

Pairing Users in Social Media via Processing Meta-data from Conversational Files

Meghna Chaudhary, Ravi Sharma, and Sriram Chellappan^(⊠)

Department of Computer Science and Engineering, University of South Florida, Tampa 33620, USA

{meghna1,ravis}@mail.usf.edu, sriramc@usf.edu

Abstract. Massive amounts of data today are being generated from users engaging on social media. Despite knowing that whatever they post on social media can be viewed, downloaded and analyzed by unauthorized entities, a large number of people are still willing to compromise their privacy today. On the other hand though, this trend may change. Improved awareness on protecting content on social media, coupled with governments creating and enforcing data protection laws, mean that in the near future, users may become increasingly protective of what they share. Furthermore, new laws could limit what data social media companies can use without explicit consent from users. In this paper, we present and address a relatively new problem in privacy-preserved mining of social media logs. Specifically, the problem here is the feasibility of deriving the topology of network communications (i.e., match senders and receivers in a social network), but with only meta-data of conversational files that are shared by users, after anonymizing all identities and content. More explicitly, if users are willing to share only (a) whether a message was sent or received, (b) the temporal ordering of messages and (c) the length of each message (after anonymizing everything else, including usernames from their social media logs), how can the underlying topology of sender-receiver patterns be generated. To address this problem, we present a Dynamic Time Warping based solution that models the meta-data as a time series sequence. We present a formal algorithm and interesting results in multiple scenarios wherein users may or may not delete content arbitrarily before sharing. Our performance results are very favorable when applied in the context of Twitter. Towards the end of the paper, we also present interesting practical applications of our problem and solutions. To the best of our knowledge, the problem we address and the solution we propose are unique, and could provide important future perspectives on learning from privacy-preserving mining of social media logs.

Keywords: Social media · Privacy · Big-data · Meta-data · Dynamic Time Warping

1 Introduction

As of today, social media is a major platform that citizens across the globe choose to communicate over. Such communications span many forms including bi-directional one-on-one messaging among peers; posting content to be viewed by a larger group; following posts from favorite personalities; advertising products to customers; reviewing products and services; and so much more. It is a fact that mining such data is a billion dollar business with so many players in the market now. One of the downsides of this scenario is compromising the privacy of common citizens.

As of today, a vast majority of users on social media do not care about privacy, and even if they do, they still actively communicate over many social media platforms. In fact, this is true even amongst the more educated citizens. However, this trend may change in the future. There are numerous reports now wherein common citizens are becoming victims because the content they post/share on social media has been accessed by third-party entities beyond the scope/knowledge of the victim. These include people losing their jobs, denied admissions to universities, being charged with crimes and fines, losing custody of children, and so much more. With increasing awareness of such reports, citizens of the future are likely to be increasingly aware of privacy violations. In parallel, governments across the world are also monitoring un-restricted access of social media content by data mining agencies, and newer laws and regulations are being generated today, one prime example being The EU General Data Protection Regulation (GDPR), that came into force in 2018.

Given these developments, there are now urgent efforts in the academia and the industry to investigate the paradigm of privacy-preserved learning from bigdata. Essentially, the issue at hand is how can we learn meaningful information from data (whether generated from social media or not), which still preserving the privacy of the data as intended by users. There are number of studies currently looking at this paradigm from perspectives of AI [1–3], Crypto [4,5], and Social Sciences [6,7]. In this paper, we make contributions towards this paradigm. Specifically, the scenario we address is one where users wish to protect identities (i.e., their own username and those of peers) and content of what they share on social media (texts, images, multi-media etc.), but are willing to share meta-data of such content - for example, whether a message was sent or received; the temporal ordering of messages; and the size of the message (e.g., number of characters in the text). In this scenario, our specific problem is identifying the network/communication topology - that is pairing the sender and receiver among multiple conversational files from multiple users containing only meta-data.

We propose a Dynamic Time Warping (DTW) based approach to our problem. Essentially, our solution models the meta-data sequences from each conversational file as a temporal sequence, and then uses DTW technique to find the best similarity match among multiple conversational files. We evaluate the technique using a limited sample of Twitter users/logs (47 users and 128 conversational files), and performance results are very favorable, and our technique is also scalable. What is important is that our proposed method gracefully degrades when users can choose to delete content (and hence the underlying meta-data) within conversational files before sharing. To the best of our knowledge, our problem in this paper is unique and has not been addressed before. There are important practical ramifications of this problem from multiple perspectives including one related to cyber-abuse, which we elaborate on towards the end of the paper.

