### Scalable Uncertainty Management 01 – Introduction

Rainer Gemulla

April 20, 2012

Information & Knowledge Management Circa 1988

Databases SQL Datalog

Knowledge bases First-order logic Free text Information retrieval NLP



### Information & Knowledge Management Today

Structured	Ir	nformation	Unstructured
Sei F	<b>mantic Web</b> RDF DWL	Informatio Sens	on Extraction
Knowledge k First-order l	ogic	Deep Web	
SQL Datalog	Semi-Se XML	tructured Info.	Free text Information retrieval NLP
Databases	Web Service SOAP WSDL	es Hyperto HTML	ext -



SUM is about managing large amounts of uncertain data.

### Outline





### Sources of uncertainty

Certain data	Uncertain data	
The temparature is 25.634589 °C.	Sensor reported $25 \pm 1 ^{\circ}C$ .	Precision of devices
Bob works for Yahoo.	Bob works for Yahoo or Microsoft.	Lack of information
MPII is located in <mark>Saarbrücken</mark> .	MPII is located in Saarland.	Coarse-grained information
Mary sighted a finch.	Mary sighted either a finch (80%) or a sparrow (20%).	Ambiguity
lt will rain in Saarbrücken tomorrow.	There is a 60% chance of rain in Saarbrücken tomorrow.	Uncertainty about future
John's age is <mark>23</mark> .	John's age is in [20,30].	Anonymization
Paul is married to Amy.	Paul is married to Amy. Amy is married to Frank.	Inconsistent data

### Where does uncertainty arise?

Everywhere!

- Information extraction (D5 research)
- Sensor networks
- Business intelligence & predictive analytics
- Forecasting
- Scientific data management
- Privacy preserving data mining
- Data integration
- Data deduplication
- Social network analysis

## Entity disambiguation (AIDA)

#### Disambiguate each mention of an entity in a piece of text.

Disambiguation Method:	Input Type:TEXT	Run Information Graph Removal St	eps
prior prior+sim prior+sim+coherence (graph) Parameters: (default should be OK)	(Alexandria) Alexandria is an ancient city on the (Mediterranean Sea) Mediterranean . It was famous for its lighthouse, one of the seven	O: Alexandria (solved by local sim. onl     37: Mediterranean	y)
Similarity Impact 0.9	world wonders.	Candidate Entity	ME
Ambiguity degree 5		Meditemanean_Sea	0.446960
		Battle_of_the_Mediterranean	0.17493
Coherence threshold:		Mediterranean_Basin	0.016009
0.9		Mediterranean_Fleet	0.11436:
		Yom_Kippur_War	0.087808
Mention Extraction:		Napoleonic_Wars	0.418878
Stanford NER Manual		University_of_the_Mediterranean	1.875323
You can manually tag the mentions by putting them between [] and [].		Meditemanean_sea_\u0028oceanography\u0029	0.003496
HTML Tables are automatically disambiguated in the manual mode.		Meditemanean_diet	0.00818
🛃 🗋 🖪 🖌 🗓 44 📰 🗃 📰 Styles 🔹 Para		Meditemanean_naval_engagements_during_World_War_I	0.01274
¥ 42 C3 C3 C3 C3 A3 45 E E E = ● (**** A * ****		Southern_Europe	0.03332(
✓ = = = = = = = = = = = = = = = = = = =		Meditemanean_race	0.01892:
		1991_Mediterranean_Games	0.001450
Alexandria is an ancient city		Meditemanean_Region\u002c_Turkey	6.581360
on the Mediterranean. It was		Israeli_coastal_plain	9.286973
famous for its lighthouse, one		Ecology_of_California	0.006114
of the seven world wonders.		Mediterranean_pass	0.003953
		Meditemanean_Squadron_\u0028United_States\u0029	0.002118

- Find web pages concerning "The King of Rock'n'Roll" (entity search)
- How much fuzz about "Santorum" in each month of 2012? (*entity tracking*)

### Text segmentation

# Segment a piece of text into fields. E.g., "52-A Goregaon West Mumbai 400 062".



- Send a promotion to customers in West Mumbai.
- Find all papers containing YAGO in the title (faceted search)

