### Advanced Data Management (CSCI 490/680)

### Graph Data

Dr. David Koop





### Discussions

- Please post at least once on the discussion board in Blackboard with a question, answer, or discussion point about one of the lectures
- Feedback is useful—I hope the lectures are clear, but I am pretty sure there are still places where I can clarify things better
- You may also post questions about the assignment there if you believe they are relevant to all students









### <u>Assignment 4</u>



#### D. Koop, CSCI 490/680, Spring 2020

- COVID-19 data
- Data Integration
  - Population
  - Temperature
- Data Fusion:
  - Our World in Data
  - Johns Hopkins
  - Wikipedia
- Questions?





## Test 2

- Online on Blackboard (webcourses.niu.edu)
- Thursday, April 9 from 3:30-4:45pm
- If you have conflicts, let me know as soon as possible
- Format:
  - Some multiple choice
  - More short answer/free response
- Focus on topics since the first test
- More details this week





### What is Data?



Less than 0 Change for U.S.: 32,712,033

D. Koop,

- here

de matin I chlorent genin blin nent and delition & MIT it S'en . It fills to philis in the addition of the the men for addition to the the star

new house is o live in ey want him plastering. He para gitti.

{much money went} Has a tractor.

#### Date: July 1980 Place:Sakaltutan Zafor:

Household now Zafor and wife; Nazif Unal and wife and youngest son, still a boy. They run two dolmuß; one with a driver from Süleymanti. Goes in and out once a day. He gets 8,000 a month. Zafor then said, keskin deoil. { not sharp - i.e.? not profitable} I said he did very well on 8,000 TL with only two journeys a day. Nazif Unal has "bought" a Durak (dolmuß stop} from Belediye and works all day in Kayseri.

http://onlineqda.hud.ac.uk/Intro\_QDA/Examples\_of\_Qualitative\_Data.php

Pisa Griffin











### What is data?

- "Data are representations of observations, objects, or other entities used as evidence of phenomena for the purposes of research or scholarship." [C. L. Borgman]
- Data can be digital but can also be physical (e.g. sculptures)
- Semantics are important (e.g. temperature to engineer and biologist)
- Grey Data: surveys, student records—think about privacy









## Sharing Data

- Required/encouraged by universities, funding agencies, publishers
- used to support the arguments." [C. L. Borgman]
- Questions:
  - How is data maintained? Who is responsible?
  - What is the process for curating data?
  - How long should data be kept?
  - How should data collection and curation be acknowledged?

# "Publications are arguments made by authors, and data are the evidence





### Data Curation Lifecycle

#### The DCC Curation Lifecycle Model













## Sequential Actions in Data Curation

- Create or Receive: Create/receive data and make sure metadata exists
- preservation
- Preservation Action: Data cleaning, validation (ensure that data remains) authentic, reliable and usable)
- Store: Store the data in a secure manner adhering to relevant standards Access, Use and Reuse: Make sure is accessible to users and reusers

 Conceptualize: Plan creation of data—capture method and storage options. Appraise and Select: Evaluate data and select for long-term curation and

• Ingest: Transfer data to an archive, repository, data centre or other custodian

Transform: Create new data from the original (migrate formats, subsets, etc.)







## 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







### Findable: DataCite Workflow









## Accessible: DOI to Landing Page with Metadata



### Document citing the data

D. Koop, CSCI 490/680, Spring 2020

Repository housing the data

Data store











### Interoperable: Standard vocabularies

					fairsharing.	org/standards/		_	Ċ				
FAIRsharing.org standards, databases, policies						<b>Q</b> Search all of FAIRs	haring	Standard	s Databases P	Policies Collec	ctions Add/Clair	n Content Stats L	og in or Regist
		Search	Standards	<b>Q</b> Search		Sea	arch	Reset	Advanced				
						Sh	owingree	cords <b>1 - 50</b> o	f 1384.				
View as Table View as Grid Sort by			« 1 2 3	3 4 5	6 7 8 9	10 11 12 13	14	15 16	17 18 19 20	21 22	23 24 25	26 27 28 »	
Name 🗳	Registry	Name	Abbreviation	Type Sub	oject	Domain	Та	axonomy	Related Database	Related Standard	Related Policy	In Collection/Recommendation	n Status
Recommended Records	Å	ABA Adult Mouse Brain	ABA	Standard	Neuroscience	Brain     Gene Expression     Brain Imaging		Mus musculus	NeuroMorpho.Org	None	None	None	ß
Recommended         Associated Publication?         No Publication       Has Publication         Claimed?         No Maintainer       Has Maintainer		Access to Biological Collection Data	ABCD	Standard	Biodiversity & Biology Life Science	None		🖉 Ali	GBIF ALA IPT - GBIF Australia Repository GBIF Spain IPT - GBIF Spain Repository Canadensys IPT - GBIF Canadensys Repository SiB Colombia IPT - GBIF Colombia Repository Plus 1 more	ABCDDNA ABCDEFG	None	TDWG Biodivensity Information Standard	8
Uncertain     Deprecated     In development     Ready		Access to Biological Collection Databases Extended for Geosciences	ABCDEFG	Standard 🥔	Earth Science Geology Paleontology Soil Science	None	6	🗲 All	GeoCASe Data Portal	XML ABCD	None	None	R
Standard Type         Terminology Artifact       771         Model/Format       405         Reporting Guideline       163	Ð	Access to Biological Collection Data DNA extension	ABCDDNA	Standard	Biodiversity Ø Biology	<ul> <li>DNA Sequence Data</li> <li>Experiment Metadata</li> <li>Sequence</li> <li>Deoxyribonucleic Acid</li> <li>Polymerase Chain Reaction</li> <li>Plus 1 more</li> </ul>		🖉 All	GenBank	MOD-CO ABCD	None	TDWG Biodivensity Information Standards	
Metric 30 Identifier Schema 15	Ð	.ACE format	ACE format	Standard 🥏	Life Science	DNA Sequence Data     O     Deoxyribonucleic Acid	ontig Genome	🕈 Ali	None	None	None	None	ß
Show More	ന്ന	AdaLab-meta	ADALAB-META	Standard None	e	None		🖉 Ali	None	None	None	None	R
Domains	đ,	AdaLab ontology	ADALAB	Standard None	9	None		🛷 Ali	None	None	None	None	R
Report   141     Data Transformation   134	E	Adverse Drug Reaction Markup Language	EU-ADR ML	Standard None	3	Adverse Reaction     Electronic Health Record     Disesse    Orug	Farget	Homo sapiens	None	XML	None	None	0
Chemical Entity 131 Phenotype 88 Show More	đ	Adverse Event Reporting entelogy	AERO	Standard 🥔	Biomedical Science Health Science & Medicine Ontology And Terminology	Adverse Reaction     Electronic Health Record		Homo sapiens	ОВО	IAO OGMS OBI	None	None	0

