

Entity Linkage

Shobeir Fakhraei
University of Southern California

Slides based on material provided by Lise Getoor, Ashwin
Machanavajjhala, Craig Knoblock, Jose-Luis Ambite, and others.

Where are we?

- Information Extraction
 - Semi-structured: Web, HTML
 - Unstructured: Text, Ads, Tweets, ...
- Entity Linkage, Data Cleaning, Normalization
- Logical Data Integration
 - Mediators, Query Rewriting
 - Warehouse, Logical Data Exchange
- Automatic Source Modeling/Learning Schema Mappings
- Semantic Web
 - RDF, SPARQL, OWL, Linked Data
- Advanced Topics
 - Geospatial Data Integration, Knowledge Graphs

Information Extraction

As a family
of techniques:

Information Extraction =
segmentation + classification + association + clustering

October 14, 2002, 4:00 a.m. PT

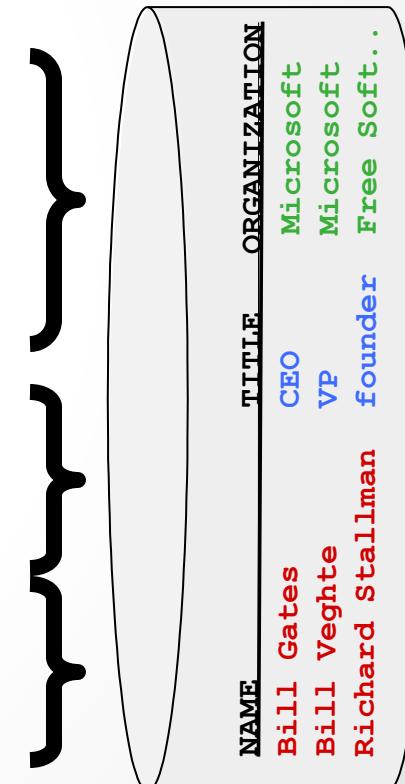
For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."

Richard Stallman, founder of the Free Software Foundation, countered saying...

- * Microsoft Corporation
CEO
Bill Gates
- * Microsoft
Gates
- * Microsoft
Bill Veghte
- * Microsoft
VP
- Richard Stallman**
founder
Free Software Foundation



IE form Ungrammatical Text

Post:

\$25 winning bid at
holiday inn sel. univ. ctr.

Reference Set:

Holiday Inn Select	University Center
Hyatt Regency	Downtown

Ref_hotelName

Ref_hotelArea

Record Linkage

\$25 winning bid at
holiday inn sel. univ. ctr.

Holiday Inn Select | University Center

“\$25”, “winning”, “bid”, ...

Extraction

\$25 winning bid ... <price> \$25 </price> <hotelName> holiday inn
sel.</hotelName> <hotelArea> univ. ctr. </hotelArea>
<Ref_hotelName> Holiday Inn Select </Ref_hotelName>
<Ref_hotelArea> University Center </Ref_hotelArea>

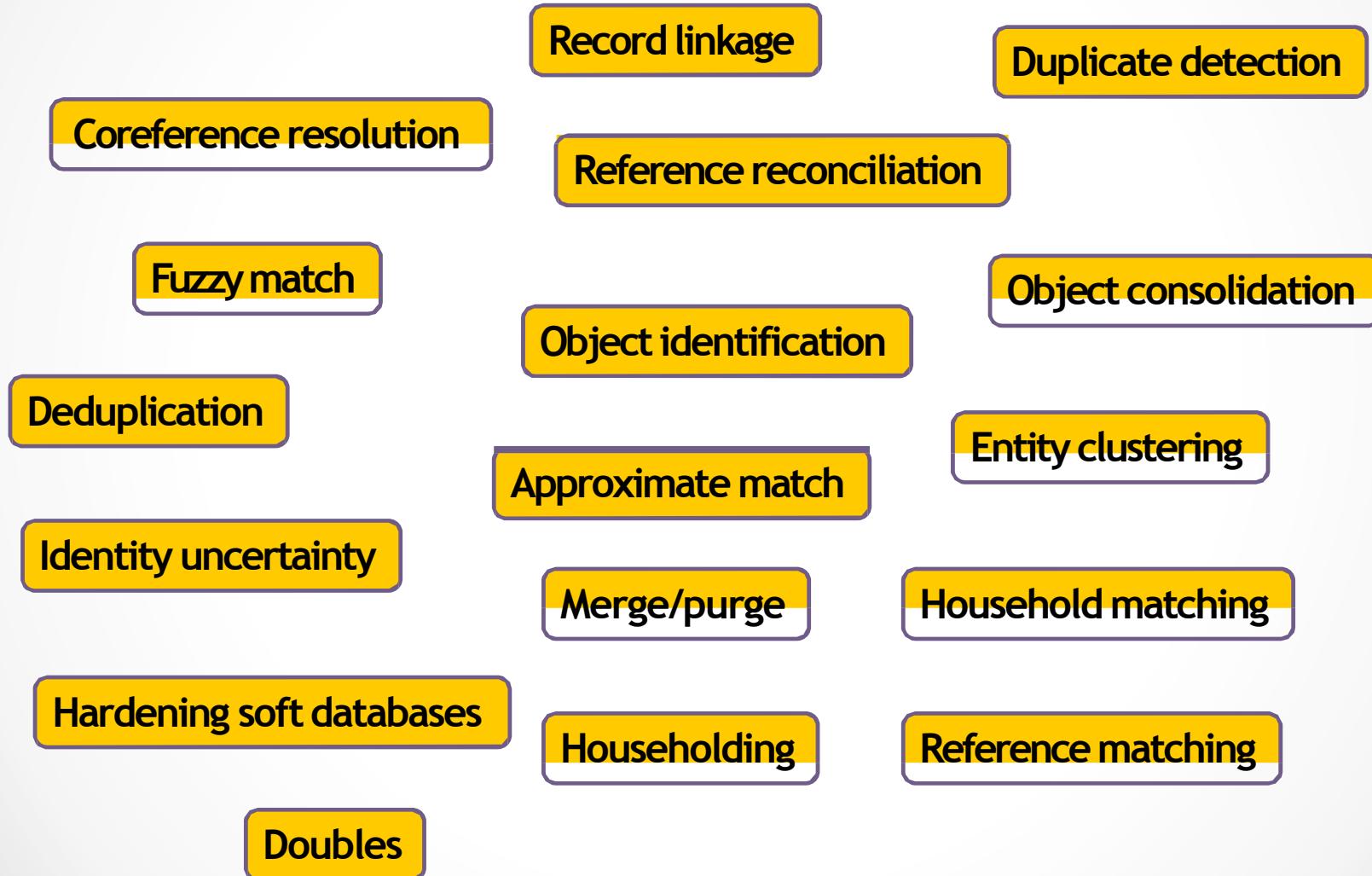
What is Entity Linkage?

Problem of identifying and linking/grouping different manifestations of the same real world object.

Examples of manifestations and objects:

- Different ways of addressing (names, email addresses, FaceBook accounts) the same person in text.
- Web pages with differing descriptions of the same business.
- Different photos of the same object.
- ...

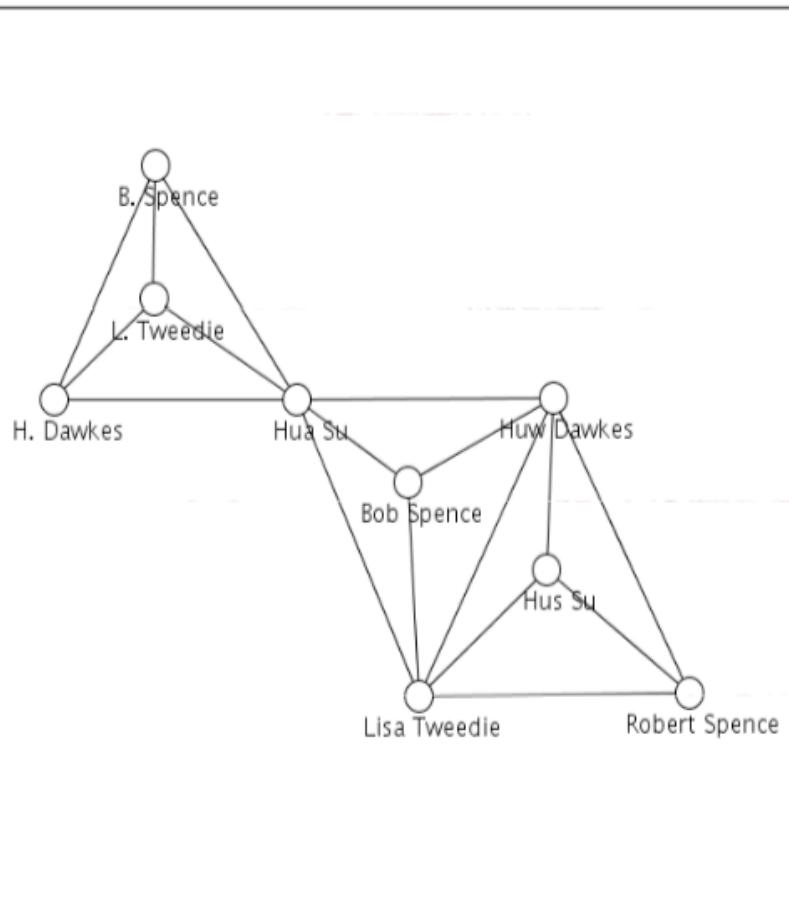
Ironically, Entity Linkage has many duplicate names



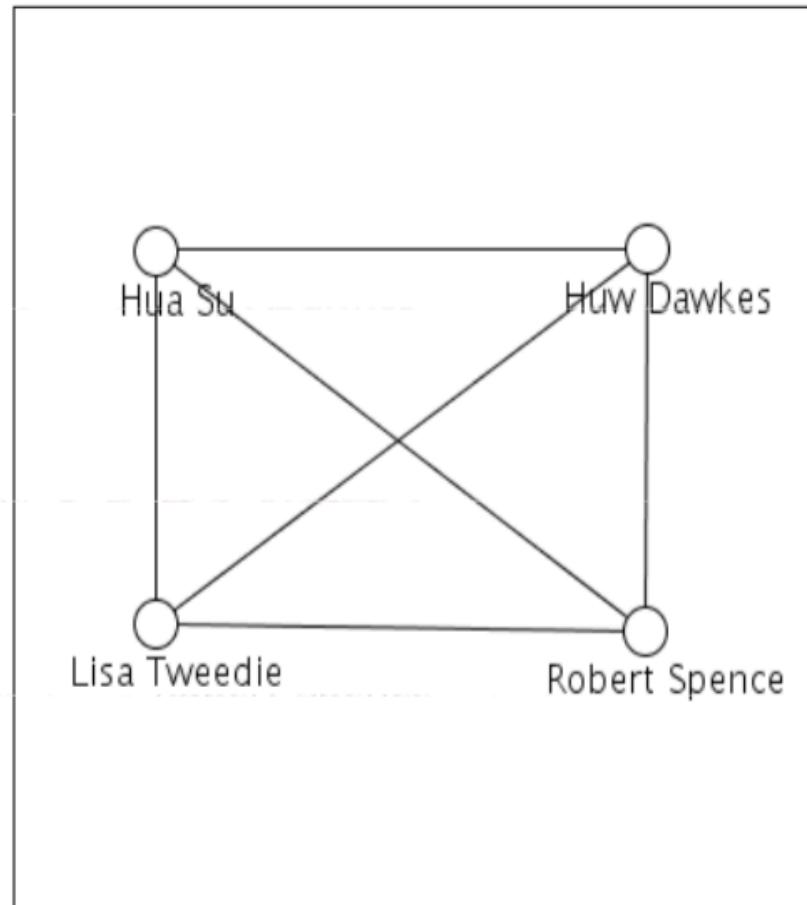
EL Motivating Examples

- *Linking Census Records*
- *Public Health*
- *Web search*
- *Comparison shopping*
- *Counter-terrorism*
- *Spam detection*
- *Machine Reading*
- ...

