

Web Based Linkage

Ergin Elmacioglu¹ ergin@psu.edu
Min-Yen Kan² kanmy@comp.nus.edu
Dongwon Lee¹ dongwon@psu.edu
Yi Zhang³ yiz@soe.ucsc.edu

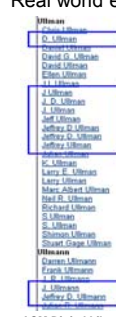
¹The Pennsylvania State University, USA
²National University of Singapore, Singapore
³University of California, Santa Cruz, USA






Motivation

- Real world examples of the linkage problem





ACM Digital Library IMDb & Wikipedia

Name Linkage

- Problem Definition:**
 - The process of detecting and merging duplicate *named entities* that represent the same real-world *object*
- Other real world examples**
 - Electronic devices (Apple iPod Nano 4GB vs. 4GB iPod nano)
 - Automobile models (Honda Fix vs. Honda Jazz)
 - Companies (T-Fal vs. Tefal)
 - Person names
 - Customs (Jane Doe vs. Doe, Jane)
 - Marriage (Carol Dusseau vs. Carol Arpac-Dusseau)
 - Misc. (Sean Engelson vs. Shlomo Argamon)

WIDM 2007 3

Name Linkage using Collective Knowledge

- There are many effective solutions if enough and discriminative contents are available
 - E.g., Contents based record linkage techniques
- Becomes challenging when
 - Not enough content to compare
 - Content does not help to identify

- Proposal: use external knowledge**
 - Ask people what they think*
 - Collective knowledge of people from the Web

WIDM 2007 4

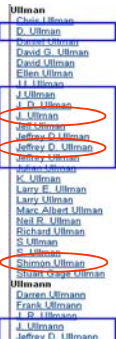
Name Linkage using Collective Knowledge

- Main Idea:** Use the Web as a collective knowledge of people in solving the Name Linkage Problem
- Hypothesis:**

If an entity **e1** is a duplicate of another entity **e2**, and if **e1** frequently appears together with information **I** on the Web, then **e2** may appear frequently with **I** on the Web, too.

WIDM 2007 5

Small Test



- Search results from Google:**
 - "Jeffrey D. Ullman" 384,000 pages 45%
 - "Jeffrey D. Ullman" + "aho" 174,000 pages
 - "J. Ullman" 124,000 pages
 - "J. Ullman" + "aho" 41,000 pages 33%
 - "Shimon Ullman" 27,300 pages
 - "Shimon Ullman" + "aho" 66 pages 0%

WIDM 2007 6

Web Based Linkage: Overview

WIDM 2007 7

Step 1. Select Representative Data

- What to select
 - A single token "aho"
 - A key phrase "stanford professor"
 - A sentence or more?
- How to select
 - Assess importance
 - $tf, tf*idf, \text{latent topic models}, \dots$
- How many to select
 - 1, 2, ... n
- Where to select from?
 - Contents of canonical entity, variant, both

WIDM 2007 8

Step 2. Acquire Knowledge from Web

- How to form the query?
 - Single information "I" (the most important data piece)
 - "Jeffrey D. Ullman" AND "Aho"
 - Multiple information "I₁", "I₂", "I₃", ... (the most k important data pieces)
 - Conjunction: "Jeffrey D. Ullman" AND "Aho" AND "database" AND "vldb"...
 - Disjunction: "Jeffrey D. Ullman" AND ("Aho" OR "database" OR "vldb"...
 - Hybrid: "Jeffrey D. Ullman" AND "Aho" AND ("database" OR "vldb"...

WIDM 2007 9

Step 3. Interpret the Collective Knowledge

For entities e_c, e_i and information t_c

- Page Counts

$sim(e_c, t_c)$	Jeffrey D. Ullman J. Ullman portal.acm.org infolab.stanford.edu en.wikipedia.org theory.lcs.mit.edu	Jeffrey D. Ullman Shimon. Ullman portal.acm.org	t_c
	= 4/16	= 1/19	
- URLs

$sim(e_c, t_c)$	Jeffrey D. Ullman J. Ullman	Jeffrey D. Ullman Shimon Ullman	t_c
	= 1/(174,000 - 41,000)	= 1/(174,000 - 66)	

