BIG Data is not flat

# Data is multi-modal, multi-relational, spatio-temporal, multi-media



shorthand:  **Graph Data**

# NEED: ML* for Graphs

*: Machine Learning

# ML for Graphs

Pattern #1: Collective Classification

Pattern #2: Link Prediction

Pattern #3: Entity Resolution

# ML for Graphs

Pattern #1: Collective Classification – inferring labels of nodes in graph

Pattern #2: Link Prediction

Pattern #3: Entity Resolution

# ML for Graphs

Pattern #1: Collective Classification – inferring labels of nodes in graph

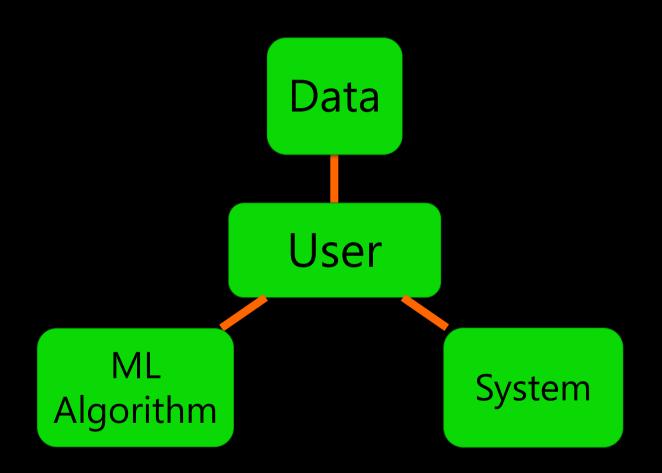Pattern #2: Link Prediction – inferring the existence of edges in graph

Pattern #3: Entity Resolution

# ML for Graphs

Pattern #1: Collective Classification – inferring labels of nodes in graph

Pattern #2: Link Prediction – inferring the existence of edges in graph

Pattern #3: Entity Resolution – clustering nodes that refer to the same underlying entity