EL and Network Analysis



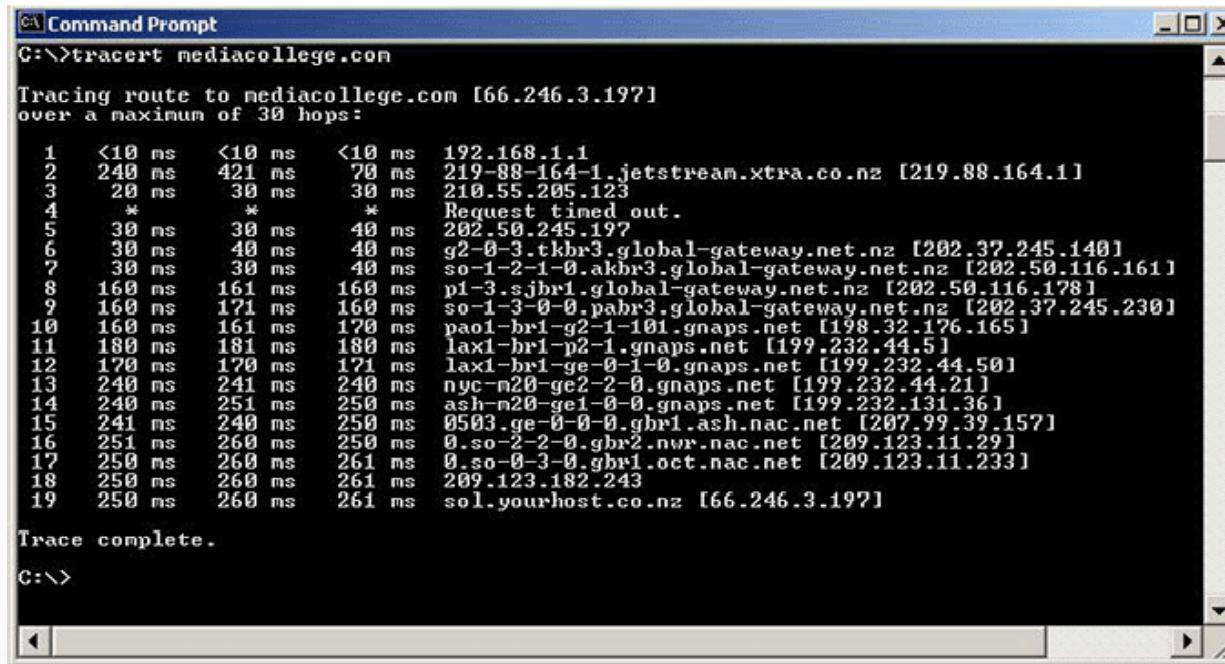
before



after

Motivation: Network Science

- Measuring the topology of the internet ... using **traceroute**



The screenshot shows a Windows Command Prompt window titled "Command Prompt". The command "tracert mediacollege.com" is run, tracing the route to the IP address 66.246.3.197 over a maximum of 30 hops. The output shows the path through various routers and ISPs, with some hops timing out. The window has a standard Windows title bar and scroll bars.

```
C:\>tracert mediacollege.com
Tracing route to mediacollege.com [66.246.3.197]
over a maximum of 30 hops:
 1  <10 ms    <10 ms    <10 ms  192.168.1.1
 2  240 ms    421 ms    70 ms   219-88-164-1.jetstream.xtra.co.nz [219.88.164.1]
 3  20 ms     30 ms    30 ms   210.55.205.123
 4  *          *          * Request timed out.
 5  30 ms     30 ms    40 ms   202.50.245.197
 6  30 ms     40 ms    40 ms   g2-0-3.tkbx3.global-gateway.net.nz [202.37.245.140]
 7  30 ms     30 ms    40 ms   so-1-2-1-0.akbr3.global-gateway.net.nz [202.50.116.161]
 8  160 ms    161 ms    160 ms   p1-3.sjbr1.global-gateway.net.nz [202.50.116.178]
 9  160 ms    171 ms    160 ms   so-1-3-0-0.pabr3.global-gateway.net.nz [202.37.245.230]
10  160 ms    161 ms    170 ms   pa01-br1-g2-1-101.gnaps.net [198.32.176.165]
11  180 ms    181 ms    180 ms   lax1-br1-p2-1.gnaps.net [199.232.44.5]
12  170 ms    170 ms    171 ms   lax1-br1-ge-0-1-0.gnaps.net [199.232.44.50]
13  240 ms    241 ms    240 ms   nyc-m20-ge2-2-0.gnaps.net [199.232.44.21]
14  240 ms    251 ms    250 ms   ash-m20-ge1-0-0.gnaps.net [199.232.131.36]
15  241 ms    240 ms    250 ms   0503.ge-0-0-0.gbr1.ash.nac.net [207.99.39.157]
16  251 ms    260 ms    250 ms   0.so-2-2-0.gbr2.nvr.nac.net [209.123.11.29]
17  250 ms    260 ms    261 ms   0.so-0-3-0.gbr1.oct.nac.net [209.123.11.233]
18  250 ms    260 ms    261 ms   209.123.182.243
19  250 ms    260 ms    261 ms   sol.yourhost.co.nz [66.246.3.197]

Trace complete.
C:\>
```

IP Aliasing Problem

[Willinger et al. 2009]

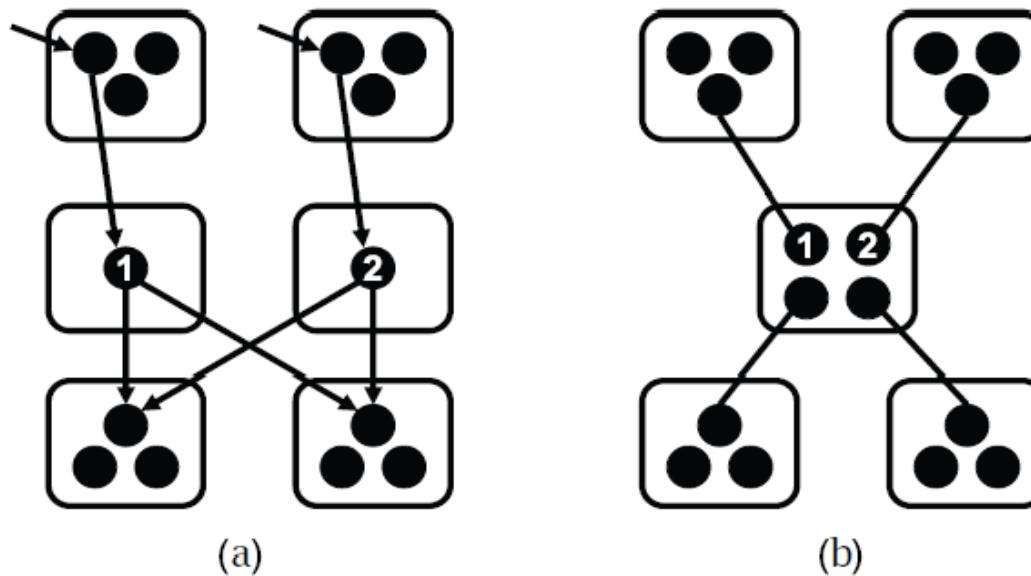


Figure 2. The IP alias resolution problem.
Paraphrasing Fig. 4 of [50], traceroute does not list routers (boxes) along paths but IP addresses of input interfaces (circles), and alias resolution refers to the correct mapping of interfaces to routers to reveal the actual topology. In the case where interfaces 1 and 2 are aliases, (b) depicts the actual topology while (a) yields an “inflated” topology with more routers and links than the real one.

IP Aliasing Problem

[Willinger et al. 2009]

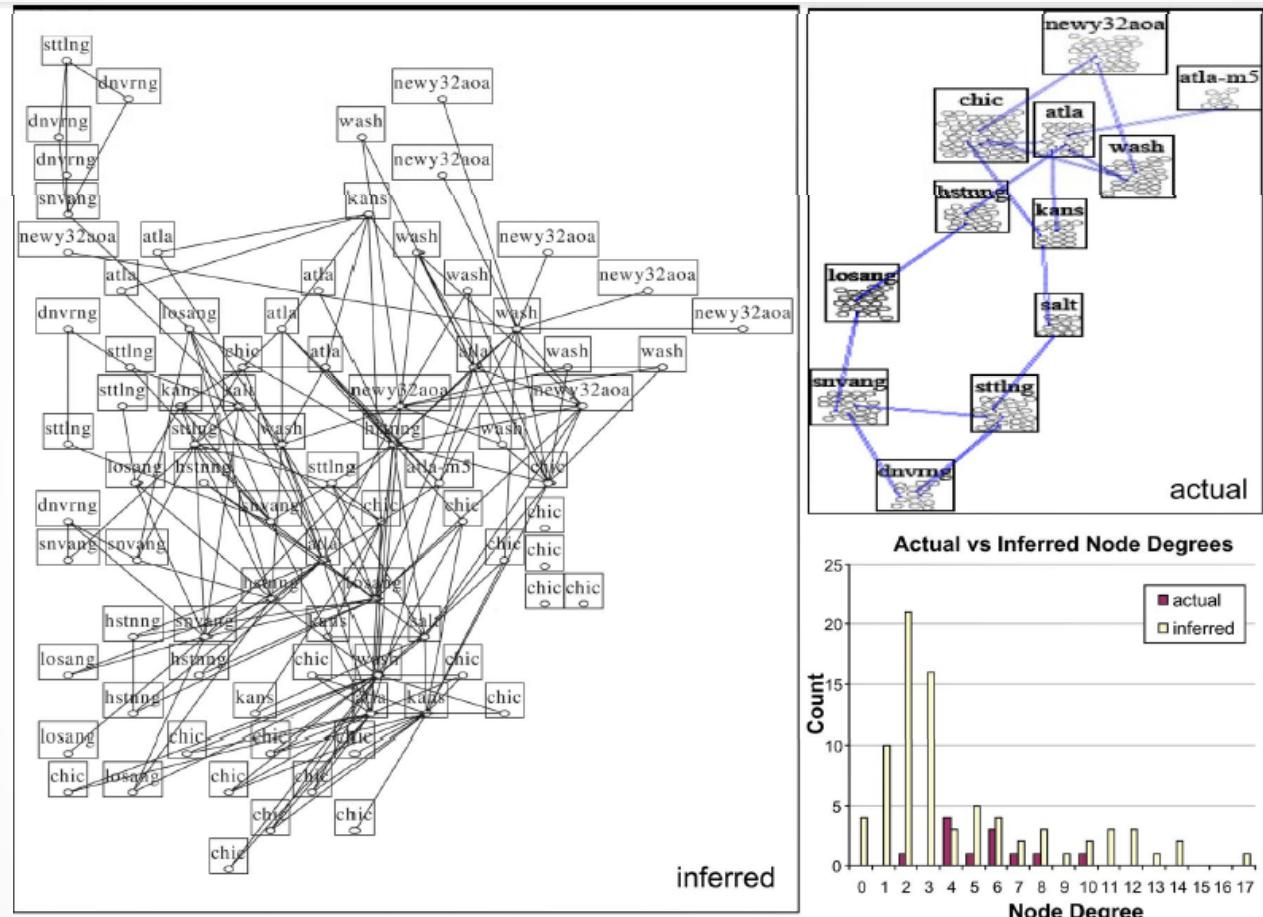


Figure 3. The IP alias resolution problem in practice. This is re-produced from [48] and shows a comparison between the Abilene/Internet2 topology inferred by Rocketfuel (left) and the actual topology (top right). Rectangles represent routers with interior ovals denoting interfaces. The histograms of the corresponding node degrees are shown in the bottom right plot. © 2008 ACM,

Example: Google Contacts

Google Search for anyone +Sean All Contacts Starred Sean Purcell Keri Andersen hikingfan@gmail.com 650-555-7646 Anissa Lee anissa.v.lee@gmail.com 646-555-8226

Frequently contacted Nick Siegel origamifan26@gmail.com

Nicole Bonsavage peanutbutterfan17@gmail.com

James corgicrazy@gmail.com

Tom Green tgreen198@gmail.com

Lindsay Carter yogafan011@gmail.com

Michael Potts bikingfan22@gmail.com

Merrill Alexander hockeyfanatic30@gmail.com

All contacts (99) Anissa Lee anissa.v.lee@gmail.com James Carroway corgicrazy@gmail.com

Find duplicates Merge all (36)

Keri Andersen

Keri Andersen hikingfan@gmail.com Home

Keri Andersen 650-555-7646 Mobile

James Carroway

James Carroway corgicrazy@gmail.com 555-234-7654 Home Mobile

James Carroway surfingfan6@gmail.com Home

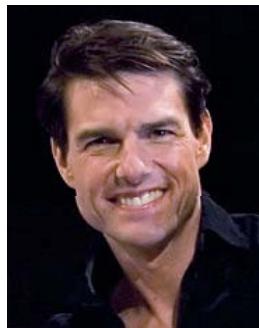
Dismiss Merge

Dismiss Merge

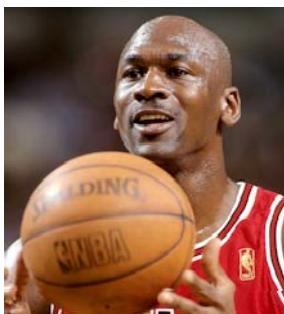
Traditional Challenges in ER

- Name/Attribute ambiguity

Thomas Cruise



Michael Jordan



Traditional Challenges in ER

- Name/Attribute ambiguity
- Errors due to data entry



↓	C1	C2
	Total Cholesterol_1	Total Cholesterol_2
682	214.4	214.4
683	184.4	184.4
684	183.5	183.5
685	240.7	240.7
686	215.1	215.1
687	198.6	198.6
688	2800.0	280.0
689	210.8	210.8
690	182.5	182.5
691	192.6	192.6

Traditional Challenges in ER

- Name/Attribute ambiguity
- Errors due to data entry
- Missing Values

Exhibit 2: Examples of variables that are set to unknown values

Administrative dates: set to 0101YY, 010199, 999999

Date of Birth 0101YY, 1506YY, 3006YY, 0107YY, 1507YY, 0101YEAR

Names: set to spaces, NK, UNKNOWN, or ZZZZ
BABY, MALE, FEMALE, TWIN, TRIPLET, INFANT

Other variables: set to 9, 99, 9999, -1
NK (Not Known)
NA (Not applicable)
NC (Not coded)
U (Unknown)

[Gill et al; Univ of Oxford2003]

Traditional Challenges in ER

- Name/Attribute ambiguity
 - Errors due to data entry
 - Missing Values
 - Changing Attributes
-
- Data formatting

