

Data Mining Concepts & Tasks

Duen Horng (Polo) Chau
Georgia Tech

CSE6242 / CX4242

Sept 9, 2014

Partly based on materials by
Professors Guy Lebanon, Jeffrey Heer, John Stasko, Christos Faloutsos

Last Time

Collection

Cleaning

Integration

Analysis

Visualization

Presentation

Dissemination

Data Cleaning

- Google Refine, Data Wrangler

Data Integration

- Many examples: Google knowledge graph, Facebook Graph Search, Freebase, Feldspar, Kayak, Apple Siri, etc.

Continuing with

Data Integration

Freebase

(a graph of entities)

“...a large collaborative knowledge base consisting of metadata composed mainly by its **community members**...”

Wikipedia.

Crowd-sourcing Approaches: Freebase

The screenshot shows the Freebase website interface. At the top, there is a navigation bar with the Freebase logo, a search bar labeled "Find topics...", and links for "Data", "Schema", "Apps", "Docs", and "Sign In or Sign Up". Below the navigation bar is a blue banner with the text "An entity graph of people, places and things, built by a community that loves open data." and a notice: "Notice: the Freebase Privacy Policy has been updated to the Google Privacy Policy."

The main content area is divided into several sections:

- Featured Data:** A table showing various categories with their member counts, activity trends, and statistics. The categories and their data are as follows:

Category	Members	Last Week Activity	Facts	Topics	Top User
Music	100+	2M	38M	11M	[User Icon]
TV	35	329K	10M	1M	[User Icon]
Film	79	69K	6M	877K	[User Icon]
People	100+	38K	7M	2M	[User Icon]
Business	100+	7K	1M	704K	[User Icon]
Books	46	838	29M	6M	[User Icon]
Location	74	800	8M	1M	[User Icon]
Government	62	511	422K	139K	[User Icon]
- Arts & Entertainment**
- Commons**
- Products & Services**
- Science & Technology**
- Society**
- Special Interests**
- Sports**
- System**
- Time & Space**
- Transportation**
- All**

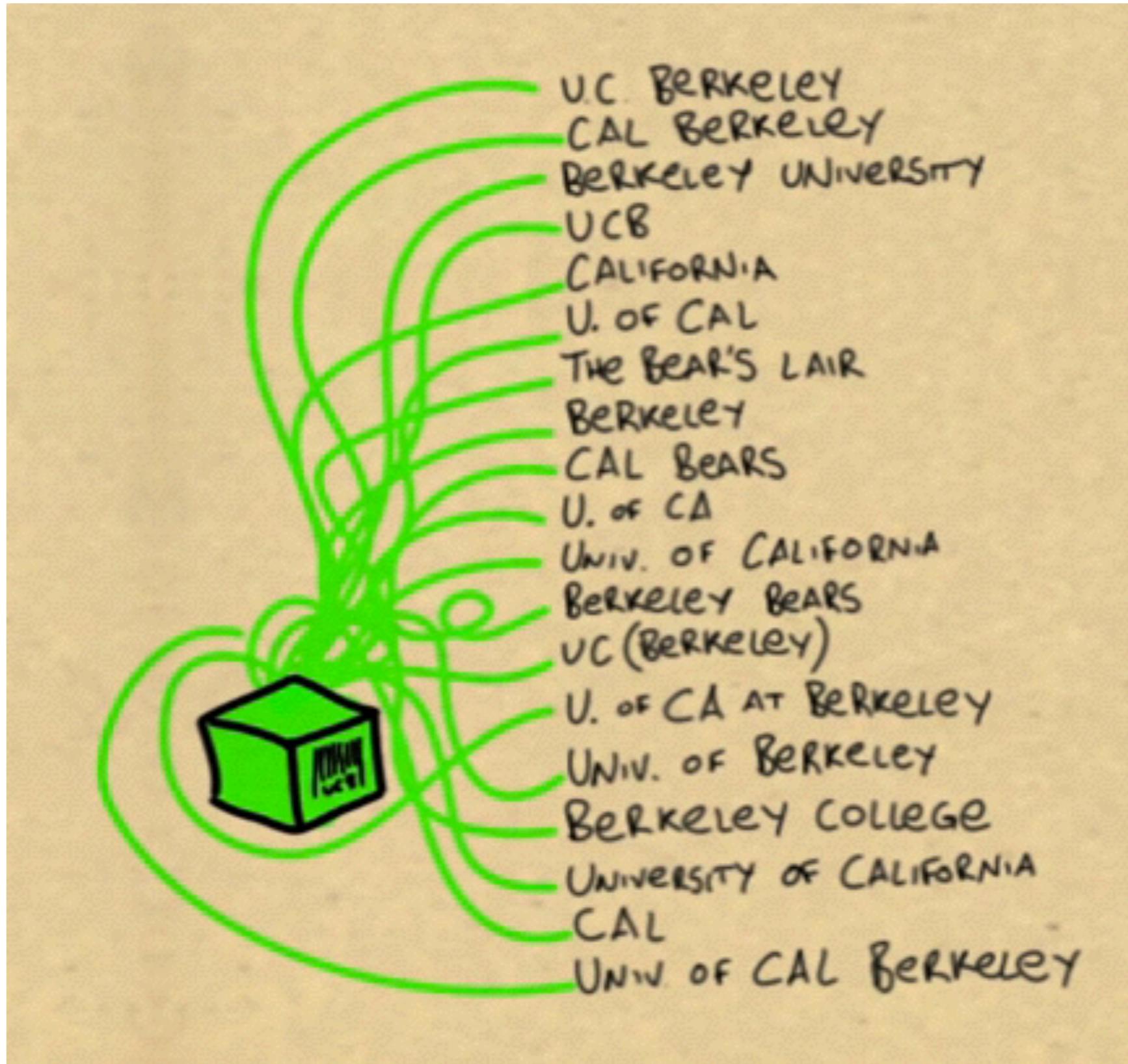
On the right side, there are three promotional boxes:

