



**Politecnico  
di Torino**

MASTER'S DEGREE IN PHYSICS OF COMPLEX SYSTEMS

---

**Data-driven characterization of viral  
events on social networks: sustainability  
issues in the palm oil production**

---

*Author:*  
Elena Candellone

*Supervisors:*  
Yamir Moreno  
Alberto Aleta  
Henrique Ferraz De Arruda

*Master's Thesis realized in collaboration with*  
ISI Foundation

October 2022

POLITECNICO DI TORINO

## *Abstract*

Master's Degree in Physics of Complex Systems

**Data-driven characterization of viral events on social networks: sustainability issues in the palm oil production**

by Elena Candellone

Palm oil is the most widely used vegetable oil in the world. However, its production and consumption have generated heated debates over the past few decades due to its environmental impact. This thesis studies the debate on palm oil on the popular social network Twitter since 2006. Using Natural Language Processing and Network Science tools, we analyze the most important viral events related to palm oil. We identify Opinion Drivers' role in the debate and how most debates are short-lived. Indeed, by studying the interevent time distributions of specific hashtags, we see that even the most far-reaching viral events are quickly forgotten. Furthermore, most viral events are described by similar characteristics, showing an underlying universality that goes beyond the specific topics. All in all, our results show that the public debate on Twitter is limited to a few countries and mainly centered around the leading actors of public opinion. Thus, rather than considering this debate intrinsic to the public, it should be regarded as mainly driven by a few organizations.

## Acknowledgements

*A Trieste*, le bottiglie risuonano quando spira la Bora.

*A Torino*, la bellezza di sentirsi a casa.

*A Parigi*, il ritmo pulsa incessante e frenetico.

*A Yamir*, por inspirarme y guiarme con sabiduría.

*A Alberto*, el mejor cicerone de Zaragoza y maestro del 'biased random walk'.

*A Henrique*, por sua imensa gentileza e paciência.

*A Ariadna*, el professional del disseny gràfic.

*Al minitaglio*, miglior invenzione dell'essere umano.

*A Toccalmatto*, tra una Zona Cesarini e una Madame Satan.

*A PCS*, per i compleanni, matrimoni, battesimi e feste di laurea.

*A Pietro*, per le improvvisate scene del crimine e le pizze all'ananas sotto la pioggia.

*A Frappa*, per le infinite chiacchierate e le birre fatte in casa.

*A Salvo*, per i volantini passati e quelli che verranno.

*A Greivin y Diana*, por los chilaquiles y las pizzas al tegamino en compañía.

*A Sam e Claudio*, per tutti i tappi che il barattolo può ancora contenere.

*Alla GANG*, per avermi riaccolta e amata.

*A Lucia*, per avermi insegnato a *resistere*.

*A Mathias e Tunnus*, per aver provato insieme a rimarginare le ferite.

*Ad Ari*, per sempre nei nostri cuori, non ti dimenticheremo.

*A Nonno Sandro*, per avermi insegnato a fischiare  
e per esser fiero di me, ovunque vada.

*A Nonna Angela*, per avermi ispirata con i racconti dei vostri viaggi,  
cercando nell'oblio il filo dei tuoi ricordi.

*A Sofia*, per tutti i treni che ancora prenderemo insieme.

*A Mamma e Papà*, per supportarmi e sopportarmi, sempre.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>ii</b>
<b>Introduction</b>	<b>1</b>
<b>1 Twitter dataset: an overview</b>	<b>3</b>
1.1 Basics of Twitter . . . . .	3
1.2 Collecting data on Twitter . . . . .	3
<b>2 Exploratory Data Analysis</b>	<b>6</b>
2.1 Time evolution of the number of Tweets . . . . .	6
2.2 Keywords and Hashtags . . . . .	8
2.3 Popular Users . . . . .	8
2.4 Geographical distribution . . . . .	11
<b>3 NLP tools: Sentiment Analysis and Topic Modeling</b>	<b>12</b>
3.1 Sentiment Analysis . . . . .	12
3.1.1 Model description and general results . . . . .	12
3.1.2 Sentiment dynamics . . . . .	13
3.1.3 Popularity and sentiment . . . . .	15
3.1.4 Sentiment analysis and Opinion Mining: are they synonyms? . . . . .	16
3.2 Topic Modeling . . . . .	16
<b>4 Network science</b>	<b>20</b>
4.1 Building blocks of Network Science . . . . .	20
4.2 Community detection . . . . .	21
4.3 Real-data networks . . . . .	23
4.3.1 Co-mentions network . . . . .	23
4.3.2 Co-hashtags network . . . . .	24
4.3.3 User-hashtag network . . . . .	25
4.4 Degree distribution properties: an overview . . . . .	26
4.5 Temporal networks: evolution of modularity and nestedness . . . . .	27
<b>5 Interevent time distribution</b>	<b>31</b>
5.1 Heavy-tailed distributions of IET . . . . .	32
5.2 Time maps: a tool for multiple timescales visualization . . . . .	33
5.3 Viral events and Negative sentiment . . . . .	34
<b>6 Cascade size distribution</b>	<b>38</b>
6.1 Hashtags cascade . . . . .	38
6.2 Sentiment cascade . . . . .	40



<b>7</b>	<b>Virality phase diagram</b>	<b>42</b>
7.1	Palm oil dataset . . . . .	43
7.2	Virality phase diagram: comparison among oils . . . . .	44
<b>8</b>	<b>Conclusion</b>	<b>47</b>
<b>A</b>	<b>Coconut and Olive: comparison among datasets</b>	<b>49</b>
A.1	Evolution of the interest . . . . .	49
A.2	Most popular hashtags: topic detection . . . . .	51
A.3	Sentiment analysis . . . . .	52
A.4	IET distribution . . . . .	53
A.4.1	Coconut dataset . . . . .	53
A.4.2	Olive dataset . . . . .	53
A.5	CS distribution . . . . .	55
A.5.1	Coconut dataset . . . . .	55
A.5.2	Olive dataset . . . . .	55
	<b>Bibliography</b>	<b>58</b>
	<b>Funding</b>	<b>62</b>

# List of Figures

2.1	Histogram representation of the number of Tweets containing the words <i>palm oil</i> per year. A peak of interest was detected in 2018. . . . .	7
2.2	Histogram representation of the number of Tweets containing the words <i>palm oil</i> , restricted to November 2018. The interest increased on the 9th of November when the Tweet in Fig. 2.3 was published. . . . .	7
2.3	Screenshot of Iceland foods' Tweet about the banned TV spot. The Tweet received more than 14,700 likes and 8,000 reTweets. . . . .	7
2.4	Histogram representation of the most popular hashtags in English. <i>palmoil</i> , <i>palm</i> , <i>oil</i> were removed, sustainability-related hashtags are predominant. . . . .	8
2.5	Wordcloud of the most common keywords in the corpus. Keywords related to sustainability are widely used. . . . .	9
2.6	Geographical distribution of Tweets from the <b>All-languages</b> dataset. The UK, the US, Malaysia, Nigeria, and Indonesia are the countries that are Tweeting the most about the palm oil topic. . . . .	11
3.1	Pie chart representation of the text classification task results. A prevalence of <i>Neutral</i> labels is highlighted, followed by the <i>Negative</i> labels. Just a few Tweets are classified as <i>Positive</i> . . . . .	14
3.2	Time evolution of sentiment. The plot represents the percentage of labeled Tweets published in a given year. <i>Positive</i> sentiment remains low and stable through time, while <i>Negative</i> sentiment has a peak around 2018. . . . .	14
3.3	Pie chart representation of the text classification task results for the <b>OD</b> . The percentage of <i>Positive</i> labels is higher compared to the <b>GP</b> case (Fig. 3.1). . . . .	15
3.4	Comparison of sentiment dynamics between the <b>GP</b> and the <b>OD</b> . Proportionally, the <b>OD</b> show a stronger growth in the <i>Positive</i> sentiment than the <b>GP</b> . . . . .	16
3.5	Wordcloud of the most common keywords for the two labeled subsets of the main corpus. In both cases, the keywords refer to sustainability issues. . . . .	17
3.6	Results of <b>LSA</b> using the <b>Tf-idf</b> vectorization. The high-dimensional occurrence matrix <b>M</b> is projected along two directions, representing the main distinctions into topics. Each point in the plot represents a Tweet of the dataset, labeled as the oil it refers to. . . . .	19
4.1	Degree distribution for the co-mentions network: the blue and the red dotted lines represent the fitted power-law and truncated power-law, respectively, while the green scatter plot represents the empirical data in logarithmic binning. . . . .	24

4.2	Degree distribution for the co-hashtags network: the blue and the red dotted lines represent the fitted power-law and truncated power-law, respectively, while the green scatter plot represents the empirical data in logarithmic binning. . . . .	25
4.3	Degree distribution for the user-hashtag network: the blue and the red dotted lines represent the fitted power-law and truncated power-law, respectively, while the green scatter plot represents the empirical data in logarithmic binning. . . . .	26
4.4	Graphical visualization of the regimes identified by the average degree and the scaling exponent of the degree distribution. On the y-axis, the average degree is normalized by the natural logarithm of the number of nodes, assuming that none of these networks has an average degree smaller than 1. Therefore the <i>subcritical regime</i> is not shown. . . . .	28
4.5	(a)-(b): Normalized time evolution of modularity and nestedness for temporal networks, realized on time windows of three years. The values are given as percentages of the maximum value to display the change in time, more than the absolute values. (c)-(d): Time evolutions of modularity and nestedness. Here, values are absolute to display the change in magnitude in the different networks. . . . .	30
5.1	IET distributions of the ten most used hashtags. The red and blue dashed lines represent the fitted power laws and truncated power laws, while the scatter plots represent the data we are fitting. . . . .	34
5.2	(a): graphical representation of the time intervals between Tweets containing a given hashtag. (b): graphical representation of the time map construction. Each event has the time before the event as an x-coordinate and the time after the event as a y-coordinate. . . . .	35
5.3	Time maps of the most common hashtags. The plot is in log-log scale and each event $e_i$ has coordinates $(t_{i-1}, t_i)$ for $i = 1, \dots, n - 1$ . The heat map colors represent in which year this event happened. . . . .	35
5.4	Time evolution of the IET and sentiment's rolling averages for the most common hashtags. The IET is normalized, such that $r_i \in [0, 1]$ and the weighted sentiment is $\langle l(t) \rangle \in [0, 2]$ . We suggest a strong correlation in the "viral hashtags" between shorter IETs and a drop in the weighted sentiment. . . . .	37
6.1	Time evolution of the cascade size for the hashtag <i>palmoil</i> . The cascade size $CS(h)$ is calculated as the number of Tweets per day containing a given hashtag $h$ . . . . .	39
6.2	Hashtags cascade distributions of the most common hashtags. The red and blue lines represent the fitted power-law and truncated power-law distributions. The scatter plot represents the empirical data (in logarithmic binning) that we are fitting. The red and blue curves overlap for each hashtag. . . . .	40
6.3	Sentiment cascades. The scatter plots represent the empirical data fitted by heavy-tailed distribution, while the lines represent the distributions that best fit the data. <i>Negative</i> sentiment is more viral than the <i>Positive</i> one, while the <i>Neutral</i> sentiment is out of the virality regime. . . . .	41

7.1	Virality phase diagram determined by the scaling exponents of the <b>IET</b> and <b>CS</b> distributions. The regions are defined by the mathematical properties of the distributions, which depend on the scaling exponents: for $\alpha < 2$ , all the moments are diverging, while for $2 < \alpha < 3$ , the first moment (average) is finite and the higher moments diverge. Finally, for $\alpha > 3$ , the second moment is also finite. . . . .	44
7.2	Virality phase diagram composed of the most common hashtags of the three datasets ( <b>Olive</b> , <b>Palm</b> , <b>Coconut</b> ). Different markers and colors identify the belonging to the different datasets. We can notice that only <b>Palm</b> hashtags are found in the <i>viral region III</i> , while most of the <b>Olive</b> hashtags are found as non-viral. . . . .	46
7.3	Virality phase diagram composed of the most common hashtags of the three datasets ( <b>Olive</b> , <b>Palm</b> , <b>Coconut</b> ). Different markers identify the belonging to the different datasets, while the prevalent sentiment determines the colors. We can notice that the hashtags in region <b>III</b> are prevalently <i>Negative</i> , while neutrality predominates in the less viral regions ( <b>IV-V-VI</b> ). . . . .	46
A.1	Histogram representation of the number of Tweets per year containing the words <i>coconut oil</i> and <i>olive oil</i> respectively. . . . .	50
A.2	The relative growth is calculated as the number of Tweets in a given year divided by the number of Tweets in a fixed reference year. Here the results for $\text{ref\_year} = 2007$ are shown for the three datasets and an average of the relative growths of the most common words in English. . . . .	50
A.3	Histogram representation of the most common hashtags for the two datasets. . . . .	51
A.4	Pie chart representation of the text classification task results. . . . .	52
A.5	Time evolution of sentiment. The plot represents the percentage of labeled Tweets published in a given year. . . . .	52
A.6	<b>IET</b> distributions of the ten most used hashtags in the <b>Coconut</b> dataset. The red and blue dashed lines represent the fitted power laws and truncated power laws, while the scatter plots represent the data we are fitting. . . . .	53
A.7	<b>IET</b> distributions of the ten most used hashtags in the <b>Olive</b> dataset. The red and blue dashed lines represent the fitted power laws and truncated power laws, while the scatter plots represent the data we are fitting. . . . .	54
A.8	Hashtags cascade distributions of the most common hashtags ( <b>Coconut</b> dataset). The red and blue lines represent the fitted power-law and truncated power-law distributions. At the same time, the scatter plot is the empirical data (in logarithmic binning) that we are fitting. . . . .	56
A.9	Hashtags cascade distributions of the most common hashtags. The red and blue lines represent the fitted power-law and truncated power-law distributions. At the same time, the scatter plot is the empirical data (in logarithmic binning) that we are fitting. . . . .	57

# List of Tables

1.1	Number of Tweets and Users for each dataset obtained from the keyword-based Data Collection process. . . . .	5
2.1	Potential opinion drivers ( <i>OD</i> ), defined as the Users with a <i>verified</i> Twitter account and at least a Tweet with more than 10,000 reactions from other Users. We notice that just five of the potential <i>OD</i> posted more than ten Tweets related to the topic. . . . .	10
2.2	Number of geo-tagged Tweets per country, considering the General dataset. We can suggest a strong interest in the topic from the countries related to palm oil production, as well as from the UK and the US. . . . .	11
3.1	Average sentiment of the <i>OD</i> . . . . .	16
4.1	Nodes of the co-mentions graph (largest connected component) with the highest degree. . . . .	24
4.2	Nodes of the co-hashtags graph (largest connected component) with the highest degree. . . . .	25
4.3	Nodes of the user-hashtag graph (largest connected component) with the highest degree. . . . .	26
4.4	Comparison of the three networks. $\ln N$ is the logarithm of the total number of nodes in the networks, while the fitting procedure is done on the largest connected component ( <i>GC</i> ) of each graph, the $p$ -value measures the significance of choice between candidate distributions, the degree cut-off is the lower bound chosen by minimizing the Kolmogorov-Smirnov distance $D$ , $\alpha$ is the scaling exponent and $\sigma$ is its standard error. . . . .	28
5.1	The fit results: $n$ is the number of bins considered (the number of scatter points in Fig. 5.1). $p$ is the significance value of the comparison between candidate distributions, $\tau_{min}$ is the lower bound of the scaling range [48], $D$ is the Kolmogorov-Smirnov distance, $\alpha$ is the scaling parameter, $\sigma$ is the standard error. <sup>1</sup> . . . . .	33
6.1	The fit results: $n$ is the number of bins considered (the number of scatter points in Fig. 6.2). $p$ is the significance value of the comparison between candidate distributions, $CS_{min}(h)$ is the lower bound of the scaling range [48], $D$ is the Kolmogorov-Smirnov distance, $\alpha$ is the scaling parameter, $\sigma$ is the standard error. . . . .	39
6.2	The fit results, where $n$ is the number of bins considered (the number of scatter points in Fig. 6.3). $p$ is the significance value of the comparison between candidate distributions, $CS_{min}(s)$ is the lower bound of the scaling range [48], $D$ is the Kolmogorov-Smirnov distance, $\alpha$ is the scaling parameter, $\sigma$ is the standard error. . . . .	41

A.1	Results of the fit for the <b>Coconut</b> dataset: $n$ is the number of bins considered (the number of scatter points in Fig. A.6). $p$ is the significance value of the comparison between candidate distributions, $\tau_{min}$ is the lower bound of the scaling range [48], $D$ is the Kolmogorov-Smirnov distance, $\alpha$ is the scaling parameter, $\sigma$ is the standard error. . . . .	54
A.2	Results of the fit for the <b>Olive</b> dataset: $n$ is the number of bins considered (the number of scatter points in Fig. A.6), $p$ is the significance value of the comparison between candidate distributions, $\tau_{min}$ is the lower bound of the scaling range [48], $D$ is the Kolmogorov-Smirnov distance, $\alpha$ is the scaling parameter, $\sigma$ is the standard error. . . . .	55
A.3	Results of the fit for the <b>Coconut</b> dataset: $n$ is the number of bins considered (the number of scatter points in Fig.A.8), $p$ is the significance value of the comparison between candidate distributions, $CS_{min}(h)$ is the lower bound of the scaling range [48], $D$ is the Kolmogorov-Smirnov distance, $\alpha$ is the scaling parameter, $\sigma$ is the standard error. . . . .	56
A.4	Results of the fit for the <b>Olive</b> dataset: $n$ is the number of bins considered (the number of scatter points in Fig.A.9), $p$ is the significance value of the comparison between candidate distributions, $CS_{min}(h)$ is the lower bound of the scaling range [48], $D$ is the Kolmogorov-Smirnov distance, $\alpha$ is the scaling parameter, $\sigma$ is the standard error. . . . .	57

# List of Abbreviations

All-languages	Dataset containing the keyword <i>palm oil</i> in all its translations
Coconut	Dataset with Tweets containing the keyword <i>coconut oil</i>
CS	Cascade Size
EDA	Exploratory Data Analysis
GC	Giant Component
GP	General Public
IET	InterEvent Time
KS	Kolmogorov-Smirnov
LSA	Latent Semantic Analysis
OD	Opinion Drivers
Olive	Dataset with Tweets containing the keyword <i>olive oil</i>
Palm	Dataset with Tweets containing the keyword <i>palm oil</i>
SVD	Singular Value Decomposition
Tf-idf	Term frequency-inverse document frequency

*Per chi viaggia in direzione ostinata e contraria...*



# Introduction

Collective phenomena permeate the reality surrounding us [1]: in nature, flocks of birds produce recognizable patterns [2] and fireflies flash following a collective rhythm [3]. Complexity science seeks to explain the emergence of collective phenomena [4], using quantitative tools derived from statistical mechanics and beyond. The need to study diametrically opposed phenomena has led complexity science to be characterized by strong interdisciplinarity. Over the years, the possibility of modeling complex systems by studying their collective behavior has begun to fascinate physicists, mathematicians, biologists, economists, and many others [5], forming a modern and ever-evolving field of research.

The advent of this new way of studying reality has also opened up a unique perspective in sociology. It is now possible to study human interactions [6], comparing them to other systems present in nature, outlining models [7] that universally describe the emergence of these collective phenomena. A further urge for a fusion of the two fields of research, such as complexity science and sociology, has been pushed recently by social media's sudden emergence and development. It is a tool that allows individuals to interact through new mechanisms, sharing news, opinions, and multimedia content. Complexity scientists began to study social media behavior to be able to quantitatively understand the emergence of collective phenomena in the opinion dynamics [8, 9], and how information (and misinformation) propagates on the social network structure.

Viral phenomena [10] emerge on social media every day, massively bringing new information to our screens that sooner or later will be forgotten. Many social and political campaigns (carried out by activists and non-governmental organizations) or marketing campaigns (carried out by companies) have found fertile ground in social media. It is easier to attract the attention of a large number of people to one's issue, creating intense interest and debate about it. However, due to the large amount of information we perceive through social media, the attention devoted to each topic we are exposed to is ephemeral, temporary [11].

In order to characterize these behaviors, we considered a specific case study. A topic that has generated sudden media interest and heated debates over the past few decades is the production and consumption of palm oil. Palm oil is the most widely used vegetable oil in the world (more than 35% of world production) [12]. However, it is controversial due to its production's strong environmental impact [13] generated. In particular, the growth in palm oil production, due to a substantial increase in global demand, has caused a devastating impact on the ecosystems of producing countries, leading to deforestation and loss of biodiversity [14], as well as conflicts over land ownership [15]. For decades, it has been an intensively debated topic within specialist circles. The general public has witnessed a viral phenomenon concerning palm oil. In the second half of the 2010s, many sociopolitical campaigns emerged to support more sustainable palm oil production. These campaigns were disseminated through traditional and modern media, such as newspapers and online social media. The convenience of extracting data from digital media allows us to focus on studying it. Mainly, we analyzed a globally popular social network, Twitter.

In this thesis, we characterize the emergence of viral events and, more generally, user attitudes concerning the specific case study of palm oil, using tools and techniques derived from different branches of complexity science. The work is organized as follows: Chapter 1 is a brief description of the Twitter social network and the *Data Collection* process; Chapter 2 conducts an exploration of the dataset, analyzing the main features and the overall trends; Chapter 3 provides an overview of the dataset from the *Natural Language Processing* point of view, implementing the *Sentiment Analysis* and the *Topic Detection* approaches; Chapter 4 analyzes the underlying social structures from a *Network Science* perspective; Chapters 5-6 employ the statistics of two quantities, the *interevent time* and the *cascade size*, to mathematically characterize the importance of an event; finally, Chapter 7 provides a general pipeline to compare and characterize the emergence of viral events around a topic, with a comparison among vegetable oils.

Leveraging strategies from different research fields allowed us to obtain an overall view of the case study and to be able to characterize the dynamics of social media interactions on sustainability issues.

## Chapter 1

# Twitter dataset: an overview

This chapter aims to describe the basic properties of Twitter, the social network we considered for this analysis, and briefly explain the Data Collection process output. Indeed, it is the fundamental starting point for Social Network Analysis.

### 1.1 Basics of Twitter

*Twitter* is a social network founded on the 26th of March, 2006, by Jack Dorsey, Noah Glass, Biz Stone, and Evan Williams. It works as a microblogging service where Users can post, like, reply, and reTweet (i.e., publish on their profile a Tweet created by someone else). Every User can follow another User, i.e., receive updates from a given User on their Twitter homepage. The action of following is not necessarily reciprocal. The followed User can choose to follow back or not. From the User's profile page, it is possible to access four main sections: Tweets, Tweets & Replies, Media, and Likes. Therefore, there is a complete timeline of the User's interactions with the Twitter social network, excluding private messages among Users (these are called Direct Messages). There are specific rules concerning the structure of a Tweet: it can contain text, images, hyperlinks, GIFs, and, more importantly for our concerns, hashtags, and mentions. In the following, a definition of these two metadata structures is given.

**Definition.** Hashtags are pieces of text starting with #, without spaces between the words, that can be used as an identification of a particular movement or topic.

It is possible to access the most recent or pertinent Tweets containing a specific hashtag from the search engine on Twitter. Moreover, a Twitter section dedicated to Trending Topics collects the most popular hashtags at that moment.

**Definition.** Mentions are hyperlinks starting with @ and followed by the username of a Twitter community member, sending back to the mentioned profile.

Mentions are structurally similar to hashtags. The only differences are found in the starting symbol (@ instead of #) and the restriction in the text content to the usernames of other social network members. Mentions are a way to invoke publicly another User of the community. We are particularly interested in these two metadata structures because it is possible to perform a quantitative analysis of the hashtags and mentions' use on social media.

### 1.2 Collecting data on Twitter

The best way to obtain a large dataset from a social network is to collect data from an API (Application Programming Interface). In this context, after requesting an

Academic Permission<sup>1</sup>, we invoked the Twitter API v2, from which it is possible to access a limited number of Tweets per month<sup>2</sup>.

