Advanced Data Management (CSCI 490/680)

Review

Dr. David Koop





Data systems rely on algorithms

DATA SYSTEMS ALGORITHMS















Data structures define performance



D. Koop, CSCI 680/490, Spring 2021



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 680/490, Spring 2021









7

Model to Predict Data's Location on Disk

Frequency Distribution

Cumulative Distribution Function (CDF)



MacMenamin







Challenges



D. Koop, CSCI 680/490, Spring 2021

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 680/490, Spring 2021

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





Sorting

(a) CDF Model Pre-Sorts



D. Koop, CSCI 680/490, Spring 2021

(b) Compact & local sort









12

Sorting

(a) CDF Model Pre-Sorts



D. Koop, CSCI 680/490, Spring 2021

(b) Compact & local sort















More...























<u>Assignment 5</u>

- Four parts
 - Loading Data
 - Spatial Analysis
 - Graph Analysis
 - Temporal Analysis
- Due tomorrow
- Questions?





Final Exam

- Monday, April 26, 4:00-5:50pm, Online (Blackboard)
- Similar format
- More comprehensive (questions from topics covered in Test 1 & 2)
- machine learning

• Will also have questions from temporal data, provenance, reproducibility,





Questions?





D. Koop, CSCI 680/490, Spring 2021

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 680/490, Spring 2021

[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 680/490, Spring 2021

• An interactive, configurable programming









NumPy arrays and slicing



D. Koop, CSCI 680/490, Spring 2021

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

→ Networks





 \rightarrow Multidimensional Table



→ Trees



D. Koop, CSCI 680/490, Spring 2021

→ Geometry (Spatial)







Northern Illinois University









Categorial, Ordinal, and Quantitative

Α	В	(C	S	Т	U
Order ID	Order Date	Order Priori	ty	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low		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		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 Spec	fied	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Spec	ana	atitativa	0.6	6/6/05
70	12/18/06	5-Low	yuai	illative	0.59	12/23/06
70	12/18/06	5-Low	ordi	nal	0.82	12/23/06
96	4/17/05	2-High			0.55	4/19/05
97	1/29/06	3-Medium	cate	gorical	0.38	1/30/06
129	11/19/08	5-Low	cute	5011041	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
	1 / = / 0 0	a				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 680/490, Spring 2021





Pass





How do data scientists spend their time?



D. Koop, CSCI 680/490, Spring 2021

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

\heartsuit	cu	stomer analysis > customer ~ Random		
070		≡ ∽ ∽ ₽ _₿ ∙ ₽	• R ₂ ¹ • + • • € ∃ A •	= -{} ч; ч; ч;
D)		Preview		
		# IMSI ~	○ CONTRACT_END ∨	CONTRACT_START
~			Log 0010 Dec 0016	lan 0000. Dec 0014
(\mathbf{P})		3101 - 310.261	Jan 2013 - Dec 2016	Jan 2000 - Dec 2014
	•	3101/0220812/21	0/4/10	10/6/12
	-	310100900/00/00	3/28/15	10/0/13
		310005432840230	5/20/15	2/14/01
		310026939721905	9/11/15	9/18/10
		310026015466952	8/27/15	3/13/06
		310170484724861	1/16/16	5/11/04
		310170765640471	05-Jul-2011	9/11/06
		310260310245556	12/24/15	3/28/01
		310150834295817	3/6/15	7/26/00
		310160464252516	9/25/15	4/4/04
		310120438750772	4/30/16	9/8/04
		310260195729676	1/16/15	1/3/04
		310026261822880	8/13/13	11/23/08
		310005667082048	8/4/16	10/22/14
		310170836020164	1/22/15	10/19/14
		310160772267782	11/21/15	12/28/14
		310170116249240	27-Sep-2011	2/9/09
		310026110612337	5/29/15	3/29/05
		310260681676970	11/17/16	5/21/07
		310004436630316	9/15/16	7/24/11
		310120423699542	2/27/15	6/29/11
		310120773194729	4/28/16	6/15/04
		310030295859214	2/7/15	3/24/12
?		310012150088547	13-Jan-2009	12/10/05
		310120387060694	10/1/16	10/25/11
D	O	19 Columns 20,000 Row	vs 8 Data Types	











Foofah: Programming by Example









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	

D. Koop, CSCI 680/490, Spring 2021

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









- ×

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

	tre	atmenta t	reatmentb				
John	Smith		2				
Jane	Doe	16	11				1.
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	liay L	ງລເລ	

Transpose











MultiIndex Row Access and Slicing

- df.loc[("Boston", 2007)] or sometimes df.loc["Boston", 2007]
- Remember that loc uses the index values, iloc uses integers
- Note: df.iloc[0] gets the first row, not df.iloc[0,0]
- Can get a subset of the data using partial indices
 - df.loc["Boston"] returns both 2007 and 2008 data
- What about slicing?
 - df.loc["Boston":"Cleveland"] \rightarrow ERROR! (Need sorted data)
 - df = df.sort index()
 - df.loc["Boston":"Cleveland"] → inclusive! - df.loc[(slice("Boston", "Cleveland"), 2007),:]







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 680/490, Spring 2021

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 680/490, Spring 2021

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 680/490, Spring 2021

Northern Illinois University

NIU











Google Dataset Search Overview









Data Curation

The DCC Curation











Computational Data Citation (MODIS)



D. Koop, CSCI 680/490, Spring 2021





45

FAIR Principles

- computers
- Accessible: Users need to know how data can be accessed, possibly including authentication and authorization
- Interoperable: Can be integrated with other data, and can interoperate with applications or workflows for analysis, storage, and processing
- Reusable: Optimize the reuse of data. Metadata and data should be welldescribed so they can be replicated and/or combined in different settings

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









Parallel DB Architecture: Shared Nothing



D. Koop, CSCI 680/490, Spring 2021



[Hellerstein et al., Architecture of a Database System]







Column Stores



Each column has a file or segment on disk

D. Koop, CSCI 680/490, Spring 2021

	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













Spanner: Google's NewSQL Cloud Database



D. Koop, CSCI 680/490, Spring 2021

- 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





51

Graph Databases focus on relationships

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







Interactive Exploration of Spatial Data













Interactive Exploration of Spatial Data













Spatial Data: Minimize Latency









Spatial Data: Precompute Optimized Storage









Spatial Data: Prefetching

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













Split-Apply-Combine



D. Koop, CSCI 680/490, Spring 2021

[W. McKinney, Python for Data Analysis]































































Provenance











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!













Using Provenance



D. Koop, CSCI 680/490, Spring 2021









61

Reproducibility













Machine Learning and Databases















Questions?







Final Exam

- Monday, April 26, 4:00-5:50pm, Online (Blackboard)
- Similar format
- More comprehensive (questions from topics covered in Test 1 & 2)
- machine learning

• Will also have questions from temporal data, provenance, reproducibility,





