# JADT 2024
# Mots comptés, textes déchiffrés

Tome 1

Anne Dister et Dominique Longrée (éd.)

**UCLouvain**

Couverture : Cédrick Fairon

# Table des matières

# Decoding Disinformation: A Comprehensive Analysis of Fake News

Dario Sacco[1], Massimo Aria[2], Sara Moccia[3]

[1]University of Naples Federico II – dario.sacco@unina.it

[2] University of Naples Federico II – massimo.aria@unina.it

[3] University of Naples Federico II – sara.moccia@unina.it

## Abstract

In the contemporary digital landscape, the proliferation of inaccurate or manipulated news and information has attained a disconcerting magnitude. The dissemination of such misinformation has prompted an examination of causal elements, notably including the concept of citizen-limited rationality. This term denotes the cognitive constraints imposed by a dearth of information, cognitive limitations, and temporal constraints. This study is bifurcated into two primary objectives. Firstly, it aims to delineate the characteristics of fake news through statistical analysis. Secondly, it seeks to discern determinants influencing the propagation of misinformation regarding individuals' awareness of fake news. To accomplish these aims, a comprehensive two-step methodology was implemented. In the initial phase, unstructured texts were extracted from the web pages of prominent fact-checking sites utilizing web scraping techniques for the investigation of fake news. The second phase involved the identification and measurement of determinants in the dissemination of misinformation through structural equation models. Employing a topic-modeling model, the investigation probed whether fake news exhibits recurrent associations of words forming distinct topics. The subsequent analysis delved into the domains and degrees of misinformation, elucidating factors contributing to its dissemination. Common topics within fake news served a dual purpose: firstly, to gauge a latent variable about the extent of misinformation, utilizing topics to chart recurring themes in fake news upon which respondents were invited to articulate their perspectives regarding veracity or falsity; and secondly, to scrutinize, via path analysis with partial least squares structural equation modeling estimation, the variables deemed determinants of the misinformation phenomenon.

**Keywords:** Fake News, Disinformation, Fact-Cheking, Topic Modeling, PLS-SEM.

## 1. Introduction

In the vast landscape of issues surrounding communication, one phenomenon that emerges with particular relevance is fake news, which is the subject of in-depth study in the context of sustainable communication. The breadth of the field of responsible communication is well known (Wang et al., 2019), but here we focus on the study of false information, selecting it because of its pervasive spread and the significant impacts it generates.

Analysis of this phenomenon reveals the presence of crucial causal factors, including the so-called "bounded rationality" of citizens, a cognitive limitation brought about by limited information, cognitive, and time resources (Veltri, 2018; Piazza and Croce, 2019). Another relevant element discussed in the literature is "homophily" (Quattrociocchi and Vicini, 2016), a phenomenon whereby individuals tend to replicate and reinforce their behaviors through identification with similar individuals. This produces a systematic vulnerability of evaluations: citizens may fall into errors of judgment as they select information based on what confirms their preexisting beliefs or what is considered reliable because it is supported by many people.

The present study was carried out with a twofold analysis of misinformation: descriptive analysis of fake news and analysis of factors influencing the spread of misinformation. The

work was thus structured in two well-defined phases: the first phase focused on the extraction of text from sources performing fact-checking activities. This approach taken involved the systematic collection of unstructured data, focusing specifically on fact-checking articles, which following a rigorous process of verifying the facts stated and the data used within a given text attributable to fake news, an article is drafted that reports the headline and content of the fake news commented on by experts who have carried out the fact-checking activity. This method of analysis allowed Topic Modeling to be applied to identify the presence of recurring themes in fake news. The next stage involved the use of a survey to investigate the degree of misinformation, as well as the factors contributing to the spread of misinformation. Common topics identified in the topics were used to assess the degree of misinformation, inviting respondents to express opinions on the truthfulness or falsity of these topics. At the same time, these topics were coded to reflect specific themes that act as determinants of misinformation.

