

SIRITVIS: Social Interaction Research Insights Topic Visualisation

Sagar Narwade¹, Gillian Kant², Benjamin Säfken¹, and Benjamin Leiding¹

¹ Technische Universität Clausthal, Clausthal-Zellerfeld, Germany ² Georg-August-Universität Göttingen, Göttingen, Germany

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

SIRITVIS is a powerful text analysis tool that has been carefully designed to analyse data from Reddit, X (Twitter), Instagram, and any external text data sources. It utilises the latest tools to automatically extract information from these sources, clean the data, and optimise topic modeling more efficiently. One feature of SIRITVIS is its ability to map social media posts on a global scale, connecting geographical locations with trending topics, thus providing valuable insights into the worldwide trends and conversations shaping our digital landscape. Additionally, it offers sentiment analysis on mapped data using the well-regarded VADER tool (Hutto & Gilbert, 2014). This software is a valuable resource for the scientific community, offering deep insights into public discussions on various social platforms. It can be used for a wide range of research purposes, including the analysis of recent discussions on global geopolitical issues. Installing SIRITVIS is straightforward, as it can be accomplished using the pip package manager. Comprehensive installation instructions can typically be found in the dedicated repository of the package.¹

Statement of Need

The surge of social media platforms has revolutionised the way we communicate, share information, and express viewpoints on various subjects. Among these platforms, Reddit, Instagram, and X (Twitter) stand out as notable sources of public discourse. Examining text data from these platforms offers valuable insights into public sentiments, preferences, and trending discussions, benefiting fields such as marketing, politics, and disaster management.

Dealing with the colossal volume of unstructured text data from social media is challenging due to its dynamic nature and sheer size. To address this challenge, we introduce SIRITVIS, a text analysis package designed to simplify the analysis of social media data. The package employs advanced neural topic models developed by AVITM (Srivastava & Sutton, 2017) and other popular topic models, including Latent Dirichlet Allocation (LDA) (Blei et al., 2003), Neural Latent Dirichlet Allocation (NeuralLDA) (Srivastava & Sutton, 2017), Product Latent Dirichlet Allocation (ProdLDA) (Srivastava & Sutton, 2017), and Contextualised Topic Models (CTM) (Bianchi et al., 2021). These models automatically identify and extract top topics in an unsupervised manner, enabling users to explore large text datasets and discover hidden patterns for meaningful insights.

SIRITVIS encompasses a comprehensive suite of functionalities designed for data extraction from social media platforms, including X (Twitter), Reddit, and Instagram, facilitated through the utilisation of respective application programming interfaces (APIs). This extraction process is followed by a meticulous data preprocessing phase, wherein valuable information is extracted

¹<https://github.com/CodeEagle22/SIRITVIS/>

41 from raw data, and superfluous elements are expunged, utilising advanced natural language
42 processing (NLP) techniques. The resultant processed data is subsequently employed for topic
43 modeling, with the flexibility for users to fine-tune hyperparameters according to their domain
44 expertise.

45 In addition, the package includes an evaluation module, which allows for the assessment
46 of trained models using a variety of metrics tailored to the user's specific requirements.
47 Furthermore, the software offers the capability to analyse and retrieve the most proficiently
48 trained models, further enhancing its utility and practicality for scientific research and analysis.

49 To enrich the analysis and comprehension of textual data, SIRITVIS includes two powerful
50 data visualisation tools: PyLDAvis ([Sievert & Shirley, n.d.](#)) and Word Cloud ([Mueller, 2023](#)).
51 The useful data visualisation tool PyLDAvis, created by Sievert and Shirley, improves the
52 interpretation of topic models from textual data. In order to give users a deeper and more
53 intuitive understanding of the latent themes contained in the text corpus, it creates dynamic
54 and interactive visualisations that assist users in exploring the connections between subjects
55 and their associated keywords (see figure 1). On the other hand, the Word Cloud tool offers
56 an engaging and straightforward means of visually representing the most frequently occurring
57 words in the dataset, simplifying the identification of crucial keywords and patterns at a glance
58 (see figure 2). Together, these tools provide users with a comprehensive and user-friendly
59 approach to uncover and communicate the hidden patterns and insights within their textual
60 data.

61 SIRITVIS is renowned for its remarkable ability to map the spatial distribution of tweets and
62 Instagram posts onto a global map, linking each geographical location with its top trending
63 topics and their respective frequencies (see figure 3). Moreover, it excels in the color-coding of
64 these locations based on the sentiments expressed in each post, providing an accurate count
65 of positive, negative, and neutral posts (see figure 4). Furthermore, SIRITVIS facilitates a
66 user-friendly exploration of specific keywords and visualises their occurrences on the world map.
67 This spatial insight contributes significantly to an enhanced understanding of public discussions
68 and lends invaluable support to data-driven decision-making across diverse domains.

69 Comparing and Contrasting Available Toolsets

70 In recent years, the field of text data analysis from social media has witnessed remarkable
71 advancements, offering researchers and practitioners an array of toolkits and approaches to
72 delve into the wealth of online content. Within this dynamic landscape, it becomes imperative
73 to discern the distinctive features of our research, encapsulated in the SIRITVIS framework, as
74 it stands apart from existing related work.