# What about Interaction?

# What's different about graphs?

Unit of Interaction

Context

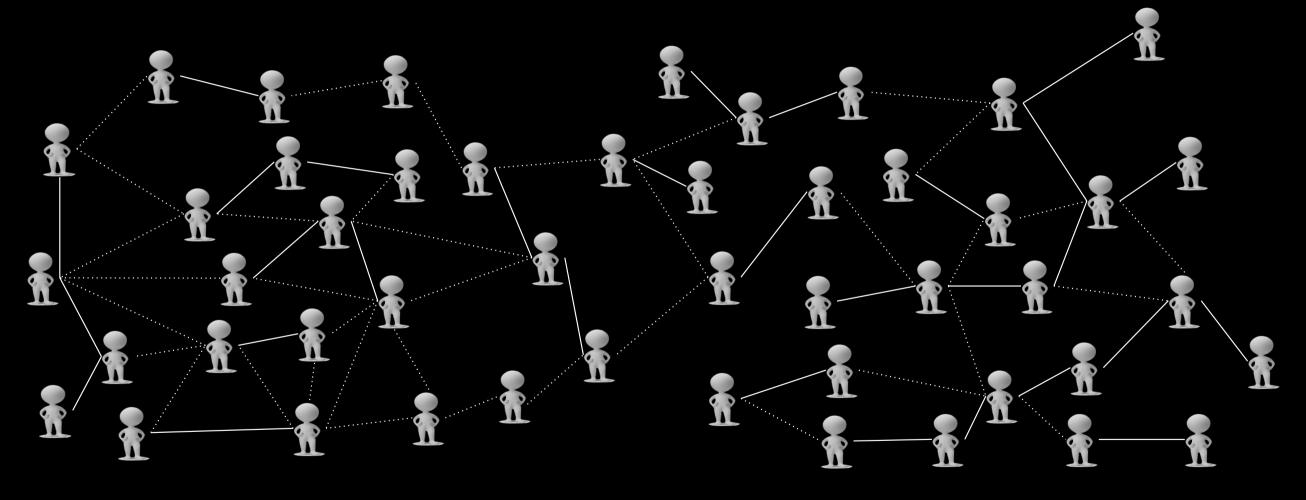Comparison

# What's different about graphs?

Unit of Interaction
Context
Comparison

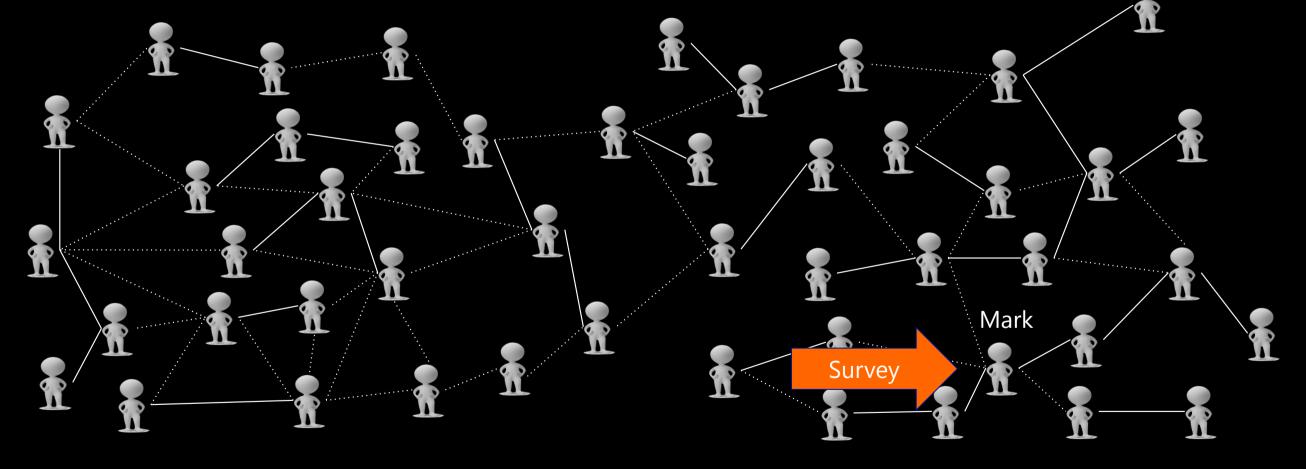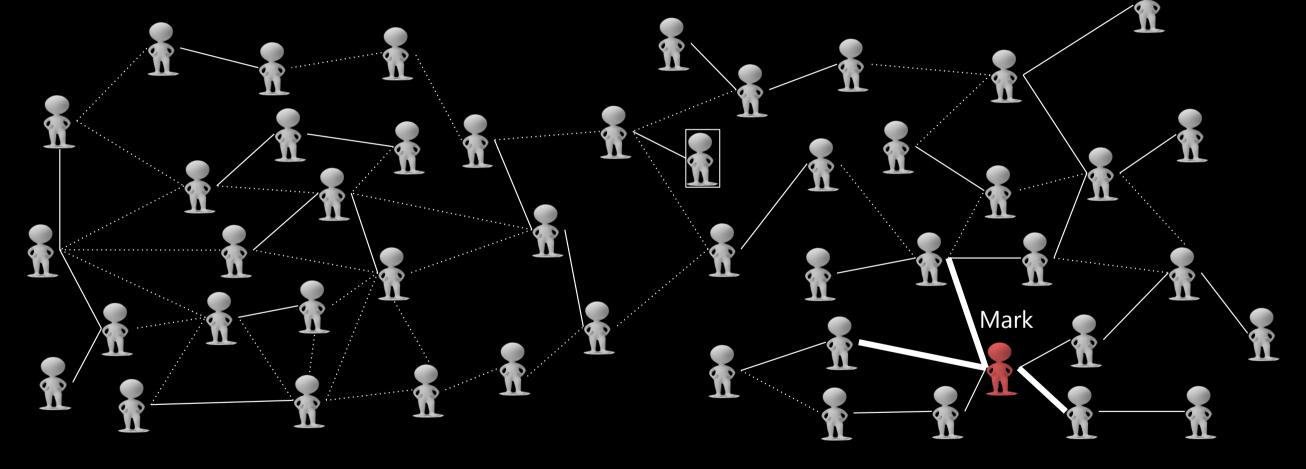# Nugget: active surveying – acquire label *and* neighbors

Sharara & Getoor IJCAI 2011;  Namata et al., MLG 2012

Most previous work assumes that only the labels are unobserved (i.e., a <u>fully</u> observed network)

Label: Positive Neutral Negative

# Network structure also often only partially observed



**Label:** 🟩 **Positive**   🟦 **Neutral**   🟥 **Negative**

Survey: Acquire the label and ego-network of a node
e.g., personal interview, targeted information gathering

