Advanced Data Management (CSCI 640/490)

Review

Dr. David Koop





Data systems rely on algorithms

DATA SYSTEMS ALGORITHMS















Data structures define performance



D. Koop, CSCI 640/490, Spring 2023



register = this room caches = this city memory = nearby city disk = Pluto

Jim Gray, Turing Award 1998











Tradeoffs in each structure







"Traditional" Database Research











Learned Data Structures and Algorithms









B-Tree



D. Koop, CSCI 640/490, Spring 2023

7

Model to Predict Data's Location on Disk

Frequency Distribution

Cumulative Distribution Function (CDF)

MacMenamin

date	7-01-01	7-01-02	7-01-02	.7-01-03	.7-01-03	7-01-04	7-01-04	.7-01-05	7-01-05	7-01-06	7-01-07	7-01-09	7-01-09	7-01-09	7-01-10	7-01-10	7-01-11	.7-01-12	.7-01-13	7-01-14	7-01-15	7-01-16	7-01-17	7-01-18	7-01-19	7-01-20	7-01-21	7-01-22	7-01-22	7-01-22	7-01-23
da	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017	2017

Challenges

D. Koop, CSCI 640/490, Spring 2023

Frameworks are not designed for nano-second execution

ML+System Co-Design

Recursive Model Index (RMI)

2-Stage RMI with Linear Model $pos_0 = a_0 + b_0 * key$ $pos_1 = m_1[pos_0].a + m_1[pos_0].b * key$ $record = local-search(key, pos_1)$

Sandwiched Bloom Filter

D. Koop, CSCI 640/490, Spring 2023

[M. Mitzenmacher, 2018 via T. Kraska, 2019]

Sorting

(a) CDF Model Pre-Sorts

D. Koop, CSCI 640/490, Spring 2023

(b) Compact & local sort

12

Sorting

(a) CDF Model Pre-Sorts

D. Koop, CSCI 640/490, Spring 2023

(b) Compact & local sort

More...

Query Optimization

Final Exam

- Wednesday, May 10, **8:00**-9:50am, PM 253
- Similar format
- More comprehensive (questions from topics covered in Test 1 & 2)
- Will also have questions from graph/spatial/temporal data, provenance, reproducibility, machine learning

Questions?

D. Koop, CSCI 640/490, Spring 2023

Review

What's involved in dealing with data?

Data	Data	Data	Data	Data
Acquisition	Analysis	Curation	Storage	Usage
 Structured data Unstructured data Event processing Sensor networks Protocols Real-time Data streams Multimodality 	 Stream mining Semantic analysis Machine learning Information extraction Linked Data Data discovery 'Whole world' semantics Ecosystems Community data analysis Cross-sectorial data analysis 	 Data Quality Trust / Provenance Annotation Data validation Human-Data Interaction Top-down/Bottom- up Community / Crowd Human Computation Curation at scale Incentivisation Automation Interoperability 	 In-Memory DBs NoSQL DBs NewSQL DBs Cloud storage Query Interfaces Scalability and Performance Data Models Consistency, Availability, Partition-tolerance Security and Privacy Standardization 	 Decision support Prediction In-use analytics Simulation Exploration Visualisation Modeling Control Domain-specific usage

D. Koop, CSCI 640/490, Spring 2023

[Big Data Value Chain, Curry et al., 2014]

Python!

- Just assign expressions to variables, no typing
 - a = 12
 - a = "abc"
 - b = a + "de"
- Functions defined using def, called using parenthesis:
 - def hello(name1="Joe", name2="Jane"): print(f"Hello {name1} and {name2}") hello(name2="Mary")
- Always indent blocks (if-else-elif, while, for, etc.):

Python Containers

- List: [1, "abc", 12.34]
- Tuple: (1, "abc", 12.34)
- Indexing/Slicing:
 - x[0], x[:-1], x[1:2], x[::2]
- Set: {1, "abc", 12.34}
- Dictionary: {'x': 1, 'y': "abc", 'z': 12.34}
- Mutable vs. Immutable
- Stored by reference
- You cannot index/slice an iterator (d.values() [-1] doesn't work)