## 2. Theoretical and social framework

The research focuses on the analysis of a peculiar manifestation of disinformation, namely so-called fake news, selected as an object of study in the broader context of sustainable communication. The notion of sustainability, which is wide in scope and applicability in numerous fields, is circumscribed here to the social sphere, with a specific focus on the domain of communication. For example, the UK document on "Sustainable Communities" (Odpm, 2003) directly interrogates such highly relevant concepts as the quality of life, social inclusion and security, planning geared toward collective well-being, equality of opportunity, and ensuring quality services for the entire population. The sustainable communication discourse finds further enrichment in the work of Vinthagen (2013). His analysis is intertwined with our investigation of fake news, serving as a fundamental conceptual bridge between the dynamics of misinformation and broader issues of equity and social resilience. The elements of social crises of the present as polarization, increasing urban poverty, conflict, urban violence, and terrorism. Demands arise for the revision of traditional methods of sustainability, including sustainable forms of communication that not only inform but also promote accountability, equity, and resilience. Fake news poses a direct threat to these principles in that it can distort reality, create information inequalities, and undermine trust in communication but also governance systems that are essential to a democratic society.

Studies by Beck (1992), Eizenberg and Jabareen (2017), and Jabareen (2015) highlight risk as a fundamental organizing principle in contemporary societies. Risk takes a prominent role in the concept of sustainability, which must be rethought to ensure the long-term security of human and nonhuman communities. Ziemann's (2007) definition of communication as the foundation of social organization and as the vehicle through which "the internal is externalized" becomes particularly pertinent. The communication process is not only a manifestation of individual consciousness, but also an expression of social attitudes, intentions, thoughts, and values. Sustainable communication, therefore, takes the form of a commitment to a critical consciousness regarding the problems and risks inherent in the relationships between humans and those with the environment, values, and social norms. With this in mind, the battle against fake news becomes an essential element of sustainable communication: fake news not only distorts reality and fuels social injustice, but also hinders the formation of a critical and informed collective consciousness about the real risks facing our societies.

The works of Veltri G. (2018) and Wang et al., (2019) share a common interest in the dynamics of disinformation dissemination in social media and its implications on society. For Veltri, social media and the echo chamber phenomenon affect individuals' rationality, limiting their

ability to process reliable information. The digital environment fosters the formation of information bubbles in which fake news can proliferate undisturbed. Wang et al. focus on the specific context of health-related misinformation in social media, highlighting the potentially serious consequences that fake news can have on public health. Similar to Veltri, they identify the psychological characteristics of individuals and the dynamics of social networks as the main factors that facilitate the spread of misinformation. Both studies emphasize the importance of sustainable communication as an antidote to misinformation. In the context of fake news, sustainable communication involves enhancing the critical capacity of citizens and promoting accurate and verifiable information.

## 3. Research Proposal and Methodology

### *3.1. Research Proposal*

In today's context of the increasing prevalence of fake news, the need to understand and analyze the dynamics underlying its proliferation becomes imperative. To address these issues, the analysis was divided into two aspects: descriptive analysis of fake news (Jain and Kasbe, 2018) and analysis of the determinants in the dissemination of disinformation about the degree of knowledge of fake news. A combined two-step strategy was used to conduct these analyses: 1) concerning fake news, unstructured texts were extracted from web pages using fact-checking sites through the development of web scraping programs; 2) based on the topics extracted from the data analysis collected in step 1, the determinants in the dissemination of misinformation were studied and a survey was conducted by administering questionnaires. Facta.news database was used to extrapolate fact-ckeeping articles. Facta is a fact-checking site, that conducts timely verification of facts and sources, is a member of the International Fact-Checking Network (Ifcn), a signatory of its Code of Principles, and is a registered newspaper with the Court of Milan (No. 56 of March 8, 2021). Once the data were extracted and preprocessed, analysis was carried out by topic modeling. The objective was to determine the existence of patterns of frequently associated words that form distinct topics within fake news, following reference studies in the field (Englmeier, 2021; Hosseini et al., 2023; Choi and Ko, 2021). Topic modeling was selected as an analysis tool because of its ability to reveal the main themes present in large textual datasets. This methodology takes advantage of recent developments in information technology, enabling automatic or semiautomatic processing of texts to extract meaning from natural language. This approach is crucial for overcoming the limitations imposed by the lack of prior knowledge and the need for a significant time commitment that would be required by manual analysis. (Misuraca M. and Spano M., 2020). Next, the analysis continues through the use of a survey to investigate the areas and levels of misinformation and the factors that contribute to the spread of misinformation. The topics common to fake news were used for a twofold purpose: 1) to measure a latent variable related to the degree of misinformation; in this regard, topics will be used to map recurring trends in fake news on which respondents will be asked to express their opinions regarding the truthfulness or falsity of the same; and 2) to trace the topics back to a particular theme that reflects a specific determinant of misinformation itself. This step will be carried out Partial Least Squares Structural Equation Modeling (PLS-SEM) estimation algorithm (Ciavolino et al., 2022).

