

# Literature Presentation

## Community Detection through Representation learning in Evolving Heterogenous Networks

Egor Dmitriev

Utrecht University

18 November 2021

- 1 Introduction & Surveys
- 2 Classical ML Methods
- 3 Deep Learning Based Methods
- 4 Deep Community Detection
- 5 Conclusion

## Introduction & Surveys

# Surveys

- [Rossetti and Cazabet, 2018] and [Dakiche et al., 2019]
- Categorize methods on problem of tracking community evolution
- Introduce a common definition for evolving communities
- Compile events/properties for analysis of evolving communities

# Evolving Communities

- Communities in real world:
  - disjoint (students belonging to different disciplines in an institute)
  - overlapping (person having membership in different social groups on Facebook)
  - hierarchical (cells in human body form tissues that in turn form organs and so on)
- Depend on underlying networks:
  - Time-series of static networks (Snapshots)
  - Real time a stream of edges (Temporal networks)
- Evaluated on synthetic (generated) communities
  - Usually based on a quality score:
    - Normalized Mutual Information score (NMI)
    - Modularity
    - etc.

- Quality function to evaluate algorithms **favors the ones that are designed to optimize it**

## Changes in Evolving Communities

- Operations that define community changes

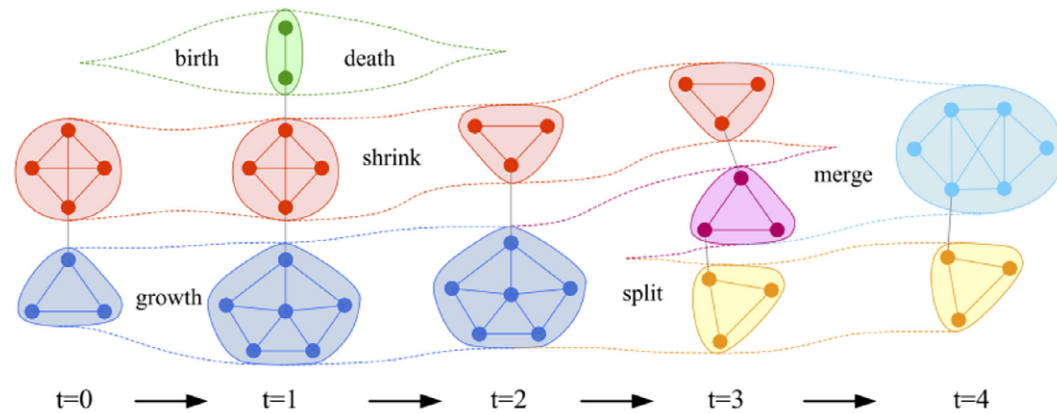


Fig. 3. Community evolution in a dynamic network (Shang, Liu, Li, Xie, & Wu, 2016).

## Challenges & Uses

- deciding if an ***element composed of several entities*** at a given instant is the same or not as another one composed of some—or even none—of such entities at a later point in time is ***necessarily arbitrary and cannot be answered unambiguously***
  - Main issues encountered by dynamic community detection approaches is the **instability of solutions**
  - Use cases:
    - forecasting emerging market trends in online retail networks
    - characterizing functions of unknown proteins
    - real-time partitioning of web-pages with different topics
    - predicting the emergence and lifetimes of overlapping communities in online social networks
- deciding if an ***element composed of several entities*** at a given instant is the same or not as another one composed of some—or even none—of such entities at a later point in time is ***necessarily arbitrary and cannot be answered unambiguously***
  - Main issues encountered by dynamic community detection approaches is the **instability of solutions**
  - Use cases:
    - forecasting emerging market trends in online retail networks
    - characterizing functions of unknown proteins
    - real-time partitioning of web-pages with different topics
    - predicting the emergence and lifetimes of overlapping communities in online social networks

# Classical ML Methods



## Louvain method (LM)

- Popular clustering algorithm
- Complexity is  $n \log n$
- Can be applied to weighted graphs
- Does not require a priori knowledge of the number of partitions

