# 1. Subgraph Analysis

# 2. Propagation Methods

# 3. **Latent Factor Models**
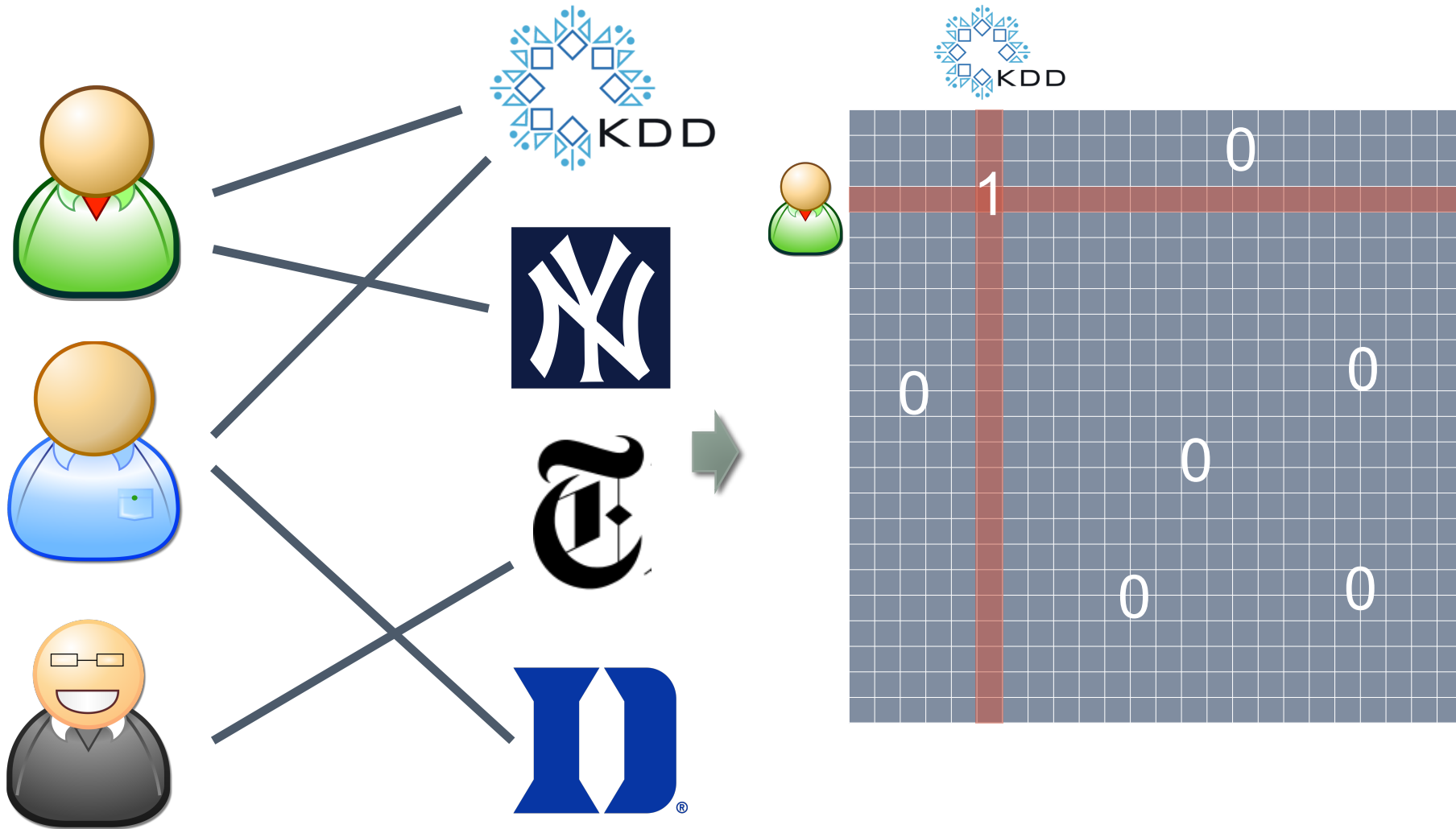
## a) Background

## b) Normal Behavior

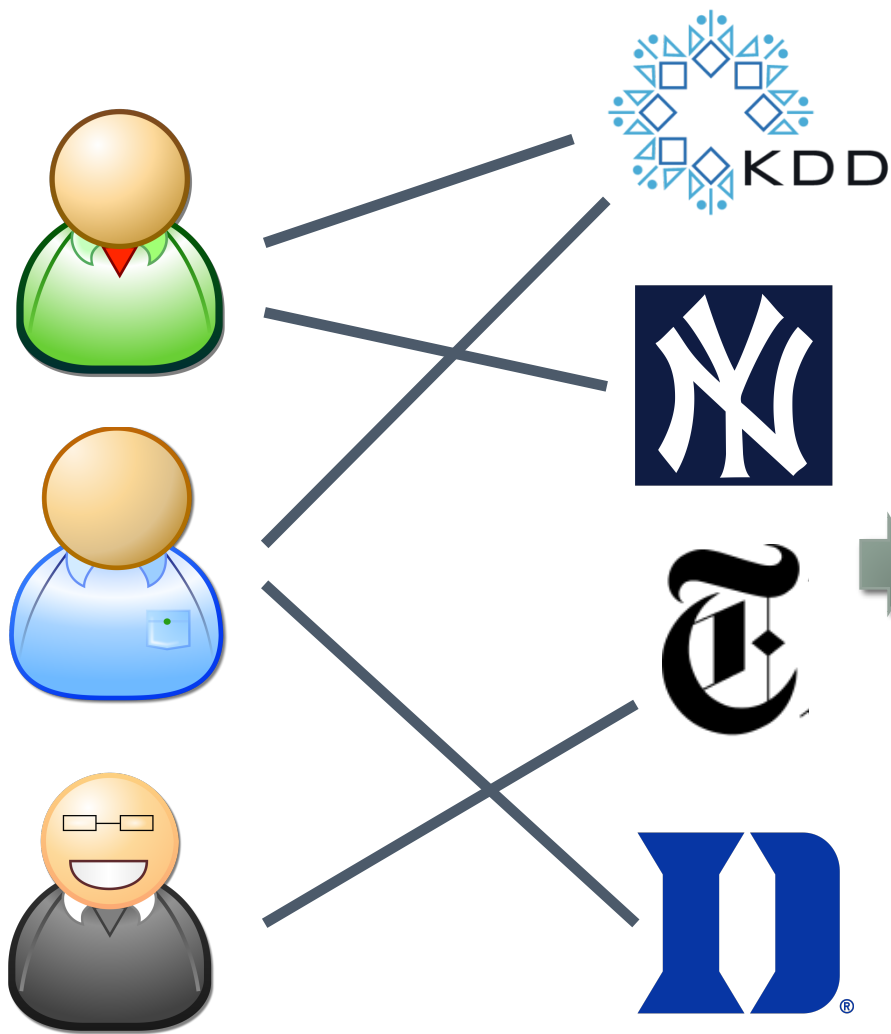## c) Abnormal Behavior
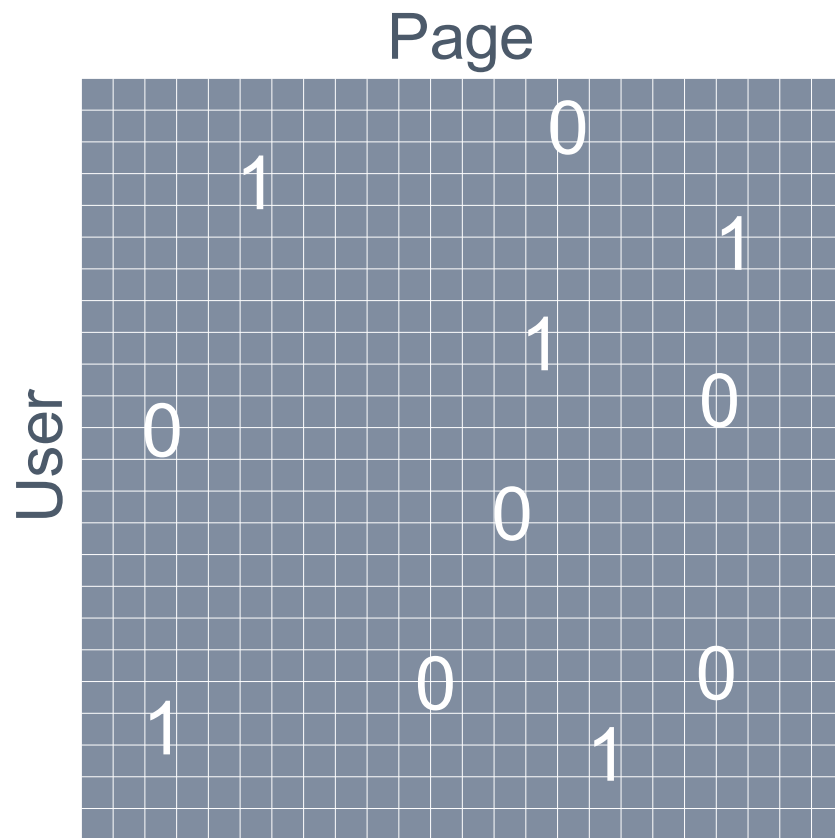
# Matrix Modeling

# Matrix Modeling

# Matrix Modeling

Matrix *M*

## HITS

Authoritativeness $\vec{v}$ is first eigenvector of $M^{\mathrm{T}}M$

$$\vec{v} = cM^{\mathrm{T}}M\vec{v}$$

Hubness $\vec{u}$ is first eigenvector of $MM^{\mathrm{T}}$

$$\vec{u} = cMM^{\mathrm{T}}\vec{u}$$

Page

User

|   |   |   | 0 |   |   |
|---|---|---|---|---|---|
| 1 |   |   |   |   |   |
|   |   |   |   |   | 1 |
|   |   |   | 1 |   |   |
| 0 |   |   |   |   | 0 |
|   |   |   | 0 |   |   |
|   |   |   | 0 |   | 0 |
| 1 |   |   |   | 1 |   |

# Matrix Modeling

Matrix *M*

## HITS

Authoritativeness $\vec{v}$ is first eigenvector of $M^{\mathrm{T}}M$

$$\vec{v} = cM^{\mathrm{T}}M\vec{v}$$

Hubness $\vec{u}$ is first eigenvector of $MM^{\mathrm{T}}$

$$\vec{u} = cMM^{\mathrm{T}}\vec{u}$$

Page

User

| | | | | | 0 | | | |
| | 1 | | | | | | | 1 |
| | | | | | 1 | | | |
| 0 | | | | | | | 0 | |
| | | | | 0 | | | | |
| | | | | | | 0 | | 0 |
| 1 | | | | | | 1 | | |

What about the other eigenvectors?

# Matrix Modeling
# Singular Value Decomposition



$$U\Sigma V^{\mathrm{T}} \approx M$$

# Matrix Modeling
# Singular Value Decomposition



$$U\Sigma V^{\mathrm{T}} \approx M$$

Hubness $\vec{u}$

Authoritativeness $\vec{v}$

# Matrix Modeling
# Singular Value Decomposition



$$U\Sigma V^{\mathrm{T}} \approx M$$

Hubness $\vec{u}$

Authoritativeness $\vec{v}$

$\Sigma$ contains normalization for $\vec{u}$ and $\vec{v}$

# Matrix Factorization

What does each eigenvector capture?
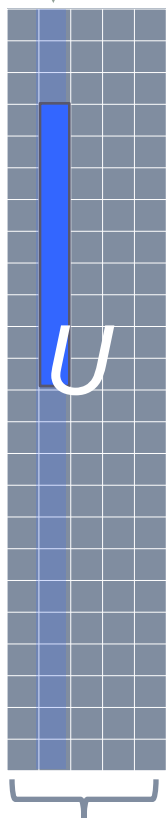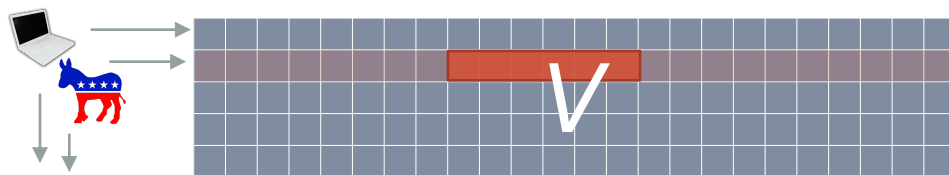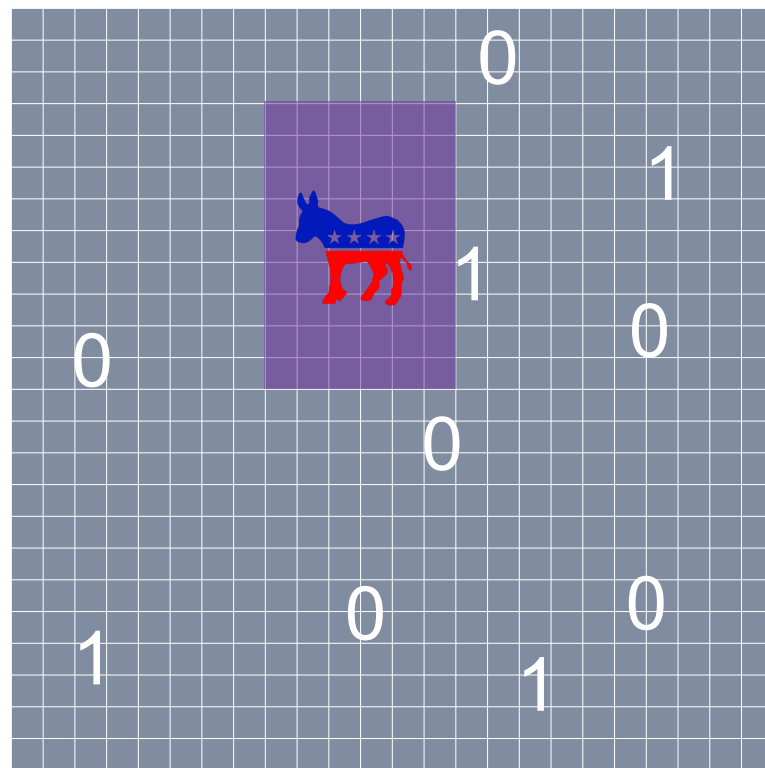


$$UV^{\mathrm{T}} \approx M$$

Each factor captures a dense block in the matrix

Topics

# Matrix Factorization

What does each eigenvector capture?



