# Understanding Illegal Drug Transactions over Darknet Forums

Yuntian He
he.1773@osu.edu
The Ohio State University
Columbus, Ohio, USA

Pranav Maneriker
maneriker.1@osu.edu
The Ohio State University
Columbus, Ohio, USA

Srinivasan Parthasarathy
srini@cse.ohio-state.edu
The Ohio State University
Columbus, Ohio, USA

## ABSTRACT

Illegal drug transactions have been a challenging issue in the broader Southern California area for decades. In recent years, such transactions have shifted to darknet forums, which enforce the anonymity of traders through encrypted communication and increase the difficulty of drug control. In this work, we aim to leverage machine learning and data science techniques to address this public health concern in Southern California and across the United States. Specifically, we emphasize the importance of exploring these forums and highlight the challenges of research in this field. In addition, we introduce several widely used darknet datasets and present several research problems including stylometry-based authorship attribution and trust signaling strategies on darknet platforms.

## CCS CONCEPTS

• **Computing methodologies → Machine learning**.

## KEYWORDS

Machine Learning, Darknets, Author Identification, AI4Science, AI4Health

## 1 INTRODUCTION

**Drug Transactions in Southern California** Illegal transactions of drugs have posed significant challenges to the local community in Southern California. Los Angeles and San Diego are labeled by the Office of National Drug Control Policy as High-Intensity Drug Trafficking Areas (HIDTA) [7] in the United States. It reveals that disrupting drug trafficking organizations (DTOs) and criminal groups in Southern California supply drugs to distributors both within this region and throughout the country. Now this region is one of the most significant distribution centers in the United States for drugs including heroin, cocaine, and marijuana.

The presence of Southern California in national or international drug markets can be attributed to the following reasons. First, the geographic location on the U.S.-Mexico border and the developed

transportation routes provides convenience for drug storage and distribution [4]. Second, the rural desert regions in Southern California are ideal for drug manufacturing. In addition, this area has a large population with cultural and economic diversity, which lead to a huge drug market. Since drug use and transactions cause severe concerns in the local community, it is necessary to understand and control illicit drug transactions to effectively enhance public health and safety.

**Darknet Forums.** Recently, most drug transactions have shifted to darkness, which are online forums that leverage digital encryption techniques to enable anonymous transactions between different parties [6]. They are hosted on Tor networks, where the identities of users including IP and geographic location are hidden. Therefore, investigating darknets for drug control is significantly challenging.

The advancement of machine learning approaches has achieved state-of-the-art performance in numerous graph-related tasks. However, they find limited use on darknet forums for various reasons. Unlike other eCommerce platforms and social networks, darknets have limited useful features leading to low utility performance. Moreover, the volume of darknet forum data is typically small. A naive option is to apply models pretrained on other networks to darknet forums. Here we present an illustrative comparison in the context of author identification (see Figure 1). The problem aims to identify the authorship based on the stylometry. We adopt LUAR [9] as the model for authorship attribution and train it on traditional network data. Our test data include four distinct datasets: two darknet forums (namely, 'The Hub' and 'Dread') and two non-darknet datasets. The results demonstrate that simply applying trained models to darknet forums fails to achieve ideal performance, which motivates our study of darknet forum analysis.

## 2 DATASET DESCRIPTION

To aid our study of drug markets, we selected forum and review-based datasets from different sources on the darkweb.

**Forum Data.** 'Dread' and 'The Hub' forums have been central to discussions on the Darkweb [2]. We sourced our data from the CrimeBB collection of forum datasets released by the Cambridge Cybercrime group [8]. As our focus is on user identification, we grouped together data corresponding to each user (author) on these forums. We divided each dataset into three temporal splits, with 70% of the posts used to train author identification models, and a midway split for each author in the remaining data to provide query and target sets for each author that are used to evaluate the trained models. We filtered the data to ensure that we had at least 2 posts per author in each split. Following this filtering step, the data on 'The Hub' spanned a period starting January 2014 until August 2019, while 'Dread' spanned a period starting February 2018 to January 2020.
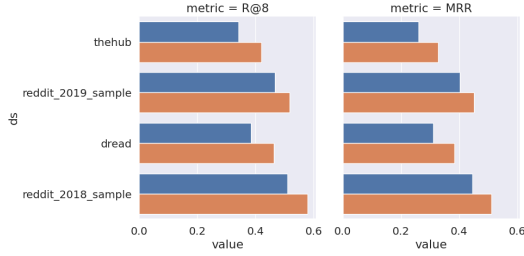
**Figure 1: Author identification results with a model pre-trained on one year of Reddit data. Performance degrades for non-Reddit datasets ('The Hub' and 'Dread').**



**Figure 2: Frequent combinations of markets for different vendors reviewed on Kilos.**

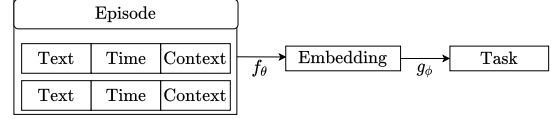**Review Data.** The Kilos dataset is a part of the Darknet market archives [1] and contains 235K reviews scraped from 6 darknet markets. Kilos was originally constructed as a search engine for Vendors' listings and reviews. Multiple vendors with reviews on Kilos are active across multiple markets, which makes this dataset suitable for studying patterns of trust signaling when vendors migrate across markets. Figure 2 shows market collections with the highest number of vendors operating on both markets in each pair.