- Google Refine:** An open source power tool to fix, discover, experiment, connect and customize your data. [Learn more »](#)
- What is Freebase?:** Learn what an entity graph is, what kind of information it contains, and why you should add your data! [Learn More »](#)
- Freebase for Developers:** A list of features including a powerful queryable API, JavaScript-based hosting framework, and libraries for other languages. [Learn More »](#)

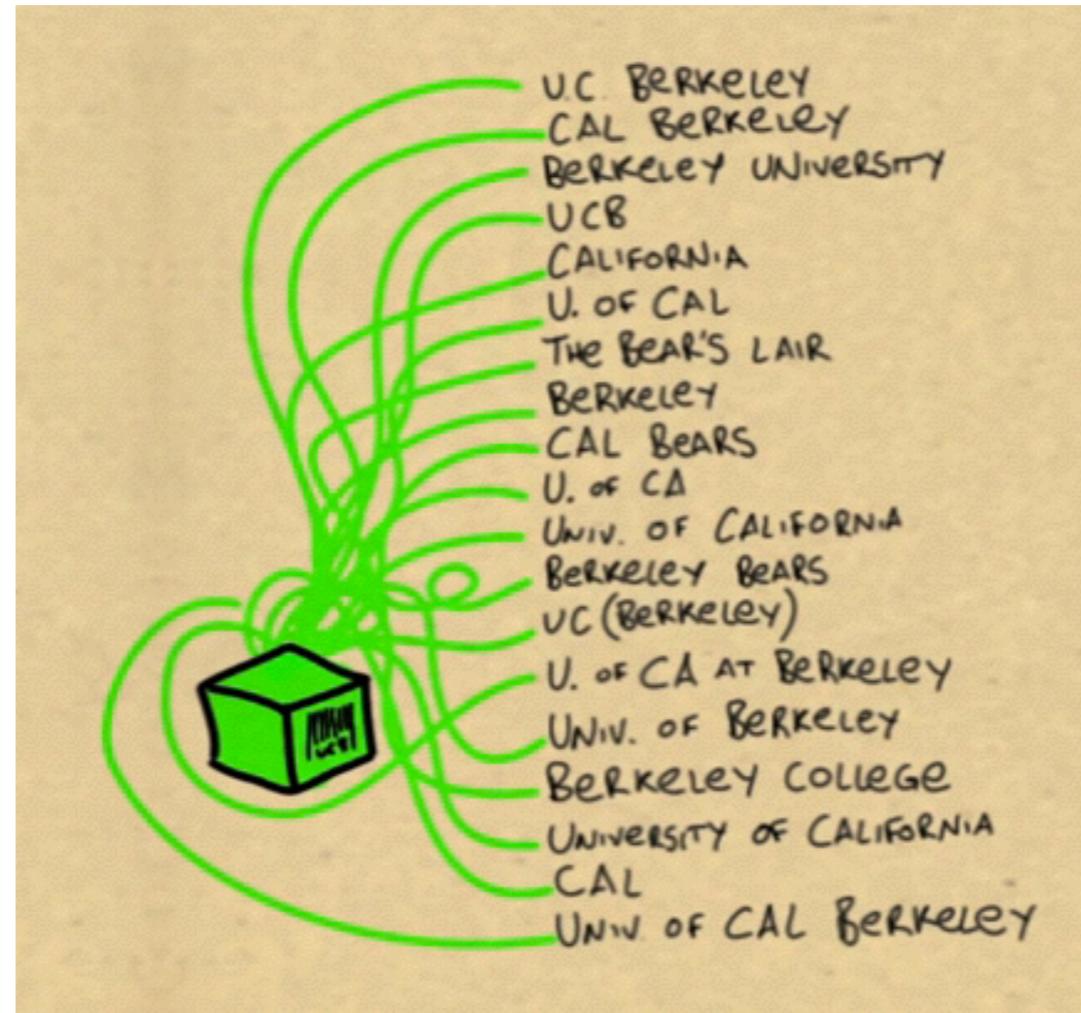
At the bottom, there are four community-related sections:

- Join the Community:** Help the Freebase community create an entity graph of people, places and things, or put it to work for you! [Sign In or Sign Up](#)
- Blog:** Latest posts: "The Freebase blog is moving to Google+" by masouras on May 29, and "Google Refine (previously Freebase Gridworks) 2.0 announced" by skud on November 10.
- Wiki:** Recent changes: [Empty]
- Discussion List:** Enter your email address to join the discussion: [Subscribe](#) [Read the List Archive](#)

What do we **need** before we can even integrate datasets/tables/schemas?



