## Advanced Data Management (CSCI 640/490)

Data Integration

Dr. David Koop





# Data Cleaning Types

- How can statistical techniques improve efficiency or reliability of data cleaning? (Data Cleaning with Statistics)
  - Example: Trifacta
  - Two tasks: Error Detection & Data Repairing
- How how can we improve the reliability of statistical analytics with data cleaning? (Data Cleaning for Statistics)
  - Example: SampleClean
- Task: Do statistics and clean along the way Similar questions if we substitute machine learning for statistics









## Misconceptions about Data Cleaning

- The end goal of data cleaning is clean data
- Data cleaning is a sequential operation
- Data cleaning is performed by one person
- Data quality is easy to evaluate













## Classifying Data Quality Problems



### D. Koop, CSCI 640/490, Spring 2023



Northern Illinois University



## Dirty and Cleaned Data

## (a) Dirty Data

id	title	pub_year	citation _count
<b>t</b> 1	CrowdDB	11	18
<b>t</b> 2	TinyDB	2005	1569
<b>t</b> 3	YFilter	Feb, 2002	298
<b>t</b> 4	Aqua		106
<b>t</b> 5	DataSpace	2008	107
<b>t</b> 6	CrowdER	2012	1
<b>t</b> 7	Online Aggr.	1997	687
•••	•••	•••	•••
<b>t</b> 10000	YFilter - ICDE	2002	298

### D. Koop, CSCI 640/490, Spring 2023

## (b) Cleaned Sample

id	title	pub_year	citation _count	#dup
<b>t</b> 1	CrowdDB	2011	144	2
<b>t</b> 2	TinyDB	2005	1569	1
<b>t</b> 3	YFilter	2002	298	2
<b>t</b> 4	Aqua	1999	106	1
<b>t</b> 5	DataSpace	2008	107	1
<b>t</b> 6	CrowdER	2012	34	1
<b>t</b> 7	Online Aggr.	1997	687	3

[J. Wang et al., 2014]









## SampleClean Framework



D. Koop, CSCI 640/490, Spring 2023





6

## Comparing the Two Approaches







## HoloClean

- quantitative methods:
  - Qualitative: use integrity constraints or external data sources
  - Quantitative: use statistics of the data
- Driven by probabilistic inference. Users only need to provide a dataset to be cleaned and describe high-level domain specific signals.
- Can scale to large real-world dirty datasets and perform automatic repairs with high accuracy

A holistic data cleaning framework that combines qualitative methods with









## HoloClean

### Input

	Dataset to be cleaned						
	DBAName	Address	City	State	Zip		
t1	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		
t2	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60609		
t3	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60609		
t4	Johnnyo's	3465 S Morgan ST	Cicago	(IL	60608		

### **Denial Constraints**

- c1: DBAName  $\rightarrow$  Zip
- c2: Zip  $\rightarrow$  City, State
- c3: City, State, Address  $\rightarrow$  Zip

### Matching Dependencies

m1:  $Zip = Ext_Zip \rightarrow City = Ext_City$ m2:  $Zip = Ext_Zip \rightarrow State = Ext_State$ m3: City = Ext\_City  $\land$  State = Ext\_State  $\land$  $\land$  Address = Ext\_Address  $\rightarrow$  Zip = Ext\_Zip

### External Information

Ext_Address	Ext_City	Ext_State	Ext_Zip
3465 S Morgan ST	Chicago	IL	60608
1208 N Wells ST	Chicago	IL	60 <mark>61</mark> 0
259 E Erie ST	Chicago	IL	60611
2806 W Cermak Rd	Chicago	IL	60623

### D. Koop, CSCI 640/490, Spring 2023



## Output

	Proposed Cleaned Dataset						
	DBAName	Address	City	State	Zip		
t1	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		
t2	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		
t3	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		
t4	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		

### **Marginal Distribution** of Cell Assignments

Cell	Possible Values	Probability
10 71	60608	0.84
t2.ZIP	60609	0.16
14 01	Chicago	0.95
t4.City	Cicago	0.05
	John Veliotis Sr.	0.99
t4.DBAName	Johnnyo's	0.01













