. Topic 4 is likely connected to Ukrainian victims of the war – "work", "rescue", "kill", "believe" and "power". Topic 5 is related to injuries, having "body", "many", "wound" and "brain" as the highest-scoring words. In this topic, there is also mention of the city of "vuhledar". Topic 6 has "war" as the highest-scoring word and is related to the manifestation of "peace" and "help". On topic 7 there are references to "russia" and "sanction", but also "artist" and "joke", terms related to Volodymyr Zelenskyy, the former comedian. Topic 8 is related to the military force of Russia. Among the highest-scoring words are "russian", "war", "military" and "destroy". Finally, topic 9 is also related to the support of Ukraine – "stop", "stand", "provide", "support". There are also mentions of the crimes committed during the conflict – "mass" and "murder".

## 5.4. Models

We proceeded with our prediction phase after completing the necessary pre-processing steps and performing Sentiment Analysis and Topic Modelling on the Twitter data collected from 2022 and 2023. During this phase, we employed various ML models, leveraging the capabilities of the *sklearn* library in Python (<https://pypi.org/project/scikit-learn/>). The classification models used in our analysis included DT, LR, NB, AdaBoost, and XGBoost. Overall, these models demonstrated excellent performance, with XGBoost having the best results.

Our analysis focused on predicting sentiments related to Anti-Russia and Anti-Ukraine, disregarding tweets classified as Neutral. To facilitate the classification, we assigned numerical categories, mapping 'Anti-Ukraine' sentiment to 0 and 'Anti-Russian' sentiment to 1.

In this study, we implemented a rigorous evaluation process to enhance the reliability and performance of our sentiment classification model. Firstly, we divided the dataset into two distinct subsets: a training set and a test set. The training set was used to train the model, allowing it to learn patterns and relationships within the data. Subsequently, the model's performance was assessed on the test set, which served as unseen data, providing an objective measure of its generalization capability. To further validate the model's effectiveness and to enhance the model's performance, we employed cross-validation on the NB algorithm and hyperparameter tuning on the XGBoost algorithm.

In the analysis of the models employed, we relied on four key evaluation metrics to interpret and assess the results: Accuracy, Precision, Recall, and F1 Score. As described by Hutto & Gilbert (2014), Accuracy represents the overall performance of the model. Precision measures the positive classifications made by the model, and it is calculated by dividing the number of true positive classifications by the total number of elements labelled as belonging to the positive class. Recall quantifies the model's ability to correctly identify positive instances out of all the actual positive instances and it is calculated by dividing the number of true positive classifications by the total number of elements known to belong to the positive class. F1 score is the mean of precision and recall. It provides a balanced measure between precision and recall and represents the overall accuracy of the model's performance.

### 5.4.1. Decision Tree

Among the employed ML models, the DT is one of the most used algorithms for classification problems. As a supervised ML approach, it learns from labelled data to discern patterns and successfully classify sentiments of new, untagged tweets. In order to boost the performance of the Sentiment Analysis system, we incorporated stemming and lemmatization techniques during the pre-processing stage. For feature extraction, we employed TF-IDF, a technique that measures the significance of a term amongst a collection of documents, and BoW.

The following table contains the values of Accuracy, Precision, Recall and F1 score for the labels 'Anti-Russia' and 'Anti-Ukraine' for each dataset:

**Table 5.9: Performance Metrics for Decision Tree**

|                         |              | <i>2022 dataset</i> | <i>2023 dataset</i> |
|-------------------------|--------------|---------------------|---------------------|
| <b><i>Accuracy</i></b>  |              | 82%                 | 88%                 |
| <b><i>Precision</i></b> | Anti-Ukraine | 83%                 | 88%                 |
|                         | Anti-Russia  | 82%                 | 89%                 |
| <b><i>Recall</i></b>    | Anti-Ukraine | 83%                 | 89%                 |
|                         | Anti-Russia  | 81%                 | 88%                 |
| <b><i>F1 Score</i></b>  | Anti-Ukraine | 83%                 | 89%                 |
|                         | Anti-Russia  | 82%                 | 88%                 |

By examining the results presented in Table 5.9, we conclude that this model exhibited promising results, achieving an accuracy of 82% for the 2022 and 88% for the 2023 datasets. These findings highlight the model's effectiveness in classifying sentiments within the analysed Twitter data. The value of precision tells us 83% and 82% of the tweets classified as 'Anti-Ukraine' and 'Anti-Russia' for the 2022 dataset, respectively, are true positive classifications among the total number of elements predicted as positive. For the 2023 dataset, 88% of the tweets classified as 'Anti-Ukraine' and 89% of the tweets classified as 'Anti-Russia' are true positive classifications among the total number of elements predicted as positive. The value of recall shows us that, over the total number of elements known to belong to the positive class, for the 2022 dataset 83% ('Anti-Ukraine' sentiment) and 81% ('Anti-Russia' sentiment) are true positive classifications. For the 2023 dataset, 89% ('Anti-Ukraine' sentiment) and 88% ('Anti-Russia' sentiment) are true positive classifications. The value of F1 score, which represents the overall accuracy of the model, is 83% for the tweets classified as 'Anti-Ukraine' and 82% for the tweets classified as 'Anti-Russia' for the 2022 dataset. For the 2023 dataset, the F1 score is 89% for the tweets classified as 'Anti-Ukraine' and 88% for the tweets classified as 'Anti-Russia'.

In summary, the DT model exhibited superior performance for the 2023 dataset compared to the 2022 dataset. This model demonstrates higher precision in classifying Anti-Ukraine sentiment within the 2022 dataset and Anti-Russia sentiment within the 2023 dataset. Regarding recall and F1 score metrics, it excels in classifying Anti-Ukraine sentiment.

#### **5.4.2. Logistic Regression**

Next, we applied the LR algorithm, "one of the best-known and widely used methods" (Sultan et al., 2023) for classification. For this analysis, we utilized the stemmed text and the complete datasets, including all the records. To facilitate text analysis, we employed BoW, which analyses text and documents based on word count, along with Count Vectorizer to convert the text into numerical data.

In Table 5.10, we can observe the metrics that evaluate the performance of the LR algorithm.