The rest of the paper is organized as follows. Section 2 discusses about important related work. Section 3 presents information about our data source and the formal problem statement. Section 4 elucidates the methodology followed for matching of conversational files. In Sect. 5, we discuss performance evaluations. Section 6 discusses about the practical relevance of our contributions in this paper. Finally, we conclude the paper in Sect. 7.

2 Related Work

We now present a brief overview of important work related to this paper.

2.1 Identifying Potential Friends in Social Networks

In [8], a study is conducted in order to determine effectiveness of textual features and network trends to make recommendations for friends on the Twitter platform. In this study, 200 most recent tweets from 100,000 users are collected and analyzed. Also, another 200 recent tweets from 10 friends and 10 followers for each user is also collected. Different information sources include posts by a user, posts by user's mentions, and friends and followers are processed, and a model integrating Bag of Words and Principal Component Analysis is designed to identify potential friends. Another model that considers network-level metrics is also used, wherein two users are considered to be friends if their social connections (i.e., those they follow and those that follow them) share similar content. Based on this network structure also, friends are predicted for each user. The paper argues that network level structures are better suited for friends recommendation compared to purely textual based features. While using pure network level properties does preserve privacy to a certain degree, the fast that user names are shared and processed raises privacy issues.

In [9], another study is conducted wherein pictures and contact tags of 10,000 Flickr users are collected, wherein the 10,000 users belong to 2000 social groups. Features from these picture and contact tags are processed to see how similarities as computed from the features match the actual underlying topology. The paper demonstrates that picture and contact tag features can model the underlying social topology very well. Based on these results, a recommendation system is also proposed for friends matching. But processing images, can have serious privacy implications.

2.2 Computational Techniques to Preserve and Compromise Privacy

There are a number of studies now that design novel encryption techniques to search from encrypted data [4,5]. The basic idea is to ensure that legitimate users have a degree of access to meaning of the data even under encryption, while adversaries do not derive any meaningful information. The limitations though are that even with state-of-the-art encryption today, the quality and quantity of information gleaned from encrypted data is still minimal, and not enough for many applications. There are also other recent papers in the realm of [1–3] that investigate how much information do machine learning models remember after training, to the point where information about ground truth data used to train the models can be recovered. These avenues of research demonstrate serious privacy breaches from the perspective of exposing data used to train algorithms today. There are also a number of studies also in the realm of using the privacy-preserving and spatially compact Bloom Filters to store and retrieve records based on similarities. Most of these studies [10,11] and [12], are in medical contexts though, and only for data searching operations.

3 Data Source and Formal Problem Formulation

The source of our dataset is Twitter. For this study, we extracted tweets from a group of 47 socially connected users and obtained 3200 most recent tweets per user. Conversational connections in terms of user names and their mentions were collected, and tweets which do not represent a conversation were discarded. The period of data collection was from January 2019 to March 2019.

The process of anonymization is a little more complex and is elaborated below. First, for each of the 47 users whose data we collect, we anonymize their identities as $User_1$ to $User_{47}$. For example, consider $User_1$. Let $User_1$ engage in communications with 10 other contacts in the timeframe of our data collection. Now when $User_1$ wishes to share his/her history, with privacy expectations, $User_1$ will not reveal the identities of these contacts. What $User_1$ will instead do is anonymize their identities arbitrarily before sharing. As such in our data collection process, every contact of $User_1$ is anonymized as $User_1^{c1}$, $User_1^{c2}$,... $User_1^{c10}$ for ten such contacts. Naturally, to protect the content shared between $User_1$ and any of its contacts, the records of corresponding communications shared are only whether a message was sent or received; the temporal ordering; and the overall length of the message/content. The same process is followed for every other user, and for every contact of that user in our study.