### Relation extraction (NELL / Yago2)

#### Extract structured relations from the web.

instanc <del>e</del>	iteration	date learned	confidence
dried_squash_seeds is a <u>nut</u>	225	28-mar-2011	99.5 🖓
sinnett_thorn_mountain_cave is a cave	225	28-mar-2011	99.7 🍰
vail_road is a street	224	26-mar-2011	98.4 🍰
narold_macmillan is a <u>scientist</u>	225	28-mar-2011	96.6 🖓
132207 is a ZIP code	224	26-mar-2011	99.4 🍰
wday_tv_collaborates with bbc_news	224	26-mar-2011	96.9 🖓
imes controls friedman	227	03-apr-2011	96.9 🗳
support_personnel is a profession that is a kind of professionals	224	26-mar-2011	96.9 🍰
<u>bc_news</u> is a newspaper in the city washingtondc	224	26-mar-2011	99.2 🎝
witter operates the website twitter_com	225	28-mar-2011	100.0 🍰

street(98.4%)

- CPL @219 (98.4%) on 13-mar-2011 [ "ramp onto \_" "second right onto \_" "first traffic light onto \_" "bear left onto \_" "off ramp onto \_" "traffic light onto \_" ] using vail\_road
- CPL @66 (87.5%) on 17-mar-2010 [ 'first traffic light onto \_' 'off ramp onto \_' 'bear left onto \_' ] using vail\_road

- What is known about Albert Einstein? (fact search)
- Who has won a Nobel Prize and is born in Ulm? (question answering)

### Reasoning with uncertainty (URDF)



### Google Squared (discontinued)

#### Find and describe items of a given category.

G	oog e squa	labs			Square it Ad	ld to th	iis Square			
con	nedy movies									
	ltem Name	🔻 Release Date 💌 🗙	Genre	VX	Director	<b>V</b> X	Country	VX	Language	V
×	The Mask	29 July 1994	Comedy		Chuck Russell		USA		English	
×	Shrek 2	19 May 2004	Adventure	۲	Chuck Russell Directed by for T www.freebase.com	Гће М <u>m</u> - <u>a</u> l	ask II 10 sources	3		
×	Scary Movie	7 July 2000	Comedy	Oth	er possible values Dean Koontz Lo Author for The M www.freebase.com	w con Iask <u>m</u> - <u>a</u>	fidence II4 sources »			
×	Role Models	7 November 2008	Comedy	0	Bob Engelman Producer for The www.infibeam.com	Low c e Mas m · a	onfidence k I3 sources »			
×	Road Trip	19 May 2000	Comedy	0	Timothy Bond Low confidence Mask Comedy Movies and read product reviews Direc Timothy Bond, Release Date: September 07, 2004. Add list. Price Al www.shopping.com - all 2 sources.x			s <b>Director</b> 2004. Add to	:	
X	Old School	21 February 2003	Cornedy	Se	arch for more value Toda Phillips	<u>es »</u>	USA		English	_

- Directors that directed at least one comedy movie?
- Birthplaces of directors of comedy movies with a budget of over \$20M?

### Information integration

	Op.Sy	stem	CustId	Name	C	ity	State	J
C	1		$C_1$	John	San Fr	ancisco	CA	
Samer	2		$C_2$	Johnny	San	Jose	CA	}vvr
	1		$C_3$	Jack	San Fr	ancisco	CA	Í
	1		$C_4$	William	San Fr	ancisco	CA	Í
	2		$C_5$	Bill	San	Jose	CA	Í
			(	a) Custome	r Data			
							_	
		Op.S	System	TransID	CustID	Sales		
			1	$Tr_1$	$C_1$	\$15	]	
			1	$Tr_2$	$C_1$	\$5		
			2	$Tr_3$	$C_2$	\$30		
			2	$Tr_4$	$C_2$	\$20		
			1	$Tr_5$	$C_3$	\$30		
			1	$Tr_6$	$C_4$	\$90		
			2	$Tr_7$	$C_5$	\$25		
			2	$Tr_8$	$C_5$	\$15		
			(b	) Transactio	on Data		-	

Which one?

#### Example

• Turnover in San Francisco? And in California? (OLAP)

Sismanis et al., ICDE09.

### Predictive analytics



- What is the effect of changing the price on future sales?
- What is the risk associated with my portfolio?

### RFID & moving objects



- How many people are attending John's lecture?
- Where are choke points when moving items through my storage facility?

### Statistical & uncertain rules

#### **Smoking and Quitting in Groups**

Researchers studying a network of 12,067 people found that smokers and nonsmokers tended to cluster in groups of close friends and family members. As more people quit over the decades, remaining groups of smokers were increasingly pushed to the periphery of the social network.



