CS 541: Advanced Data Management

Dong Deng

1

Self Introduction



Dong Deng

- Assistant Professor at Rutgers University
- Senior Scientist at Inception Institute of AI
- Postdoc from MIT Database Group
- Ph.D. from Tsinghua University

10+ years of research experience in the **data management** field

Course Object

- Introducing students the cutting-edge research on data management and data curation (preparation)
- Training students to master basic skills for being a researcher

Course Setup

- Grading
 - Paper Review: 25% (write review for 1 paper each week)
 - Paper Presentation: 25% (35mins + 15mins QA, 1 or 2 presentations depends on enrollment)
 - Paper Discussion: 10% (ask at least 10 questions during paper discussions)
 - Programming Assignments: $2 \times 10\% = 20\%$
 - Exam: 20% (based on contents in the lectures)
- TA

– Yanshi Luo yanshi.luo@rutgers.edu

Policy

- Don't be Late
 - Everyone has a budget of 2 days to be used on assignments
 - Once it is used up, 20% penalty per day for each late day
- Don't Cheat
 - We will do plagiarism check
 - If you got caught, you will fail this course

If you are struggling, let us know!

Course Overview

Data Science



Data Curation (Preparation)



Data Curation (Preparation)



The Landscape of Data Curation



Data Curation Example



Data Curation Example



Data Processing Pipeline

• What you think you do?

• What you really do?



Data

Analysis/

Modeling

Predictive

Model

Course Topics

- Data Discovery
- Data Wrangling
- Data Transformation
- Data Standardization
- Data Cleaning
- Data Integration
- Data Visualization

They are all interrelated

Data Discovery

• Merck has approximated 4000 Oracle databases, a large data lake, and uncountable individual files. Data scientists spend >90% of their time finding datasets relevant to their task at hand.

- Many data lakes:
 - US Government open data: www.data.gov
 - Climate Data Library: <u>http://iridl.ldeo.columbia.edu/index.htmln</u>

Data Wrangling

• The process of transforming and mapping data from one "raw" data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.

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	

. . .

	Tel	Fax
Niles C.	(800)645-8397	(907)586-7252
Jean H.	(918)781-4600	(918)781-4604
Frank K.	(615)564-6500	(615)564-6701

Data Wrangling

• The process of transforming and mapping data from one "raw" data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.

Gri	id Columns			Find column	Filters ~ < Recipe	Add Step	
		Source		Preview	Choose e trep	oformation	
	transaction_date	✓ #.# ticket_pr	ce ~ #.# disc	ount #.# column1	choose a train	sionation	
					derive		
		1 I I I		- I C	Constant of America	- luma with the second of a fermina	
	010 - Dec 2010	0 - 30	0.00 - 0.20	0 - 30	Creates a new o	column with the result of a formula	
	/00/19	0.01	0	0.01			
	/08/19	29.99	0	29.99	Formula		
	/09/18	29.99	0	29.99			
	/10/18	29.99	0	29.99	ticket_price * (1 - discount)		
	/11/17	29.99	0	29.99			
	/12/17	29.99	0	29.99			
	/01/05	0.01	0	0.01	Group by		
	/01/05	9.99	0	9.99			
	/02/04	9.99	0	9.99	Choose cord		_
	/03/06	9.99	0.05	9.49049999999999			
	/04/05	9.99	0.05	9.49049999999999	Order by		Tritocta
	/05/05	9.99	0	9.99	order by		IIIICLA
	/06/04	9.99	0	9.99	Edit formul	la	
	/07/04	9.99	0	9.99			
	/08/03	9.99	0	9.99			
	/09/02	9.99	0	9.99	New column n	name	
	/10/02	9.99	0	9.99	Construction of the second sec		
	/11/01	9.99	0.05	9.49049999999999	String		

Data Standardization



https://unsupervisedmethods.com/nips-accepted-papers-stats-26f124843aa0

Data Standardization

"It strikes me as a little comical that it required such a kludgey effort to pull these numbers together..."

"I wrote a bunch of manual patterns to map names to canonical versions ("UCL" to "University College London" and "Google Inc" to "Google"), although it is likely that I still missed some cases..."

"I've created a number of rules to map together alternative organization names and misspellings..."

IBM Research IBM Research, NY IBM Research, USA IBM T J Watson Research Center IBM T. J. Watson Research IBM T.J Watson Research Center IBM T.J. Watson Research Center IBM T.J. Watson Research Center IBM Thomas J. Watson Research Center IBM TJ Watson Research Center

.