One thing is very important to note here, and this stems from the way each user's data is independently anonymized before processing. Let us for now assume that two users anonymized as $User_2$ and $User_3$ are indeed conversing with each other. Due to the way we independently anonymize the contact for each user in our study, it could be the case that $User_3$ is identified as $User_2^{c1}$ in the dataset corresponding to $User_2$; and $User_2$ is identified as $User_3^{c5}$ in the

dataset corresponding to User₃. That is, there is no linkage among these contacts due to independence of anonymization under our data collection process. The problem we address here is how to still derive the social/network topology by identifying that User₂ and User₃ are indeed communicating with each other using only the metadata collected from each user independently, and especially when certain logs of data could be deleted arbitrarily by each user, irrespective of whether or not the other party deletes them.

Figure 1 presents a snapshot of the dataset collected in the above manner. In the dataset illustrated, User₁ has 10 connections labeled from User₁^{c1} to User₁^{c10}; User₂ has 7 connections labeled from User₂^{c1} to User₂^{c7}; User₃ has 5 connections labeled from $User_3^{c1}$ to $User_3^{c5}$; $User_5$ has 8 connections labeled from $User_5^{c1}$ to $User_5^{c8}$; and so on. Finally, $User_{46}$ has 9 connections labeled from $User_{46}^{c1}$ to $User_{46}^{c9}$ and User₄₇ has 8 connections labeled from User₄₇^{c1} to User₄₇^{c8}. The corresponding sent/received patterns and length of messages for these users are also presented in Fig. 1. Given this dataset, our problem is to determine the overall social/network topology. Essentially, our goal is to identify (based on alignment of sent/received patterns and lengths in the dataset) that User₁ is a social contact of User₃; User₂ is a social contact of User₃; User₃ is a social contact of User₅; User₅ is a social contact of User₄₆; User₄₆ is a social contact of User₄₇ and so on. While this is easier to do if all messages among all contacts are retained and shared, this will not be the case in reality. Users can arbitrarily choose to delete messages in their conversations with one or more or contacts, irrespective of whether or not the same message is deleted by the other party. Under such situations, which is common in practice, the problem of establishing social connections becomes much harder with potential for false negatives and false positives as well.

4 Our Dynamic Time Warping Framework

In this paper, we employ a Dynamic Time Warping (DTW) approach to solve our problem. DTW is an ideal technique for our problem, since, the calculation of the distance (or similarity) metric between time-series sequences of datasets using DTW can overcome problems related to comparing short patterns of data and class imbalances. The technique is also independent of certain non-linear variations in the time dimension [13–16].

To apply DTW for our problem, we need to encode our dataset as a timeseries sequence. This is doable in our case because the sent/received patterns over time for any user are essentially temporal orderings, to which the length of a corresponding message can be easily appended. With this temporal ordering in place across all conversational files for all users, the problem now is find the best alignment between encoded files in order to generate the connection between users, and hence the overall social topology. As we present below, the DTW technique involves determination of a warping path between temporal sequences, from which the path that optimally minimizes the corresponding distance is the ideal one.

Note that, there are certain critical steps that need to be executed to compute the warping path, and hence the distance metric. These are presented below.

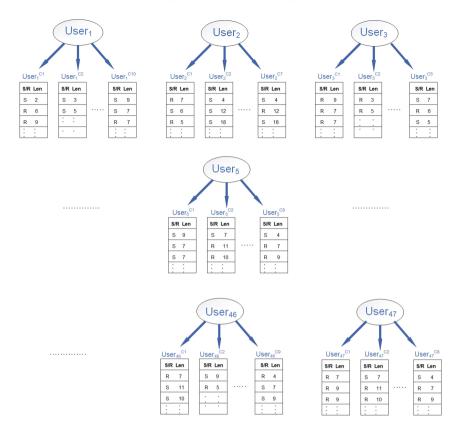


Fig. 1. A Snapshot of our dataset for analysis. Our problem is to determine social connections even in the presence of arbitrary message drops

- 1. **The boundary case**: We ensure that the starting points and the ending points of each pair of sequences are identified and matched as such. This ensures that all data points in the entire sequences are compared in determining the similarity metric, and hence prevents data loss.
- 2. **Ensuring monotonicity**: When comparing sequences of data points, any set of points in time in one sequence that are already aligned with points in another sequence are not used to evaluate for matching with later points in time. In order words, comparison of similarities between points in multiple sequences is monotonically increasing.
- 3. The step size: In computing the distance metric across a warping path, we ensure that every point within the neighborhood of a data point is considered for distance measurement, and as such jumps across data points are not allowed [17].

