# Record Linkage

Everything Data

CompSci 290.01 Spring 2014

**DUKE**
COMPUTER SCIENCE

# Announcements (Wed. Jan. 28)

- **Homework #3** will be posted by tomorrow morning
  - Due midnight Sunday

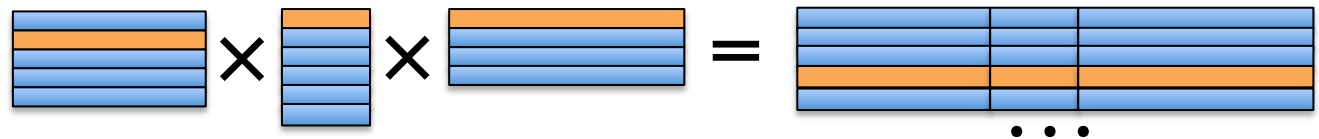# Recap: *Querying Relational Databases in SQL*

**SELECT** *columns or expressions*  5. Compute one output row for each "wide row"

(or for each group of them if query has grouping/aggregation)

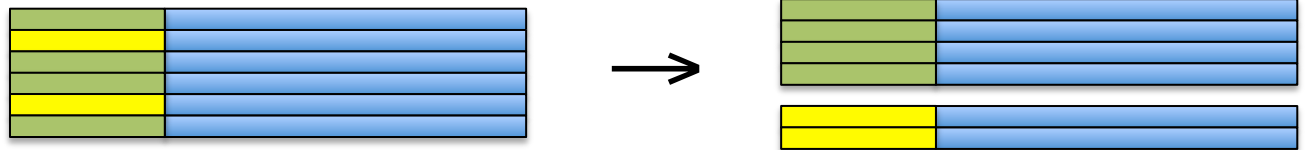**FROM** *tables*  1. Generate all combinations of rows, one from each table; each combination forms a "wide row"



**WHERE** *conditions*  2. Filter—keep only "wide rows" satisfying *conditions*

**GROUP BY** *columns*  3. Group—"wide rows" with matching values for *columns* go into the same group



**ORDER BY** *output columns*;  4. Sort the output rows

# Problem

- Forbes magazine article: "Wall Street's favorite senators"

# Problem

- Forbes magazine article: "Wall Street's favorite senators"

```
Chris,Dodd,Democrat,CT,35.7,9161489
Richard,Shelby,Republican,AL,33.4,2542878
Charles,Schumer,Democrat,NY,32.8,3255362
Tom,Carper,Democrat,DE,32.5,1453446
Mike,Crapo,Republican,ID,32.2,946531
Bob,Bennett,Republican,UT,32.3,1078302
Jack,Reed,Democrat,RI,31.5,1280500
Tim,Johnson,Democrat,SD,29.1,1396308
Mike,Enzi,Republican,WY,25.1,564100
Joe,Liebermen,Independent,CT,25,7878838
```

- What are their ages?

# Solution

- Join with the persons table (from govtrack)