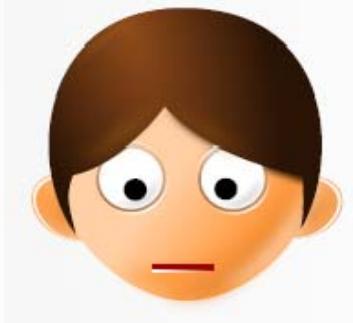


- Abbreviations / Data Truncation

ABSTRACT PROBLEM STATEMENT

Abstract Problem Statement

Real World



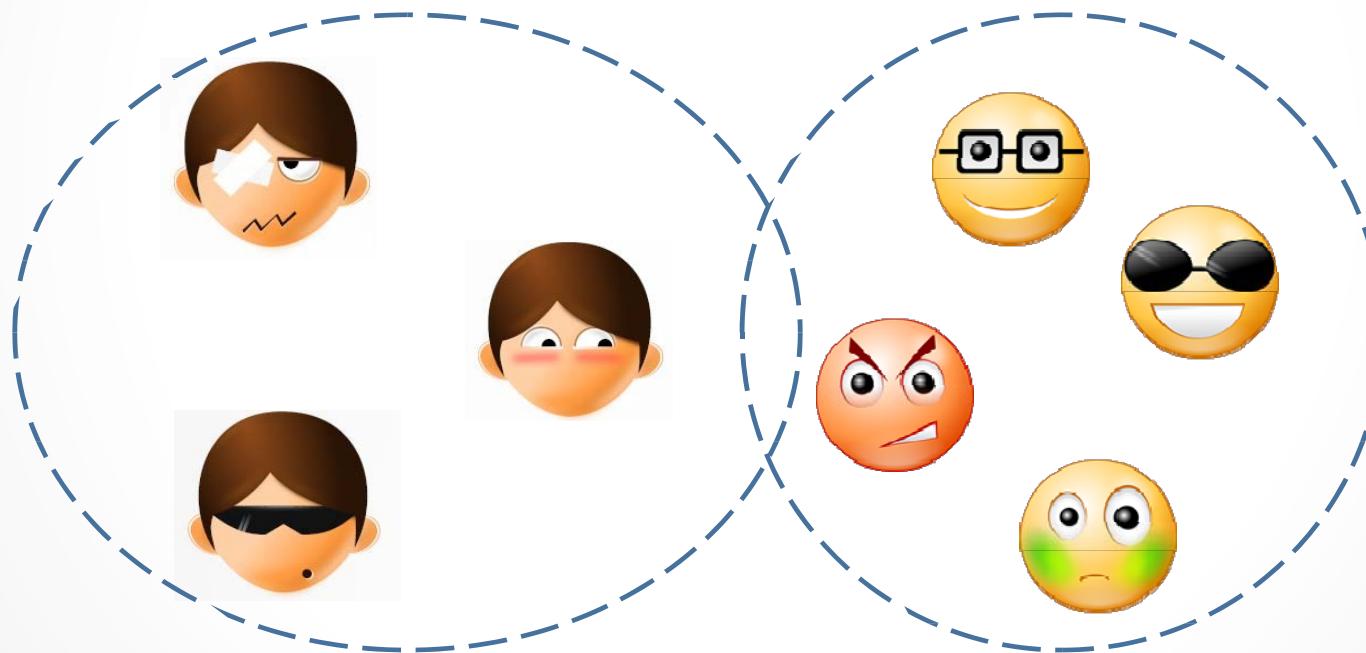
Digital World



Records /
Mentions

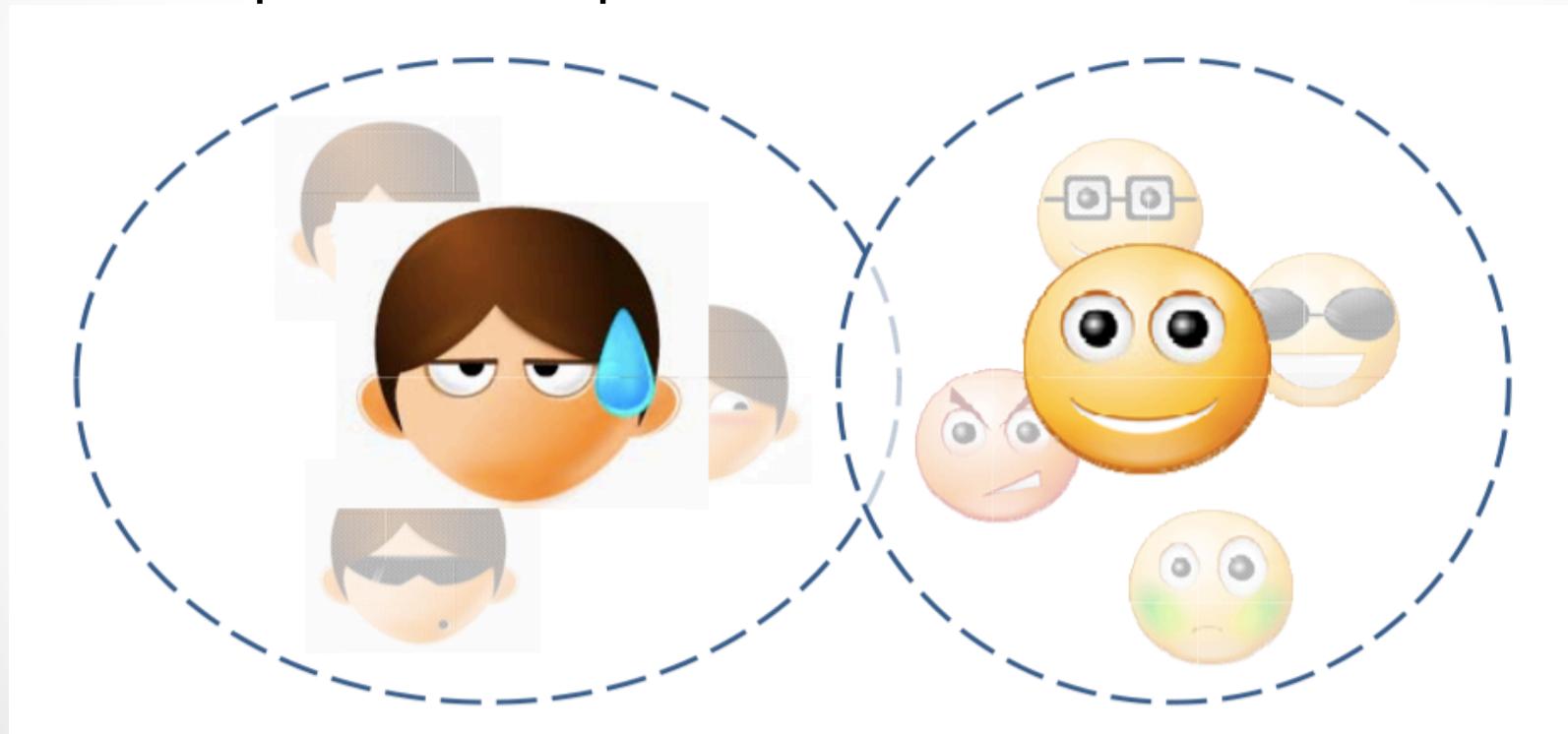
Deduplication Problem Statement

- Cluster the records/mentions that correspond to same entity



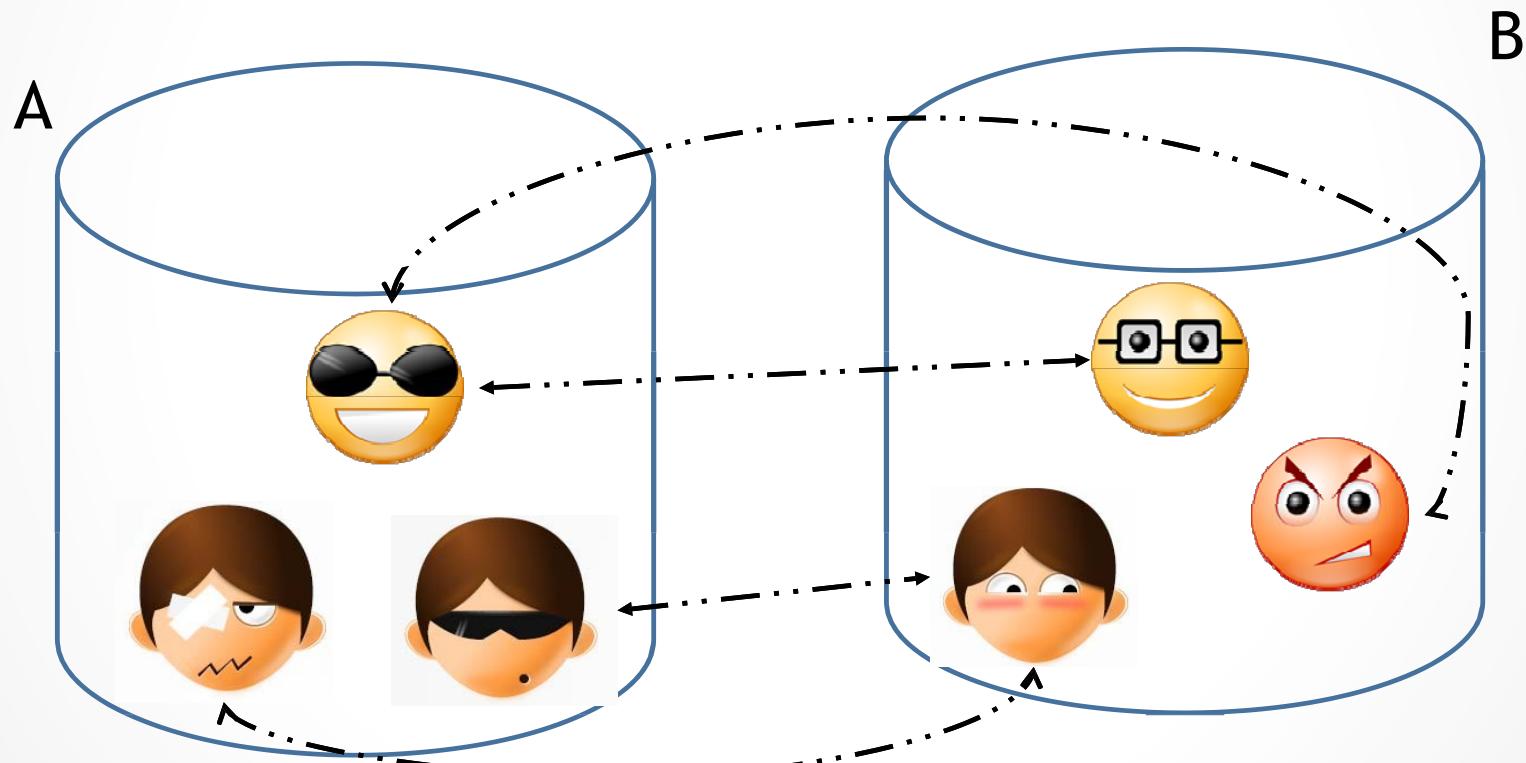
Deduplication Problem Statement

- Cluster the records/mentions that correspond to same entity
 - Compute cluster representative