As a general rule, the data is stored in JSON files<sup>3</sup> containing the following information (when applicable):

- *created\_at*: creation date of the Tweet;
- *entities*: set of all the special characteristics of a Tweet (if there are mentioned Users, hashtags, etc.);
- *lang*: the language of the Tweet;
- *id*: unique numerical ID identifying the Tweet;
- *source*: app or browser from which the Tweet is created;
- *text*: text corpus of the Tweet;
- *public\_metrics*: engagement metrics (number of reTweets, mentions, replies, likes, quotes);
- *author\_id*: unique numerical ID of the Tweet's author;
- *conversation\_id*: numerical ID of the original Tweet (i.e., the *root* of replies' and mentions' *tree*);
- *referenced\_Tweets*: list of referenced Tweets;
- *attachments*: list of the Tweet's attachments (media files, polls, etc.);
- *in\_reply\_to\_user\_id*: User's numerical ID at which this User is replying;
- *geo*: geolocalization details. Geolocalization is optional.

Furthermore, two additional types of JSON files are obtained from the Data Collection process. They present a deeper description through metadata of the Users involved in the Tweet (author, mentioned Users) and the geographical coordinates. These files help complete an overall description of the Tweets.

We built three different datasets by requesting to the API all the Tweets containing specific keywords created from Twitter's foundation to the 31st of December, 2021. The three datasets, further referred to as **Palm**, **Olive**, **Coconut** datasets, are obtained by requesting all the Tweets containing respectively the keywords *palm oil*, *olive oil*, and *coconut oil*. Due to the time limitations imposed by the Academic Permission (i.e., it is possible to request a limited number of Tweets per month), we decided to limit our search queries to the previously mentioned instead of trying alternative wordings. In this way, the Tweets containing, for instance, just *palm* or *oil* are not considered. The same reasoning is applied to bad-spelled words or the absence of white space (*palmoil*). Table 1.1 shows the number of Tweets and Users composition of each dataset. The **Olive** dataset is the biggest among the three, followed by the **Coconut** and the **Palm**. We have also extracted the data about other oils for the preliminary analysis: *canola*, *peanut*, *soybean*, and *sunflower*. However, we decided

<sup>1</sup>More information about the Academic Permission requirements and advantages can be found here [developer.twitter.com/en/products/twitter-api/academic-research](https://developer.twitter.com/en/products/twitter-api/academic-research).

<sup>2</sup>The limit is set to 10 million Tweets per months, with a maximum streaming rate of 50 server requests every 15 minutes.

<sup>3</sup>Acronym of JavaScript Object Notation, it is a standard data format based on JavaScript.

Oil	Number of Tweets	Number of Users
<b>Palm</b>	3,771,073	1,210,537
<b>Olive</b>	6,677,619	2,821,484
<b>Coconut</b>	5,269,203	2,111,945

TABLE 1.1: Number of Tweets and Users for each dataset obtained from the keyword-based Data Collection process.

to consider just the three oils mentioned above, as the number of Tweets regarding the other oils is negligible by two orders of magnitude. There is another reason to investigate the opinions around the three oils further. They are considered the most dangerous oils for the environment because of the number of species threatened by their production [16]. Furthermore, another dataset is created: the previous datasets are only composed of Tweets in English due to the keyword-based approach in the Data Collection process. Therefore, we decided to look for the translations of *palm oil* (ex. *olio di palma* in Italian) in all the other languages: this dataset is composed of 6,512,541 Tweets and 2,396,991 Users. The latter dataset (further referred to as **All-languages** dataset) is useful for analyzing the geographical distribution of the Tweets regarding palm oil (Section 2.4).

As a general rule, we decided to focus on the **Palm** dataset, recalling the other datasets' results only when a comparison is required. This is the case of Section 3.2, where we perform the *Topic Modeling* on the **Palm, Olive, Coconut** datasets. Similar reasoning is valid for Chapter 7, where we compare the virality phase diagrams of the three datasets. Moreover, Appendix A contains further details about the **Olive** and **Coconut** datasets.

## Chapter 2

# Exploratory Data Analysis

In order to unveil the main features of the dataset, the first step is to use the so-called *Exploratory Data Analysis (EDA)*. It is a general term that defines the investigation of the dataset through visualization techniques. Our interest is focused on capturing the dataset's size and time evolution. We focus on disclosing eventual viral events by looking at the time evolution in the number of Tweets and detecting the topics of interest by analyzing the hashtags and keywords. Other helpful information is retrieved by looking at the geographical distribution of the Users, which can give important insights into where the public campaigns are created and carried on with more robust engagement.

In this chapter, first, we analyze the time evolution of the number of Tweets to detect the emergence of viral events (Section 2.1). Next, we explore the most common hashtags and keywords to highlight the most interesting discussion topics (Section 2.2). Then, we identify the most popular Users and the ones that drove public opinion throughout social campaigns (Section 2.3). Finally, we analyze the geographical distribution of the Tweets to establish a correlation between the countries that are producing palm oil and the countries that are tweeting the most (section 2.4).

## 2.1 Time evolution of the number of Tweets

Exploring the dataset, we tracked the activity of 1,210,537 Users. As shown in Fig. 2.1, there has been a substantial increment in the interest in palm oil from the early stages, with a prominent peak in 2018. To further estimate the significance of this event, in Appendix A.1, we analyze the relative growth of this dataset, comparing it to the general growth of the social network over the years. Moreover, focusing the attention on this peak of interest, we highlighted a particular month, November 2018, where we hypothesize that a viral event happened. As shown in Fig. 2.2, there is a strong and sudden growth (from 1,000 to 31,562 Tweets) between the 8th and the 9th of November. The viral event that brought attention to palm oil was led by the Twitter account of the company Iceland Foods, which published a Christmas TV spot against palm oil, which the UK TV channels banned.

In Fig. 2.3, we can see the Tweet post that received more than 14,700 likes and 8,000 reTweets, indicating the growing interest in the topic. This campaign, supported by Greenpeace, generated more awareness about the sustainability of products such as palm oil, which we further investigate in the following chapters.

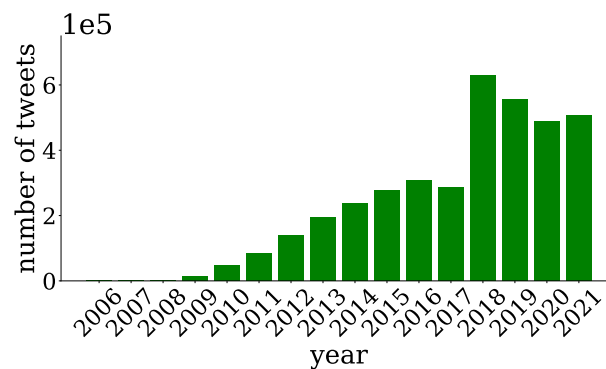


FIGURE 2.1: Histogram representation of the number of Tweets containing the words *palm oil* per year. A peak of interest was detected in 2018.

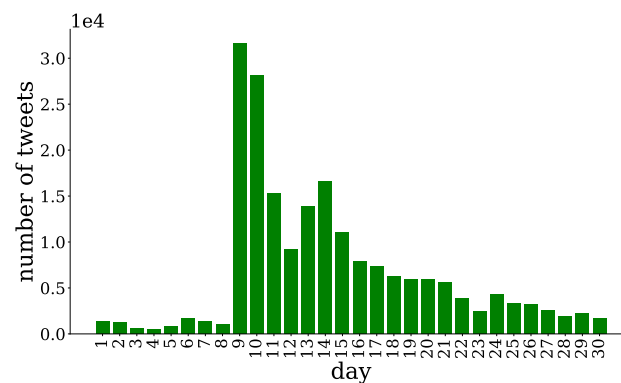


FIGURE 2.2: Histogram representation of the number of Tweets containing the words *palm oil*, restricted to November 2018. The interest increased on the 9th of November when the Tweet in Fig. 2.3 was published.



FIGURE 2.3: Screenshot of [Iceland foods' Tweet](#) about the banned TV spot. The Tweet received more than 14,700 likes and 8,000 reTweets.

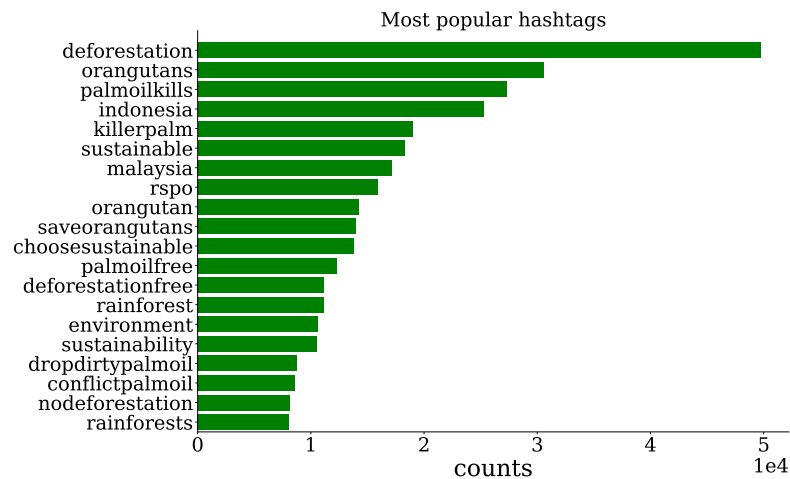


FIGURE 2.4: Histogram representation of the most popular hashtags in English. *palmoil*, *palm*, *oil* were removed, sustainability-related hashtags are predominant.

## 2.2 Keywords and Hashtags

From a quantitative point of view, analyzing the dataset's keywords and the most used hashtags is quite significant. Usually, a petition or a public campaign is supported by a specific hashtag. For example, in Fig. 2.3, the starters of this campaign tried to propose the use of the hashtag *#NoPalmOilChristmas*. We could assume that it is not a *successful* hashtag (i.e., it is not present in the list of the most common hashtags) because it is peculiar to this Christmas TV spot, while the most used hashtags are quite more general. We extracted the hashtags from each Tweet, following the definition given in Section 1.1. In this section, we neglect the hashtags *palmoil*, *palm*, *oil* to analyze other aspects not directly related to the words used to create the dataset. As shown in Fig. 2.4, the most popular hashtags are related to palm oil sustainability issues (*deforestation*, *orangutans*, *Indonesia*, *sustainable*), supporting the idea that most of the public campaigns against palm oil were carried on to stop the deforestation and preserve the rainforest microclimate. The orangutans were the leading actors of the Iceland foods' TV spot (and the inhabitants of rainforests); therefore, the substantial use of that hashtag is expected and well-motivated.

The *wordcloud*<sup>1</sup> package is another way to inspect graphically the most common keywords. It is an NLP tool that performs text *pre-processing* (i.e., it removes the so-called stop-words, then it extracts the single words from the corpus, checking if there are errors in the spelling) and evaluates the frequency of each keyword. Then, the graphical representation of each word is sized proportionally to its frequency in the dataset. Fig. 2.5 represents the *wordcloud* for our dataset. As previously highlighted in the hashtags' analysis, it is possible to notice a significant prevalence of keywords related to sustainability issues.

## 2.3 Popular Users

Another interesting perspective for the dataset characterization is the identification of the potential *opinion drivers* (OD), that we define as follows:

<sup>1</sup>[http://amueller.github.io/word\\_cloud/](http://amueller.github.io/word_cloud/)



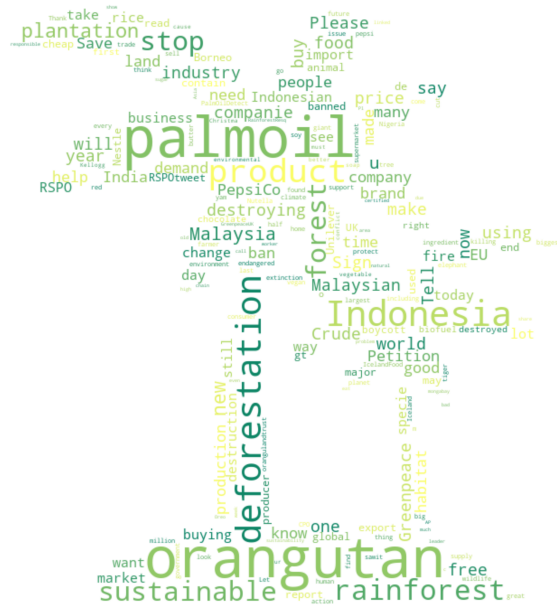


FIGURE 2.5: Wordcloud of the most common keywords in the corpus. Keywords related to sustainability are widely used.

- they shall own a *verified* account (i.e., Twitter assures that the account is authentic and that it respects some guidelines<sup>2</sup>);
- the sum of the number of retweets, replies, likes, and quotes is bigger than 10,000 for at least one of their Tweets in the dataset.

We identified 36 Users with these characteristics, shown in Table 2.1. Notice that only five of these Users created more than ten Tweets related to the palm oil topic. This result can be explained by the fact that some opinion drivers mentioned palm oil just in a few Tweets. Probably, they are not running a public campaign, or they are not involved in some organizations, even if they are respecting the constraints we fixed. Therefore, one way to further restrict our analysis is to add another constraint. Here, we considered only the Users with more than 10 Tweets related to palm oil. In this way, the *OD* are restricted to the following Users:

- **Iceland Foods**: the British supermarket chain that proposed the TV spot;
- **Greenpeace UK**: the British branch of the international NGO that relates to environmental issues;
- **Andreas Harsono**: an Indonesian human rights activist;
- **AJ+**: the Al-Jazeera's social media publisher;
- **Business Insider**: an American financial news website.

We were able to identify and constrain the leading actors of the debate, noticing that only a few verified accounts were deeply involved in the discussion. In the following chapters, we will further analyze the similarities and differences between the General Public and this small group of Opinion Drivers. Was this public campaign effective? Was the General Public deeply involved in the discussion?

<sup>2</sup>More information about the Twitter guidelines for the verified accounts can be found here <https://help.twitter.com/en/managing-your-account/about-twitter-verified-accounts>

Name	Self-reported location	Number of Tweets
Iceland Foods	In a UK town near you	776
Greenpeace UK	UK	517
Andreas Harsono	Indonesia	130
AJ+	-	53
Business Insider	New York, NY	49
Carl Franzen	Pittsburgh, PA	6
Major Gaurav Arya (Retd)	India	6
Sema	Texas	5
IcelandFoodsIreland	Ireland	4
Peter McGuire	Dublin	4
Elizabeth Cotignola	Chicago/Montreal	3
merry jerry	Nashville, TN	3
Anshul Saxena	India	2
Bret Von Dehl	New Jersey, USA	2
Saket Gokhale	Goa/New Delhi/Mumbai	2
jordan	Washington, D.C.	2
Conor McDonnell	Los Angeles	2
Gabi de Ferrer	London, England	2
Imran Khan	Pakistan	1
Michael Bonfiglio	-	1
Phil Nolan	Brooklyn (He/Him)	1
Rani Timekey Baker	probably Portland, Oregon	1
bleep	-	1
Deterministic Optimism	-	1
Ramon on Zoom 1/29 7p est	Lorain, OH	1
James Coleman	South San Francisco, CA	1
Nicole Schuman, M.A.	Astoria, N.Y.	1
David Klion	Brooklyn	1
Joel Birch	Sunshine Coast	1
Delphine Rivet	-	1
Drew Holden	Washington, DC	1
Omar Sakr	Sydney, New South Wales	1
"teen suicide" the band	Around	1
demy	-	1
Rania Khalek	Lebanon	1
Nathan Bernard	Maine, USA	1

TABLE 2.1: Potential opinion drivers (*OD*), defined as the Users with a *verified* Twitter account and at least a Tweet with more than 10,000 reactions from other Users. We notice that just five of the potential *OD* posted more than ten Tweets related to the topic.

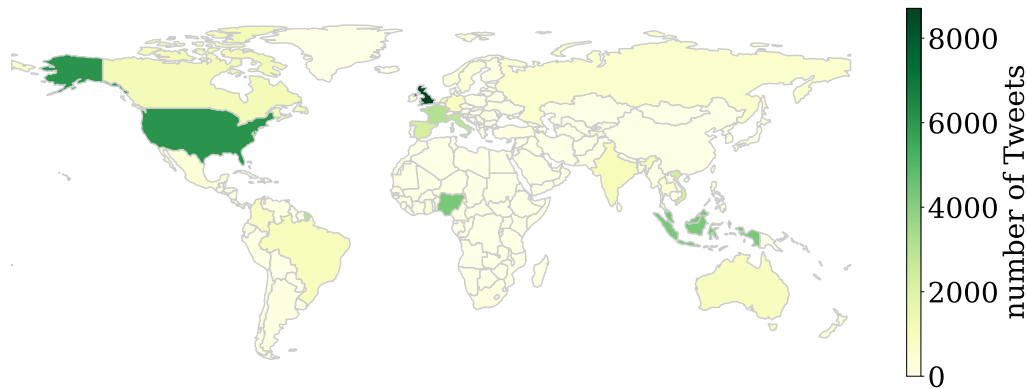


FIGURE 2.6: Geographical distribution of Tweets from the **All-languages** dataset. The UK, the US, Malaysia, Nigeria, and Indonesia are the countries that are Tweeting the most about the palm oil topic.

Country	Number of Tweets
United Kingdom	8,706
United States of America	6,099
Malaysia	4,372
Nigeria	4,334
Indonesia	4,301

TABLE 2.2: Number of geo-tagged Tweets per country, considering the General dataset. We can suggest a strong interest in the topic from the countries related to palm oil production, as well as from the UK and the US.

## 2.4 Geographical distribution

Using the **All-languages** dataset (not restricted to the Tweets in English), it is possible to highlight which countries are the most involved in the palm oil topic. From the Twitter API, it is possible to extract a database with all the tagged locations, 54,522 (only 0.84% of the dataset). The percentage of Tweets with active localization is generally low. However, we assume that the other Tweets are distributed with the same proportions in the different locations, considering the reachability limits of the social network (it is banned from some countries, the localization is not available in others, etc.).

In Fig. 2.6, the geographical distribution is shown through a heat map. As expected, the most active countries are the UK and the US. At the same time, a peculiar high activity is highlighted in Malaysia, Nigeria, and Indonesia (Table 2.2 displays the most active countries and the respective number of Tweets.). According to [Forbes](#), Indonesia produced 58% of the total amount of palm oil in 2019, followed by Malaysia (26%), Thailand (4%), Colombia (2%) and Nigeria (1%). We can conclude that the countries more directly involved in the sustainability issues of palm oil are also the most active on social media, together with the countries where the public campaigns were carried out (the UK and the US)<sup>3</sup>. In the following, we focus on the English datasets, as we showed that this is the primary source of Tweets.

<sup>3</sup>Further investigations could lead in the direction of understanding if these results are influenced by the total volume of Tweets for each country.

## Chapter 3

# NLP tools: Sentiment Analysis and Topic Modeling

### 3.1 Sentiment Analysis

*Sentiment Analysis* (also known as Opinion Mining) is a Machine Learning text classification task that assigns a degree of polarity to a text. It spans from *Positive* to *Negative* labels, detected through the examination of the text corpus, the structure of the sentences, and the use of punctuation signs. For example, a sentence like “I love palm oil” will probably be detected as *Positive*, while the sentence “I hate palm oil” will be detected as *Negative*. We employed a multilingual pre-trained model, described in Section 3.1.1, to assign to each Tweet the labels *Positive*, *Neutral*, and *Negative*. A score was assigned to each label: it represents the probability that the considered Tweet has a given sentiment. For further analysis, we consider each Tweet as labeled by the one with the highest score among the three.

In the following, we briefly describe the model, highlighting the general results of the classification task (Section 3.1.1). To investigate the sentiment heterogeneity, we also study the time evolution of the sentiment in Section 3.1.2. This analysis can lead to important insights into how viral events are related to sentiment. As reported in [17], we expect a correlation between negative sentiments and attention on a specific topic. Then, we analyze some interesting results followed from sentiment analysis of the so-called *opinion drivers* (OD) (Section 3.1.3). They are the most active Users that influence public opinion by producing Tweets with massive engagement (many mentions, reTweets, and replies). Furthermore, in Section 3.1.4, we study the empirical relationship between the detected sentiment and the Users’ opinions.

#### 3.1.1 Model description and general results

*Twitter-xlm-roberta-base-sentiment*<sup>1</sup> is a multilingual model pre-trained on  $\sim 198M$  Tweets in 30 languages and fine-tuned for sentiment analysis<sup>2</sup>. This pre-trained model best suits our aim based on the evaluation benchmark performed in [21]. It is trained on a corpus of Tweets, while XLM-RoBERTa [22] (the more general model) is trained on a larger, less specific corpus. Furthermore, Barbieri et al. [21] showed that a more Twitter-specific model is better for capturing the sentiment in the case of

<sup>1</sup>[huggingface.co/cardiffnlp/twitter-xlm-roberta-base-sentiment](https://huggingface.co/cardiffnlp/twitter-xlm-roberta-base-sentiment)

<sup>2</sup>It is possible to find the same model fine-tuned for other specific tasks, for example irony detection [18], emotion recognition [19] or offensive language identification [20]. They can be useful to analyze further different aspects of the debate, but we decided to restrict our analysis to a three-label system, the sentiment analysis. We also decided only to consider the English dataset for the reasons presented in Section 2.4.

Twitter datasets. It is due to the particular Tweet structure: constraint in the maximum length, presence of hashtags and hyperlinks, and more primitive sentences.

The first step needed to obtain a suitable text corpus is the *preprocessing*. For this specific model, the only request was to substitute a username mentioned with @ with just @user and a hyperlink starting with http with just http. In this way, the model recognizes them as negligible information. Furthermore, we performed the *Tokenization*: it transforms a text corpus into a set of tokens. Depending on the type of Tokenization, we are going to obtain either tokens made of words (ex. *shorter*), letters (ex. *s-h-o-r-t-e-r*), or sub-word tokens (ex. *short-er*). The Tokenizer used for this method is of the first type. Then, running the classification task on each Tweet of the English dataset, we obtain as output

$$scores(t) = \{s_t(i)\}_{i=0,1,2}, \quad (3.1)$$

where  $s_t(i)$  are respectively the scores assigned to the three labels<sup>3</sup>. Hence, we attribute the label  $l$  to the Tweet  $t$ , such as

$$l(t) \in \operatorname{argmax}_{i \in \{0,1,2\}} scores(t), \quad (3.2)$$

i.e., we choose the label with the highest score. The **Palm dataset** classification produces the results shown in Fig. 3.1.

We notice that many Tweets are labeled as *Neutral*. We summarized the possible reasons that can lead to high neutrality into three principal ones:

- *Balance of sentiments*: the presence of *Negative* and *Positive* words in the same Tweet could lead to a balance of opposite stances. Therefore, the algorithm does not detect a particular attitude in these Tweets;
- *Tweet's brevity constraints*: due to the peculiar structure of the Tweets, it is more difficult to detect an attitude from a short text with mainly basic sentences than from a corpus made by other types of text;
- *Sentiment  $\neq$  Opinion*: it is crucial to focus on the difference between sentiment and opinions. It is possible to express a *Negative* opinion using a *Positive* attitude, for example, by expressing joy for the success of a public campaign against palm oil. The expression "Opinion Mining" can be misleading in this context because we are far from detecting opinions using the Sentiment Analysis task.

In the following, we analyze how sentiment changes over time, how it is related to a User's popularity, and finally, prove the missing relation between sentiment and opinions in our context.

### 3.1.2 Sentiment dynamics

As shown in Fig. 3.2, where the time evolution of labeled Tweets is plotted, the Tweets detected as *Positive* remain stable around 10-15 % of the total. In contrast, the percentage of *Negative* Tweets increases over time. It is interesting to appreciate that the highest percentage of *Negative* Tweets was detected in 2018 when the biggest viral event happened. This is evidence of the relationship between *Negative* sentiment and viral events, confirming the results obtained in [17]. This aspect is further investigated in Sections 5.3 and 6.2.

<sup>3</sup>0 is *Negative*, 1 is *Neutral*, 2 is *Positive*.