Algorithm Description: We now present our DTW algorithm for the problem above. But first, we present the encoding technique for our ground truth data

to facilitate discussions. Consider a sequence of messages from any arbitrary User (say User_x) to another arbitrary user (say User_y) as [S, R, S, S, R, ..., S, R]. This sequence will be encoded as [1, 0, 1, 1, 0, ..., 1, 0] during processing. If the corresponding message lengths are [10, 12, 4, 5, 6, ..., 12, 23] are considered, the encoded sequence now becomes [101, 120, 41, 51, 60, ..., 121, 230] that integrates both message lengths and sent/received status.

We now present a simple example of how to execute the DTW technique to find similarities for our problem statement. We present the example for the case of considering sent/received patterns only, without considering message lengths, but the technique is straightforward to integrate lengths also. Consider a sent/received sequence denoted as $T_1 = [1, 0, 1, 0]$ (read from bottom to top in the left in green in Fig. 2). Consider another sequence $T_2 = [0, 0, 1]$ (in red in the bottom in Fig. 2). Note that two sequences are not the same even in the number of entries they have. This can happen in our case, since users are allowed to arbitrarily delete content before sharing, and our technique will accommodate this case. In the DTW technique, for this example, we first compute the distance matrix $Dist_{Matrix}$ of dimensions 4×3 , where each entry in the matrix is computed using the following equation (where ED stands for Euclidean Distance):

$$DTW_Dist[i, j] = ED[T_1[i], T_2[j]] + min(DTW_Dist[i - 1, j], DTW_Dist[i - 1, j - 1], DTW_Dist[i, j - 1])$$
(1)

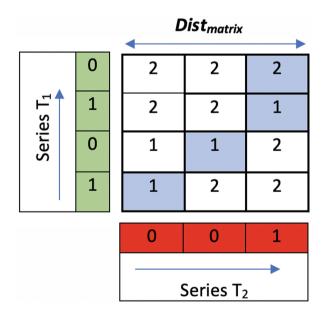


Fig. 2. Distance matrix computation for two series using DTW (Color figure online)

Finally, we traverse the optimal warping path from the end point to the starting point between the two paths (which in the distance matrix is from topright cell to the bottom-left cell) by choosing the adjacent cell with the minimum distance. In Fig. 2, this is shown in blue shaded color. The sum of the entries in these cells is the overall distance computed, and denotes the similarity between the two time sequences, which in this case is 5.

It is easy to infer that the above procedure can be expanded to include sequences encoded with message lengths. Furthermore, for two sequences that are exactly the same, the distance computed will be 0. Finally, we point out that in order to compare a sequence of messages shared by one user with a sequence of messages shared by another user for similarity matching, we invert the sequence of messages in the other user prior to determining the warping path. This ensures that a message sent by the first user is compared with a message received message by the other user; and that a message received by the first user is matched with a message sent by the other user. Algorithm 1 presents the formal sequence of steps in our solution to compute similarity across message sequences. The complexity here is $O(f^2 \times m^2)$ where f denotes the number of conversation files to compare, and f denotes the maximum length among sequences to compare. Our technique to implement DTW was built using Python. Note that Python provides open-source Just-in-Time (JIT) compiler Numba to generate faster machine code which helps accelerate computation.

5 Results

5.1 Overview

We present results of our technique to match users that are communicating with each other, based on processing only meta-data logs of conversational files. Before that we present some preliminaries. First, of all we point out that the total number of Twitter users in our experiment was 47. The total number of conversational files was 128. The 3200 most recent tweets per user were obtained, and the period of data collection was from January to March 2019. On an average, each conversational file contained anywhere from 1,500 to 3,000 messages. We only processed textual messages between users. As mentioned above, all user identities were anonymized. Only sent/received status of messages, their temporal ordering and lengths were processed. No actual textual content was processed. We considered two classes of features for determining the connections between users. The first one only included the sent/received patterns (without considering the lengths). The second includes the lengths of the corresponding messages along with the sent/received patterns.