What do we **need** before we can even integrate datasets/tables/schemas?



You need an ID for every unique entity/item/object/thing... Easy?

What do we **need** before we can even integrate datasets/tables/schemas?

person_id	name	state_id
1	Smith	111
2	Johnson	222
3	Obama	222

+

state_id	state_name
111	GA
222	NY
333	CA



person_id	name	state
1	Smith	GA
2	Johnson	NY
3	Obama	NY

Entity Resolution

(A hard problem in data integration)

Polo Chau

P. Chau

Duen Horng Chau

Duen Chau

D. Chau

**Why is Entity Resolution
so Important?**

Related: [iphone 5](#) [iphone 4](#) [iphone unlocked](#) [iphone 3gs](#) [iphone verizon](#) [iphone 5c](#) [iphone 4 unlocked](#) [iphone 3](#) [samsung galaxy s3](#) ...

Include description

Categories

- Cell Phones & Accessories (2,653,244)
- Cell Phone Accessories (2,414,030)
- Cell Phones & Smartphones
- Other (144,703)
- Replacement Parts & Tools (56,276)
- Wholesale Lots (4,886)
- More ▾
- See all categories

Features see all

Contract see all

Condition see all

- New (3,232)
- New other (see details) (1,831)
- Manufacturer refurbished (559)
- Seller refurbished (1,563)
- Used (19,256)
- For parts or not working (8,281)

Price

- Under \$25
 - \$25 - \$50
 - \$50 - \$100
 - Over \$100
- \$ to \$ >>

Format see all

- All Listings (32,722)
- Auction (15,207)
- Buy It Now (23,263)

Item Location see all

- Default
- Within
- of >>
- US Only
- North America
- Worldwide

Delivery Options see all

- Free shipping
- Free in-store pickup

All Listings Auction Buy It Now

Sort: Best Match View:

All > Cell Phones & Accessories > Cell Phones & Smartphones

iphone 32,722 listings [Follow this search](#)

Find Your iPhone

<p>iPhone 5s</p>  <ul style="list-style-type: none"> • 4" Retina Display • True Tone Flash • Slo-mo Video • Touch ID <p>Shop iPhone 5s</p>	<p>iPhone 5c</p>  <ul style="list-style-type: none"> • 4" Retina Display • 8 MP Camera • 1080p HD Video • Multi Color Opt. <p>Shop iPhone 5c</p>	<p>iPhone 5</p>  <ul style="list-style-type: none"> • 4" Retina Display • 8 MP Camera • 1080p HD Video • Face Time, Siri <p>Shop iPhone 5</p>
--	--	---



Apple iPhone 4S A1387, Sprint, 16GB, White, Clean ESN

\$54.00

2 bids

2m left (Today 7:53AM)



Apple iPhone 4 - 8GB - (Verizon) Smartphone - Black or White - Good

USA SELLER *** WARRANTY *** ACCESSORIES INCLUDED

\$79.88

Buy It Now

Free shipping

2053+ Watchers

Save \$10 for every 3 items you buy

FAST 'N FREE

Get it on or before Sat, Sep. 13



Apple iPhone 4 - 16GB - Black (Verizon) Smartphone 7.1.2 MC676LL/A Clean ESN

\$79.00

0 bids

\$125.00

Buy It Now

2m left (Today 7:54AM)

Popular on eBay



Apple iPhone 4 - 8GB - Verizon Straight Talk...

\$109.95

Buy It Now
Free shipping



U Apple iPhone 4 - 8GB - Black (Verizon)...

\$78.95

Buy It Now
Free shipping



U Apple iPhone 4 - 8GB - White (Verizo...

\$79.95

Buy It Now
Free shipping

D-Dupe

Interactive Data Deduplication and Integration
TVCG 2008

University of Maryland

Bilgic, Licamele, Getoor, Kang, Shneiderman

<http://linqs.cs.umd.edu/basilic/web/Publications/2008/kang:tvcg08/kang-tvcg08.pdf>

<http://www.cs.umd.edu/projects/linqs/ddupe/> (skip to 0:55)

Search Potential Duplicate Pairs by Similarity Metric

Potential Duplicate Pairs Similarity Metric

Similarity	Left Node	Right Node
0.982	Elizabeth Churchill	Elizabeth F. Churchill
0.981	Kristian Simsarian	Kristian T. Simsarian
0.981	Gregg Vanderheiden	Gregg C. Vanderheiden
0.981	Christine Neuwirth	Christine M. Neuwirth
0.981	George W. Fitzmaurice	George Fitzmaurice
0.981	Catherine R. Marshall	Catherine C. Marshall
0.980	Pamela K. Schraedley	Pamela Schraedley
0.980	Katherine M. Everitt	Katherine Everitt
0.980	Mja Van Der Wege	Mja M. Van Der Wege
0.980	Elizabeth Veinott	Elizabeth S. Veinott
0.979	Timothy Bickmore	Timothy W. Bickmore