$$UV^{\mathrm{T}} \approx M$$

Each factor captures a dense block in the matrix

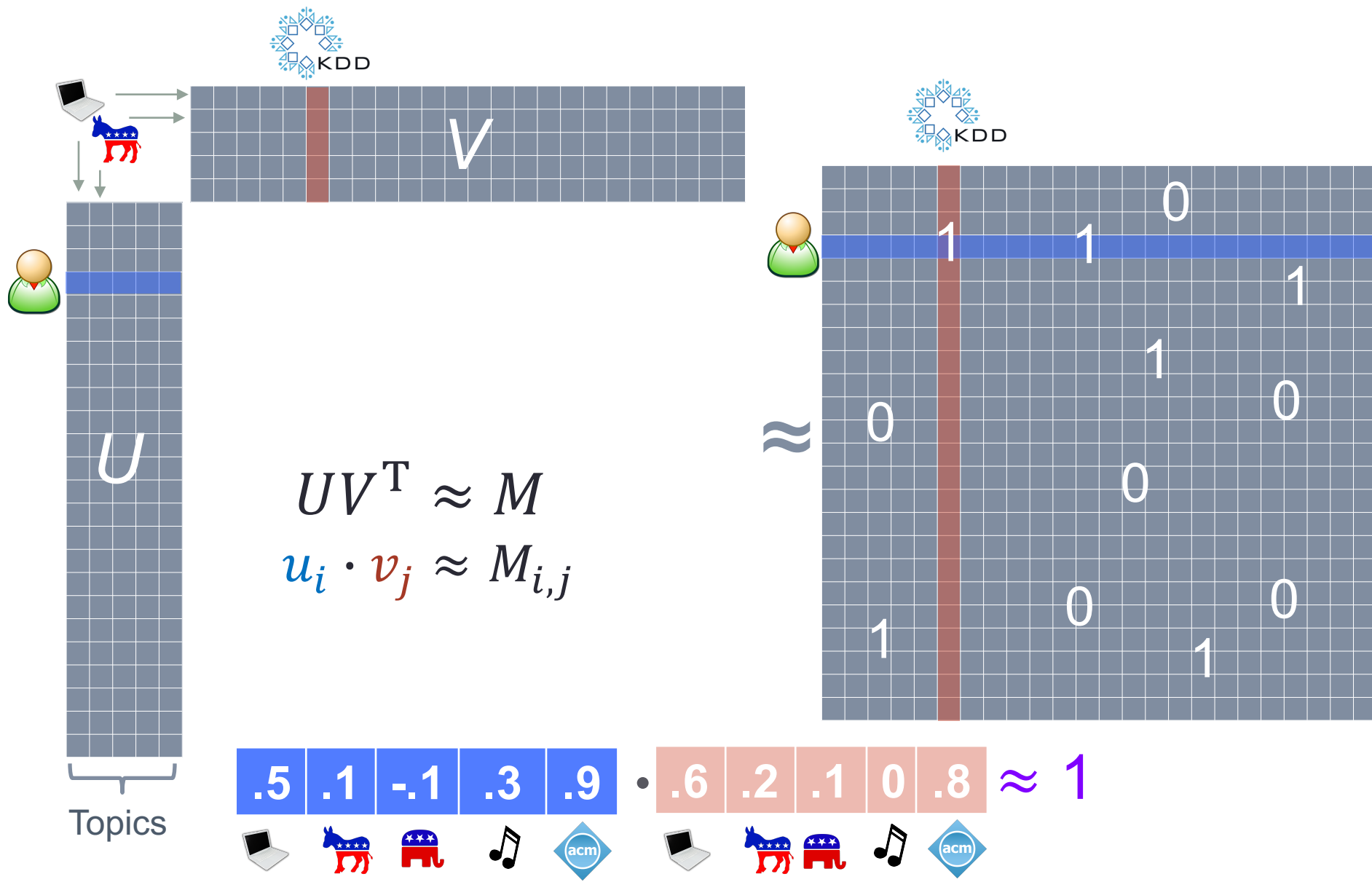Topics

# Matrix Factorization

$$UV^{\mathrm{T}} \approx M$$

$$u_i \cdot v_j \approx M_{i,j}$$

$V$

$U$

$\approx$

# Matrix Factorization



$$UV^{\mathrm{T}} \approx M$$

$$u_i \cdot v_j \approx M_{i,j}$$

Topics

# Matrix Factorization



$$UV^{\mathrm{T}} \approx M$$

$$u_i \cdot v_j \approx M_{i,j}$$

Topics

| .5 | .1 | -.1 | .3 | .9 | • | .6 | .2 | .1 | 0 | .8 | ≈ 1 |

# 1. Subgraph Analysis

# 2. Propagation Methods

# **3. Latent Factor Models**

## a) Background

## b) Normal Behavior

## c) Abnormal Behavior

# Matrix Completion

$$V$$

$$U$$

$$\approx$$

Movies

Users

| | | | | | | | ? | | |
| 1 | | | | | | | | 2 | |
| | ? | | | 1 | | | | | |
| ? | | | | | | | ? | | |
| | | | | ? | | | | | |
| | | ? | | | | | ? | |
| 5 | | | | | | 4 | | |

Recommendation systems
predict missing entries

# Matrix Completion



$V$

$U$

Can't find singular vectors with missing entries. Instead,

$$\min_{U,V} \sum_{(i,j)\in M} (M_{i,j} - \vec{u}_i \cdot \vec{v}_j)^2$$

$\approx$

1    ?

2

?    1

?    ?

?

?    ?

5    ?    4

# Matrix Completion



$V$

$U$

Genres

$\approx$

Can't find singular vectors with missing entries. Instead,

$$\min_{U,V} \sum_{(i,j) \in M} (M_{i,j} - \vec{u}_i \cdot \vec{v}_j)^2$$

# Matrix Completion



$V$

$U$

Can't find singular vectors with missing entries. Instead,

$$\min_{U,V} \sum_{(i,j)\in M} (M_{i,j} - \vec{u}_i \cdot \vec{v}_j)^2$$

Genres

| 1.2 | -.1 | .5 | .8 | -.5 |

$\cdot$

| .6 | .8 | 0 | 0 | .1 |

$\approx 1$

$\approx$

# Matrix Completion



Can't find singular vectors
with missing entries.  Instead,

$$\min_{U,V} \sum_{(i,j)\in M} (M_{i,j} - \widehat{M}_{i,j})^2$$

$$\widehat{M}_{i,j} = \vec{u}_i \cdot \vec{v}_j$$

Genres

# Adding Latent Factors



$$\min_{U,V} \sum_{(i,j)\in M} (M_{i,j} - \widehat{M}_{i,j})^2$$

Consider additional factors:
- Dataset mean $\mu$
- Row (user) baseline $b_i$
- Column (movie) baseline $b_j$

$$\widehat{M}_{i,j} = \mu + b_i + b_j + \vec{u}_i \cdot \vec{v}_j$$

# Adding Latent Factors

What if we know the time of the rating
(time of the edge being created)?

# Adding Latent Factors

### Mean Rating by Date (Netflix)



Mean Score — Time (days)

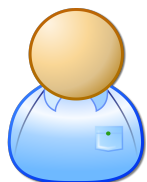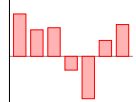Collaborative Filtering with Temporal Dynamics
Yehuda Koren
*KDD* 2009

# Adding Latent Factors



Mean Rating by Movie Age (Netflix)

# Adding Latent Factors

$$\min_{U,V} \sum_{(i,j)\in M} (M_{i,j} - \widehat{M}_{i,j})^2$$

$V$

$U$

Time factors:
- Column (movie)- time baseline $b_{j,\mathrm{Bin}(t)}$
- Row (user)-time baseline *function* $b_i(t)$

$$\widehat{M}_{i,j} = \mu + b_i + b_j + \vec{u}_i \cdot \vec{v}_j \\ + b_{j,\mathrm{Bin}(t)} + b_i(t)$$

$\approx$

| 5 | | | | ? | |
|---|---|---|---|---|---|
| | ? | | 1 | | 2 |
| ? | | | | | ? |
| | | ? | | | |
| 5 | | | ? | | ? |
| | | | | 4 | |

Collaborative Filtering with Temporal Dynamics
Yehuda Koren
*KDD* 2009

# Bayesian Modeling

# Bayesian Modeling

$\mu_U$

$\sim$

Sample user factors from
Normal distribution

Bayesian Probabilistic Matrix Factorization
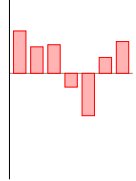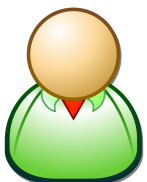Ruslan Salakhutdinov and Andriy Mnih
*ICML* 2008
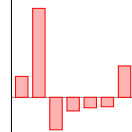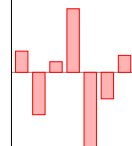
# Bayesian Modeling



$\mu_U$

$\sim$

Sample user factors from
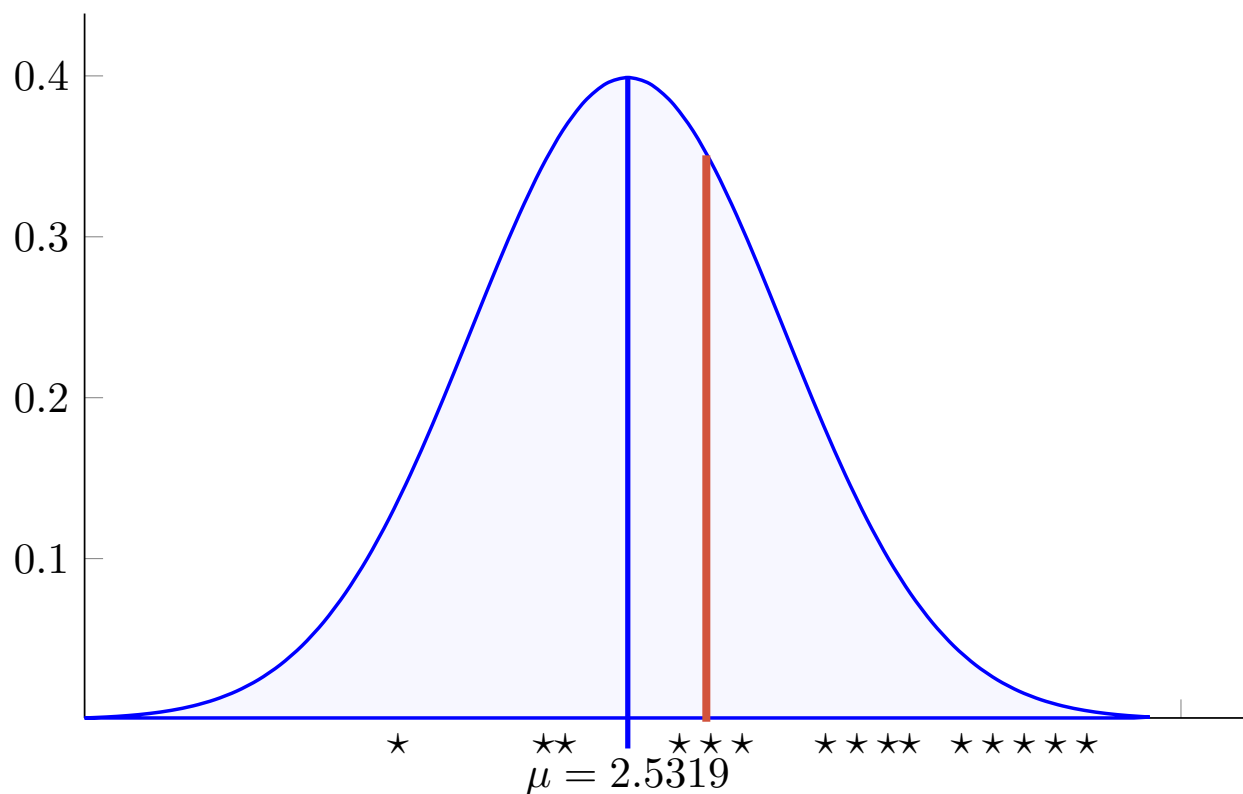Normal distribution

Update mean based on
user factors

Bayesian Probabilistic Matrix Factorization
Ruslan Salakhutdinov and Andriy Mnih
*ICML* 2008

# Bayesian Modeling



Similarly sample movie factors

# Bayesian Modeling



$$p\big(M_{i,j}\big|U,V\big) = \mathcal{N}\big(M_{i,j}\big|\vec{u}_i \cdot \vec{v}_j, \sigma^2\big)$$