- But there is no key to join on …

# Record Linkage

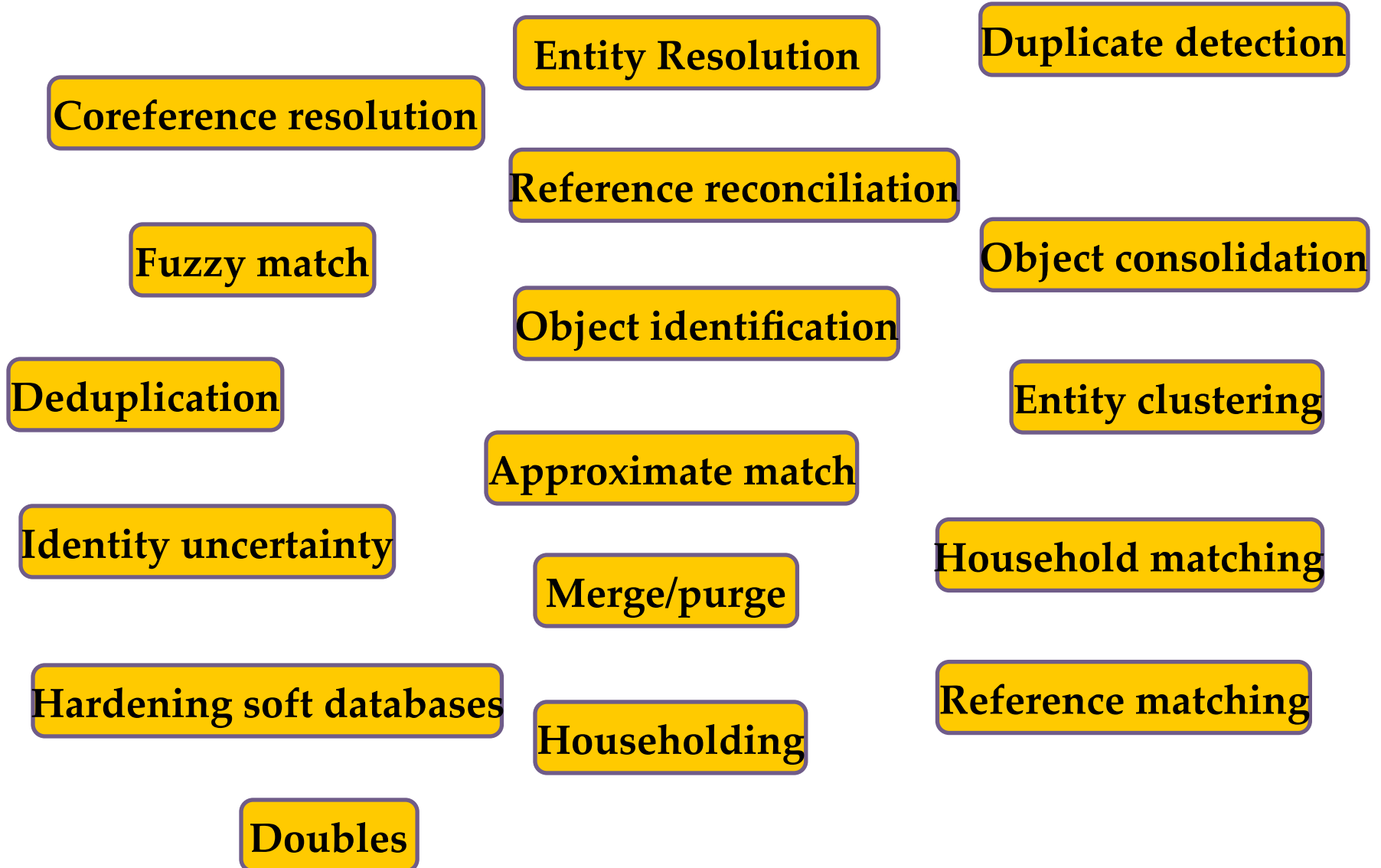- Problem of finding duplicate entities across different sources *(or even within a single dataset).*



Article | Talk

Read | Edit | View history | Search

## Record linkage

From Wikipedia, the free encyclopedia
(Redirected from Entity resolution)

**Record linkage** (RL) refers to the task of finding records in a data set that refer to the same entity across different data sources (e.g., data files, books, websites, databases). Record linkage is necessary when joining data sets based on entities that may or may not share a common identifier (e.g., database key, URI, National identification number), as may