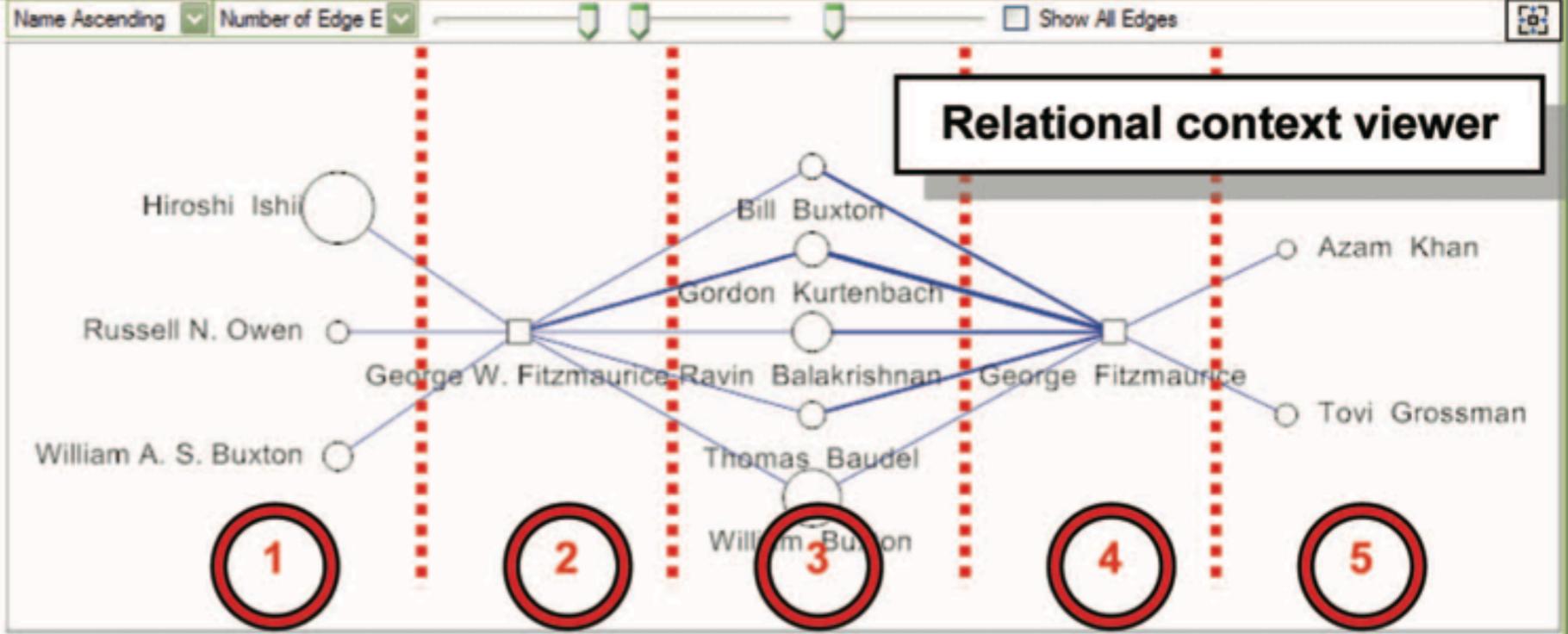
Search Algorithm: Blocking Algorithm - Sample Clustering By Nam

Search Potential Duplicates: Both Within and Across Data Source

Number of Potential Duplicate Pairs (1 ~ 300): 200

Search Potential Duplicate Pairs

Potential duplicate viewer



Potential Duplicates Viewer

person_id	full_name	last_name	first_name	middle_name	suffix	affiliation
P95459	George W. Fitzmaurice	Fitzmaurice	George	W.		
P95460	George Fitzmaurice	Fitzmaurice	George			Alias/wavefront, Toronto, Ontario, Canada and University

Merge Duplicates Mark Distinct

Search Nodes by Keywords

Search

person_id	full_name	last_name	first_name	mid

Search Potential Duplicates of Selected Node

Node Detail Viewer (10 items)

