# Large Graph Mining

## *Christos Faloutsos*

## CMU

# Thank you!

- Hillol Kargupta

# Outline

- **Problem definition / Motivation**
- Static & dynamic laws; generators
- Tools: CenterPiece graphs; fraud detection
- Conclusions

C. Faloutsos                                      3

# Motivation

Data mining: ~ find patterns (rules, outliers)

- Problem#1: How do real graphs look like?
- Problem#2: How do they evolve?
- Problem#3: How to generate realistic graphs

TOOLS

- Problem#4: Who is the 'master-mind'?
- Problem#5: Fraud detection

# Problem#1: Joint work with

Dr. Deepayan Chakrabarti
(CMU/Yahoo R.L.)

# Graphs - why should we care?



Internet Map
[lumeta.com]



Food Web
[Martinez '91]



Friendship Network
[Moody '01]



Protein Interactions
[genomebiology.com]

# Graphs - why should we care?

- IR: bi-partite graphs (doc-terms)

$D_1$      $T_1$

...    ...

$D_N$      $T_M$

- web: hyper-text graph

- ... and more:

# Graphs - why should we care?

- network of companies & board-of-directors members

- 'viral' marketing

- web-log ('blog') news propagation

- computer network security: email/IP traffic and anomaly detection

- ....

C. Faloutsos

# Problem #1 - network and graph mining



- How does the Internet look like?

- How does the web look like?

- What is 'normal'/'abnormal'?

- which patterns/laws hold?

C. Faloutsos

# Graph mining

- Are real graphs random?

C. Faloutsos

# Laws and patterns

- Are real graphs random?

- A: NO!!
  - Diameter
  - in- and out- degree distributions
  - other (surprising) patterns

C. Faloutsos

# Solution#1

- Power law in the degree distribution [SIGCOMM99]

**internet domains**

**att.com**

log(degree)

**ibm.com**

-0.82

log(rank)

# Solution#1': Eigen Exponent *E*

Eigenvalue



Exponent = slope

$E = -0.48$

May 2001

Rank of decreasing eigenvalue

- A2: power law in the eigenvalues of the adjacency matrix

# But:

How about graphs from other domains?

C. Faloutsos

# More power laws:

- web hit counts [w/ A. Montgomery]

### Web Site Traffic

log(count)

Zipf

``ebay''

log(in-degree)

users

sites

# epinions.com

- who-trusts-whom [Richardson + Domingos, KDD 2001]

count



trusts-2000-people user

(out) degree

# Motivation

Data mining: ~ find patterns (rules, outliers)

✓ Problem#1: How do real graphs look like?

• Problem#2: How do they evolve?

• Problem#3: How to generate realistic graphs

TOOLS

• Problem#4: Who is the 'master-mind'?

• Problem#5: Fraud detection

# Problem#2: Time evolution

- with Jure Leskovec (CMU/MLD)



- and Jon Kleinberg (Cornell – sabb. @ CMU)

# Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:
  - diameter ~ O(log N)
  - diameter ~ O(log log N)
- What is happening in real data?

C. Faloutsos

# Evolution of the Diameter

- Prior work on Power Law graphs hints at **slowly growing diameter**:
  - diameter ~ O (log N )
  - diameter ~ O(log log N)
- What is happening in real data?
- Diameter **shrinks** over time

C. Faloutsos        20

# Diameter – ArXiv citation graph

- Citations among physics papers
- 1992 –2003
- One graph per year

diameter



time [years]

# Diameter – "Autonomous Systems"

- Graph of Internet
- One graph per day
- 1997 – 2000

C. Faloutsos

# Diameter – "Affiliation Network"

- Graph of collaborations in physics – authors linked to papers
- 10 years of data



diameter

Effective diameter vs time [years]

Legend:
- Full graph
- Post '95 subgraph
- Post '95 subgraph, no past

# Diameter – "Patents"

- Patent citation network
- 25 years of data

# Temporal Evolution of the Graphs

- N(t) … nodes at time t
- E(t) … edges at time t
- Suppose that

  $$N(t+1) = 2 * N(t)$$

- Q: what is your guess for

  $$E(t+1) =? 2 * E(t)$$

# Temporal Evolution of the Graphs

- N(t) … nodes at time t
- E(t) … edges at time t
- Suppose that

  $$N(t+1) = 2 * N(t)$$

- Q: what is your guess for

  $$E(t+1) = ? 2 * E(t)$$

- A: over-doubled!