Note that in our problem formulation, users can choose to arbitrarily delete messages from one or more conversational files before sharing. As such in our dataset, we provision for that. Specifically, we define two new parameters denoted as α and β . Here α denotes the percentage of files whose content a user opts to delete before sharing. In other words, when $\alpha=0$, the user does not modify any file; and when $\alpha=100\%$, the user chooses to modify every file shared. The next parameter β denotes the percentage of messages within a file that a user chooses to delete. Here again, when $\beta=0$ for a particular file, the user does

Algorithm 1. DTW to Compute Similarity

```
\overline{\text{User Files} = U_{files}}
Source User File Series = X
Receiver User File Series = Y
Length of Source User File Series = m
Length of Receiver User File Series = n
Array DTW_{Dist}[0...,0...n]
Euclidean Distance = ED
Similarity score = similarity_{score}
Input: List X and Y of patterns in two conversation files
Output: The similarity score (measure of distance)
1: for p \leftarrow 1, number of user files -1 do
2:
       for q \leftarrow p + 1, number of user files do
3:
           X = U_{files}[p]
           Y = U_{files}[q]
4:
5:
           cost=0
6:
           function DTW((X[1...m], Y[1,...n]))
 7:
               for i \leftarrow 1, m do
8:
                   for j \leftarrow 1, n do
9:
                       DTW\_Dist[i, j]
                                                      ED[X[i], Y[j]] + min(DTW\_Dist[i -
    [1, j], DTW\_Dist[i - 1, j - 1], DTW\_Dist[i, j - 1]
10:
                   end for
11:
               end for
12:
               path \leftarrow [m, n]
13:
               y \leftarrow m
14:
               z \leftarrow n
15:
                while z >= 0 and y >= 0 do
16:
                   if z = 0 then
17:
                       z \leftarrow z - 1
18:
                   else if y = 0 then
19:
                       y \leftarrow y - 1
20:
                   else
                       if DTW\_Dist[y, z-1] = min(DTW\_Dist[y, z-1], DTW\_Dist[y-1])
21:
    [1, z], DTW\_Dist[y - 1, z - 1]) then
22:
                           z \leftarrow z - 1
23:
                       else if DTW\_Dist[y-1,z] = min(DTW\_Dist[y,z-1], DTW\_Dist[y-1])
    1, z, DTW_Dist[y - 1, z - 1]) then
24:
                           y \leftarrow y - 1
25:
                       else
26:
                           y \leftarrow y - 1
27:
                           z \leftarrow z - 1
28:
                       end if
29:
                   end if
30:
                end while
31:
                Add \ path[z,y] \ to \ path \ array
32:
                for [z, y] \in path do
33:
                   cost \leftarrow cost + ED[y, z]
34:
                   return cost
35:
                end for
36:
            end function
37:
            similarity_{score} = DTW(X[1...m], Y[1,...n])
38:
            print \ similarity_{score}
39:
        end for
40: end for
```

not delete any conversation in that file; and when $\beta=100\%$ for a file, the user chooses to delete every conversation in that file. In presentation of results below, we vary α and β from 0 to 90%. Note that the files chosen for deletion and the messages chosen to be deleted within a file for each user are random. Also, note that content is deleted independently for each user, irrespective of whether or not the corresponding content is deleted in the communication files of the other party. This is the most practical scenario, and hence we evaluate our system under this scenario.

5.2 Metrics and Results

We employ three standard measures to evaluate our algorithm for our problem. These are Precision, Recall and the ROC curve.

Precision and Recall: In Fig. 3(a) and (b), we plot the Precision of our system for varying values of α and β . Here, Fig. 3(a) is the plot for the case where only sent/received patterns are considered (without considering message lengths). Figure 3(b) is the plot for the case where sent/received patterns and message lengths are both considered. Matching of conversational files across users is done for those pair of files that have the lowest similarity scores. It is straightforward to calculate True Positives, False Positives and False Negatives from this, which will directly yield the computation of Precision and Recall.

First off, we see from that when α and β are low, the Precision and Recall are high for both feature sets. This is straightforward, since matches among communicating users are more reliable when there are very few message drops. As α and β increase, we see lower precision and recall values. We also see that Precision is more sensitive to β than to α . This is because, even when more conversational files are chosen to be deleted (i.e., higher values of α), the fact that only few messages within each file are deleted (i.e., lower values of β) enables our system to match correct conversational files better. On the other hand, even when α is lower, the situation when β is higher, means that more messages within files are deleted, which makes correct matches harder. This is why our performance metrics are more sensitive to β than to α . Finally, we see that Precision and Recall are better in Figs. 3(b) and 4(b) than in Figs. 3(a) and 4(a), since the inclusion of message lengths along with sent/received patterns in Figs. 3(b) and 4(b) results in better performance than the case with only considering sent/received patterns as done in Figs. 3(a) and 4(a).