Bayesian Probabilistic Matrix Factorization
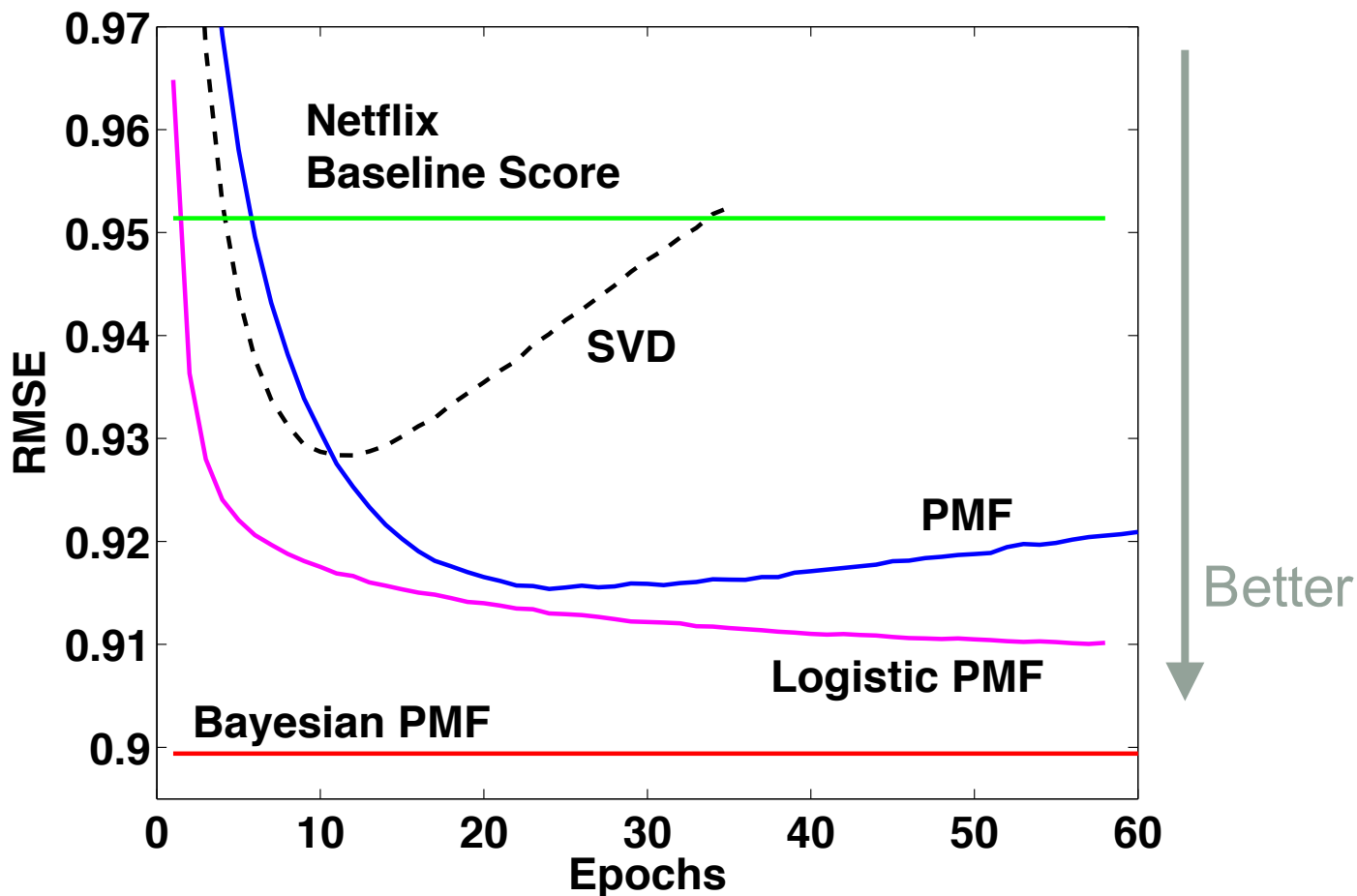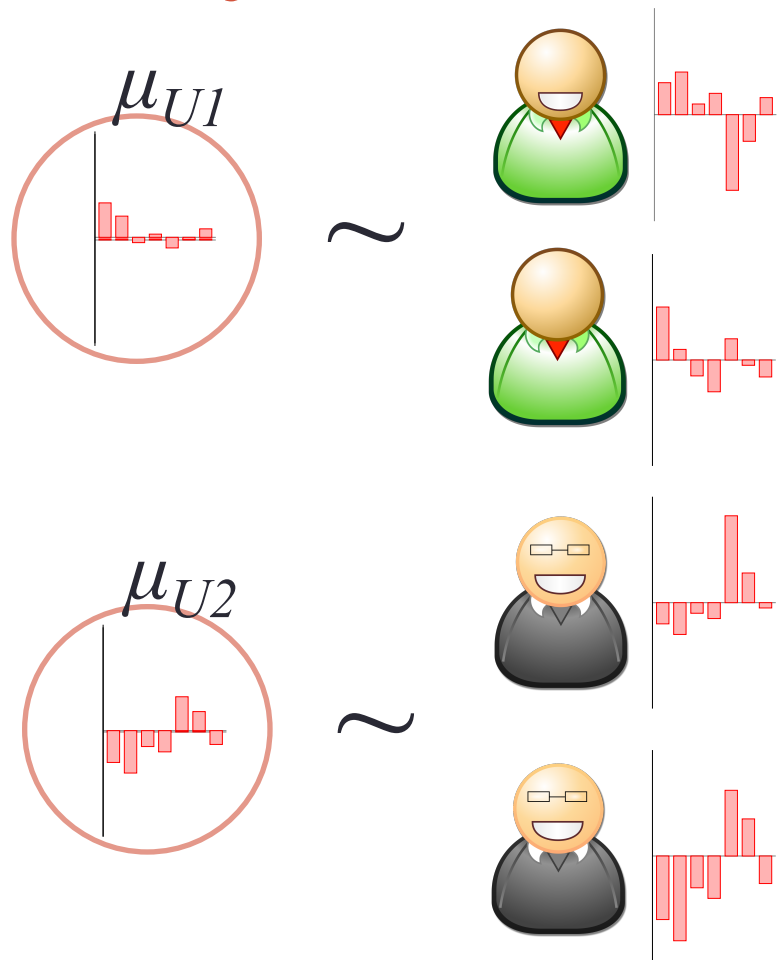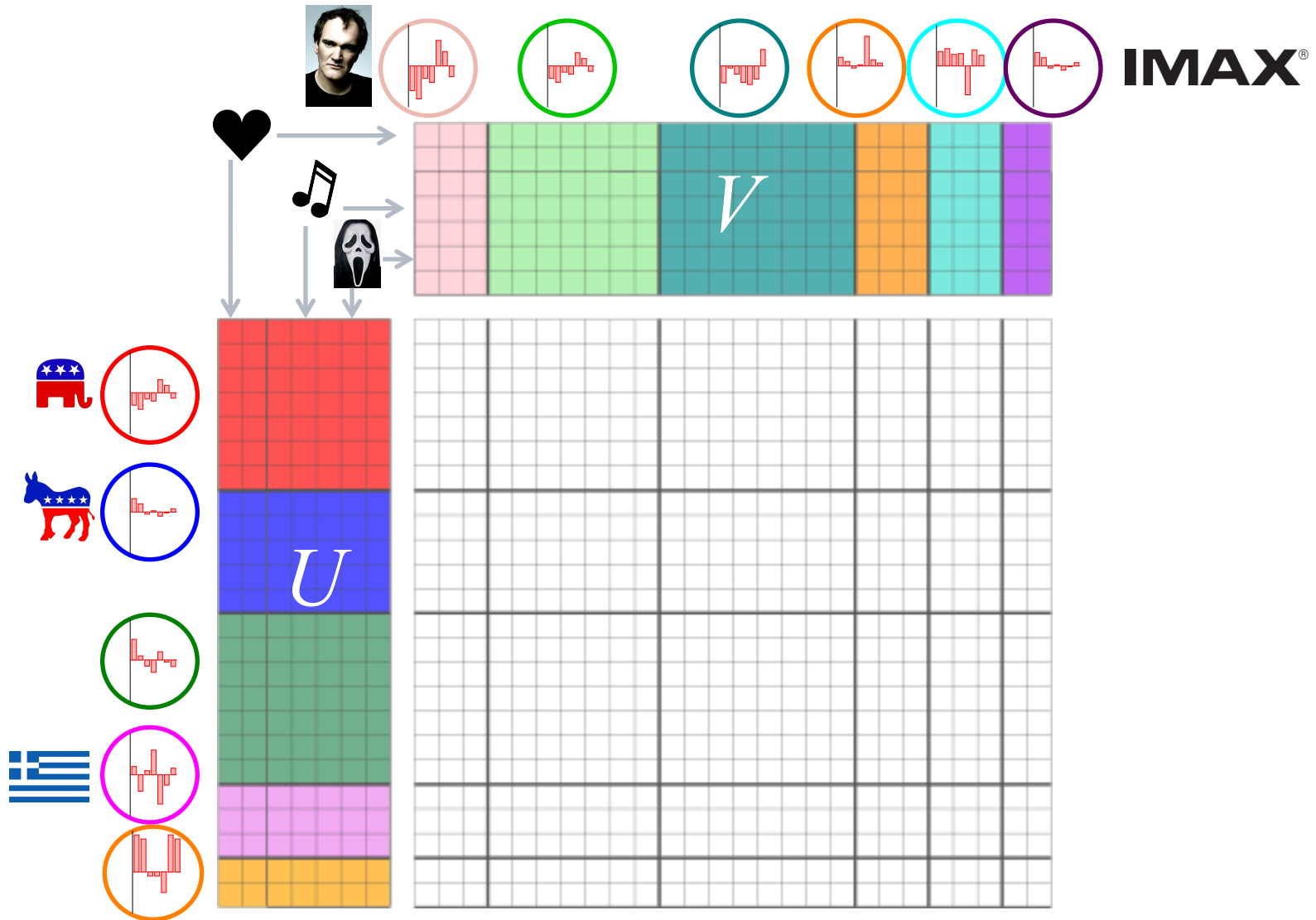Ruslan Salakhutdinov and Andriy Mnih
*ICML* 2008

# Bayesian Modeling



Bayesian Probabilistic Matrix Factorization
Ruslan Salakhutdinov and Andriy Mnih
*ICML* 2008

# Bayesian Modeling with Co-Clustering

$\mu_{U1}$
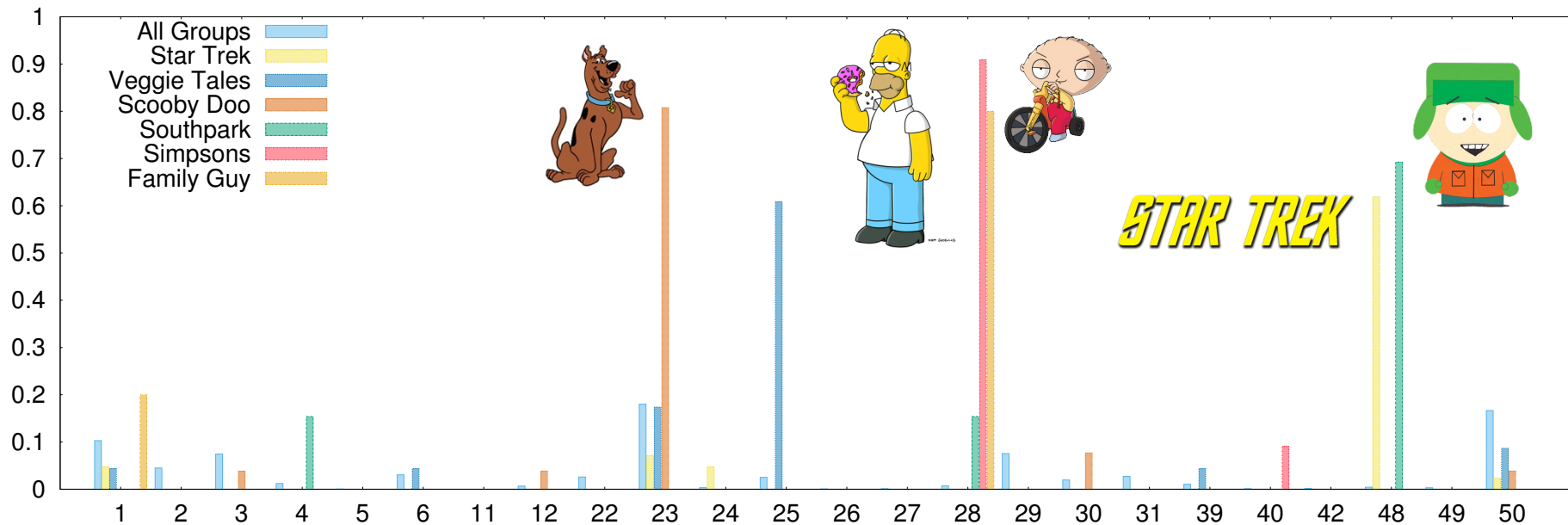
$\sim$

$\mu_{U2}$

$\sim$

Cluster users
with similar factors

CoBaFi: Collaborative Bayesian Filtering
Alex Beutel, Kenton Murray,
Christos Faloutsos Alex Smola
*WWW* 2014

# Bayesian Modeling with Co-Clustering



| | Cluster 28 | Cluster 30 | Cluster 48 |
|---|---|---|---|
| | Simpsons | Scooby Doo | Star Trek |
| | Family Guy | Spy Kids | Back to the Future |
| | Monty Python | Stuart Little | Southpark |
| | Curb your Enthusiasm | Dr. Dolittle | Lord of the Rings |
| | The Twilight Zone | Lion King | Harry Potter |
| | Arrested Development | Agent Cody Banks | The X-Files |

CoBaFi: Collaborative Bayesian Filtering
Alex Beutel, Kenton Murray,
Christos Faloutsos Alex Smola
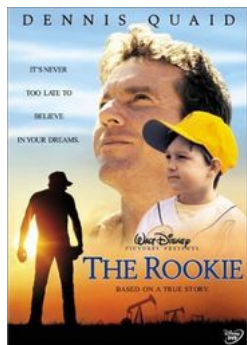*WWW* 2014

# Online Rating Models



Typically fit a Gaussian - Minimize RMSE



| Data | Normal CF |

CoBaFi: Collaborative Bayesian Filtering
Alex Beutel, Kenton Murray,
Christos Faloutsos Alex Smola
*WWW* 2014

# Online Rating Models



Typically fit a Gaussian - Minimize RMSE
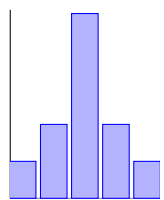


Data    Normal CF    CoBaFi

CoBaFi: Collaborative Bayesian Filtering
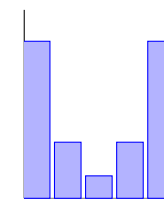Alex Beutel, Kenton Murray,
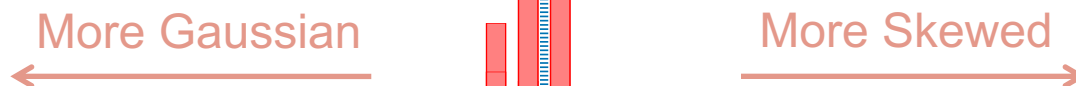Christos Faloutsos Alex Smola
*WWW* 2014

# Shape of Netflix reviews

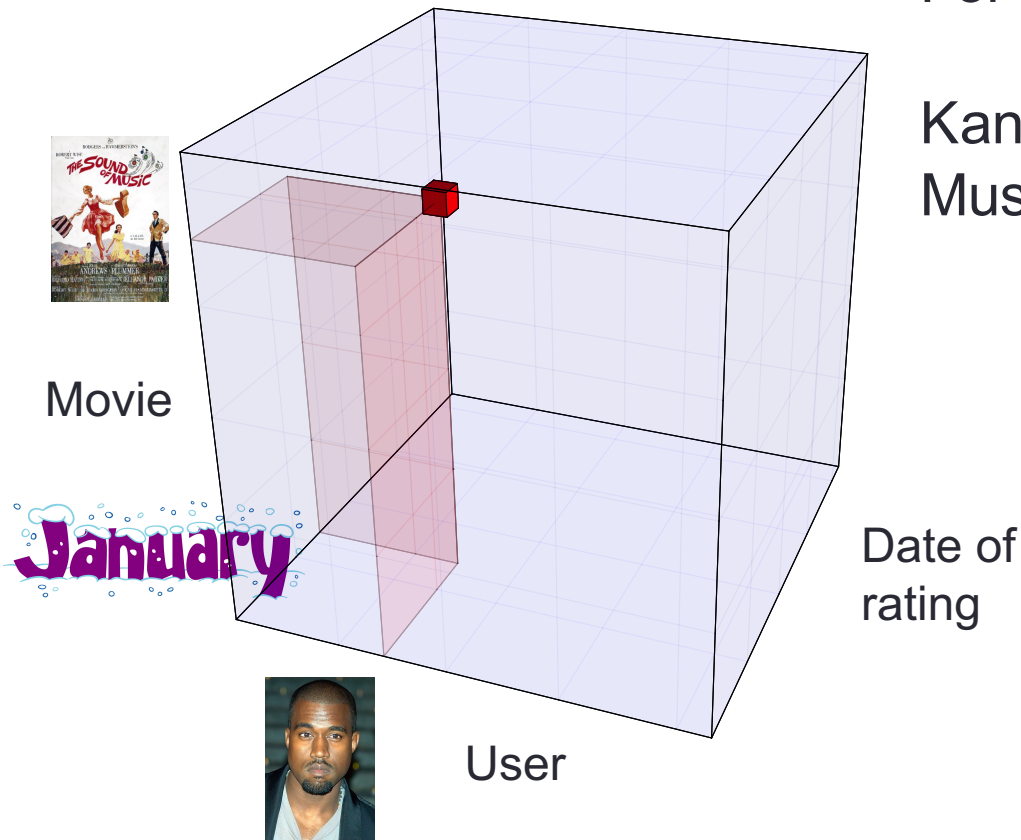| Most Gaussian | Most skewed |
|---|---|
| The Rookie | The O.C. Season 2 |
| The Fan | Samurai X: Trust and Betrayal |
| Cadet Kelly | Aqua Teen Hunger Force: Vol. 2 |
| Money Train | Sealab 2001: Season 1 |
| Alice Doesn't Live Here | Aqua Teen Hunger Force: Vol. 2 |
| Sea of Love | Gilmore Girls: Season 3 |
| Boiling Point | Felicity: Season 4 |