  – But obeying the ``Densification Power Law''

# Densification – Physics Citations

- Citations among physics papers
- 2003:
  - 29,555 papers, 352,807 citations



$E(t)$

$N(t)$

# Densification – Physics Citations

- Citations among physics papers
- 2003:
  - 29,555 papers, 352,807 citations

$E(t)$



Apr 2003

1.69

Jan 1993

- Edges
- $= 0.0113 \ x^{1.69} \ R^2 = 1.0$

Number of nodes

$N(t)$

# Densification – Physics Citations

- Citations among physics papers
- 2003:
  - 29,555 papers, 352,807 citations

$E(t)$

Apr 2003

1.69

1: tree

Jan 1993

Number of edges

- Edges
- $= 0.0113 \, x^{1.69} \, R^2 = 1.0$

Number of nodes   $N(t)$

# Densification – Physics Citations

- Citations among physics papers
- 2003:
  - 29,555 papers, 352,807 citations

$E(t)$

clique: 2

1.69

Apr 2003

Jan 1993

* Edges

— $= 0.0113 \; x^{1.69} \; R^2 = 1.0$

Number of edges

Number of nodes

$N(t)$

C. Faloutsos

# Densification – Patent Citations

- Citations among patents granted
- 1999
  - 2.9 million nodes
  - 16.5 million edges
- Each year is a datapoint



E(t)

1.66

N(t)

# Densification – Autonomous Systems

- Graph of Internet
- 2000
  - 6,000 nodes
  - 26,000 edges
- One graph per day



$E(t)$

Number of edges

$10^{4.4}$
$10^{4.3}$
$10^{4.2}$
$10^{4.1}$

1.18

- Edges
— $= 0.87 \ x^{1.18} \ R^2 = 1.00$

$10^{3.5}$   $10^{3.6}$   $10^{3.7}$   $10^{3.8}$
Number of nodes

$N(t)$

# Densification – Affiliation Network

- Authors linked to their publications
- 2002
  - 60,000 nodes
    - 20,000 authors
    - 38,000 papers
  - 133,000 edges



$E(t)$

Number of edges

1.15

* Edges
— = 0.4255 $x^{1.15}$ $R^2$=1.0

Number of nodes

$N(t)$

# Motivation

Data mining: ~ find patterns (rules, outliers)

✓ Problem#1: How do real graphs look like?

✓ Problem#2: How do they evolve?

• Problem#3: How to generate realistic graphs

TOOLS

• Problem#4: Who is the 'master-mind'?

• Problem#5: Fraud detection

# Problem#3: Generation

- Given a growing graph with count of nodes $N_1$, $N_2$, ...

- Generate a realistic sequence of graphs that will obey all the patterns

C. Faloutsos

# Problem Definition

- Given a growing graph with count of nodes $N_1$, $N_2$, ...

- Generate a realistic sequence of graphs that will obey all the patterns
  - Static Patterns

    Power Law Degree Distribution

    Power Law eigenvalue and eigenvector distribution

    Small Diameter
  - Dynamic Patterns

    Growth Power Law

    Shrinking/Stabilizing Diameters

# Problem Definition

- Given a growing graph with count of nodes $N_1, N_2, ...$

- Generate a realistic sequence of graphs that will obey all the patterns

- **Idea: Self-similarity**

  – Leads to power laws

  – Communities within communities

  – ...

C. Faloutsos

# Kronecker Product – a Graph



Adjacency matrix

# Kronecker Product – a Graph

- Continuing multiplying with $G_1$ we obtain $G_4$ and so on …



$G_4$ adjacency matrix

C. Faloutsos

# Kronecker Product – a Graph

- Continuing multiplying with $G_1$ we obtain $G_4$ and so on …



$G_4$ adjacency matrix

C. Faloutsos

# Kronecker Product – a Graph

- Continuing multiplying with $G_1$ we obtain $G_4$ and so on …



$G_4$ adjacency matrix

C. Faloutsos

# Properties:

- ## We can PROVE that
  - Degree distribution is multinomial ~ power law
  - Diameter: constant
  - Eigenvalue distribution: multinomial
  - First eigenvector: multinomial
- ## See [Leskovec+, PKDD'05] for proofs

# Problem Definition

- Given a growing graph with nodes $N_1, N_2, ...$

- Generate a realistic sequence of graphs that will obey all the patterns
  - Static Patterns
    - ✓ Power Law Degree Distribution
    - ✓ Power Law eigenvalue and eigenvector distribution
    - ✓ Small Diameter
  - Dynamic Patterns
    - ✓ Growth Power Law
    - ✓ Shrinking/Stabilizing Diameters

- First and only generator for which we can **prove** all these properties

# (Q: how to fit the parm's?)

A:

- Stochastic version of Kronecker graphs +
- Max likelihood  +
- Metropolis sampling
- [Leskovec+, ICML'07]

C. Faloutsos

# Experiments on real AS graph

### Degree distribution



### Hop plot



### Adjacency matrix eigen values



### Network value



C. Faloutsos

# Conclusions

- Kronecker graphs have:
  - All the static properties
    - ✓ Heavy tailed degree distributions
    - ✓ Small diameter
    - ✓ Multinomial eigenvalues and eigenvectors
  - All the temporal properties
    - ✓ Densification Power Law
    - ✓ Shrinking/Stabilizing Diameters
  - We can formally prove these results

# Motivation

Data mining: ~ find patterns (rules, outliers)

✓ Problem#1: How do real graphs look like?

✓ Problem#2: How do they evolve?