**Table 5.10: Performance Metrics for Logistic Regression**

|                         |              | <i>2022 dataset</i> | <i>2023 dataset</i> |
|-------------------------|--------------|---------------------|---------------------|
| <b><i>Accuracy</i></b>  |              | 61%                 | 59%                 |
| <b><i>Precision</i></b> | Anti-Ukraine | 59%                 | 56%                 |
|                         | Anti-Russia  | 67%                 | 72%                 |
| <b><i>Recall</i></b>    | Anti-Ukraine | 83%                 | 90%                 |
|                         | Anti-Russia  | 37%                 | 28%                 |
| <b><i>F1 Score</i></b>  | Anti-Ukraine | 69%                 | 69%                 |
|                         | Anti-Russia  | 47%                 | 40%                 |

Through the examination of Table 5.10, the LR algorithm exhibited lower accuracies than the DT model. It achieved an accuracy of 61% for the 2022 dataset and 59% for the 2023 dataset. Regarding precision, 59% of the tweets classified as 'Anti-Ukraine' and 67% of the tweets classified as 'Anti-Russia' were true positive classifications. Similarly, for the 2023 dataset, 56% of the tweets labelled as 'Anti-Ukraine' and 72% labelled as 'Anti-Russia' were accurately classified as true positives. For recall metric, on the 2022 dataset, the model successfully identified 83% of 'Anti-Ukraine' sentiment and 37% of 'Anti-Russia' sentiment as true positives. In the case of the 2023 dataset, the model achieved 90% recall for 'Anti-Ukraine' sentiment and 28% recall for 'Anti-Russia' sentiment. On the 2022 dataset, the F1 score is 69% for 'Anti-Ukraine' tweets and 47% for 'Anti-Russia' tweets. Meanwhile, for the 2023 dataset, the F1 score reaches 69% for 'Anti-Ukraine' tweets and 40% for 'Anti-Russia' tweets.

In conclusion, the results obtained with this model exhibited a significant discrepancy. In contrast to other models, LR demonstrated higher accuracy for the 2022 dataset. When considering precision, Anti-Russia sentiment outperformed on both datasets. However, the precision values were superior on the 2022 dataset for the Anti-Ukraine sentiment. Moreover, the recall values favoured the Anti-Ukraine sentiment by a substantial margin, extending to the F1 score.

### **5.4.3. Naïve Bayes**

The NB model is "regarded as one of the most effective and efficient inductive learning algorithms, and it has been implemented as an efficient classifier in several different research projects concerning social media" (Sultan et al., 2023). In our study, we employed the entire dataset as an input and utilized the tokens generated during the pre-processing phase for classification. Like the LR model, we used Count Vectorizer to convert textual data into numerical representations. To eliminate any potential bias and validate the model's efficiency, we employed CV on this algorithm. CV involves partitioning the data into multiple subsets, and performing iterations where each subset serves as both training and validation data. The accuracy values obtained with CV for both datasets are presented in Tables 5.11 and 5.12.

**Table 5.11: Cross-validation values for the 2022 dataset**

|                          |            |
|--------------------------|------------|
| <i>First fold of CV</i>  | 0.83446281 |
| <i>Second fold of CV</i> | 0.82281092 |
| <i>Third fold of CV</i>  | 0.81515906 |
| <i>Fourth fold of CV</i> | 0.81018581 |
| <i>Fifth fold of CV</i>  | 0.80888944 |

**Table 5.12: Cross-validation values for the 2023 dataset**

|                          |            |
|--------------------------|------------|
| <i>First fold of CV</i>  | 0.81376761 |
| <i>Second fold of CV</i> | 0.8212173  |
| <i>Third fold of CV</i>  | 0.82373239 |
| <i>Fourth fold of CV</i> | 0.86010563 |
| <i>Fifth fold of CV</i>  | 0.70747841 |

In Table 5.11 and 5.12, we present the accuracy scores obtained for each fold during CV procedure for both datasets. The methodology employed involved a 5-fold cross-validation approach, which means the datasets were divided into 5 subsets or folds. The NB model was trained and evaluated 5 times, each time using a different fold as the test set and the remaining folds as the training set. The accuracy scores obtained were good and provided an estimate of how well this classification model was likely to perform on unseen data. Nonetheless, it is noteworthy that the 2023 dataset exhibited superior performance metrics when compared to the 2022 dataset in the final model assessment.

Table 5.13 contains the values for Accuracy, Precision, Recall and F1 Score for the final NB model.

**Table 5.13: Performance Metrics for Naïve Bayes**

|                         |              | <i>2022 dataset</i> | <i>2023 dataset</i> |
|-------------------------|--------------|---------------------|---------------------|
| <b><i>Accuracy</i></b>  |              | 83%                 | 87%                 |
| <b><i>Precision</i></b> | Anti-Ukraine | 80%                 | 85%                 |
|                         | Anti-Russia  | 87%                 | 91%                 |
| <b><i>Recall</i></b>    | Anti-Ukraine | 90%                 | 92%                 |
|                         | Anti-Russia  | 75%                 | 82%                 |
| <b><i>F1 Score</i></b>  | Anti-Ukraine | 84%                 | 88%                 |
|                         | Anti-Russia  | 80%                 | 86%                 |

The NB model exhibited strong performance in sentiment classification, achieving an accuracy of 83% for the 2022 dataset and 87% for the 2023 dataset. Regarding precision, in the context of the 2022 dataset, the model accurately identified 59% of tweets labelled as 'Anti-Ukraine' and 67% of tweets labelled as 'Anti-Russia' as true positive classifications among those predicted as positive. Similarly, for the 2023 dataset, the model achieved true positive rates of 56% for 'Anti-Ukraine' tweets and 72% for 'Anti-Russia' tweets. Concerning recall, for the 2022 dataset, the model successfully identified 90% of the 'Anti-Ukraine' sentiment instances and 75% of the 'Anti-Russia' sentiment

instances as true positives. In the case of the 2023 dataset, the model achieved even higher recall rates, correctly capturing 92% of the 'Anti-Ukraine' sentiment instances and 82% of the 'Anti-Russia' sentiment instances. The F1 score for the 2022 dataset stands at 84% for 'Anti-Ukraine' tweets and 80% for 'Anti-Russia' tweets, indicating a good overall accuracy in sentiment classification. For the 2023 dataset, the F1 score reaches 88% for 'Anti-Ukraine' tweets and 86% for 'Anti-Russia' tweets, further confirming the model's effectiveness in accurately classifying sentiments.