Comprehensions

- List Comprehensions:
 - squares = $[i^{*2} \text{ for i in range}(10)]$
- Dictionary Comprehensions:
 - squares = {i: i^*2 for i in range(10) }
- Set Comprehensions:
 - squares = $\{i^{*2} \text{ for } i \text{ in range}(10)\}$
- Comprehensions allow filters:
 - squares = [i**2 for i in range(10) if i % 2 == 0]

JupyterLab

- environment Supports many activities including notebooks • Runs in your web browser • Notebooks: IUDVter - Originally designed for Python - Supports other languages, too - Displays results (even interactive maps) inline - You decide how to divide code into executable cells
 - Shift+Enter to execute a cell

D. Koop, CSCI 640/490, Spring 2023

• An interactive, configurable programming

Relational Algebra

- Six basic operators
 - select: σ
 - project:
 - union: U
 - set difference: -
 - Cartesian product: x
 - rename: p

D. Koop, CSCI 640/490, Spring 2023

Definition: A procedural language consisting of a set of operations that take one or two relations as input and produce a new relation as their result.

Components of SQL

- Data Manipulation Language (DML): provides the ability to query and modify tuples in the database.
- An SQL relation is defined using the create table command:
- A typical SQL query has the form: select A_1, A_2, \ldots, A_n **from** *l*₁, *l*₂, ..., *l*_m where *P*

 Data Definition Language (DDL): the specification of information about relations, including schema, types, integrity constraints, indices, storage

information from the database and to insert tuples into, delete tuples from,

create table $r(A_1 D_1, A_2 D_2, ..., A_n D_n, (C_1), ..., (C_k))$

- A_i is an **attribute**
- D_i is the **data type**
- ri represents a relation
- *P* is a **predicate**

NumPy arrays and slicing

D. Koop, CSCI 640/490, Spring 2023

Expression	Shape
arr[:2, 1:]	(2, 2)
arr[2] arr[2, :] arr[2:, :]	(3,) (3,) (1, 3)
arr[:, :2]	(3, 2)
arr[1, :2] arr[1:2, :2]	(2,) (1, 2)

[W. McKinney, Python for Data Analysis]

Boolean Indexing

- names == 'Bob' gives back booleans that represent the element-wise comparison with the array names
- Boolean arrays can be used to index into another array:
 - data[names == 'Bob']
- Can even mix and match with integer slicing
- Can do boolean operations (&, |) between arrays (just like addition, subtraction)
 - data[(names == 'Bob') | (names == 'Will')]
- Note: or and and do not work with arrays
- We can set values too! data [data < 0] = 0

What is Data?

→ Tables

 \rightarrow Multidimensional Table

D. Koop, CSCI 640/490, Spring 2023

→ Geometry (Spatial)

Northern Illinois University

Categorial, Ordinal, and Quantitative

Α	B	(2	S	Т	U
Order ID	Order Date	Order Priorit	tv	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-1 ow		Large Box	0.8	10/21/06
6	2/21/08	4-Not Speci	fied	Small Pack	0.55	2/22/08
32	7/16/07	2-High	neu	Small Pack	0.79	7/17/07
32	7/16/07	2-High		Jumbo Box	0.72	7/17/07
32	7/16/07	2-High		Medium Box	0.6	7/18/07
32	7/16/07	2-High		Medium Box	0.65	7/18/07
35	10/23/07	4-Not Speci	fied	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Speci	fied	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent		Small Box	0.55	11/3/07
65	3/18/07	1-Urgent		Small Pack	0.49	3/19/07
66	1/20/05	5-Low		Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Speci	fied	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Spec		atitativo	0.6	6/6/05
70	12/18/06	5-Low	quai	Illalive	0.59	12/23/06
70	12/18/06	5-Low	ordi	nol	0.82	12/23/06
96	4/17/05	2-High	UIUI	11a1	0.55	4/19/05
97	1/29/06	3-Medium	cate	onrical	0.38	1/30/06
129	11/19/08	5-Low	cate	Surran	0.37	11/28/08
130	5/8/08	2-High		Small Box	0.37	5/9/08
130	5/8/08	2-High		Medium Box	0.38	5/10/08
130	5/8/08	2-High		Small Box	0.6	5/11/08
132	6/11/06	3-Medium		Medium Box	0.6	6/12/06
132	6/11/06	3-Medium		Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Speci	fied	Large Box	0.82	5/3/08
135	10/21/07	4-Not Speci	fied	Small Pack	0.64	10/23/07
166	9/12/07	2-High		Small Box	0.55	9/14/07
193	8/8/06	1-Urgent		Medium Box	0.57	8/10/06
194	4/5/08	3-Medium		Wrap Bag	0.42	4/7/08
101	1 / = / 0 0	A 11 11				1 (2) (0.0

