#### Advanced Data Management (CSCI 640/490)

#### Data Fusion

Dr. David Koop





## Tidy Data Principles

- **Tidy Data**: Codd's 3rd Normal Form (Databases)
  - 1. Each variable forms a column
  - 2. Each observation forms a row
  - 3. Each type of observational unit forms a table (DataFrame)
- Other structures are messy data

D. Koop, CSCI 640/490, Spring 2024







### Problem: Variables stored in both rows & columns

Mexico Weather, Global Historical Climatology Network

id	year	month	element	d1	d2	d3	d4	d5	d6	d7	d8
MX17004	2010	1	tmax								
MX17004	2010	1	tmin								
MX17004	2010	2	tmax		27.3	24.1					
MX17004	2010	2	tmin		14.4	14.4					
MX17004	2010	3	tmax					32.1			
MX17004	2010	3	tmin					14.2			
MX17004	2010	4	tmax								
MX17004	2010	4	tmin								
MX17004	2010	5	tmax								
MX17004	2010	5	tmin								

#### D. Koop, CSCI 640/490, Spring 2024







### Problem: Variables stored in both rows & columns

Mexico Weather, Global Historical Climatology Network

id	year	month	element	d1	d2	d3	d4	d5	d6	d7	d8
MX17004	2010	1	tmax								
MX17004	2010	1	tmin								
MX17004	2010	2	tmax		27.3	24.1					
MX17004	2010	2	tmin		14.4	14.4					
MX17004	2010	3	tmax					32.1			
MX17004	2010	3	tmin					14.2			
MX17004	2010	4	tmax								
MX17004	2010	4	$\operatorname{tmin}$								
MX17004	2010	5	tmax								
MX17004	2010	5	$\operatorname{tmin}$								

Variable in columns: day; Variable in rows: tmax/tmin











## Solution: Melting + Pivot

id	date	element	value	id	date	tmax	$\operatorname{tmin}$
MX17004	2010-01-30	tmax	27.8	MX17004	2010-01-30	27.8	14.5
MX17004	2010-01-30	$\operatorname{tmin}$	14.5	MX17004	2010-02-02	27.3	14.4
MX17004	2010-02-02	tmax	27.3	MX17004	2010-02-03	24.1	14.4
MX17004	2010-02-02	$\operatorname{tmin}$	14.4	MX17004	2010-02-11	29.7	13.4
MX17004	2010-02-03	tmax	24.1	MX17004	2010-02-23	29.9	10.7
MX17004	2010-02-03	$\operatorname{tmin}$	14.4	MX17004	2010-03-05	32.1	14.2
MX17004	2010-02-11	tmax	29.7	MX17004	2010-03-10	34.5	16.8
MX17004	2010-02-11	$\operatorname{tmin}$	13.4	MX17004	2010-03-16	31.1	17.6
MX17004	2010-02-23	tmax	29.9	MX17004	2010-04-27	36.3	16.7
MX17004	2010-02-23	tmin	10.7	MX17004	2010-05-27	33.2	18.2

(a) Molten data

D. Koop, CSCI 640/490, Spring 2024

(b) Tidy data

[H. Wickham, 2014]



## Getting Lost in Transformations

Bureau of I.A.	
Regional Director	Numbers
Niles C.	Tel: (800)645-8397
	Fax: (907)586-7252
Jean H.	Tel: (918)781-4600
	Fax: (918)781-4604
Frank K.	Tel: $(615)564-6500$
	Fax: (615)564-6701



Intermediate Table

D. Koop, CSCI 640/490, Spring 2024

	Tel	Fax
Niles C.	(800)645-8397	
		(615)564-6701
Jean H.	(918)781-4600	
Frank K.	(615)564-6500	

#### Problem Table

		Tel	Fax
<b>_</b>	Niles C.	(800)645-8397	(907)586-7252
	Jean H.	(918)781-4600	(918)781-4604
l+	Frank K.	(615)564-6500	(615)564-6701

#### **Desired Solution**





# Foofah: Input, Output, and Transformations

Fax:(918)781-4604



Raw Data:

- A grid of values, i.e., spreadsheets "Somewhat" structured - must have some regular structure or is automatically generated.