Mark

Survey

Label:  ■ Positive  ■ Neutral  ■ Negative

Survey: Acquire the label and ego-network of a node
e.g., personal interview, targeted information gathering

Mark

Label:   Positive   Neutral   Negative

# Survey: Acquire the label and ego-network of a node
## e.g., personal interview, targeted information gathering



Catherine

Survey

**Label:** Positive | Neutral | Negative

Survey: Acquire the label and ego-network of a node
e.g., personal interview, targeted information gathering

Catherine

Survey

Label:  Positive  Neutral  Negative

% Reduction in Required Responders
Active Survey vs Random

# What's different about graphs?

Unit of Interaction

Context

Comparison

# Context

too little: single node

too much: whole graph

just right: relational context
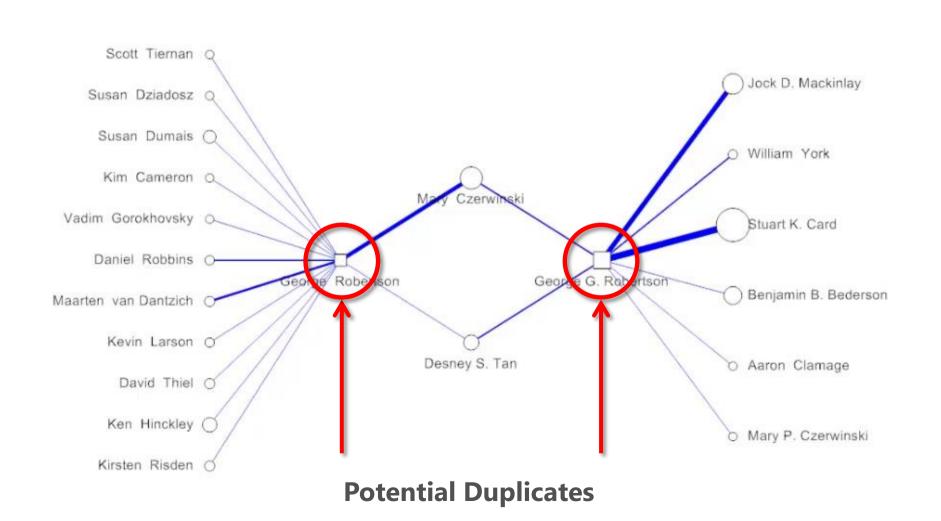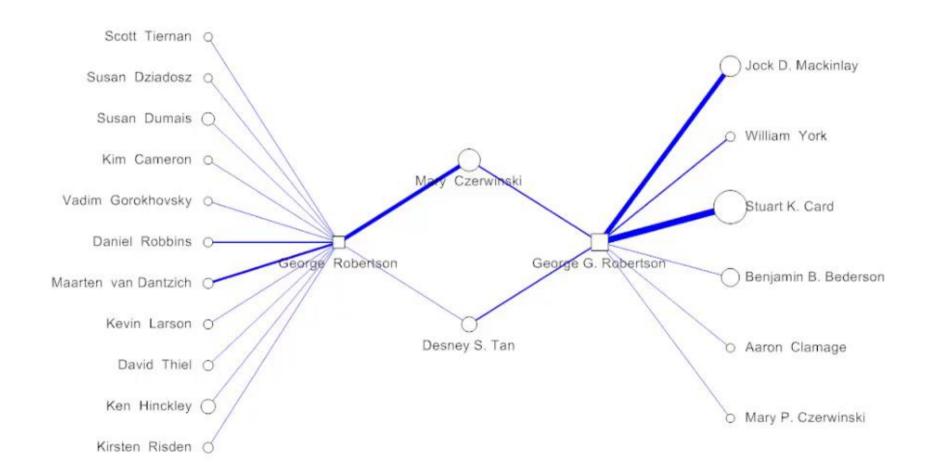
# D-Dupe: Interactive Entity Resolution Tool



Potential Duplicate Viewer

Relational Context Viewer

Data Detail Viewer

Kang, Getoor, Shneiderman, Bilgic, Licamele, TVCG 2008
http://www.cs.umd.edu/projects/linqs/ddupe