75 Although alternatives such as TLocVis ([Kant et al., 2020](#)), TweetViz ([Stojanovski et al., 2014](#))
76 and Twitmo ([Abuchmueller, n.d.](#)) have their merits, SIRITVIS sets itself apart by providing
77 exceptional flexibility, usability, and comprehensiveness. Its extensive array of advanced topic
78 models, alongside a distinctive capability for pinpointing geographical information and seamless
79 integration with pyLDAvis ([Sievert & Shirley, n.d.](#)) for visualising outcomes, empowers
80 researchers to extract profound insights from social media text data.

81 What sets SIRITVIS apart is its comprehensive suite of evaluation metrics, facilitated by the
82 octis tool ([Terragni et al., 2021](#)). These metrics cover important aspects such as topic diversity,
83 accuracy, inverted RBO, coherence, and Jaccard similarity, ensuring that the Topic Models
84 generated by SIRITVIS are not only reliable, but also imbued with substantive meaning. This
85 robust evaluation framework is a hallmark of the research, emphasizing the toolkit's dedication
86 to producing meaningful results.

87 Furthermore, SIRITVIS places a strong emphasis on user-friendliness and accessibility. Its
88 intuitive interface and detailed documentation cater to both novice and experienced users,
89 making the toolkit approachable and easy to navigate. Additionally, SIRITVIS is designed to
90 accommodate various data sources and formats, ensuring that researchers can adapt it to their

Selected Topic: **0** Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric: $\lambda = 1$ 0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)

Topic 8 is highlighted in red. The map shows the relative positions of 9 topics (1-9) based on their semantic similarity. Topic 8 is distinct from the others.

Top-30 Most Relevant Terms for Topic 8 (7.3% of tokens)

Term	Overall term frequency	Estimated term frequency within the selected topic
crypto	100	40
new	40	20
easy	20	10
special	15	10
year	10	10
make	10	10
experience	10	10
people	10	10
companies	10	10
cryptocurrencies	10	10
world	10	10
ethereum	70	10
time	10	10
network	10	10
ecosystem	10	10
bank	10	10
claim	10	10
use	10	10
early	10	10
bitcoin	40	10
join	30	10
industry	10	10
data	10	10
bni	10	10
love	10	10
market	10	10
launch	10	10
building	10	10
https	100	10
defi	10	10

1. $\text{saliency}(\text{term } w) = \text{frequency}(w) \cdot \left[\sum_t p(t | w) \cdot \log(p(t | w) / p(t)) \right]$ for topics t ; see Chuang et. al (2012)
 2. $\text{relevance}(\text{term } w | \text{topic } t) = \lambda \cdot p(w | t) + (1 - \lambda) \cdot p(w | t) / p(w)$; see Sievert & Shirley (2014)

[illegible]