Re	Inspection Type	Inspection Date	Zip	State	City	Address	Risk	Facility Type	License #	AKA Name	DBA Name	Inspection ID	
Not F	Licens	01/13/2020	60607.0	IL 6	CHICAGO	210 N CARPENTER ST	All	NaN	2709319.0	UNCOOKED LLC	UNCOOKED LLC	2356580	0
	License Re-Inspection	01/13/2020	60602.0	IL 6	CHICAGO	33 N LA SALLE ST	Risk 1 (High)	Restaurant	2689550.0	MOJO 33 NORTH LASALLE LLC	MOJO 33 NORTH LASALLE LLC	2356551	1
Not F	Licens	01/10/2020	60618.0	IL 6	CHICAGO	2949 W BELMONT AVE	Risk 1 (High)	NaN	2708992.0	LA BIZNAGA #2	LA BIZNAGA #2	2356492	2
	Canvas	01/09/2020	60641.0	IL 6	CHICAGO	4920 W IRVING PARK RD	Risk 1 (High)	Restaurant	1617900.0	LAS TABLAS	LAS TABLAS	2356432	3
	Canvas	01/09/2020	60643.0	IL 6	CHICAGO	9613 S WESTERN AVE	Risk 1 (High)	Restaurant	2074456.0	GIORDANO'S OF BEVERLY	GIORDANO'S OF BEVERLY	2356423	4
													•••
	Suspected Food Poisoning	02/18/2010	60604.0	IL 6	CHICAGO	77 W JACKSON BLVD	Risk 1 (High)	Restaurant	1801495.0	PANDA EXPRESS #236	PANDA EXPRESS #236	112321	199687
	Complain	02/08/2010	60615.0	IL 6	CHICAGO	1453 E HYDE PARK BLVD	Risk 1 (High)	Restaurant	81030.0	UNCLE JOE'S	KENNYS RIBS & CHICKEN	74300	199688
	License Re-Inspection	01/28/2010	60630.0	IL 6	CHICAGO	5527-5531 N Milwaukee AVE	Risk 1 (High)	Restaurant	2016764.0	Cafe Marbella	Cafe Marbella	70314	199689
	TASK FORCE LIQUOR 147	02/18/2010	60649.0	IL 6	CHICAGO	7544 S STONY ISLAND AVE	Risk 3 (Low)	Grocery Store	2004292.0	WALGREENS # 07876	WALGREENS # 07876	78309	199690
	License Re-Inspection	01/12/2010	60641.0	IL 6	CHICAGO	4908 W Irving Park RD	Risk 1 (High)	Restaurant	2013419.0	YSABEL'S GRILL ASIAN CUISINE	YSABEL'S FILIPINO CUISINE	150209	199691

199692 rows × 17 columns

• Data Frames are tables with many database-like operations Index shared across all columns # just the beginning of the dataset of Teach Select, project, merge (join), and more Read and write many file formats

D. Koop, CSCI 640/490, Spring 2023

Pass

How do data scientists spend their time?

D. Koop, CSCI 640/490, Spring 2023

What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

[CrowdFlower Data Science Report, 2016]

Data Wrangling

- Automated Transformation Suggestions
- Editable Natural Langua
 - Transform Script Import Export Fill Bangladesh by copying value Split data repeatedly on **newline** into **rows** above

averaging Fill Bangladesh by ✓ copying interpolating values from above

Fill Bangladesh by averaging t values from above

- Visual Transformation Pl
- Transformation History

split	# split1	# split2	# split3	# split4
	2004	2004	2004	2003
ATE	Participation Rate 2004	Mean SAT I Verbal	Mean SAT I Math	Participation Ro
v York	87	497	510	82
necticut	85	515	515	84
ssachusetts	85	518	523	82
v Jersey	83	501	514	85
v Hampshire	80	522	521	75
	77	489	476	77
ine	76	505	501	70
ınsylvania	74	501	502	73
laware	73	500	499	73
orgia	73	494	493	66
split	# fold	Abc fold1	# value	
v York	2004	Participation Rate 2004	87	
v York	2004	Mean SAT I Verbal	497	
v York	2004	Mean SAT I Math	510	
v York	2003	Participation Rate 2003	82	
v York	2003	Mean SAT I Verbal	496	
v York	2003	Mean SAT I Math	510	
necticut	2004	Participation Rate 2004	85	
necticut	2004	Mean SAT I Verbal	515	
nnecticut	2004	Mean SAT I Math	515	
nnecticut	2003	Participation Rate 2003	84	
nnecticut	2003	Mean SAT I Verbal	512	
	2002	Mark CAT T Math		