### 3.2. Methodology and work-flow

### 3.2.1. Text collection and and preprocessing

Information found on the Web comprises a diverse range of data, including quantitative and qualitative structural, semi-structural, and unstructured data. Such data manifest themselves in various forms, such as Web pages, HTML tables, Web databases, tweets, blog posts, and video content (Watson, 2014). The acquisition and organization of this type of data, with an emphasis on text, was accomplished through the use of a web scraping program implemented through the Python programming language. To perform the web scraping operation, two specific libraries were adopted, namely *"BeautifulSoup"* and *"selenium"*. The fact-checking articles included in the dataset were selected based on a period between December 1, 2020 and December 1, 2023. The selection of these articles was conducted by considering their subject matter, with a particular preference for topical issues. Subsequently, the following were excluded from the dataset: duplicate articles, articles containing only commentary, and articles that incorporated misleading or fraudulent information.

The process began by collecting 1450 fake articles through web scraping techniques. The first step, related to the pretreatment of the collected texts, was based on a Text Normalization process, which involved the removal of elements such as URLs, emoji, IP addresses, numbers, and dates. Next, the text was subjected to an annotation process that included, tokenization, lemmatization, and tagging to universal parts of speech (UPOS). This annotation was accomplished through the use of pre-trained models derived from Universal Dependencies (Straka M. and Straková J., 2017). The programming language R was used to perform these operations with the help of the libraries "*Udpipe"*, "*tidyverse"*, and *"tall"*.

The lemmatized dataset presented a total of 674,478 rows, corresponding to the number of tokens in the texts. The next step involved studying the distribution of conditional lemmas concerning UPOS, focusing particularly on ADJ, PROPN, NOUN, and VERB. At this stage, a customized stopword list was developed. Multi-words (n-grams) were created and only the UPOSs of interest, namely ADJ, PROPN, NOUN, VERB, and the Multi-Worlds, were dissected. Unwanted words were removed through the use of the previously created stop words list, while hapaxes, i.e., terms present only once in the corpus, were removed. This process produced a dataset consisting of 158,419 rows.

### 3.2.2. Topic Modeling

We used the Latent Dirichlet Allocation (LDA) model to identify topics. In the field of NLP and under the "bag-of-words" assumption, LDA is a generative probabilistic topic model proposed by Blei et al (2003), aiming to uncover latent topics within a collection of documents. LDA views documents as mixtures of topics, with each topic represented as a distribution over a fixed vocabulary. The model employs conjugate Dirichlet priors to generate topic and document distributions, where each document is characterized by topic proportions and each topic by word distributions.

The optimal number of Topics for LDA analysis was determined using two different approaches developed by Cao et al. (2009) and Deveaud et al. (2014) Figure 2, respectively. Following the method proposed by Cao et al. (2009), the optimal number of topics is selected by minimizing the average cosine similarity between topics, while the metric of Deveaud et al. (2014), by maximizing the average distance between topic distribution pairs estimated by LDA.