# Ironically, Record Linkage has many names

**Entity Resolution**
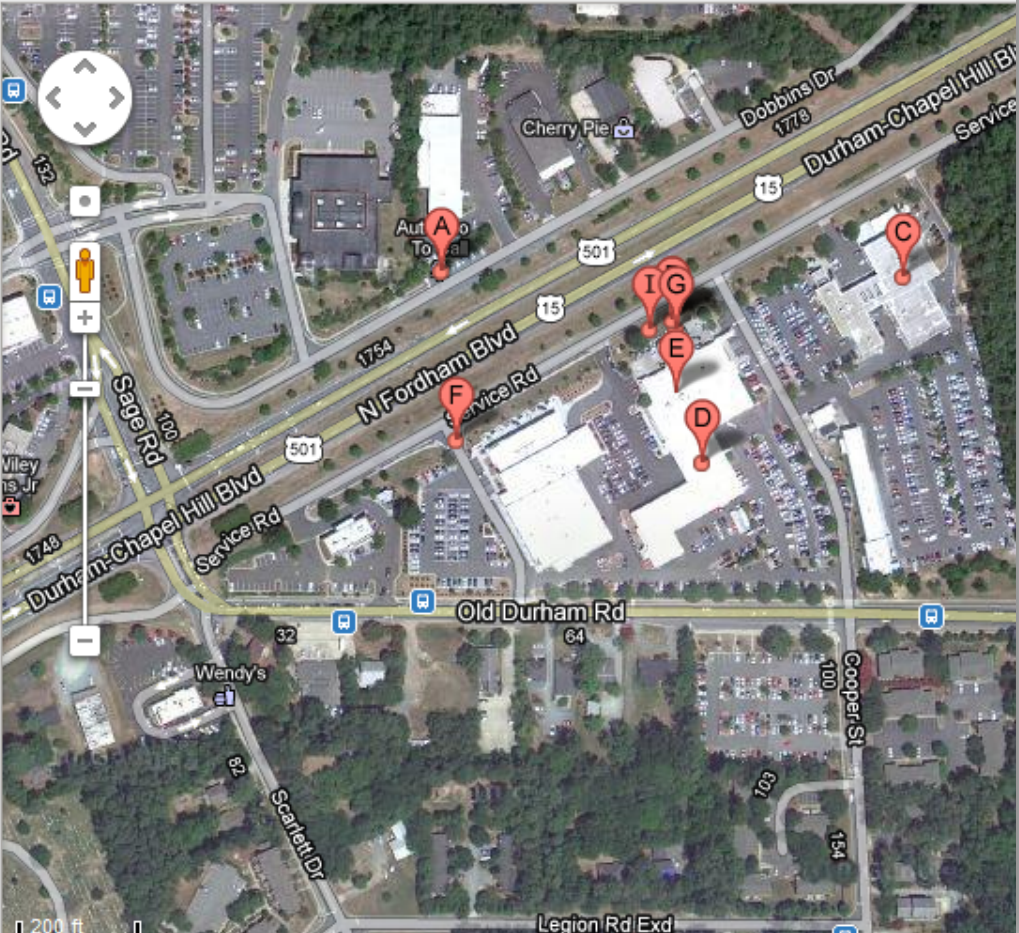
**Duplicate detection**

**Coreference resolution**

**Reference reconciliation**

**Fuzzy match**

**Object consolidation**

**Object identification**

**Deduplication**

**Entity clustering**

**Approximate match**

**Identity uncertainty**

**Household matching**

**Merge/purge**

**Hardening soft databases**

**Reference matching**

**Householding**

**Doubles**

# Motivating Example 1: Web

# Motivating Example 1: Web

# Motivating Example 1: Web

# Motivating Example 2: Network Science

- Measuring the topology of the internet … using `traceroute`

# IP Aliasing Problem  [Willinger et al. 2009]



(a)          (b)

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,

# And many many more examples

- *Linking Census Records*
- *Public Health*
- *Medical records*
- *Web search – query disambiguation*
- *Comparison shopping*
- *Maintaining customer databases*
- *Law enforcement and Counter-terrorism*
- *Scientific data*
- *Genealogical data*
- *Bibliographic data*

# Opportunity

# Back to our example

- Join with the persons table (from govtrack)

- But there is no key to join on …

- What about (firstname, lastname)?