In a nutshell, the NB model outperformed on the 2023 dataset in contrast to the 2022 dataset. Concerning precision, this model delivered superior results when classifying Anti-Russia sentiment compared to Anti-Ukraine sentiment. However, the situation was reversed when evaluating recall and F1 score metrics.

#### 5.4.4. AdaBoost

The AdaBoost model, also called Adaptive Boosting, is a boosting technique and “one of the most promising, fast convergence, and easy to implement machine learning algorithm” (Wang, 2012). This ensemble method “creates a set of poor learners by maintaining a collection of weights over training data and adjusts them after each weak learning cycle adaptively” (Wang, 2012).

Our study leveraged the sampled datasets and utilized tokenized and lemmatized text as pre-processing techniques for the AdaBoost model. Like the DT, we employed the TF-IDF vectorizer for feature extraction, which helps transform the text data into numerical features suitable for ML models. Regarding the performance evaluation, the AdaBoost model demonstrated very good results in sentiment classification. Table 5.14 shows the performance metrics used to evaluate the model's effectiveness in sentiment classification.

**Table 5.14: Performance Metrics for AdaBoost**

|                  |              | 2022 dataset | 2023 dataset |
|------------------|--------------|--------------|--------------|
| <b>Accuracy</b>  |              | 84%          | 89%          |
| <b>Precision</b> | Anti-Ukraine | 85%          | 90%          |
|                  | Anti-Russia  | 83%          | 88%          |
| <b>Recall</b>    | Anti-Ukraine | 84%          | 88%          |
|                  | Anti-Russia  | 84%          | 89%          |
| <b>F1 Score</b>  | Anti-Ukraine | 85%          | 89%          |
|                  | Anti-Russia  | 84%          | 88%          |

For the 2022 dataset, the model achieved an accuracy of 84%, while for the 2023 dataset, it achieved an even higher accuracy of 89%. In terms of precision, when considering the 2022 dataset, the model exhibited an ability to accurately identify 85% of tweets labelled as 'Anti-Ukraine' and 83% of tweets labelled as 'Anti-Russia' as true positive classifications among those predicted as positive. Similarly, for the 2023 dataset, the model achieved higher true positive rates, correctly classifying 90% of 'Anti-Ukraine'

tweets and 88% of 'Anti-Russia' tweets. Regarding recall, for the 2022 dataset, the model demonstrated strong performance in capturing 84% of the instances associated with 'Anti-Ukraine' sentiment and 84% of those linked to 'Anti-Russia' sentiment as true positive classifications. For the 2023 dataset, the model achieved even higher recall rates, successfully capturing 88% of 'Anti-Ukraine' sentiment instances and 89% of 'Anti-Russia' sentiment instances. Lastly, for the 2022 dataset, the F1 score reached 85% for 'Anti-Ukraine' tweets and 84% for 'Anti-Russia' tweets, indicating a good overall accuracy in sentiment classification. For the 2023 dataset, the F1 score reached 89% for 'Anti-Ukraine' tweets and 88% for 'Anti-Russia' tweets, further confirming the model's effectiveness in accurately classifying sentiments.

In summary, the AdaBoost model demonstrated superior performance on the 2023 dataset compared to the 2022 dataset. The results produced by this model show minimal discrepancy when classifying both sentiments.

### 5.4.5. XGBoost

The XGBoost algorithm, short for Extreme Gradient Boosting, is a powerful and robust ML algorithm. This model is “an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable” (Zhang & Zhan, 2017). As described by Zhang & Zhan (2017), XGBoost leverages the combination of numerous simple trees with low accuracy to construct a more accurate and robust model.

Our analysis employed the XGBoost algorithm using the sampled datasets and the “cleaned\_text” column as inputs. Furthermore, we conducted hyperparameter tuning. Hyperparameters are adjustable parameters that influence the model's behaviour during training and help us identify the best configuration that yields the best performance for the XGBoost algorithm.

Table 5.15 displays the values of the performance metrics used to evaluate the model's success in sentiment classification:

**Table 5.15: Performance Metrics for XGBoost**

|                         |              | <i>2022 dataset</i> | <i>2023 dataset</i> |
|-------------------------|--------------|---------------------|---------------------|
| <b><i>Accuracy</i></b>  |              | 90%                 | 93%                 |
| <b><i>Precision</i></b> | Anti-Ukraine | 91%                 | 94%                 |
|                         | Anti-Russia  | 89%                 | 91%                 |
| <b><i>Recall</i></b>    | Anti-Ukraine | 90%                 | 92%                 |
|                         | Anti-Russia  | 90%                 | 94%                 |
| <b><i>F1 Score</i></b>  | Anti-Ukraine | 90%                 | 93%                 |
|                         | Anti-Russia  | 90%                 | 92%                 |

The XGBoost model achieved high accuracy, reaching 90% for the 2022 and 93% for the 2023 datasets. In terms of precision, for the 2022 dataset, the model demonstrated the ability to accurately classify 91% of tweets labelled as 'Anti-Ukraine' and 89% of tweets labelled as 'Anti-Russia' as true positive instances among those predicted as

positive. Furthermore, for the 2023 dataset, the model achieved even higher true positive rates, correctly identifying 94% of 'Anti-Ukraine' tweets and 91% of 'Anti-Russia' tweets. As for recall, the model's performance was commendable in the 2022 dataset, successfully capturing 90% of instances associated with 'Anti-Ukraine' sentiment and 90% of those linked to 'Anti-Russia' sentiment as true positive classifications. On the 2023 dataset, the model achieved even higher recall rates, accurately capturing 92% of 'Anti-Ukraine' sentiment instances and 94% of 'Anti-Russia' sentiment instances. The F1 score, which combines precision and recall, was calculated for both datasets, reached 90% for 'Anti-Ukraine' tweets and 90% for 'Anti-Russia' tweets, indicating a strong overall accuracy in sentiment classification. For the 2023 dataset, the F1 score reached 93% for 'Anti-Ukraine' tweets and 92% for 'Anti-Russia' tweets, further confirming the model's effectiveness in accurately classifying sentiments.

In review, the XGBoost model exhibited better performance on the 2023 dataset in comparison to the 2022 dataset. Like the AdaBoost model, the results obtained with this model exhibit a relatively small degree of variation.