## 3 RESEARCH ON DARKNET

**Stylometry-based Authorship Attribution.** Since one malicious vendor may operate multiple accounts, it is important to find these accounts for a target user. Text stylometry represents a user's writing style. It has been utilized as a key signature for authorship attribution in social networks. Applying machine learning models with statistical or tokenized features can achieve promising performance. However, traditional stylometry-based approaches rely on long text corpora while text information in darknet forums is usually shorter. To this end, we need to leverage additional information for authorship attribution.

Figure 3 illustrates our proposed solution SYSML [5]. The model is fed with an episode (set) of user posts as input and generates a vector representation from textual, temporal, and contextual information. The learned representation is used by a task learner to make predictions of user identity. Specifically, the text embedding module encodes tokenized text data and feeds the encodings to several sliding window filters and pooling layers. A fully connected layer finally generates the text embedding of a given episode. Then



**Figure 3: The Workflow of SYSML.**

the time embedding module encodes the day of the week for each post as a vector representation. Moreover, the contextual embedding module uses the relationship between different entities to generate embeddings. It builds a heterogeneous graph consisting of users, subforums, threads, and posts. Then it leverages a random-walk-based embedding approach [3] to learn node embeddings with specific metapaths. The learned representations capture the context of interactions in the darknet forum. The final representation of an episode is learned from all these embeddings by a transformer pooling layer.

We evaluate SYSML on multiple datasets. The results show the efficacy of author attribution in darknets with a lift of 2.5X on Mean Reciprocal Rank and 2X on Recall@10 compared with our baselines. Please refer to [5] for full results.

**Understanding Trust Signaling.** Prior to exchanging drugs, vendors need to first build trust through text posted in darknet forums. Therefore it is necessary to study how trust is formed in this highly anonymized environment. Specifically, we emphasize (1) how vendors express trustworthiness to their potential buyers, (2) how buyers evaluate the trustworthiness of vendors, and (3) how vendors react to threats to their trustworthiness.

We select two representative darknet forums (namely, 'The Hub' and 'Dread') from CrimeBB [8]. We also collect the Kilos dataset [1] which contains multiple drug transactions. First, we extract 93 Kilos vendors from 'The Hub' and 586 from 'Dread', respectively, by matching users who have the same username. We assume these accounts on darknet forums belong to corresponding vendors and then extract their activity history. These text data showcase how vendors convey their trustworthiness. In future work, we plan to analyze the transaction data to evaluate their trust signaling strategies.

From the remaining data, we collect posts that mention one or more vendors. This includes 3,414 posts on 'The Hub' and 14,049 posts on 'Dread'. Those posts reflect how buyers review the vendors. We plan to apply natural language processing models and perform sentiment analysis to evaluate the quality of reviews.

## 4 CONCLUSION

In this work, we study the issue of illicit drug transactions over darknet forums to address public health concerns in Southern California. We introduce darknet datasets and showcase how we apply data science and machine learning methods to these datasets for two distinct research problems.

## ACKNOWLEDGMENTS

# REFERENCES

[1] Gwern Branwen, Nicolas Christin, David Décary-Hétu, Rasmus Munksgaard Andersen, StExo, El Presidente, Anonymous, Daryl Lau, Delyan Kratunov Sohhlz, Vince Cakic, Van Buskirk, Whom, Michael McKenna, and Sigi Goode. 2015. Dark Net Market archives, 2011-2015. https://www.gwern.net/DNM-archives. https://www.gwern.net/DNM-archives

[2] Selina Y Cho and Joss Wright. 2019. Into the Dark: A Case Study of Banned Darknet Drug Forums. In *International Conference on Social Informatics*. Springer, 109–127.

[3] Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. 2017. metapath2vec: Scalable representation learning for heterogeneous networks. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. 135–144.

[4] Nathan P Jones. 2016. Pangas, Trickery, Intimidation, and Drug Trafficking in California. Small Wars Journal.

[5] Pranav Maneriker, Yuntian He, and Srinivasan Parthasarathy. 2021. SYSML: StYlometry with Structure and Multitask Learning: Implications for Darknet Forum Migrant Analysis. *Empirical Methods in Natural Language Processing* (2021).

[6] James Martin. 2014. *Drugs on the dark net: How cryptomarkets are transforming the global trade in illicit drugs*. Springer.

[7] Executive Office of the President Office of National Drug Control Policy. 2022. *The Report of High Intensity Drug Trafficking Areas Program*. https://www.whitehouse.gov/wp-content/uploads/2022/12/HIDTA-Annual-Report-to-Congress-2022.pdf

[8] Sergio Pastrana, Daniel R Thomas, Alice Hutchings, and Richard Clayton. 2018. Crimebb: Enabling cybercrime research on underground forums at scale. In *Proceedings of the 2018 World Wide Web Conference*. 1845–1854.

[9] Rafael A Rivera-Soto, Olivia Elizabeth Miano, Juanita Ordonez, Barry Y Chen, Aleem Khan, Marcus Bishop, and Nicholas Andrews. 2021. Learning universal authorship representations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 913–919.