*3.2.3. Path-Analysis using PLS-SEM*

After careful topic analysis and study, ten fake news items were selected using the highest theta parameters (θ; document topic distributions) of the LDA model, finalized to a battery of items to be submitted in the interview. Among these fake news items was a particular statement that constituted the false part. Seven of these statements were submitted to the respondents as selected. For the remaining three, the corresponding true news item was entered. The identification of the true news was possible because many fake news stories are based on events that happened. Both real news and fake news were submitted to the respondents to keep the reader's concentration high. The ten selected news items were used as a battery of items to measure a latent variable designed to gauge the degree of misinformation among survey participants. These items were coded using a six-level scaling method (1 to 6) in which participants made a judgment about the truthfulness or falsity of the information. With a label of 1, participants considered the published news completely false, while with a value of 6, they believed it to be true. An even-step scale was adopted, excluding the middle step representing neutrality, to obtain a polarized opinion from the participants.

Accordingly, a pilot survey was conducted by administering a questionnaire to a sample of students at Federico II University, with a request to distribute the questionnaire further, given the limited time for the survey. A total of 235 questionnaires were collected. In addition to the construct related to misinformation, seven causal relationships believed to be determinants of misinformation were hypothesized. The variables are "Status Quo," "Altruism," "Information Sharing," "News Overload," "Passive Information," "Vaccine Information," and "Social Media Information." The constructs and questionnaires are validated in the literature (Apuke O. D. and Omar B. 2021); Sterie, L. G., et al., 2023). The choice to use the construct "Vaccine Information" is due to the vast majority of fake news in the social health field. These variables take the role of exogenous in the PLS-SEM model and were assessed through a five-step item battery. The variable "Disinformation" is considered an endogenous variable in the model. The PLS-SEM estimation algorithm produced the following results, which are shown in the table along with the respective significance estimated by bootstrapping.

## 4. Text analysis

The word cloud (Figure 1) was generated using the extracted lemmas, allowing for a concise and representative visualization of the main concepts present in the analyzed corpus.



*Figure 1: World Cloud of Fake News*

## 4.1. Topic Modeling resoults



*Figure 2: K Choice using Cao-Juan and Daveaud metrics.*

By increasing the lemma number, the two metrics suggested using a topic number between 14 and 16. We evaluated the various topics described by the topics and opted for a lemma number of 150 and a topic number of 5. In general, a topic model with an excessive number of topics will have many overlaps and words or synonyms within a topic, making it difficult to distinguish between different topics or concepts. The choice of setting K equal to 5 was also derived from two reasons, the first related to the purpose of the research, to search for macro-topics that fake news deals with, and the second related to the lexical complexity of the texts, which by observing the Types/Tokens Ratio (TTR) index equal to 12.75 we observe that the texts use a relatively small number of terms.

The following extracted Topics are given.

Topic 1 (Figure 3) appears to cover topics related to politics and public order, with a focus on events involving public figures and government institutions. Related articles on the topic discuss various events involving political figures and government institutions, including controversial arrests by police, protests at public events involving politicians such as Mario Draghi, and false or unsubstantiated news stories regarding alleged misconduct by politicians such as Nancy Pelosi and Ursula von der Leyen. In addition, the salaries of senior public officials are discussed, and political protests and activism are also discussed, with instances of demonstrations against institutions such as the European Central Bank, highlighting the role of opposition and criticism in the area of economic and financial policies. We define this topic as "Chronicle and Politics".

Topic 2 (Figure 4) includes a range of topics revolving around issues affecting Italy and international policy contexts, with a focus on public health, the food industry, international organizations, and social issues. In the articles related to it, there are discussions regarding authorizations by the Italian Ministry of Health. There are references regarding Italian companies and alleged collaborations with the government. Food additives are discussed, with misleading information regarding their dangerousness and regulation in Europe, with reference also to the McDonald's fast food chain. Food flavorings, citric acid, and food in general are discussed. Much mention is also made of WHO (World Health Organization) and its president Tedros Adhanom Ghebreyesus.