## 5.5. Conclusion

Among the models evaluated, XGBoost exhibited the highest performance, achieving accuracy rates of 90% for the 2022 dataset and 93% for the 2023 dataset. On the other hand, LR demonstrated the lowest performance, with accuracy rates of 61% for the 2022 dataset and 59% for the 2023 dataset. The higher accuracy of the XGBoost model indicates its superior ability to classify sentiments accurately. The lower accuracy of the LR model suggests limitations in capturing the complexities of sentiment in the given datasets.

In our analysis, it becomes evident that superior performance metrics are obtained when working with the 2023 dataset in contrast to the dataset from 2022. This observed disparity in performance can be attributed to the difference in the volume of available tweets between the two respective datasets. Concerning the performance of the employed ML models, it exhibited distinctive patterns. Notably, LR and NB models demonstrated superior precision when classifying Anti-Russia sentiment, while the recall and F1 score values were better when dealing with the classification of Anti-Ukraine sentiment. The DT model obtained more favourable performance metrics in the classification of Anti-Ukraine sentiment, which can be attributed, in part, to the prevalence of Anti-Ukraine sentiment in both datasets. Finally, AdaBoost and XGBoost models had better results in terms of overall performance metrics, displaying a consistent and small degree of performance variation between the two sentiments under consideration.

Overall, these results provide valuable insights into the performance of various ML models for Sentiment Analysis on Twitter data, highlighting the strengths and weaknesses of each model in classifying Anti-Russian or Anti-Ukrainian sentiments.

## Chapter 6

### Conclusions

In retrospect, the analysis conducted throughout this thesis aimed to fulfill the goal of analysing how Twitter users' sentiment regarding the conflict evolved over six months between August 2022 and February 2023. This chapter focuses on the key findings that have emerged from our investigation.

Our research revealed interesting patterns in Twitter activity surrounding the conflict. In our analysis of the 2022 dataset, we observed that the day with the highest volume of tweets was August 24th, 2022. On that day, Ukraine marked the six-month milestone since the conflict's start. Additionally, August 24<sup>th</sup> is celebrated as Ukraine's Independence Day, contributing to the heightened Twitter activity. In the 2023 dataset, the day with most tweets occurred on February 27<sup>th</sup>, 2023. This date coincides with the marking of the first year since the conflict's initiation. The significance of this milestone increased online discussions, making February 27<sup>th</sup> the day with the highest tweet volume in 2023.

In examining the geographical origin of tweets concerning the conflict, Ukraine and the United States emerged as the top two countries of Twitter activity. However, a closer look at the data revealed a shift in these trends between 2022 and 2023. In 2022, Ukraine took the lead as the country with the highest number of tweets. This surge in Ukrainian Twitter activity can be directly linked to the significant events, namely the Independence Day and the six-month milestone since the outbreak of the war. Conversely, in 2023, the United States emerged as the primary location of tweet origin. This shift was tied to an important event - President Biden's visit to Ukraine, commemorating the first anniversary of the conflict. His visit gained substantial national and international attention and impacted the online discourse.

One of the most intriguing findings of our study is the prevalence of Anti-Ukraine sentiment. This consistent trend is evident in both datasets. As a result, in response to the primary objective of this investigation, we can conclude that the sentiment of Twitter users remains persistently Anti-Ukraine throughout the six-month analysis period in both datasets.

While the difference between both sentiments is not substantial, it holds significance for several reasons. Contrary to the narrative presented in existing literature, which highlights substantial support from Western countries, including the USA, EU, and Canada, for Ukraine since the outbreak of the conflict, our discoveries reveal a different perspective. Notably, these Western countries were among the locations with the highest Twitter activity, besides Ukraine itself. However, it is essential to recognize that our analysis uncovered tweets originating from ambiguous locations such as "Earth", "Everywhere", and "Global". These vague locations may introduce noise into the Sentiment Analysis. Furthermore, our findings deviate from a recent analysis conducted

by bruegel.org, where the authors concluded that "European public opinion remains supportive of Ukraine" (Demertzis et al., 2023).

Regarding the performance of the classification models, when analysing the accuracy values, the models exhibited better results when classifying data from the 2023 dataset, outperforming their behaviour on the 2022 dataset. This disparity in performance can be partly attributed to the difference in the number of tweets between the two datasets. The 2023 dataset has a larger volume of tweets compared to its 2022 counterpart. The increased data volume provided the models with more extensive training and testing sets, contributing to improved accuracy and performance metrics. However, this does not happen with the LR. The consistently poor performance exhibited by this particular model in both the 2022 and 2023 datasets, as reflected in the confusion matrices, strongly suggests that this model is not the best suited for our data.

The DT model demonstrated superior precision, recall, and F1 score when classifying Anti-Ukraine sentiment. This model exhibited a better capacity for accurately identifying tweets expressing sentiments against Ukraine. In contrast, the LR and NB models displayed distinct performance patterns. These algorithms excelled in precision when classifying Anti-Russia sentiment but exhibited stronger recall and F1 score results when classifying Anti-Ukraine sentiment. This discovery contradicts the conventional Sentiment Analysis pattern, where models specialize in recognizing particular sentiment categories. The AdaBoost and XGBoost models demonstrated less differentiation between Anti-Ukraine and Anti-Russia sentiments, with no significant precision, recall, or F1 score disparity.

Among the models evaluated, XGBoost emerged as the standout performer. This model consistently delivered high accuracy, achieving a remarkable 90% accuracy rate for the 2022 dataset and an even more impressive 93% for the 2023 dataset. Comparing our findings to the existing literature, it is evident that our research achieved more robust performance across multiple metrics. While LR and NB were reported as the best-performing models in prior studies, our analysis, particularly the XGBoost model, outperformed these reference points regarding accuracy, precision, recall and F1 score.

The main contribution of this research, setting it apart from other studies, lies in introducing distinct sentiment classes, namely 'Anti-Ukraine' and 'Anti-Russia'. This approach adds a better understanding of public sentiments in the context of the conflict.