User Input:

 Sample from raw data • Transformed view of the sample

Tel:(800)645-839



Program to synthesize: A loop-free Potter's Wheel [2] program

**Transformations Targeted:** 1. Layout transformation



#### D. Koop, CSCI 640/490, Spring 2024















## AutoSuggest

- Goals:
  - Automate "Complex" Data Preparation steps - Focus on frame transformations (not per-cell transformations)

  - Learn from Jupyter Notebooks
  - Use interactive methods to help users select from top-k options
- Two Types of Predictions:
  - Single-Operator Prediction: Given two tables and an operation, decide how to best apply the operation (what are the parameters)
  - Next-Operator Prediction: Given all operations performed so far, predict the next one









## **Pivot/Unpivot Prediction**

- Pivot is hard to get right
  - Index
  - Header
  - Aggregation Function
  - Aggregation Columns
- Use GroupBy Prediction
- Look for NULLs and use affiinity
- Affinity-Maximizing Pivot Table
- Unpivot requires compatibility

Sector	Ticker	Company	Year	Quarter	Market Cap	Revenu
Aerospace	AJRD	AEROJET ROCKETD	2006	Q1	1442.67	472.07
Aerospace	AJRD	AEROJET ROCKETD	2006	Q2	1514.80	489.22
Aerospace	BA	BOEING CO	2006	Q1	343.41	210.6
Utilities	YORW	YORK WATER CO	2008	Q4	600.19	271.73

Sector	Ticker	Company	2006	2007	2008
Aerospace	AJRD	AEROJET ROCKETD	6218.09	6342.45	7088.62
	ATRO	ASTRONICS CORP	1050.97	1071.99	1198.11
Business Services	HHS	HARTE-HANKS INC	2473.75	2523.22	2820.07
	NCMI	NATL CINEMEDIA	856.92	874.06	976.89
Consumer Staples	YTEN	TIELD10 BIOSCI	533.13	543.79	607.77
Utilities	YORW	YORK WATER CO	1902.37	1940.42	2168.70

Ticker	Company	Year	Aerospace	<b>Business Services</b>	 Utilit
AJRD	AEROJET ROCKETD	2006	6218.09	NULL	 NU
AJRD	AEROJET ROCKETD	2007	6342.45	NULL	 NU
AJRD	AEROJET ROCKETD	2008	7088.62	NULL	 NU
ATRO	ASTRONICS CORP	2006	1050.97	NULL	 NU
HHS	HARTE-HANKS INC	2006	NULL	2473.75	 NU
YORW	YORK WATER CO	2008	NULL	NULL	 2168

















#### Data Integration

select title, startTime from Movie, Plays where Movie.title=Plays.movie AND location="New York" AND director="Ava DuVernay"

Sources S1 and S3 are relevant, sources S4 and S5 are irrelevant, and source S2 is relevant but possibly redundant.



D. Koop, CSCI 640/490, Spring 2024

**Movie**: Title, director, year, genre Actors: title, actor **Plays**: movie, location, startTime **Reviews**: title, rating, description

S3	S4	S5
emas in NYC:	Cinemas in SF:	Reviews:
nema, title,	location, movie,	title, date
startTime	startingTime	grade, review









### Data Integration

- Lots of data sources, how do we answer questions where we need to access data from more than one?
- Schema matching
- Problem of heterogeneity
- Al-Complete problem: difficulty is the same as making computers as intelligent as people
- Two techniques:
  - Mediation
  - Data Warehouses





## Data Warehouses: Offline Replication

- Determine physical schema
- Define a database with this schema
- Define procedural mappings in an "ETL tool" to import the data and clean it.
- Periodically copy all of the data from the data sources
  - Note that the sources and the warehouse are basically independent at this point









#### Virtual Data Warehouses



D. Koop, CSCI 640/490, Spring 2024







#### Integrated Schema Example









## Why is Data Integration Hard?

- Systems-level reasons:
  - Managing different platforms
  - SQL across multiple systems is not so simple
  - Distributed query processing
- Logical reasons:
  - Schema (and data) heterogeneity
- 'Social' reasons:
  - Locating and capturing relevant data in the enterprise.
  - Convincing people to share (data fieldoms)
    - Security, privacy and performance implications







## Reading Quiz





#### <u>Assignment 3</u>

- Met Art Data
- Use OpenRefine & Pandas (no loops)





#### Data Fusion





## Record Linkage Motivation

- Often data from different sources need to be integrated and linked
  - To allow data analyses that are impossible on individual databases
  - To improve data quality
  - To enrich data with additional information
- Lack of unique entity identifiers means that linking is often based on personal information
- confidentiality is vital
- privacy concerns

When databases are linked across organisations, maintaining privacy and

• The linking of databases is challenged by **data quality**, **database size**, and







## Motivating Example

- Preventing the outbreak of epidemics requires monitoring of occurrences of unusual patterns of symptoms, ideally in real time
- Data from many different sources will need to be collected (including travel and immigration records; doctors, emergency and hospital admissions; drug purchases; social network and location data; and possibly even animal health data)