- Does John smoke? (social network analysis)
- "Mississippi" most often refers to the state of Mississippi. (*entity disambiguation*)

### Anonymized data

	1	Non-Sen	sitive	Sensitive
	Zip Code	Age	Nationality	Condition
1	1305*	$\leq 40$	*	Heart Disease
4	$1305^{*}$	$\leq 40$	*	Viral Infection
9	$1305^{*}$	$\leq 40$	*	Cancer
10	1305*	$\leq 40$	*	Cancer
5	1485*	> 40	*	Cancer
6	$1485^{*}$	> 40	*	Heart Disease
7	1485*	> 40	*	Viral Infection
8	1485*	> 40	*	Viral Infection
2	1306*	$\leq 40$	*	Heart Disease
3	1306*	$\leq 40$	*	Viral Infection
11	1306*	$\leq 40$	*	Cancer
12	1306*	$\leq 40$	*	Cancer

#### Example

• Medical research, trend analysis, allocation of public funds, ....



1 Uncertainty in the Real World



### How to deal with uncertainty? (1)

Clean it (then deny it)!

- E.g., data warehouse systems
- Advantages
  - Lots of expertise and tools for cleaning data
  - Can be stored and queried in traditional DBMS
- Disadvantages
  - Loss of information
  - No risk assessment
  - High expense of cleaning
  - New data may "break" the clean database

#### • Important, but not covered in this lecture!

Customers					
	Sys	Cust	Name	City	State
Samal	1	$C_1$	John	SFO	CA
Samei	2	$C_2$	Johnny	SJ	CA
	1	C <sub>3</sub>	Jak	SFO	CA

#### CleanedCustomers

Cust	Name	City	State
C <sub>12</sub>	Johnny	SFO	CA
C <sub>3</sub>	Jak	SFO	CA

### How to deal with uncertainty? (2)

Manage it!

### Approach I: Incomplete databases

• A data integration scenario

Transactions

Sys	TransID	Cust	Sales
1	$T_1$	C1	\$15
1	$T_2$	$C_1$	\$5
2	T <sub>3</sub>	$C_2$	\$30
1	$T_4$	C <sub>3</sub>	\$30

#### Resolving entities via an incomplete database

ResolvedCustomersEntNameCityStateE1John || JohnnySFO || SJCAE2JakSFOCA

#### ResolvedTransactions

TransID	Ent	Sales
$T_1$	$E_1$	\$15
T <sub>2</sub>	$E_1$	\$5
T <sub>3</sub>	$E_1$	\$30
$T_4$	$E_2$	\$30

#### Some query results

Sales by city

City	Sum(Sales)	Status
SFO	\$30-\$80	guaranteed
SJ	\$50	non-guaranteed

#### Sales by state

State	Sum(Sales)	Status
CA	\$80	guaranteed

### Approach II: Probabilistic databases

• Bird watcher's observations

9	Sightin	gs	
	Name	Bird	Species
$t_1$ :	Mary	Bird-1	Finch: 0.8    Toucan: 0.2
<i>t</i> <sub>2</sub> :	Susan	Bird-2	Nightingale: 0.65    Toucan: 0.35
t3:	Paul	Bird-3	Humming bird: 0.55    Toucan: 0.45

• Which species exist in the park?

#### ObservedSpecies



• Observe: Cleaning up data by most likely choice would miss Toucan!

### Approach III: Probabilistic graphical models

- Anna and Bob are friends. Anna smokes, but does not have cancer. What do we know about Bob?
- Uncertain knowledge





### How to deal with uncertainty? (2)

Manage it!

- Advantages
  - No or little loss of information
  - Uncertainty might be resolved more accurately at query time
  - Risk assessment is possible
  - Less upfront effort
  - Arrival of new data handled gracefully
- Disadvantages
  - Increased cost of data processing
  - Active research area with lots of open issues (and interesting results)
  - No commercial DBMS systems available!
- This lecture!

#### Course overview

- Modelling uncertainty
  - Incomplete databases
  - Probabilistic databases
  - Probabilistic graphical models for relational data
- Managing uncertain data
  - Languages (relational algebra, datalog, relational calculus)
  - Provenance
  - Algorithms
  - Complexity
  - Approximation techniques
  - Systems
- Applications
  - Information extraction, sensor networks, business intelligence & predictive analytics, forecasting, scientific data management, privacy preserving data mining, data integration, data deduplication, social network analysis, ...

### Suggested reading

- Charu C. Aggarwal (Ed.) Managing and Mining Uncertain Data (Chapter 1) Springer, 2009.
- Daphne Koller, Nir Friedman *Probabilistic Graphical Models: Principles and Techniques* (Chapter 1) The MIT Press, 2009
- Dan Suciu, Dan Olteanu, Christopher Ré, Christoph Koch *Probabilistic Databases* (Chapter 1) Morgan & Claypool, 2011
- Charu C. Aggarwal, Philip S. Yu
   A Survey of Uncertain Data Algorithms and Applications
   IEEE Transactions of Knowledge and Data Engineering, 21(5),
   pp. 609–623, May 2009