ROC Curves: In Figs. 5 and 6, we plot the ROC curves for a range of α and β values. We once again get very good ROC curves when α and β are low, and they get progressively poor when α and β increase. We also see superior ROC curves when message lengths are integrated with sent/received patterns (Fig. 6) compared to considering only sent/received patterns (Fig. 5).

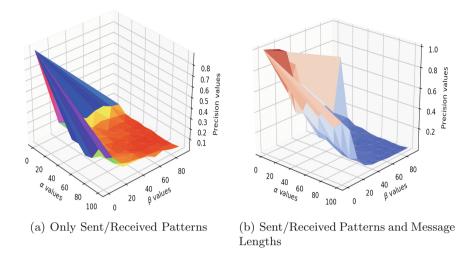


Fig. 3. Precision plots for two classes of features, and for varying α and β values

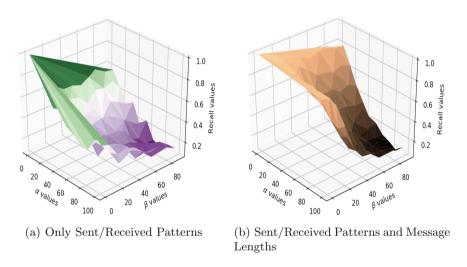


Fig. 4. Recall plots for two classes of features, and for varying α and β values

5.3 Summary

To summarize, we believe that our technique is satisfactory despite limited users and limited datasets. The fact that our technique is sensitive to message drops is reasonable, but for relatively smaller drops, the performance is still acceptable. We believe if system deployers have a fair idea of their user profile, and some knowledge of message drop parameters, the right thresholds can be chosen for better accuracies. Furthermore, with larger scale datasets, more advanced machine learning techniques could also be developed and performance improved.

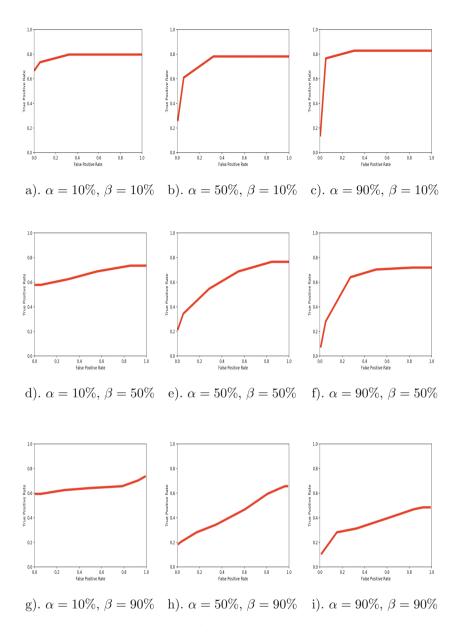


Fig. 5. ROC curves with only sent/received patterns for various α and β values

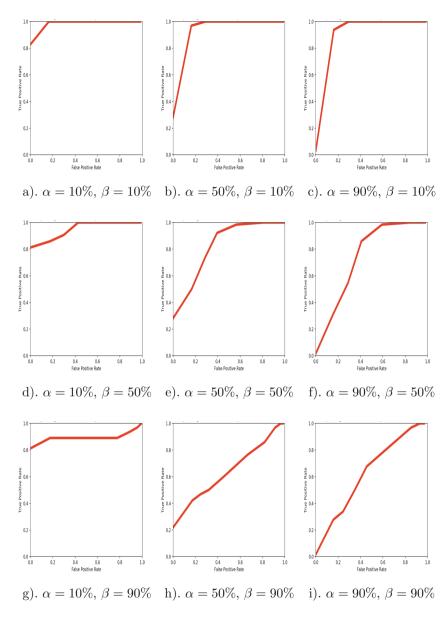


Fig. 6. ROC curves with sent/received patterns and message lengths for various α and β values

