

# MidTerm Progress Report

**Stavroula Afroditi Christopoulou**

ID: 11412984

Masters: Data Science

Uva Information Studies

**Supervisor**

dr.Stevan Rudinac

## Abstract

The remarkable emergence of numerous supporters of extreme doctrines, triggered the elaboration of this thesis project, which is based on the Political Extremism content of social media sources. Through this research, social, interpersonal, political and national level findings are expected to be revealed, under the prism of Multimodal Sentiment Analysis. .

## 1 Multimodal Sentiment Analysis on Political Extremism Content

The main research question of this thesis is Whether it is possible to estimate "extremeness" of ideologies in extremist groups, based on automatic, multi-modal sentiment analysis. Followed by the subquestions whether it is possible to identify which of the posts are highly positive/negative and whether women extremists can express extremeness in the same levels that men extremists do. For the purposes of this research, the big amount of data is processed with DAS-4 supercomputer, and the research findings are approached with Vader Sentiment analysis which is a Python valence based package tool that analyses people's opinions and gives positive, neutral, negative and compound scores of sentiment.

At this initial stage of research, data have been gathered and preprocessed, while extraction and localisation of interesting relationships between the data have been accomplished. Main focus constitutes the creation of a clear data structure, so that the decomposition, selection and deeper investigation of data categories and relationships between the data would be accurate and manipulable for investigation.

Research is based on a data set that includes sensitive information of extremist content from

extremist users that participate on a social networking website. It consists of Json files divided into entities (categories, content of text, entities, publication date, replied or quoted, topics, user) and users (avatars, category, profile picture, topics, user, user id), and the main focus was the creation of access to any category of the data that seems interesting and inspiring for further exploration.

Starting with the data cleaning, I created access to the categories and topics on both files, users and entities, and subsequently deeper access and extraction of the most active users, most active categories and most active topics, as well as to the content of the posts per year, category or season (divided in quarters) or by topic. All this material after sentiment analysis process, will reveal interesting relationships which is the most important purpose of this research.

According to the preliminary research, some interesting findings support the basic research question that extremeness can be estimated through automatic multimodal sentiment analysis, however, the most interesting part is the serendipitous findings that have been revealed and opened the way to a series of zestful new subquestions. Following the path that the metadata and available data have revealed, it was challenging to research which were the most popular extremist users and which entities they prefer (as well as the total number of their posts), in order to investigate the entities with the highest traffic and the levels of sentiment at these specific posts.

A step closer to the answer of my research questions have been taken by examining specific countries with different political regimes (:Italy, South Africa, Canada, Britain and Russia) and comparing the sentiment analysis results, trying to reveal possible correlation between "extreme" elation and political status quo. Furthermore, I extracted the yearly activity from all the content to

compare the amount of increase or decrease in use, based on important political social events that have been held and could be correlated with extremism.

Dividing the content of users posts in 4 seasons, and extracting the sentiment scores in each period of the year using the available temporal metadata, it was very interesting to inspect whether seasonality affects the temperament, behaviour and extremeness of users expression.

Finally, one of the basic aforementioned sub-questions (women compared to men extremeness), was explored by extracting a special category of content that is used only by female users of this networking website, opposed to an averaged sample of general population members that was extracted from the dataset.

This preliminary phase of the research, brings constantly into light numerous interesting correlations and relationships between the data, in a never ending procedure of searching and inspecting.

Having as groundwork this first stage of research, the appropriate literature, practical tools and exploration instinct, the research questions are expected to be answered and reveal interesting relationships between the data, creating a solid and complete thesis project.

### 1.1 Some indicative results

The 5 most active users and the number of their posts:

[(u'revision', 19992), (u'ADAMANT', 11668), (u'kazan188', 8990), (u'junkers88', 8616), (u'eyzwydopen', 7955)]

As well as the Vader sentiment analysis on the content of these posts:

The positive average in revision is: 0.022688936426964644 The negative average in revision is: 0.037019856288738495 The neutral average in revision is: 0.9399663759530368 The compound average in revision is: -0.1309411222642758

The positive average in ADAMANT is: 0.09557305630026798 The negative average in ADAMANT is: 0.07370670241286878 The neutral average in ADAMANT is: 0.8307202412868613 The compound average in ADAMANT is: 0.08663483914209089

The positive average in kazan188 is: 0.01923228048964795 The negative average in kazan188 is: 0.02856354843584701 The neutral average in kazan188 is: 0.952049418165329

The compound average in kazan188 is: -0.07671949523953422

The positive average in junkers88 is: 0.07600817691878851 The negative average in junkers88 is: 0.08967589667348064 The neutral average in junkers88 is: 0.8337576658613618 The compound average in junkers88 is: -0.1487218175060397

The positive average in eyzwydopen is: 0.08077365079365094 The negative average in eyzwydopen is: 0.11226920634920624 The neutral average in eyzwydopen is: 0.8069688888888878 The compound average in eyzwydopen is: -0.2374127301587304

The 5 most active topics and the number of posts are: 'Pic thread!', 11162, 'Ron Paul Revolution with 2,500,000 views', 10920, 'Tales of the Holocaust', 6883, '-The Hellenic thread', 6144, 'Addressing Filth.', 3736

To find the general population Sentiment and compare it with "For Stormfront Ladies only", which is an extremists' women category, I aggregated three categories with high number of posts (Lounge, Politics Continuing Crises, Strategy and Tactics) and found the mean sentiment:

The aggregated positive average is: 0.11310993616524327 The aggregated negative average is: 0.08380836725449455 The aggregated neutral average is: 0.8008926626930194 The aggregated compound average is: 0.08856957581687702 And for women, Stormfront Ladies only category: The positive average in For Stormfront Ladies Only is: 0.1553962861363155 The negative average in For Stormfront Ladies Only is: 0.06480544418653662 The neutral average in For Stormfront Ladies Only is: 0.7793409298726913 The compound average in For Stormfront Ladies Only is: 0.3100517619751059

Which shows that women are expressing more positively, confirming the social stereotypes.

All code and results, are in my git hub account (<https://github.com/StavroulaChristopoulou/ThesisProject.git>)