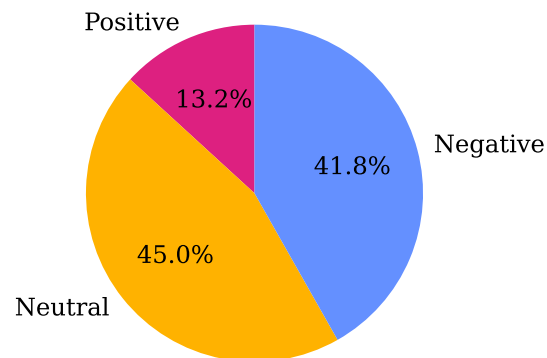


FIGURE 3.1: Pie chart representation of the text classification task results. A prevalence of *Neutral* labels is highlighted, followed by the *Negative* labels. Just a few Tweets are classified as *Positive*.

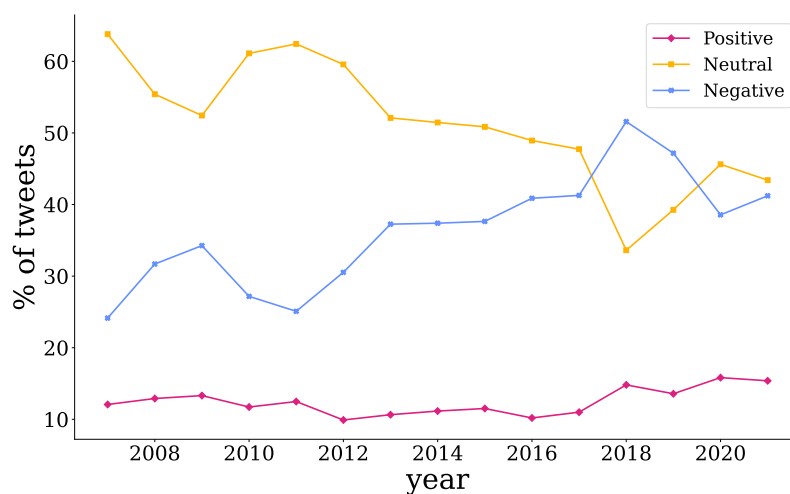


FIGURE 3.2: Time evolution of sentiment. The plot represents the percentage of labeled Tweets published in a given year. *Positive* sentiment remains low and stable through time, while *Negative* sentiment has a peak around 2018.

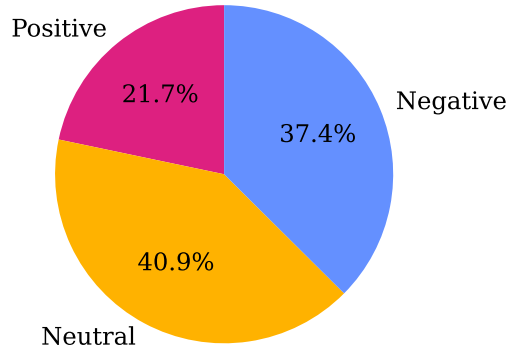


FIGURE 3.3: Pie chart representation of the text classification task results for the **OD**. The percentage of *Positive* labels is higher compared to the **GP** case (Fig. 3.1).

### 3.1.3 Popularity and sentiment

In this section, we are interested in analyzing the relationship between sentiment and popularity and how the *opinion drivers* differ from the rest of the Users. In detail, we extracted the statistics on the sentiment of the five Users identified as **OD**; the labels are distributed as shown in Fig. 3.3. It is possible to highlight a stronger presence of *Positive* attitudes among the **OD**, compared to the General Public (further referred as **GP**). To further investigate this phenomenon, we calculated the average sentiment for the **OD** (see Table 3.1). The average is computed as

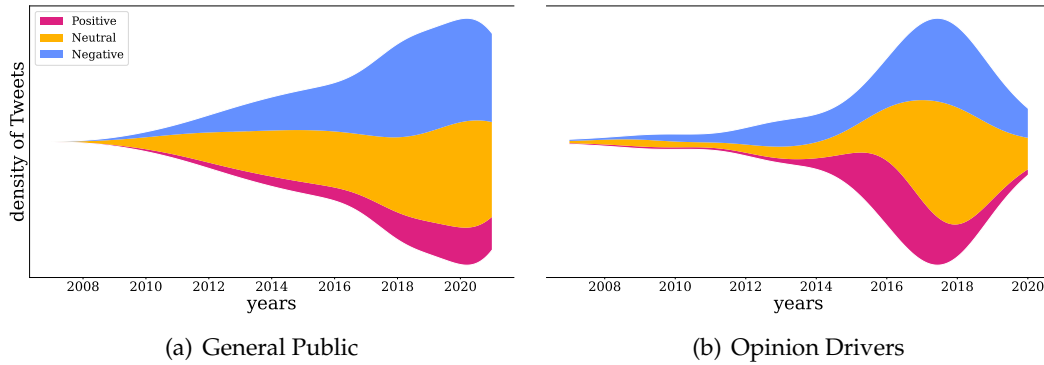
$$\langle l(u) \rangle = \frac{1}{N_u} \sum_{t \in u} l(t),$$

where  $u$  is a given User,  $N_u$  is the number of Tweets in the dataset created by the User  $u$ , and  $l(t)$  is the sentiment assigned to the Tweet  $t$ .

To complete the framework, it is also helpful to compare the time evolution of the **OD** sentiments with the **GP**. In Fig. 3.4, the Gaussian smoothing of the sentiment dynamics is plotted. The two plots differ significantly: after 2015, the **OD** *Positive* Tweets start to grow. The average sentiments do not detect this behavior. One reason for this higher *Positive* attitude could be identified in seemingly counterintuitive reasoning: we expect the **OD** to be the strongest ones trying to transmit anger and interest in their social campaign. Nonetheless, it is possible to state that the **OD** are not using the most engaging words, as they are less polarized than the **GP**. Considering the Iceland Food campaign (Fig. 2.3), the company did not need to use strong words such as *kill* because they are still maintaining a *Positive* attitude by using words such as *enjoy*. The anger and the strong polarization are left to the Users that comment on this post and spread the campaign among their followers.

User $u$	Number of Tweets	Average entiment $\langle l(u) \rangle$
IcelandFoods	776	<i>Neutral</i>
GreenpeaceUK	517	<i>Neutral</i>
andreasharsono	130	<i>Negative</i>
ajplus	53	<i>Negative</i>
BusinessInsider	49	<i>Negative</i>

TABLE 3.1: Average sentiment of the OD.

FIGURE 3.4: Comparison of sentiment dynamics between the GP and the OD. Proportionally, the OD show a stronger growth in the *Positive* sentiment than the GP.

### 3.1.4 Sentiment analysis and Opinion Mining: are they synonyms?

In this subsection, we present evidence of the discrepancy in the definitions of *sentiment* and *opinion*. *Ad absurdum*, let us state that sentiment and opinion are overlapped definitions. Then, we expect a one-to-one correspondence between the "extreme" labels (i.e., *Positive* and *Negative*) and the two opinions (i.e. being *for* or *against* palm oil). To see if this statement is correct, we used the classification obtained from Sentiment Analysis to divide the keywords into two subsets (the keywords related to *Positive*-labeled and *Negative*-labeled Tweets<sup>4</sup>).

As shown in Fig. 3.5, the two subsets present a strong overlap. The same keywords present a high frequency in both the *Positive* and *Negative* Tweets. Users are concerned about sustainability, deforestation, and preservation of the microclimate for the life of the orangutans, and they are using both *Positive* and *Negative* attitudes. The only visible difference is the absence of the keywords *palmoilkill* and *killerpalm* in the *Positive* Tweets. This result indicates that our model is working correctly, identifying *Negative* words, even without correlations between Sentiment and Opinions.

## 3.2 Topic Modeling

The so-called *Topic Modeling* is an important branch of Natural Language Processing, which aims at discerning and detecting different topics in a text corpus. Both supervised and unsupervised strategies can be implemented to achieve this aim. Here, we chose an unsupervised approach based on *Bag-of-Words* [23] (BOW). Specifically, BOW divides the text into multisets of words and assigns the frequency of each

<sup>4</sup>For the sake of this analysis, we neglect the *Neutral* labels.



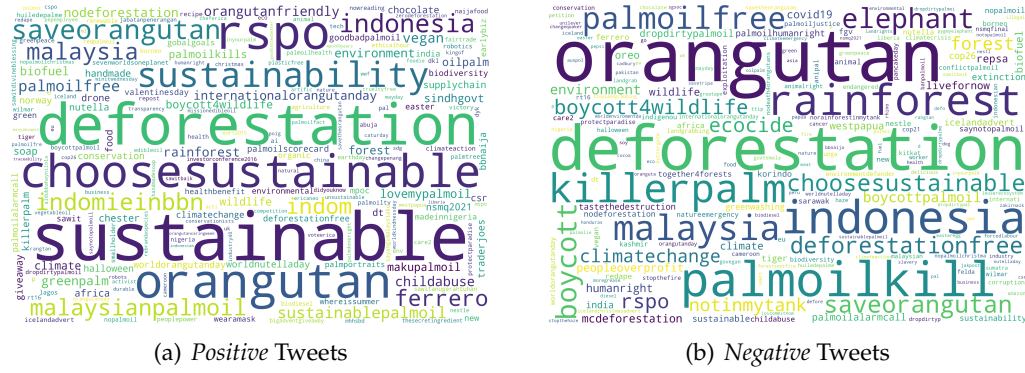


FIGURE 3.5: Wordcloud of the most common keywords for the two labeled subsets of the main corpus. In both cases, the keywords refer to sustainability issues.

word. Therefore, an occurrence matrix (also known as *document-term matrix*) is created as follows

$$M_{i,j} = w_{i,j}, \quad (3.3)$$

where  $w_{i,j}$  is the frequency of the word  $j$  in the document (here it is a Tweet)  $i$ . To account for the importance of a word for the tweet and also the entire corpus, instead of using the word counts, we assign weights computed with the so-called **Tf-idf** (*term frequency-inverse document frequency*) [24]. The product of two elements gives the Tf-idf weights as

$$\text{tf-idf}_{i,j} = \text{tf}_{i,j} \times \text{idf}_j, \quad (3.4)$$

where  $\text{tf}$  and  $\text{idf}$  are the *term frequency* and *inverse document frequency*, respectively. The first element,  $\text{tf}$ , is given by

$$\text{tf}_{i,j} = \frac{n_{i,j}}{|d_i|}, \quad (3.5)$$

where  $n_{i,j}$  is the number of occurrences of the word  $j$  in the Tweet  $i$ , and  $d_i$  is the length of the Tweet  $i$ . There are many distinct possibilities to calculate  $\text{idf}$ , and here we adopt the implementation of Scikit-learn [25], as follows

$$\text{idf}_j = \log_{10} \frac{|D|}{1 + |\{d \in D : j \in d\}|}, \quad (3.6)$$

where  $|D|$  is the size of the dataset (number of tweets), and  $|\{d \in D : j \in d\}|$  is the number of Tweets containing the word  $j$ . Eq. 3.6 is a penalty used to reduce the importance of the words present in many Tweets. Consequently, it focuses on the words that characterize this particular text. After assigning to each word a weight using Eq. 3.4, it is possible to create the Tf-idf matrix  $\mathbf{M}'$ . Then, using the *Latent Semantic Analysis* (also known as Truncated Singular Value Decomposition) [26], it is possible to project the high-dimensional occurrence matrix into two dimensions, still preserving the similarity among Tweets. The *Latent Semantic Analysis*, further referred to as **LSA**, is a technique that performs the Singular Value Decomposition (SVD) of the Tf-idf matrix  $\mathbf{M}'$  as

$$\mathbf{M}' = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad (3.7)$$

where  $\mathbf{U}$  and  $\mathbf{V}$  are orthogonal matrices and  $\mathbf{\Sigma}$  is diagonal. Notice that **LSA** can be calculated for both  $\mathbf{M}'$  or  $\mathbf{M}$ , but here we choose the Tf-idf matrix,  $\mathbf{M}'$ . As we are interested in projecting our data in a lower  $d$ -dimensional space ( $d = 2$  for the sake of visualization), the first  $d$  rows of  $\mathbf{U}$  and  $\mathbf{V}^\top$  represent the coordinates of the single words and the Tweets in the  $d$ -dimensional space.

More schematically, here we summarize the entire procedure:

1. *text preprocessing*: We converted all the words into *lowercase*, then the *Tokenization* is applied (the **RegexpTokenizer** from the Natural Language Toolkit [27]). From the list of tokens, *stopwords* were removed<sup>5</sup>. Finally, the *Lemmatization* is performed, using the **WordNet** lemmatization [28]<sup>6</sup>;
2. *tf-idf vectorization*: using the package *TfidfVectorizer* from scikit-learn [25], we set as parameters *max\_features* = 3,000, which limits the vocabulary (i.e., the dimensions of matrix  $\mathbf{U}$ ), to the top  $n$  words in terms of **tf-idf** frequency, as well as *min\_df* = 10 and *max\_df* = 0.7, which are cut-offs in the number of words. If a word appears in the corpus less than *min\_df* times, it is ignored from the vocabulary. Similarly, *max\_df* is a cut-off in the frequency of words. If a word occurs with a higher frequency than *max\_df*, it is ignored from the vocabulary. This method returns a sparse matrix (the Tf-idf matrix  $\mathbf{M}'$ ) with dimensions  $15,717,895 \times 3,000$ , where the size of the dataset is  $|D| = 15,717,895$ <sup>7</sup>;
3. *Latent Semantic Analysis*: using the package *TruncatedSVD* from scikit-learn [25] we set as parameters *n\_components* = 2, which is the dimensionality of the space in which we want to project the data, as well as *algorithm* = *randomized*, which determines the type of **SVD** solver used. We chose the default one, further described in [29]. This method returns the matrices  $\mathbf{U}$  and  $\mathbf{V}$ , with dimensions respectively  $3,000 \times 2$  and  $2 \times 15,717,895$  (recall that  $|D| = 15,717,895$  and *max\_features* = 3,000);
4. *results visualization*: matrix  $\mathbf{V}^\top$  gives the coordinates in a 2-dimensional space of all the Tweets. In the case of **LSA**, the two main directions (called *topic 1* and *topic 2* in the plot) represent the directions along which there are the highest variances. Every Tweet is labeled following the segmentation into **Palm**, **Coconut**, **Olive** datasets.

Here, we focused our analysis on the visualization of the tweets using **LSA**. Fig. 3.6 shows the results obtained. We can notice that the three oils are grouped differently. Specifically, all the Tweets related to palm oil are found close to the origin. In contrast, the Tweets about olive and coconut oil spread along the positive and negative directions on the y-axis. Using this analysis strategy, we could observe different topics around which the Users debate and estimate the topic similarities among datasets. We can state that the **Olive** and **Coconut** dataset presents a wider variety of topics. In contrast, the Tweets devoted to **Palm** are highly monothematic (i.e., the debates are predominantly related to sustainability and deforestation).

<sup>5</sup>The list of *stopwords* is provided by the Natural Language Toolkit [27].

<sup>6</sup>The **WordNet** lemmatization [28] consists of searching a given word in the **WordNet** database, and, subsequently, substituting the word with the correspondent lemma. If the word is not found in the dataset, the lemmatization process returns the original word.

<sup>7</sup>The dataset contains all the Tweets of **Palm**, **Coconut**, **Olive**.

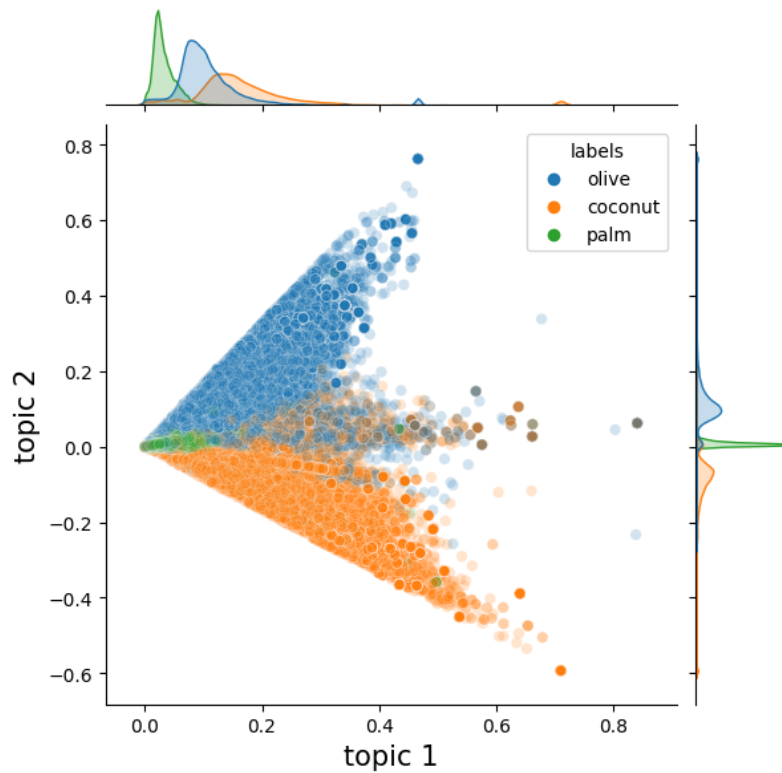


FIGURE 3.6: Results of **LSA** using the **Tf-idf** vectorization. The high-dimensional occurrence matrix  $\mathbf{M}$  is projected along two directions, representing the main distinctions into topics. Each point in the plot represents a Tweet of the dataset, labeled as the oil it refers to.

## Chapter 4

# Network science

### 4.1 Building blocks of Network Science

Complexity science is based on the need to characterize the collective phenomena arising from complex systems. Among the branches of complexity science, a fundamental approach is the one developed by *Network Science*, which aims to recognize patterns in the interactions among agents, studying the network structure underlying these real systems. Following the interdisciplinary spirit of complexity science, there are many applications of Network Science throughout the different fields of research [5] (e.g., physics, math, biology, computer science, and many others). Scholars have developed mathematical and statistical tools to analyze the networks and to underline a common universality *leitmotif* throughout the disciplines. Networks (or graphs) are composed of nodes and edges, where edges create connections between nodes. In the following, we describe some Network Science basics that can be useful for our particular context. Taking a graph  $G = (V, E)$ , where  $V$  is the set of nodes, and  $E$  is the set of edges<sup>1</sup>, we can define the following quantities [30]:

- a graph is *undirected* when there is not a precise direction of the edges, while it is *directed* if the edges are oriented;
- the *degree*  $k_i$  of the node  $i$  is the number of edges that are connecting the node to other nodes (we distinguish between *in-degree* and *out-degree* in the case of directed graphs, which are respectively the number of edges ongoing and outgoing from a given node);
- the *degree distribution* provides the probability that a randomly chosen node has a given degree. In many real networks, it is characterized by a power-law distribution, as we are going to explain further in the next section;
- a *bipartite network* is a graph whose nodes can be divided into two disjoint subsets, and the edges connect nodes of different subsets uniquely;
- a *weighted network* is a graph whose edges are weighted, while an *unweighted network* has all the edges with the same weight ( $w_{ij} = 1 \quad \forall (ij) \in E$ );
- the *adjacency matrix* is a network representation in which each matrix element is given by

$$A_{ij} = \begin{cases} w_{ij} & \text{if } (ij) \in E \\ 0 & \text{if } (ij) \notin E \end{cases}$$

where  $w_{ij}$  is the weight of the edge  $(i, j)$ ;

---

<sup>1</sup>A set of edges of an undirected graph is defined as  $E \subseteq \{\{i, j\} \mid i, j \in V \text{ and } i \neq j\}$ . It differs from the definition in the case of a directed graph; in the former case, it is an unordered set, while in the latter, it is ordered.

- a *connected component* is a subset of connected nodes (i.e., starting from each node in the subset, it is possible to reach with a path all the other nodes in the same subset). In the case of directed graphs, we distinguish between a *strongly* and *weakly connected* graph. The first is defined as a graph where it is possible to reach any node from every starting point, considering the directionality of the edges. The latter definition is similar, but every node is reachable from each starting point only in the undirected version of the same graph.

Since the well-known Millgram experiment [31], Network Science has been applied to social systems and, after the rise of the Internet, to social media. It is possible to identify different entities that can concur for the role of nodes and edges, in particular, we extracted three types of networks from our Twitter dataset, characterized as follows:

- *co-mentions* network: nodes represent Users, and the edges are created between Users mentioned in the same Tweet. The edge weights are proportional to the frequency of the co-mentions in the dataset. The idea behind the construction of this network is to bind together those Users who should, in some ways, discuss together around a topic;
- *co-hashtags* network: nodes represent hashtags, and the edges are created between hashtags used in the same Tweet. The edge weights are proportional to the frequency of the co-occurrence of hashtags. We constructed this network to highlight the relationships among the hashtags;
- *user-hashtag* network: nodes represent Users and hashtags, and the edges are created between Users and the hashtags used in their Tweets. The edge weights are proportional to the frequency of each event. This network underlines the relationships between Users and the topics they discuss.

It is possible to consider two different perspectives: first, the *static* perspective, where we study the time-independent graphs, taking into account the nodes and the edges created through time as immutable. Afterward, we study the *temporal* version of these graphs. In the latter, we are more concerned about time evolution, the generative processes, and how the burstiness of human activities can influence the network structure and its main features.

## 4.2 Community detection

The type of real networks we are interested in often present some organization in substructures, the so-called *communities*<sup>2</sup>. *Community detection* is a task that is actively in development, and even the definition of what a community is can vary. As Fortunato et al. state [32], there are two main approaches for defining a community: the *classical* approach sees the communities as denser subgraphs well separated from each other; therefore, it is more focused on counting the internal and external edges of a community, while the *modern* approach is more focused on considering the probabilities that a node has of sharing edges with the nodes in a community; this perspective assumes that the probability of being connected to a node within the same community is higher than the probability of being connected to an external node. In the latter approach, the procedure aims to find a preferential linking

<sup>2</sup>In other contexts, they are called *clusters* or *modules*.

pattern [32]. Therefore, it is necessary to introduce a generative model<sup>3</sup> to compute the probabilities of having edges between nodes. The most used generative model for community detection is the Stochastic Block Model [34], which is a widely used algorithm that divides a network into subgroups called *blocks*. Take  $N$  nodes and  $M$  blocks: the probability of having an edge between the nodes  $i$  and  $j$

$$P(\langle ij \rangle \in E) = p_{g_i, g_j}$$

depends on the  $g_i = 1, \dots, M$  which is the group membership of the node  $i$ . The stochastic block matrix is a  $M \times M$  matrix made by the elements  $p_{g_i, g_j}$ , and it is possible to extract further insights into the network's community structure. However, it is only one of the many possible approaches to tackle this challenging task.

We decided to stick to a more classical approach, using an algorithm based on modularity optimization. For the sake of further analysis, we define two quantities:

- *Modularity* [35]

$$Q = \frac{1}{2L} \sum_{ij} \left( A_{ij} - \frac{k_i k_j}{2L} \right) \delta(C_i, C_j)$$

where we sum over all the pairs of vertices  $i$  and  $j$  the difference between the actual adjacency matrix elements  $A_{ij}$  and the degree-preserving null model's expected number of edges between nodes  $i$  and  $j$ , which is  $p_{ij} = \frac{k_i k_j}{2L}$ <sup>4</sup>. Modularity is a measure of how much a network is well-divided in communities with respect to its randomized version and  $Q \in [-1/2, 1]$ . The higher the modularity, the better the partition,  $Q = 0$  means that there is not any partition, and  $Q < 0$  means that every node is in a different partition;

- *Nestedness* is given by the maximum eigenvalue of the adjacency matrix of the network  $\lambda_{max}$  [36]. It is a measure of the nested interaction structure and the self-organization patterns.

A widely-used community detection algorithm is the so-called *Louvain algorithm* [37]. It is a greedy optimization based on modularity, described in detail in the following.

1. Assign a different community to each node;
2. Calculate the modularity change<sup>5</sup> by moving the community of node  $i$  in each of the communities of its neighbors  $j$ ;
3. Move the node  $i$  in the community, which gives the largest gain. Otherwise, if the change in modularity is smaller or equal to zero, do not move the node  $i$ ;

<sup>3</sup>A generative model describes how a particular type of network is generated. Another example of a generative model is the Barabasi-Albert model [33], which produces a network using preferential attachment mechanisms.