## Chronicle and Politics

| Word | Beta Probability (β) |
|------|------|
| woman | 0,074 |
| colleague | 0,061 |
| public | 0,055 |
| government | 0,053 |
| president | 0,053 |
| state | 0,041 |
| law | 0,040 |
| spokesman | 0,040 |
| police | 0,035 |
| euros | 0,034 |
| body | 0,032 |
| press | 0,029 |
| politician | 0,027 |
| protest | 0,026 |

*Figure 3: Topic 1*

## Institutions, Agreements, Government and Markets

| Word | Beta Probability (β) |
|------|------|
| Italian | 0,091 |
| country | 0,055 |
| Italy | 0,054 |
| organization | 0,047 |
| European | 0,041 |
| milion | 0,035 |
| firm | 0,035 |
| wordwide | 0,033 |
| health | 0,033 |
| WHO | 0,030 |
| company | 0,030 |
| circular | 0,029 |
| pandemic | 0,029 |
| ministry | 0,029 |

*Figure 4: Topic 2*

## Vaccine and Covid 19

| Word | Beta Probability (β) |
|------|------|
| vaccine | 0,235 |
| Covid | 0,198 |
| vaccination | 0,055 |
| doctor | 0,047 |
| virus | 0,046 |
| disease | 0,040 |
| Pfizer | 0,037 |
| system | 0,031 |
| risk | 0,028 |
| coronavirus | 0,027 |
| Sars | 0,027 |
| patient | 0,026 |
| adverse | 0,025 |
| infection | 0,022 |

*Figure 5: Topic 3*

## Adverse Effects

| Word | Beta Probability (β) |
|------|------|
| data | 0,104 |
| death | 0,054 |
| vaccinate | 0,037 |
| decease | 0,034 |
| population | 0,033 |
| hospital | 0,032 |
| age | 0,029 |
| high | 0,028 |
| effect | 0,027 |
| increase | 0,027 |
| vaccinated | 0,026 |
| climatic | 0,024 |
| journal | 0,024 |
| level | 0,023 |

*Figure 6: Topic 4*

## Geopolitical Conflicts

| Word | Beta Probability (β) |
|------|------|
| Ukraine | 0,083 |
| United States | 0,072 |
| child | 0,068 |
| Ukrainian | 0,046 |
| city | 0,044 |
| Russian | 0,041 |
| American | 0,039 |
| military | 0,039 |
| artificial | 0,028 |
| flag | 0,028 |
| war | 0,028 |
| road | 0,027 |
| canal | 0,026 |
| strike | 0,026 |

*Figure 7: Topic 5*

Other topics are about Italy and the prices of many products in national and international markets. We define this topic as "Institutions, Agreements, Government and Markets".

Topic 3 (Figure 5) seems to focus mainly on the topic of COVID-19 vaccination and related issues. The role of physicians and the health care system in the administration of Pfizer and mRNA vaccines is discussed, with attention to potential risks and adverse effects related to them. This is a topic that relates directly to the management of the coronavirus pandemic, with an analysis of diseases caused by the virus. In addition, there appears to be a focus on the safety

and efficacy of vaccines, as well as the evaluation of cases of infection and patients who develop serious complications following vaccine administration. This topic reflects the relevance and topicality of the discussion on anti-Covid vaccines, fueling the complexity of the issues involved and marking the risks associated with it. We term this topic "Vaccine and Covid 19."

Topic 4 (Figure 6) seems to focus on data analysis of adverse health and environmental effects. Reports circulate related to deaths and vaccination with a range of related factors. Italian hospitals and data inherent in deaths are mentioned. Climate change is also mentioned as a possible influential variable of health and environmental risk. In addition, reference is made to the role of scientific journals in providing data and conclusions on this topic. The topic suggests the importance of understanding the complex interactions between vaccination, public health, environmental factors, and other elements to draw meaningful conclusions about mortality dynamics. Of course, fake news stories about adverse effects and deaths leave a lot of room for vaccine-related fake news. Let us call such a topic "Adverse Effects."