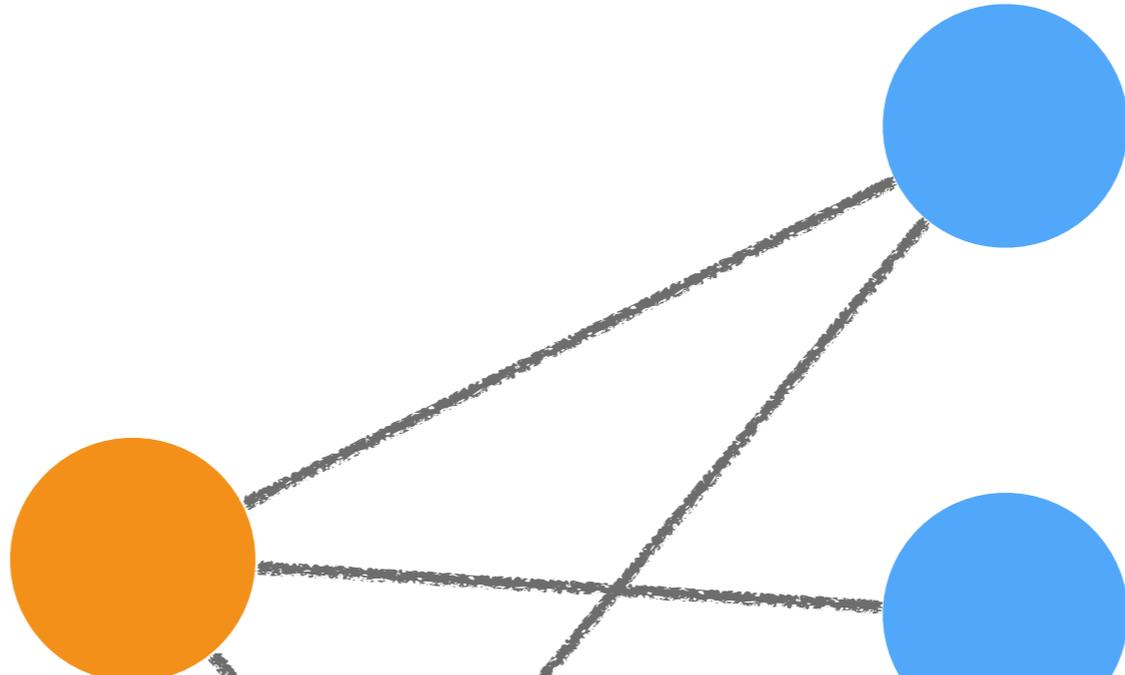
person_id	full_name	last_name	first_name	mid
P110925	Hiroshi Ishii	Ishii	Hiroshi	
P298693	William A. S. Buxton	Buxton	William	A. S.
P250512	Russell N. Owen	Owen	Russell	N.
P284951	Tovi Grossman	Grossman	Tovi	
P23365	Azam Khan	Khan	Azam	

Edge Data

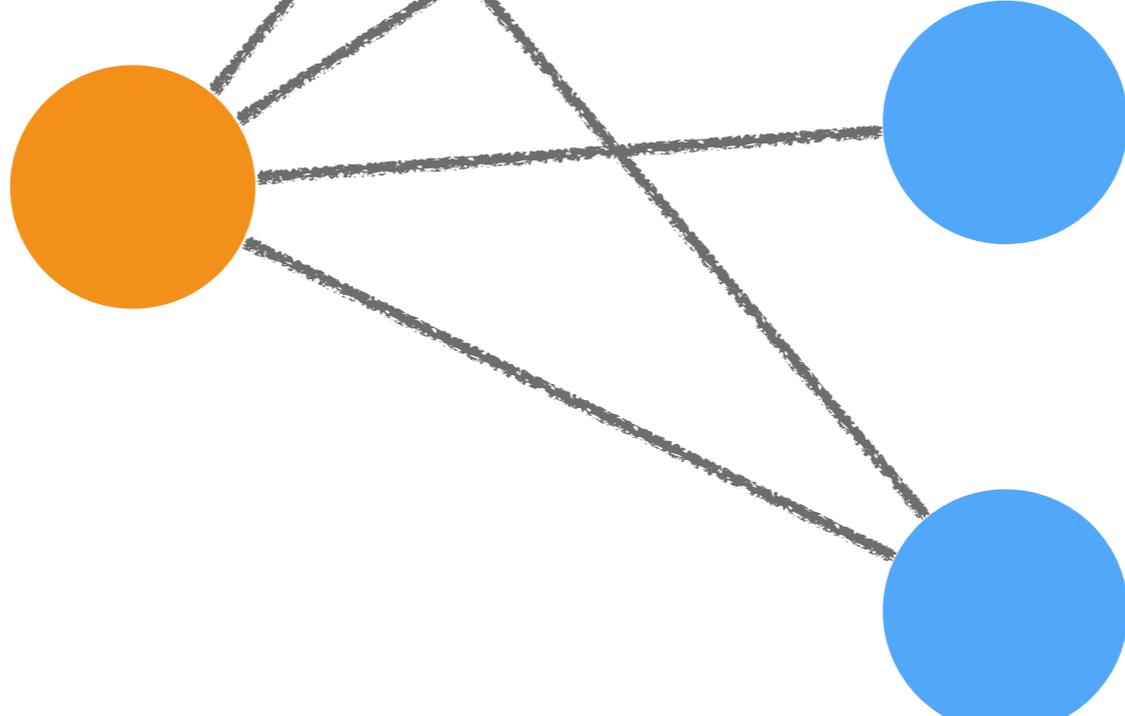
article	
223964	Brooks
303047	The Hotbox
503398	Creating principal 3D curves with digital tape drawing
303033	An exploration into supporting artwork orientation in the user i
258578	An emotional evaluation of orasable user interfaces