Inaccurate



Inconsistent

FlightView		FlightAware AAL119 (Track inbound flight)	Orbitz
FLIGHT TRACKER	Aircraft Origin	(meh.air) (11.8100) American Atrines "American" Boeing 737-800 (twin-jet) (B738/Q - track or photos) Terminal A / Gate 32 / Newark Liberty Intl (KEWR - tr	American Airlines # 119 Leg 1: In Transit
Departure Airport: Scheduled Time: 6:15 PM, Dec 08 Takeoff Time: 6:53 PM, Dec 08 Terminal - Gate: Terminal A - 32 ArrivalStatus: In Air	Destination Route Date Duration Progress	Terminal 4 / Gate 42B / Los Angeles Inti (KLAX - <u>track</u> <u>Other flights between these airports</u> ZDMAZ Q42 BTERX Q400 AIR 780 VHP 780 MCI 724 SLN 7102 ALS 744 RSK 7 <u>Decodo</u> 2011年 12月 08日 (Thursday) 5 hours 43 minutes 20 minutes left 5 hours 23 minutes	Departs: Newark (EWR) <u>View real-time airpo</u> Gate: 32 Scheduled Estimated Actual 6:22p Dec 8 - Dec 8
Airport: Scheduled Time: 9:40 PM, Dec 08 9:42 PM, Dec 08 Estimated Time: Track This Flight Live! Time Remaining: 25 min Terminal - Gate: Terminal 4 - 42B Baggage Claim: 4	Status Distance Fare Cabin Departure Arrival	En Route (2,284 sm down; 168 sm to go) Direct: 2,451 sm Planned: 2,458 \$51.99 to \$3,561.11; average: \$241.96 (airline insight) First: Dinner / Economy: Food for sale Scheduled 7-day Average Actual/Estimated 06:15PM EST 07:08PM EST 06:53PM EST 08:33PM PST 09:17PM PST 09:36PM PST	Arrives: Los Angeles (LAX) <u>View real-time ai</u> Gate: 42B Scheduled Estimated Actual 9:54p Dec 8 9:47p Dec 8

Incomplete

Country \$	UN R/P 10% ^[4] \$	UN R/P 20% ^[5] \$	World Bank Gini (%) ^[6]	WB Gini (year) ≎	CIA R/P 10% ^[7]	Year	CIA Gini (%) ^[8]	CIA Gini (year)	GPI Gini (%) ^[9] \$
Z Seychelles			65.8	2007					\frown
Comoros			64.3	2004					
Mamibia	106.6	56.1	63.9	2004	129.0	2003	59.7	2010	
South Africa	33.1	17.9	63.1	2009	31.9	2000	65.0	2005	
Botswana	43.0	20.4	61.0	1994			63	1993	
Haiti	54.4	26.6	59.2	2001	68.1	2001	59.2	2001	
Angola			58.6	2000					62.0
Honduras	59.4	17.2	57.0	2009	35.2	2003	57.7	2007	



- Address errors caused 6.8 billion undelivered mails in 2013
- Estimated \$1.5 billion spent on processing
- At least \$3.4 billion wasted postage

DATA



Bad Data Costs the U.S. \$3 Trillion Per Year

by Thomas C. Redman

SEPTEMBER 22, 2016

Data Integration





Comparison Shopping

"GE estimates they would save \$100M/year by integrating orders for better pricing" FORTUNE

Supplier Mastering, Customer Mastering, Collaborative Research etc.

Data Integration





Comparison Shopping

"GE estimates they would save \$100M/year by integrating orders for better pricing" **FORTUNE**

Will be increasingly more important

- ➤ data sharing becomes pervasive
- translation of legacy data

Data Integration is Critical



Who has More Data Sources?

- Large manufacturing enterprise
 - Has 325 procurement systems
 - Estimates they would save \$100M/year by "most favored nation status"
- Large drug company
 - Has 10,000 bench scientists
 - Wants to integrate their "electronic lab notebooks"
- Large auto company
 - Wants to integrate customer databases In Europe In 40 languages

A Closer Look at Data Integration

Entity Resolution: Find Duplicate Records

Entity Consolidation: *Merge Duplicate Records*

ID	Name	Address	Telephone
P1	Mary Lee	9 St, Wisconsin	(718) 453-0681
P2	M. Lee	9th St, WI	7184530681
P3	James Smith	3rd E Ave, CA	212-213-2888x264
P4	J. Smith	3 E Avenue, CA	(212) 213-2888
P5	Lee, Mary	9 Street, WI	+1-718-453-0681
P6	Smith, James	5th Street, WA	+1-212-213-2888





Schema Mapping

Entity Resolution



Similarity Measures



Set Similarity Measures

		Unnormalized	Score	Normalized Score	
	Sequence	Edit Distar	nce	Edit Similarity	
	Set	Overlap Size		Jaccard/Cosine Similarity	
Overlap($(A,B) = A \cap B $	3	A		
Jaccard	$(A,B) = \frac{ A \cap B }{ A \cup B }$	0.6			
Cosine($(A,B) = \frac{ A \cap B }{\sqrt{ A B }}$	0.75	B		

Set Similarity Measures

	Unnormalized	Normalized
Sequence	Edit Distance	Edit Similarity
Set	Overlap Size	Jaccard/Cosine Similarity

 $Overlap(A,B) = |A \cap B| \qquad 2$

$$Jaccard(A,B) = \frac{|A \cap B|}{|A \cup B|} \qquad 0.66 \qquad \begin{array}{c} David Smith \rightarrow A = \{David, Smith\} \\ Smith R. David \rightarrow B = \{David, R., Smith\} \\ Cosine(A,B) = \frac{|A \cap B|}{\sqrt{|A||B|}} \qquad 0.82 \qquad 32 \end{array}$$