# Stars

Movies

More Gaussian ← → More Skewed

# Stars

TV Shows

CoBaFi: Collaborative Bayesian Filtering
Alex Beutel, Kenton Murray,
Christos Faloutsos Alex Smola
*WWW* 2014

# What is a tensor?

- Tensors are used for structured data > 2 dimensions
- Think of as a 3D-matrix
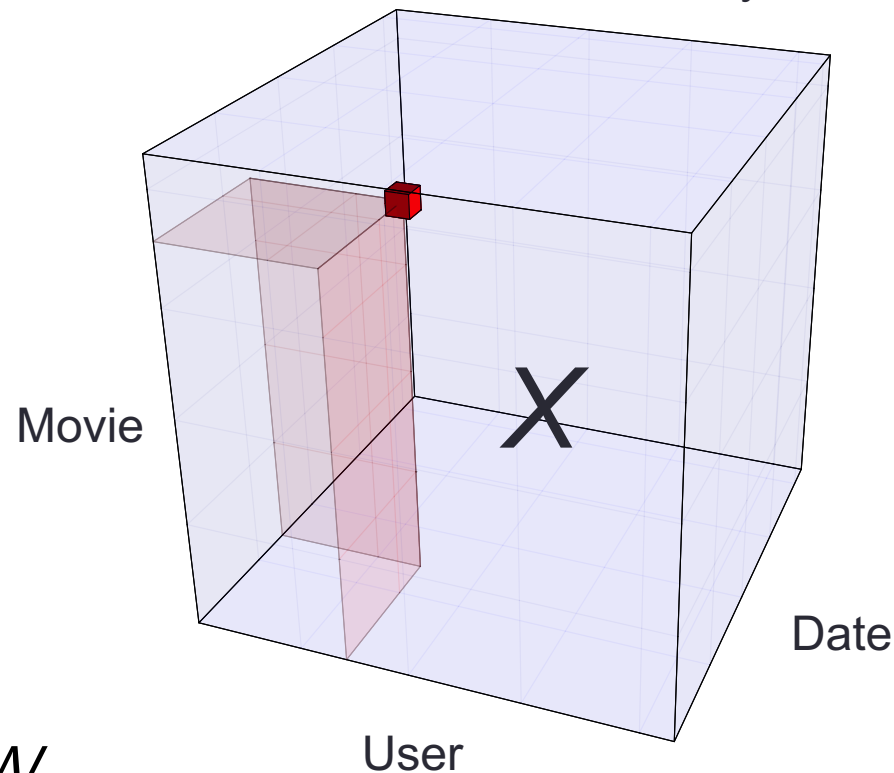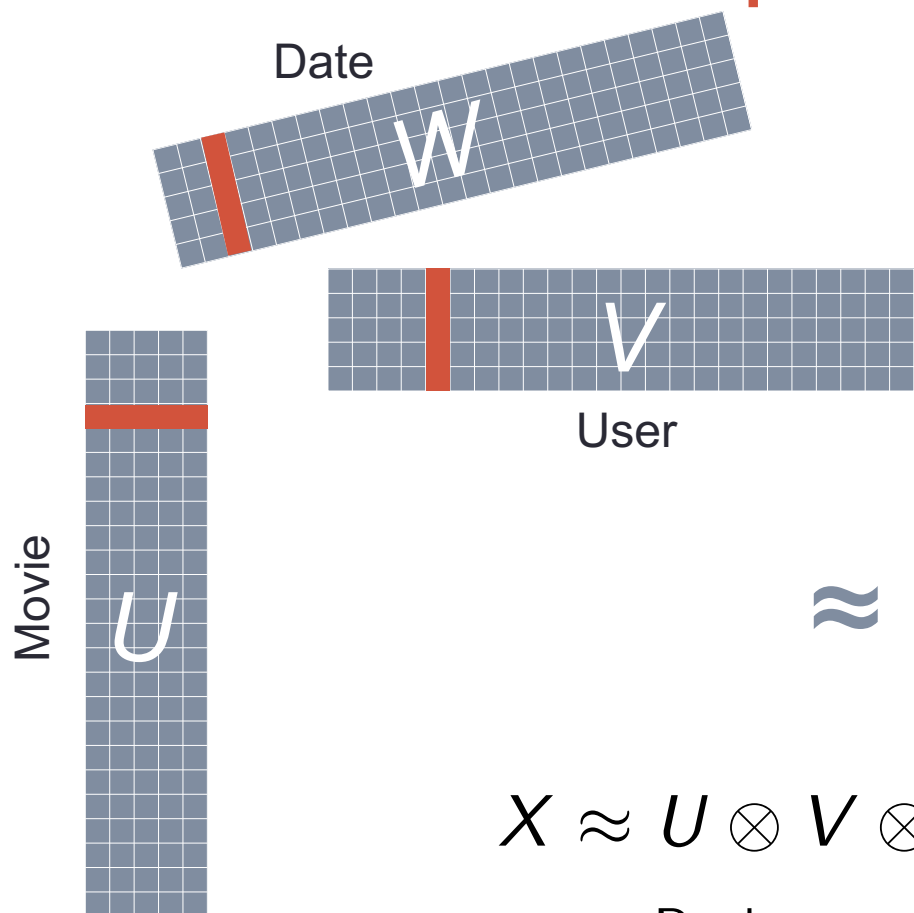


Movie

January

User

Date of rating

For example:

Kanye West rated The Sound of Music five stars last January.

# Tensor Decomposition

Kanye West rated The Sound of
Music five stars last January.

Date

W

V

User

Movie

U

≈

$$X \approx U \otimes V \otimes W$$

$$X_{i,j,k} \approx \sum_{r=1}^{\text{Rank}} U_{i,r} V_{j,r} W_{k,r}$$
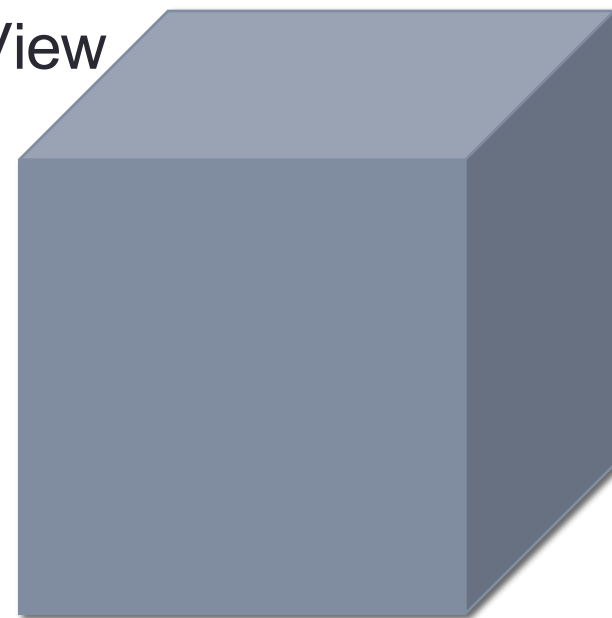
Movie

X

Date

User

# Graph Clustering with Tensors

Multiple possible views
of the DBLP network:
1. Who-cites-whom
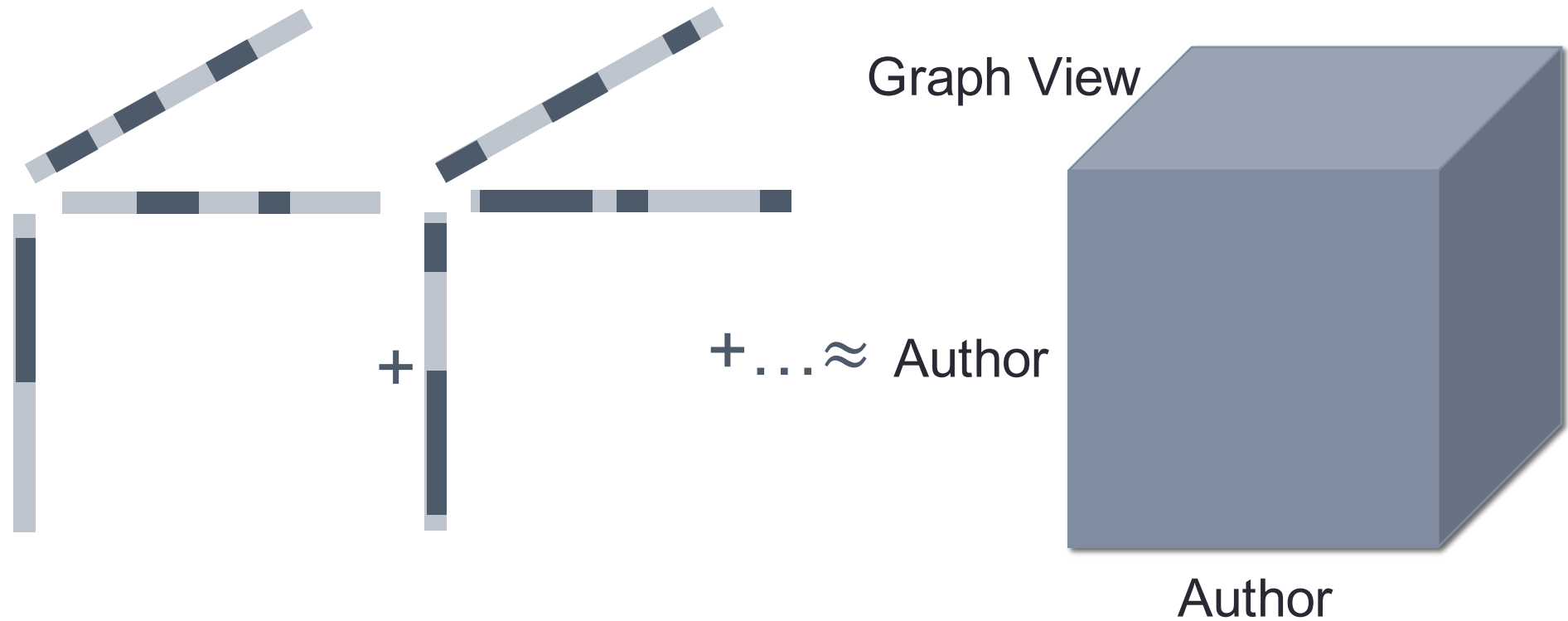2. Co-authorship
3. Using same words in title

Graph View

Author

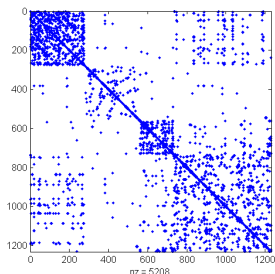Author

# Graph Clustering with Tensors



Graph View

$+ \dots \approx$ Author

Author

## Sparse Tensor Factorization

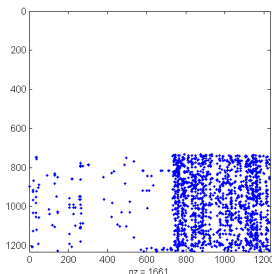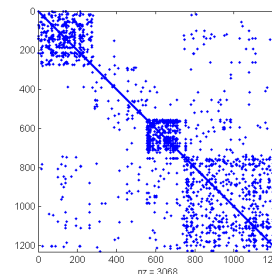# Graph Clustering with Tensors

DBLP-1
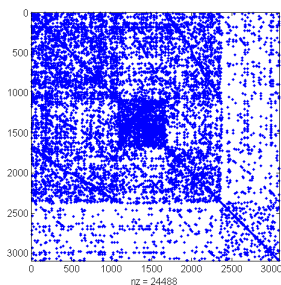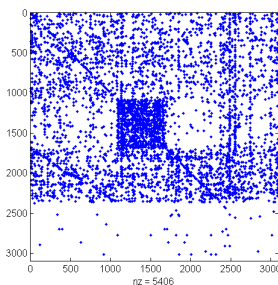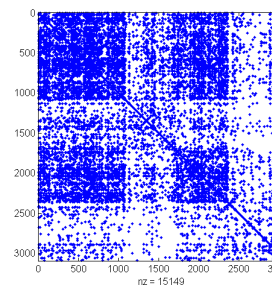


(a) citation    (b) co-auth.    (c) co-term
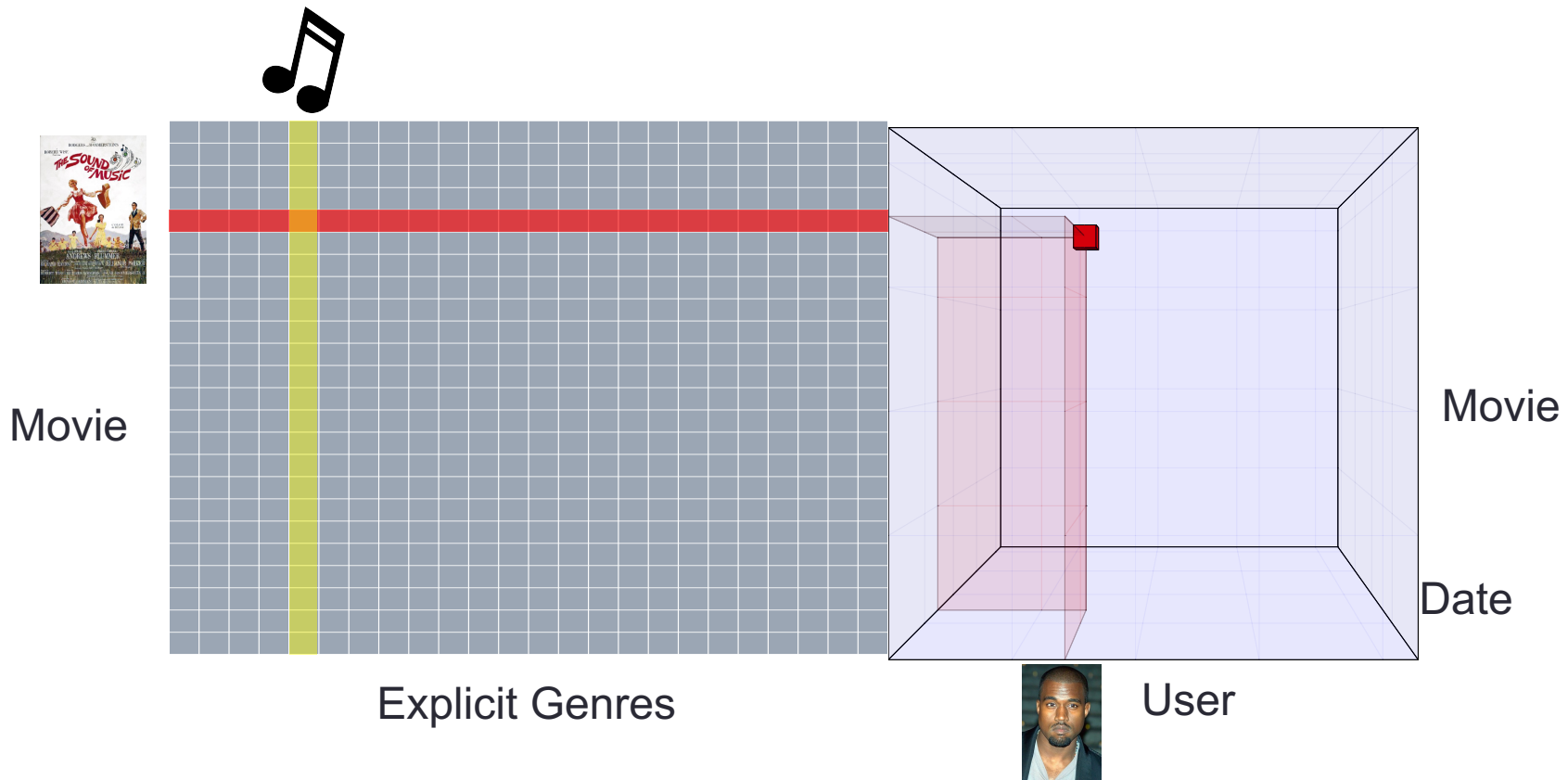
DBLP-2