WIDM 2007 10

Step 3. Interpret the Collective Knowledge

- Web Page Contents
 - Use top-k returned Web pages for each entity
 - Represent each set by a **Virtual Document**
 - Some heuristics
 - D (m): Top m ($\leq k$) documents are concatenated
 - T (all, n): Top n tokens with the highest weight from all top-k web pages
 - Snippet (m): Snippets of top m ($\leq k$) web pages
- Probabilistic Language Model: KL-divergence
 - $sim(e_c, e_i) = \text{doc_sim}(v\text{doc}(e_c), v\text{doc}(e_i))$

WIDM 2007 11

Virtual Document Creation

WIDM 2007 12

Experimental Validation

- Tested against
 - ACM, ArXiv, IMDB
- Variations tested
 - Step 1: single token, top-k tokens
 - Step 2: conjunctive query only
 - Step 3
 - Page count, URL
 - Virtual Document: 10 heuristics and 2 language models
- Search Engines
 - Google, MS Live Search

Presented Next

WIDM 2007 13

Experimental Validation

ACM data set:

- 43 authors
 - 14.2 citations/author
- 21 candidates/block
 - 3.1 citations/candidate
- 1.8 name variants/block
 - 6.7 citations/variant

Jeffrey D. Ullman

Foto N. Afrati, Chen Li, Jeffrey D. Ullman Using views to generate efficient evaluation plans for queries. *J. Comput. Syst. Sci.* 73: 703-724 (2007)

foto, n, afrati, chen, li, ...
 coauthor

using, views, generate, efficient, ...
 title

j, comput, syst, sci, 73, ...
 venue

WIDM 2007 14

Experimental Validation

IMDB data set:

- 50 actors
 - 24 titles/entity
- 20 candidates/block
 - 24 titles/candidate
- 1 name variant/block (a.k.a. name)
 - 23.5 titles/variant

Christina Aguilera

Moulin Rouge
What Women Want
The Bold and the Beautiful
The Next Best Thing
.....

moulin, rouge, bold, beautiful, ...
 Christina Aguilera

what_women_want, next, best, thing, ...
 Xtina (a.k.a.)

WIDM 2007 15

Scalability

- Not scalable:
 - A large number of Web accesses
 - Network traffic, load of search engine and web sites

Running times of the experiments on the ACM data set using the *tf* scheme on the title attribute

Methods	recall	Running time	# of Web accesses
Baseline - jaccard on the current content	0.495	0.19 min	0
Google - cosine using <i>D(3)</i> and Google	0.653	290.44 min	7951
Google - jaccard using <i>D(3)</i> and Google	0.711	289.32 min	7951
Google - lang. model 1 using <i>D(3)</i> and Google	0.572	289.29 min	7951
Google - lang. model 2 using <i>D(3)</i> and Google	0.756	289.47 min	7951
Google - cosine using <i>D(3)</i> and the local snapshot	0.462	5.53 min	0
Google - jaccard using <i>D(3)</i> and the local snapshot	0.445	6.38 min	0
Google - lang. model 1 using <i>D(3)</i> and the local snapshot	0.457	5.37 min	0
Google - lang. model 2 using <i>D(3)</i> and the local snapshot	0.494	5.39 min	0

- Solutions:
 - A better blocking scheme
 - Local snapshot of the Web
 - Stanford WebBase Project
 - ~100 million web pages from >50,000 sites including many .edu domains
 - Downloaded the half of the data & filtered
 - Local snapshot containing 3.5 million relevant pages

WIDM 2007 16

Related Work

- Abundant research on related problems
 - DB: approximate join, merge/purge, record linkage
 - DL: citation matching, author name disambiguation
 - AI: identity uncertainty
 - LIS: name authority control
- In a nutshell, existing approaches often do:
 - For two entities, e_1 and e_2 , capture their information in data structures, $D(e_1)$ and $D(e_2)$
 - Measure the distance or similarity between data structures: $dist(D(e_1), D(e_2)) = d$
 - Determine for matching:
 - If $d < threshold$, then e_1 and e_2 are matching entities
- Work well for common applications
- Ours performs better when
 - Entities **lack** useful information

WIDM 2007 17

Conclusion & Future Work

- The Name Linkage Problem
 - incomplete, noisy, non-descriptive data
 - usage of the Web to get additional information for the linkage process
- Future Work
 - A formal framework
 - Extension to "name disambiguation"
 - Experimentation on more & (larger) data sets
 - Scalability

WIDM 2007 18