TDE: Transform Data by Example

С	D
Transaction Date	output
Wed, 12 Jan 2011	2011-01-12-Wednesday
Thu, 15 Sep 2011	2011-09-15-Thursday
Mon, 17 Sep 2012	
2010-Nov-30 11:10:41	
2011-Jan-11 02:27:21	
2011-Jan-12	
2010-Dec-24	
9/22/2011	
7/11/2012	
2/12/2012	

С	D	Transform Data by Example
Transaction Date	output	≡
Wed, 12 Jan 2011	2011-01-12-Wednesday	Show Instruct Get Transformations
Thu, 15 Sep 2011	2011-09-15-Thursday	
Mon, 17 Sep 2012	2012-09-17-Monday	System.DateTime Parse(System.String)
2010-Nov-30 11:10:41	2010-11-30-Tuesday	System.Convert ToDateTime(System.String)
2011-Jan-11 02:27:21	2011-01-11-Tuesday	
2011-Jan-12	2011-01-12-Wednesday	DateFormat.Program Parse(System.String)
2010-Dec-24	2010-12-24-Friday	
9/22/2011	2011-09-22-Thursday	
7/11/2012	2012-07-11-Wednesday	
2/12/2012	2012-02-12-Sunday	© Microsoft Privacy Terms Feedback

Transform by Pattern: Automating Unify/Repair

Auto-Unify

S-timestamp 🗖	S-phone 🗸	S-coordinates 🖉
2019-12-23	(425) 882-8080	(38°57'N, 95°15'W)
2019-12-24	(425) 882-8080	(38°61'N, 95°21'W)
2019-12-23	(206) 876-1800	(39°19'N, 95°18'W)
2019-12-24	(206) 876-1800	(39°26'N, 95°23'W)
2019-12-23	(206) 903-8010	(39°42'N, 96°38'W)
R-timestamp 🚽	R-phone 🖉	R-coordinates 🖉
Nov. 16 2019	650-853-1300	N37°31' W122°14'
Nov. 17 2019	650-853-1300	N37°18' W122°19'
Nov. 16 2019	425-421-1225	N37°48' W122°17'
Nov. 17 2019	425-421-1225	N37°60' W123°08'
Nov. 16 2019	650-253-0827	N37°01' W123°72'

• Auto-Repair

Year	Artist	Issue Price (BU)			
1989	John Mardon	\$16.25			
1990	D.J. Craig	\$16.75			
1991	D.J. Craig	\$16.75			
1992	Karsten Smith	17.50			
1993	Stewart Sherwood	\$17.50			
1994	lan D. Sparkes	\$17.95			
(b) EN-Wiki: Currency values					

D. Koop, CSCI 640/490, Spring 2023

Women's	Time ♦
winner	
Anikó	2:31:24
Kálovics	2.01.21
Lenah	2.22.02
Cheruiyot	2.21.02
Lenah	2.33 11
Cheruiyot	2.33.44
Emily	2.28 12
Kimuria	2.20.42
lana Ekimat	2.32.08
	2.32.00
\mathbf{c}	EN_
\cup_{j}	

wiki:time

#	Original air date ^[1]
12	March 23, 2008
13	March 30, 2008
14	April 6, 2008
15	13 April 2008
16	20 April 2008
(d) E	EN-Wiki: Date

TBP: Learning from Tables

1	Name	#	Born	Died		
-	Washington, George	USA President (1)	02/22/1732	12/14/1799		
	Adams, John	USA President (2), VP (1)	10/30/1735	07/04/1826		
	Jefferson, Thomas	USA President (3), VP (2)	04/13/1743	07/04/1826		
	Madison, James	USA President (4)	03/16/1751	06/28/1836		
	Monroe, James	USA President (5)	04/28/1758	07/04/1851		