Topic 5 (Figure 7) appears to be about several elements related to a geopolitical situation involving Ukraine, the United States, and Russia, with a focus on military involvement and international tensions. Ukraine is mentioned along with the United States and Russia, highlighting the complexity of relations between them. Children are also mentioned, which would emphasize facts and situations arising from the conflict. The term "artificial" refers to technologies used in the context of war and related agreements between these conflicting countries. The mention of flags, wars, and attacks suggests an environment of tension and conflict. The presence of "Israel" indicates the underscoring of news proliferation concerning world conflicts. In summary, topic 5 seems to focus on conflict dynamics, international politics, and security, with particular reference to Ukraine, the United States, and Russia. We call this topic "Geopolitical Conflicts".

## 5. PLS – SEM resoult

| Exogenous Variables | Endogenous Variable | Bootstrap Mean | Bootstrap SD | T Stat. | P-value | 2,5% CI | 97,5% CI |
|---|---|---|---|---|---|---|---|
| Status quo | disinformation | 0,134 | 0,087 | 1,529 | 0,0638 | -0,038 | 0,317 |
| Altruism | disinformation | -0.025 | 0,082 | -0,524 | 0,3005 | -0,175 | 0,153 |
| Information sharing | disinformation | -0,059 | 0,088 | -0,73 | 0,233 | -0,2229 | 0,105 |
| News overload | disinformation | 0,165*** | 0,053 | 2,598 | 0,005 | 0,065 | 0,265 |
| Passive information | disinformation | 0,134** | 0,074 | 1,821 | 0,035 | -0,017 | 0,274 |
| Social media | disinformation | 0,291*** | 0,095 | 3.251 | 0,0007 | 0,098 | 0,462 |
| vaccines | disinformation | -0,116** | 0,058 | -1,687 | 0,0464 | -0,23 | -0,007 |

*Table 1:PLS-SEM model (\*\*\* indicates a 1% level of statistical significance, while \*\* indicates a 5% level of significance)*

## 6. Discussion

The topic analysis showed that fake news constitutes a pervasive and insidious phenomenon, likely to affect a wide range of topics. Although it was possible to identify macro-topics related to the disinformation produced, the diversification of articles producing the topics appears to be extremely broad, with some news stories so subtly distorted from reality that it is difficult to discern veracity without careful verification. In political and public policy contexts, false

information can be used to spread non-truths about public figures and government institutions to damage their reputations or influence public opinion about specific political events. In the field of public health and the food industry, fake news can be conveyed to disseminate misinformation about health risks associated with certain foods and chemicals or to generate confusion about government policies, such as on nutrition, or the regulation of markets. This phenomenon may interfere with the market, changing expectations, or with the opinions of those exposed. Relative to anti-COVID vaccination, false information can fuel misinformation, raising doubts about their safety and efficacy and consequently contributing to the spread of conspiracy theories and vaccine rejection. In the context of adverse health and environmental effects, fake news can be exploited to generate panic or mistrust related to alleged risks associated with certain health claims or certain environmental phenomena. Finally, in the context of geopolitical conflicts, fake news can be employed to manipulate national and international public opinion by supporting specific political narratives or justifying military actions. In sum, fake news poses a threat to democracy and social welfare because it undermines trust in information and fact-based public debate. It is essential to promote greater critical awareness and improve media literacy skills to increase citizen empowerment and limit the harmful effect of false information on society.

The results of the PLS-SEM model show that being subjected to information overload from different and varied sources, passively receiving information, and being informed on platforms such as social media have a positive impact on misinformation. Thus, those who rely on social media or are exposed to an excessive amount of information, even passively, from various sources are more likely to be misinformed. In particular, social media is associated with the virality of unverified news and the manipulation of information, thus contributing to a distorted perception of facts. Information overload, on the other hand, generates cognitive overload that hinders the ability to process information critically, leading individuals to confusion.

Addressing this phenomenon requires a combination of factors such as individual education, regulation of digital platforms, and collective efforts to promote a culture of verification and critical thinking. Information regarding vaccines, contrary to expectations, hurts misinformation, as it enables users to recognize the importance and benefits of prevention.

Just as correct information enables citizens to orient themselves toward choices that are supported by scientific evidence and therefore safe, in the same way, techniques for biasing certain factors that are generative of explanatory constructs of the phenomenon allow the advantage of a reliable reading of reality for both planning and promoting virtuous behaviors that adhere to the most evident social needs.