Figure 2: Word Cloud

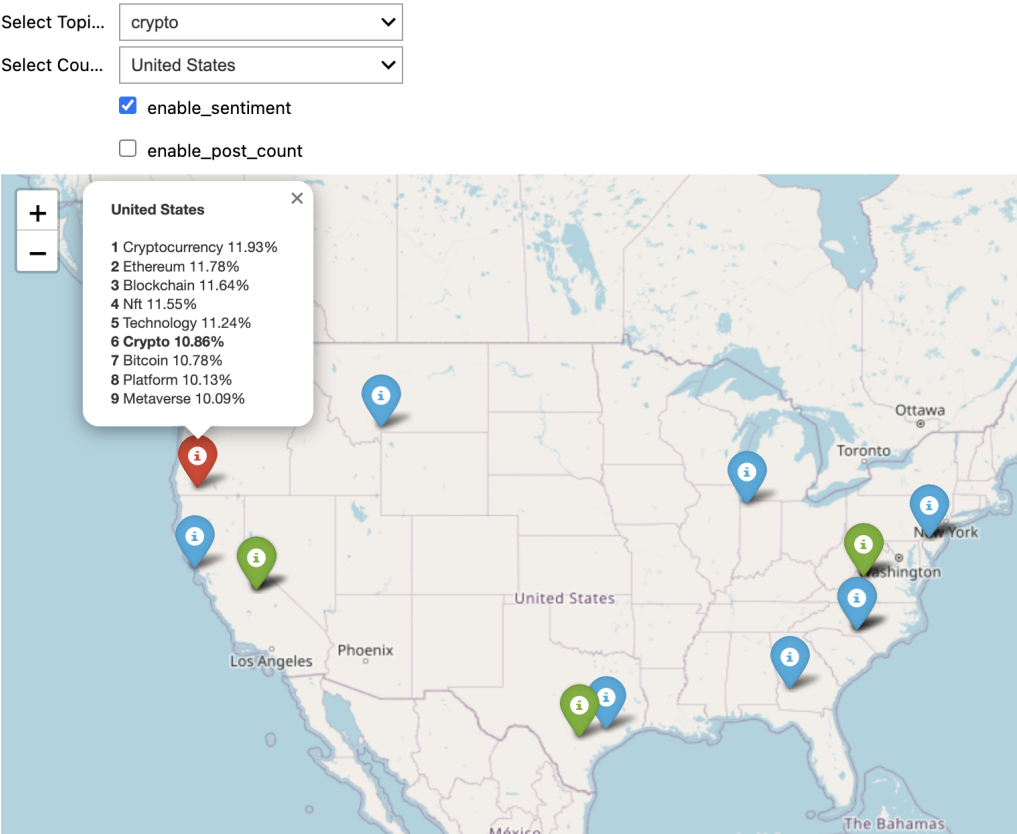


Figure 3: Topic Mapper Frequency Count

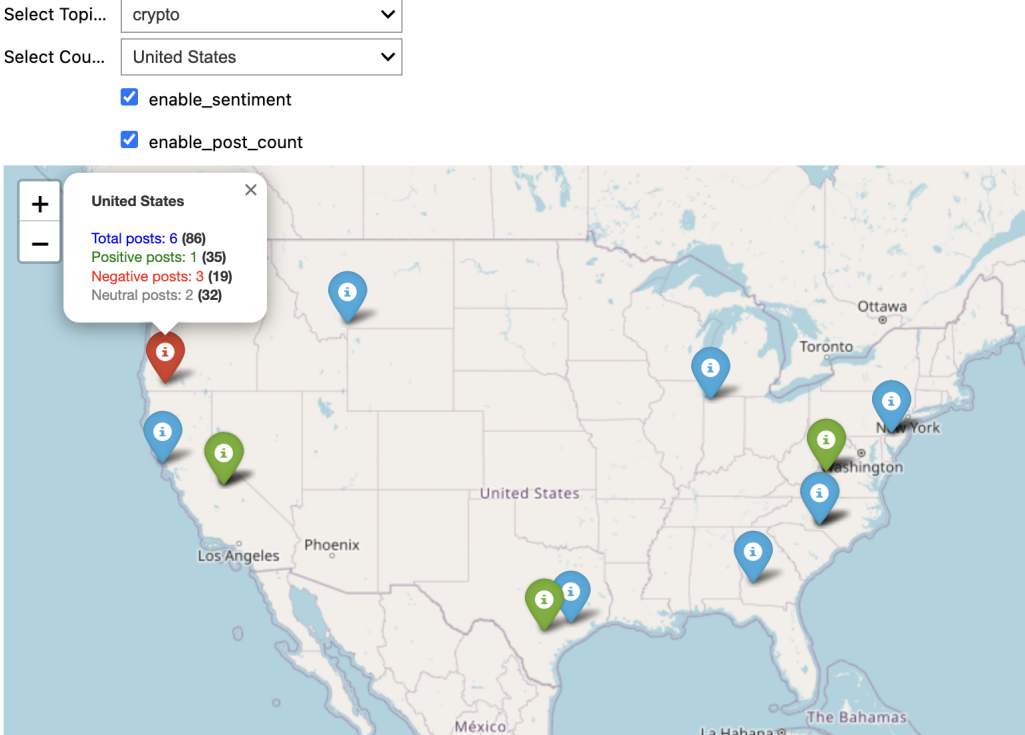


Figure 4: Topic Mapper Sentiment Count

References

- Abuchmueller. (n.d.). *GitHub - abuchmueller/Twitmo: Collect Twitter data and create topic models with R*. <https://github.com/abuchmueller/Twitmo>
- Bianchi, F., Terragni, S., Hovy, D., Nozza, D., & Fersini, E. (2021). Cross-lingual contextualized topic models with zero-shot learning. *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 1676–1683. <https://doi.org/10.18653/v1/2021.eacl-main.143>
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- Hutto, C. J., & Gilbert, E. (2014). VADER: A parsimonious rule-based model for sentiment analysis of social media text. *Eighth International Conference on Weblogs and Social Media (ICWSM-14)*.
- Kant, G., Weisser, C., & Säfken, B. (2020). TTLocVis: A twitter topic location visualization package. *Journal of Open Source Software*, 5(54), 2507. <https://doi.org/10.21105/joss.02507>
- Mueller, A. C. (2023). *Wordcloud* (Version 1.9.1). <https://github.com/amueller/wordcloud>
- Sievert, C., & Shirley, K. (n.d.). *LDAvis: A method for visualizing and interpreting topics* (pp. 63–70). <https://nlp.stanford.edu/events/illvi2014/papers/sievert-illvi2014.pdf>
- Srivastava, A., & Sutton, C. (2017). *Autoencoding variational inference for topic models*. <https://arxiv.org/abs/1703.01488>
- Stojanovski, D., Dimitrovski, I., & Madjarov, G. (2014). Tweetviz: Twitter data visualization. *Proceedings of the Data Mining and Data Warehouses*, 1–4.
- Terragni, S., Fersini, E., Galuzzi, B. G., Tropeano, P., & Candelieri, A. (2021). OCTIS: Comparing and optimizing topic models is simple! *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, 263–270. <https://www.aclweb.org/anthology/2021.eacl-demos.31>