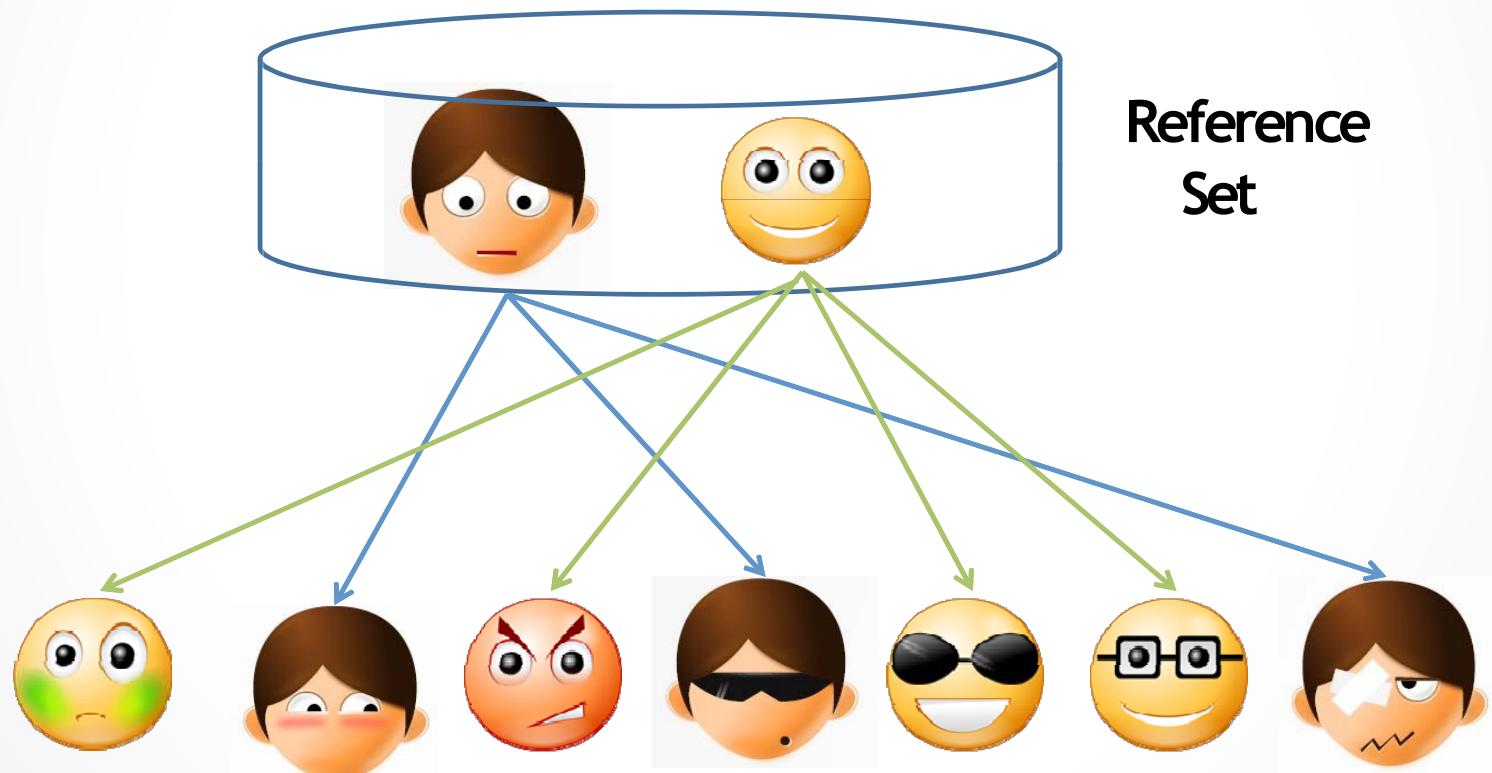
Record Linkage Problem Statement

- Link records that match across databases



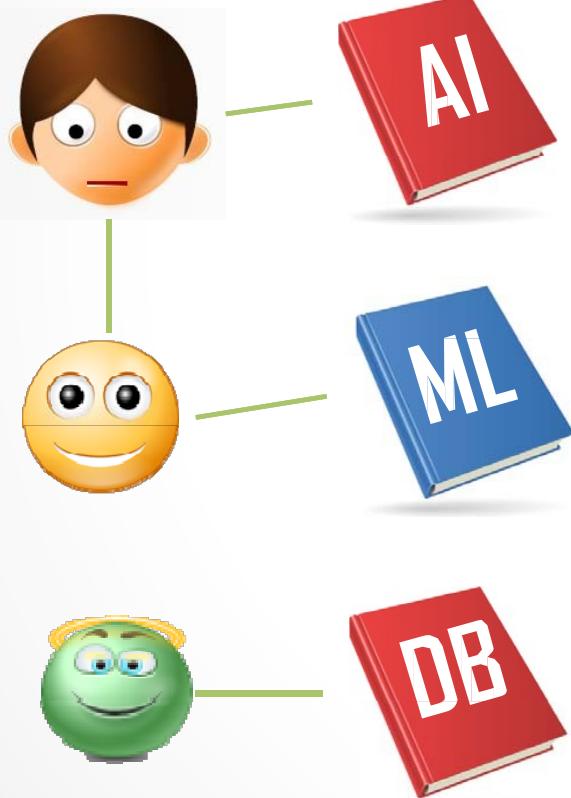
Reference Matching Problem

- Match noisy records to clean records in a reference table

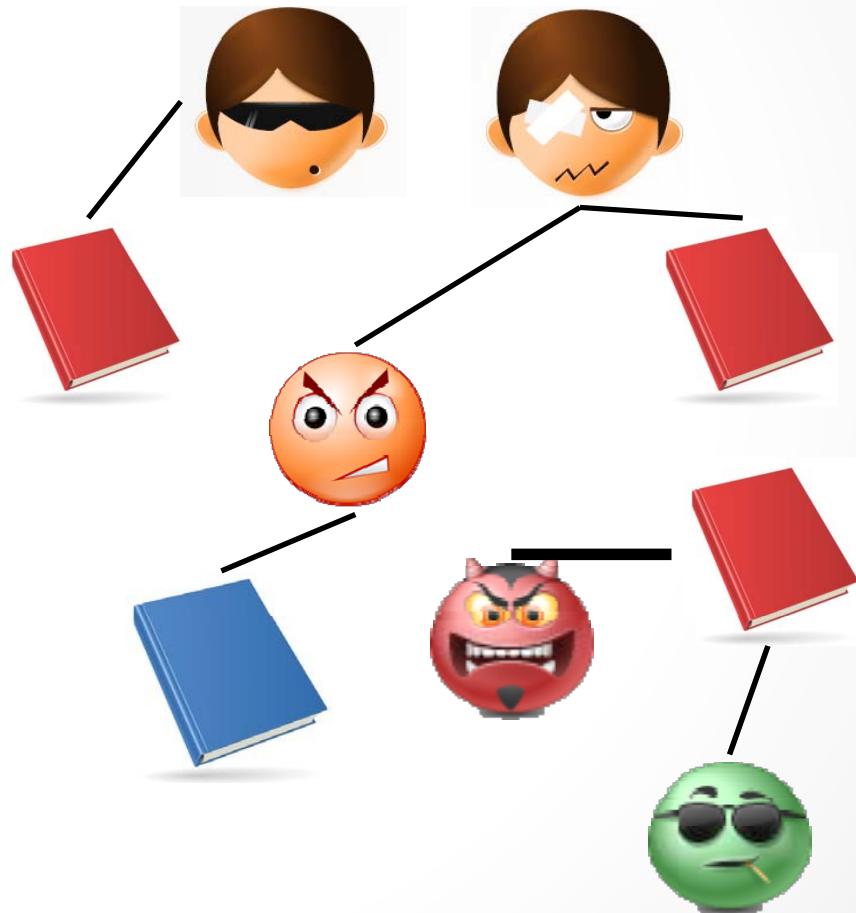


Abstract Problem Statement

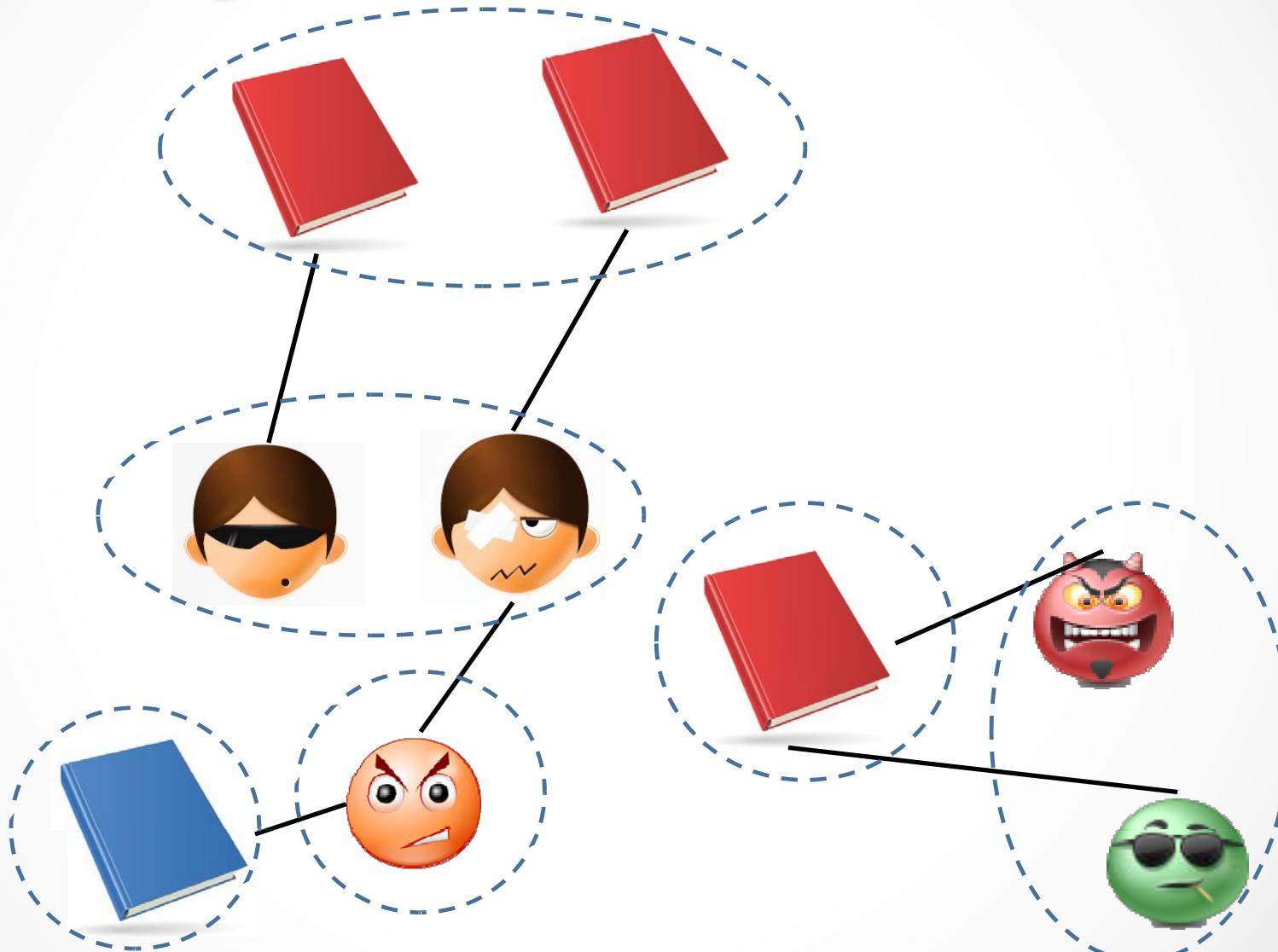
Real World



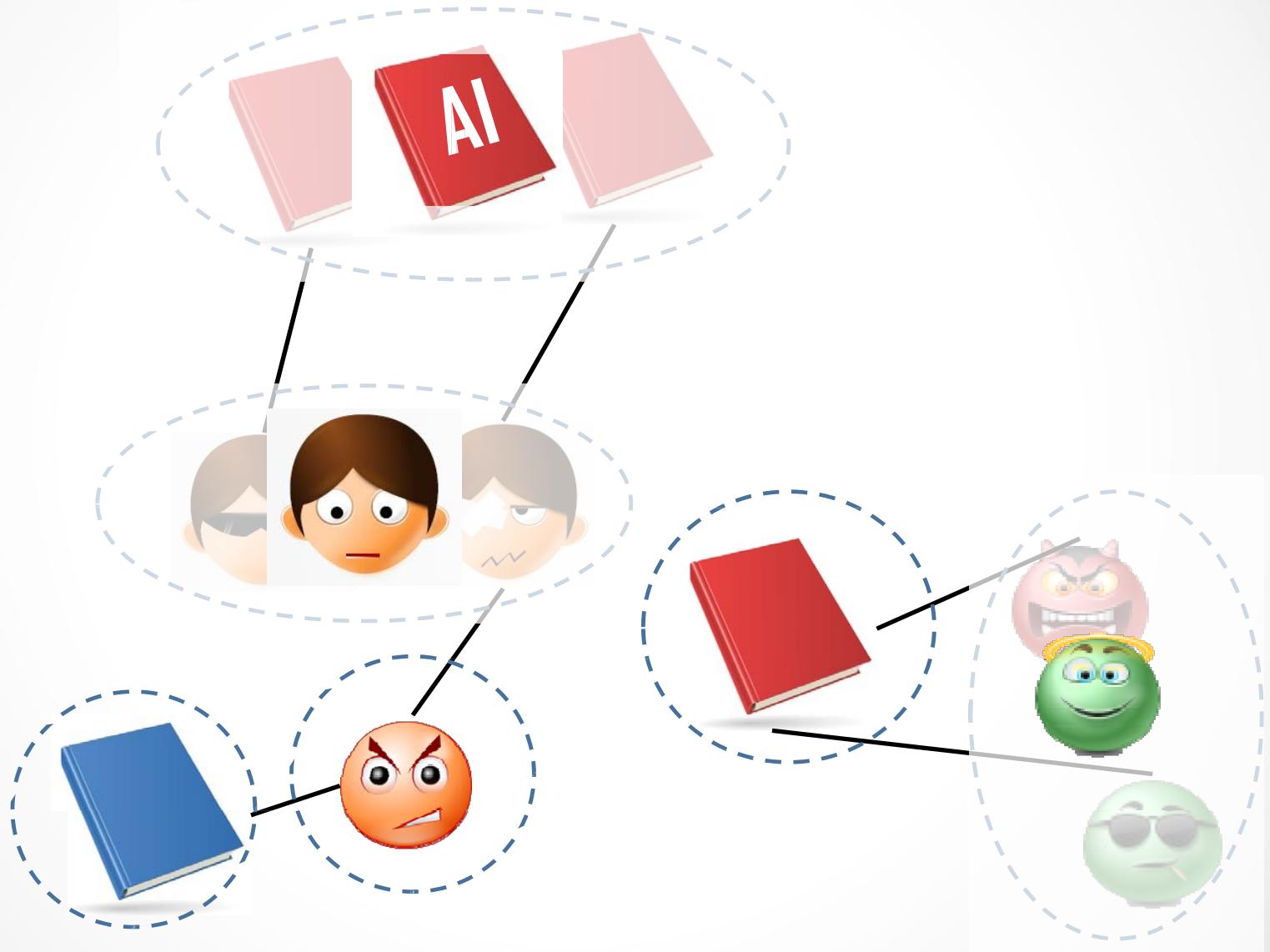
Digital World



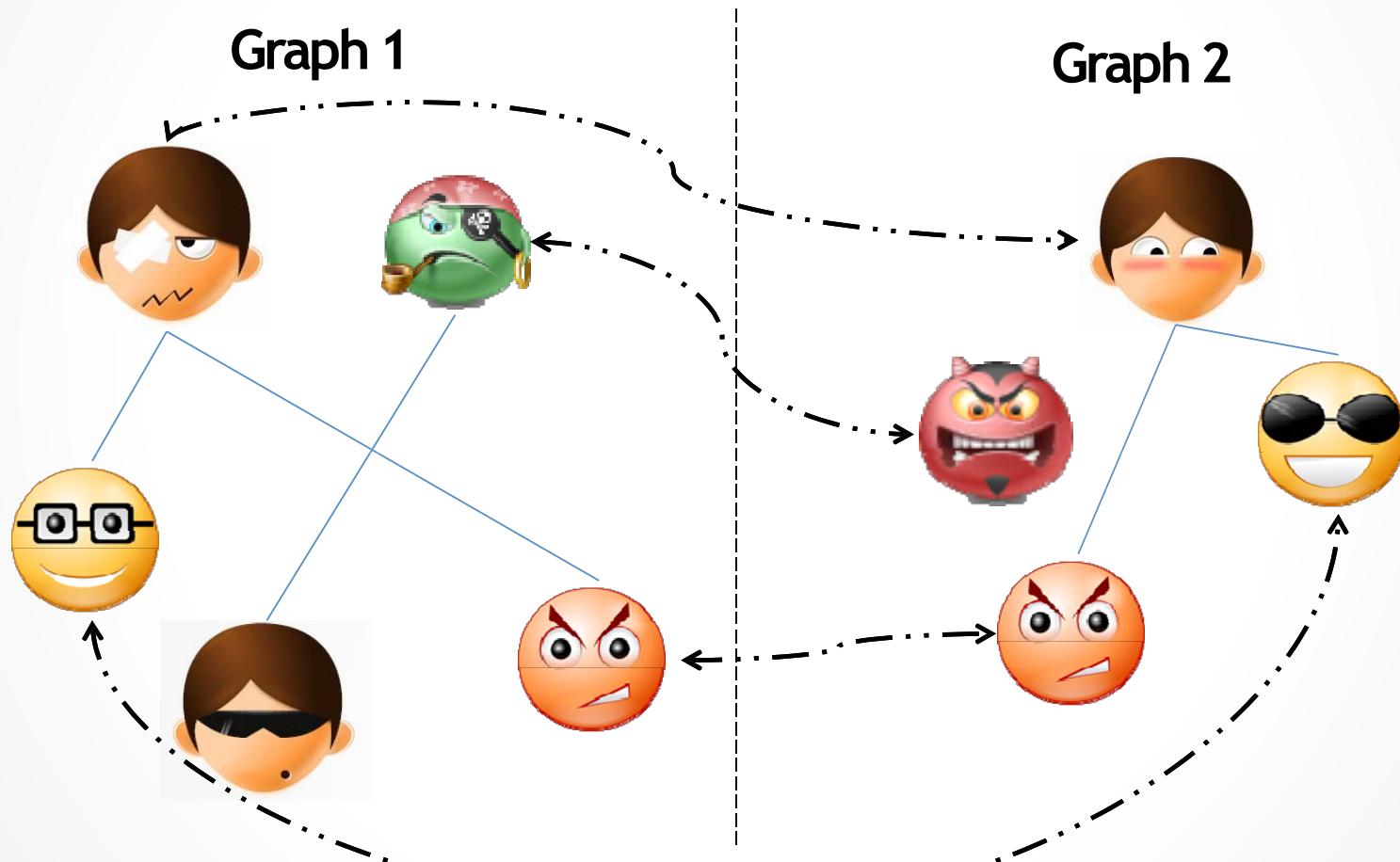
Deduplication Problem Statement



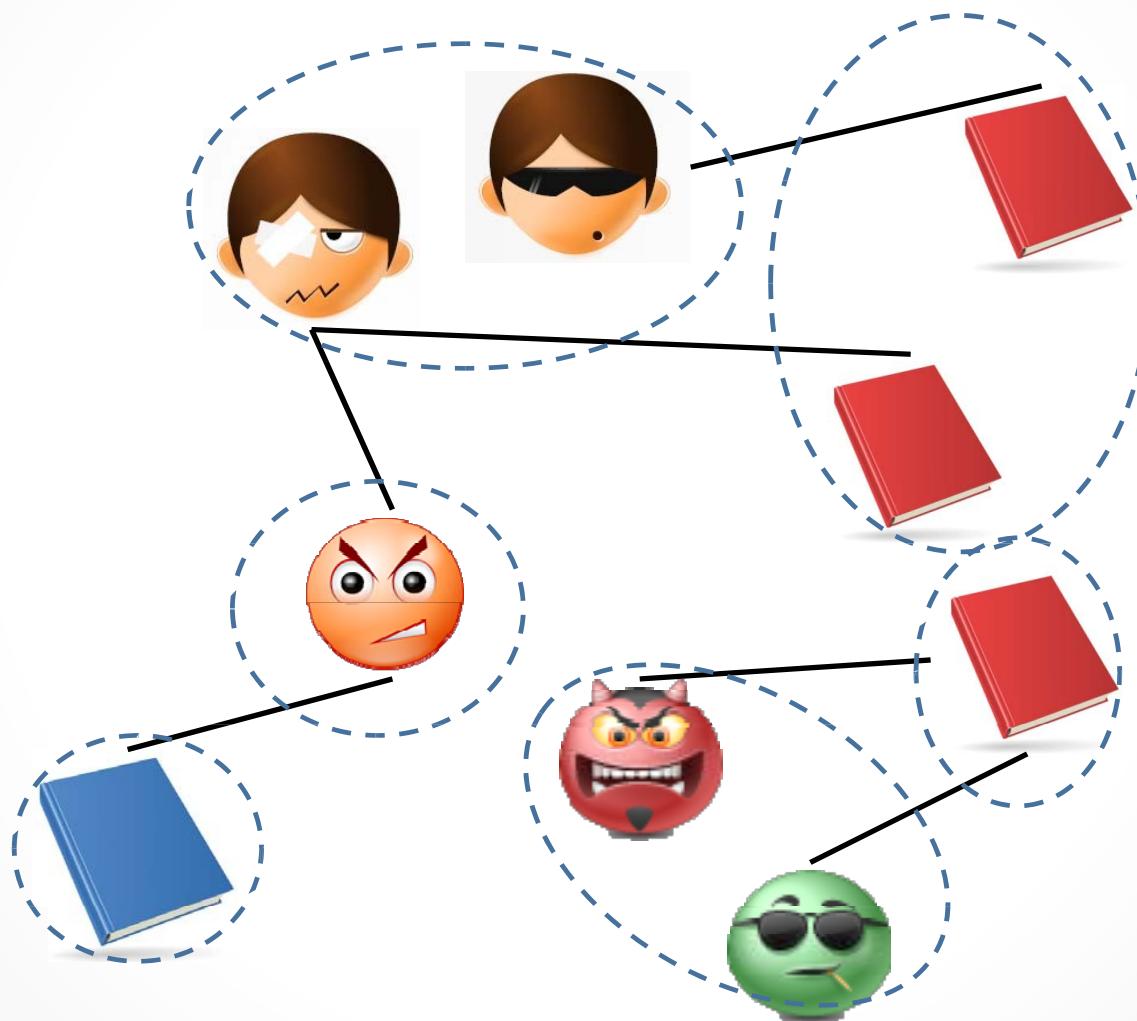
Deduplication with Canonicalization



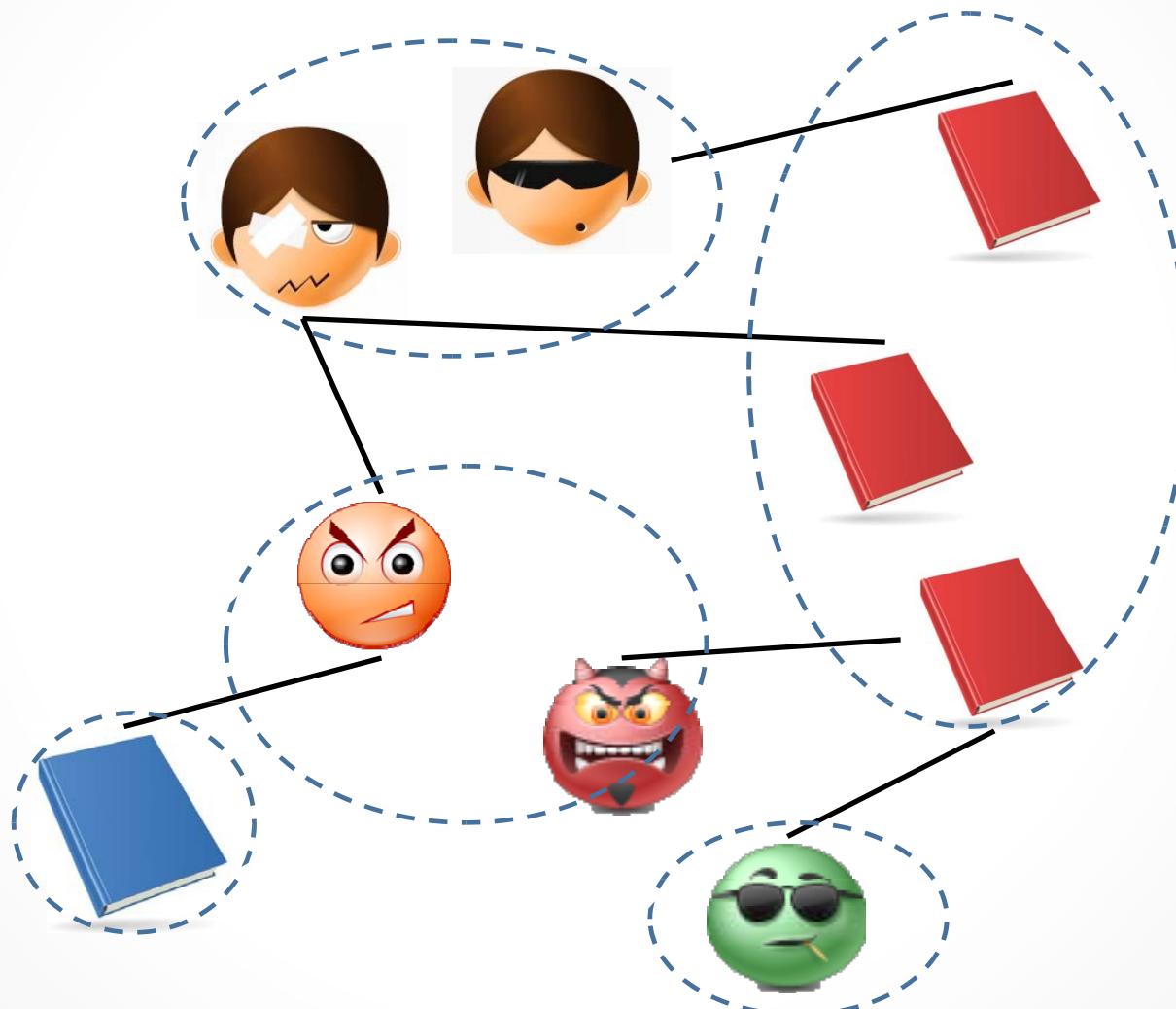
Graph Alignment



Relationships are crucial



Relationships are crucial



Problem Definition (textbook)

- Given two relational tables X and Y with identical schemas
 - assume each tuple in X and Y describes an entity (e.g., person)
- We say tuple x from X matches tuple y from Y if they refer to the same real-world entity
 - (x,y) is called a match
- Goal: find all matches between X and Y

Example

Table X

	Name	Phone	City	State
X ₁	Dave Smith	(608) 395 9462	Madison	WI
X ₂	Joe Wilson	(408) 123 4265	San Jose	CA
X ₃	Dan Smith	(608) 256 1212	Middleton	WI

(a)

Table Y

	Name	Phone	City	State
Y ₁	David D. Smith	395 9426	Madison	WI
Y ₂	Daniel W. Smith	256 1212	Madison	WI

(b)

Matches

(x₁, y₁)
(x₃, y₂)

(c)

■ Other variations

- Tables X and Y have different schemas
- Match tuples within a single table X
- The data is not relational, e.g., XML or RDF

Why is This Different than String Matching?

- In theory, can treat each tuple as a string by concatenating the fields, then apply string matching techniques
- But doing so makes it hard to apply sophisticated techniques and domain-specific knowledge
- E.g., consider matching tuples that describe persons
 - suppose we know that in this domain two tuples match if the names and phone match exactly
 - this knowledge is hard to encode if we use string matching
 - it is better to keep the fields apart

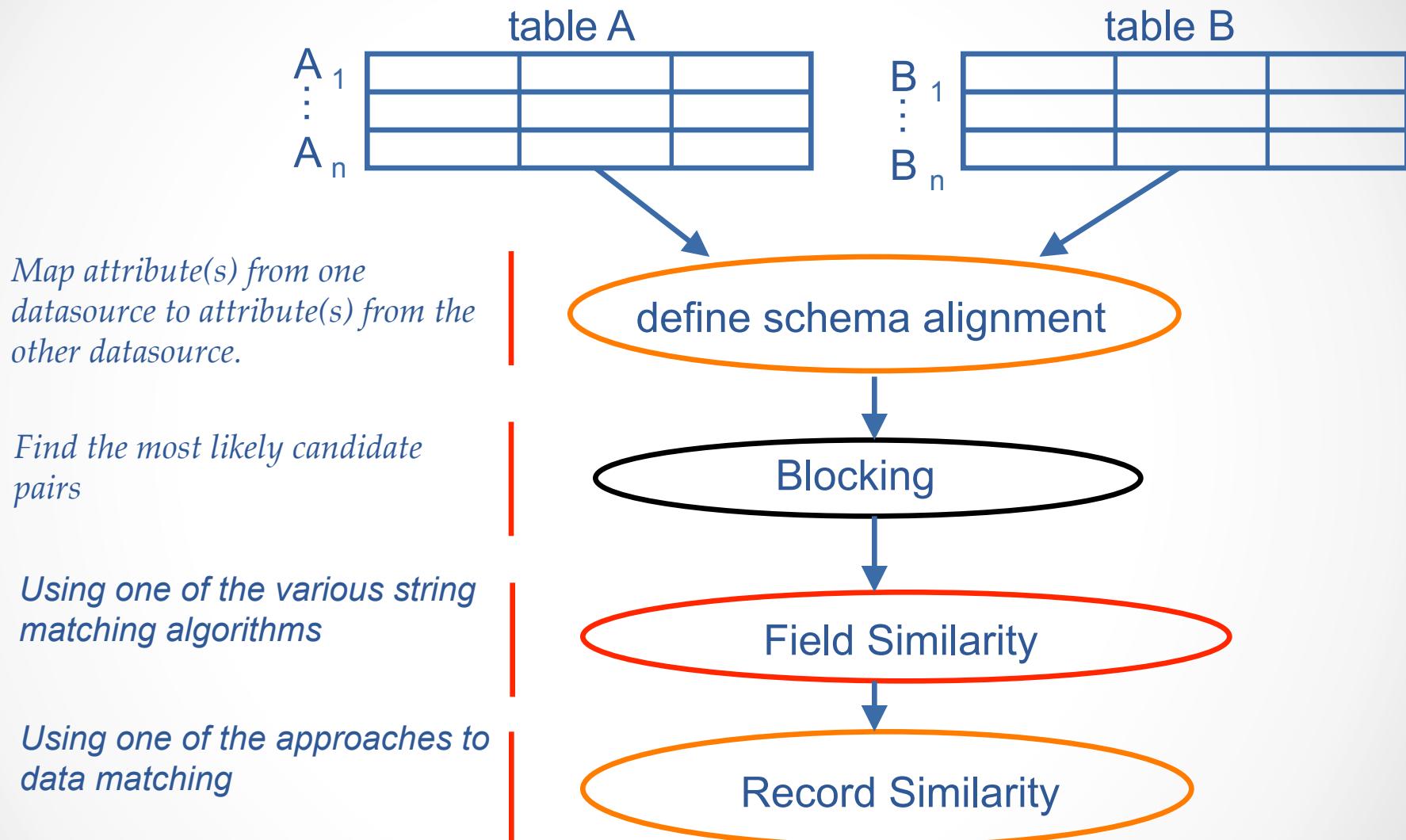
Challenges

- Same as in string matching
- How to match accurately?
 - difficult due to variations in formatting conventions, use of abbreviations, shortening, different naming conventions, omissions, nicknames, and errors in data
 - several common approaches: rule-based, learning-based, clustering, probabilistic, collective
- How to scale up to large data sets
 - again many approaches have been developed, as we will discuss