✓ Problem#3: How to generate realistic graphs

TOOLS

➡ • Problem#4: Who is the 'master-mind'?

• Problem#5: Fraud detection

# Problem#4: MasterMind – 'CePS'

- w/ Hanghang Tong, KDD 2006

- htong <at> cs.cmu.edu

# **Center-Piece Subgraph(Ceps)**

- **Given** Q query nodes
- **Find** Center-piece ($\leq b$ )

- App.
  - Social Networks
  - Law Inforcement, …

- Idea:
  - Proximity -> random walk with restarts

# Case Study: AND query

R. Agrawal

Jiawei Han

V. Vapnik

M. Jordan

C. Faloutsos

# Case Study: AND query

# Case Study: AND query



H.V. Jagadish — 15 — R. Agrawal

H.V. Jagadish — 10 — Laks V.S. Lakshmanan

Laks V.S. Lakshmanan — 13 — Jiawei Han

10

R. Agrawal — 1 — Mannila

Jiawei Han — 1

R. Agrawal — 2

Mannila

Christos Faloutsos — 1 — Smyth

Smyth — 1

1

1

V. Vapnik

Christos Faloutsos — 1

1

M. Jordan

V. Vapnik — 4 — Corinna Cortes

Corinna Cortes — 6 — Pregibon

# Conclusions



- Q1:How to measure the importance?

- A1: RWR+K_SoftAnd

- Q2:How to do it efficiently?

- A2:Graph Partition (Fast CePS)
  - ~90% quality
  - 150x speedup (ICDM'06)

# Motivation

Data mining: ~ find patterns (rules, outliers)

✓ Problem#1: How do real graphs look like?

✓ Problem#2: How do they evolve?

✓ Problem#3: How to generate realistic graphs

TOOLS

✓ Problem#4: Who is the 'master-mind'?

• Problem#5: Fraud detection

C. Faloutsos

# E-bay Fraud detection

w/ Polo Chau &
Shashank Pandit, CMU

# E-bay Fraud detection

- lines: positive feedbacks
- would you buy from him/her?



C. Faloutsos

# E-bay Fraud detection

- lines: positive feedbacks
- would you buy from  him/her?

- or him/her?

C. Faloutsos

# E-bay Fraud detection - NetProbe

C. Faloutsos

# OVERALL CONCLUSIONS

- Graphs pose a wealth of fascinating problems

- self-similarity and power laws work, when textbook methods fail!

- New patterns (shrinking diameter!)

- New generator: Kronecker

# Promising directions

- Reaching out
  - Sociology, epidemiology; physics, ++…
  - Computer networks, security, intrusion det.
  - Num. analysis (tensors)

time

IP-source

IP-destination

# Promising directions – cont'd

- Scaling up, to Gb/Tb/Pb

  – Storage Systems

  – Parallelism (hadoop/map-reduce)

C. Faloutsos

# E.g.: self-* system @ CMU

- >200 nodes
- 40 racks of computing equipment
- 774kw of power.
- target: 1 PetaByte
- goal: self-correcting, self-securing, self-monitoring, self-...

C. Faloutsos

# DM for Tera- and Peta-bytes

Two-way street:

<- DM can use such infrastructures to find patterns

-> DM can help such infrastructures become self-healing, self-adjusting, 'self-*'

# References

- Hanghang Tong, Christos Faloutsos, and Jia-Yu Pan *Fast Random Walk with Restart and Its Applications* ICDM 2006, Hong Kong.

- Hanghang Tong, Christos Faloutsos *Center-Piece Subgraphs: Problem Definition and Fast Solutions,* KDD 2006, Philadelphia, PA

# References

- Jure Leskovec, Jon Kleinberg and Christos Faloutsos *Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations* KDD 2005, Chicago, IL. ("Best Research Paper" award).

- Jure Leskovec, Deepayan Chakrabarti, Jon Kleinberg, Christos Faloutsos *Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication* (ECML/PKDD 2005), Porto, Portugal, 2005.

# References

- Jure Leskovec and  Christos Faloutsos, *Scalable Modeling of Real Graphs using Kronecker Multiplication*, ICML 2007, Corvallis, OR, USA

- Shashank Pandit, Duen Horng (Polo) Chau, Samuel Wang and Christos Faloutsos *NetProbe: A Fast and Scalable System for Fraud Detection in Online Auction Networks* WWW 2007, Banff, Alberta, Canada, May 8-12, 2007.

- Jimeng Sun, Dacheng Tao, Christos Faloutsos *Beyond Streams and Graphs: Dynamic Tensor Analysis,* KDD 2006, Philadelphia, PA

# References

- Jimeng Sun, Yinglian Xie, Hui Zhang, Christos Faloutsos. *Less is More: Compact Matrix Decomposition for Large Sparse Graphs*, SDM, Minneapolis, Minnesota, Apr 2007. [pdf]

# THANK YOU!

Contact info:

www. cs.cmu.edu /~christos

(w/ papers, datasets, code, etc)