







EagleMine: vision-guided large graph mining

<u>Wenjie Feng</u>⁺ Shenghua Liu⁺ Christos Faloutsos^{*} Bryan Hooi^{*} Huawei Shen⁺ Xueqi Cheng⁺

*Institute of Computing Technology ICT, CAS*Computer Science Dept., Carnegie Mellon University





Human Healthcare



Human Healthcare





Human Healthcare



Chest CT scan





Human Healthcare



Chest CT scan





[Liver Cancer]



Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]



Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]

How about very graph diagnosis?



Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]

How about very graph diagnosis?



Social Network [Facebook]



Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]

How about very graph diagnosis?



View spaces





Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]

How about very graph diagnosis?



Social Network [Facebook]



Anomaly Patterns



Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]

How about very graph diagnosis?





10E-4 0.02



Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]

How about very graph diagnosis?





Human Healthcare



Chest CT scan



Cancer / Tumor



[Liver Cancer]

How about very graph diagnosis?



suspicious users:

182x: "best*"

178x: "black*"

105x: "blue*"

51x: "coolboy*" **223x:** "18-year-old*"

6% (deleted next year)



Outline

- Introduction
- Overview <<</p>
- Proposed Method: EagleMine
- Experiments
- Conclusion



 Our EagleMine: Automatically summarize the node distribution in graph correlated feature spaces.





Detected micro-clusters







EagleMine achieves shortest MDL.

Quantitative



STING





DBSCAN













EagleMine achieves shortest MDL.

Quantitative



STING



manually tuned para.

DBSCAN



Watershed Qualitative





 Our EagleMine: Detect explainable anomaly pattern more effective and scalable for running linearly in # of nodes .





Outline

- Introduction
- Overview
- Proposed Method: EagleMine <<</p>
 - Problem definition <
 - Algorithm
- Experiments
- Conclusion

Graph-Theoretic Scagnostics

graph-theoretic summaries of high-dimensional scattered point data.



Graph-Theoretic Scagnostics

 graph-theoretic summaries of high-dimensional scattered point data.



scagnostics. [L. Wikinson et al, INFOVIS'05], [J. W. Tukey et al NCGA'85]



9



Graph diagnosis ('CT' scans)

Graph nodes distribution in correlated feature spaces.



Graph diagnosis ('CT' scans)

Graph nodes distribution in correlated feature spaces. E.g. Degree (in- / out- degree); # Triangle; Coreness; PageRank: (undirected homogeneous graph $\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$) spectral vectors: $\mathbf{A} \approx \sigma \cdot \mathbf{u} \cdot \mathbf{v}^T$ (\mathbf{u} : hubness / \mathbf{v} : authority), etc.



Graph diagnosis ('CT' scans)

Graph nodes distribution in correlated feature spaces.

E.g. Degree (in- / out- degree); # Triangle; Coreness;

PageRank: (undirected homogeneous graph $\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$) spectral vectors: $\mathbf{A} \approx \sigma \cdot \mathbf{u} \cdot \mathbf{v}^T$ (\mathbf{u} : hubness / \mathbf{v} : authority), etc.





Graph diagnosis ('CT' scans)

Graph nodes distribution in correlated feature spaces. E.g. Degree (in- / out- degree); # Triangle; Coreness; PageRank: (undirected homogeneous graph $Ax = \lambda x$)

spectral vectors: $\mathbf{A} \approx \sigma \cdot \mathbf{u} \cdot \mathbf{v}^T$ (u: hubness / v: authority), etc.





High-dimension spaces



Graph diagnosis ('CT' scans)

Graph nodes distribution in correlated feature spaces.

E.g. Degree (in- / out- degree); # Triangle; Coreness;

PageRank: (undirected homogeneous graph $\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$) spectral vectors: $\mathbf{A} \approx \sigma \cdot \mathbf{u} \cdot \mathbf{v}^T$ (\mathbf{u} : hubness / \mathbf{v} : authority), etc.