General Approach to Record Linkage

1. Identification of candidate pairs (blocking)
 - Comparing all possible record pairs would be computationally wasteful
2. Compute Field Similarity
 - String similarity between individual fields is computed
3. Compute Record Similarity
 - Field similarities are combined into a total record similarity estimate

Overview



Rule-based Matching

...

Rule-based Matching

- The developer writes rules that specify when two tuples match
 - typically after examining many matching and non-matching tuple pairs, using a development set of tuple pairs
 - rules are then tested and refined, using the same development set or a test set
- Many types of rules exist, we will consider
 - linearly weighted combination of individual similarity scores
 - logistic regression combination
 - more complex rules

Linearly Weighted Combination Rules

- Compute the sim score between tuples x and y as a linearly weighted combination of individual sim scores
 - $\text{sim}(x,y) = \sum_{i=1}^n \alpha_i * \text{sim}_i(x, y)$
 - n is number of attributes in each table
 - $\text{sim}_i(x,y)$ is a sim score between the i -th attributes of x and y
 - $\alpha_i \in [0,1]$ is a pre-specified weight that indicates the importance of the i -th attribute to $\text{sim}(x,y)$, such that $\sum_{i=1}^n \alpha_i = 1$
- We declare x and y matched if $\text{sim}(x,y) \geq \beta$ for a pre-specified β , and not matched otherwise
 - in another variation: declare x and y matched if $\text{sim}(x,y) \geq \beta$, not matched if $\text{sim}(x,y) \leq \gamma$, and subject to human review otherwise; we do not consider this in the chapter

Example

Table X

	Name	Phone	City	State
X ₁	Dave Smith	(608) 395 9462	Madison	WI
X ₂	Joe Wilson	(408) 123 4265	San Jose	CA
X ₃	Dan Smith	(608) 256 1212	Middleton	WI

(a)

Table Y

	Name	Phone	City	State
Y ₁	David D. Smith	395 9426	Madison	WI
Y ₂	Daniel W. Smith	256 1212	Madison	WI

(b)

Matches
 (x_1, y_1)
 (x_3, y_2)

(c)

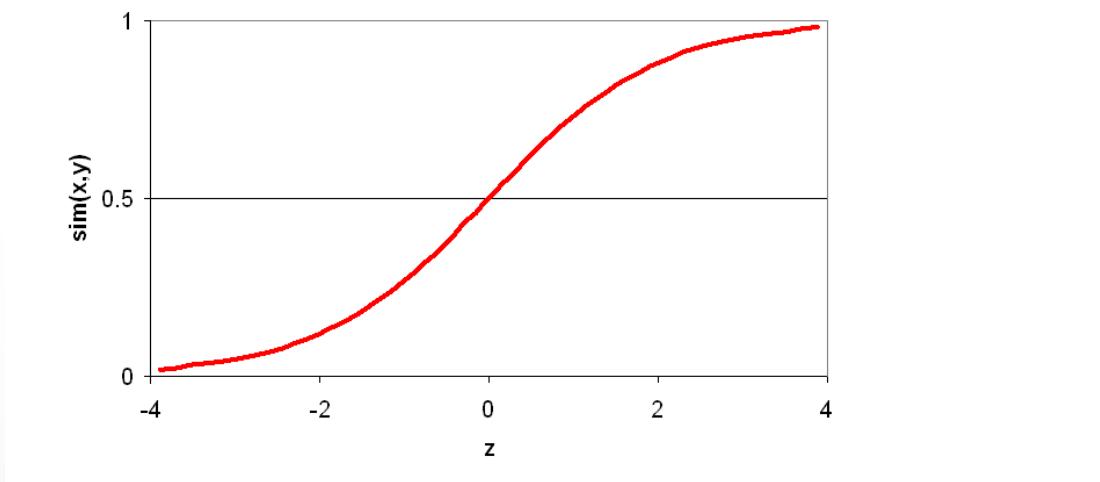
- $\text{sim}(x,y) = 0.3s_{\text{name}}(x,y) + 0.3s_{\text{phone}}(x,y) + 0.1s_{\text{city}}(x,y) + 0.3s_{\text{state}}(x,y)$
 - $s_{\text{name}}(x,y)$: based on Jaro-Winkler
 - $s_{\text{phone}}(x,y)$: based on edit distance between x's phone (after removing area code) and y's phone
 - $s_{\text{city}}(x,y)$: based on edit distance
 - $s_{\text{state}}(x,y)$: based on exact match; yes → 1, no → 0