Several recommendations can enhance the results of further studies. In the future, it would be beneficial to incorporate DL models to improve accuracy and robustness in NLP tasks, especially with complex and multilingual datasets. Expanding the analysis to integrate languages beyond English would add a more comprehensive understanding of global sentiment and diverse perspectives surrounding the conflict. Increasing the dataset's size and time period would improve the depth and reliability of Sentiment Analysis, allowing for more comprehensive insights into sentiment trends. Lastly, fine-tuning VADER model to better scrutinize both sentiment categories, namely 'Anti-Ukraine' and 'Anti-Russia', and adapting sentiment lexicons to capture the nuances of the conflict would enrich the analysis and enhance its relevance.

## References

- Aggarwal, C. (2018). *Machine Learning for Text*. doi: 10.1007/978-3-319-73531-3
- Al-Hassan, A., & Al-Dossari, H. (2022). Detection of hate speech in Arabic tweets using deep learning. *Multimedia Systems*, 28, 1963-1974 <https://doi.org/10.1007/s00530-020-00742-w>
- Amores, J., Blanco-Herrero, D., Sánchez-Holgado, P., & Frías-Vázquez, M. (2021). Detecting ideological hatred on Twitter. Development and evaluation of a political ideology hate speech detector in tweets in Spanish. *Cuadernos.Info*, 49, 98-123. <https://doi-org/10.7764/cdi.49.27817>
- Balahur, A. (2013). Sentiment Analysis in Social Media Texts. Proceedings of the 4th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, Georgia, 120-128. <https://aclanthology.org/W13-1617>
- Beaumont, P., Koshiw, I. & Harding, L. (2023, February 20th). *Joe Biden visits Kyiv in major show of support for Ukraine*. The Guardian. <https://www.theguardian.com/world/2023/feb/20/joe-biden-visits-ukraine-kyiv> (Accessed on July 6th, 2023)
- Blei, D., Ng, A., & Jordan, M. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3, 993-1022. doi: 10.5555/944919.944937
- Bobichev, V., Kanishcheva, O., & Cherednichenko, O. (2017). Sentiment analysis in the Ukrainian and Russian news. 2017 IEEE First Ukraine Conference on Electrical and Computer Engineering (UKRCON), Kyiv, 1050-1055. doi: 10.1109/UKRCON.2017.8100410
- Bowden, M. (2023, January 27th). *Canada and France follow United States and Germany in pledging tank supplies for Ukraine*. Euronews. <https://www.euronews.com/2023/01/27/canada-and-france-follow-united-states-and-germany-in-pledging-tank-supplies-for-ukraine> (Accessed on July 6th, 2023)
- Cahyana, N., Saifullah, S., Fauziah, Y., Aribowo, A., & Drezewski, R. (2022). Semi-supervised Text Annotation for Hate Speech Detection using K-Nearest Neighbors and Term Frequency-Inverse Document Frequency. *International Journal of Advanced Computer Science and Applications*, 13(10), 147-151. <https://doi.org/10.25139/ijair.v4i2.5267>
- Castells, M. (2008) The New Public Sphere Global Civil Society, Communications Networks, and Global Governance. *The Annals of the American Academy*, 616. <https://doi.org/10.1177/0002716207311877>



Castells, M. (2022). The Network Society Revisited. *American Behavioral Scientist*, 0, 1-7. doi:10.117700027642221092803

Corazza, M., Menini, S., Cabrio, E., Tonelli, S., & Villata, S. (2020). A Multilingual Evaluation for Online Hate Speech Detection. *ACM Transactions on Internet Technology*, 20(2), 1-21 <https://doi.org/10.1145/3377323>

D'Anieri, P. (2023). Conclusion: From Cold War to Hot War. In P. D'Anieri (Ed.). *Ukraine and Russia From Civilized Divorce to Uncivil War* (pp. 308-333). doi: 10.1017/9781009315555

D'Anieri, P. (2023). New World Order? 1989–1993. In P. D'Anieri (Ed.). *Ukraine and Russia From Civilized Divorce to Uncivil War* (pp. 29-65). doi: 10.1017/9781009315555

D'Anieri, P. (2023). Reform and Reversal, 2004–2010. In P. D'Anieri (Ed.). *Ukraine and Russia From Civilized Divorce to Uncivil War* (pp. 135-169). doi: 10.1017/9781009315555

D'Anieri, P. (2023). The Conflict Smolder, 2015-2021. In P. D'Anieri (Ed.). *Ukraine and Russia From Civilized Divorce to Uncivil War* (pp. 243-272). doi: 10.1017/9781009315555

D'Anieri, P. (2023). The Sources of Conflict over Ukraine. In P. D'Anieri (Ed.). *Ukraine and Russia From Civilized Divorce to Uncivil War* (pp. 1-28). doi: 10.1017/9781009315555

D'Anieri, P. (2023). War. In P. D'Anieri (Ed.). *Ukraine and Russia From Civilized Divorce to Uncivil War* (pp. 272-308). doi: 10.1017/9781009315555

Davidson, T., Warmsley, D., Macy, M., & Weber, I. (2017) Automated Hate Speech Detection and the Problem of Offensive Language. *arXiv Preprint arXiv:1703.04009*

Defersha, N., & Tune, K. (2021). Detection of Hate Speech Text in Afan Oromo Social Media using Machine Learning Approach. *Indian Journal of Science and Technology*, 14(31), 2567-2578 <https://doi.org/10.17485/IJST/v14i31.1019>

Del Vigna, F., Cimino, A., Dell'Orletta, F., Petrocchi, M., & Tesconi, M. (2017). Hate me, hate me not: Hate speech detection on Facebook. *Proceedings of the First Italian Conference on Cybersecurity*, Veneza, Itália, 1816, 86-95. <http://ceur-ws.org/Vol-1816/paper-09.pdf>

Demertzis, M., Grand, C., & Moffat, L. (2023, June 5<sup>th</sup>). *European public opinion remains supportive of Ukraine*. bruegel. <https://www.bruegel.org/analysis/european-public-opinion-remains-supportive-ukraine> (Accessed on September 11th, 2023)

Drus, Z., & Khalid, H. (2019). Sentiment Analysis in Social Media and Its Application: Systematic Literature Review. *Procedia Computer Science*, 161, 707-714. doi: 10.1016/j.procs.2019.11.174

Garcia, M., & Cunanan-Yabut, A. (2022). Public Sentiment and Emotion Analyses of Twitter Data on the 2022 Russian Invasion of Ukraine. In *2022 9<sup>th</sup> International Conference on Information Technology, Computer, and Electrical Engineering (ICITACEE)*, 242-247. Doi:10.1109/icitacee55701.2022.9924136