Data detail viewer

Polo



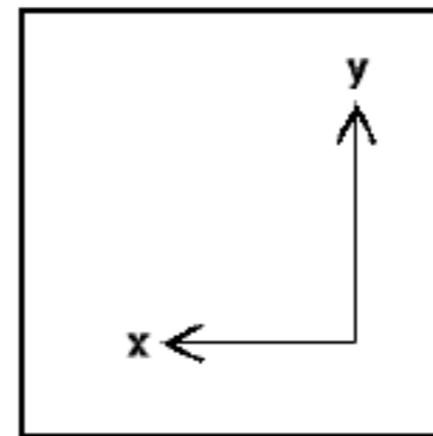
Poalo



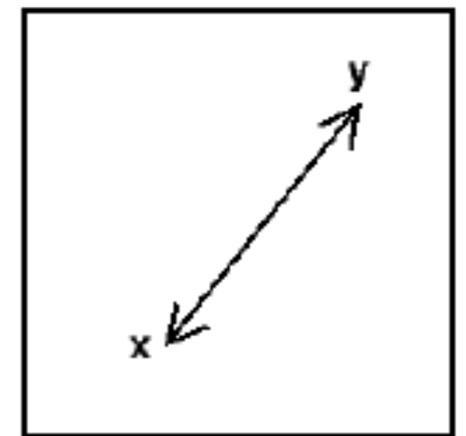
Numerous **similarity** functions

Excellent read: <http://infolab.stanford.edu/~ullman/mmds/ch3a.pdf>

- **Euclidean distance**
Euclidean norm / L2 norm
- **Manhattan distance**
- **Jaccard Similarity**
e.g., overlap of nodes' #neighbors



Manhattan



Euclidean

- **Jaccard Similarity**
e.g., overlap of nodes' #neighbors

Jaccard similarity of sets S and T is $|S \cap T| / |S \cup T|$

- **String edit distance**
e.g., “Polo Chau” vs “Polo Chan”
- **Many more...**

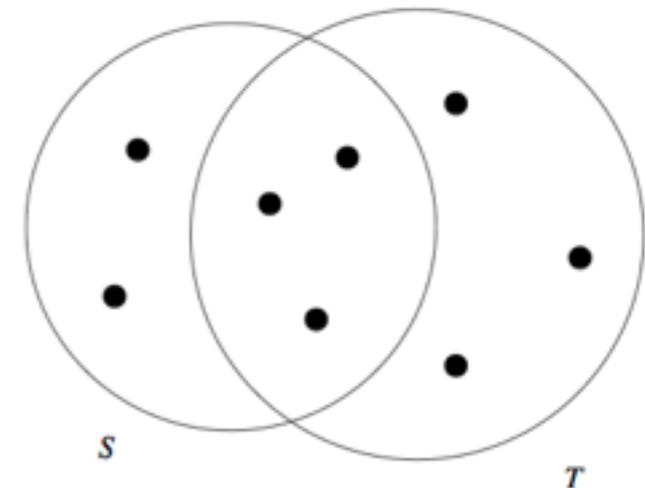


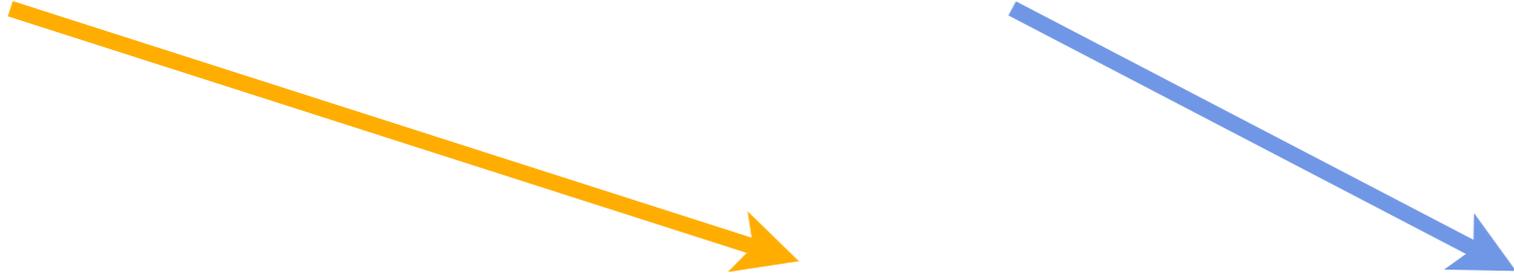
Figure 3.1: Two sets with Jaccard similarity 3/8

Core components: **Similarity functions**

Determine how two entities are similar.