## Pandas Operations

- Filtering out missing data (dropna):
  - Can choose rows or columns
- Filling in missing data (fillna):
- Finding problems in data (statistics, special values)
- Filtering duplicates: (drop\_duplicates, unique)
- Cleaning data: (map, replace, clamping data)

special values) ates, unique) oping data)





## <u>Assignment 3</u>

- Salary Data
- Use Pandas (not loops)
- Part 5: use melt/pivot or a similar high-level operation

### D. Koop, CSCI 640/490, Spring 2023

# • Part 2: CSCI 640 students need to do (b), CSCI 490 students can choose





# Outline

- Combining Data
- Data Integration
- Data Matching (Entity Resolution)
- Data Fusion
- Data Fusion Techniques
  - Integrating Conflicting Data: The Role of Source Dependence, X. L. Dong et al., 2009
  - Quiz at the beginning of class





## Example: Football Game Data

- Data about football games, teams, & players
  - Game is between two Teams
  - Each Team has Players
- For each game, we could specify every player and all of their information... why is this bad?





## Example: Football Game Data

- Data about football games, teams, & players
  - Game is between two Teams
  - Each Team has Players
- For each game, we could specify evaluate player and all of their information... this bad?
- Normalization: reduce redundancy, keep information that doesn't change separate
- 3 Relations: Team, Player, Game
- Each relation only encodes the data specific to what it represents

### Player

very	
why	is

ld	Name	Height		Weight		
Team						
ld	Name	Wins		Losses		
Game						
ld	Location			Date		





## Example: Football Game Data

- Have each game store the id of the home team and the id of the away team (one-toone)
- Have each player store the id of the team he plays on (many-to-one)

• What happens if a player plays on 2+ teams?







## How does this relate to pandas?

- DataFrames in pandas are ~relations (tables)
- We may wish to normalize data in a similar manner in pandas
- However, operating on 2+ DataFrames at the same time can be unwieldy, can we merge them together?
- Two potential operations:
  - Have football game data (just the Game table) from 2013, 2014, and 2015 and wish to merge the data into one data frame
  - Have football game data and wish to find the average temperature of the cities where the games were played





15

## Concatenation

- Take two data frames with the same columns and add more rows
- pd.concat([data-frame-1, data-frame-2, ...])
- Default is to add rows (axis=0), but can also add columns (axis=1)
- Can also concatenate Series into a data frame.
- concat preserves the index so this can be confusing if you have two default indices (0,1,2,3...)—they will appear twice
  - Use ignore\_index=True to get a 0,1,2...

## e columns and add more rows ta-frame-2, ...])





## Merges (aka Joins)

- Example: Football game data merged with temperature data

### Game

Id	Location	Date	Home	Away
0	Boston	9/2	1	15
1	Boston	9/9	1	7
2	Cleveland	9/16	12	1
3	San Diego	9/23	21	1

## No data for San Diego-

D. Koop, CSCI 640/490, Spring 2023

# Need to merge data from one DataFrame with data from another DataFrame

## Weather

wld	City	Date	Temp
0	Boston	9/2	72
1	Boston	9/3	68
7	Boston	9/9	75
21	Boston	9/23	54
			•••
36	Cleveland	9/16	81







# Merges (aka Joins)

- Want to join the two tables based on the location and date Location and date are the keys for the join
- What happens when we have missing data?
- Merges are **ordered**: there is a left and a right side
- Four types of joins:
  - Inner: intersection of keys (match on both sides)
  - Outer: union of keys (if there is no match on other side, still include with NaN to indicate missing data)
  - Left: always have rows from left table (no unmatched right data) Right: like left, but with no unmatched left data







## Inner Strategy

## Merged

Id	Location	Date	Home	Away	Temp	wld
0	Boston	9/2	1	15	72	0
1	Boston	9/9	1	7	75	7
2	Cleveland	9/16	12	1	81	36

No San Diego entry





## Outer Strategy

### Merged

Id	Location	Date	Home	Away	Temp	wld
0	Boston	9/2	1	15	72	0
NaN	Boston	9/3	NaN	NaN	68	1
1	Boston	9/9	1	7	75	7
NaN	Boston	9/10	NaN	NaN	76	8
NaN	Cleveland	9/2	NaN	NaN	61	22
2	Cleveland	9/16	12	1	81	36
3	San Diego	9/23	21	1	NaN	NaN