# Attempt 1:

SELECT w.*, date_part('year', current_date) - date_part('year', p.birthday) AS age

FROM wallst w, persons p

WHERE w.first_name = p.first_name

and w.last_name = p.last_name;

# Problems

- Join condition is too specific
  - Nicknames used instead of real first names

# Attempt 2:

- Join on Last name + Age < 100 (senator must be alive)

SELECT w.*, date_part('year', current_date) - date_part('year', p.birthday) AS age

FROM wallst w, persons p

WHERE w.lastname = p.last_name and date_part('year', current_date) - date_part('year', p.birthday) < 100;

# Problem:

- Join condition is too inclusive
  - Many individuals share the same last name.

| Surname | Approx # | Rank |
|---------|----------|------|
| Smith | 2.4 M | 1 |
| Johnson | 1.8 M | 2 |
| Williams | 1.5 M | 3 |
| Brown | 1.4 M | 4 |
| Jones | 1.4 M | 5 |

http://www.census.gov/genealogy/www/data/2000surnames/

# "Where is Joe Liebermen ?"

- Spelling mistake
  - Liebermen vs Lieberman

- Need an approximate matching condition!

# Levenshtein (or edit) distance

- The minimum number of character **edit** operations needed to turn one string into the other.

LIEBERMAN

LIEBERMEN

– Substitute A to E. Edit distance = 1

# Levenshtein (or edit) distance

- Distance between two string s and t is the shortest sequence of **edit commands** that transform s to t.

- Commands:

  Costs can be different

  – Copy character from s to t        (cost = 0)
  – Delete a character from s        (cost = 1)
  – Insert a character into t        (cost = 1)
  – Substitute one character for another (cost = 1)

# Levenshtein (or edit) distance

Ashwin Machanavajjhala

Aswhin Maachanavajhala

# Levenshtein (or edit) distance

String s: Ashwin MaGchanavajjhala

sub        ins        del

String t: Aswhin MaachanavajGhala

Total cost: 4

# Computing the edit distance

|   |   | A | S | W | A | N |
|---|---|---|---|---|---|---|
|   | 0 | 1 |   |   |   |   |
| A | 1 |   |   |   |   |   |
| S |   |   |   |   |   |   |
| W |   |   |   |   |   |   |
| H |   |   |   |   |   |   |
| I |   |   |   |   |   |   |
| N |   |   |   |   |   |   |

Cost of changing "**G**" → "A"

Cost of changing "ASWH" → "AS"

# Computing the edit distance

|   |   | A | S | W | A | N |
|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 |   |   |   |
| A | 1 | 0 | 1 |   |   |   |
| S | 2 | 1 | 0 |   |   |   |
| W | 3 | 2 |   |   |   |   |
| H |   |   |   |   |   |   |
| I |   |   |   |   |   |   |
| N |   |   |   |   |   |   |

Cost of changing "ASW" → "AS":

Minimum of:
- Cost of "AS" → "AS" + 1 (delete W)
- Cost of "ASW" → "A" + 1 (insert S)
- Cost of "AS" → "A" + 1 (substitute W with S)

# Computing the edit distance

|   |   | A | S | W | A | N |
|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 |
| A | 1 | 0 | 1 | 2 | 3 | 4 |
| S | 2 | 1 | 0 | 1 | 2 | 3 |
| W | 3 | 2 | 1 | 0 | 1 | 2 |
| H | 4 | 3 | 2 | 1 | 1 | 2 |
| I | 5 | 4 | 3 | 2 | 2 | 2 |
| N | 6 | 5 | 4 | 3 | 3 | ? |

# Computing the edit distance

|   |   | A | S | W | A | N |
|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 |
| A | 1 | 0 | 1 | 2 | 3 | 4 |
| S | 2 | 1 | 0 | 1 | 2 | 3 |
| W | 3 | 2 | 1 | 0 | 1 | 2 |
| H | 4 | 3 | 2 | 1 | 1 | 2 |
| I | 5 | 4 | 3 | 2 | 2 | 2 |
| N | 6 | 5 | 4 | 3 | 3 | 2 |

Remember the minimum in each step and retrace your path.