(a) citation    (b) co-auth.    (c) co-term

# Graph Clustering with Tensors

| Dataset | Baseline | GraphFuse |
|---------|----------|-----------|
| DBLP-1  | 0.12     | **0.30**  |
| DBLP-2  | 0.08     | **0.12**  |

## Modeling Accuracy

# Coupled Matrix + Tensor Decomposition



Movie

Explicit Genres

Movie

Date

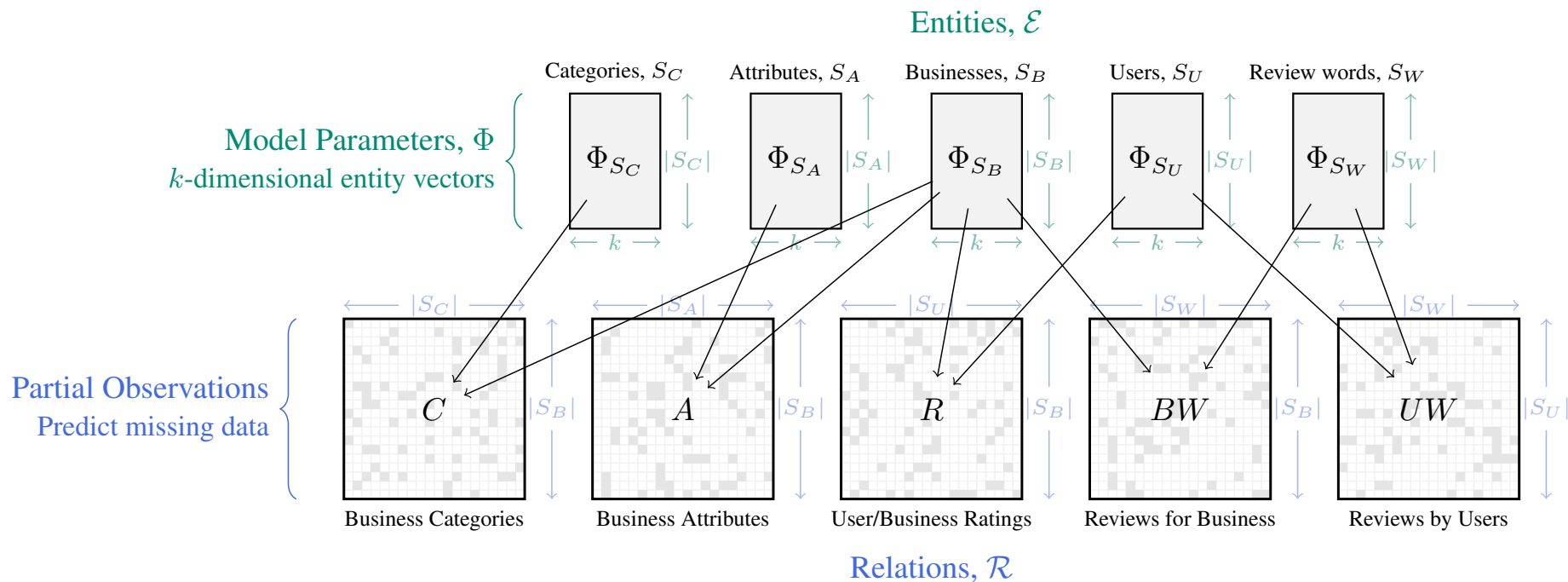User

# Coupled Matrix + Tensor Decomposition



$$X \approx U \otimes V \otimes W$$

$$Y \approx UA^{\top}$$

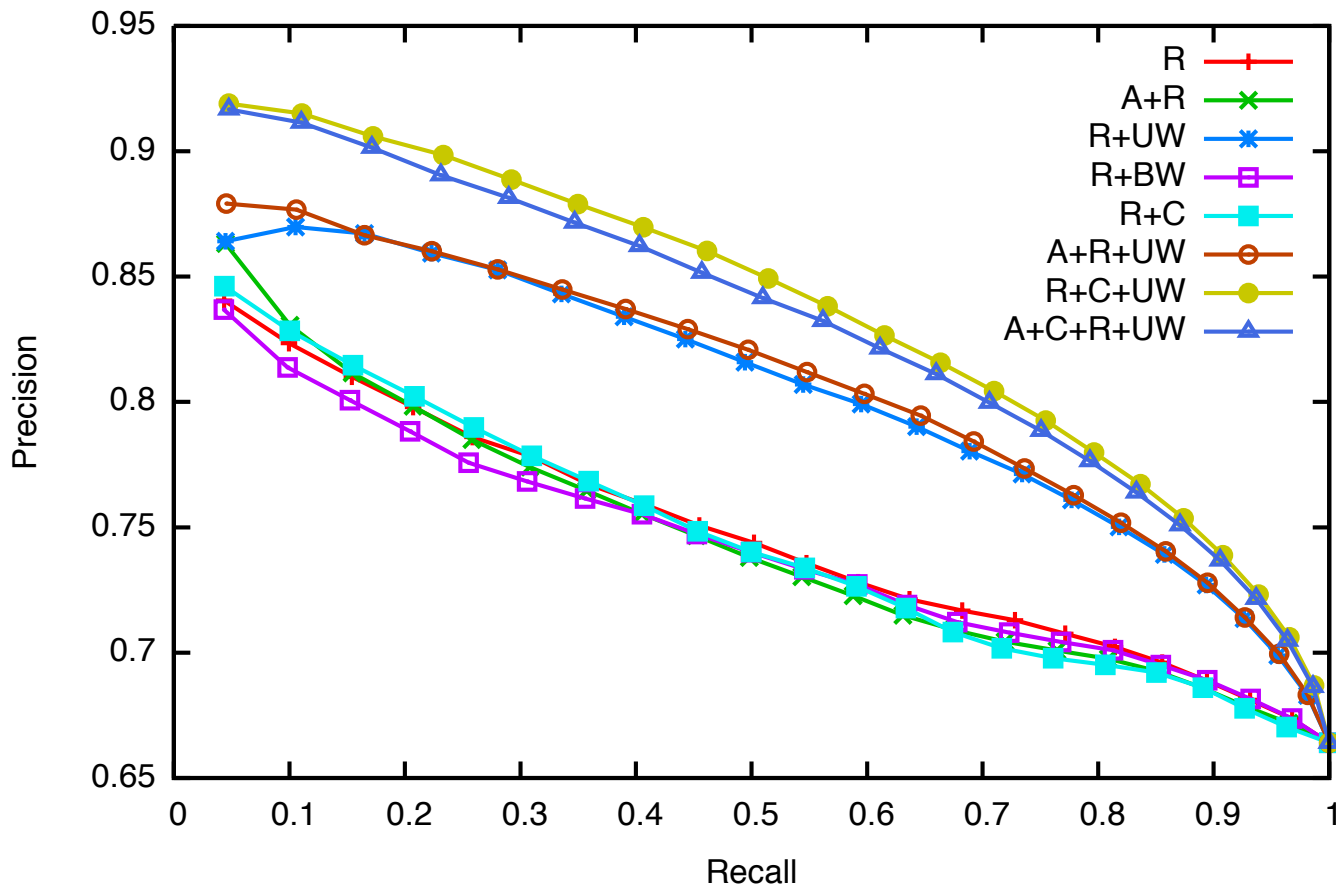$$\min_{U,V,W,A} \|X - U \otimes V \otimes W\|_F^2 + \|Y - UV^{\top}\|_F^2$$

# Joint Factorization



Entities, $\mathcal{E}$

Model Parameters, $\Phi$
$k$-dimensional entity vectors

Categories, $S_C$  Attributes, $S_A$  Businesses, $S_B$  Users, $S_U$  Review words, $S_W$

$\Phi_{S_C}$ $|S_C|$  $\Phi_{S_A}$ $|S_A|$  $\Phi_{S_B}$ $|S_B|$  $\Phi_{S_U}$ $|S_U|$  $\Phi_{S_W}$ $|S_W|$

$\leftarrow k \rightarrow$

Partial Observations
Predict missing data

$|S_C|$  $|S_A|$  $|S_U|$  $|S_W|$  $|S_W|$

$C$ $|S_B|$  $A$ $|S_B|$  $R$ $|S_B|$  $BW$ $|S_B|$  $UW$ $|S_U|$

Business Categories  Business Attributes  User/Business Ratings  Reviews for Business  Reviews by Users

Relations, $\mathcal{R}$

# Joint Factorization



PR Curve (Ratings)

Most valuable:
1. Ratings
2. Review text
3. Business
   Categories

Collective Factorization for Relational Data:
An Evaluation on the Yelp Datasets
Nitish Gupta, Sameer Singh