## Acknowledgement

## Bibliographie

Apuke O. D. and Omar B. (2021). Fake news and COVID-19: modelling the predictors of fake news sharing among social media users. *Telematics and Informatics*, 56, 101475.

Beck U. (1992). *Risk Society: Towards a New Modernity.* London: Sage.

Blei D. M., Ng A. Y. and Jordan M. I. (2003). Latent dirichlet allocation. *Journal of machine Learning research*, 3 (Jan), 993-1022.

Cao J., Xia T., Li J., Zhang Y. and Tang S. (2009). A density-based method for adaptive LDA model selection. *Neurocomputing*, 72 (7-9), 1775-1781.

Choi H. and Ko Y. (2021, October). Using topic modeling and adversarial neural networks for fake news video detection. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2950-2954.

Ciavolino E., Aria M., Cheah J. H. and Roldán J. L. (2022). A tale of PLS structural equation modelling: episode I—a bibliometrix citation analysis. *Social Indicators Research*, 164, 1323-1348.

Deveaud R., SanJuan E. and Bellot P. (2014). Accurate and effective latent concept modeling for ad hoc information retrieval. *Document numérique*, 17 (1), 61-84.

Hosseini M., Sabet A. J., He S. and Aguiar D. (2023). Interpretable fake news detection with topic and deep variational models. *Online Social Networks and Media*, 36, 100249.

Eizenberg E. and Jabareen Y. (2017). Social sustainability: A new conceptual framework. *Sustainability*, 9 (1), 68.

Jain A. and Kasbe A. (2018, February). Fake news detection. In *2018 IEEE International Students' Conference on Electrical, Electronics and Computer Science (SCEECS)*, 1-5.

Jabareen Y. (2015). *The Risk City: Cities Countering Climate Change: Emerging Planning Theories and Practices around the World.* New York: Springer.

Misuraca M. and Spano M. (2020). Unsupervised analytic strategies to explore large document collections. In *Text Analytics: Advances and Challenges*. New-York: Springer, 17-28.

Office of the Deputee Prime Minister. (2003). *Sustainable communities: building for the future.* TSO, London. OECD, 2003 Social Issues in the Provision and Pricing of Water Services.

Quattrociocchi W. and Vicini A. (2016). *Misinformation: Guida alla società dell'informazione e della credulità*. Milano: FrancoAngeli.

Piazza T. and Croce M. (2019). Epistemologia delle fake news. *Sistemi intelligenti*, 31 (3), 439-468.

Straka M. and Straková J. (2017). Tokenizing, pos tagging, lemmatizing and parsing ud 2.0 with udpipe. In *Proceedings of the CoNLL 2017 shared task: Multilingual parsing from raw text to universal dependencies*, Vancouver, 88-99.

Veltri G. (2018). La tempesta perfetta: social media, fake news e la razionalità limitata del cittadino. *Media Education*, 9 (1), 36-56.

Vinthagen S. (2013). Ten theses on why we need a "Social Science Panel on Climate Change". *ACME: An International Journal for Critical Geographies*, 12 (1), 155-176.

Wang Y., McKee M., Torbica A. and Stuckler D. (2019). Systematic literature review on the spread of health-related misinformation on social media. *Social science & medicine*, 240, 112552.

Watson H. J. (2014). Tutorial: Big Data Analytics: Concepts, Technologies, and Applications. *Communications of the Association for Information Systems*, 34 (1), 1247-1268.

Ziemann A. (2007). Kommunikation der Nachhaltigkeit. Eine kommunikationstheoretische Fundierung. In G. Michelsen & J. Godemann (Eds), *Handbuch Nachhaltigkeitskommunikation. Grundlagen und Praxis.* Munich: Oekom, 123–133.

Sterie L. G., Sitar-Tăut D. A. and Mican D. (2023, May). Analyzing the Antecedents of Fake News Sharing in Online Social Networks. In *International Conference on Informatics in Economy*. Singapore, 149-158.