<sup>4</sup> $C_i$  is the community of node  $i$ ,  $L$  is the total number of edges in the graph.  $A_{ij}$  is the adjacency matrix element that gives information about the connection between nodes  $i$  and  $j$ , while  $k_i$  is the degree of node  $i$ .

<sup>5</sup>The change in modularity determined by moving the node  $i$  in the community of node  $j$  is given by [37]

$$\Delta Q = \left[ \frac{\Sigma_{in} + k_{i,in}}{2L} - \left( \frac{\Sigma_{tot} + k_i}{2L} \right)^2 \right] - \left[ \frac{\Sigma_{in}}{2L} - \left( \frac{\Sigma_{tot}}{2L} \right)^2 - \left( \frac{k_i}{2L} \right)^2 \right],$$

where  $\Sigma_{in}$  is the sum of the edge weights internal to the community of node  $j$ ,  $\Sigma_{tot}$  is the weighted sum of the edges between incoming in the community of  $j$ .  $k_i$  is the degree of node  $i$ , while  $k_{i,in}$  is the number of edges from node  $i$  to nodes of the same community of  $j$  (included  $j$ ).



4. Repeat 2 and 3 until there is no gain in modularity anymore;
5. Create the induced graph (i.e., the graph made by the communities as nodes and edges within communities as self-loops);
6. Repeat the optimization on the induced graph until there is no gain in modularity.

Modularity optimization is an NP-hard task, but this algorithm is a widely-used approximation, even if modularity has some resolution limitations. Furthermore, it is strongly dependent on the size (number of nodes and edges) of the network [32]. Still, we can extract some valuable insights from this algorithm's application. Specifically, in Section 4.5, we analyze the correlation between the self-organization patterns and the modular structure of the network, further characterizing the meaning of a viral event by studying the structure of the temporal networks.

### 4.3 Real-data networks

In this section, we describe the three networks and compare them with the expected results from the literature. It was shown in [38] that heavy-tailed degree distributions characterize real networks due to the presence of hubs, which are nodes with a very high degree compared to the other nodes. For example, the friendship networks would be characterized by more popular people (i.e., with more friends) and less famous people, usually the majority. Therefore, there is a high presence of nodes with a small degree and an exceptional presence of nodes with a high degree. For this reason, we fit the degree distribution of our networks with heavy-tailed distributions (power-law and truncated power-law), and we compare the scaling exponents and the average degrees with the literature. This way, we can identify which regime our networks can be classified. We further describe the different regimes in Section 4.4.

#### 4.3.1 Co-mentions network

The co-mentions network is built as follows. We collected all the Tweets with two or more mentions (i.e., the usernames present in the Tweet corpus that follows the @) and added them to the list of nodes. An edge is created between all the usernames present in the same Tweet<sup>6</sup>, if the edge was already present, the weight of that edge is increased by one. The network comprises 218,762 nodes and 1,079,032 edges, the average degree is 9.9, and the maximum degree is 7,459. As we can see from the strong discrepancy between these values, we expect a strongly heterogeneous graph, different from the random network approximation [39]. The number of connected components in the network is 20,812, and the largest connected component comprises 164,384 nodes and 1,020,456 edges (respectively, 75% and 95% of the total number of nodes and edges). In the following, we only consider the largest connected component, as it includes a high percentage of the total size of the graph.

Table 4.1 shows the nodes with highest degree. These Users are mentioned many times by the Users, and they include companies (*nestle*, *icelandfoods*, *nutellaglobal*), NGOs (*wwf*, *greenpeace*), and pages specifically related to palm oil (*rspoTweets*, *palm-choice*, *orangutans*, *palmoildetect*). It suggests that the act of mentioning someone aims

<sup>6</sup>If there are more than two mentions, we created an edge for all the possible combinations.

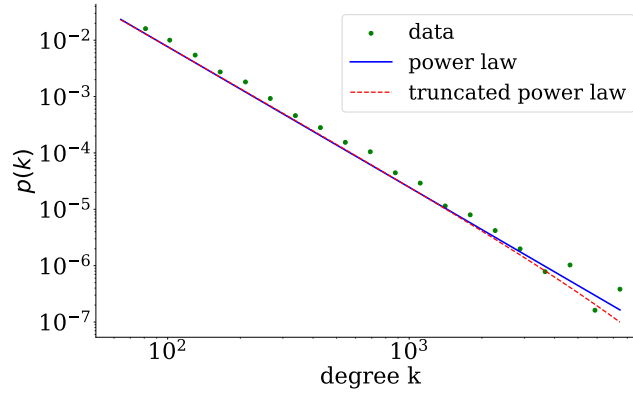


FIGURE 4.1: Degree distribution for the co-mentions network: the blue and the red dotted lines represent the fitted power-law and truncated power-law, respectively, while the green scatter plot represents the empirical data in logarithmic binning.

Username (node)	Number of mentions (degree)
orangulandtrust	7,459
rspoTweets	6,968
palmchoice	6,061
wwf	4,979
greenpeace	4,205
nestle	4,181
icelandfoods	4,128
orangutans	3,939
nutellaglobal	3,758
palmoildetect	3,554

TABLE 4.1: Nodes of the co-mentions graph (largest connected component) with the highest degree.

to require the attention of that User, which can be useful to ask for new sustainability policies for the companies or to support a social campaign launched by an NGO. Then, it is possible to analyze the degree distribution of the largest connected component. As shown in Fig. 4.1, the degree distribution is well-fitted by heavy-tailed distributions, in particular, a truncated power-law with exponent  $\alpha = (2.49 \pm 0.02)$  is the best fit.

### 4.3.2 Co-hashtags network

The co-hashtags network is built similarly to the co-mentions network, but instead of considering Users, we consider hashtags. The network comprises 56,328 nodes and 327,894 edges, the average degree is 11.6, and the maximum degree is 16,106. We expect the presence of hubs in the network again. The number of connected components is 3,021, and the largest connected component counts 48,941 nodes and 321,621 edges (respectively 87% and 98% of the whole graph's nodes and edges). We further investigate the largest connected component again, as it is a good sample of the whole graph.

Table 4.2 shows the nodes with the highest degree. The hashtags are very similar to the most common ones analyzed in previous chapters, even if we measure the



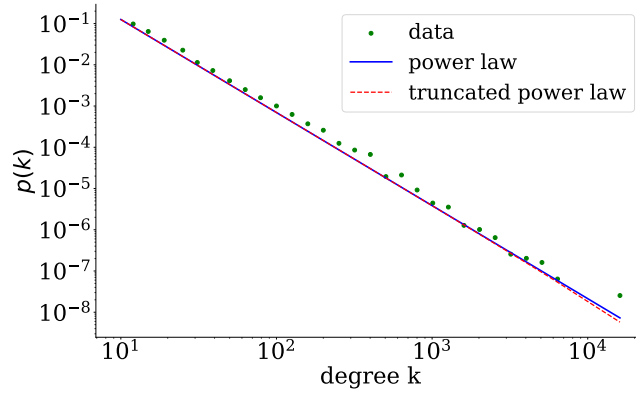


FIGURE 4.2: Degree distribution for the co-hashtags network: the blue and the red dotted lines represent the fitted power-law and truncated power-law, respectively, while the green scatter plot represents the empirical data in logarithmic binning.

Hashtag (node)	Co-occurrence (degree)
palmoil	16,106
palm	6,015
oil	4,723
deforestation	4,087
indonesia	3,805
malaysia	3,257
sustainability	2,672
sustainable	2,630
environment	2,511
orangutans	2,317

TABLE 4.2: Nodes of the co-hashtags graph (largest connected component) with the highest degree.

co-occurrence in this case. Once again, it is possible to stress the Users' attention on the environmental issues, as all these hashtags are related to deforestation and sustainability. Then, we can consider the degree distribution of the largest connected component, as shown in Fig. 4.2, which is well described by heavy-tailed distribution. In particular, a truncated power law with scaling exponent  $\alpha = (2.25 \pm 0.01)$  is the best fit.

### 4.3.3 User-hashtag network

The last network is constructed in the following way. We consider only the Tweets where at least one hashtag is present. We create two subsets of nodes made by the usernames related to these Tweets and the hashtags used. Then, the edges are created by connecting each User with the hashtags they used. Therefore, the graph is bipartite by construction. It comprises 405,564 nodes and 1,010,582 edges, the average degree is 5.0, and the maximum degree is 91,166. This graph is again a good candidate for having a heavy-tailed distribution due to the presence of hubs. The number of connected components is 12,860, and the largest connected component comprises 371,493 nodes and 988,858 edges (92% and 98% of the whole graph's nodes and edges).

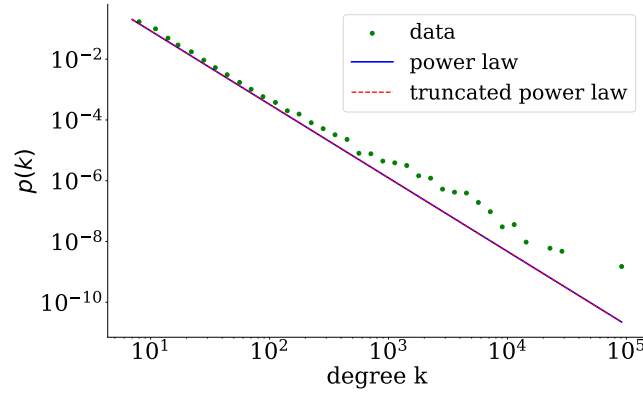


FIGURE 4.3: Degree distribution for the user-hashtag network: the blue and the red dotted lines represent the fitted power-law and truncated power-law, respectively, while the green scatter plot represents the empirical data in logarithmic binning.

Node	Degree
palmoil	91,166
deforestation	27,498
orangutans	18,725
indonesia	14,295
malaysia	10,464
orangutan	10,245
palm	9,556
sustainable	7,895
rainforest	7,844
rspo	7,027

TABLE 4.3: Nodes of the user-hashtag graph (largest connected component) with the highest degree.

Table 4.3 shows the nodes with the highest degree. It is possible to notice that there are no Users among them. It is more probable that a popular hashtag is used many times than a User using many hashtags, considering the limited number of Tweets that the Users are posting concerning this topic. The degree distribution of the largest connected component is well-fitted by a power-law, as shown in Fig. 4.3, with scaling exponent  $\alpha = (2.420 \pm 0.008)$ .

## 4.4 Degree distribution properties: an overview

Table 4.4 shows a summary of all the features extracted from the giant components of these three networks. They differ significantly in size; therefore, it is not easy to measure size-dependent quantities. For this reason, in the next section, we further investigate the *temporal* networks to appreciate how size influences the network structure. We now compare the degree distributions. The first two networks are well described by truncated power laws, while a power law better describes the last one. As the p-values suggest, there is a preference for the truncated power law for the first network. In contrast, for the latter two networks, the p-values are  $p \geq 0.5$ . Therefore, both distributions are acceptable, and there is no preference for one of the

two. Furthermore, the null hypothesis is accepted at 95% of accuracy in all three cases, as the KS distances are smaller than the critical values from the p-value table.

One possibility to distinguish the topological regimes is to analyze the average degree  $\langle k \rangle$  of the network [30]:

- $0 < \langle k \rangle < 1$  (*subcritical regime*): few links, characterized by small connected component, absence of hubs;
- $\langle k \rangle = 1$  (*critical point*): separates the regimes, emergence of a largest connected component with hubs;
- $\langle k \rangle > 1$  (*supercritical regime*): coexistence of isolated, small connected components and a giant component;
- $\langle k \rangle > \ln N$  (*connected regime*): all the nodes are part of the largest connected component, the graph is fully connected.

Co-mentions and User-Hashtag networks are part of the *supercritical regime*, while the Co-hashtags network is in the *connected regime*. These results are compatible with the literature [5]. It was shown that many real networks are either supercritical or fully connected due to their heterogeneous structure and the presence of the hubs.

Another way of classifying the networks is by analyzing the properties of the scaling exponent of the heavy-tailed distribution. We can identify some regimes [30] again:

- $\alpha \leq 2$  (*anomalous regime*): the average degree and the second moment diverge, while the largest degree<sup>7</sup> would be bigger than the total number of nodes in the network. Therefore it is not physically possible;
- $2 < \alpha < 3$  (*scale-free regime*): the average degree is finite, while the second moment diverges;
- $\alpha > 3$  (*random network regime*): both the average degree and the second moment are finite.

All three networks we are considering are in the *scale-free regime*, and a summary of the two properties is shown in Fig. 4.4. Some discrepancies can be identified even if the networks show the characteristics we expect from real networks. The bipartite network is significantly less connected than the others, which is motivated by the construction constraints (i.e., it is impossible to connect a node with another node in the same subset). Another interesting difference is the higher connectedness of the Co-hashtags network than the Co-mentions one. This can be explained by noticing a more substantial variety of the User involved, while the hashtags used for this particular phenomenon are limited. In the following, we study the time evolution of two measures of the network structure, modularity and nestedness, by considering the temporal version of these networks.

## 4.5 Temporal networks: evolution of modularity and nestedness

We present the analysis of the *temporal networks* generated from the empirical dataset, focusing our attention on the time evolution of modularity and nestedness, defined

<sup>7</sup>The size of the largest degree is given by  $k_{max} = N^{1/(\alpha-1)}$

	Co-mentions	Co-hashtags	User-hashtag
Average degree	9.9	11.6	5.0
$\ln N$	12.3	10.9	12.9
Number of nodes (GC)	164,384	48,941	371,493
Number of edges (GC)	1,020,456	321,621	988,858
Best distribution	truncated power-law	truncated power-law	power-law
$p$ -value	0.04	0.5	0.9
degree cutoff	64	10	7
$D$	0.01	0.02	0.01
$\alpha$	<b>2.49</b>	<b>2.25</b>	<b>2.420</b>
$\sigma$	0.02	0.01	0.008

TABLE 4.4: Comparison of the three networks.  $\ln N$  is the logarithm of the total number of nodes in the networks, while the fitting procedure is done on the largest connected component (GC) of each graph, the  $p$ -value measures the significance of choice between candidate distributions, the degree cut-off is the lower bound chosen by minimizing the Kolmogorov-Smirnov distance  $D$ ,  $\alpha$  is the scaling exponent and  $\sigma$  is its standard error.

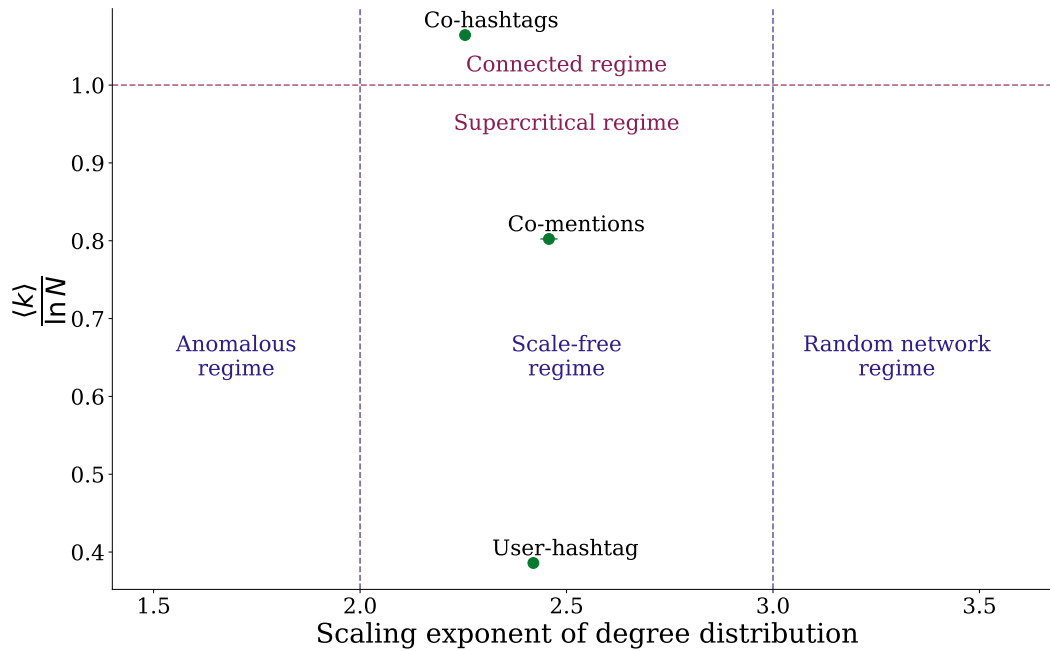


FIGURE 4.4: Graphical visualization of the regimes identified by the average degree and the scaling exponent of the degree distribution. On the y-axis, the average degree is normalized by the natural logarithm of the number of nodes, assuming that none of these networks has an average degree smaller than 1. Therefore the *subcritical regime* is not shown.

in Section 4.2. We focused on the Co-hashtags and the User-hashtag networks because the Co-mentions network's physical interpretation of these measures is unclear. In detail, to generate the temporal networks, we proceeded as follows. We identified sliding windows of three years, then extracted a frame of each network in these sliding windows, i.e., considering only the Tweets in that period. Furthermore, to limit the side effects due to the network size (modularity is strongly dependent on the average degree [32]), we decided to fix the number of nodes per time window. In particular, for the Co-hashtags network, we considered the 1000 most used hashtags per time window, while for the User-hashtag network, we considered the 500 most active Users and the 500 most used hashtags per time window. We are only interested in keeping track of the global properties of the temporal networks. Therefore, we are not interested in considering the evolution of specific communities, and the approximation to the 1000 nodes with the highest degree is the most suitable for our scopes. Modularity relaxes towards values below 0.2 for both networks, with a common relaxation time of five-time windows. As shown in Fig. 4.5, we can observe that the two networks exhibit a similar trend for both measures. While modularity decreases over time, nestedness shows the opposite behavior and keeps growing. In contrast, nestedness' starting values are below 20, while the growth is more robust in the case of the Co-hashtags network. We could sum up the physical interpretation of these behaviors in the following reasoning:

- The divergent trends of the two measures suggest that the emergence of viral events is strongly related to non-modular and well-nested structures. This is in agreement with previous studies [36], which suggests that a well-defined community structure does not allow a significant spread of information. In contrast, the spread of information is favored by the self-adaptation of internal structures and hierarchies due to a higher nestedness;
- In the case of the Co-hashtags network, the diversity of ideas, characterized by the presence of well-defined communities, is missing through time. The initial variety collapses through stable values, and no innovation is brought from the hashtags' point of view;
- Finally, the strong decrease in modularity in the case of the User-hashtag network suggests that, even in the presence of new Users interested in the topic, there is a lack of new ideas. The diversity among communities is relaxing to an asymptotic value.

More generally, our results suggest that viral events such as palm oil campaigns can be described as the advent of highly-nested structures. These structures minimize competition among individuals, as shown in [40], converging to a common point of view by significantly reducing the diversity of ideas and constraining the leading actors of public opinion.

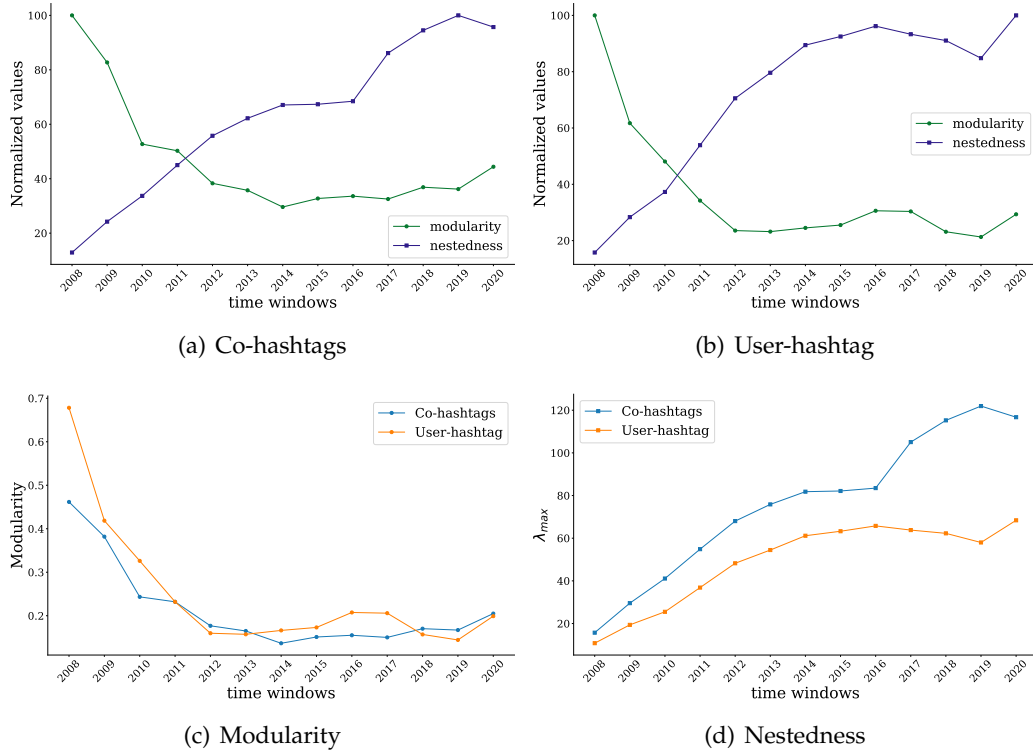


FIGURE 4.5: (a)-(b): Normalized time evolution of modularity and nestedness for temporal networks, realized on time windows of three years. The values are given as percentages of the maximum value to display the change in time, more than the absolute values. (c)-(d): Time evolutions of modularity and nestedness. Here, values are absolute to display the change in magnitude in the different networks.

## Chapter 5

# Interevent time distribution

An interesting branch of Probability Theory is dedicated to the study of waiting time distributions, which we address as Interevent Time (**IET**) Distribution [41]. Classically, the Poisson process [42] was considered the best description for the waiting time distribution. Specifically, it is a memoryless process with a probability density function of the **IETs** in the exponential form

$$p(\tau) = \lambda e^{-\lambda\tau} \quad (5.1)$$

where  $\tau$  is the interevent time and  $\lambda = \frac{1}{\langle\tau\rangle}$  is the rate of activity. It is broadly accepted that in many human activities, the **IET** distribution strongly deviates from the Poissonian approximation (Vasquez et al. [43] first proved it). Human actions are far from being performed at a constant rate, as a Poisson process would predict. Instead, they are characterized by burstiness and jagged peaks of activity. As we confirm with the experimental results, human activities are better described by heavy-tailed distributions [11]. We considered two possible candidate distributions: the *power-law*, with probability density function<sup>1</sup>

$$p(\tau) \propto \tau^{-\alpha} \quad (5.2)$$

where  $\alpha$  is the scaling parameter [44], and the *truncated power-law*<sup>2</sup>

$$p(\tau) \propto \tau^{-\alpha} e^{-\lambda\tau} \quad (5.3)$$

which is a power law with an exponential cut-off [44]. There are many examples of human activities in a virtual environment with heavy-tailed **IETs** distributions: email activity patterns [43], contagion processes, financial markets evolution [45], web browsing [46], opinions on social networks [47].

In our context, the **IET** is defined as the time interval between two Tweets containing the same hashtag. In the following sections, we analyze the **IET** distributions of the most common hashtags, fitting them with heavy-tailed distributions (Section 5.1). Second, we visually explore the emergence of viral events through the time maps (Section 5.2). Finally, we compare the evolution of the **IET** with the time evolution of sentiment (Section 5.3).