# Edit Distance Variants

- Needleman-Munch
  - Different costs for each operation

- Affine Gap distance
  - John Reed vs John Francis "Jack" Reed
  - Consecutive inserts cost less than the first insert.

# Back to our example … Attempt 3

SELECT w.firstname, w.lastname, w.state, w.party, p.first_name, p.last_name, date_part('year', current_date) - date_part('year', p.birthday) AS age

FROM  wallst w, persons p

WHERE levenshtein(w.lastname, p.last_name) <= 1 and date_part('year', current_date) - date_part('year', p.birthday) < 100;

# Jaccard Distance

- Useful similarity function for sets
  - *(and for… long strings).*
- Let A and B be two sets
  - Words in two documents
  - Friends lists of two individuals

$$\text{Jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

# Jaccard similarity for names

- Use character trigrams

LIEBERMAN =  {GGL, GLI, LIE, IEB, EBE,
BER, ERM, RMA,MAN, ANG, NGG}
LIEBERMEN =  {GGL, GLI, LIE, IEB, EBE,
BER, ERM, RMA,MEN, ENG, NGG}

Jaccard(s,t) = 9/13 = 0.69

# Attempt 4:

SELECT w.firstname, w.lastname, w.state, w.party, p.first_name, p.last_name, date_part('year', current_date) - date_part('year', p.birthday) AS age

FROM  wallst w, persons p

WHERE similarity(w.lastname, p.last_name) >= 0.5 and date_part('year', current_date) - date_part('year', p.birthday) < 100;

# Translation / Substitution Tables

- Strings that are usually used interchangeably
  - New York vs Big Apple
  - Thomas vs Tom
  - Robert vs Bob

# Attempt 5

select w.firstname, w.lastname, w.state, p.first_name, p.last_name, date_part('year', current_date) - date_part('year', p.birthday) AS age

from wallst w, persons p

where levenshtein(w.lastname, p.last_name) <= 1 and date_part('year', current_date) - date_part('year', p.birthday) < 100

and (w.firstname = p.first_name or w.firstname IN (select n.nickname from nicknames n where n.firstname = p.first_name));

# Almost there …

- Tim matches both Timothy and Tim
  - Can fix it by matching on STATE
  - *Homework exercise* ☺

# Summary of Similarity Methods

**Easiest and most efficient**

- Equality on a boolean predicate
- Edit distance
  - Levenstein, Affine
- Set similarity
  - Jaccard
- Vector Based
  - Cosine similarity, TFIDF

- Translation-based
- Numeric distance between values
- Phonetic Similarity
  - Soundex, Metaphone
- Other
  - Jaro-Winkler, Soft-TFIDF, Monge-Elkan

# Summary of Similarity Methods

**Handle Typographical errors**

**Useful for abbreviations, alternate names.**

- Equality on a boolean predicate
- Edit distance
  - Levenstein, Affine
- Set similarity
  - Jaccard
- Vector Based
  - Cosine similarity, TFIDF

**Good for Text (reviews/ tweets), sets, class membership, …**

- Translation-based
- Numeric distance between values
- Phonetic Similarity
  - Soundex, Metaphone
- Other
  - Jaro-Winkler, Soft-TFIDF, Monge-Elkan

**Good for Names**

# Evaluating Record Linkage

- Hard to get all the matches to be exactly correct in real world problems
  - As we saw in real examples

- Need to quantify how good the matching is.

# Property Testing

- Consider a universe U of <span style="color:red">objects</span>
  - Documents (in web search)
  - Pairs of records (in record linkage)

- Suppose you want to identify a subset M in U that satisfies a specific <span style="color:red">property</span>
  - Relevance to a query (in web search)
  - Do the records match (in record linkage)

# Property Testing

- Consider a universe U of objects
- Suppose you want to identify a subset M in U that satisfies a specific property

- Let A be an (imperfect) algorithm that guesses whether or not an element in U satisfies the property
  - Let $M_A$ be the subset of objects that A identifies as satisfying the property.