Pros and Cons

- Pros
 - conceptually simple, easy to implement
 - can learn weights α_i from training data
- Cons
 - an increase \pm in the value of any s_i will cause a linear increase $\alpha_i * \pm$ in the value of s
 - in certain scenarios this is not desirable, thereafter a certain threshold an increase in s_i should count less (i.e., “diminishing returns” should kick in)
 - e.g., if $s_{name}(x,y)$ is already 0.95 then the two names already very closely match so any increase in $s_{name}(x,y)$ should contribute only minimally

Logistic Regression Rules

- address the diminishing-returns problem
- $\text{sim}(x,y) = 1 / (1 + e^{-z})$, where $z = \sum_{i=1}^n \alpha_i * \text{sim}_i(x, y)$
 - here α_i are not constrained to be in $[0,1]$ and sum to 1
 - so z goes from $-\infty$ to $+\infty$, in which case $\text{sim}(x,y)$ gradually increases, but minimally so after z has exceeded a certain value
→ ensuring diminishing returns



Logistic Regression Rules

- Are also very useful in situations where
 - there are many “signals” (e.g., 10-20) that can contribute to whether two tuples match
 - we don’t need all of these signals to “fire” in order to conclude that the tuples match
 - as long as a reasonable number of them fire, we have sufficient confidence
- Logistic regression is a natural fit for such cases
- Hence is quite popular as a first matching method to try

More Complex Rules

- Appropriate when we want to encode more complex matching knowledge
 - e.g., two persons match if names match approximately and either phones match exactly or addresses match exactly
 1. If $s_{name}(x,y) < 0.8$ then return “not matched”
 2. Otherwise if $e_{phone}(x,y) = \text{true}$ then return “matched”
 3. Otherwise if $e_{city}(x,y) = \text{true}$ and $e_{state}(x,y) = \text{true}$ then return “matched”
 4. Otherwise return “not matched”

Pros and Cons of Rule-Based Approaches

- Pros
 - easy to start, conceptually relatively easy to understand, implement, debug
 - typically run fast
 - can encode complex matching knowledge
- Cons
 - can be labor intensive, it takes a lot of time to write good rules
 - can be difficult to set appropriate weights
 - in certain cases it is not even clear how to write rules
 - learning-based approaches address these issues

Learning-based Matching

...

Learning-based Matching

- Here we consider supervised learning
 - learn a matching model M from training data, then apply M to match new tuple pairs
 - will consider unsupervised learning later

Learning-based Matching

- The training:
 1. start with training data: $T = \{(x_1, y_1, l_1), \dots (x_n, y_n, l_n)\}$, where each (x_i, y_i) is a tuple pair and l_i is a label: “yes” if x_i matches y_i and “no” otherwise

Learning-based Matching

- The training:
 1. start with training data: $T = \{(x_1, y_1, l_1), \dots (x_n, y_n, l_n)\}$, where each (x_i, y_i) is a tuple pair and l_i is a label: “yes” if x_i matches y_i and “no” otherwise
 2. define a set of features f_1, \dots, f_m , each quantifying one aspect of the domain judged possibly relevant to matching the tuples

Learning-based Matching

- The training:
 1. start with training data: $T = \{(x_1, y_1, l_1), \dots (x_n, y_n, l_n)\}$, where each (x_i, y_i) is a tuple pair and l_i is a label: “yes” if x_i matches y_i and “no” otherwise
 2. define a set of features f_1, \dots, f_m , each quantifying one aspect of the domain judged possibly relevant to matching the tuples
 3. convert each training example (x_i, y_i, l_i) in T to a pair $(\langle f_1(x_i, y_i), \dots, f_m(x_i, y_i) \rangle, c_i)$
 1. $v_i = \langle f_1(x_i, y_i), \dots, f_m(x_i, y_i) \rangle$ is a feature vector that encodes (x_i, y_i) in terms of the features
 2. c_i is an appropriately transformed version of label l_i (e.g., yes/no or 1/0, depending on what matching model we want to learn)

Learning-based Matching

- The training:
 1. start with training data: $T = \{(x_1, y_1, l_1), \dots, (x_n, y_n, l_n)\}$, where each (x_i, y_i) is a tuple pair and l_i is a label: “yes” if x_i matches y_i and “no” otherwise
 2. define a set of features f_1, \dots, f_m , each quantifying one aspect of the domain judged possibly relevant to matching the tuples
 3. convert each training example (x_i, y_i, l_i) in T to a pair $(\langle f_1(x_i, y_i), \dots, f_m(x_i, y_i) \rangle, c_i)$
 1. $v_i = \langle f_1(x_i, y_i), \dots, f_m(x_i, y_i) \rangle$ is a feature vector that encodes (x_i, y_i) in terms of the features
 2. c_i is an appropriately transformed version of label l_i (e.g., yes/no or 1/0, depending on what matching model we want to learn)
 4. thus T is transformed into $T' = \{(v_1, c_1), \dots, (v_n, c_n)\}$

Learning-based Matching

- The training:

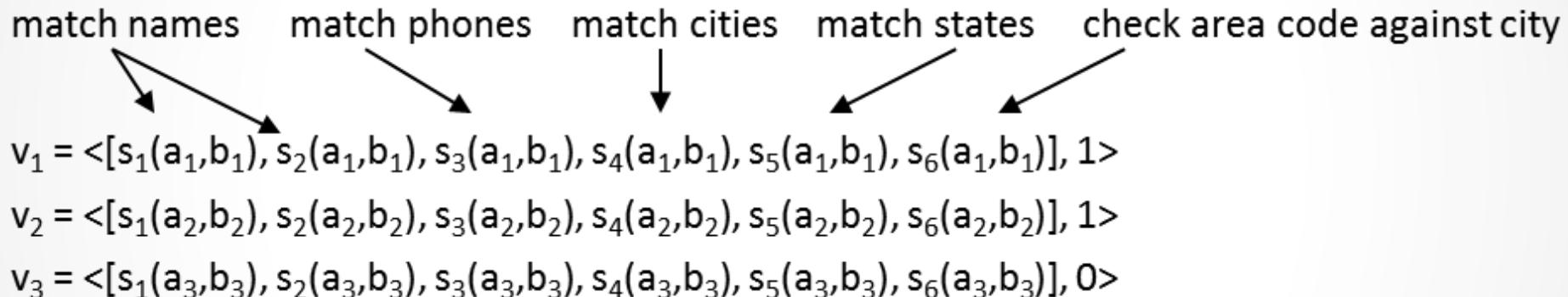
1. start with training data: $T = \{(x_1, y_1, l_1), \dots (x_n, y_n, l_n)\}$, where each (x_i, y_i) is a tuple pair and l_i is a label: “yes” if x_i matches y_i and “no” otherwise
2. define a set of features f_1, \dots, f_m , each quantifying one aspect of the domain judged possibly relevant to matching the tuples
3. convert each training example (x_i, y_i, l_i) in T to a pair $(\langle f_1(x_i, y_i), \dots, f_m(x_i, y_i) \rangle, c_i)$
 1. $v_i = \langle f_1(x_i, y_i), \dots, f_m(x_i, y_i) \rangle$ is a feature vector that encodes (x_i, y_i) in terms of the features
 2. c_i is an appropriately transformed version of label l_i (e.g., yes/no or 1/0, depending on what matching model we want to learn)
4. thus T is transformed into $T' = \{(v_1, c_1), \dots, (v_n, c_n)\}$
5. apply a learning algorithm (e.g. decision trees, SVMs) to T' to learn a matching model M

Learning-based Matching

- Applying model M to match new tuple pairs
 - given pair (x,y) , transform it into a feature vector
$$v = \langle f_1(x,y), \dots, f_m(x,y) \rangle$$
 - apply M to v to predict whether x matches y

Example: Learning a Linearly Weighted Rule

< a_1 = (Mike Williams, (425) 247 4893, Seattle, WA), b_1 = (M. Williams, 247 4893, Redmond, WA), yes>
< a_2 = (Richard Pike, (414) 256 1257, Milwaukee, WI), b_2 = (R. Pike, 256 1237, Milwaukee, WI), yes>
< a_3 = (Jane McCain, (206) 111 4215, Renton, WA), b_3 = (J. M. McCain, 112 5200, Renton, WA), no>



- s_1 and s_2 use Jaro-Winkler and edit distance
- s_3 uses edit distance (ignoring area code of a)
- s_4 and s_5 return 1 if exact match, 0 otherwise
- s_6 encodes a heuristic constraint

Example: Learning a Linearly Weighted Rule

- Goal: learn rule $s(a,b) = \sum_{i=1}^6 \alpha_i s_i(a, b)$
- Perform a least-squares linear regression on training data

$$v_1 = \langle [s_1(a_1, b_1), s_2(a_1, b_1), s_3(a_1, b_1), s_4(a_1, b_1), s_5(a_1, b_1), s_6(a_1, b_1)], 1 \rangle$$

$$v_2 = \langle [s_1(a_2, b_2), s_2(a_2, b_2), s_3(a_2, b_2), s_4(a_2, b_2), s_5(a_2, b_2), s_6(a_2, b_2)], 1 \rangle$$

$$v_3 = \langle [s_1(a_3, b_3), s_2(a_3, b_3), s_3(a_3, b_3), s_4(a_3, b_3), s_5(a_3, b_3), s_6(a_3, b_3)], 0 \rangle$$

to find weights α_i that minimizes the squared error

$$\sum_{i=1}^3 (c_i - \sum_{j=1}^6 \alpha_j s_j(v_i))^2$$

- c_i is the label associated with v_i
- $s_j(v_i)$ is the value of the j-th element of v_i

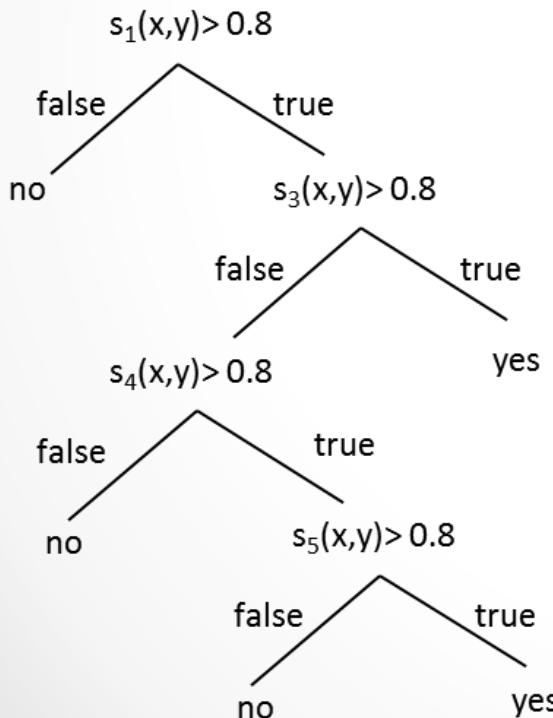
Example: Learning a Decision Tree

match names match phones match cities match states check area code against city

$v_1 = \langle [s_1(a_1, b_1), s_2(a_1, b_1), s_3(a_1, b_1), s_4(a_1, b_1), s_5(a_1, b_1), s_6(a_1, b_1)], \text{yes} \rangle$

$v_2 = \langle [s_1(a_2, b_2), s_2(a_2, b_2), s_3(a_2, b_2), s_4(a_2, b_2), s_5(a_2, b_2), s_6(a_2, b_2)], \text{yes} \rangle$

$v_3 = \langle [s_1(a_3, b_3), s_2(a_3, b_3), s_3(a_3, b_3), s_4(a_3, b_3), s_5(a_3, b_3), s_6(a_3, b_3)], \text{no} \rangle$



Now the labels are yes/no,
not 1/0

The Pros and Cons of Learning-based Approach

- Pros compared to rule-based approaches
 - in rule-based approaches must manually decide if a particular feature is useful → labor intensive and limit the number of features we can consider
 - learning-based ones can automatically examine a large number of features
 - learning-based approaches can construct very complex rules
- Cons
 - still require training examples, in many cases a large number of them, which can be hard to obtain
 - clustering can address this problem

Matching by Clustering

...