# Nugget: Relational Context



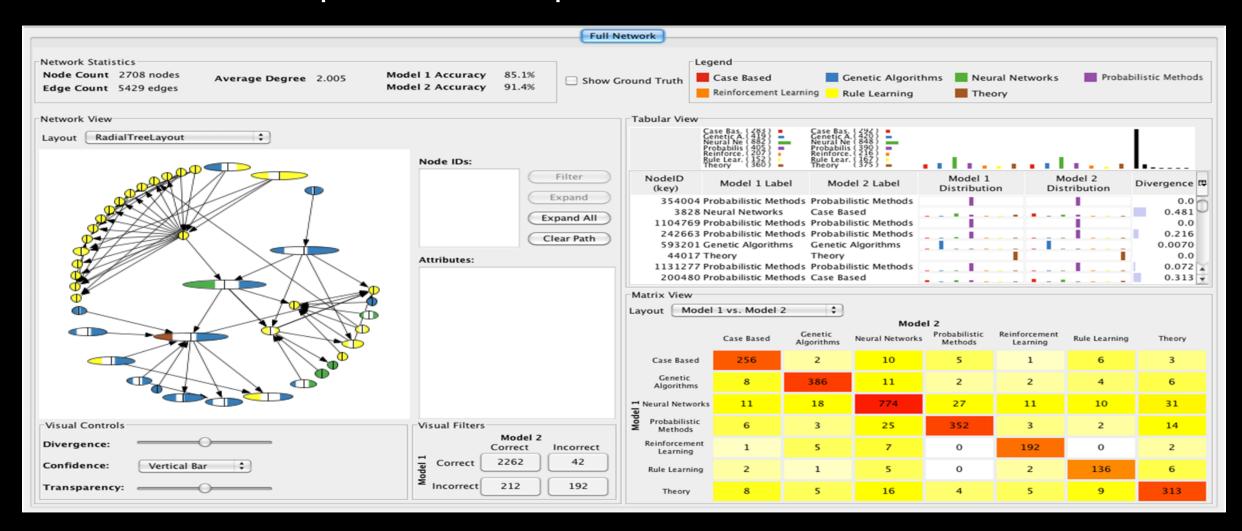Potential Duplicates

# Nugget: Relational Context

# What's different about graphs?

Unit of Interaction
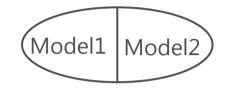
Context

Comparison
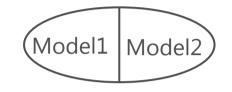
# Comparing ML Algorithms

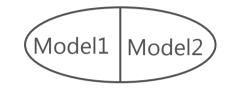Flat Data: confusion matrix
Graph Data: ?

# G-Pare: Graph Comparison



Sharara, Sopan, Namata, Getoor, VAST 2011
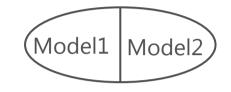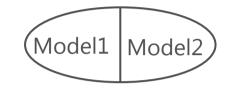http://www.cs.umd.edu/projects/linqs/gpare

# Nugget: Node Visualization

# Nugget: Node Visualization



| Color Coding | Predicted Label | Neutral / Positive | Agree | Disagree |
|---|---|---|---|---|

# Nugget: Node Visualization



| | | |
|---|---|---|
| **Color Coding** | Predicted Label | Neutral / Positive — Agree — Disagree |
| **Fill Area** | Prediction Confidence | High Confidence — Moderate Confidence — Low Confidence |

# Nugget: Node Visualization



| | | |
|---|---|---|
| **Color Coding** | Predicted Label | ▪ Neutral ▪ Positive / Agree / Disagree |
| **Fill Area** | Prediction Confidence | High Confidence / Moderate Confidence / Low Confidence |
| **Eccentricity** | KL-Divergence | |

# Nugget: Node Visualization



| Color Coding | Predicted Label |  |
| Fill Area | Prediction Confidence |  |
| Eccentricity | KL-Divergence |  |
| Border Highlighting | Ground Truth (Prediction Accuracy) |  |

# Nugget: Node Visualization



Neutral
Positive

- Model 1 prediction: "Positive"
  Model 2 prediction: "Neutral"

- Model 1 is more confident in its prediction than Model 2

- Distributions of the two models vary significantly

- Model 1's prediction matches the ground truth

# Finding regions of disagreement

# GrDB: Putting it all together, first steps...



Eldin Moustafa, Miao, Deshpande, Getoor, SIGMOD Demo 2013
http://www.cs.umd.edu/projects/linqs/grdb

# Closing

**State-of-the-Art:** interaction unit, context and comparison important

**Challenges:** interaction/ML for complex tasks involving graphs is hard

**Opportunities:** creating common abstractions that work for both interaction for ML and ML for interaction