# Property Testing

## Real World

Crying Wolf!

| Algorithm Guess | Satisfies P | Doesn't Satisfy P | |
|---|---|---|---|
| Satisfies P | *True positives (TP)* | *False positives (FP)* | $M_A$ |
| Doesn't satisfy P | *False negatives (FN)* | *True negatives (TN)* | $U - M_A$ |
| | M | U - M | |

# Venn diagram view

True negatives
(TN)

True positives
(TP)

U

M  M_A

False negatives
(FN)

False positives
(FP)

# Error: Precision / Recall

$$\text{Precision} = TP/(TP + FP)$$
$$= |M \cap M_A|/|M_A|$$

*fraction of answers returned by A that are correct*

$$\text{Recall} = TP/(TP + FN)$$
$$= |M \cap M_A|/|M|$$

*fraction of correct answers that are returned by A*

# Error: F-measure

$$\text{Precision} = \frac{|M \cap M_A|}{|M_A|}$$

$$\text{Recall} = \frac{|M \cap M_A|}{|M|}$$

$$\text{F1 score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

# Example

- M:

| firstname | lastname | state | first_name | last_name | age |
|-----------|----------|-------|------------|-----------|-----|
| Bob | Bennett | UT | Robert | Bennett | 81 |
| Tom | Carper | DE | Thomas | Carper | 67 |
| Mike | Crapo | ID | Michael | Crapo | 63 |
| Chris | Dodd | CT | Christopher | Dodd | 70 |
| Mike | Enzi | WY | Michael | Enzi | 70 |
| Tim | Johnson | SD | Tim | Johnson | 68 |
| Joe | Liebermen | CT | Joseph | Lieberman | 72 |
| Jack | Reed | RI | John | Reed | 65 |
| Charles | Schumer | NY | Charles | Schumer | 64 |
| Richard | Shelby | AL | Richard | Shelby | 80 |

(10 rows)

# Example:

Algorithm A:

select * from wallst w, persons p

where w.lastname = p.last_name and

date_part('year', current_date) - date_part('year', p.birthday) < 100

and (w.firstname = p.first_name or w.firstname IN (select n.nickname from nicknames n where n.firstname = p.first_name));

Exact match on last name

Age < 100

First name is same or a nickname

# Example

- $M_A$:

| firstname | lastname | state | first_name | last_name | age |
|-----------|----------|-------|------------|-----------|-----|
| Bob | Bennett | UT | Robert | Bennett | 81 |
| Charles | Schumer | NY | Charles | Schumer | 64 |
| Chris | Dodd | CT | Christopher | Dodd | 70 |
| Jack | Reed | RI | John | Reed | 65 |
| Mike | Crapo | ID | Michael | Crapo | 63 |
| Mike | Enzi | WY | Michael | Enzi | 70 |
| Richard | Shelby | AL | Richard | Shelby | 80 |
| Tim | Johnson | SD | Timothy | Johnson | 68 |
| Tim | Johnson | SD | Tim | Johnson | 68 |
| Tom | Carper | DE | Thomas | Carper | 67 |

(10 rows)