# 1. Subgraph Analysis
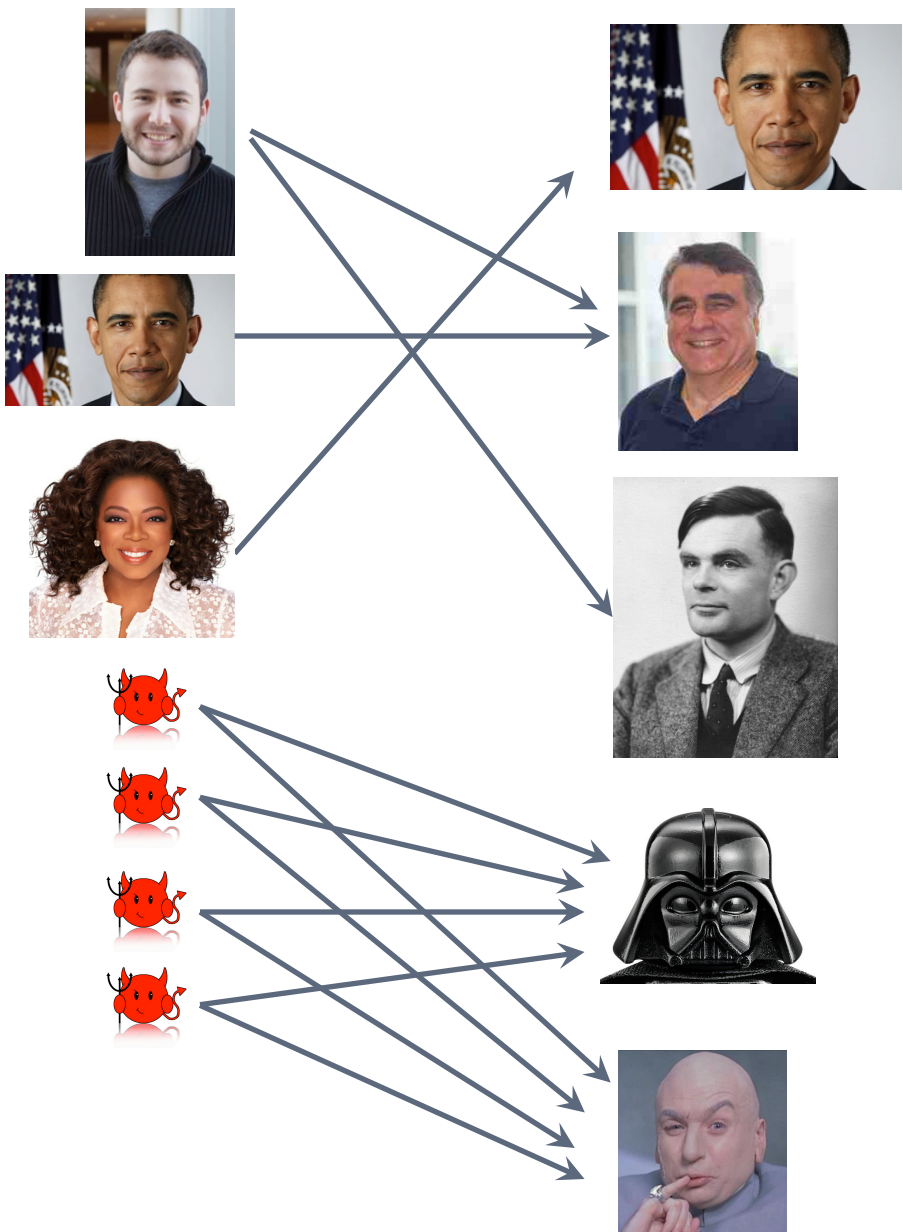
# 2. Propagation Methods

# 3. Latent Factor Models
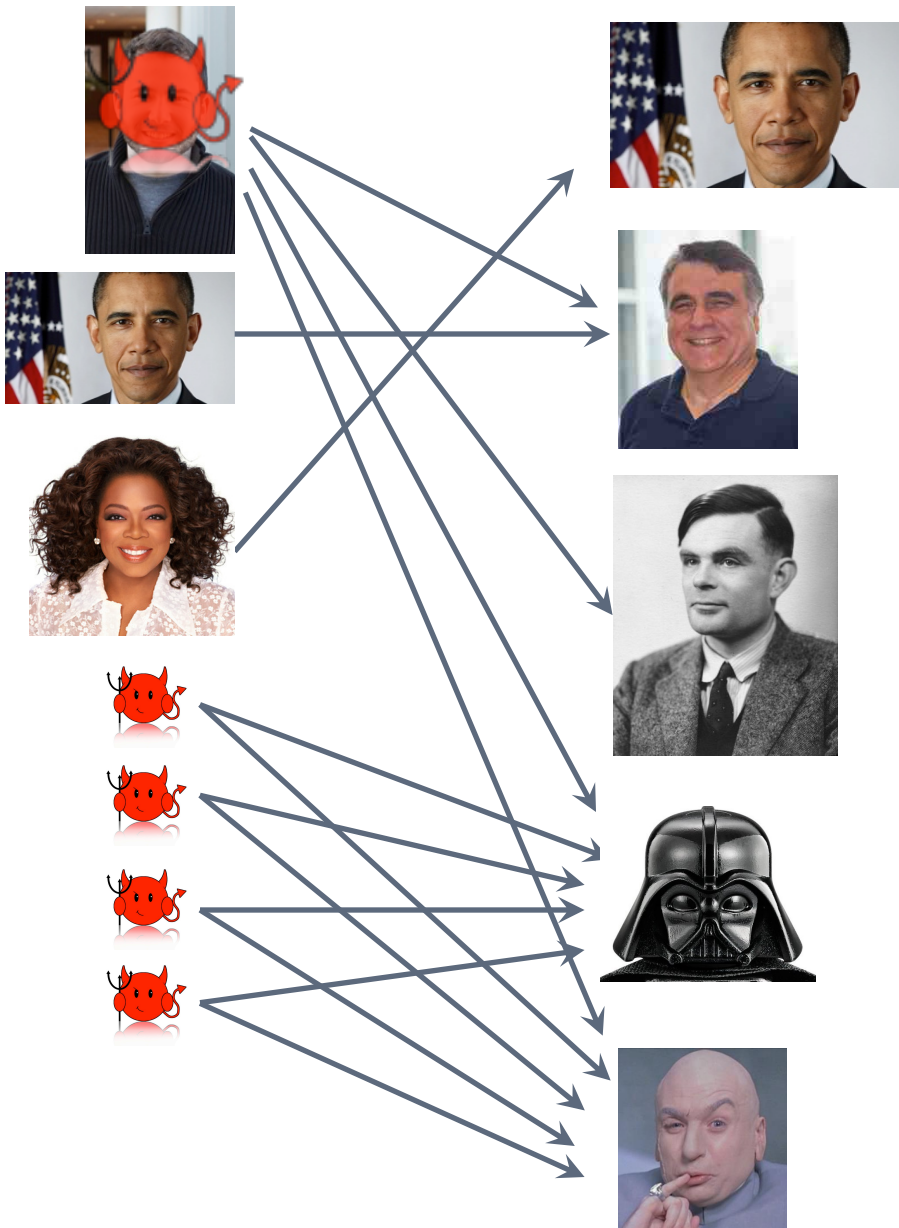
## a) Background

## b) Normal Behavior
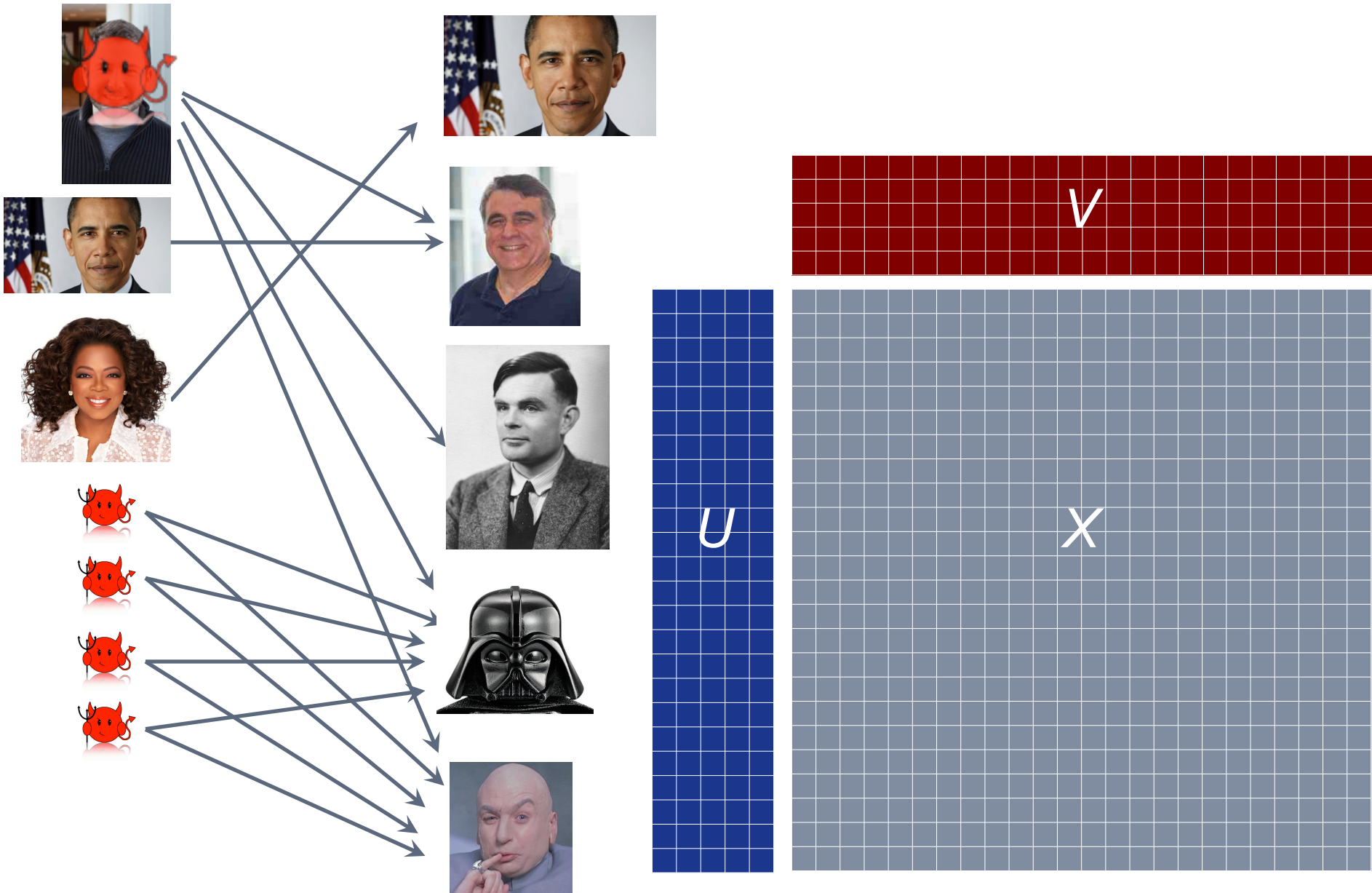
## c) Abnormal Behavior
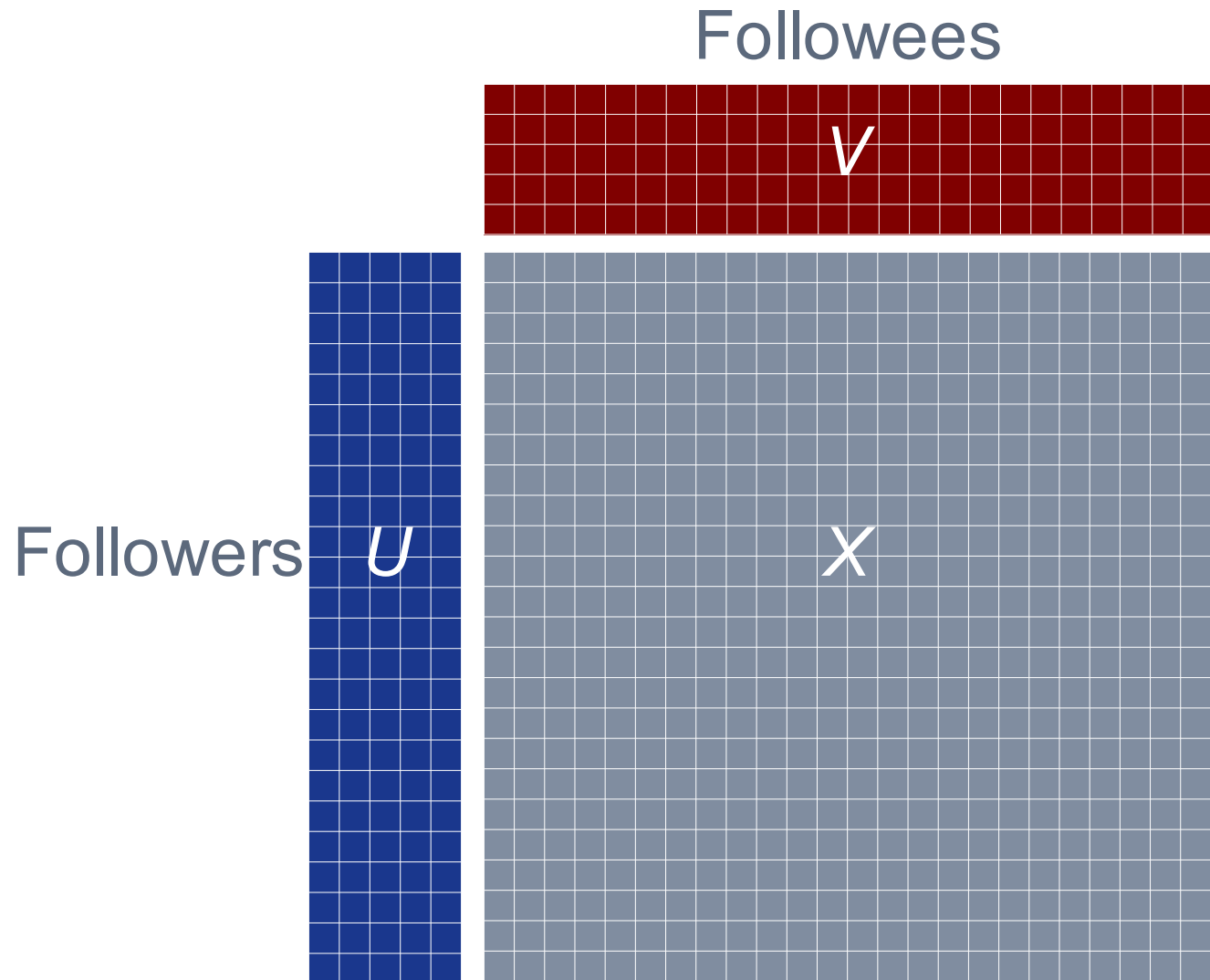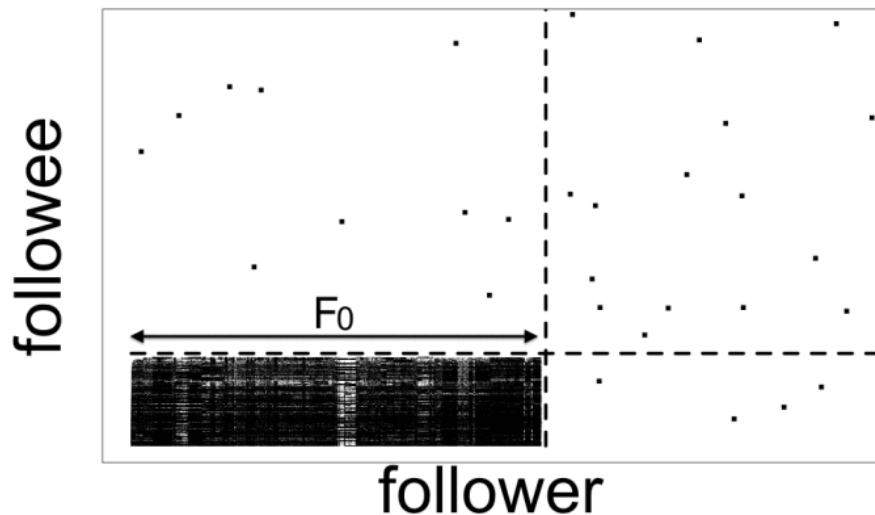
# Fraud Detection

# Fraud Detection

# Fraud within a factorization

# Fraud within a factorization

Followees

Followers

*V*

*U*

*X*

# Fraud within a factorization

Followees

$V$

| 1.5 | 1 | -0.5 | -2 | 1 |

Followers $U$

$X$

| ? | ? | ? | ? | ? |
| ? | ? | ? | ? | ? |
| ? | ? | ? | ? | ? |

# Fraud within a factorization

# Fraud within a factorization

$u_1$          $u_3$

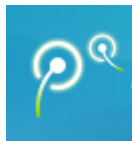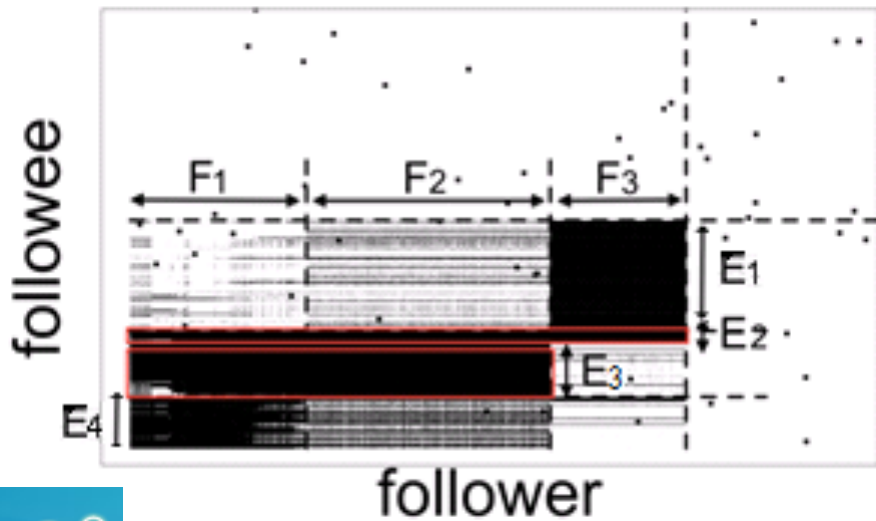| $\varepsilon$ | $\varepsilon$ | 1 | $\varepsilon$ | $\varepsilon$ |