6 Discussions on Practical Relevance of This Paper

We now present some important perspectives on the practical impact of our work in this paper. First off, with massive scale participation of citizens on social media platforms, there are billion dollar industries that focus on mining social media data for profit. Unfortunately, the policies of such companies in terms of what they process, what they share, who they sell results to are not transparent at all. However, as mentioned earlier in the paper, these trends are likely to change with increasing privacy awareness among multiple governmental agencies, civil liberties organizations, and the educated public. As such, in the near future, organizations will have to contend with people imposing limits on what companies can and cannot use about data from the public. Now, of course, different users may have different privacy expectations, and naturally, there will be a notion of adaptive privacy requirements and data sharing in such cases. The dataset we generate in this paper (and, in our opinion) is highly privacy preserving, since all user identities are anonymized, and only meta-data of content is processed (only sent/received patterns and lengths of messages).

Now a question may arise regarding the practical utility of our problem statement - which is to determine entities that communicate with each other based on processing the meta-data, and hence the overall social topology. We present perspectives below. The first application relates to cyber-abuse, especially as it pertains to young people. Currently, the process of research on cyber-abuse primarily involves adult researchers looking at content of messages and then indicating whether or not a message (or maybe an image) constitutes abuse. This is fundamentally flawed, since the perspectives of the actual victim - in this case children are not solicited. It is common knowledge that unique slang that children use, context of communications (as relates to events in schools or playgrounds), code-mixing etc. make it very hard for adults (that are absent from the social contexts of younger victims) to decipher the emotional impact of messages. This issue is actually well studied in area called digital divide [18–20].

On the other hand though, requesting content directly from victims (again, children in this case) is also problematic because of IRB regulations, and the high risk it entails when sensitive data from an already sensitive population is analyzed. This issue significantly impedes the possibility of getting robust ground truth data, which is important for research with realistic outcomes. We believe that a system like ours can mitigate these shortcomings. We are currently designing smart-phone apps where users (of any age), can willingly assent to share meta-data of their communication logs, wherein the meta-data will only be lengths of messages, sizes of files, results after performing sentiment analysis on the messages (within the device). All user identities will be anonymized. Subsequently, and if there is larger scale adoption of our system, we could use results from this paper to derive the social/network topologies of young people. In addition, if the user is also willing to mark certain messages as abusive, then the meta-data of those messages, along with the preceding and succeeding messages can be analyzed to derive signatures of abuse. If these results can be mapped back to the derived network topologies, we could perform research with significant impact to answer questions like (a) feasibility of early warning of cyber-abuse from meta-data alone; (b) identify victims and abusers in the topology, and apply graph theoretic results to understand topology evolution; (c) use metrics like graph centrality to see which nodes are more significant; (d) model how the graph enables the dissemination of abusive content across various nodes and so much more. We strongly believe that results from this paper, coupled with the insights mentioned above have significant impact to cyber-abuse research in the near future. Naturally, the impact of this research can also extend to other domains where users are willing to share meta-data for critical applications like - modeling the efficiency of an office environment as it pertains to social communications among employees; model and analyze privacy-preserving topologies in the realm of doctor-patient or nurse-patient interactions for better dissemination of health related information etc. Designing systems for these applications, and furthering research in analyzing meta-data in these unique contexts is part of our future work.

7 Conclusions

In this paper, we presented a unique problem in the realm of privacy-preserved mining of social media logs to derive network topologies. The source of data for our study was Twitter, and included 47 users. The novelty of our study is the significant privacy accorded to the dataset during analysis, wherein all user identities were anonymized, and only message lengths of content was processed, and not the actual content. We presented a Dynamic Time Warping based algorithm for our problem, and presented interesting results on the accuracy of deriving network topology from meta-data alone. Towards, the end of the paper, we presented some practical impact of work accomplished in this paper. With increasing privacy awareness across the globe, coupled with newer privacy laws coming into effect, we believe our work in this paper is timely and relevant, and can create new societal scale and privacy-preserving big-data applications.

Acknowledgment. This work was supported in part by US National Science Foundation (Grant # 1718071). Any opinions, findings and conclusions are those of the authors alone, and do not reflect views of the funding agency.