# Example

$$\text{Precision} = \frac{|M \cap M_A|}{|M_A|}$$

$$= \frac{9}{10} = 0.9$$

$$\text{Recall} = \frac{|M \cap M_A|}{|M|}$$

$$= \frac{9}{10} = 0.9$$

$$\text{F1 score} = 2 \, \frac{0.9 \times 0.9}{0.9 + 0.9} = 0.9$$

# Summary

- Many interesting data analyses require reasoning across different datasets

- May not have access to keys that uniquely identify individual rows in both datasets
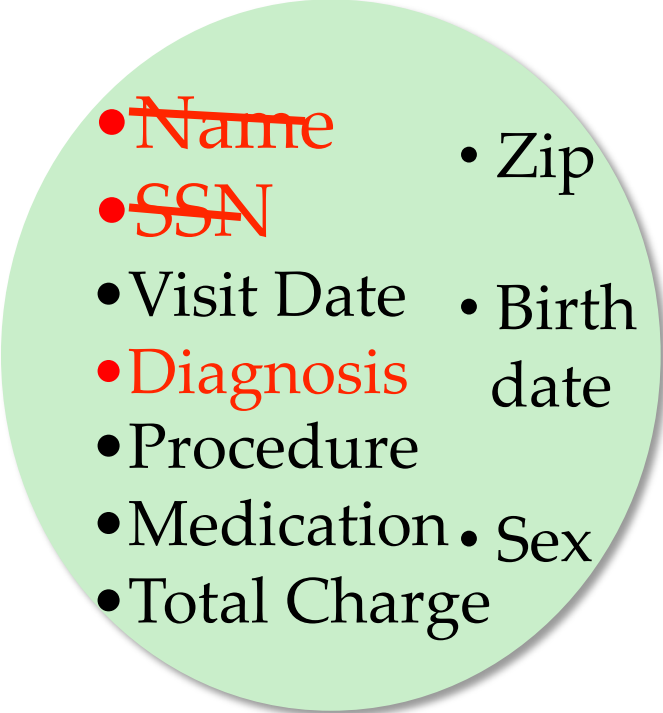
# Summary

- Use combinations of attributes that are approximate keys (or **quasi-identifiers**)

- Use similarity measures for fuzzy or approximate matching
  - **Levenshtein** or **Edit** distance
  - **Jaccard** Similarity

- Use translation tables

# Summary

- Record Linkage is rarely perfect
  - Missing attributes
  - Messy data errors
  - …

- **Precision/Recall** is used to measure the quality of linkage.
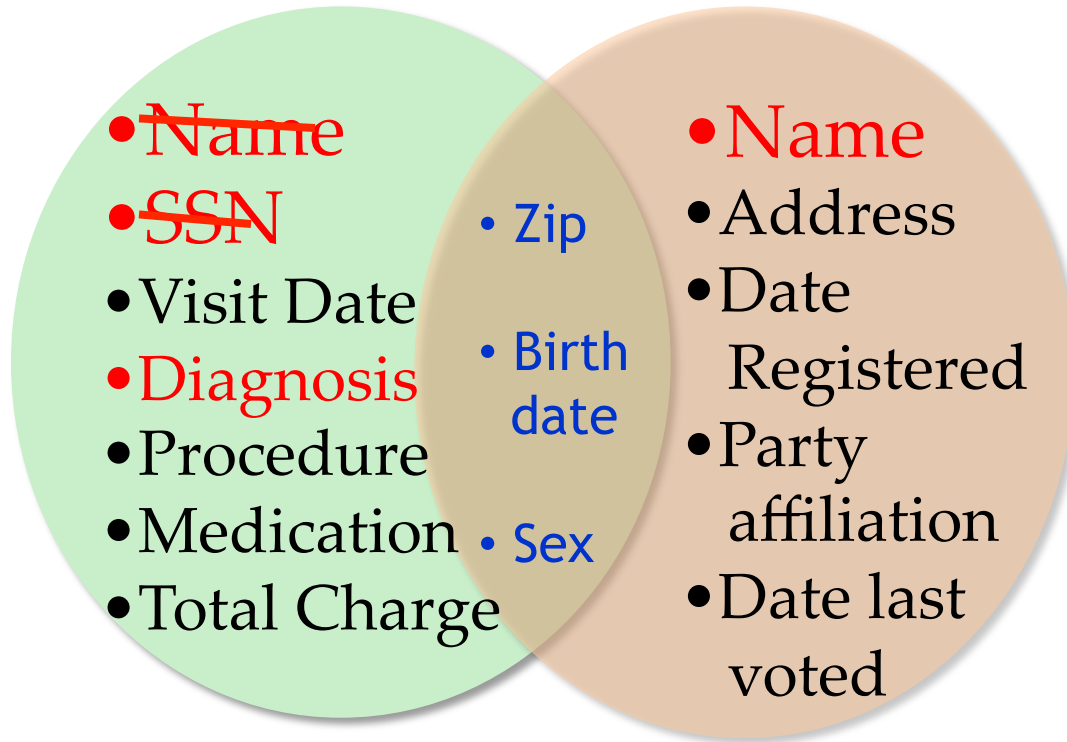
# The Ugly side of Record Linkage
[Sweeney IJUFKS 2002]