Gatopoulos, D., & Varentytsia, I. (2022, August 24th). *Russia rocket kills 22 at Ukraine train station in Independence Day attack, Zelenskyy says*. PBS NewsHour. <https://www.pbs.org/newshour/world/russian-rocket-kills-15-at-ukraine-train-station-in-independence-day-attack-zelenskyy-says> (Accessed on July 5th, 2023)

Gonçalves, P., Araújo, M., Benevenuto, F., & Cha, M. (2013). Comparing and Combining Sentiment Analysis Methods. Proceedings of the first ACM conference on Online social networks, Boston, 27-38. doi: 10.48550/arXiv.1406.0032

Hajibabae, P., Malekzadeh, M., Ahmadi, M., Heidari, M., Esmaeilzadeh, A., Abdolazimi, R., & Jones, J. (2022). Offensive Language Detection on Social Media Based on Text Classification. In *2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC)*, 92-98. DOI: 10.1109/CCWC54503.2022.9720804

Hasan, A., Sharma, T., Khan, A., & Al-Abyadh, M. (2022) Analysing Hate Speech against Migrants and Women through Tweets Using Ensembled Deep Learning Model. *Hindawi*, 2022, 1-8. <https://doi-org/10.1155/2022/8153791>

Hasan, M., Islam, L., Jahan, I., Meem, S., & Rahman, R. (2023). Natural Language Processing and Sentiment Analysis on Bangla Social Media Comments on Russia-Ukraine War Using Transformers. *Vietnam Journal of Computer Science*. doi: 10.1142/S2196888823500021

Hauter, J. (2021). Forensic conflict studies: Making sense of war in the social media age. *Media, War & Conflict*, 0, 1-20. doi: 10.1177/17506352211037325

Hickman, L., Thapa, S., Tay, L., Cao, M., & Srinivasan, P. (2022). Text Preprocessing for Text Mining in Organizational Research: Review and Recommendations. *Organizational Research Methods*, 25(1), 114-146. doi: 10.1177/1094428120971683

Hussein, D. (2018). A survey on sentiment analysis challenges. *Journal of King Saud University – Engineering Sciences*, 30, 330-338. doi: 10.1016/j.jksues.2016.04.002

Hutto, C., & Gilbert, E. (2014) VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. *Proceedings of the Eight International AAAI Conference on Weblogs and Social Media*, 216-225.

<http://www.aaai.org/ocs/index.php/ICWSM/ICWSM14/paper/view/8109%5Cnhttp://comp.social.gate>

Jelodar, H., Wang, Y., Yuan, C., Feng, X., Jiang, X., Li, Y., & Zhao, L. (2019). Latent Dirichlet allocation (LDA) and topic modelling: models, applications, a survey. *Multimedia Tools and Applications*, 78, 15169-15211. doi: 10.1007/s11042-018-6894-4

Kiilu, K., Okeyo, G., Rimiru, R., & Ogada, K. (2018). Using Naïve Bayes Algorithm in detection of Hate Tweets. *International Journal of Scientific and Research Publications*, 8(3), 99-107 DOI: 10.29322/IJSRP.8.3.2018.p7517

Koroutchev, R. (2023). The Explosive Ukrainian Migration due to the Russian armed Conflict in 2022: the case of Bulgaria. *Journal of Liberty and International Affairs*, 9(1), 303-311. doi: 10.47305/JLIA2391309k

MacAvaney, S., Yao, H., Yang, E., Russell, K., Goharian, N., & Frieder, O. (2019). Hate Speech detection: Challenges and solutions. *PLoS ONE*, 14(8), 1-16. <https://doi.org/10.1371/journal.pone.0221152>

Makhortykh, M., & Sydorova, M. (2017). Social media and visual framing of the conflict in Eastern Ukraine. *Media, War & Conflict*, 10(3), 359-381. doi: 10.1177/1750635217702539

Marchellim, G., & Ruldeviyani, Y. (2021). Sentiment analysis of hate speech as an information tool to prevent riots and environmental damage. *IOP Conf. Series: Earth and Environmental Science*, 700, 1-10. doi: 10.1088/1755-1315/700/1/012024

Maruf, A., Ziyad, Z., Haque, M., & Khanam, F. (2022). Emotion Detection from Text and Sentiment Analysis of Ukraine Russia War using Machine Learning Technique. *International Journal of Advanced Computer Science and Applications*, 13(20), 868-882. doi: 10.14569/IJACSA.2022.01312101

McMahon, R. (2014). Ukraine in Crisis. *Council on Foreign Relations*, 1-8

Mutanga, R., Naicker, N., & Olugbara, O. (2022). Detecting Hate Speech on Twitter Network Using Ensemble Machine Learning. *International Journal of Advanced Computer Science and Applications*, 13(3), 331-339 <https://doi.org/10.14569/IJACSA.2022.0130341>

Nayak, A., Kanive, A., Chandavekar, N., & Ramasamy, B. (2016). Survey on Pre-Processing Techniques for Text Mining. *International Journal of Advanced Trends in Computer Science and Engineering*, 5(6), 16875-16879. doi: 10.18535/ijecs/v5i6.25

Noonan, J. (2023). Ukraine Conflict as a Case of the Political Contradictions of Contemporary Imperialism. *International Critical Thought*, 13(1), 1-21. doi: 10.1080/21598282.2022.2163416

O'Connell, M. (2017). The Crisis in Ukraine 2014-. *International Law and the Use of Force: A Case-Based Approach*, Olivier Corten and Tom Ruys, eds, Oxford University Press, Forthcoming, Notre Dame Law School Legal Studies Research Abstract No. 1720, 1-22

Odey, S., & Basse, S. (2022). Ukrainian foreign policy toward Russia between 1991 and 2004: the start of the conflict. *Journal of Liberty and International Affairs*, 8(2), 346-361. doi: 10.47305/JLIA2282346a

Öztürk, N., & Ayvaz, S. (2018). Sentiment analysis on Twitter: A text mining approach to the Syrian refugee crisis. *Telematics and Informatics*, 35, 136-147. <https://doi.org/10.1016/j.tele.2017.10.006>

Prasetyo, V. & Samudra, A. (2022). Hate speech content detection system on Twitter using K-nearest neighbor method. *AIP Conference Proceedings 2470*, 050001, 1-10. <https://doi.org/10.1063/5.0080185>