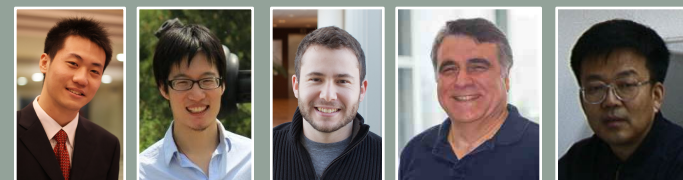# Fraud within a factorization

# Fraud within a factorization



Inferring Strange Behavior from Connectivity Pattern in Social Networks
Meng Jiang, Peng Cui, Alex Beutel,
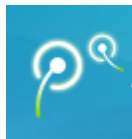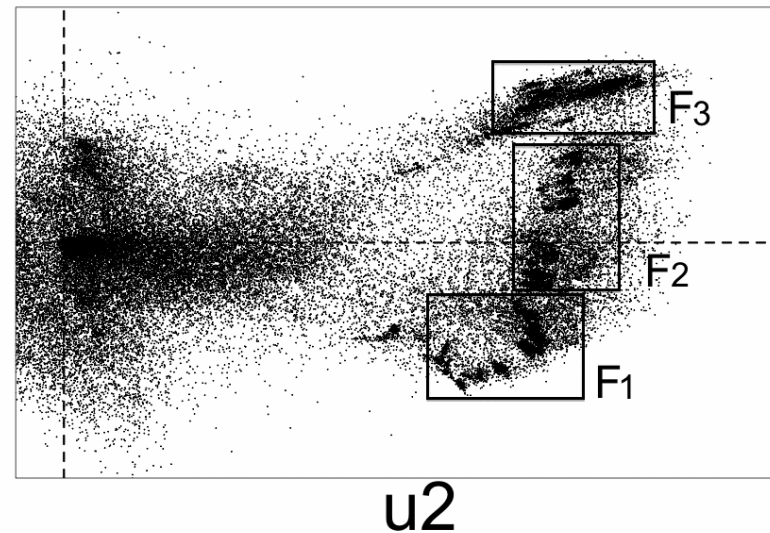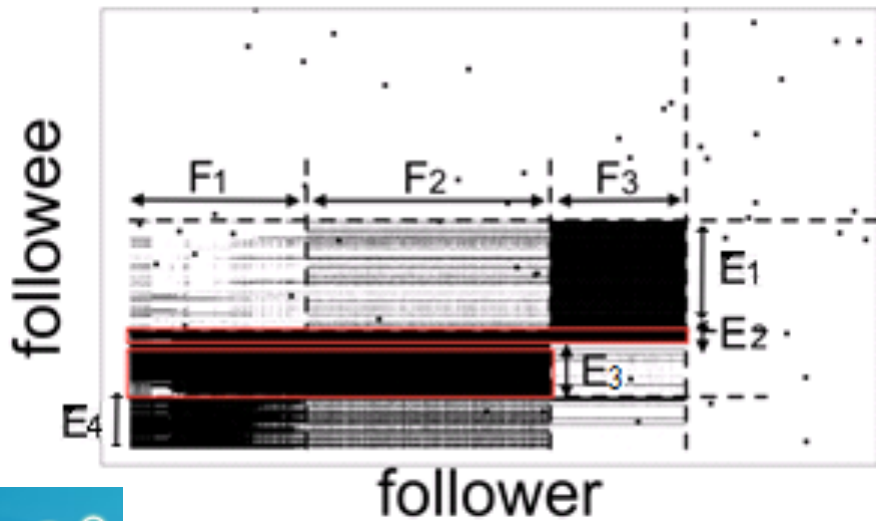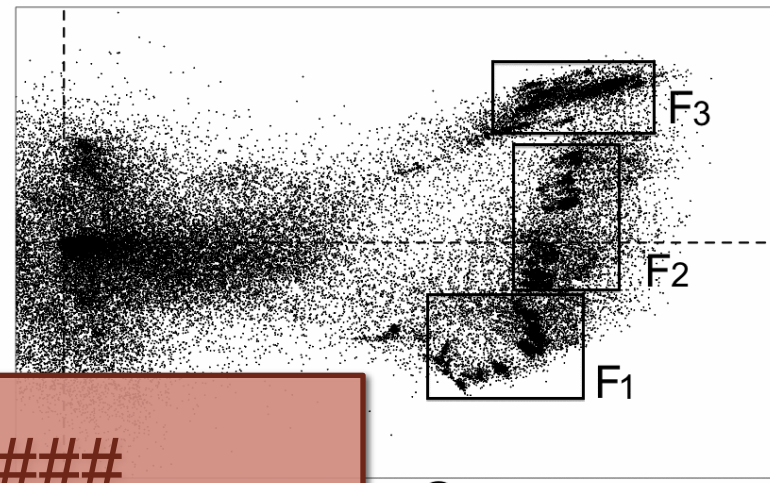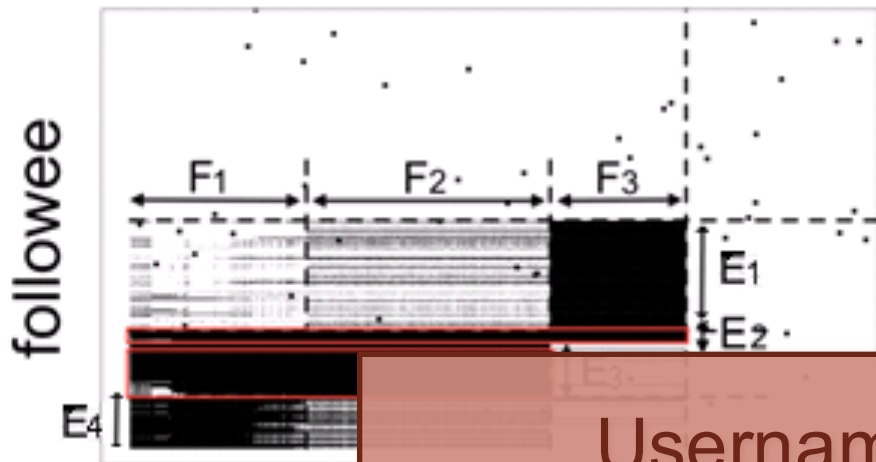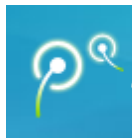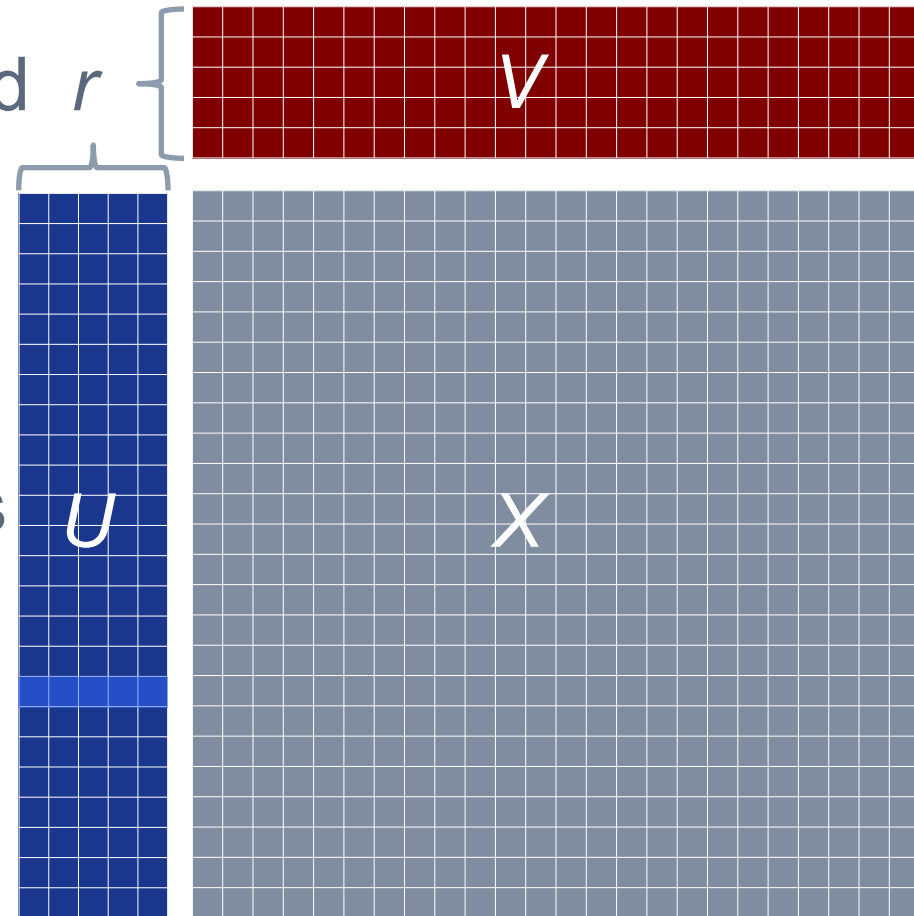Christos Faloutsos, Shiqiang Yang.
*PAKDD*, 2014

# Fraud within a factorization



Username: a#####
Birthday: January 1st

# Complementary Fraud Detection

Followees

Limited $r$

$V$

Followers $U$

$X$

? ? ? ? ?

# Complementary Fraud Detection

960 fraudsters
safely following
960 customers



Singular Value (Attack Size) vs Number of components (k)

Honest users

Fraudsters

Honest objects          Customers

Spotting Suspicious Link Behavior with fBox: An Adversarial Perspective
Neil Shah, Alex Beutel, Brian Gallagher,
Christos Faloutsos
*ICDM*, 2014.

# Complementary Fraud Detection

Followees

Limited $r$

$V$

Followers

$U$

$X$

0 | 0 | 0 | 0 | 0

11111111
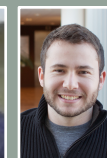
# Complementary Fraud Detection



Spotting Suspicious Link Behavior with fBox: An Adversarial Perspective
Neil Shah, Alex Beutel, Brian Gallagher,
Christos Faloutsos
*ICDM*, 2014.

# Complementary Fraud Detection



Spotting Suspicious Link Behavior with fBox: An Adversarial Perspective
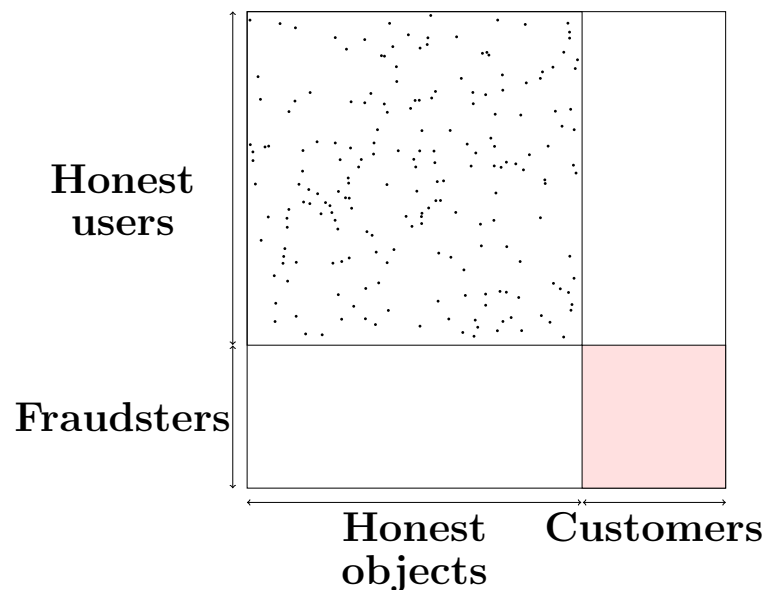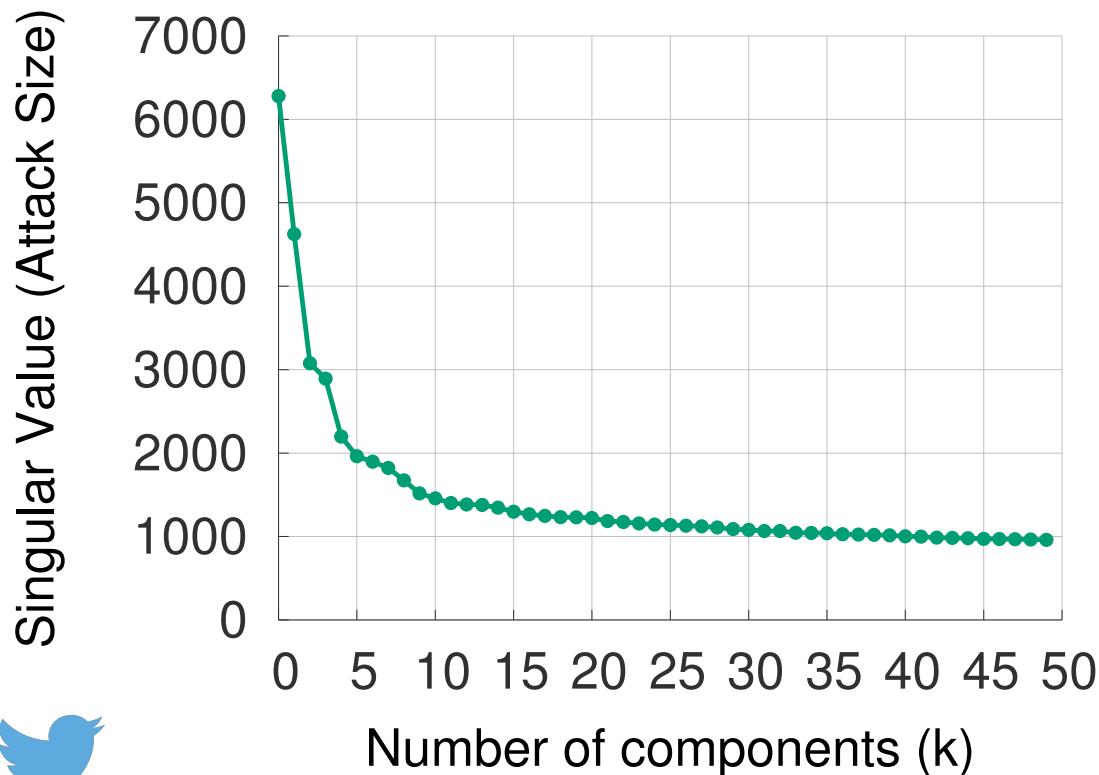Neil Shah, Alex Beutel, Brian Gallagher,
Christos Faloutsos
*ICDM,* 2014.
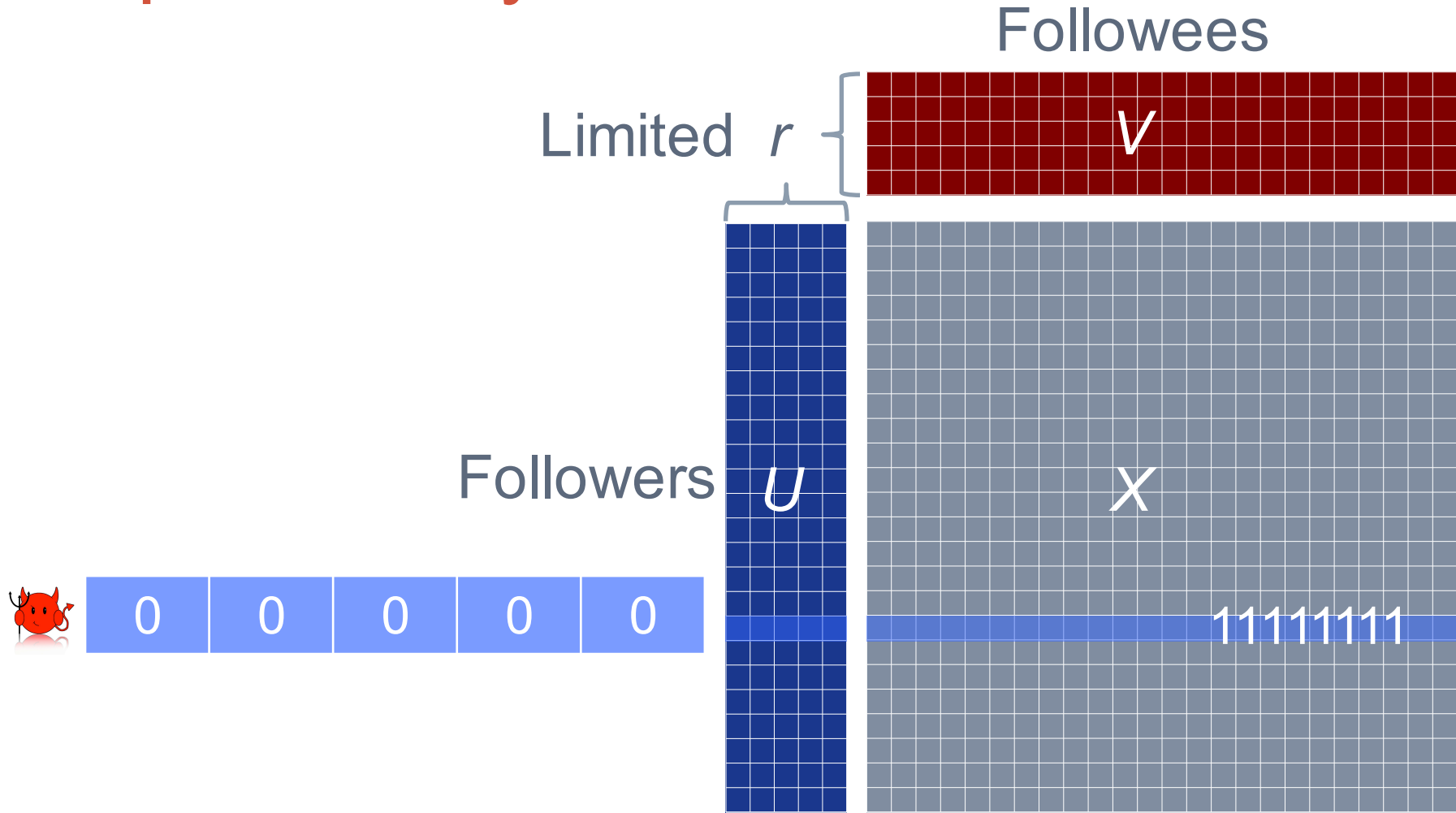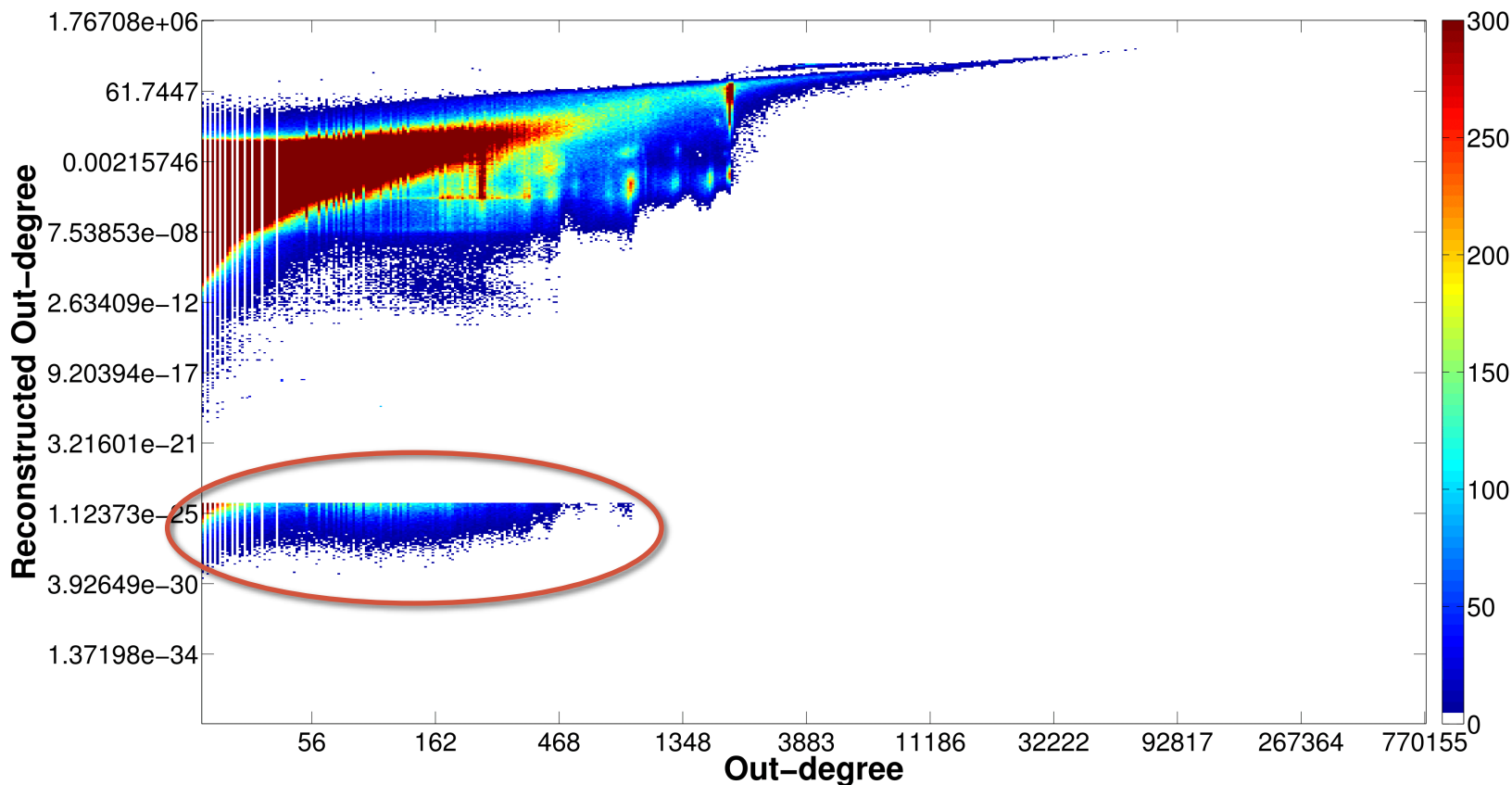
# Complementary Fraud Detection