High-dimension spaces



• Vision-guided diagnosis & mining for large graph:



Vision-guided diagnosis & mining for large graph:
Given:



Vision-guided diagnosis & mining for large graph: Given:

- Graph $\mathcal{G}\!=(V,E)$, $\boldsymbol{\mathsf{A}}$ is the adjacency matrix
- Node features histogram $\ensuremath{\mathcal{H}}$



Vision-guided diagnosis & mining for large graph: Given:

- Graph $\mathcal{G}\!=(V,E)$, $\boldsymbol{\mathsf{A}}$ is the adjacency matrix
- Node features histogram $\ensuremath{\mathcal{H}}$





Vision-guided diagnosis & mining for large graph: Given:

- Graph $\mathcal{G}\!=(V,E)$, $\boldsymbol{\mathsf{A}}$ is the adjacency matrix
- Node features histogram $\ensuremath{\mathcal{H}}$

Goal (automatically):



Vision-guided diagnosis & mining for large graph: Given:

- Graph $\mathcal{G} = (\mathbf{V}, \mathbf{E})$, $\boldsymbol{\mathsf{A}}$ is the adjacency matrix
- Node features histogram $\ensuremath{\mathcal{H}}$

Goal (automatically):

- I. recognize and monitor node groups;
- 2. *summary* graph nodes in feature space & *identify* suspicious micro-clusters





Vision-guided diagnosis & mining for large graph: Given:

- Graph $\mathcal{G}\!=(V,E)$, $\boldsymbol{\mathsf{A}}$ is the adjacency matrix
- Node features histogram $\ensuremath{\mathcal{H}}$

Goal (automatically):

- I. recognize and monitor node groups;
- 2. *summary* graph nodes in feature space & *identify* suspicious micro-clusters





Vision-guided diagnosis & mining for large graph: Given:

- Graph $\mathcal{G} = (\mathbf{V}, \mathbf{E})$, $\boldsymbol{\mathsf{A}}$ is the adjacency matrix
- Node features histogram $\ensuremath{\mathcal{H}}$

Goal (automatically):

- I. recognize and monitor node groups;
- 2. summary graph nodes in feature space & identify suspicious micro-clusters

Solution: EagleMine algorithm








Outline

- Introduction
- Overview
- Proposed Method: EagleMine <<</p>
 - Problem definition
 - Algorithm <
- Experiments
- Conclusion





Vision-guided diagnosis & mining for large graph:



- Vision-guided diagnosis & mining for large graph:
 - I. Recognize node groups with Hierarchical tree structure for ${\cal H}$

WATERLEVELTREE algorithm



- Vision-guided diagnosis & mining for large graph:
 - I. Recognize node groups with Hierarchical tree structure for ${\cal H}$

WATERLEVELTREE algorithm

 2. Search the tree and get summarization of the histogram: TREEXPLORE algorithm



- Vision-guided diagnosis & mining for large graph:
 - I. Recognize node groups with Hierarchical tree structure for ${\cal H}$

WATERLEVELTREE algorithm

 2. Search the tree and get summarization of the histogram: TREEXPLORE algorithm

Vocabulary-based summarization model

Hypothesis test selection



- **Vision-guided** diagnosis & mining for large graph:
 - I. Recognize node groups with Hierarchical tree structure for ${\cal H}$

WATERLEVELTREE algorithm

 2. Search the tree and get summarization of the histogram: TREEXPLORE algorithm

Vocabulary-based summarization model

Hypothesis test selection



1. Recognize micro-clusters WATERLEVELTREE



- 1. Recognize micro-clusters WATERLEVELTREE
 - I. Build waterlevel tree;



- 1. Recognize micro-clusters WATERLEVELTREE
 - I. Build waterlevel tree;





- 1. Recognize micro-clusters WATERLEVELTREE
 - I. Build waterlevel tree;