Putri, T., Sriadhi, S., Sari, R., Rahmadani, R., & Hutahaean, H. (2020). A comparison of classification algorithms for hate speech detection. *IOP Conference Series: Materials Science and Engineering*, 830, 1-6. doi: 10.1088/1757-899X/830/3/032006

Ratten V. (2022). The Ukraine/Russia conflict: Geopolitical and international business strategies. *Thunderbird International Business Review*, 1-7. doi: 10.1002/tie.22319

Ruwandika, N., & Weerasinghe, A. (2018). Identification of Hate Speech in Social Media. *2018 International Conference on Advances in ICT for Emerging Regions*, 273-278. DOI: 10.1109/ICTER.2018.8615517

Schwarz, P. (2022, December 21). *Former German Chancellor Merkel admits the Minsk agreement was merely to buy time for Ukraine's arms build-up*. World Socialist Web Site. <https://www.wsws.org/en/articles/2022/12/22/ffci-d22.html> (Accessed on July 4, 2023)

Šerstka, A. (2021). *Big Data: A New Perspective on Conflict Resolution* (Master's thesis, Charles University). Retrieved from <https://dspace.cuni.cz/handle/20.500.11956/124611>, 1-67

Smart, B., Watt, J., Benedetti, S., Mitchell, L., Roughan, M. (2022). #IStandWithPutin Versus #IStandWithUkraine: The Interaction of Bots and Humans in Discussion of the Russia/Ukraine War. In: *Hopfgartner, F., Jaidka, K., Mayr, P., Jose, J., Breitsohl,*

J. (eds) *Social Informatics. SocInfo 2022. Lecture Notes in Computer Science*, 13618, 1-14. [https://doi.org/10.1007/978-3-031-19097-1\\_3](https://doi.org/10.1007/978-3-031-19097-1_3)

Sultan, D., Toktarova, A., Zhumadillayeva, A., Aldeshov, S., Mussiraliyeva, S., Beissenova, G., Tursynbayev, A., Baenova, G., & Imanbayeva, A. (2023). Cyberbullying-related Hate Speech Detection Using Shallow-to-deep Learning. *Computers, Materials & Continua*, 74(1), 2115-2131 DOI: 10.32604/cmc.2023.032993

Turki, T., & Roy, S. (2022). Novel Hate Speech Detection Using Word Cloud Visualization and Ensemble Learning Coupled with Count Vectorizer. *Applied Sciences*, 12, 1-13. <https://doi.org/10.3390/app12136611>

Underwood, N., Wood, D., & Hurt, A. (2022, August 29th). *Russia-Ukraine war: A weekly recap and look ahead (August 29th)*. National Radio Public. <https://www.npr.org/2022/08/29/1118820771/russia-ukraine-war-a-weekly-recap-and-look-ahead-aug-29> (Accessed on July 5th, 2023)

Underwood, N., Wood, D., & Hurt, A. (2022, September 5th). *Russia-Ukraine war: A weekly recap and look ahead (September 5th)*. National Radio Public. <https://www.npr.org/2022/09/05/1119926980/russia-ukraine-war-a-weekly-recap-and-look-ahead-sept-5> (Accessed on July 5th, 2023)

Underwood, N., Wood, D., & Hurt, A. (2023, February 27th). *Latest in Ukraine: As Russia's war starts its second year, eyes are on China (February 27th)*. National Radio Public. <https://www.npr.org/2023/02/27/1159636226/russia-ukraine-war-latest-news-updates-feb-27> (Accessed on July 5th, 2023)

Vijayarani, S., Ilamathi, J., & Nithya (2015). Preprocessing Techniques for Text Mining – An Overview. *International Journal of Computer Science & Communication Networks*, 5(1), 7-16.

Wadhvani, G., Varshney, P., Gupta, A., & Kumar, S., (2023). Sentiment Analysis and Comprehensive Evaluation of Supervised Machine Learning Models Using Twitter Data on Russia-Ukraine War. *SN Computer Science*, 4(346), 1-11. doi: 10.1007/s42979-023-01790-5

Wang, R. (2012). AdaBoost for Feature Selection, Classification and Its Relation with SVM, A Review. *Physics Procedia: 2012 International Conference on Solid State Devices and Materials Science*, China, 25, 800-807. doi: 10.1016/j.phpro.2012.03.160

Wankhade, M., Rao, A., & Kulkarni, C. (2022). A survey on sentiment analysis methods, applications, and challenges. *Artificial Intelligence Review*, 55, 5731-5780. doi: 10.1007/s10462-022-10144-1

Yue, L., Chen, W., Li, X., Zuo, W., & Yin, M. (2019). A survey of sentiment analysis in social media. *Knowledge and Information Systems*, 60, 617-663. doi: 10.1007/s10115-018-1236-4

Zhang, L., & Zhan, C. (2017). Machine Learning in Rock Facies Classification: An Application of XGBoost. *International Geophysical Conference*, Qingdao, China, 1371-1374. doi: 10.1190/IGC2017-351

Zhu, Y., Haq, E., Lee, L., Tyson, G., & Hui, P. (2022). A Reddit Dataset for the Russo-Ukrainian Conflict in 2022, 1-7. *arXiv:2206.05107*



## Appendix Index

|  |    |
|--|----|
| Appendix 1: WordCloud for 2022 dataset.....                                      | 60 |
| Appendix 2: WordCloud for 2023 dataset.....                                      | 60 |
| Appendix 3: List of variables after data cleaning:.....                          | 60 |
| Appendix 4: WordCloud for the ‘Anti-Ukraine’ sentiment on the 2022 dataset ..... | 61 |
| Appendix 5: WordCloud for the ‘Anti-Russia’ sentiment on the 2022 dataset.....   | 61 |
| Appendix 6: WordCloud for the ‘Anti-Ukraine’ sentiment on the 2023 dataset ..... | 62 |
| Appendix 7: WordCloud for the ‘Anti-Russia’ sentiment on the 2023 dataset.....   | 62 |
| Appendix 8: WordCloud for August 24 <sup>th</sup> , 2022 .....                   | 62 |
| Appendix 9: WordCloud for August 30 <sup>th</sup> , 2022 .....                   | 63 |
| Appendix 10: WordCloud for September 6 <sup>th</sup> , 2022.....                 | 63 |
| Appendix 11: WordCloud for September 9 <sup>th</sup> , 2022.....                 | 63 |
| Appendix 12: WordCloud for February 7th, 2023 .....                              | 64 |
| Appendix 13: WordCloud for February 21st, 2023.....                              | 64 |
| Appendix 14: WordCloud for February 25th, 2023 .....                             | 65 |
| Appendix 15: WordCloud for February 27 <sup>th</sup> , 2023 .....                | 65 |