<sup>1</sup>The normalization constant is given by  $C = (\alpha - 1)\tau_{min}^{\alpha-1}$

<sup>2</sup> $C = \lambda^{1-\alpha} / \Gamma(1 - \alpha, \lambda\tau_{min})$

## 5.1 Heavy-tailed distributions of IET

To perform this analysis, we selected the ten most used hashtags (in descending order): *palmoil*, *deforestation*, *orangutans*, *palmoilkills*, *indonesia*, *palm*, *killerpalm*, *sustainable*, *malaysia*, *rspo*. We maintained the hashtags *palmoil* and *palm* because they can be seen as a measure of the general behavior of the dataset. As previously stated, the **IET** is the time interval between two consecutive Tweets containing a given hashtag. In practice, a smaller **IET** corresponds to a more substantial interest in that hashtag. Considering the most common hashtags, we expect two types of behaviors. (i) Either they are used frequently and uniformly throughout the dataset, characterized by a Poisson process, or (ii) there is a burst peak of interest, and afterward, they are forgotten. Heavy-tailed distributions characterize the latter case. Therefore, we expect to find this bursty behavior.

Using the *power-law* package [48], which allows visualizing, fitting, and comparing heavy-tailed distributions, we obtained the following results for the most common hashtags (shown in Fig. 5.1). Here we describe the simple procedure done to determine the best fit:

1. Compute the log-likelihood ratio  $R$  between two candidate distributions (power-law and truncated power-law) and their significance value  $p$ .
2. The sign of  $R$  establishes which distribution is more likely to fit the data. At the same time,  $p$  quantifies how accurate the decision is ( $p \geq 0.5$  indicates that neither distribution is better than the other).
3. After choosing the best distribution, it is necessary to choose the optimal value of  $\tau_{min}$ <sup>3</sup>, which is computed by minimizing the Kolmogorov-Smirnov distance (**KS**)<sup>4</sup> between the data and the fit.

As shown in Table 5.1, the distributions of seven hashtags are better described by truncated power-laws, while power-laws describe just three. Nevertheless, the  $p$ -values of the comparison between candidate curves are exceptionally high (neglecting the “killerpalm” case), which means that both the truncated power law and the power law are good fits for these empirical distributions, without a strong preference for one of the two hypotheses. Furthermore, the null hypothesis is accepted at 95% of accuracy in all cases, and it has been done by comparing the KS distance with the critical values from the P-value table<sup>5</sup>. All the values of the scaling exponents are bigger than 1, as expected. Distributions with  $\alpha < 1$  are not normalizable. Neglecting the “killerpalm” case again, we can notice that  $2 < \alpha < 3$ . This is the *scale-free* regime, where the first moment of the distribution is finite, while the higher moments diverge. We should consider the cases “rspo” and “indonesia” because their scaling exponent is compatible with  $\alpha = 3$ , which is the critical point between two different regimes. For  $\alpha > 3$ , the second moment of the distribution does not diverge.

<sup>3</sup> $\tau_{min}$  is the lower bound for the scaling range: it has to be fixed because the power-law is not defined for  $\tau = 0$ .

<sup>4</sup>**KS** statistics is defined as [44]

$$D = \max_{\tau \geq \tau_{min}} |S(\tau) - P(\tau)|$$

where  $S(\tau)$  is the CDF for the empirical data, while  $P(\tau)$  is the CDF for the heavy-tailed distribution we are considering as hypothesis (*null model*).

<sup>5</sup>**P-value table** can be found at this URL.



Hashtag	Best distribution	$n$	$p$ -value	$\tau_{min}$	$D$	$\alpha$	$\sigma$
killerpalm	truncated power-law	60	0.1	0.001	0.02	1.428	0.004
palmoilkills	truncated power-law	17	0.4	86	0.03	2.4	0.1
deforestation	truncated power-law	18	0.5	12	0.02	2.47	0.04
sustainable	truncated power-law	15	0.6	40	0.03	2.52	0.07
palmoil	power-law	27	0.9	1	0.02	2.55	0.01
orangutans	truncated power-law	14	0.4	36	0.03	2.64	0.07
malaysia	power-law	17	0.9	47	0.04	2.80	0.08
palm	truncated power-law	19	0.6	16	0.03	2.80	0.05
rspo	power-law	16	0.9	56	0.04	2.9	0.1
indonesia	truncated power-law	15	0.9	62	0.03	2.9	0.1

TABLE 5.1: The fit results:  $n$  is the number of bins considered (the number of scatter points in Fig. 5.1).  $p$  is the significance value of the comparison between candidate distributions,  $\tau_{min}$  is the lower bound of the scaling range [48],  $D$  is the Kolmogorov-Smirnov distance,  $\alpha$  is the scaling parameter,  $\sigma$  is the standard error.<sup>6</sup>

The scaling exponents of the IET distributions can be interpreted as a measure of *memory*, in which hashtags that are described by distributions with bigger scaling exponents are forgotten faster than the others. They present a lower probability of showing bigger IETs than the other hashtags. It is strong evidence that, even if human activities are characterized by burstiness and ease of forgetting, their behavior strongly depends on the attention given to a particular social or political issue. Looking at two specific cases, while “indonesia” is a topic that presents a lower peak of interest but is more uniformly considered in time, “killerpalm” was probably used in a specific public campaign and forgotten afterward. We further analyze these behaviors in the following. We compare the scaling exponents with those related to the cascade distributions to develop a more general framework and characterize the different regimes.

## 5.2 Time maps: a tool for multiple timescales visualization

We are interested in visualizing IETs that differ significantly. An insightful way to visualize discrete events on different timescales is by using time maps [49].

Consider a set of events<sup>7</sup>  $E = \{e_i \text{ for } i = 1, \dots, n\}$ . The IETs are the time intervals between the events: denote them as  $IET = \{t_i \text{ for } i = 0, \dots, n-1\}$ . The time map is realized by representing each event  $e_i$  as a point in the plane, with coordinates  $e_i = (t_{i-1}, t_i)$  for  $i = 1, \dots, n-1$ . A graphical representation of the time map construction is shown in Fig. 5.2. A time map can be divided into four regimes:

- low  $x$ , low  $y$  coordinates (*viral event region*): the events happen one next to the other, the time intervals are small before and after any event in this region, and the IETs are quite uniform;

<sup>6</sup> $\sigma = \frac{\alpha-1}{\sqrt{n}} + O(1/n)$

<sup>7</sup>In this paragraph, we intend as an event a Tweet that contains a particular hashtag.

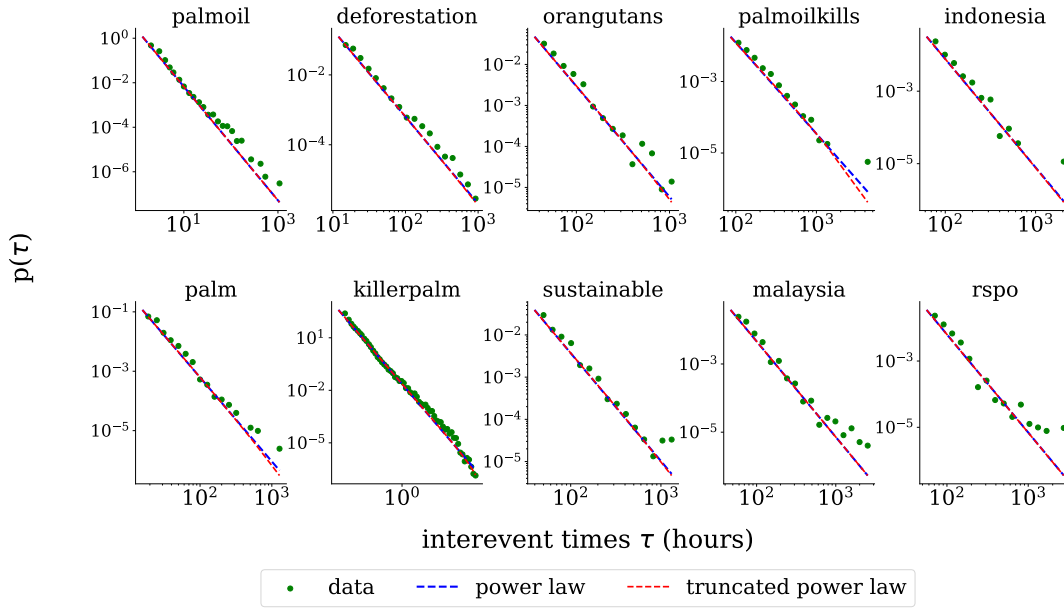


FIGURE 5.1: IET distributions of the ten most used hashtags. The red and blue dashed lines represent the fitted power laws and truncated power laws, while the scatter plots represent the data we are fitting.

- low  $x$ , high  $y$  coordinates (*pre-viral event region*): the time interval before the event is bigger than the one after the event; the interest in the topic is increasing, and we are possibly approaching a viral event;
- high  $x$ , low  $y$  coordinates (*post-viral event region*): the time interval before the event is smaller than the one after the event; the topic is becoming less interesting, a viral event has possibly just happened;
- high  $x$ , high  $y$  coordinates (*relaxed region*): the time intervals are big before and after the event, with uniform values. The interest in the topic is low and constant.

As shown in Fig. 5.3, these are the time maps for the most common hashtags. We can visualize the differences in the behavior of the hashtags. *orangutans*, *indonesia*, *malaysia*, *rspo* have a peak of interest around 2018, when the Iceland Foods campaign spread out; *palmoilkills* had spread just in 2016, while it is not mainly present in other years; *palmoil* and *deforestation* are largely present in the whole dataset. Therefore, we cannot identify a precise peak of interest, as the scatter points strongly overlap. *Palm* and *sustainable* do not present strong evidence of viral events, as they are vaguely used throughout the dataset. *killerpalm* has a peculiar behavior because it was mainly used during a probable viral event in 2016. It was forgotten in the next years and used again around 2021.

### 5.3 Viral events and Negative sentiment

In Section 3.1.2, we showed that it is possible to relate the emergence of a viral event to sentiment dynamics. More specifically, the growth of attention on a topic is followed by an increase in the *Negative* attitude. Another way to highlight the relation between a change in sentiment and the emergence of viral events is by analyzing the time evolution of sentiment and the IETs. Next, we define two quantities:

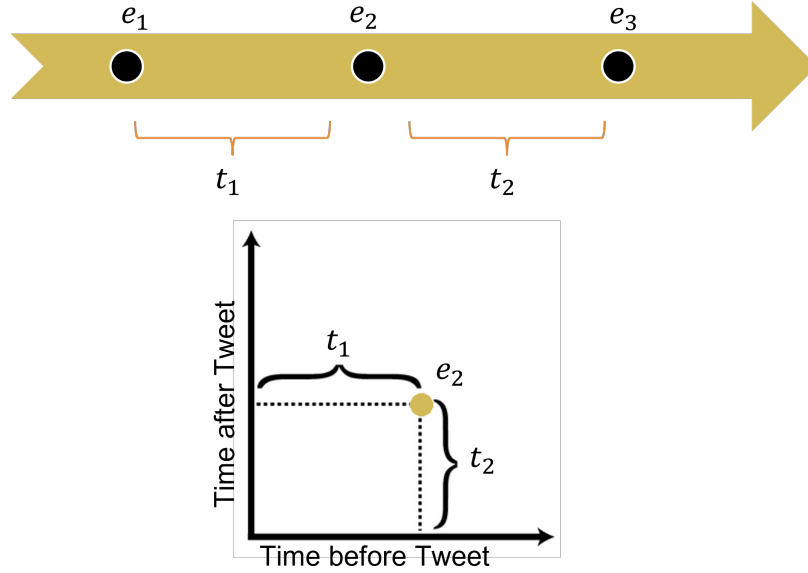


FIGURE 5.2: (a): graphical representation of the time intervals between Tweets containing a given hashtag. (b): graphical representation of the time map construction. Each event has the time before the event as an x-coordinate and the time after the event as a y-coordinate.

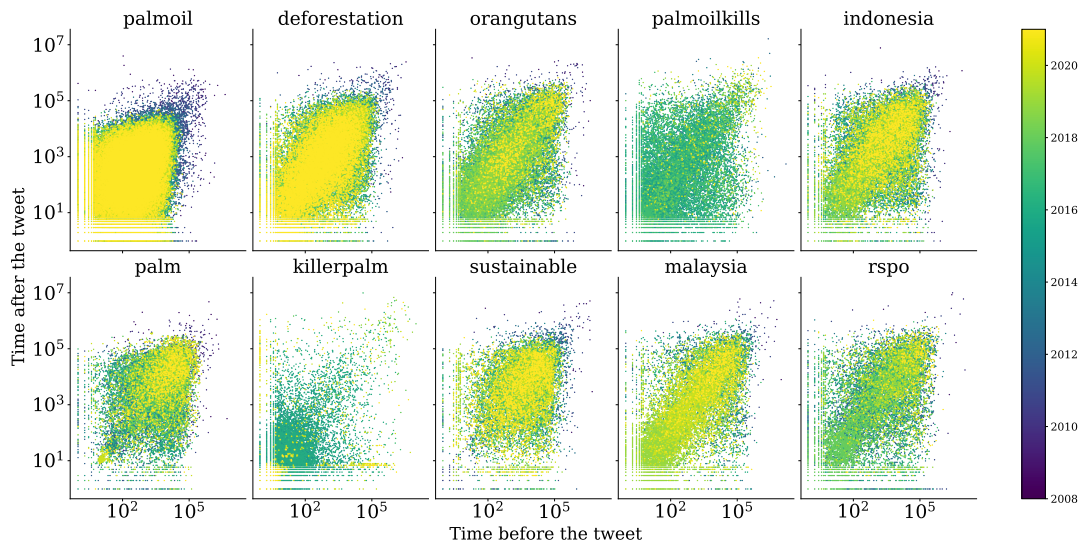


FIGURE 5.3: Time maps of the most common hashtags. The plot is in log-log scale and each event  $e_i$  has coordinates  $(t_{i-1}, t_i)$  for  $i = 1, \dots, n - 1$ . The heat map colors represent in which year this event happened.

- *weighted sentiment*:

$$\langle l(t) \rangle = \frac{\sum_{i=0}^2 s_i \cdot i}{\sum_{i=0}^2 s_i},$$

where  $i \in [0, 2]$  is the sentiment label and  $s_i$  the related score;

- *IET rolling average*: a set  $\{r_i(w)\}$  for  $i = 1, \dots, n - 1$  where

$$r_i(w) = \frac{1}{w} \sum_{j=i}^{i+w-1} t_j$$

is the rolling average over a window of size  $w$ ,  $t_j$  is the  $j$ -th **IET** and  $n$  is the total number of **IETs**.

As shown in Fig. 5.4, it is possible to compare the time evolution of the **IET** rolling average and the weighted sentiment. Every hashtag presents a different behavior, as the weighted sentiment firmly varies. From the hashtag that can be defined as *mainstream*<sup>8</sup>, no further insights can be found. They never present an intense viral event, which would be characterized by a sudden decrease in the **IET** rolling average. Therefore, these cases have no strong correlation between **IET** and sentiment. Nevertheless, for the so-called *viral hashtags*<sup>9</sup> an interesting behavior emerges. It is possible to observe that there are parts in which the **IET** remains close to zero for a long number of Tweets, and these parts correspond to a sudden decrease of the weighted sentiment. It is interesting to show that viral events regarding sustainability and social issues are strongly related to a *Negative* attitude and strong, *Negative* tones.

---

<sup>8</sup>palmoil, deforestation, palm

<sup>9</sup>indonesia, sustainable, malaysia, rspo

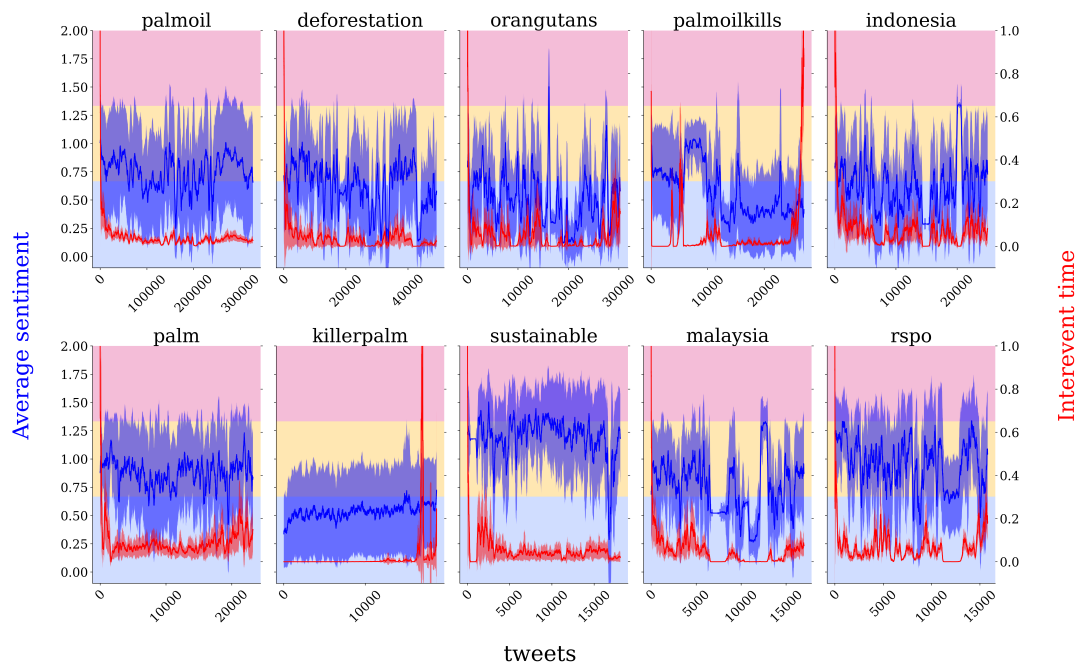


FIGURE 5.4: Time evolution of the **IET** and sentiment's rolling averages for the most common hashtags. The **IET** is normalized, such that  $r_i \in [0, 1]$  and the weighted sentiment is  $\langle l(t) \rangle \in [0, 2]$ . We suggest a strong correlation in the "viral hashtags" between shorter **IET**s and a drop in the weighted sentiment.

## Chapter 6

# Cascade size distribution

*Cascade size statistics* is a topic broadly studied in many research fields. It can be related to many phenomena, such as neuronal firings mechanisms [50], earthquakes [51], fractures in porous media [52], magnetic systems [53], spreading of diseases [54] and information [55] among individuals. It is possible to generically define a cascade as the series of events generated as consequences of a first significant event. There are many possible definitions of a cascade in practice, as this concept applies to various topics. Many features of the cascades were studied through the years. For example, the (average) cascade shape [56], whose symmetry strongly depends on the nature of the phenomenon (i.e., if it is Poissonian or not), and several models, aimed to describe the dynamics of the cascade phenomenon, were developed. The first, paradigmatic example is the threshold model developed by Watts [57], which initiated further developments, including models considering the underlying network structure [58].

In our particular context, it is possible to define the cascade size (CS) for two different quantities:

- *hashtags cascade*  $CS(h)$ : number of Tweets in a time window containing a given hashtag  $h$ ;
- *sentiment cascade*  $CS(s)$ : number of Tweets in a time window labelled by a given sentiment  $s$ .

In the following, we analyze the cascade size distributions by fitting heavy-tailed distributions on the empirical data. With the first definition (*hashtags cascade*), we compare the scaling exponents of the cascade distributions for the most common hashtags (Section 6.1), completing the framework of the IET distributions and providing another tool to interpret virality. Using the second definition (*sentiment cascade*), we analyze the virality of the different sentiments, bringing another evidence of the relationship between viral events and negative sentiment (Section 6.2).

## 6.1 Hashtags cascade

Looking at the time evolution of the cascade size (Fig. 6.1 represents the time evolution of cascade size in the case of the hashtag *palmoil*), it is possible to observe that there are sudden increases in the cascade size, indicating the emergence of viral events. We analyzed the CS distributions similarly to the IET distributions, using the fitting procedure described in Section 5.1, as shown in Fig. 6.2. Table 6.1 shows the fitting parameters resulting from the analysis. The best distribution is the power law for all the hashtags, except for *killerpalm* and *orangutans*. Still, the significance value  $p$  is bigger than 0.5 for all of them, except for *orangutans*. Therefore, there is not a strong preference between the two candidate distributions, and they both describe

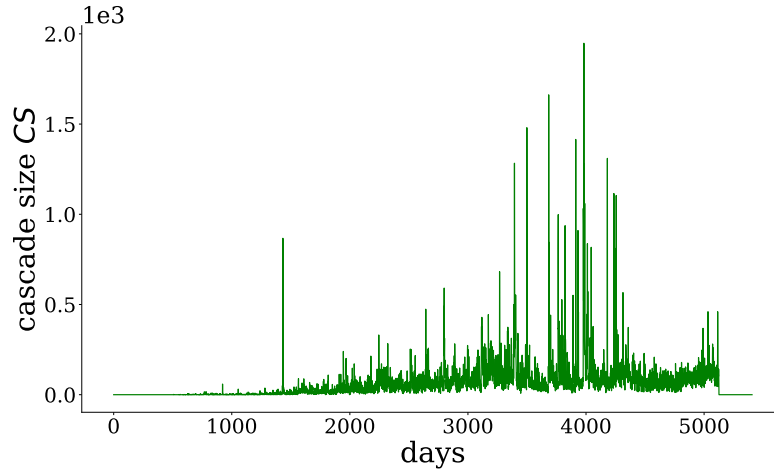


FIGURE 6.1: Time evolution of the cascade size for the hashtag *palmoil*. The cascade size  $CS(h)$  is calculated as the number of Tweets per day containing a given hashtag  $h$ .

Hashtag ( $h$ )	Best distribution	$n$	$p$ -value	$CS_{min}(h)$	$D$	$\alpha$	$\sigma$
killerpalm	truncated power-law	24	0.8	19	0.08	2.01	0.09
orangutans	truncated power-law	17	0.5	14	0.03	2.26	0.06
deforestation	power-law	19	0.9	20	0.02	2.57	0.07
malaysia	power-law	23	0.8	7	0.02	2.61	0.07
palmoilkills	power-law	17	0.9	56	0.04	2.8	0.2
indonesia	power-law	16	0.9	15	0.03	2.93	0.09
rspo	power-law	18	0.9	14	0.04	3.1	0.1
palmoil	power-law	12	0.9	103	0.02	3.16	0.07
sustainable	power-law	16	0.9	18	0.03	3.9	0.3
palm	power-law	11	0.9	17	0.03	4.3	0.2

TABLE 6.1: The fit results:  $n$  is the number of bins considered (the number of scatter points in Fig. 6.2).  $p$  is the significance value of the comparison between candidate distributions,  $CS_{min}(h)$  is the lower bound of the scaling range [48],  $D$  is the Kolmogorov-Smirnov distance,  $\alpha$  is the scaling parameter,  $\sigma$  is the standard error.

well our empirical data. It is confirmed by the Kolmogorov-Smirnov test, which is below the critical value in all the cases. Therefore, the null hypothesis is accepted at 95% of accuracy. The scaling exponents are comprised of different ranges, indicating different behaviors of these distributions. When  $2 < \alpha < 3$ , the average is finite while the second moment is not. We observe a significant difference in the cascade sizes; there are few events with a bigger size. This is less evident for the distributions with  $\alpha > 3$ , where the first and the second moment are both finite. In this case, there is less discrepancy in the size of the different events, and there are fewer events with a bigger size (or more events with a bigger size, still, there is less discrepancy in the sizes of the events). The scaling exponent  $\alpha$  is a measure of the intensity of a viral event and the lower it is, the higher the virality of a given hashtag. In Chapter 7, we fully characterize the virality of the most common hashtags.