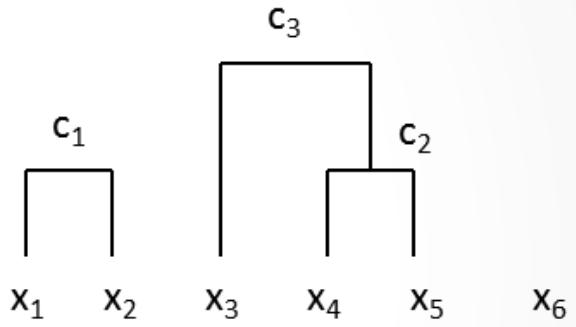
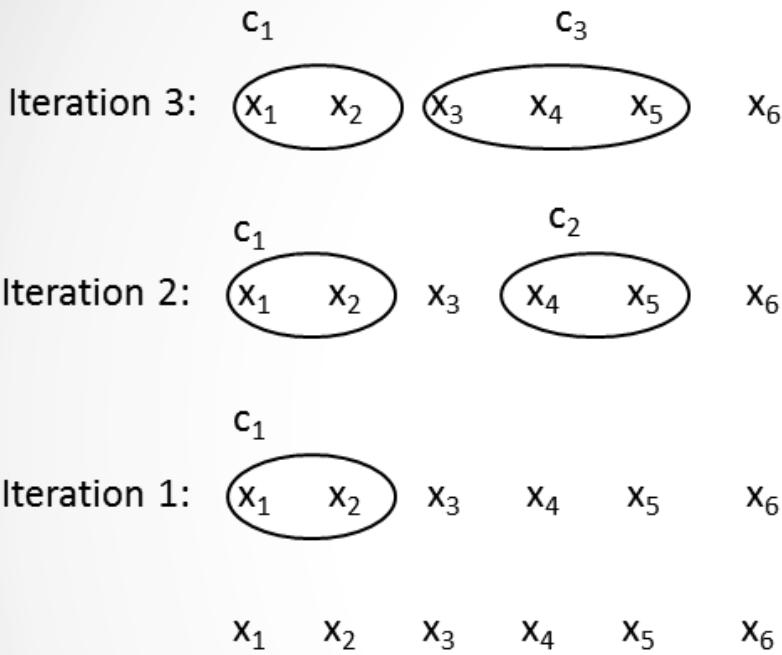
Matching by Clustering

- Many common clustering techniques have been used
 - agglomerative hierarchical clustering (AHC), k-means, graph-theoretic, ...
 - here we focus on AHC, a simple yet very commonly used one

Matching by Clustering

- AHC
 - partitions a given set of tuples D into a set of clusters
all tuples in a cluster refer to the same real-world entity,
tuples in different clusters refer to different entities
 - begins by putting each tuple in D into a single cluster
 - iteratively merges the two most similar clusters
 - stops when a desired number of clusters has been reached, or until the similarity between two closest clusters falls below a pre-specified threshold

Example



- $\text{sim}(x,y) = 0.3s_{\text{name}}(x,y) + 0.3s_{\text{phone}}(x,y) + 0.1s_{\text{city}}(x,y) + 0.3s_{\text{state}}(x,y)$

Computing a Similarity Score between Two Clusters

- Let c and d be two clusters

- Single link:

$$s(c,d) = \min_{x_i \in c, y_j \in d} \text{sim}(x_i, y_j)$$

- Complete link:

$$s(c,d) = \max_{x_i \in c, y_j \in d} \text{sim}(x_i, y_j)$$

- Average link:

$$s(c,d) = [\sum_{x_i \in c, y_j \in d} \text{sim}(x_i, y_j)] / [\# \text{ of } (x_i, y_j) \text{ pairs}]$$

Computing a Similarity Score between Two Clusters

- Canonical tuple:
 - create a canonical tuple that represents each cluster
 - sim between c and d is the sim between their canonical tuples
 - canonical tuple is created from attribute values of the tuples
 - e.g.,
“Mike Williams” & “M. J. Williams” → “Mike J. Williams”
(425) 247 4893 & 247 4893 → (425) 247 4893

Key Ideas underlying the Clustering Approach

- View matching tuples as the problem of constructing entities (i.e., clusters)
- The process is iterative
 - leverage what we have known so far to build better entities
- In each iteration merge all matching tuples within a cluster to build an “entity profile”, then use it to match other tuples → merging then exploiting the merged information to help matching

Scaling Up Data Matching

...

Scaling up Data Matching

- Two goals:
 - minimize # of tuple pairs to be matched
(Blocking)
 - minimize time it takes to match each pair

Blocking: Motivation

- Naïve pairwise: $|R|^2$ pairwise comparisons
 - 1000 business listings each from 1,000 different cities across the world
 - 1 trillion comparisons
 - 11.6 days (if each comparison is 1 μ s)
- Mentions from different cities are unlikely to be matches
 - Blocking Criterion: City
 - 1 billion comparisons
 - 16 minutes (if each comparison is 1 μ s)

Blocking: Motivation

- Mentions from different cities are unlikely to be matches
 - May miss potential matches

Get directions My places

www.bankofamerica.com/

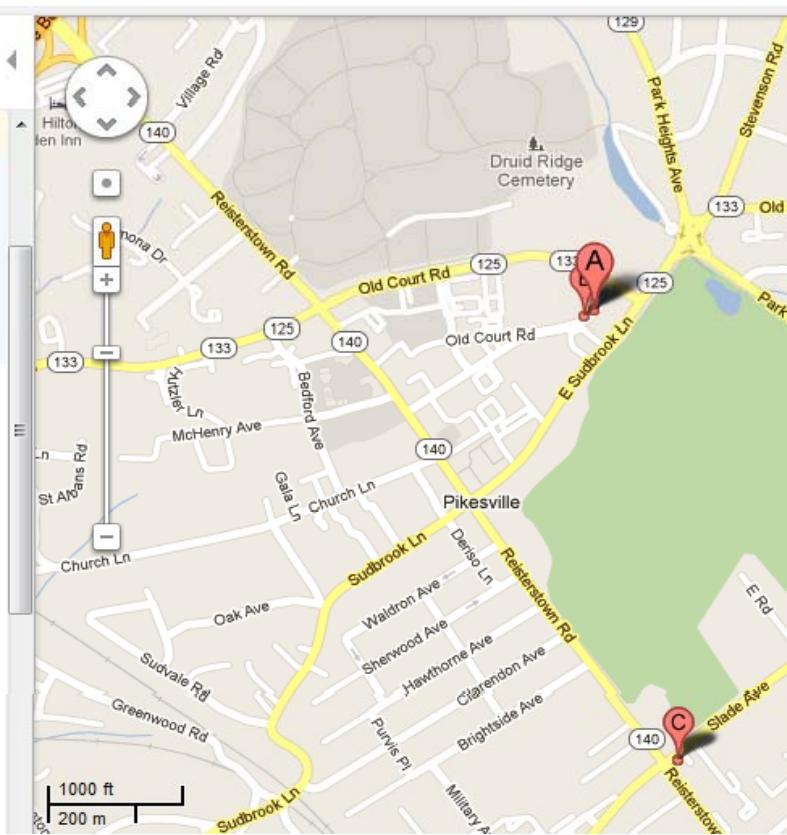
A Bank of America
3621 Old Court Road, Baltimore, MD
(410) 484-8511 ·
[locators.bankofamerica.com](#)
"Bank Of America is my husband's bank. Absolutely no problems ever and has ..." - [insiderpages.com](#)

B Bank of America
108 Old Court Road, Pikesville, MD
(410) 484-4301 ·
[locators.bankofamerica.com](#)

C Bank of America
25 Slade Avenue, Pikesville, MD
(410) 653-6482 ·
[locators.bankofamerica.com](#)

D ATM (Bank of America)
3621 Old Court Road, Baltimore, MD
M&T Bank

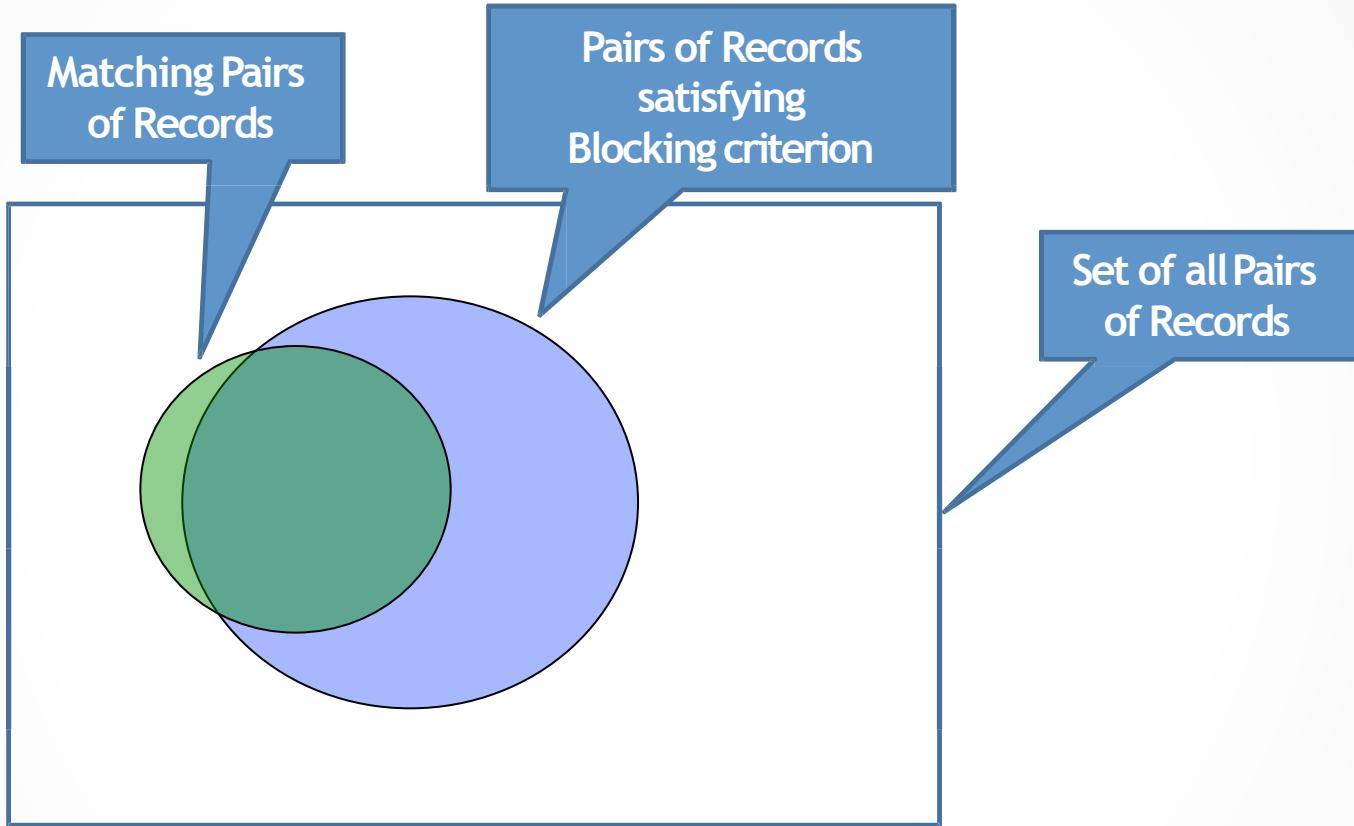
E ATM (Bank of America)
25 Slade Avenue, Pikesville, MD
(410) 653-6482 ·
[locators.bankofamerica.com](#)



The map displays a geographic area with several roads and landmarks. Key features include:

- Old Court Rd:** A major road running through the area, with segments labeled "Old Court Rd" and "Reisterstown Rd".
- Slade Ave:** A street running parallel to Old Court Rd.
- Reisterstown Rd:** A road connecting the northern and southern parts of the area.
- Sudbrook Ln:** A road branching off Old Court Rd.
- Village Rd:** A road in the northern part of the area.
- Druid Ridge Cemetery:** Located in the northern part of the map.
- Hilton Garden Inn:** Located near the intersection of Old Court Rd and Reisterstown Rd.
- Park Heights Ave:** A road in the eastern part of the map.
- Stevenson Rd:** A road in the far east.
- Local Streets:** Numerous smaller streets like Bedford Ave, McHenry Ave, Gaia Ln, S Church Ln, and Oak Ave.
- Landmarks:** The map includes a "Druid Ridge Cemetery" and a "Hilton Garden Inn".
- Scale:** A scale bar at the bottom left indicates distances of 1000 ft and 200 m.

Blocking: Motivation



Scaling up Data Matching

- Blocking:
 - hashing
 - sorting
 - indexing
 - canopies
 - using representatives
 - combining the techniques

Scaling up Data Matching

■ Simple Hashing

- hash tuples into buckets, match only tuples within each bucket
- e.g., hash house listings by zipcode, then match within each zip

Scaling up Data Matching

- Sorting
 - use a key to sort tuples, then scan the sorted list and match each tuple with only the previous **(w-1)** tuples, where w is a pre-specified window size
 - key should be strongly “discriminative”: brings together tuples that are likely to match, and pushes apart tuples that are not
 - example keys: soc sec, student ID, last name, soundex value of last name
 - employs a stronger heuristic than hashing: also requires that tuples likely to match be within a window of size w

Scaling up Data Matching

- Indexing
 - index tuples such that given any tuple a , can use the index to quickly locate a relatively small set of tuples that are likely to match a
e.g., inverted index on names

Scaling up Data Matching

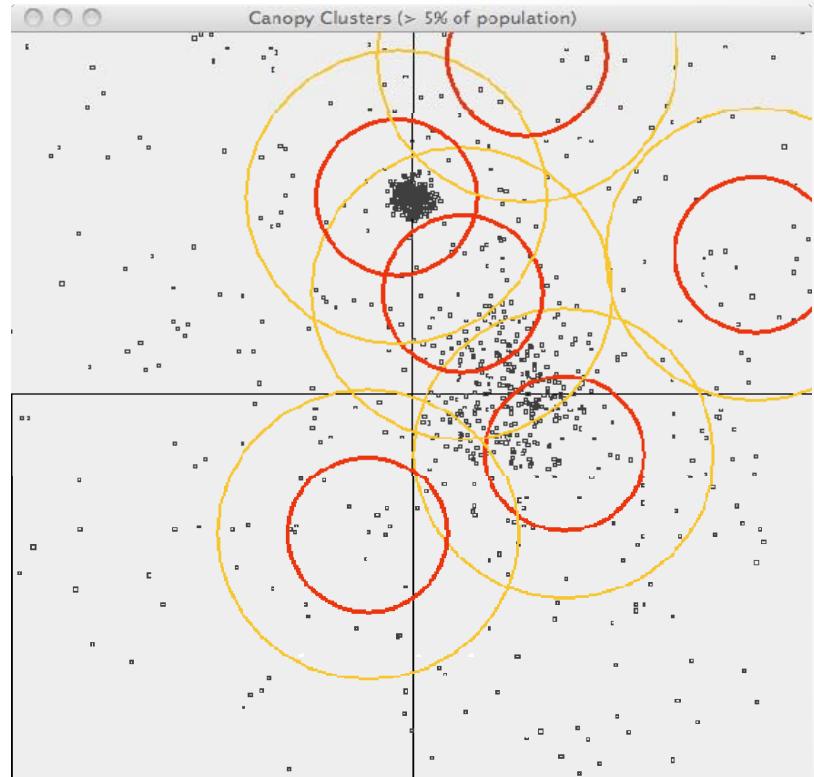
- Canopies
 - use a computationally cheap sim measure to quickly group tuples into overlapping clusters called canopies (or umbrella sets)
 - use a different (far more expensive) sim measure to match tuples within each canopy
 - e.g., Strings that differ by at least 3 in length can not have an edit distance smaller than 3.

Canopy Clustering [McCallum et al KDD'00]

Input: Mentions M ,
 $d(x,y)$, a distance metric,
thresholds $T_1 > T_2$

Algorithm:

1. Pick a random element x from M
2. Create new canopy C_x using
mentions y s.t. $d(x,y) < T_1$
3. Delete all mentions y from M
s.t. $d(x,y) < T_2$
4. Return to Step 1 if M is not empty.