#### [P. Christen, 2019], image: [Pharexia, Wikipedia]



Northern Illinois University







#### Record Linkage

P. Christen





### Record Linkage Process













## Record Linkage Techniques

- Deterministic matching
  - Rule-based matching (complex to build and maintain)
- Probabilistic record linkage [Fellegi and Sunter, 1969]
  - Use available attributes for linking (often personal information, like names, addresses, dates of birth, etc.)
  - Calculate match weights for attributes
- "Computer science" approaches
  - Based on machine learning, data mining, database, or information retrieval techniques
  - Supervised classification: Requires training data (true matches) - Unsupervised: Clustering, collective, and graph based











## Record Linkage/Entity Resolution Recipe

- Problem: Link references to the same entity
- Short Answers:
  - Random Forest with attribute similarity features
  - Deep Learning to handle text and noise
  - End-to-end solutions still being worked on

D. Koop, CSCI 640/490, Spring 2024

#### [X. L. Dong and T. Rekatsinas, 2018]









### Data Integration and Data Fusion

- Data Integration: focus on integrating data from different sources • When sources are orthogonal, no problems
- What happens when two sources provide the same type of information and they conflict?
- Data Fusion: create a single object while resolving conflicting values









#### Data Fusion — Resolving Data Conflicts in Integration

X. L. Dong and F. Naumann





## Data Fusion Summary

- Conflict resolution strategies
- "Truth-discovery" techniques
  - Accuracy
  - Freshness
  - Dependence
- Fusion Issues
  - Accuracy
  - Efficiency
  - Usability
  - How fusion fits with the rest of data integration?







#### Data Conflicts



D. Koop, CSCI 640/490, Spring 2024

#### [L. Dong and F. Naumann, 2009]



Northern Illinois University







## Information Integration









## Information Integration









### Data Fusion

- Problem: Given a duplicate, create a single object representation while resolving conflicting data values.
- Difficulties:
  - Null values: Subsumption and complementation
  - Contradictions in data values
  - process
  - Metadata: Preferences, recency, correctness
  - Lineage: Keep original values and their origin
  - Implementation in DBMS: SQL, extended SQL, UDFs, etc.

- Uncertainty & truth: Discover the true value and model uncertainty in this







## Conflict Resolution Strategies









#### Integrating Conflicting Data: The Role of Source Dependence

X. L. Dong, L. Berti-Equille, and D. Srivastava





#### Discussion

- What is the paper's main contribution?
- Do you buy the argument? Any issues with the experiments?
- Can you think of any scenarios where the proposed technique will fail?
- Questions?







#### Example Problem













#### Example Problem

	SI	S2	<b>S3</b>
Stonebraker	MIT	Berkeley	MIT
Dewitt	MSR	MSR	UWisc
Bernstein	MSR	MSR	MSR
Carey	UCI	AT&T	BEA
Halevy	Google	Google	UW











### Naive Voting Works

	SI	S2	S3
Stonebraker	MIT	Berkeley	MIT
Dewitt	MSR	MSR	UWisc
Bernstein	MSR	MSR	MSR
Carey	UCI	AT&T	BEA
Halevy	Google	Google	UW











### Naive Voting Only Works if Data Sources are Independent













## Naive Voting Only Works if Data Sources are Independent

	SI	S2	S3	S4	S5
Stonebraker	MIT	Berkeley	MIT	MIT	MS
Dewitt	MSR	MSR	UWisc	UWisc	UWisc
Bernstein	MSR	MSR	MSR	MSR	MSR
Carey	UCI	AT&T	BEA	BEA	BEA
Halevy	Google	Google	UW	UW	UW











### S4 and S5 copy from S3

	SI	S2	S3	S4	S5
Stonebraker	MIT	Berkeley	MIT	MIT	MS
Dewitt	MSR	MSR	UWisc	UWisc	UWisc
Bernstein	MSR	MSR	MSR	MSR	MSR
Carey	UCI	AT&T	BEA	BEA	BEA
Halevy	Google	Google	UW	UW	UW











### S4 and S5 copy from S3

	SI	S2	S3	S4	S5
Stonebraker	MIT	Berkeley	MIT	MIT	MS
Dewitt	MSR	MSR	UWisc	UV/isc	UVVisc
Bernstein	MSR	MSR	MSR	MSR	MSR
Carey	UCI	AT&T	BEA	BEA	BEA
Halevy	Google	Google	UW	UW	UW











	SI	S2	S3	S4	S5
Stonebraker	MIT	Berkeley	MIT	MIT	MS
Dewitt	MSR	MSR	UWisc	UWisc	UWisc
Bernstein	MSR	MSR	MSR	MSR	MSR
Carey	UCI	AT&T	BEA	BEA	BEA
Halevy	Google	Google	UW	UW	UW









	SI	S2	S3	S4	S5	
Stonebraker	MIT	Berkeley	MIT	MIT	MS	
Dewitt	MSR	MSR	UWisc	UWisc	UWisc	
Bernstein	MSR	MSR	MSR	MSR	MSR	
Carey	UCI	AT&T	BEA	BEA	BEA	
Halevy	Google	Google	UW	UW	UW	