References

- Melis, L., Song, C., De Cristofaro, E., Shmatikov, V.: Exploiting unintended feature leakage in collaborative learning. arXiv preprint arXiv:1805.04049 (2018)
- Hunt, T., Song, C., Shokri, R., Shmatikov, V., Witchel, E.: Chiron: privacy-preserving machine learning as a service. arXiv preprint arXiv:1803.05961 (2018)
- Song, C., Ristenpart, T., Shmatikov, V.: Machine learning models that remember too much. In: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, pp. 587–601. ACM (2017)
- Bost, R., Minaud, B., Ohrimenko, O.: Forward and backward private searchable encryption from constrained cryptographic primitives. In: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, pp. 1465– 1482. ACM (2017)
- Demertzis, I., Papamanthou, C.: Fast searchable encryption with tunable locality. In: Proceedings of the 2017 ACM International Conference on Management of Data, pp. 1053–1067. ACM (2017)

- Jung, A.R.: The influence of perceived Ad relevance on social media advertising: an empirical examination of a mediating role of privacy concern. Comput. Hum. Behav. 70, 303–309 (2017)
- Tsay-Vogel, M., Shanahan, J., Signorielli, N.: Social media cultivating perceptions of privacy: a 5-year analysis of privacy attitudes and self-disclosure behaviors among facebook users. New Media Soc. 20(1), 141–161 (2018)
- 8. Benton, A., Arora, R., Dredze, M.: Learning multiview embeddings of twitter users. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp. 14–19 (2016)
- Huang, S., Zhang, J., Wang, L., Hua, X.: Social friend recommendation based on multiple network correlation. IEEE Trans. Multimedia 18(2), 287–299 (2016). https://doi.org/10.1109/TMM.2015.2510333
- Vatsalan, D., Christen, P.: Privacy-preserving matching of similar patients. J. Biomed. Inform. 59, 285–298 (2016). https://doi.org/10.1016/j.jbi.2015.12.004. http://www.sciencedirect.com/science/article/pii/S1532046415002841
- Randall, S.M., Ferrante, A.M., Boyd, J.H., Bauer, J.K., Semmens, J.B.: Privacy-preserving record linkage on large real world datasets. J. Biomed. Inform. 50, 205–212 (2014). https://doi.org/10.1016/j.jbi.2013.12.003. http://www.sciencedirect.com/science/article/pii/S1532046413001949. Special Issue on Informatics Methods in Medical Privacy
- 12. Chi, Y., Hong, J., Jurek, A., Liu, W., O'Reilly, D.: Privacy preserving record linkage in the presence of missing values. Inf. Syst. **71**, 199–210 (2017). https://doi.org/10.1016/j.is.2017.07.001. http://www.sciencedirect.com/science/article/pii/S030643791630504X
- Fulcher, B.D., Jones, N.S.: Highly comparative feature-based time-series classification. IEEE Trans. Knowl. Data Eng. 26(12), 3026–3037 (2014)
- SerrA, J., Arcos, J.L.: An empirical evaluation of similarity measures for time series classification. Knowl.-Based Syst. 67, 305–314 (2014). https://doi. org/10.1016/j.knosys.2014.04.035. http://www.sciencedirect.com/science/article/ pii/S0950705114001658
- Bellman, R., Kalaba, R.: On adaptive control processes. IRE Trans. Autom. Control. 4(2), 1–9 (1959)
- 16. Myers, C., Rabiner, L., Rosenberg, A.: Performance tradeoffs in dynamic time warping algorithms for isolated word recognition. IEEE Trans. Acoust. Speech Signal Process. **28**(6), 623–635 (1980)
- 17. Senin, P.: Dynamic time warping algorithm review. Inf. Comput. Sci. **855**(1–23), 40 (2008). Department University of Hawaii at Manoa Honolulu, USA
- 18. Chassiakos, Y.L.R., Radesky, J., Christakis, D., Moreno, M.A., Cross, C., et al.: Children and adolescents and digital media. Pediatrics 138(5), e20162593 (2016)
- 19. Ballano, S., Uribe, A.C., Munté-Ramos, R.A.: Young users and the digital divide: readers, participants or creators on internet? (2014)
- Miller, J.L., Paciga, K.A., Danby, S., Beaudoin-Ryan, L., Kaldor, T.: Looking beyond swiping and tapping: review of design and methodologies for researching young children's use of digital technologies. Cyberpsychology: J. Psychosoc. Res. Cyberspace 11(3), 6 (2017)