## Left Strategy

## Merged

Id	Location	Date	Home	Away	Temp	wld
0	Boston	9/2	1	15	72	0
1	Boston	9/9	1	7	75	7
2	Cleveland	9/16	12	1	81	36
3	San Diego	9/23	21	1	NaN	NaN





# Right Strategy

## Merged

Id	Location	Date	Home	Away	Temp	wld
0	Boston	9/2	1	15	72	0
NaN	Boston	9/3	NaN	NaN	68	1
1	Boston	9/9	1	7	75	7
NaN	Boston	9/10	NaN	NaN	76	8
NaN	Cleveland	9/2	NaN	NaN	61	22
2	Cleveland	9/16	12	1	81	36
					•••	

## No San Diego entry









# Data Merging in Pandas

- pd.merge(left, right, ...) or left.merge(right, ...)
- Default merge: join on matching column names
- Better: specify the column name(s) to join on via on kwarg
  - If column names differ, use left on and right on
  - Multiple keys: use a list
- not being joined on
- Can also merge using the index by setting left index Or right index to True

• how kwarg specifies type of join ("inner", "outer", "left", "right") • Can add suffixes to column names when they appear in both tables, but are







# Merge Arguments

Argument	Description		
left	DataFrame to be merged on the left sid		
right	DataFrame to be merged on the right sight		
how	One of 'inner', 'outer', 'left',		
ΟΠ	Column names to join on. Must be found given, will use the intersection of the col		
left_on	Columns in left DataFrame to use as jo		
right_on	Analogous to left_on for left DataF		
left_index	Use row index in left as its join key (or		
right_index	Analogous to left_index.		
sort	Sort merged data lexicographically by joi some cases on large datasets).		
suffixes	Tuple of string values to append to colun 'data' in both DataFrame objects, wo		
сору	If False, avoid copying data into result copies.		
indicator	Adds a special column _merge that ind 'right_only', or 'both' based o		

### D. Koop, CSCI 640/490, Spring 2023

le.

```
or 'right'; defaults to 'inner'.
```

in both DataFrame objects. If not specified and no other join keys lumn names in left and right as the join keys.

oin keys.

-rame.

<sup>•</sup> keys, if a MultiIndex).

in keys; True by default (disable to get better performance in

nn names in case of overlap; defaults to  $('_x', '_y')$  (e.g., if uld appear as 'datax' and 'datay' in result). ting data structure in some exceptional cases; by default always

icates the source of each row; values will be 'left\_only', the origin of the joined data in each row. [W. McKinney, Python for Data Analysis]









# Outline

- Combining Data
- Data Integration
- Data Matching (Entity Resolution)
- Data Fusion
- Data Fusion Techniques
  - Integrating Conflicting Data: The Role of Source Dependence, X. L. Dong et al., 2009
  - Quiz at the beginning of class







## Introduction to Data Integration

A. Doan, A. Halevy, and Z. Ives





## Data Integration

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

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 2023

**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 Matching & Data Fusion

- <u>Google Thinks I'm Dead</u> (<u>I know otherwise.</u>) [R. Abrams, NYTimes, 2017]
- Not only Google, but also Alexa:
  - "Alexa replies that Rachel Abrams is a sprinter from the Northern Mariana Islands (which is true of someone else)."
  - "He asks if Rachel Abrams is deceased, and Alexa responds yes, citing information in the Knowledge Graph panel."









# Data Integration, Data Matching, & Data Fusion

- Data Integration: focus on integrating data from different sources Data Matching (aka Entity Resolution aka Record Linkage): want to know that two entities (often in different sources) are the same "real"
- entity
- 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









# Outline

- Combining Data
- Data Integration
- Data Matching (Entity Resolution)
- Data Fusion
- Data Fusion Techniques
  - Integrating Conflicting Data: The Role of Source Dependence, X. L. Dong et al., 2009
  - Quiz at the beginning of class







## Data Fusion — Resolving Data Conflicts in Integration

X. L. Dong and F. Naumann





# Information Integration









# Information Integration









# Outline

- Combining Data
- Data Integration
- Data Matching (Entity Resolution)
- Data Fusion
- Data Fusion Techniques
  - Integrating Conflicting Data: The Role of Source Dependence, X. L. Dong et al., 2009
  - Quiz at the beginning of class