Тз

30.	George Washington	-	57y, 10d	22.02.1732	14.12.1799	T ₄	1.	George Washington	Virginia	Feb. 22, 1732	Dec. 14, 1
31.	John Quincy Adams	Nat-Rep	57y, 7m, 20d	11.07.1767	23.02.1848	I	3.	Thomas Jefferson	Virginia	Apr. 13, 1743	July 4, 18
32.	Thomas Jefferson	Dem-Rep	57y, 10m, 18d	13.04.1743	04.07.1826	I	4	James Madison	Virginia	Mar 16 1751	June 28
33.	James Madison	Dem-Rep	57y, 11m, 15d	16.03.1751	28.06.1836	I	6	John Quiney Adama	Maaaabuaatta	hub 41 4707	Eab 02 /
34.	James Monroe	Dem-Rep	58y, 10m, 3d	28.04.1758	04.07.1831		0.	John Quincy Adams	wassachusetts	July 11, 1767	Feb. 23,

	_
	E
_	-

	Name and		State of				Age at	Age at	Т ₆	PRESIDENT	BIRTH DATE	BIRTH PLACE	DEATH DATE	LOCATION OF DEATH
	(party) ¹	Term	birth	Born	Died	Religion ²	inaug.	death		George Washington	Feb 22, 1732	Westmoreland Co., Va.	Dec 14, 1799	Mount Vernon, Va.
1.	Washington (F) ³	1789-1797	Va.	2/22/1732	12/14/1799	Episcopalian	57	67						
2.	J. Adams (F)	1797–1801	Mass.	10/30/1735	7/4/1826	Unitarian	61	90		John Adams	Oct 30, 1735	Quincy, Mass.	July 4, 1826	Quincy, Mass.

l ₂	Date of birth 🔺	President 🗢	Birthplace ¢	State [†] of birth ¢
	February 22, 1732	George Washington	Westmoreland County	Virginia†
	October 30, 1735	John Adams	Braintree	Massachusetts [†]

Tidy Data

	tr	eatmenta	treatmentb			
John Si	mith		2	_		
Jane D	oe	16	11			
Mary J	ohnson	3	1		name	Urt
				-	John Smith	a
	Initi	al Data			Jane Doe	a
					Mary Johnson	a
					John Smith	b
					Jane Doe	b
	John Smith	Jane Do	e Mary Joł	nnson	Mary Johnson	b
nenta		- 1	.6	3	Tidv r	Jata
penth	5) 1	1	1	IIUy L	Jala

	trea	atmenta t	reatmentb				
John	Smith		2				
Jane	Doe	16	11				
Mary	Johnson	3	1		name	trt	result
					John Smith	a	
	Initia	l Data			Jane Doe	a	16
					Mary Johnson	a	3
					John Smith	b	2
					Jane Doe	b	11
	John Smith	Jane Doe	Mary Joh	nson	Mary Johnson	b	1
treatmenta		16		3		$) \rightarrow + \rightarrow$	
treatmentb	2	11		1	TIQY L	ງລເລ	

Transpose











AutoSuggest

- Automate "Complex" Data Preparation steps
- Focus on frame transformations (not per-cell transformations)
- Learn from Jupyter Notebooks
- Two Types of Predictions:
 - Single-Operator Prediction
 - Next-Operator Prediction

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	210.66
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	Business Services	 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 Cleaning: SampleClean









Data Cleaning: 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	60610
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
10 710	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









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







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





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





























D. Koop, CSCI 640/490, Spring 2023

Northern Illinois University

NIU











	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	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

D. Koop, CSCI 640/490, Spring 2023

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	UWisc
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 2023

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









Source Dependence: Iteration on Truth and Sources

Source-accuracy Computation















Source Dependence: Iteration on Truth and Sources

Source-accuracy Computation

Step 3

D. Koop, CSCI 640/490, Spring 2023

Step 2 Truth Discovery

> Dependence Detection

> > Step





Northern Illinois University







NoSQL Motivation

Scalability



D. Koop, CSCI 640/490, Spring 2023

Impedance Mismatch









Column Stores



Each column has a file or segment on disk

D. Koop, CSCI 640/490, Spring 2023

	Person	Genre
oubtfire	Robin Williams	Comedy
	Roy Scheider	Horror
У	Jeff Goldblum	Horror
Magnolias	Dolly Parton	Drama
rdcage	Nathan Lane	Comedy
rokovitch	Julia Roberts	Drama
K	7	