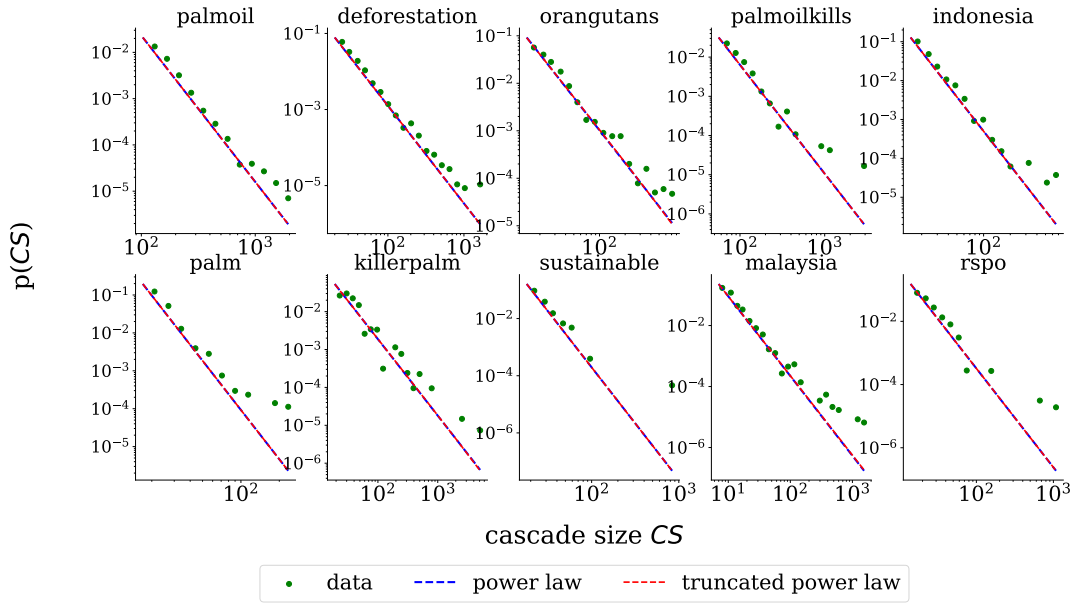


FIGURE 6.2: Hashtags cascade distributions of the most common hashtags. The red and blue lines represent the fitted power-law and truncated power-law distributions. The scatter plot represents the empirical data (in logarithmic binning) that we are fitting. The red and blue curves overlap for each hashtag.

## 6.2 Sentiment cascade

As defined in Chapter 6, the sentiment cascade is given by the number of Tweets in a time window<sup>1</sup> labeled by a given sentiment. Fig. 6.3 shows the three empirical distributions obtained and the relative results of the fitting procedure with heavy-tailed distributions. As shown in Table 6.2, the *Negative* cascade is described by a truncated power-law. It was found to better describe the distribution than the power-law due to the p-value lower than 0.5. The *Neutral* and *Positive* cascades are well described by both distributions, as the p-value is bigger than 0.5. The null hypothesis is accepted with an accuracy of 95% in all the cases. Therefore we can assume that the empirical distributions are well-described by these heavy-tail distributions above the lower bound  $CS_{min}(s)$ . Regarding the scaling exponents, we can observe that *Negative* cascade is characterized by  $2 < \alpha < 3$ . Therefore, we expect the *Negative* Tweets to be more viral than the other labeled Tweets. Furthermore, the *Positive* cascade exponent is still in the same regime of the *Negative* one, but it is close to the critical point  $\alpha = 3$ . It seems that the *Positive* Tweets are less viral than the *Negative* ones. In contrast, the *Neutral* cascade exponent is in the regime  $\alpha > 3$ . Our results suggest that the *Neutral* Tweets are not part of viral events or are not particularly relevant to the discussion.

<sup>1</sup>In our analysis, we considered a time window of one day.



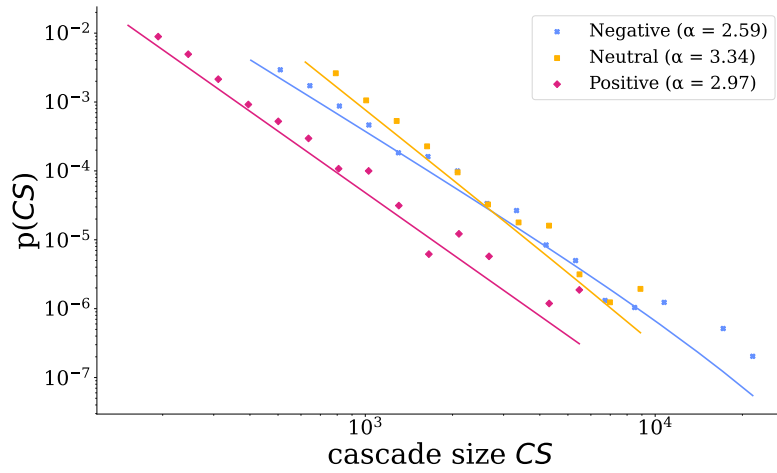


FIGURE 6.3: Sentiment cascades. The scatter plots represent the empirical data fitted by heavy-tailed distribution, while the lines represent the distributions that best fit the data. *Negative* sentiment is more viral than the *Positive* one, while the *Neutral* sentiment is out of the virality regime.

Label (s)	Best distribution	$n$	$p$ -value	$CS_{min}(s)$	$D$	$\alpha$	$\sigma$
Negative	truncated power-law	17	0.2	403	0.02	2.59	0.05
Neutral	truncated power-law	11	0.8	621	0.02	3.3	0.1
Positive	power-law	15	0.9	152	0.02	2.97	0.06

TABLE 6.2: The fit results, where  $n$  is the number of bins considered (the number of scatter points in Fig. 6.3).  $p$  is the significance value of the comparison between candidate distributions,  $CS_{min}(s)$  is the lower bound of the scaling range [48],  $D$  is the Kolmogorov-Smirnov distance,  $\alpha$  is the scaling parameter,  $\sigma$  is the standard error.

## Chapter 7

# Virality phase diagram

According to the Oxford English Dictionary<sup>1</sup>, the term *virality* represents

*"the tendency of an image, video, or piece of information to be circulated **rapidly** and **widely** from one internet user to another"*.

Therefore, a coherent analysis of the viral events need to consider the two dimensions defining virality, i.e., the speed of diffusion and the size of these events. The first task was broadly discussed through the **IET** analysis in Chapter 5, where we studied the persistence in memory of the most common hashtags, considered indicators of viral events. However, it is also necessary to tackle the second task to characterize a viral event thoroughly. For this reason, in Chapter 6, we introduced the **CS** statistics, a technique to analyze a viral event's magnitude. In this way, it is possible to determine if an exceptional interest also characterized a bursty event.

Following the definition of *virality*, it is necessary to introduce two dimensions to characterize a viral event: the interevent time (**IET**) and the cascade size (**CS**). We identify and characterize the following regimes<sup>2</sup>:

- $\alpha_{IET} < 2 \quad \wedge \quad \alpha_{CS} < 2$  *doubly unpredictable regime (Region I)*: the average and higher moments are not defined for both distributions. In this regime, we expect unpredictable behavior both on the timing and magnitude sides due to the lack of well-defined averages;
- $\alpha_{IET} < 2 \quad \wedge \quad 2 < \alpha_{CS} < 3$  *unpredictably fast, virally big regime (Region II)*: the average of the **AS** is finite, while the average of the **IET** distribution and the higher moments of both distributions are diverging. We expect the viral events to happen unpredictably fast due to the lack of a well-defined **IET** average, while the magnitude of the event is viral, meaning that we can identify an average size, but the infinite variance allows the emergence of a viral event;
- $2 < \alpha_{IET} < 3 \quad \wedge \quad 2 < \alpha_{CS} < 3$  *viral regime (Region III)*: both the averages are finite, while the higher moments are still ill-defined. This regime is characterized by the emergence of viral events that are both bursty and exceptional;
- $\alpha_{IET} > 3 \quad \wedge \quad 2 < \alpha_{CS} < 3$  *unvirally slow, virally big regime (Region IV)*: both the averages and the second moment of the **IET** distribution are finite. Therefore, the emergence of viral events is well defined from the magnitude perspective, while on the other side, the speed of viral events does not diverge significantly from the other periods;

<sup>1</sup><https://www.lexico.com/definition/virality>

<sup>2</sup>Here,  $\alpha_{IET}$  and  $\alpha_{CS}$  are respectively the scaling exponents of the **IET** distribution and the **CS** distribution. We describe only the regimes of interest, where any distribution is found, for the sake of brevity.

- $2 < \alpha_{IET} < 3 \quad \wedge \quad \alpha_{CS} > 3$  *virally fast, unvirally small regime (Region V)*: the averages are both finite, and the variance of the cascade size distribution is also. The viral events are bursty, but there is not a strong discrepancy between the sizes of a viral and an unviral event;
- $\alpha_{IET} > 3 \quad \wedge \quad \alpha_{CS} > 3$  *unviral regime (Region VI)*: both the averages and the second moments are finite for both distributions. Viral events do not characterize this regime.

In Section 7.1, we study the virality phase diagram of the **Palm** dataset, characterizing the behavior of the different hashtags. Finally, in Section 7.2, we give a more general framework by considering the phase diagram created with the three datasets, **Palm**, **Olive**, and **Coconut**.

## 7.1 Palm oil dataset

In this section, we consider the scaling exponents found in Sections 5.1 and 6.1 for the **Palm** dataset. Fig. 7.1 shows the virality phase diagram and where each distribution is located. We can observe that *killerpalm* is found at the border between regions I-II. It can be classified as an unpredictable, viral event principally from the timing perspective but partially from the magnitude point of view. This hashtag was probably used for a single viral event and was forgotten afterward. *Malaysia*, *deforestation*, *orangutans* and *palmoilkills* are located in the region III: they describe viral events as expected from the empirical and theoretical predictions [56], i.e., with heavy-tailed distribution with finite average and diverging higher moments. Furthermore, *palmoil*, *sustainable* and *palm* are located in region V, which shows something coherent with the construction of the dataset. There are moments when these hashtags are massively used. However, due to the constraints imposed to create the dataset (i.e., we chose the Tweets containing the words “palm oil”), it is expected to find the hashtags with similar words (*palmoil*, *palm*) throughout all the dataset. We expect them with a more consistent presence.

The *sustainable* hashtag found in this region V suggests once again that the palm oil topic is strongly connected to the sustainability issue, as we can observe that these hashtags are used with a similar frequency and quantity. Moreover, *rspo* and *indonesia* are borderline cases. They are compatible respectively with regions III-V-VI and III-IV-V. We can suggest that *rspo* is a viral hashtag from the timing point of view. Nevertheless, it is compatible with the *unviral* region. We suggest that this hashtag is probably present in many tweets complaining about the lack of new sustainability policies, as the RSPO is an organization that shall control palm oil production. Therefore, we expect the use of this hashtag throughout the dataset. Similarly, *indonesia* can be considered a viral hashtag characterized by viral events only in magnitude and not from the timing perspective.

From these observations, it is possible to suggest the presence of a strong heterogeneity in the virality of different hashtags, even by considering the most common ones. As these hashtags can be used to quantify human activities on social media, we could state that we behave in a burst, sometimes unpredictable way as human beings. Therefore, further quantitative studies on human behaviors might be challenging and fascinating simultaneously. To further analyze these behaviors, in the next section, we find a comparison among the different oils. The aim is to check if the debates around palm oil evolved similarly to the ones around different topics.

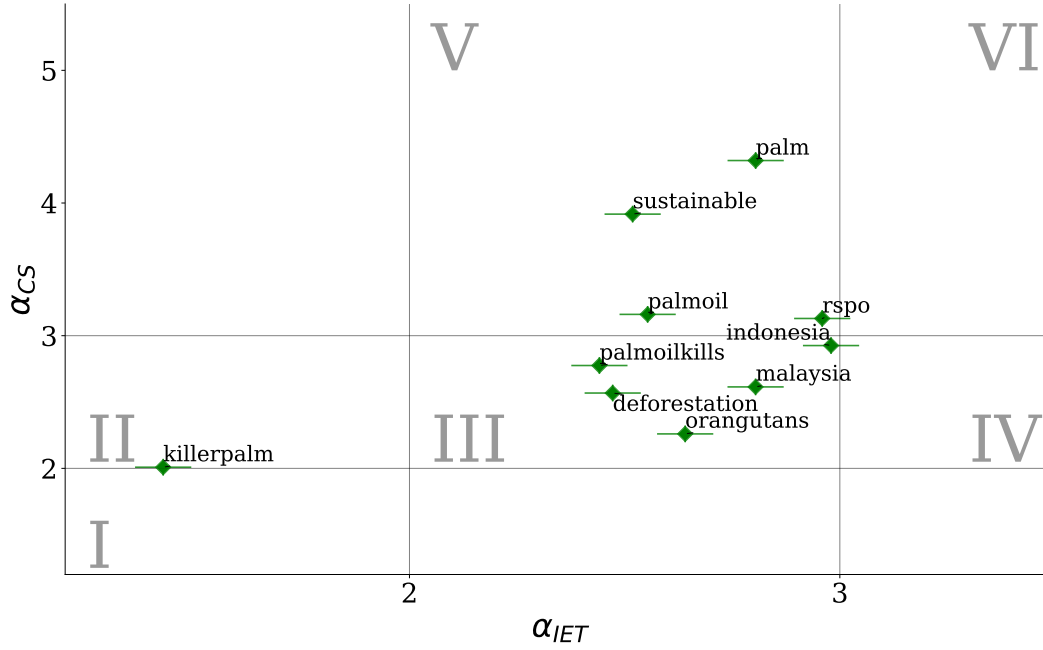


FIGURE 7.1: Virality phase diagram determined by the scaling exponents of the IET and CS distributions. The regions are defined by the mathematical properties of the distributions, which depend on the scaling exponents: for  $\alpha < 2$ , all the moments are diverging, while for  $2 < \alpha < 3$ , the first moment (average) is finite and the higher moments diverge. Finally, for  $\alpha > 3$ , the second moment is also finite.

## 7.2 Virality phase diagram: comparison among oils

In the following, we complete the results found in Section 7.1, by adding the results obtained from the **Olive** and **Coconut** datasets. Further details from IET and CS distributions for the two other datasets are discussed in Appendix A.4 and A.5. Here, two different ways of representing the viral diagram are shown. The first diagram, shown in Fig. 7.2, represents the most common hashtags of the three datasets, with different colors and markers that represent the belonging to a dataset. The second representation is shown in Fig. 7.3, where every scatter point (hashtag) is labeled by the most prevalent sentiment. In other words, for a given hashtag, we restricted the dataset to the Tweets where this hashtag was present, and we extracted which was the most common label among the three sentiments.

Some remarkable considerations can be made:

- *palm oil is the most viral topic*: the region III contains only **Palm** hashtags. This could tell us that other debates are led differently, and only the sustainability issues regarding palm oil show this balance between speed of diffusion and magnitude of the event;
- *contest-based hashtags are unpredictably fast*: the hashtags related to a challenge or a contest (i.e., *win*, *giveaway*) are strongly unpredictable. A massive engagement does not follow them, and they have restricted attention developed around the topic. The Users react to these Tweets just for personal benefit; they are not engaged in bigger social or political campaigns;
- *olive oil is the least viral*: neglecting the contest-based hashtags, all the **Olive** hashtags are found in the less (or non) viral regions (IV-V-VI). This behavior

might be a symptom of a weaker interest in this oil's social and political issues. Looking at the semantics of these hashtags, we can see that they are more related to health and food topics rather than sustainability;

- *virality is Negative*: the hashtags found in the region **III** are also the hashtags whose prevalent sentiment is the *Negative* one, neglecting *killerpalm*. Again, we can prove that the virality of an event is strongly correlated to its sentiment, in particular to a *Negative* one;
- *non-virality is Neutral*: accordingly to the results found in Section 6.2, *Neutral* hashtags are found in non-viral regions. This evidence supports again the relationship between the attitude used in a Tweet and the virality that arise from it;
- *contests are the expression of positivity*: the only *Positive* hashtags are the ones related to a contest (*win, giveaway*) or, more generally, to food and cosmetics topics.

More generally, it is possible to notice a detachment between the three oils' debates. The sustainability issues related to palm oil went more viral than the debates around olive or coconut oil. One reason could be found in looking at these other oils from the public opinion's point of view. There is a lack of demonization around them, and the customers do not feel the sustainability danger while buying these products. All the discussions about coconut are more focused on the cooking recipes and the skincare advice. Both coconut and olive oil are more related to health issues than the palm oil discussions, but these issues are seen as less problematic, as the *Neutral* sentiment suggests. It is interesting to notice how similar topics can be treated so differently, depending on the rise or not of social campaigns.

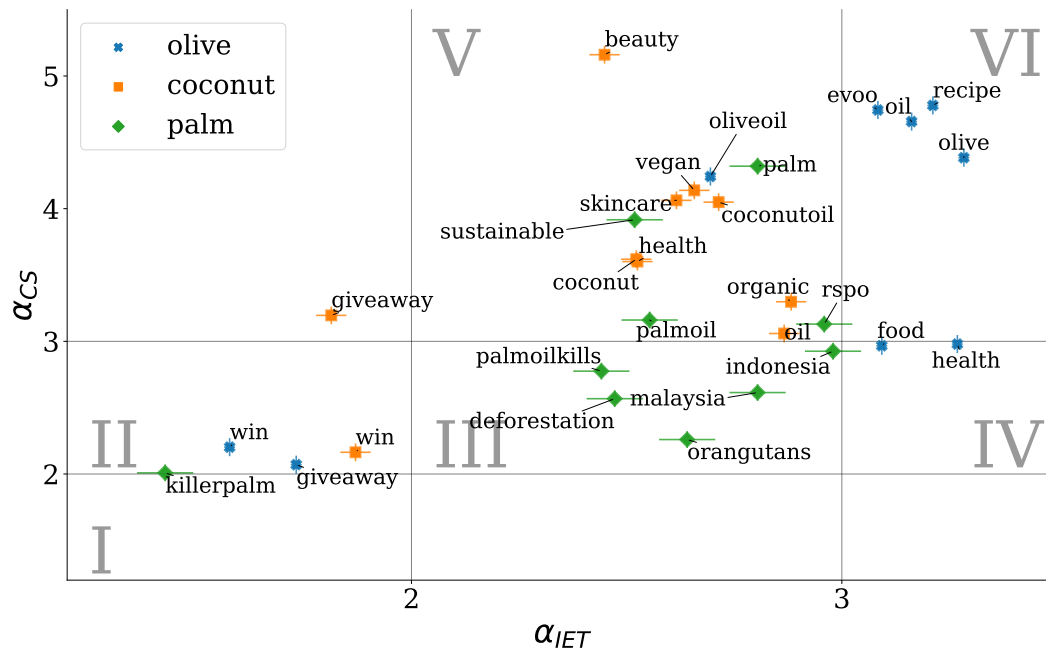


FIGURE 7.2: Virality phase diagram composed of the most common hashtags of the three datasets (**Olive**, **Palm**, **Coconut**). Different markers and colors identify the belonging to the different datasets. We can notice that only **Palm** hashtags are found in the *viral region* **III**, while most of the **Olive** hashtags are found as non-viral.

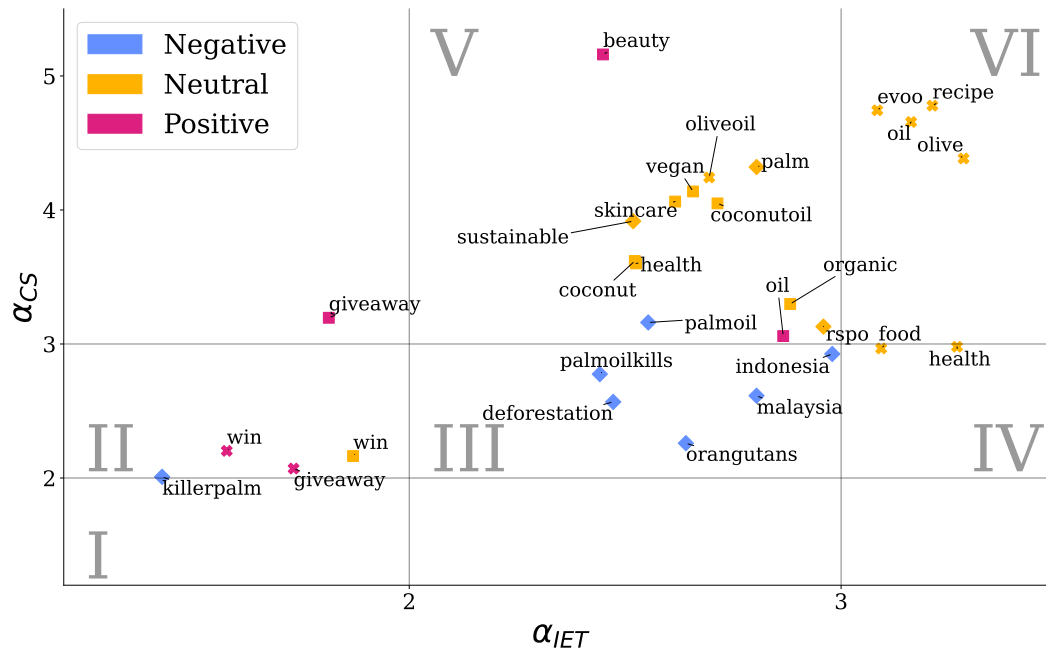


FIGURE 7.3: Virality phase diagram composed of the most common hashtags of the three datasets (**Olive**, **Palm**, **Coconut**). Different markers identify the belonging to the different datasets, while the prevalent sentiment determines the colors. We can notice that the hashtags in region **III** are prevalently *Negative*, while neutrality predominates in the less viral regions (**IV-V-VI**).

## Chapter 8

# Conclusion

In this thesis, we characterized the structure of viral events using a data-driven approach, i.e., analyzing the evolution of a public debate on social media.

In Chapter 2, we detected the emergence of viral events through the analysis of the time evolution of the number of Tweets, noticing a substantial interest in Iceland Foods' TV advertisement. We identified the *Opinion Drivers*, a mix of companies, NGOs, activists, and news websites. Meanwhile, we extracted the most used keywords and hashtags, noticing that they are mainly related to sustainability issues. Moreover, the geographical distribution of the Tweets strongly reflects the engagement of each country in the public debate.

In Chapter 3, we showed that sentiment dynamics is correlated with the emergence of viral events, particularly noticing the presence of a peak of *Negative* sentiments in the year of the Iceland Foods campaign. Then, we detected a counter-intuitive behavior, the *Opinion Drivers* use fewer *Negative* attitudes than the *General Public*. Afterward, our results indicate no correlation between our case study's sentiment classification task and opinion detection. Comparing different vegetable oils, we noticed that the debate around palm oil is strongly monothematic. The other oils present a broader range of preponderant topics.

In Chapter 4, we studied three types of networks, analyzing their structure and locating them in regimes identified by the scaling exponent of the degree distribution and the average degree. The three networks are all located in the *scale-free* regime, as expected by other experiments on real networks [30], while two of them are in the *supercritical regime* (regarding the average degree), and one is found in the *connected regime*. The latter evidence brings some valuable insights into the structure of these networks, further confirmed by Section 4.5, where we analyzed the time evolution of modularity and nestedness in temporal networks. We showed that the network structure changes significantly, and the diversity of ideas is reduced, yielding collaboration and self-organization among users.

In Chapter 5, we analyzed the *IET* distributions, quantitatively measuring the persistence in memory of the most popular hashtags. We noticed that most of these hashtags could be found in the regime where the average is finite, but the second moment diverges. It is compelling evidence of the volatility of viral events on social media. Afterward, we identified the viral events for each hashtag using an innovative visualization technique (the *time maps*). We related the sentiment to the *IET*, analyzing the change in their averages.

In Chapter 6, we analyzed, similarly to the *IET* distributions, the *CS* distributions for the most common hashtags, noticing that only half of these hashtags can be found in the *viral* regime, where the average is finite, and the higher moments are diverging. This evidence can tell us that many hashtags are broadly used along the dataset without showing significant peaks in the size of the event related to their appearance. Furthermore, we studied the *sentiment cascades*, showing that the *Negative*



Tweets tend to be more viral than the *Positive* ones, and the *neutrality* is correlated to the absence of *virality*.

In Chapter 7, we combined the results obtained from the two previous chapters. We created a more general pipeline to characterize the *virality* of an event. The two dimensions of the virality phase diagram are related to the timing and size perspectives: a viral event should be fast and widely spread. We used the mathematical properties of the heavy-tailed distribution to explore these two dimensions quantitatively. We compared the results of different vegetable oils, noticing that the debate around palm oil went more viral than the others and showing a strong correlation between *virality* and *Negative* sentiment again. Other results previously obtained were confirmed by this further analysis, such as the correlation between *Neutral* sentiment and *non-virality*.