Scaling up Data Matching

- Using representatives
 - applied during the matching process
 - assigns tuples that have been matched into groups such that those within a group match and those across groups do not
 - create a representative for each group by selecting a tuple in the group or by merging tuples in the group
 - when considering a new tuple, only match it with the representatives

Scaling up Data Matching

- Combining the techniques
 - e.g., hash houses into buckets using zip codes, then sort houses within each bucket using street names, then match them using a sliding window

Multi-Pass Blocking: Example

Census Data

<i>First Name</i>	<i>Last Name</i>	<i>Phone</i>	<i>Zip</i>
Matt	Michelson	555-5555	12345
Jane	Jones	555-1111	12345
Joe	Smith	555-0011	12345

A.I. Researchers

<i>First Name</i>	<i>Last Name</i>	<i>Phone</i>	<i>Zip</i>
Matthew	Michelson	555-5555	12345
Jim	Jones	555-1111	12345
Joe	Smeth	555-0011	12345

match

match

Rule-based Blocking: Generating Candidates

(method, attribute)



(token, last name) AND (1st letter, first name) = block-key

<i>First Name</i>	<i>Last Name</i>
Matt	Michelson
Jane	Jones

<i>First Name</i>	<i>Last Name</i>
Matthew	Michelson
Jim	Jones

(token, zip)

<i>First Name</i>	<i>Last Name</i>	<i>Zip</i>
Matt	Michelson	12345
Matt	Michelson	12345
Matt	Michelson	12345

<i>First Name</i>	<i>Last Name</i>	<i>Zip</i>
Matthew	Michelson	12345
Jim	Jones	12345
Joe	Smeth	12345

Blocking – Multi-pass

(Hernandez & Stolfo, 1998)

- Sort neighborhoods on block keys
- Multiple independent runs using keys
 - runs capture different match candidates
- Example:
 - 1st → (token, last name)
 - 2nd → (token, first name) & (token, phone)

Blocking – Multi-pass

- Terminology:
 - Each pass is a “conjunction”
 - (token, first) AND (token, phone)
 - Combine passes to form “disjunction”
 - [(token, last)] OR
 - [(token, first) AND (token, phone)]
 - Disjunctive Normal Form rules
 - form “Blocking Schemes”

Blocking Effectiveness

- Determined by rules
- Determined by choices for attributes and methods
 - (token, zip) captures all matches, but all pairs too
 - (token, first) AND (token, phone) gets half the matches, and only 1 candidate generated
 - Which is better? Why?
- How to quantify??

Blocking Effectiveness

$$\text{Reduction Ratio (RR)} = 1 - \frac{\|C\|}{(\|S\| * \|T\|)}$$

S,T are data sets;

C is the set of candidates

$$\text{Pairs Completeness (PC) [Recall]} = S_m / N_m$$

S_m = # true matches in candidates,

N_m = # true matches between S and T

Examples:

(token, last name) AND (1st letter, first name)

$$\text{RR} = 1 - 2/9 \approx 0.78$$

$$\text{PC} = 1 / 2 = 0.50$$

(token, zip)

$$\text{RR} = 1 - 9/9 = 0.0$$

$$\text{PC} = 2 / 2 = 1.0$$

Census Data			
First Name	Last Name	Phone	Zip
Matt	Michelson	555-5555	12345
Jane	Jones	555-1111	12345
Joe	Smith	555-0011	12345

A.I. Researchers			
First Name	Last Name	Phone	Zip
Matthew	Michelson	555-5555	12345
Jim	Jones	555-1111	12345
Joe	Smeth	555-0011	12345

match match

How to choose methods and attributes?

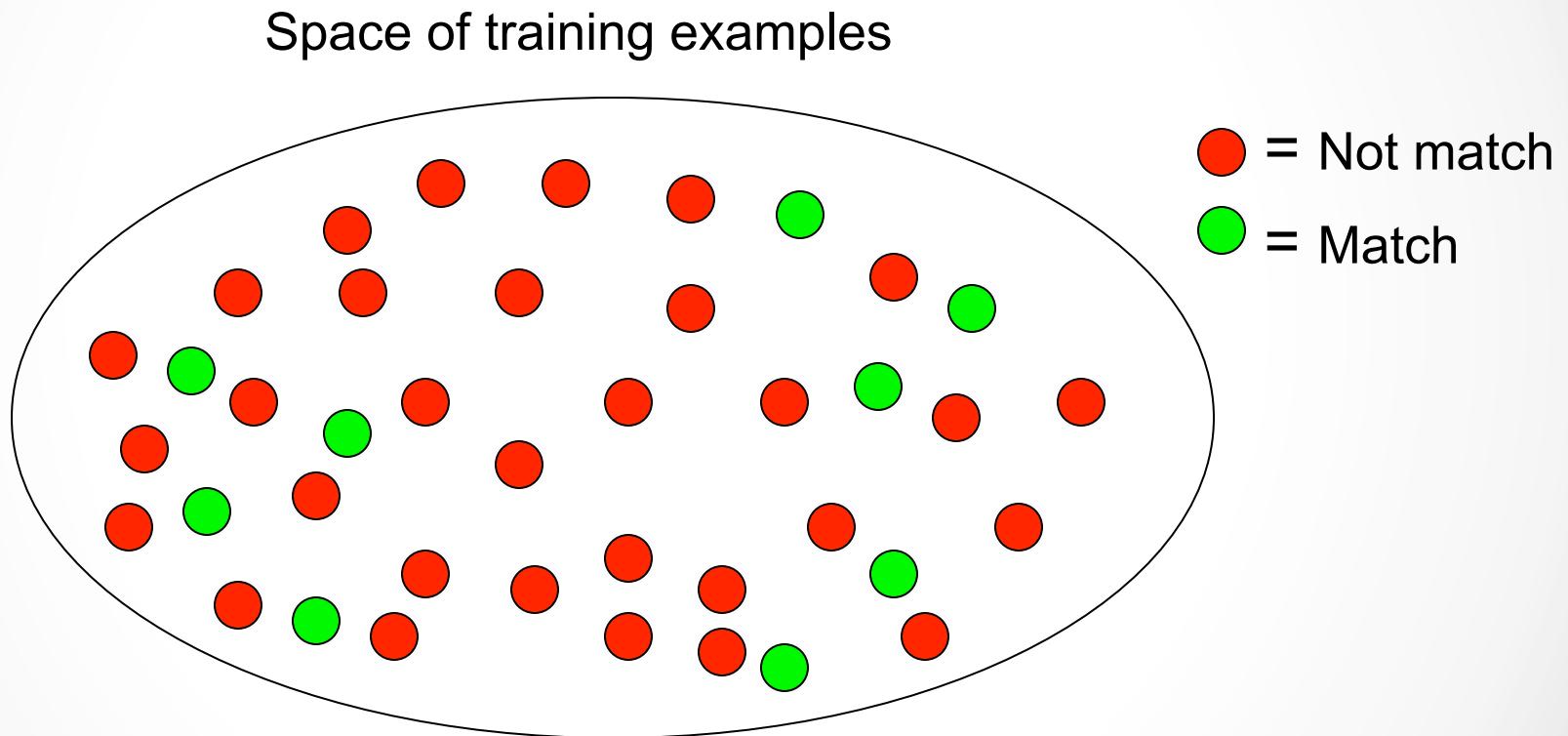
- Blocking Goals:
 - Small number of candidates (High **RR**)
 - Don't leave any true matches behind! (High **PC**)
- Previous approaches:
 - Ad-hoc by researchers or domain experts
- New Approach:
 - Blocking Scheme Learner (BSL) – modified Sequential Covering Algorithm

Michelson & Knoblock, Learning Blocking Schemes for Record Linkage, AAAI 2006

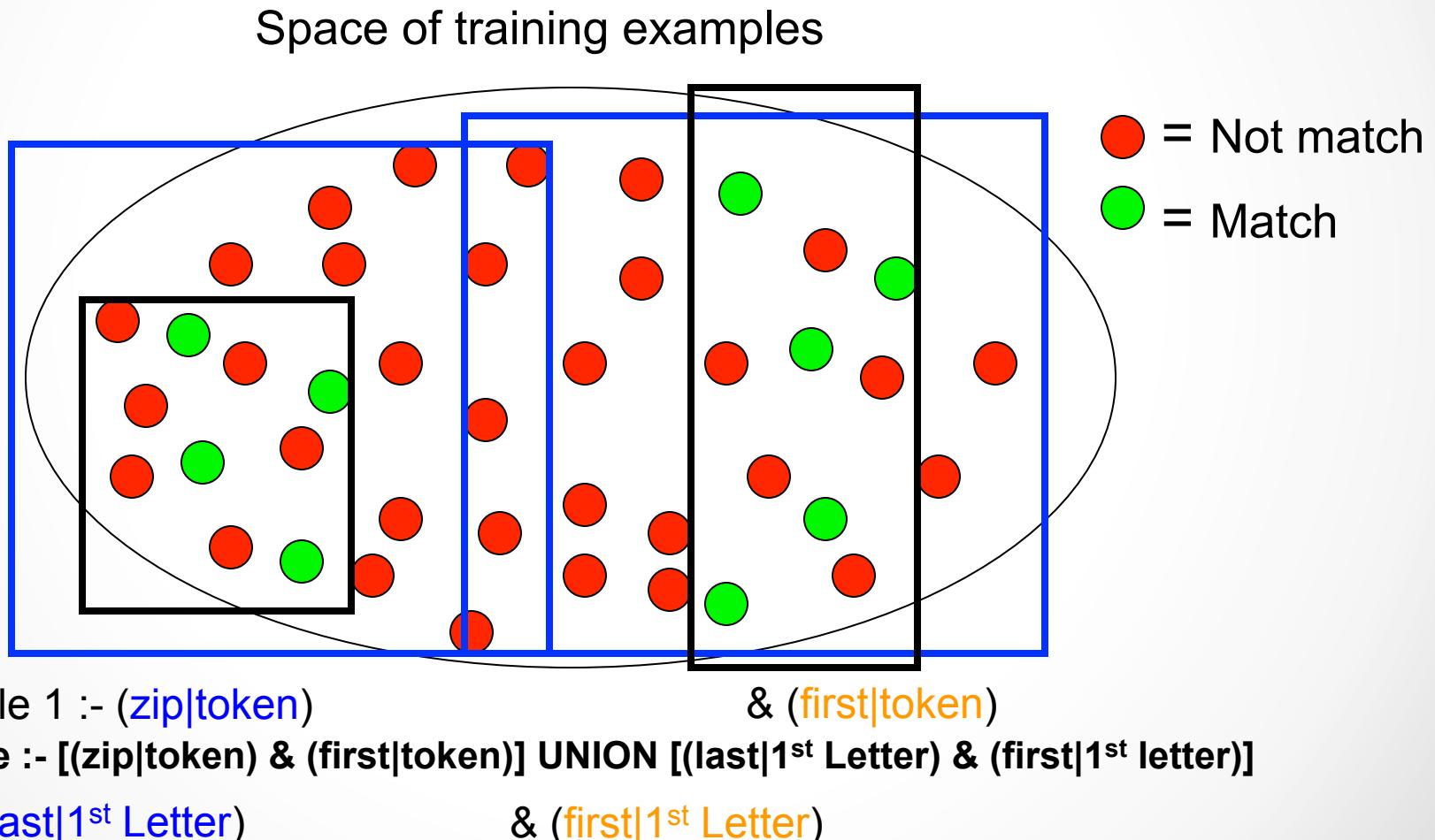
Learning Schemes – Intuition

- Learn restrictive conjunctions
 - partition the space → minimize False Positives
- Union restrictive conjunctions
 - Cover all training matches
 - Since minimized FPs, conjunctions should not contribute many FPs to the disjunction

Example to clear things up!



Example to clear things up!



SCA: propositional rules

- Multi-pass blocking = disjunction of conjunctions
- Learn conjunctions and union them together!
- Cover all training matches to maximize PC

```
SEQUENTIAL-COVERING( class, attributes, examples, threshold)
```

```
LearnedRules ← {}
```

```
Rule ← LEARN-ONE-RULE(class, attributes, examples)
```

```
While examples left to cover, do
```

```
    LearnedRules ← LearnedRules U Rule
```

```
    Examples ← Examples – {Examples covered by Rule}
```

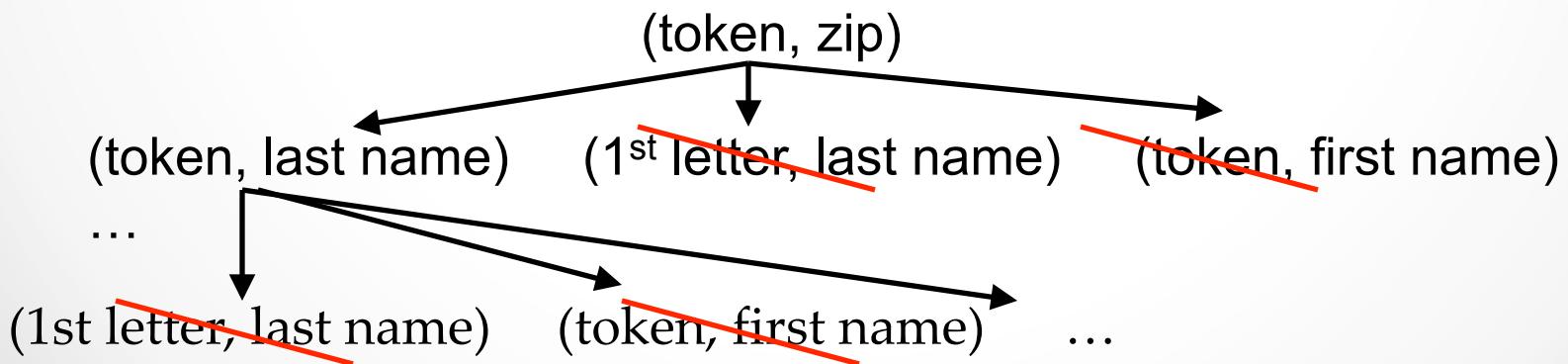
```
    Rule ← LEARN-ONE-RULE(class, attributes, examples)
```

```
    If Rule contains any previously learned rules, remove them
```

```
Return LearnedRules
```

Learn-One-Rule

- Learn conjunction that maximizes **RR**
- General-to-specific beam search
 - Keep adding/intersecting (attribute, method) pairs
 - Until can't improve **RR**
 - Must satisfy minimum **PC**



Scaling up Data Matching

- Two goals:
 - minimize # of tuple pairs to be matched
(Blocking)
 - minimize time it takes to match each pair

Scaling up Data Matching

- For the second goal of minimizing time it takes to match each pair
 - no well-established technique as yet
 - tailor depending on the application and the matching approach
 - e.g., if using a simple rule-based approach (i.e., matches individual attributes then combines their scores using weights)
Use **short circuiting**: stop the computation of the sim score if it is already so high that the tuple pair will match even if the remaining attributes do not match

Scaling up Using Parallel Processing

- Commonly done in practice
- Examples: hash tuples into buckets, then match each bucket in parallel