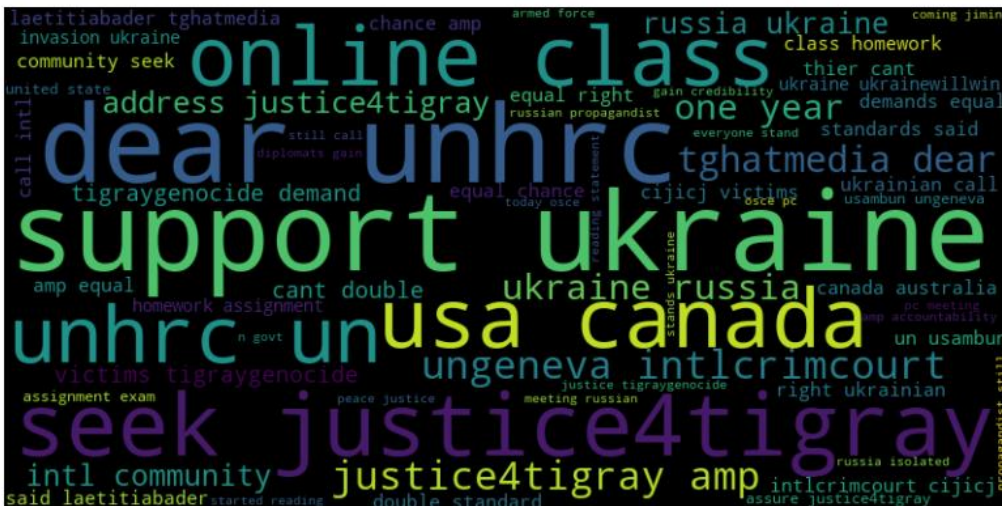




**Appendix 6: WordCloud for the ‘Anti-Ukraine’ sentiment on the 2023 dataset**



**Appendix 7: WordCloud for the ‘Anti-Russia’ sentiment on the 2023 dataset**



**Appendix 8: WordCloud for August 24<sup>th</sup>, 2022**



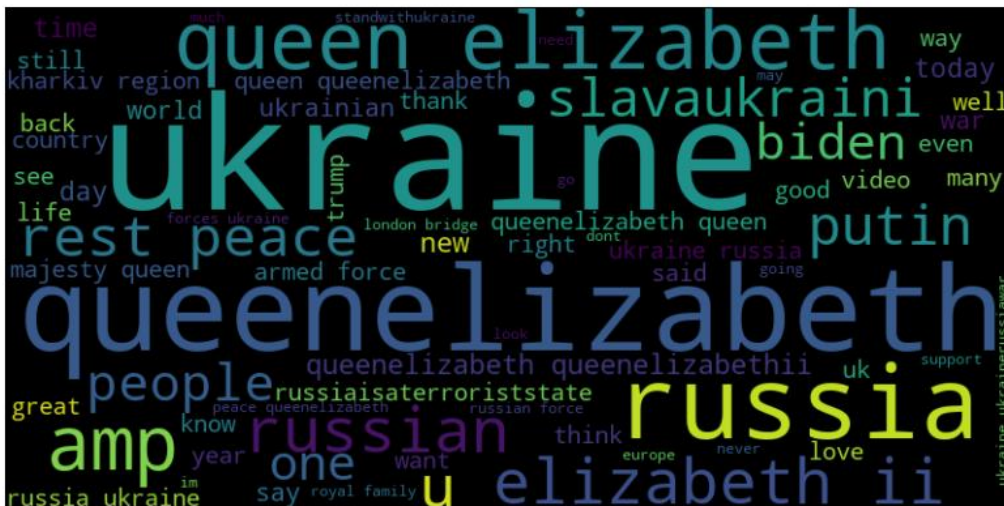
Appendix 9: WordCloud for August 30<sup>th</sup>, 2022



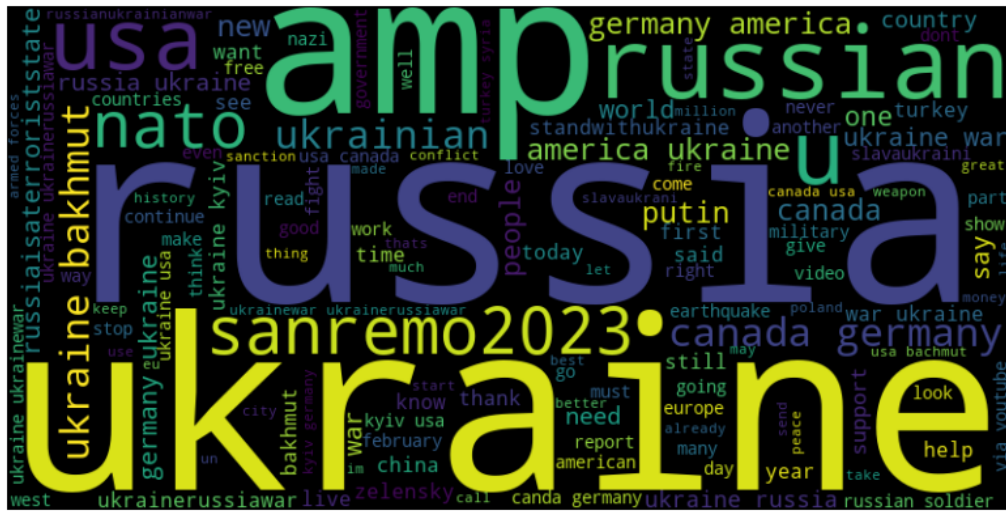
Appendix 10: WordCloud for September 6<sup>th</sup>, 2022



Appendix 11: WordCloud for September 9<sup>th</sup>, 2022



Appendix 12: WordCloud for February 7th, 2023



Appendix 13: WordCloud for February 21st, 2023