D-Dupe's approach:

Attribute similarity + **relational similarity**


$$\mathit{sim}(e_i, e_j) = (1 - \alpha) \times \mathit{sim}_A(e_i, e_j) + \alpha \times \mathit{sim}_R(e_i, e_j),$$

$$0 \leq \alpha \leq 1,$$



Similarity score for a pair of entities

Attribute similarity (a weighted sum)



$$sim_A(e_i, e_j) = \sum_{k=1}^n w_k \times sim_fun_k(e_i \cdot a_k, e_j \cdot a_k),$$
$$-1 \leq w_k \leq 1 \quad \text{and} \quad \sum_{k=1}^n |w_k| = 1,$$

Summary for data integration

Opportunities

- enable new services (Siri, padmapper)
- enable new ways to discover info
- improve existing services
- reduce redundancy
- new way to interactive with data
- promote knowledge transfer (e.g., between companies)

Data Mining Concepts & Tasks

Collection

Cleaning

Integration

Analysis

Visualization

Presentation

Dissemination

Each data-driven (business, decision-making) problem is **unique**, e.g., different goals, constraints.

Good news: many (sub)tasks that underlie these problems are **common**

Here is an **overview** of the common tasks, based on *Data Science for Business: What you need to know about data mining and data-analytic thinking*

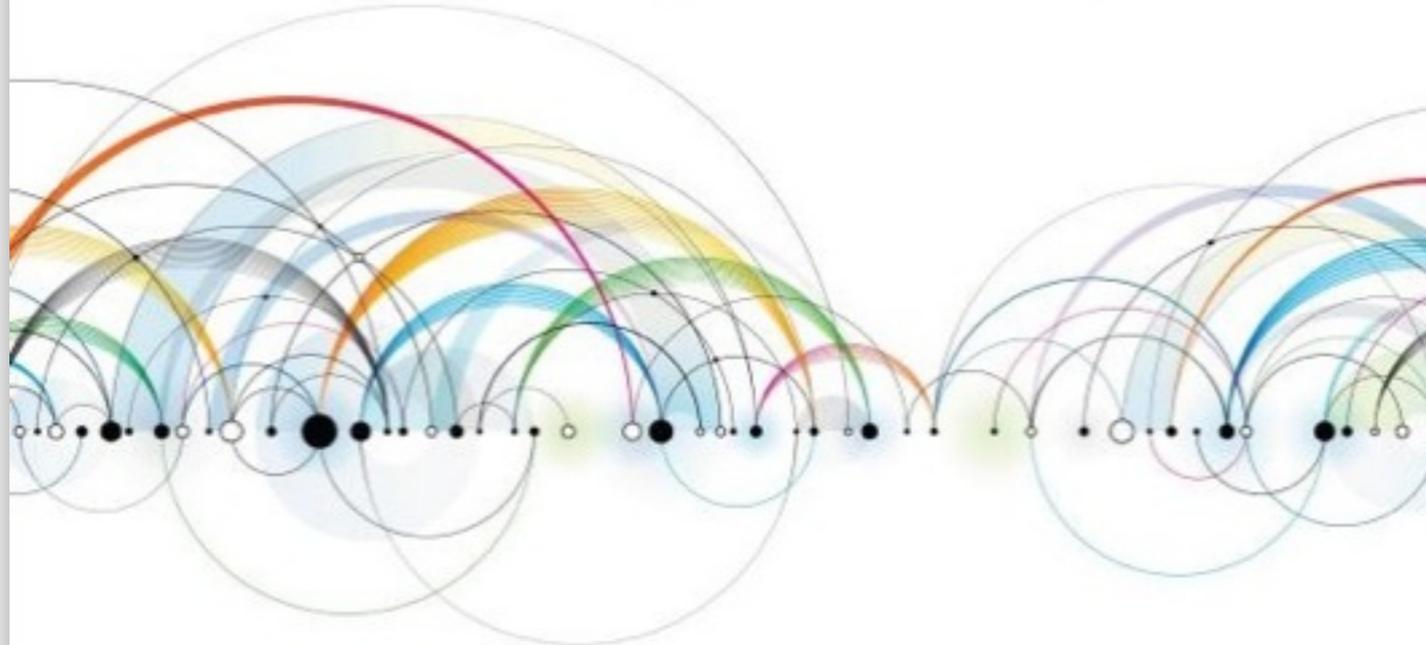
Copyrighted Material

"A must-read resource for anyone who is serious about embracing the opportunity of big data."

—Craig Vaughan, Global Vice President, SAP

Data Science *for* Business

What You Need to Know
About Data Mining and
Data-Analytic Thinking



Foster Provost & Tom Fawcett

Copyrighted Material

1. (soft) Classification, Probability Estimation (supervised learning)

Predict which of a (small) set of classes an entity belong to.

Examples: Is this app malicious or benign? Will this customer click on this ad?

More Examples?

payment transaction -> fraudulent?

news/emails -> spam?

tumor -> benign?

sentiment analysis -> +, -, neutral

weather -> rain, storm, sunny

movies genres -> action, etc.

friends -> close, acquaintance, etc.

online dating -> will work out or not?

surveillance system -> suspicious or not

2. Regression (“value estimation”) (supervised learning)

Predict the **numerical value** of some variable for an entity.

Example: how much minutes will this cellphone customer use?

Related to classification, but predict **how much**, instead of **discrete decisions** (e.g., yes, no)

More Examples?

stock prices

price of plane tickets

weather prediction

credit scores

time until machine fails (data center)

inventory management (supply chain)

population change (city, population planning)

sports stat (gambling)

3. Similarity Matching

Find similar entities (from a large dataset) based on what we know about them.