- ~~Name~~
- ~~SSN~~
- Visit Date
- Diagnosis
- Procedure
- Medication
- Total Charge
- Zip
- Birth date
- Sex

**Medical Data**

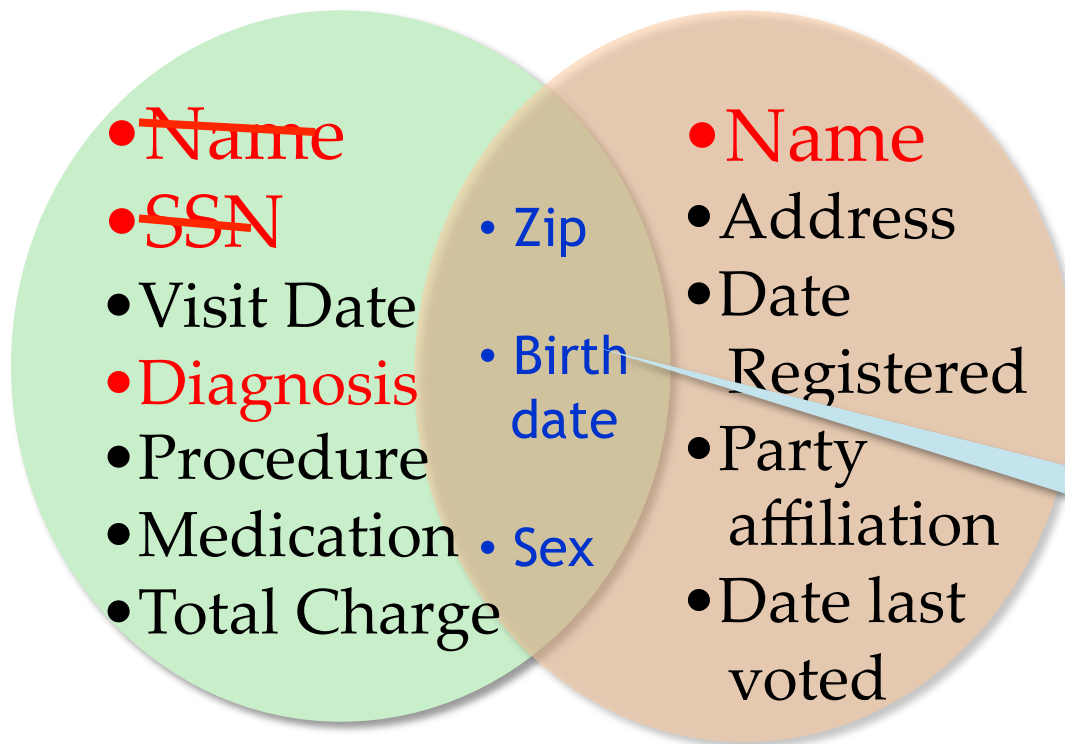# The Ugly side of Record Linkage
[Sweeney IJUFKS 2002]



**Medical Data**   **Voter List**

- Governor of MA **uniquely identified** using ZipCode, Birth Date, and Sex.

**Name linked to Diagnosis**

# The Ugly side of Record Linkage
[Sweeney IJUFKS 2002]



Medical Data — Voter List

- Name (crossed out)
- SSN (crossed out)
- Visit Date
- Diagnosis
- Procedure
- Medication
- Total Charge

- Zip
- Birth date
- Sex

- Name
- Address
- Date Registered
- Party affiliation
- Date last voted

**Quasi Identifier**

- 87 % of US population **uniquely identified** using ZipCode, Birth Date, and Sex.