	SI	S2	<b>S3</b>	<b>S4</b>	S5
Stonebraker	MIT	Berkeley	MIT	MIT	MS
Dewitt	MSR	MSR	UWisc	UWisc	UWisc
Bernstein	MSR	MSR	MSR	MSR	MSR
Carey	UCI	AT&T	BEA	BEA	BEA
Halevy	Google	Google	UW	UW	UW

D. Koop, CSCI 640/490, Spring 2024

2. With only a snapshot it is hard to decide which source is a copier.











I. Sharing common data does not in itself imply copying.

	SI	S2	S3	S4	S5		
Stonebraker	MIT	Berkeley	MIT	MIT	MS		
Dewitt	MSR	MSR	UWisc	UWisc	UV/isc		
Bernstein	MSR	MSR	MSR	MSR	M\$R		
Carey	UCI	AT&T	BEA	BEA	BEA		
Halevy	Google	Google	UW	UW			
	3. A copier can also provide or verify some data by itself, so it is inappropriate to ignore all of its data.						

#### D. Koop, CSCI 640/490, Spring 2024

2. With only a snapshot it is hard to decide which source is a copier.









### Source Dependence

- directly or transitively from a common source (can be one of S or T).
  - Independent source
  - Copier
    - copying part (or all) of data from other sources
    - may verify or revise some of the copied values
    - may add additional values
- Assumptions
  - Independent values
  - Independent copying
  - No loop copying

D. Koop, CSCI 640/490, Spring 2024

# Source dependence: two sources S and T deriving the same part of data











#### Core Case

- Conditions
  - Same source accuracy
  - Uniform false-value distribution
  - Categorical value
- highest probability to be true.

#### Proposition: W. independent "good" sources, Naïve voting selects values with





Northern Illinois University









#### deas

- If two sources share a lot of false values, they are more likely to be dependent.
- highly different from the accuracy of S1.

• S1 is more likely to copy from S2, if the accuracy of the common data is











#### Combining Accuracy and Dependence

#### Source-accuracy Computation













### Combining Accuracy and Dependence

Source-accuracy Computation

#### Step 3

D. Koop, CSCI 640/490, Spring 2024



#### Step









## The Motivating Example

	SI	S2	S3	S4	S5		
Stonebraker	MIT	Berkeley	MIT	MIT	MS		
Dewitt	MSR	MSR	UWisc	UWisc	UWisc		
Bernstein	MSR	MSR	MSR	MSR	MSR		
Carey	UCI	AT&T	BEA	BEA	BEA		
Halevy	Google	Google	UW	UW	UW		
$S_2$ $Q_2$ $Q_2$ $Q_2$ $Q_2$ $Q_3$ $S_3$ $Rnd 2$ $S_2$ $Q_4$ $Q_5$ $Q_5$ $Q_5$ $S_4$ $Q_5$ $S_5$ $S_4$ $Q_5$ $S_5$ $S_1$ $S_2$ $S_2$ $Q_4$ $Q_5$ $Q_5$ $S_5$ $S_5$ $S_5$ $S_1$ $S_1$ $S_2$ $S_1$ $S_2$ $S_2$ $S_1$ $S_2$ $S_2$ $S_3$ $S_5$ $S_1$ $S_1$ $S_2$ $S_1$ $S_2$ $S_2$ $S_2$ $S_1$ $S_2$ $S_2$ $S_2$ $S_1$ $S_2$ $S_2$ $S_1$ $S_2$ $S_2$ $S_1$ $S_2$ $S_2$ $S_2$ $S_1$ $S_2$ $S_2$ $S_2$ $S_2$ $S_1$ $S_2$ $S_2$ $S_2$ $S_2$ $S_2$ $S_1$ $S_2$ $S_$							
	Rnd 3	Rnd II	S <sub>2</sub> S <sub>4</sub>	49 49.44 .55 .55 .49.44 .55 .55 .55 .55 .55 .55	[X L Dong		













### The Motivating Example

Accuracy	SI	S2	S3	S4	S5
Round I	.52	.42	.53	.53	.53
Round 2	.63	.46	.55	.55	.55
Round 3	.71	.52	.53	.53	.37
Round 4	.79	.57	.48	.48	.31
• • •					
Round 11	.97	.61	.40	.40	.21

Value	Carey			Halevy	
Confidence	UCI	AT&T	BEA	Google	UW
Round I	1.61	1.61	2.0	2.1	2.0
Round 2	1.68	1.3	2.12	2.74	2.12
Round 3	2.12	1.47	2.24	3.59	2.24
Round 4	2.51	1.68	2.14	4.01	2.14
Round 11	4.73	2.08	1.47	6.67	1.47