The discussion around the sustainability of palm oil is an interesting and ever-present example of how an awareness-raising campaign can bring the general public's attention to hitherto little-addressed issues and how the environmental issue is becoming preponderant in public discussion. The thread that binds the various techniques used in this thesis is linked to a fundamental concept in *Probability Theory*, which is the study of large deviations. Human behavior is characterized by distributions whose mean value is defined, so we expect most events to be around that value. However, the variance is not defined in the case of heavy-tailed distributions, so some events deviate significantly from the mean value: this is the central phenomenon described by *Large Deviations Theory*. The importance of accurately characterizing large deviations is evident in many research areas, with applications to environmental issues and beyond. Therefore, the universality that characterizes heavy-tailed distributions allows us to pursue the interdisciplinarity that complexity science aims to achieve.

The overall characteristics that arose from analyzing the public debate around the palm oil topic can be summarized as follows. The public debate followed a viral event on social media. Therefore, the interest was limited in time and engagement. The countries strongly involved in the debate are the ones that produce the oil and the ones where the debates arose, while the interest remained low in the rest of the world. This behavior suggests something about the way how viral events emerge; the geographical factor is found to be relevant. Similar reasoning can be done by looking at the limited number of *Opinion Drivers*, the most involved users are the companies producing palm oil and the environmental-involved organizations. The debate is, therefore, led by a few leading actors, slightly touching the general public.

In contrast, the debate around other vegetable oils follows a different trend: the range of topics of interest is broader, and environmental issues are not central in the discussion. Healthcare and nutrition are treated differently than sustainability issues. These topics are emerging less virally, and the attitude is not as negative as in the palm oil case. These results show a substantial discrepancy in the marketing strategy carried on by the *Opinion Drivers* of the different vegetable oils. Further developments can arise from applying the pipeline created in this thesis to different sociopolitical issues.



## Appendix A

# Coconut and Olive: comparison among datasets

In the following, we include additional results from the analysis of the **Coconut** and **Olive** dataset. Respectively, the time evolution in the number of Tweets is compared (Section A.1), as well as the hashtags' statistics (Section A.2). Then we analyze the sentiment dynamics (Section A.3), and we express the results obtained by the fitting procedure of, respectively, the **IET** distribution (Section A.4) and the **CS** distribution (Section A.5).

### A.1 Evolution of the interest

In this section, we consider the time evolution in the number of Tweets for the **Coconut** and **Palm** datasets (Fig.A.1). The **Coconut** dataset presents a large peak around 2015-2016, followed by a significant decrease in interest in the topic. In contrast, the **Olive** shows a strong increase and a weaker decrease, characterized by two peaks of similar intensity around 2013 and 2020. **Coconut** behaves similarly to **Palm**, where there is a strong peak followed by a decrease, while **Olive** presents a more uniform behavior after a relaxation time.

One question that could arise is about the relationship between these particular topics' time evolution and the social network's general evolution. Do we evidence an increase in the Tweets about palm oil only because the social network was growing? To answer, as the Twitter API does not allow extracting the data about the total number of Tweets per day, we considered the following approach. We used as queries for the Data Collection process the 100 most common words in English<sup>1</sup>, and we requested the number of Tweets per day containing these words. Due to the search query constraints established by the Twitter API, only 81 words were considered valid queries (the API does not allow using stop-words as queries). We computed the relative yearly growth, i.e., the number of Tweets per year containing a given word, divided by the value obtained in a year chosen as a reference. Fig. A.2 shows the relative growth from the reference year 2017 for the three datasets and the average value for the most common words. We can see that the relative growth of the **Palm** and **Coconut** datasets is stronger than the average behavior of Twitter. In contrast, the **Olive** dataset does not grow significantly. This is another evidence of the growing interest in these topics on social media over the years.

---

<sup>1</sup>The most common words are: *take, for, who, would, will, do, like, they, day, them, his, want, get, all, how, even, say, our, come, two, have, think, about, up, no, make, out, which, of, back, that, some, see, most, than, go, because, any, new, your, year, well, with, I, also, so, use, other, these, can, know, into, not, could, after, him, on, give, her, be, there, way, now, she, only, good, one, look, over, their, what, just, then, people, us, he, time, as, when, first, work.*

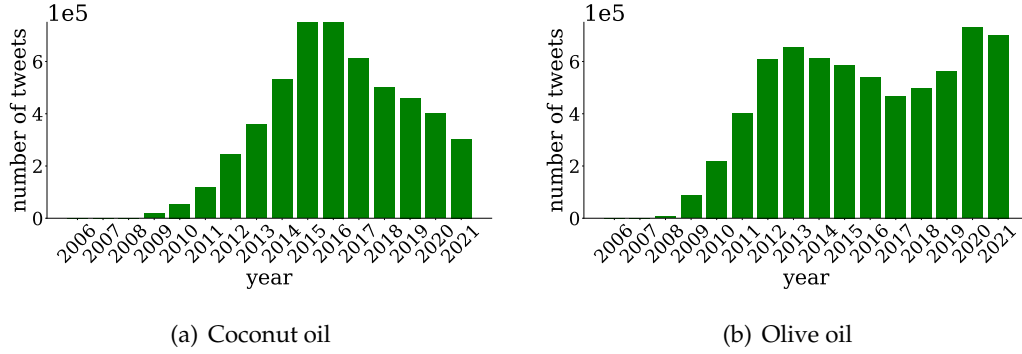


FIGURE A.1: Histogram representation of the number of Tweets per year containing the words *coconut oil* and *olive oil* respectively.

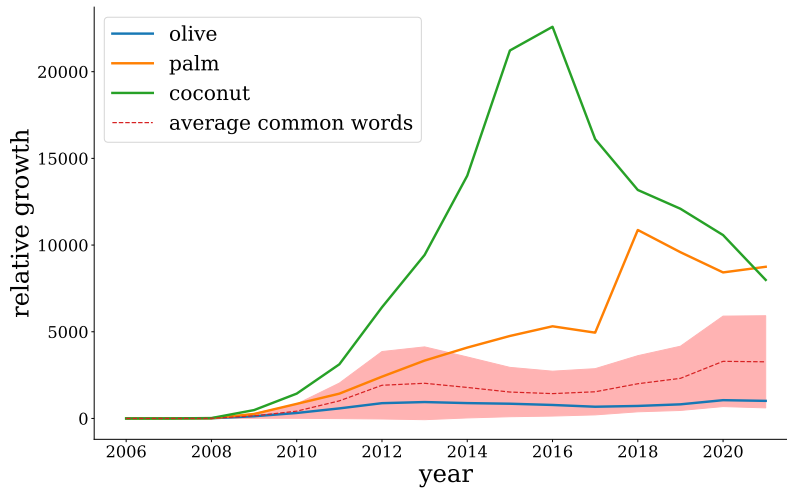


FIGURE A.2: The relative growth is calculated as the number of Tweets in a given year divided by the number of Tweets in a fixed reference year. Here the results for  $\text{ref\_year} = 2007$  are shown for the three datasets and an average of the relative growths of the most common words in English.

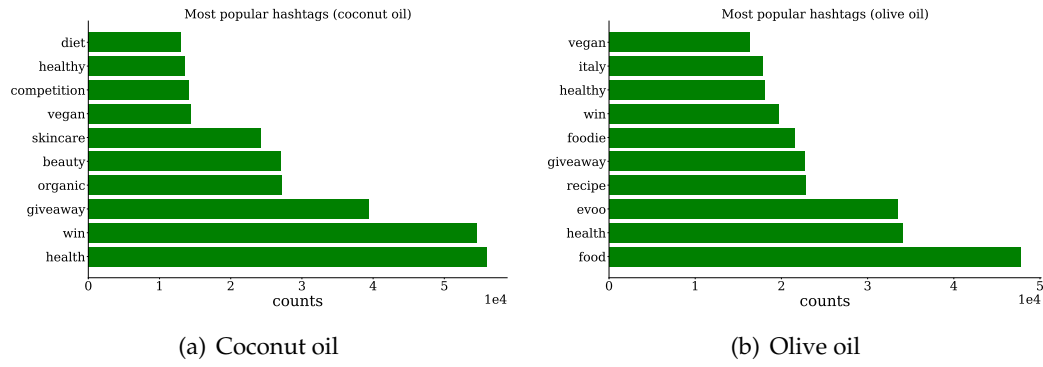


FIGURE A.3: Histogram representation of the most common hashtags for the two datasets.

## A.2 Most popular hashtags: topic detection

We extracted the hashtags from each Tweet in both datasets, following the definition given in Section 1.1. We decided to neglect the hashtags *coconutoil*, *coconut*, *oil* and *oliveoil*, *olive*, *oil* respectively, to analyze other aspects not directly related to the words used to create the dataset. As shown in Fig. A.3, the most popular hashtags are related to health and nutrition issues rather than sustainability, as in the **Palm** case. Contest-based hashtags (*giveaway*, *win*) are relevant in both datasets, while they were absent in the **Palm**. Moreover, the health and cosmetics topics are more significant in the **Coconut**, while the nutritional topic is prevalent in the **Olive**.

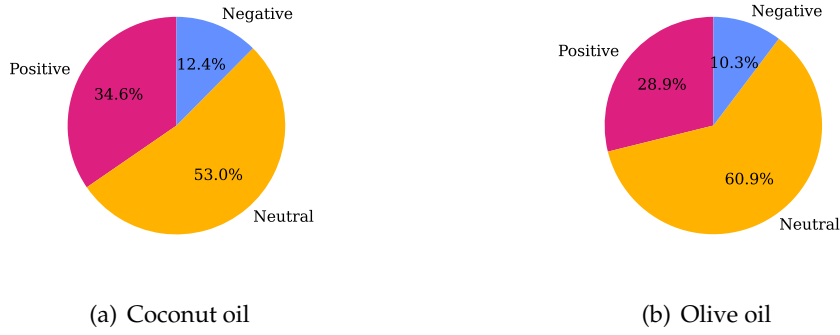


FIGURE A.4: Pie chart representation of the text classification task results.

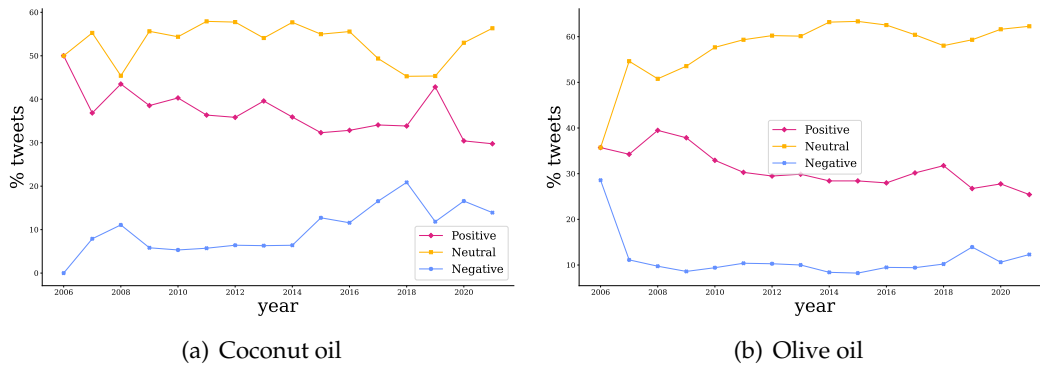


FIGURE A.5: Time evolution of sentiment. The plot represents the percentage of labeled Tweets published in a given year.

### A.3 Sentiment analysis

We performed a text classification task using the pre-trained model described in Chapter 3. Fig.A.4 shows the results for the **Coconut** and **Olive** datasets: as in the **Palm** case, *Neutral* Tweets are the most prevalent, showing here higher percentages than in the previous case. In contrast, the *Negative* sentiment shows a lower presence in these datasets, while the *Positive* one is more significant. This evidence could suggest that the **Palm** discussion evolved differently from the other two cases, where the attitude to the debates was less *Negative*. To further study the sentiment of these datasets, it is possible to show the time evolution of the sentiment. Fig.A.5 shows the percentage of Tweets labeled by a given sentiment per year. We notice that in both cases, the sentiments remain quite constant over time after an initial relaxation time. There is no strong evidence of the emergence of a *Negative* (or *Positive*) viral event; the datasets are found as quite uniform.

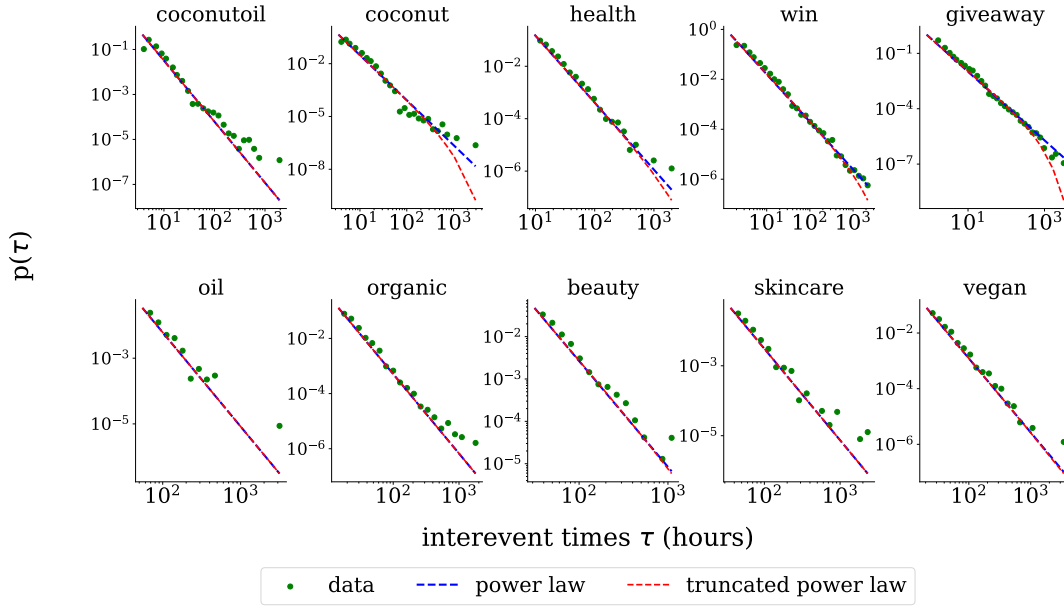


FIGURE A.6: **IET** distributions of the ten most used hashtags in the **Coconut** dataset. The red and blue dashed lines represent the fitted power laws and truncated power laws, while the scatter plots represent the data we are fitting.

## A.4 IET distribution

In this section, we express the results of the heavy-tailed fitting procedure (previously described in Section 5.1) for the **Coconut** and **Olive** datasets.

### A.4.1 Coconut dataset

As shown in Table A.1, some hashtags (*giveaway*, *win*, *beauty*, *coconut*, *health*, *vegan*) are better described by a truncated power-law, while a power-law better describes the remaining ones. In the latter case, the hashtags described power-law do not show a strong preference for it, as they present a very high p-value (i.e., bigger than 0.5). The null hypothesis is accepted at 95% of accuracy for all the hashtags. It has been done by comparing the KS distance with the critical values from the P-value table. All the scaling exponents  $\alpha$  are bigger than 1, and neglecting the contest-based hashtags (*giveaway*, *win*), they are all found in the *scale-free regime*  $2 < \alpha < 3$ . Here, the first moment is finite, while the higher moments diverge.

### A.4.2 Olive dataset

As shown in Table A.2, a few hashtags (*giveaway*, *win*, *oliveoil*) are better described by a truncated power-law, while a power-law better describes the remaining ones. Neglecting the contest-based hashtags (*giveaway*, *win*), the remaining ones do not show a strong preference for distribution, as they present a very high p-value (i.e., bigger than 0.5). The null hypothesis is accepted at 95 % of accuracy for all the hashtags. It has been done by comparing the KS distance with the critical values from the P-value table. All the scaling exponents  $\alpha$  are bigger than 1, few are found in the  $2 < \alpha < 3$  regime (i.e., *oliveoil*, *foodie*). It is possible to notice a stronger presence in the regime  $\alpha > 3$  compared to the other datasets. This means that the *Olive* dataset is less viral from the timing perspective.

Hashtag	Best distribution	$n$	$p$ -value	$\tau_{min}$	$D$	$\alpha$	$\sigma$
giveaway	truncated power-law	39,326	0.01	0.8	0.03	1.814	0.007
win	truncated power-law	54,546	0.01	1	0.02	1.87	0.01
beauty	truncated power-law	26,992	0.4	32	0.02	2.45	0.07
coconut	truncated power-law	90,807	0.04	3	0.03	2.52	0.02
health	truncated power-law	56,021	0.3	10	0.01	2.53	0.03
skincare	power-law	24,187	0.9	36	0.03	2.62	0.09
vegan	truncated power-law	14,377	0.8	21	0.01	2.66	0.05
coconutoil	power-law	139,095	0.9	4	0.02	2.71	0.03
oil	power-law	28,054	0.9	55	0.03	2.8	0.1
organic	power-law	27,117	0.9	15	0.01	2.88	0.05

TABLE A.1: Results of the fit for the **Coconut** dataset:  $n$  is the number of bins considered (the number of scatter points in Fig. A.6).  $p$  is the significance value of the comparison between candidate distributions,  $\tau_{min}$  is the lower bound of the scaling range [48],  $D$  is the Kolmogorov-Smirnov distance,  $\alpha$  is the scaling parameter,  $\sigma$  is the standard error.

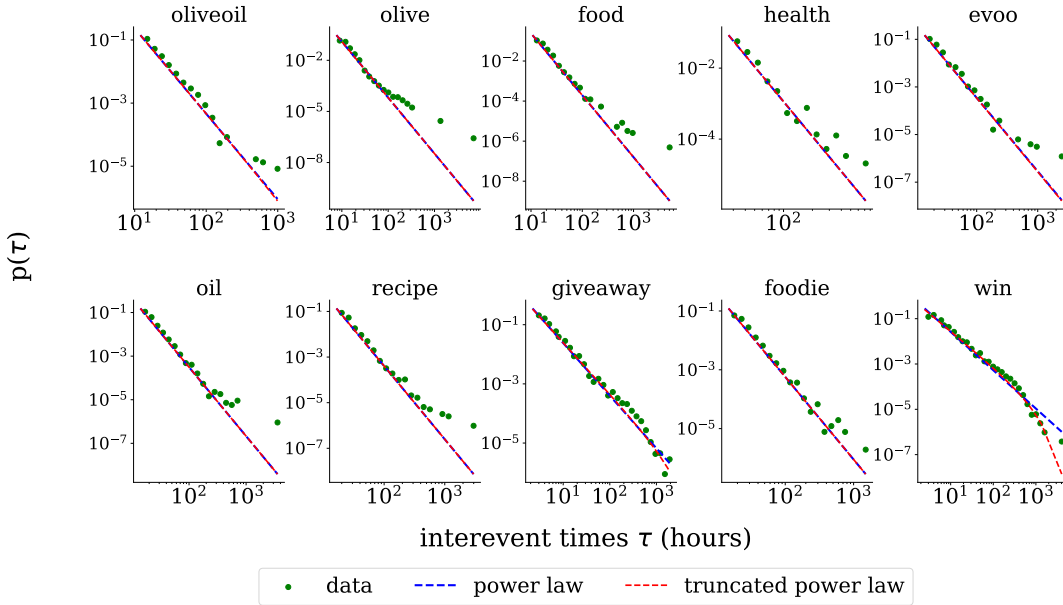


FIGURE A.7: **IET** distributions of the ten most used hashtags in the **Olive** dataset. The red and blue dashed lines represent the fitted power laws and truncated power laws, while the scatter plots represent the data we are fitting.

Hashtag	Best distribution	$n$	$p$ -value	$\tau_{min}$	$D$	$\alpha$	$\sigma$
win	truncated power-law	19,655	0.01	2	0.02	1.58	0.01
giveaway	truncated power-law	22,698	0.01	2	0.02	1.73	0.01
oliveoil	truncated power-law	89,288	0.8	12	0.02	2.69	0.07
foodie	power-law	21,510	0.9	15	0.02	2.81	0.04
evoo	power-law	33,493	0.8	14	0.02	3.08	0.05
food	power-law	47,727	0.9	11	0.02	3.09	0.05
oil	power-law	31,874	0.9	14	0.02	3.16	0.06
recipe	power-law	22,848	0.9	16	0.03	3.21	0.05
health	power-law	34,087	0.9	27	0.04	3.3	0.1
olive	power-law	58,901	0.9	8	0.02	3.28	0.05

TABLE A.2: Results of the fit for the **Olive** dataset:  $n$  is the number of bins considered (the number of scatter points in Fig. A.6),  $p$  is the significance value of the comparison between candidate distributions,  $\tau_{min}$  is the lower bound of the scaling range [48],  $D$  is the Kolmogorov-Smirnov distance,  $\alpha$  is the scaling parameter,  $\sigma$  is the standard error.

## A.5 CS distribution

In this section, we express the results of the heavy-tailed fitting procedure on the cascade size (similarly to what was done in Section 6.1) for the **Coconut** and **Olive** datasets.

### A.5.1 Coconut dataset

As shown in Table A.3, just the hashtag *win* is better described by a truncated power-law, while a power-law better describes the remaining ones. Neglecting the hashtag *win*, the remaining ones do not show a strong preference for distribution, as they present a very high  $p$ -value (i.e., bigger than 0.5). The null hypothesis is accepted at 95% of accuracy for all the hashtags. It has been done by comparing the KS distance with the critical values from the P-value table. All the scaling exponents  $\alpha$  are bigger than 1, one is found in the  $2 < \alpha < 3$  regime (i.e., *win*). It is possible to notice a stronger presence in the regime  $\alpha > 3$ , even with very high values for the scaling exponents (around 4-5). This means that the dataset is not viral at all from the magnitude perspective.

### A.5.2 Olive dataset

As shown in Table A.4, few hashtags (*giveaway*, *win*, *food*, *oliveoil*) are better described by a truncated power-law, while a power-law better describes the remaining ones. Neglecting the hashtags *giveaway*, *oliveoil*, the remaining ones do not show a strong preference for distribution, as they present a very high  $p$ -value (i.e., bigger than 0.5). The null hypothesis is accepted at 95% of accuracy for all the hashtags, and it has been done by comparing the KS distance with the critical values from the P-value table. All the scaling exponents  $\alpha$  are bigger than 1, and few are found in the  $2 < \alpha < 3$  regime (i.e., *giveaway*, *win*, *food*, *health*). It is possible to notice a stronger presence in the regime  $\alpha > 3$ , even with very high values for the scaling exponents (from 4 to even 7). This means that the dataset is not viral at all from the magnitude perspective.

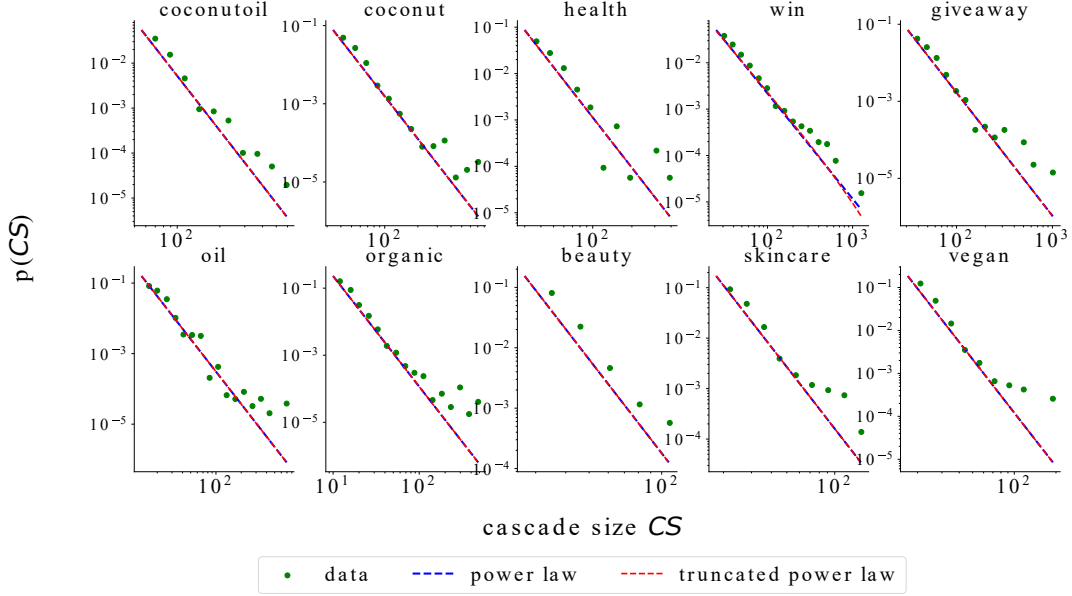


FIGURE A.8: Hashtags cascade distributions of the most common hashtags (**Coconut** dataset). The red and blue lines represent the fitted power-law and truncated power-law distributions. At the same time, the scatter plot is the empirical data (in logarithmic binning) that we are fitting.