# Independent Community Detection and Matching

- **First detect communities** at each time step and **then match** them across different time-steps
- Unmodified traditional community detection methods can be reused
- Parallelism can be used for community detection
- Community detection algorithms are unstable leading to poor matching
- Examples:
  - **Sun, Tang, Pan, and Li (2015):**
    - Applied the Louvain algorithm to find the communities.
    - Then **built a correlation matrix** to between communities in  $t$  and  $t + 1$
  - **Greene et al. (2010):**
    - Using the static algorithm MOSES to detect the communities on each snapshot.
    - Then, they described a **weighted bipartite matching to map communities**
  - [Rossetti, 2020]
    - Allows for overlapping communities using modified node labeling algorithm
    - 
    - Matching based on multiple labels in  $t$ ,  $t-1$ ,  $t+1$
- Independent Community Detection and Matching:
  - Unmodified traditional community detection methods can be reused
  - Parallelism can be used for community detection.
  - Major drawback: Community detection algorithms are unstable.
- **Sun, Tang, Pan, and Li (2015):**
  - Applied the Louvain algorithm to find the communities.
  - Then **built a correlation matrix** to between communities in  $t$  and  $t + 1$
- **Greene et al. (2010):**
  - Using the static algorithm MOSES to detect the communities on each snapshot.
  - Then, they described a **weighted bipartite matching to map communities**
- [Rossetti, 2020]
  - Allows for overlapping communities using modified node labeling algorithm
  - Matching based on multiple labels in  $t$ ,  $t-1$ ,  $t+1$
  - focuses on lowering the time complexity while at the same time increasing the partition
  - Events are detected and evaluated against ground truth
  - Provides a deterministic output

# Dependent Community Detection

- **Detect communities at time  $t$  and then use them to detect communities at time  $t + 1$ ,**
  - Reduce computational cost but do not allow parallelism
  - Examples:
    - **Gao, Luo, and Bu (2016):**
      - Evolutionary community discovery algorithm based on **leader nodes**
      - Each **community is considered as a set of follower nodes** congregating close to a potential leader
- Dependent Community Detection
    - Reduce computational cost by reusing much of the previous community
    - Traditional community detection methods are no longer directly applicable
    - Does not allow parallelism in community detection
  - **Gao, Luo, and Bu (2016):**
    - Evolutionary community discovery algorithm based on **leader nodes**
    - Each **community is considered as a set of follower nodes** congregating close to a potential leader

# Simultaneous Community Detection on All Snapshots

- **Construct a single graph** and then run a classic community detection
  - Solution for the lack of stability of the independent community detection
- Simultaneous Community Detection on All Snapshots
    - **Main advantage:** is providing a solution for the lack of stability of the independent community detection
    - Difficulty to detect complex operations such as merging and splitting

## Dynamic Community Detection on Temporal Networks (online approach)

- Update the ones previously found according to network modifications
- Problem: Modifications are done at a **local level**
- Examples:
  - **Shang et al. (2012):**
    - Update graph real-time, and **locally modify the concerned communities** in a way to increase the modularity
  - **Held and Kruse (2016):**
    - Assumption that there exist some **highly connected nodes**, called hubs, which will group people around them.
  - [Xu et al., 2020]
    - A dynamic network snapshot is totally re-partitioned once the error accumulation degree of incremental clustering exceeds a pre-defined threshold
    - Use

- Dynamic Community Detection on Temporal Networks (online approach)
  - Since the communities evolve naturally through modifications, there is, **no longer, an instability problem**
  - Advantage: low complexity of tracking communities, since changes can be incremental
  - Problem: Modifications are done at a **local level**, they can **involve drifting towards invalid communities**
- **Shang et al. (2012):**
  - Method consists in adding (or removing) each new edge as it appears (or disappears), and to **locally modify the concerned communities** in a way to increase the modularity
- **Held and Kruse (2016):**
  - Based on the assumption that there exist some highly connected nodes, called hubs, which will group people around them.
  - It is based on the assumption that there exist some highly connected nodes, called hubs, which will group people around them. So, in the first step, the proposed algorithm detects these hubs by the node degree and assigns to all non-hub elements the closest hub as a cluster label, then iteratively changes the resultant clustering by applying changes: adding or removing nodes or edges
- [Xu et al., 2020]
  - A dynamic network snapshot is totally re-partitioned once the error accumulation degree of incremental clustering exceeds a pre-defined threshold