[J. Swanhart, Introduction to Column Stores]









CAP Theorem











Cassandra: Replication and Consistency













Three Types of NewSQL Systems

- New Architectures
 - New codebase without architectural baggage of legacy systems - Examples: VoltDB, Spanner, Clustrix
- Transparent Sharding Middleware:
 - Transparent data sharding & query redirecting over cluster of single-node DBMSs
 - Examples: citusdata, ScaleArc (usually support MySQL/postgres wire)
- Database-as-a-Service:
 - Distributed architecture designed specifically for cloud-native deployment Examples: xeround, GenieDB, FathomDB (usually based on MySQL)









Spanner: Google's NewSQL Cloud Database



- Which type of system is Spanner?
 - C: consistency, which implies a single value for shared data
 - A: 100% availability, for both reads and updates
 - P: tolerance to network partitions
- Which two?
 - CA: close, but not totally available
 - So actually **CP**









Dataframe Data Model



- Typed Row/column labels
 - Labels can become data
- Indexing by label or row/column number
 - "Named notation" or "Positional notation"

- Combines parts of matrices, databases, and spreadsheets
- Ordered, but not necessarily sorted
 - Rows and columns
- No predefined schema necessary
 - Types can be induced at runtime











Differences between Databases & Dataframes



D. Koop, CSCI 640/490, Spring 2023



Incremental + inspection

Mixed types, R/C and data/metadata equiv.

600+ functions









Modin as a Way to Scale Dataframes







Data Science Jungle



D. Koop, CSCI 640/490, Spring 2023



Northern Illinois University



Magpie Goals

































































Gorilla Time Series Data Compression



D. Koop, CSCI 640/490, Spring 2023

Northern Illinois University

NIU






Graph Databases focus on relationships

- Directed, labelled, attributed multigraph
- Properties are key/value pairs that represent metadata for nodes and edges









Graph DBMS Problems

- performance
 - Slow loading speeds
 - Query speeds over magnitude slower than RDBMS
- scalability
 - Low datasize limit, typically << RAM
 - Little benefit from parallelism
- reliability
 - Loads never terminate
 - Query run out of memory or crash
 - Bugs













Interactive Exploration of Spatial Data













Interactive Exploration of Spatial Data













Visualization: Minimize Latency











Visualization: Task-Prioritized Prefetching



D. Koop, CSCI 680/490, Spring 2022





Northern Illinois University



Visualization: Prefetching

- Predict which tiles a user will need next and prefetch those
 - Use common patterns (zoom, pan)
 - Use regions of interest (ROIs)













Spatial Data: Beast Architecture



Spatial Modules



Job Monitoring and Scheduling

RDD Runtime

Storage (HDFS)

From Scratch Approach

(Spatial) User Program + RDD APIs + Job Monitoring and Scheduling + RDD Runtime + Storage + ...











Spatial Data: Partitioning/Indexing & Range Query









Data Curation

The DCC Curation











Data Curation: FAIR Principles

- computers
- Accessible: Users need to know how data can be accessed, possibly including authentication and authorization
- applications or workflows for analysis, storage, and processing

• Findable: Metadata and data should be easy to find for both humans and

• Interoperable: Can be integrated with other data, and can interoperate with

 Reusable: Optimize the reuse of data. Metadata and data should be welldescribed so they can be replicated and/or combined in different settings











Provenance



D. Koop, CSCI 640/490, Spring 2023





73

Prospective and Retrospective Provenance

- Recipe for baking a cake versus the actual process & outcome Prospective provenance is what was specified/intended
- - a workflow, script, list of steps
- Retrospective provenance is what actually happened
 - actual data, actual parameters, errors that occurred, timestamps, machine information
- **Do not need** prospective provenance to have retrospective provenance!



D. Koop, CSCI 640/490, Spring 2023







74

Reproducibility













Machine Learning and Databases















Questions?





Final Exam

- Wednesday, May 10, **8:00**-9:50am, PM 253
- Similar format
- More comprehensive (questions from topics covered in Test 1 & 2)
- Will also have questions from graph/spatial/temporal data, provenance, reproducibility, machine learning