- 1. Recognize micro-clusters WATERLEVELTREE
 - I. Build waterlevel tree;





1. Recognize micro-clusters WATERLEVELTREE







1. Recognize micro-clusters WATERLEVELTREE

0





1. Recognize micro-clusters WATERLEVELTREE

0





1. Recognize micro-clusters WATERLEVELTREE







1. Recognize micro-clusters WATERLEVELTREE







1. Recognize micro-clusters WATERLEVELTREE







1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





- 1. Recognize micro-clusters WATERLEVELTREE
 - I. **Build** waterlevel tree; II. Tree **refine**;





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE




1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE





1. Recognize micro-clusters WATERLEVELTREE









Flood fill, morphology operations [Rafael C Gonzalez et al 2007]



2. Summary node group model.



2. Summary node group model.

I. Vocabulary-based summarization model;



- 2. Summary node group model.
 - I. Vocabulary-based summarization model;
 - use statistical distribution describe node in island.



- 2. Summary node group model.
 - I. **Vocabulary**-based summarization model; use statistical distribution describe node in island.
 - Vocabulary: distributions $\mathcal Y$



- 2. Summary node group model.
 - I. **Vocabulary**-based summarization model; use statistical distribution describe node in island.
 - Vocabulary: distributions $\mathcal Y$

E.g. Gaussian, Uniform, Laplacian, Exp. etc.



- 2. Summary node group model.
 - I. Vocabulary-based summarization model;
 - use statistical distribution describe node in island.
 - Vocabulary: distributions \mathcal{Y}

E.g. Gaussian, Uniform, Laplacian, Exp. etc.





- 2. Summary node group model.
 - I. Vocabulary-based summarization model;
 - use statistical distribution describe node in island.
 - Vocabulary: distributions $\mathcal Y$
 - E.g. Gaussian, Uniform, Laplacian, Exp. etc.

Discretized, Truncated, Multivariate DTM gaussian





- 2. Summary node group model.
 - I. Vocabulary-based summarization model;
 - use statistical distribution describe node in island.
 - Vocabulary: distributions $\mathcal Y$
 - E.g. Gaussian, Uniform, Laplacian, Exp. etc.

Discretized, Truncated, Multivariate DTM gaussian





- 2. Summary node group model.
 - I. Vocabulary-based summarization model;
 - use statistical distribution describe node in island.
 - Vocabulary: distributions \mathcal{Y}

E.g. Gaussian, Uniform, Laplacian, Exp. etc.

Discretized, Truncated, Multivariate DTM gaussian





- 2. Summary node group model.
 - I. **Vocabulary**-based summarization model; use statistical distribution describe node in island.
 - Vocabulary: distributions $\mathcal Y$



- 2. Summary node group model.
 - I. **Vocabulary**-based summarization model; use statistical distribution describe node in island.
 - Vocabulary: distributions \mathcal{Y}
 - **Description**: For *C* node groups:



- 2. Summary node group model.
 - I. Vocabulary-based summarization model;
 - use statistical distribution describe node in island.
 - Vocabulary: distributions \mathcal{Y}
 - **Description**: For *C* node groups:
 - Assignments: $S = \{s_1, \cdots, s_C\}$
 - Model parameter: $\Theta = \{\theta_1, \cdots, \theta_C\}$
 - Outliers: unassigned bins $\ensuremath{\mathcal{O}}$



2. Summary node group model.



2. Summary node group model.



2. Summary node group model.





2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_{lpha} Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

II. WaterLevel Tree **Explore**;

- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

- II. WaterLevel Tree **Explore**;
- BFS search
- Determine optimal node groups with Hypothesis Test

Vocabulary Assignment: s_lpha Heuristic rule / Pearson's χ^2





2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test




2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test
- Stitch for enhancing





2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test
- Stitch for enhancing







2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test
- Stitch for enhancing







2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test
- Stitch for enhancing







2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test
- Stitch for enhancing







2. Summary node group model.

- BFS search
- Determine optimal node groups with Hypothesis Test
- Stitch for enhancing







3. Identify suspicious micro-clusters.



3. Identify suspicious micro-clusters.

Node group **suspiciousness score**:



3. Identify suspicious micro-clusters.

Node group **suspiciousness score**:

Weighted probability KL distance with the majority island.