## Deep Learning Based Methods

# User community detection via embedding of social network structure and temporal content

- [Fani et al., 2020]
- Content on the social network are often reflective of issues in the real world **topics discussed on the network constantly change** and hence users' interests towards these topics
- Social Network Connections  $\mathcal{G} = (\mathbb{U}, \mathbb{A})$ 
  - Primarily based on the **homophily** principle
  - Explicit social connection does not necessarily indicate user interest similarity
- Combine both Temporal Social Content  $\mathcal{D} = (\mathbb{U}, \mathbb{M}, \mathbb{T})$ 
  - Model communities **based on topics of interest**

- Temporal Social Content:  $\mathcal{D} = (\mathbb{U}, \mathbb{M}, \mathbb{T})$ 
  - $\mathbb{U}$ : Users,  $\mathbb{M}$ : Text content,  $\mathbb{T}$  time periods
- Social Network Graph:  $\mathcal{G} = (\mathbb{U}, \mathbb{A})$ 
  - $\mathbb{U}$ : Users/Nodes,  $\mathbb{A}$ : Edges
- Homophily: (densely connected groups of users imply a user community)

## User community detection via embedding of social network structure and temporal content

- Identify Topics using LDA and Construct Preference Time series
- Learn dense representation of users interests using CBOW model
  - Use “Regions of like-mindedness” for as scoring function
- Topological Embeddings are constructed using a Skip-Gram model
  - Use DFS based random walk over the Social Network

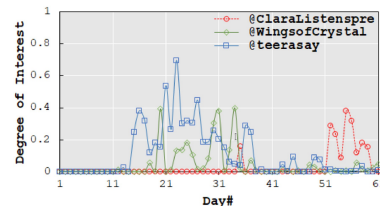


Figure 2: Different temporal behaviour of three Twitter users with respect to the War in Afghanistan topic

- **Region of like-mindedness:** Parts in  $X$  where users share interest in same topics given a threshold (for level of interest)
- BFS favours structural equivalence
- DFS in contrast, respects homophily and leads to similar (close) embeddings for densely connected users



# User community detection via embedding of social network structure and temporal content

- Finally embeddings are combined into one  $h(\mathbf{W}_{\mathcal{D}}, \mathbf{W}_{\mathcal{G}}) = \alpha \mathbf{W}_{\mathcal{D}} + (1 - \alpha) \mathbf{W}_{\mathcal{G}}$
- Community detection:
  - Construct a weighted graph:  $G = (\mathbb{U}, \mathbb{E}, w)$
  - Leverage the Louvain Method (LM)
- Remarks:
  - Users end up in one community per users
  - Consider negative sampling

- With as weights the user embedding **dot products**
-

# Vehicle Trajectory Clustering Based on Dynamic Representation Learning of Internet of Vehicles

- [Wang et al., 2020]
- Vehicle trajectory clustering aims to regroup similar vehicle trajectories together into different groups
- Challenge: As the location of vehicles is constantly changing, the vehicle social network is a dynamic network
- Social Network is constructed from trajectories:
  - Discretize vehicle positions using a Grid
  - Given vehicle  $n$  at time  $t$ , connect  $k$  closest vehicles
- Vehicle trajectory clustering aims to regroup similar vehicle trajectories together into different groups
  - Extract relevant information in order to, for instance, calculate the optimal path from one position to another, detect abnormal behavior, monitor the traffic flow in a city, and predict the next position of an object
  - The road networks of different city regions may be totally different
  - Vehicle may present totally different trajectories over different time periods of a day
  - Meanwhile, the patterns on weekdays and weekends may also differ.