93% Precision

70% of accounts missed by Twitter

Spotting Suspicious Link Behavior with fBox: An Adversarial Perspective
Neil Shah, Alex Beutel, Brian Gallagher,
Christos Faloutsos
*ICDM*, 2014.

# Practitioner's Guide

| Method | Graph Type | Node Attributes | Edge Attributes | Seed Labels |
|---|---|---|---|---|
| EigenSpokes | Directed+ | | | |
| Get-the-Scoop | Directed+ | | | |
| fBox | Directed+ | | | |
| CoBaFi | Bipartite+ | | ✔ | |
| CDOutliers | Undirected | ✔ | | |

# Detecting Fraud within Recommendation

# Detecting Fraud within Recommendation



CoBaFi: Collaborative Bayesian Filtering
Alex Beutel, Kenton Murray,
Christos Faloutsos Alex Smola
*WWW* 2014

# Clustering Fraudsters

$\mu_1$     $\mu_2$     $\mu_3$     $\mu_4$     $\mu_5$

**Naïve Spammers**

**Spam + Noise**

**Hijacked Accounts**

# Clustered Fraudsters

Clustered naïve spammers

Clustered hijacked accounts

Clustered "attacked" movies



83% are clustered together

# Outliers in Joint Factorization



Enforce $U_1 \approx U_2$ and $U_1, U_2, V_1, V_2 \geq 0$

Community Distribution Outlier Detection in Heterogeneous Information Networks
Manish Gupta, Jing Gao, and Jiawei Han
*ECML/PKDD* 2013

# Outliers in Joint Factorization

Interesting design of $X_1$ and $X_2$; see paper for details



Enforce $U_1 \approx U_2$ and $U_1, U_2, V_1, V_2 \geq 0$

# Outliers in Joint Factorization

Rows of $V_2$ represent common patterns in $X_2$ (cluster centroids)



Enforce $U_1 \approx U_2$ and $U_1, U_2, V_1, V_2 \geq 0$

Community Distribution Outlier Detection in Heterogeneous Information Networks
Manish Gupta, Jing Gao, and Jiawei Han
*ECML/PKDD* 2013

# Outliers in Joint Factorization

Rows of $V_2$ represent common patterns in $X_2$ (cluster centroids)

An anomaly is a row of $X_i$ that is *not* similar to any row in $V_i$



$V_2$

$U_2$

$X_2$

# Practitioner's Guide

| Method | Graph Type | Node Attributes | Edge Attributes | Seed Labels |
|---|---|---|---|---|
| EigenSpokes | Directed+ | | | |
| Get-the-Scoop | Directed+ | | | |
| fBox | Directed+ | | | |
| CoBaFi | Bipartite+ | | ✔ | |
| CDOutliers | Undirected | ✔ | | |

# Recap

- SVD captures communities of interest

- Bayesian methods can:

  - Handle missing values

  - Give factorization models (-> patterns, & anomalies)

- Group-outliers: spotted by CoBaFi, Get-the-Scoop, etc.

# CONCLUSION

$$\left\{ \begin{array}{l} \text{Undirected} \\ \text{Directed} \\ \text{Bipartite} \end{array} \right\}$$

✖

$$\left\{ \text{Node Attributes} \right\}$$

✖

$$\left\{ \text{Edge Attributes} \right\}$$

✖

$$\left\{ \begin{array}{l} \text{Unsupervised} \\ \text{Semi-Supervised} \end{array} \right\}$$

# Open Problems / Opportunities

**P1. Complex data:** How should we integrate data from multiple data sources?

$$\left\{ \begin{array}{c} \text{Undirected} \\ \text{Directed} \\ \text{Bipartite} \end{array} \right\}$$

$\times$

$$\left\{ \text{Node Attributes} \right\}$$

$\times$

$$\left\{ \text{Edge Attributes} \right\}$$

$\times$

$$\left\{ \begin{array}{c} \text{Unsupervised} \\ \text{Semi-Supervised} \end{array} \right\}$$

# Open Problems / Opportunities

**P2. Adversarial analysis:** Can we offer provable guarantees on detecting fraud and spam?

# Open Problems / Opportunities

**P3. Early detection:** Can we detect fraudsters before they cause significant damage?

# Summary

Local Subgraph Analysis: Patterns and Features
e.g. using ego-nets



POSTNET

http://www.sizemore.co.uk/
2005/08/i-feel-some-movies
-coming-on.html

http://instapundit.com/
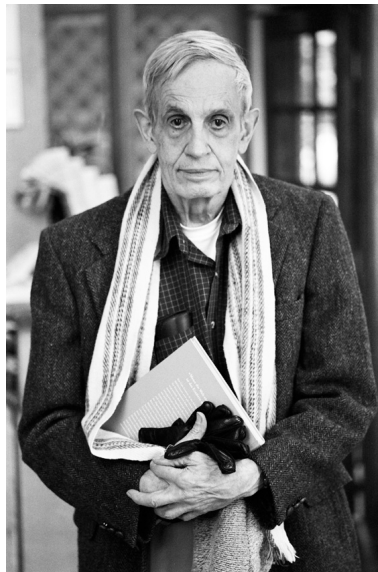archives/025235.php

| | |
|---|---|
| —— | $1.1094x + (-0.21414) = y$ |
| - - - | $1.1054x + (-0.21432) = y$ |
| – – – | $2.1054x + (-0.51535) = y$ |

|E|

|N|

# Summary

## Propagation Methods
"Guilt-by-association"
"Importance-by-association" = PageRank

# Summary

## Latent Factor Models
Find multiple communities, patterns and anomalies.

# Take Away
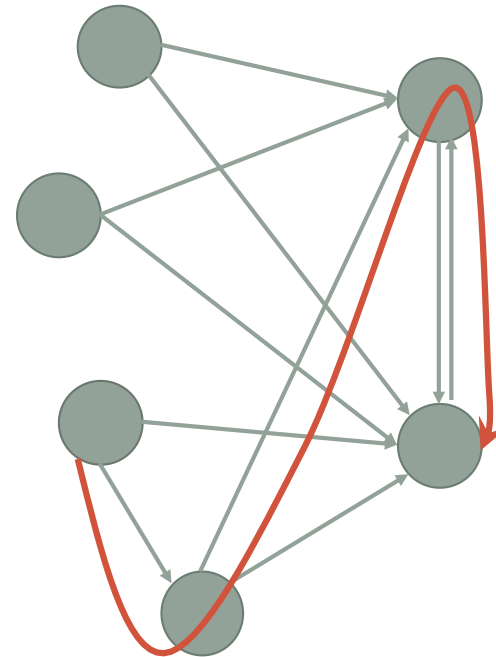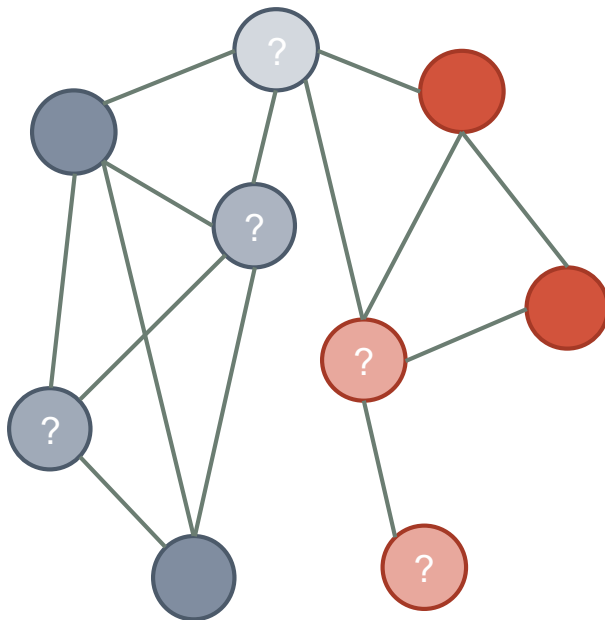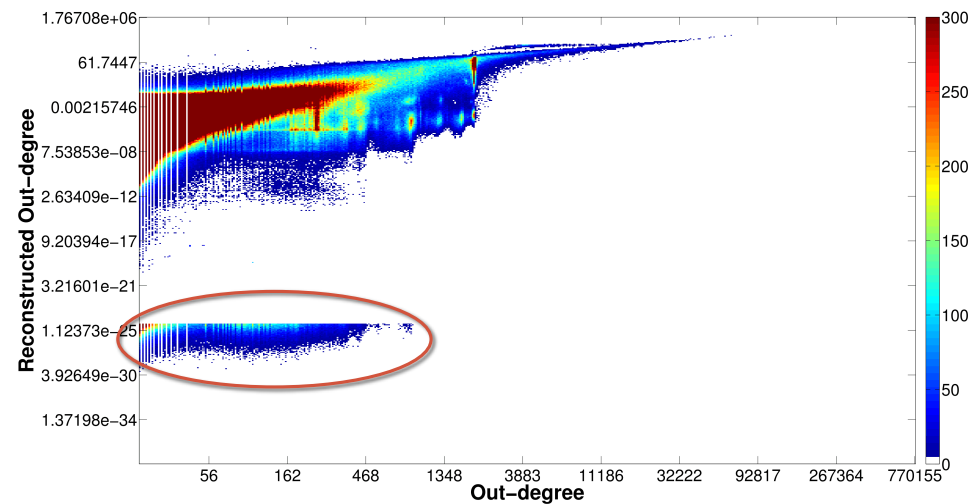
## User Modeling and Fraud Detection are two sides of the same coin.

# ODDx3 workshop TODAY 9:30-5:45

Afternoon Schedule:

- **Keynote** by Vipin Kumar

- **Panel** 'What is an Anomaly?' by Tiberio Caetano, Vipin Kumar, Tina Eliassi-Rad, Ted Senator, Jimeng Sun

- **Research talks**

  **http://outlier-analytics.org/odd15kdd/**

**ACM SIGKDD 2015 Workshop**

**ODDx3: Outlier Definition, Detection, and Description**

# Thanks again to

# Questions?

Carnegie
Mellon
University

Stony Brook
University

References and resources available at
**cs.cmu.edu/~abeutel/kdd2015_userbehavior**



KDD 2015　　A. Beutel, L. Akoglu, C. Faloutsos　46
**Pattern: Ego-net Power Law Density**

POSTNET

http://www.sizemore.co.uk/
2005/08/i-feel-some-movies
-coming-on.html

http://instapundit.com/
archives/025235.php

1.1094x + (−0.21414) = y
1.1054x + (−0.21432) = y
2.1054x + (−0.51535) = y

Oddball: Spotting anomalies in weighted graphs
Leman Akoglu, Mary McGlohon, Christos Faloutsos
*PAKDD* 2010



KDD 2015　　A. Beutel, L. Akoglu, C. Faloutsos　3
**Semi-supervised Classification**

Given a graph and
labels for some nodes,
can we learn the labels
for the other nodes?



KDD 2015　　A. Beutel, L. Akoglu, C. Faloutsos　15
**Matrix Factorization**　What does each
eigenvector capture?

$UV^{T} \approx M$

Each factor captures a
dense block in the matrix

Genres