Examples?

Online dating

recommendation systems (similar songs, movies)

image “classifier” (find all sunset images)

suggestions for online shopping

market segmentation

suggestion of friends on facebook

online advertisement

-> restaurant “classification” (italian, Chinese)

search results (google “similar” results)

search query matching



4. Clustering (unsupervised learning)

Group entities together by their similarity. (User provides # of clusters)

Examples?

factors for diseases

movie categories (genres; soft clustering)

market segmentation for targeted advertisement

social network analysis (whether people like the same thing)

geographical data (identify “neighborhood”, popular landmarks)

5. Co-occurrence grouping

(Many names: frequent itemset mining, association rule discovery, market-basket analysis)

Find associations between entities based on transactions that involve them

(e.g., bread and milk often bought together)



How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did

6. Profiling / Pattern Mining / Anomaly Detection (unsupervised)

Characterize **typical** behaviors of an entity (person, computer router, etc.) so you can find **trends** and **outliers**.

Examples?

computer instruction prediction

removing noise from experiment (data cleaning)

detect anomalies in network traffic

moneyball

weather anomalies (e.g., big storm)

google sign-in (alert)

smart security camera

embezzlement

trending articles



7. Link Prediction / Recommendation

Predict if two entities should be connected, and how strongly that link should be.

Examples?

two people on Facebook

amazon (things bought together); association-rule mining

netflix: recommend jim carey movie

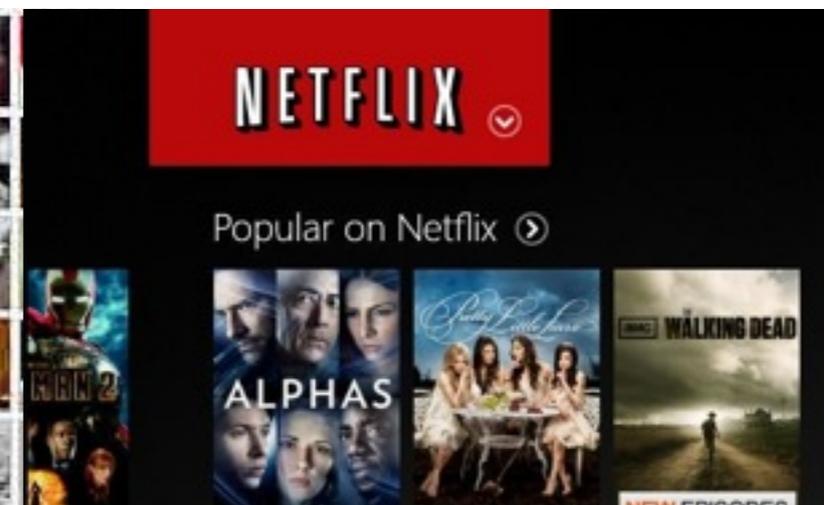
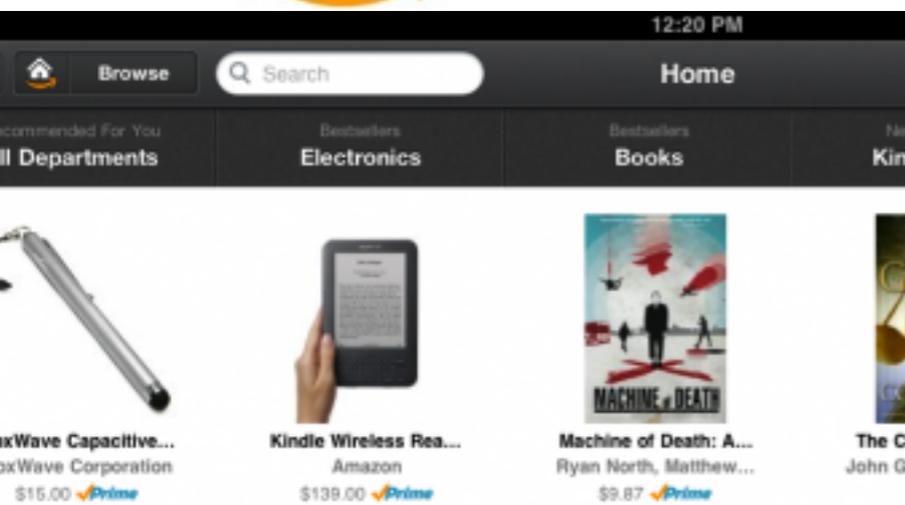
related questions on quora

top apps on apple store

crime group detection (bad guys on social network)

google search suggestions

amazon.com



8. Data reduction (“dimensionality reduction”)

Shrink a large dataset into smaller one, with as little loss of information as possible

When to do it? Examples? Why do it?

Original data is too big -> too hard to process, or take too long

2D -> 1D (many Ds -> few Ds): for visualization, for more efficient algorithms

Graph partitioning - split a large graph into smaller subgraphs

Start thinking about project

- What kind of datasets and problems do you want to solve?
- What techniques do you need?
- Will describe project requirements in next lecture