3. Identify suspicious micro-clusters.

Node group **suspiciousness score**:

Weighted probability *KL distance* with the majority island.

$$\kappa(\theta_i) = \log \bar{d}_i \cdot \sum_{\boldsymbol{b}} N_i \cdot KL(P_{\theta_i}(\boldsymbol{b}) || P_{\theta_m}(\boldsymbol{b}))$$



3. Identify suspicious micro-clusters.

Node group suspiciousness score:

Weighted probability KL distance with the majority island.

$$\kappa(\theta_i) = \log \bar{d}_i \cdot \sum_{\boldsymbol{b}} N_i \cdot KL(P_{\theta_i}(\boldsymbol{b}) || P_{\theta_m}(\boldsymbol{b}))$$

The majority island – Normal nodes



DTM Gaussian Description





Outline

- Introduction
- Overview
- Proposed Method: EagleMine
- Experiments <<</p>
- Conclusion



Data sets

Table I. Dataset statistic information

	# of nodes	# of edges	Content
Amazon rating	(2.14M, 1.23M)	5.84M	Rate
Android	(1.32M, 61.27K)	2.64M	Rate
BeerAdvocate	(33.37K, 65.91K)	1.57M	Rate
Yelp	(686K, 85.54K)	2.68M	Rate
Tagged	(2.73M, 4.65M)	150.8M	Anonymized Links
Youtube	(3.22M, 3.22M)	9.37M	Who-follow-who
Sina weibo	(2.75M, 8.08M)	50.1M	User-retweet-msg



Data sets

Table I. Dataset statistic information

	# of nodes	# of edges	Content
Amazon rating	(2.14M, 1.23M)	5.84M	Rate
Android	(1.32M, 61.27K)	2.64M	Rate
BeerAdvocate	(33.37K, 65.91K)	1.57M	Rate
Yelp	(686K, 85.54K)	2.68M	Rate
Tagged	(2.73M, 4.65M)	150.8M	Anonymized Links
Youtube	(3.22M, 3.22M)	9.37M	Who-follow-who
Sina weibo	(2.75M, 8.08M)	50.1M	User-retweet-msg

Related publications [J. McAuley, R. Pandey & J. Leskovec, KDD'15], [J. McAuley & J.
Leskovc, WWW'13], [A. Mislove, M. Marcon et al, SIGCOM'07], [S. Fakraei et al KDD'15]
Most of datasets are public available at <u>https://snap.stanford.edu/data/index.html</u>



Exp I. Quantitative evaluation of summarization

Measure: MDL (minimum Description Length)

The best model has shortest MDL.





Exp I. Quantitative evaluation of summarization

Measure: MDL (minimum Description Length)

The best model has shortest MDL.



EagleMine achieves concise summarization of graph.



Sina_weibo user: out-degree vs. hubness





Sina_weibo user: out-degree vs. hubness





Sina_weibo msg: in-degree vs. authority





Tagged friendship: #triangle vs. degree









DBSCAN







Exp 3. Anomaly detection

Sina Weibo dataset: (user-retweet-msg)

user: 2.75M # msg: 8.08M # edge: 50.1M

Time: Nov. 1-30, 2013



Suspicious Users





Exp 3. Anomaly detection

Sina Weibo dataset: (user-retweet-msg)



GetScoop [M. Jiang et al PAKDD'14], SpokEn [B.Prakash PAKDD'10], Fraudar [B. Hooi et al KDD'16]



Exp 4. Scalability test

EagleMine time complexity:





Exp 4. Scalability test

EagleMine time complexity:





Exp 5. Network pattern discovery

User-Retweet-Msg Graph





Jellyfish structure anomaly pattern



Exp 5. Network pattern discovery

User-Retweet-Msg Graph





Jellyfish structure anomaly pattern



Exp 5. Network pattern discovery

User-Retweet-Msg Graph





Jellyfish structure anomaly pattern



Outline

- Introduction
- Overview
- Proposed Method: EagleMine
- Experiments
- Conclusion <<</p>

Conclusion

EagleMine: vision-guided large graph mining

Automated summarization:

graph nodes distribution in correlated feature spaces

- *Effectiveness:* describe and detect node groups, and get visually sense-making results
- Anomaly detection: spot explainable anomaly patterns
- Scalability: runs linearly in # of graph node
- <u>Reproducible</u>: open source code & data

GitHub https://github.com/wenchieh/eaglemine







taken (s)

Linea



Reference

- **[L. Wikinson et al, 05]** Leland Wilkinson, Anushka Anand, and Robert Grossman. Graph-theoretical scagnostics. Proceedings IEEE Symposium on Information Visualisation, INFOVIS'05.
- [J. W. Tukey et al, 85] J. W. Tukey and P. A. Tukey. Computer graphics and exploratory data analysis: An introduction. Nat Computer Graphics Association.
- **[L. Vincent 91]** Luc Vincent and Pierre Soille. Watersheds in digital spaces: an efficient algorithm based on immersion simulations. PAMI'91.
- [M. Ester et al 96] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density based algorithm for discovering clusters in large spatial databases with noise. In KDD'96.
- [W. Wang et al 97] Wei Wang, Jiong Yang, Richard Muntz, et al. STING: A statistical information grid approach to spatial data mining. In VLDB'97.
- [M. Jiang et al 14] Meng Jiang, Peng Cui, Alex Beutel, Christos Faloutsos, and Shiqiang Yang. Inferring strange behavior from connectivity pattern in social networks. In PAKDD'14.
- [B.Prakash 10] B Aditya Prakash, Ashwin Sridharan, Mukund Seshadri, Sridhar Machiraju, and Christos Faloutsos. Eigenspokes: Surprising patterns and scalable community chipping in large graphs. PAKDD'10.
- [B. Hooi et al. 16] Bryan Hooi, Hyun Ah Song, Alex Beutel, Neil Shah, Kijung Shin, and Christos Faloutsos.
 2016. Fraudar: Bounding graph fraud in the face of camouflage. In SIGKDD'16.
- [J. DiCarlo et al.12] James J. DiCarlo, Davide Zoccolan, and Nicole C. Rust. 2012. How Does the Brain Solve Visual Object Recognition? Neuron'12.
- [Rafael C Gonzalez et al 2007] Rafael C Gonzalez and Richard E Woods. 2007. Image processing. Digital image processing.
- [J. McAuley, R. Pandey & J. Leskovec, 15] Julian McAuley, Rahul Pandey, and Jure Leskovec. 2015.
 Inferring Networks of Substitutable and Complementary Products. In KDD.

Reference

- [J. McAuley & J. Leskovc, 13] Julian John McAuley and Jure Leskovec. 2013. From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews. In WWW.
- [A. Mislove, M. Marcon et al, 07] Alan Mislove, Massimiliano Marcon, Krishna P. Gummadi, Peter Druschel, and Bobby Bhattacharjee. 2007. Measurement and Analysis of Online Social Networks. In SIGCOMM.
- [S. Fakraei et al. 15] Shobeir Fakhraei, James Foulds, Madhusudana Shashanka, and Lise Getoor. 2015.
 Collective Spammer Detection in Evolving Multi-Relational Social Networks. In KDD'15.











Source codes and datasets used in the paper are available at https://github.com/wenchieh/eaglemine

2018/8/16