Hashtag ( $h$ )	Best distribution	$n$	$p$ -value	$CS_{min}(h)$	$D$	$\alpha$	$\sigma$
win	truncated power-law	17	0.08	25	0.03	2.16	0.05
oil	power-law	17	0.9	13	0.03	3.06	0.09
giveaway	power-law	15	0.9	31	0.04	3.2	0.1
organic	power-law	16	0.9	10	0.02	3.30	0.09
health	power-law	11	0.9	30	0.02	3.6	0.1
coconut	power-law	13	0.9	34	0.03	3.62	0.09
coconutoil	power-law	10	0.9	56	0.04	4.0	0.1
skincare	power-law	9	0.9	18	0.03	4.1	0.2
vegan	power-law	10	0.9	17	0.03	4.1	0.3
beauty	power-law	5	0.9	27	0.05	5.2	0.4

TABLE A.3: Results of the fit for the **Coconut** dataset:  $n$  is the number of bins considered (the number of scatter points in Fig.A.8),  $p$  is the significance value of the comparison between candidate distributions,  $CS_{min}(h)$  is the lower bound of the scaling range [48],  $D$  is the Kolmogorov-Smirnov distance,  $\alpha$  is the scaling parameter,  $\sigma$  is the standard error.



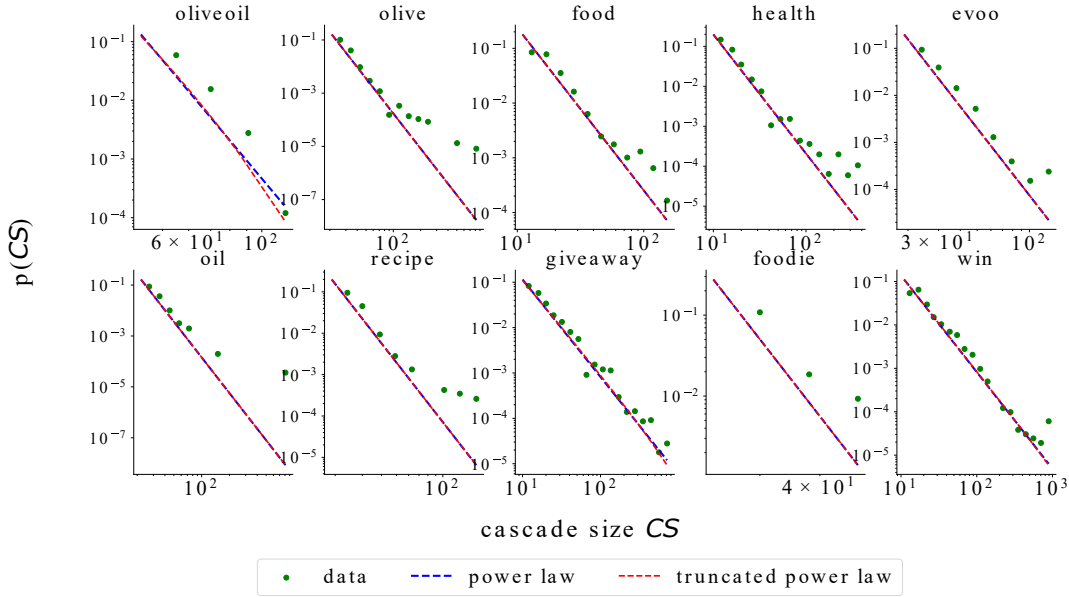


FIGURE A.9: Hashtags cascade distributions of the most common hashtags. The red and blue lines represent the fitted power-law and truncated power-law distributions. At the same time, the scatter plot is the empirical data (in logarithmic binning) that we are fitting.

Hashtag ( $h$ )	Best distribution	$n$	$p$ -value	$CS_{min}(h)$	$D$	$\alpha$	$\sigma$
giveaway	truncated power-law	18	0.1	10	0.02	2.07	0.05
win	truncated power-law	19	0.6	11	0.04	2.20	0.07
food	truncated power-law	11	0.9	11	0.03	2.97	0.06
health	power-law	15	0.9	10	0.03	2.98	0.07
oliveoil	truncated power-law	4	0.2	43	0.04	4.2	0.2
olive	power-law	15	0.9	21	0.03	4.4	0.1
oil	power-law	15	0.9	22	0.03	4.7	0.3
evoo	power-law	8	0.9	19	0.04	4.7	0.2
recipe	power-law	9	0.9	19	0.04	4.8	0.4
foodie	power-law	3	0.9	24	0.07	7.6	0.9

TABLE A.4: Results of the fit for the **Olive** dataset:  $n$  is the number of bins considered (the number of scatter points in Fig.A.9),  $p$  is the significance value of the comparison between candidate distributions,  $CS_{min}(h)$  is the lower bound of the scaling range [48],  $D$  is the Kolmogorov-Smirnov distance,  $\alpha$  is the scaling parameter,  $\sigma$  is the standard error.

# Bibliography

- [1] S. Boccaletti et al. “Complex networks: Structure and dynamics”. In: *Physics Reports* 424.4 (2006), pp. 175–308. ISSN: 0370-1573. DOI: <https://doi.org/10.1016/j.physrep.2005.10.009>. URL: <https://www.sciencedirect.com/science/article/pii/S037015730500462X>.
- [2] Andrea Cavagna et al. “Scale-free correlations in starling flocks”. In: *Proceedings of the National Academy of Sciences* 107.26 (2010), pp. 11865–11870. DOI: [10.1073/pnas.1005766107](https://doi.org/10.1073/pnas.1005766107). eprint: <https://www.pnas.org/doi/pdf/10.1073/pnas.1005766107>. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.1005766107>.
- [3] J. Buck. “Synchronous rhythmic flashing of fireflies. II”. In: *The Quarterly Review of Biology* 63.3 (Sept. 1988), pp. 265–289. ISSN: 0033-5770. DOI: [10.1086/415929](https://doi.org/10.1086/415929).
- [4] M. E. J. Newman. “Complex Systems: A Survey”. In: 79.8 (Aug. 2011), pp. 800–810. DOI: [10.1119/1.3590372](https://doi.org/10.1119/1.3590372). URL: <https://doi.org/10.1119/1.3590372>.
- [5] Mark Newman. *Networks*. Oxford university press, 2018.
- [6] Stanley Wasserman, Katherine Faust, et al. “Social network analysis: Methods and applications”. In: (1994).
- [7] Claudio Castellano, Santo Fortunato, and Vittorio Loreto. “Statistical physics of social dynamics”. In: *arXiv:0710.3256 [cond-mat, physics:physics]* (May 11, 2009). DOI: [10.1103/RevModPhys.81.591](https://doi.org/10.1103/RevModPhys.81.591). arXiv: [0710.3256](https://arxiv.org/abs/0710.3256). URL: <http://arxiv.org/abs/0710.3256> (visited on 04/29/2022).
- [8] Henrique Ferraz de Arruda et al. “Modelling how social network algorithms can influence opinion polarization”. In: *Information Sciences* 588 (Apr. 1, 2022), pp. 265–278. ISSN: 0020-0255. DOI: [10.1016/j.ins.2021.12.069](https://doi.org/10.1016/j.ins.2021.12.069). URL: <https://www.sciencedirect.com/science/article/pii/S0020025521012901> (visited on 05/30/2022).
- [9] Antonio F. Peralta, János Kertész, and Gerardo Iñiguez. “Opinion dynamics in social networks: From models to data”. In: *arXiv:2201.01322 [nlin, physics:physics]* (Jan. 10, 2022). arXiv: [2201.01322](https://arxiv.org/abs/2201.01322). URL: <http://arxiv.org/abs/2201.01322> (visited on 04/29/2022).
- [10] Sharad Goel et al. “The Structural Virality of Online Diffusion”. In: *Management Science* (July 22, 2015), p. 150722112809007. ISSN: 0025-1909, 1526-5501. DOI: [10.1287/mnsc.2015.2158](https://doi.org/10.1287/mnsc.2015.2158). URL: <http://pubsonline.informs.org/doi/10.1287/mnsc.2015.2158> (visited on 05/05/2022).
- [11] Albert-László Barabási. “The origin of bursts and heavy tails in human dynamics”. In: *Nature* 435.7039 (May 2005), pp. 207–211. ISSN: 0028-0836, 1476-4687. DOI: [10.1038/nature03459](https://doi.org/10.1038/nature03459). URL: <http://www.nature.com/articles/nature03459> (visited on 05/30/2022).

- [12] Maria Vincenza Chiriaco et al. "Palm oil's contribution to the United Nations sustainable development goals: outcomes of a review of socio-economic aspects". In: *Environmental Research Letters* 17.6 (June 2022), p. 063007. DOI: [10.1088/1748-9326/ac6e77](https://doi.org/10.1088/1748-9326/ac6e77). URL: <https://doi.org/10.1088/1748-9326/ac6e77>.
- [13] Krystof Obidzinski et al. "Environmental and social impacts of oil palm plantations and their implications for biofuel production in Indonesia". In: *Ecology and Society* 17.1 (2012).
- [14] Erik Meijaard et al. *Oil palm and biodiversity : a situation analysis by the IUCN Oil Palm Task Force*. June 2018. ISBN: 978-2-8317-1910-8 (PDF) 978-2-8317-1911-5 (print version). DOI: [10.2305/IUCN.CH.2018.11.en](https://doi.org/10.2305/IUCN.CH.2018.11.en).
- [15] John F. McCarthy. "Processes of inclusion and adverse incorporation: oil palm and agrarian change in Sumatra, Indonesia". In: *The Journal of Peasant Studies* 37.4 (2010). PMID: 20873030, pp. 821–850. DOI: [10.1080/03066150.2010.512460](https://doi.org/10.1080/03066150.2010.512460). eprint: <https://doi.org/10.1080/03066150.2010.512460>. URL: <https://doi.org/10.1080/03066150.2010.512460>.
- [16] Erik Meijaard et al. "Coconut oil, conservation and the conscientious consumer". en. In: *Curr. Biol.* 30.16 (Aug. 2020), pp. 3274–3275.
- [17] Mike Thelwall, Kevan Buckley, and Georgios Paltoglou. "Sentiment in Twitter events". In: *Journal of the American Society for Information Science and Technology* 62.2 (2011), pp. 406–418. DOI: <https://doi.org/10.1002/asi.21462>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/asi.21462>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/asi.21462>.
- [18] Cynthia Van Hee, Els Lefever, and Véronique Hoste. "SemEval-2018 Task 3: Irony Detection in English Tweets". In: Jan. 2018, pp. 39–50. DOI: [10.18653/v1/S18-1005](https://doi.org/10.18653/v1/S18-1005).
- [19] Saif Mohammad et al. "SemEval-2018 Task 1: Affect in Tweets". In: *Proceedings of The 12th International Workshop on Semantic Evaluation*. New Orleans, Louisiana: Association for Computational Linguistics, June 2018, pp. 1–17. DOI: [10.18653/v1/S18-1001](https://doi.org/10.18653/v1/S18-1001). URL: <https://aclanthology.org/S18-1001>.
- [20] Marcos Zampieri et al. "SemEval-2019 Task 6: Identifying and Categorizing Offensive Language in Social Media (OffensEval)". In: *Proceedings of the 13th International Workshop on Semantic Evaluation*. Minneapolis, Minnesota, USA: Association for Computational Linguistics, June 2019, pp. 75–86. DOI: [10.18653/v1/S19-2010](https://doi.org/10.18653/v1/S19-2010). URL: <https://aclanthology.org/S19-2010>.
- [21] Francesco Barbieri, Luis Espinosa-Anke, and Jose Camacho-Collados. "XLM-T: Multilingual Language Models in Twitter for Sentiment Analysis and Beyond". In: (2022).
- [22] Alexis Conneau et al. "Unsupervised Cross-lingual Representation Learning at Scale". In: *arXiv preprint arXiv:1911.02116* (2019).
- [23] Harry Bunt and Reinhard Muskens. *Computing Meaning*. Kluwer, 1999.
- [24] Christopher D. Manning and Hinrich Schütze. *Foundations of Statistical Natural Language Processing*. Cambridge, Massachusetts: The MIT Press, 1999. URL: <http://nlp.stanford.edu/fsnlp/>.
- [25] F. Pedregosa et al. "Scikit-learn: Machine Learning in Python". In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.

- [26] Susan T Dumais et al. "Latent semantic analysis". In: *Annu. Rev. Inf. Sci. Technol.* 38.1 (2004), pp. 188–230.
- [27] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc.", 2009.
- [28] George A. Miller. "WordNet: A Lexical Database for English". In: *Commun. ACM* 38.11 (1995), 39–41. ISSN: 0001-0782. DOI: [10.1145/219717.219748](https://doi.org/10.1145/219717.219748). URL: <https://doi.org/10.1145/219717.219748>.
- [29] Nathan Halko, Per-Gunnar Martinsson, and Joel A. Tropp. "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions". In: (2009). DOI: [10.48550/ARXIV.0909.4061](https://arxiv.org/abs/0909.4061). URL: <https://arxiv.org/abs/0909.4061>.
- [30] Albert-László Barabási and Márton Pósfai. *Network science*. Cambridge: Cambridge University Press, 2016. ISBN: 9781107076266 1107076269. URL: <http://barabasi.com/networksciencebook/>.
- [31] Stanley Milgram. "The Small-World Problem". In: *Psychology Today* 1.1 (1967), pp. 61–67.
- [32] Santo Fortunato and Darko Hric. "Community detection in networks: A user guide". In: *Physics Reports* 659 (Nov. 2016), pp. 1–44. DOI: [10.1016/j.physrep.2016.09.002](https://doi.org/10.1016/j.physrep.2016.09.002). URL: <https://doi.org/10.1016/j.physrep.2016.09.002>.
- [33] Réka Albert and Albert-László Barabási. "Statistical mechanics of complex networks". In: *Reviews of Modern Physics* 74.1 (Jan. 2002), pp. 47–97. DOI: [10.1103/revmodphys.74.47](https://doi.org/10.1103/revmodphys.74.47). URL: <https://doi.org/10.1103/revmodphys.74.47>.
- [34] Paul W. Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. "Stochastic blockmodels: First steps". In: *Social Networks* 5.2 (1983), pp. 109–137. ISSN: 0378-8733. DOI: [https://doi.org/10.1016/0378-8733\(83\)90021-7](https://doi.org/10.1016/0378-8733(83)90021-7). URL: <https://www.sciencedirect.com/science/article/pii/0378873383900217>.
- [35] M. E. J. Newman. "Modularity and community structure in networks". In: *Proceedings of the National Academy of Sciences* 103.23 (2006), pp. 8577–8582. DOI: [10.1073/pnas.0601602103](https://www.pnas.org/doi/pdf/10.1073/pnas.0601602103). eprint: <https://www.pnas.org/doi/pdf/10.1073/pnas.0601602103>. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.0601602103>.
- [36] Javier Borge-Holthoefer et al. "Emergence of consensus as a modular-to-nested transition in communication dynamics". In: *Scientific Reports* 7.1 (Jan. 30, 2017), p. 41673. ISSN: 2045-2322. DOI: [10.1038/srep41673](https://doi.org/10.1038/srep41673). URL: <https://doi.org/10.1038/srep41673>.
- [37] Vincent D Blondel et al. "Fast unfolding of communities in large networks". In: *Journal of Statistical Mechanics: Theory and Experiment* 2008.10 (Oct. 2008), P10008. DOI: [10.1088/1742-5468/2008/10/p10008](https://doi.org/10.1088/1742-5468/2008/10/p10008). URL: <https://doi.org/10.1088/1742-5468/2008/10/p10008>.
- [38] Albert-László Barabási and Réka Albert. "Emergence of Scaling in Random Networks". In: *Science* 286.5439 (1999), pp. 509–512. DOI: [10.1126/science.286.5439.509](https://www.science.org/doi/pdf/10.1126/science.286.5439.509). eprint: <https://www.science.org/doi/pdf/10.1126/science.286.5439.509>. URL: <https://www.science.org/doi/abs/10.1126/science.286.5439.509>.
- [39] Paul Erdos, Alfréd Rényi, et al. "On the evolution of random graphs". In: *Publ. Math. Inst. Hung. Acad. Sci* 5.1 (1960), pp. 17–60.

- [40] Ugo Bastolla et al. "The architecture of mutualistic networks minimizes competition and increases biodiversity". In: *Nature* 458.7241 (Apr. 1, 2009), pp. 1018–1020. ISSN: 1476-4687. DOI: [10.1038/nature07950](https://doi.org/10.1038/nature07950). URL: <https://doi.org/10.1038/nature07950>.
- [41] Elohim Reis, Aming Li, and Naoki Masuda. "Generative models of simultaneously heavy-tailed distributions of interevent times on nodes and edges". In: *Physical Review E* 102 (Nov. 2020). DOI: [10.1103/PhysRevE.102.052303](https://doi.org/10.1103/PhysRevE.102.052303).
- [42] Frank A Haight. *Handbook of the Poisson distribution*. Tech. rep. 1967.
- [43] Alexei Vazquez et al. "Impact of Non-Poissonian Activity Patterns on Spreading Processes". In: *Phys. Rev. Lett.* 98 (15 Apr. 2007), p. 158702. DOI: [10.1103/PhysRevLett.98.158702](https://link.aps.org/doi/10.1103/PhysRevLett.98.158702). URL: <https://link.aps.org/doi/10.1103/PhysRevLett.98.158702>.
- [44] Aaron Clauset, Cosma Rohilla Shalizi, and M. E. J. Newman. "Power-Law Distributions in Empirical Data". In: *SIAM Review* 51.4 (Nov. 2009), pp. 661–703. DOI: [10.1137/070710111](https://doi.org/10.1137/070710111). URL: <https://doi.org/10.1137/070710111>.
- [45] Guido Caldarelli, Matteo Marsili, and Y.-C. Zhang. "A prototype model of stock exchange". In: *Europhys. Lett.* 40.5 (1997), pp. 479–484. DOI: [10.1209/epl/i1997-00491-5](https://doi.org/10.1209/epl/i1997-00491-5). URL: <https://doi.org/10.1209/epl/i1997-00491-5>.
- [46] Z. Dezsö et al. "Dynamics of information access on the web". In: *Phys. Rev. E* 73 (6 June 2006), p. 066132. DOI: [10.1103/PhysRevE.73.066132](https://doi.org/10.1103/PhysRevE.73.066132). URL: <https://link.aps.org/doi/10.1103/PhysRevE.73.066132>.
- [47] Ye Wu et al. "Human comment dynamics in on-line social systems". In: *Physica A: Statistical Mechanics and its Applications* 389.24 (2010), pp. 5832–5837. ISSN: 0378-4371. DOI: <https://doi.org/10.1016/j.physa.2010.08.049>. URL: <https://www.sciencedirect.com/science/article/pii/S0378437110007521>.
- [48] Jeff Alstott, Ed Bullmore, and Dietmar Plenz. "powerlaw: A Python Package for Analysis of Heavy-Tailed Distributions". In: *PLOS ONE* 9.1 (Jan. 2014), pp. 1–11. DOI: [10.1371/journal.pone.0085777](https://doi.org/10.1371/journal.pone.0085777). URL: <https://doi.org/10.1371/journal.pone.0085777>.
- [49] Max C. Watson. "Time maps: A tool for visualizing many discrete events across multiple timescales". In: *2015 IEEE International Conference on Big Data (Big Data)*. 2015, pp. 793–800. DOI: [10.1109/BigData.2015.7363824](https://doi.org/10.1109/BigData.2015.7363824).
- [50] John M. Beggs and Dietmar Plenz. "Neuronal Avalanches in Neocortical Circuits". In: *Journal of Neuroscience* 23.35 (2003), pp. 11167–11177. ISSN: 0270-6474. DOI: [10.1523/JNEUROSCI.23-35-11167.2003](https://doi.org/10.1523/JNEUROSCI.23-35-11167.2003). eprint: <https://www.jneurosci.org/content/23/35/11167.full.pdf>. URL: <https://www.jneurosci.org/content/23/35/11167>.
- [51] Daniel S. Fisher et al. "Statistics of Earthquakes in Simple Models of Heterogeneous Faults". In: *Phys. Rev. Lett.* 78 (25 June 1997), pp. 4885–4888. DOI: [10.1103/PhysRevLett.78.4885](https://doi.org/10.1103/PhysRevLett.78.4885). URL: <https://link.aps.org/doi/10.1103/PhysRevLett.78.4885>.
- [52] Jordi Baró et al. "Statistical Similarity between the Compression of a Porous Material and Earthquakes". In: *Phys. Rev. Lett.* 110 (8 Feb. 2013), p. 088702. DOI: [10.1103/PhysRevLett.110.088702](https://doi.org/10.1103/PhysRevLett.110.088702). URL: <https://link.aps.org/doi/10.1103/PhysRevLett.110.088702>.

- [53] James P. Sethna et al. "Hysteresis and hierarchies: Dynamics of disorder-driven first-order phase transformations". In: *Phys. Rev. Lett.* 70 (21 May 1993), pp. 3347–3350. DOI: [10.1103/PhysRevLett.70.3347](https://doi.org/10.1103/PhysRevLett.70.3347). URL: <https://link.aps.org/doi/10.1103/PhysRevLett.70.3347>.
- [54] Oscar A. Pinto and Miguel A. Muñoz. "Quasi-Neutral Theory of Epidemic Outbreaks". In: *PLOS ONE* 6.7 (July 2011), pp. 1–7. DOI: [10.1371/journal.pone.0021946](https://doi.org/10.1371/journal.pone.0021946). URL: <https://doi.org/10.1371/journal.pone.0021946>.
- [55] Javier Borge-Holthoefer et al. "Cascading behaviour in complex socio-technical networks". In: *Journal of Complex Networks* 1.1 (Apr. 2013), pp. 3–24. ISSN: 2051-1310. DOI: [10.1093/comnet/cnt006](https://doi.org/10.1093/comnet/cnt006). eprint: <https://academic.oup.com/comnet/article-pdf/1/1/3/1370358/cnt006.pdf>. URL: <https://doi.org/10.1093/comnet/cnt006>.
- [56] James P. Gleeson and Rick Durrett. "Temporal profiles of avalanches on networks". In: *Nature Communications* 8.1 (Oct. 2017). DOI: [10.1038/s41467-017-01212-0](https://doi.org/10.1038/s41467-017-01212-0). URL: <https://doi.org/10.1038/s41467-017-01212-0>.
- [57] Duncan J. Watts. "A simple model of global cascades on random networks". In: *Proceedings of the National Academy of Sciences* 99.9 (2002), pp. 5766–5771. DOI: [10.1073/pnas.082090499](https://doi.org/10.1073/pnas.082090499). eprint: <https://www.pnas.org/doi/pdf/10.1073/pnas.082090499>. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.082090499>.
- [58] James P. Gleeson et al. "Effects of Network Structure, Competition and Memory Time on Social Spreading Phenomena". In: *Physical Review X* 6.2 (May 2016). DOI: [10.1103/PhysRevX.6.021019](https://doi.org/10.1103/PhysRevX.6.021019). URL: <https://doi.org/10.1103/PhysRevX.6.021019>.

## Funding

The author acknowledges the financial support of Soremartec S.A. and Soremartec Italia, Ferrero Group. The funders had no role in study design, data collection, and analysis, decision to publish, or preparation of the manuscript.