# Vehicle Trajectory Clustering Based on Dynamic Representation Learning of Internet of Vehicles

- Authors propose: DynWalks
  - Performs truncated random walks of length  $l$
  - **performs random walks on selected nodes**
    - Embedding vectors of other nodes remains unchanged
- Calculate vehicle representations using Skip-Gram Negative Sampling
- Finally clustering is done:
  - K-means, K-medoids, GMM
  - For each timestep

- Performs truncated random walks with length  $l$  on each selected node for  $r$  times
- By using a silding window with length  $w + 1 + w$  to slide on each random walk sequence

# Deep Community Detection

# Deep Learning for Community Detection: Progress, Challenges and Opportunities

- [Liu et al., 2020]
- Challenges:
  - An Unknown Number of Communities
  - Network Heterogeneity
  - Large-scale Networks
- [Bhatia and Rani, 2018]
  - DeCom finds the number of clusters by analyzing the structure of the network
  - Initial clusters are found using PageRank
  - Clusters are refined using GAE by optimizing for modularity
- [Cavallari et al., 2017]
  - Suboptimal: community detection is independent of its node embedding
  - Communities are seen as distributions over the embeddings
  - Introduces two step process
    - Construct initial embeddings
    - Refine embeddings and the community distributions (GMM) alternately
- An Unknown Number of Communities:
  - [Bhatia and Rani, 2018] based on random walk-based personalized PageRank. However, this type of method cannot guarantee that every node in the network is assigned to a community
- Network Heterogeneity:
  - Network heterogeneity refers to networks that contain significantly different types of entities and relationships, which means the strategies used for homogeneous networks do not necessarily work.
- Large-scale Networks:
  - Today, large-scale networks can contain millions of nodes, edges, and structural patterns and can also be highly dynamic, as networks like Facebook and Twitter demonstrate.
- GAE: Graph Auto Encoders
- [Cavallari et al., 2017]
  - uses a unified objective where embedding and community distributions are refined

## Conclusion

# Conclusion

- There are still few methods that
  - Utilize deep representation
  - Use representation vectors for clustering
  - Explore multimodal settings

## References I

## References

- S. Cavallari, V. W. Zheng, H. Cai, K. C.-C. Chang, and E. Cambria. Learning Community Embedding with Community Detection and Node Embedding on Graphs. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 377–386, Singapore Singapore, Nov. 2017. ACM. ISBN 978-1-4503-4918-5. doi: 10.1145/3132847.3132925.
- N. Dakiche, F. Benbouzid-Si Tayeb, Y. Slimani, and K. Benatchba. Tracking community evolution in social networks: A survey. *Information Processing & Management*, 56(3): 1084–1102, May 2019. ISSN 0306-4573. doi: 10.1016/j.ipm.2018.03.005.



# References II

H. Fani, E. Jiang, E. Bagheri, F. Al-Obeidat, W. Du, and M. Kargar. User community detection via embedding of social network structure and temporal content. *Information Processing & Management*, 57(2):102056, Mar. 2020. ISSN 03064573. doi: 10.1016/j.ipm.2019.102056.

F. Liu, S. Xue, J. Wu, C. Zhou, W. Hu, C. Paris, S. Nepal, J. Yang, and P. S. Yu. Deep Learning for Community Detection: Progress, Challenges and Opportunities. *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, pages 4981–4987, July 2020. doi: 10.24963/ijcai.2020/693.

G. Rossetti. ANGEL: Efficient, and effective, node-centric community discovery in static and dynamic networks. *Applied Network Science*, 5(1):26, June 2020. ISSN 2364-8228. doi: 10.1007/s41109-020-00270-6.

# References III

G. Rossetti and R. Cazabet. Community Discovery in Dynamic Networks: A Survey. *ACM Computing Surveys*, 51(2):35:1–35:37, Feb. 2018. ISSN 0360-0300. doi: 10.1145/3172867.

W. Wang, F. Xia, H. Nie, Z. Chen, Z. Gong, X. Kong, and W. Wei. Vehicle Trajectory Clustering Based on Dynamic Representation Learning of Internet of Vehicles. *IEEE Transactions on Intelligent Transportation Systems*, PP:1–10, June 2020. doi: 10.1109/TITS.2020.2995856.

Z. Xu, X. Rui, J. He, Z. Wang, and T. Hadzibeganovic. Superspreaders and superblockers based community evolution tracking in dynamic social networks. *Knowledge-Based Systems*, 192:105377, Mar. 2020. ISSN 0950-7051. doi: 10.1016/j.knosys.2019.105377.