#### D. Koop, CSCI 490/680, Spring 2020





### Reusable: Licensing

- Citation of a dataset is expected as a scholarly norm, not by law
- CC0:
  - "I hereby waive all copyright and related or neighboring rights together with all associated claims and causes of action with respect to this work to the extent possible under the law"
- CC BY: license, not a waiver as CC0
  - "You must give appropriate credit, provide a link to the license, and indicate if changes were made."
- Data Use Agreements (DUA): Used when data are restricted due to proprietary or privacy concerns.





### Reusable: Data Citation & Metrics



D. Koop, CSCI 490/680, Spring 2020





Northern Illinois University



D. Koop, CSCI 490/680, Spring 2020

### Specific Types of Data









### Graphs: Social Networks









### What is a Graph?

• An abstract representation of a set of objects where some pairs are connected by links.



Object (Vertex, Node)



Link (Edge, Arc, Relationship)







### What is a Graph?



D. Koop, CSCI 490/680, Spring 2020

 In computing, a graph is an abstract data structure that represents set objects and their relationships as vertices and edges/ links, and supports a number of graphrelated operations

- Objects (nodes): {A,B,C,D}
- Relationships (edges):

   {(D,B), (D,A), (B,C), (B,A), (C,A)}
- Operation: shortest path from  ${}_{\rm D}$  to  ${}_{\rm A}$







## Different Kinds of Graphs

- Undirected Graph
- Directed Graph
- Pseudo Graph
- Multi Graph
- Hyper Graph













## Graphs with Properties

- Each vertex or edge may have properties associated with it
- May include identifiers or classes











## Types of Graph Operations

- Connectivity Operations:
  - number of vertices/edges, in- and out-degrees of vertices
  - histogram of degrees can be useful in comparing graphs
- Path Operations: cycles, reachability, shortest path, minimum spanning tree
- Community Operations: clusters (cohesion and separation)
- Centrality Operations: degree, vulnerability, PageRank
- Pattern Matching: subgraph isomorphism
  - can use properties
  - useful in fraud/threat detection, social network suggestions







## What is a Graph Database?

- A database with an explicit graph structure
- Each node knows its adjacent nodes
- the same
- Plus an Index for lookups

D. Koop, CSCI 490/680, Spring 2020

### • As the number of nodes increases, the cost of a local step (or hop) remains













### How do Graph Databases Compare?











### Graph Databases Compared to Relational Databases

#### Optimized for aggregation



#### D. Koop, CSCI 490/680, Spring 2020

#### Optimized for connections













### Graph Databases Compared to Key-Value Stores

#### Optimized for simple look-ups



D. Koop, CSCI 490/680, Spring 2020

#### Optimized for traversing connected data











### Graph Databases Compared to Document Stores

#### Optimized for "trees" of data





D. Koop, CSCI 490/680, Spring 2020

Optimized for seeing the forest and the trees, and the branches, and the trunks











### Graph Databases

### D. Lembo and R. Rosati





## Why Graph Database Models?

- Graphs has been long ago recognized as one of the most simple, natural and intuitive knowledge representation systems
- Graph data structures allow for a natural modeling when data has graph structure
- Queries can address direct and explicitly this graph structure
- Implementation-wise, graph databases may provide special graph storage structures, and take advantage of efficient graph algorithms available for implementing specific graph operations over the data













### **Relational Model**

				_
NAME	LASTNAME	I	PERSON	PARE
George	Jones		Julia	Georg
Ana	Stone		Julia	Ana
Julia	Jones		David	James
James	Deville		David	Julia
David	Deville		Mary	James
Mary	Deville		Mary	Julia
	1		-	l

#### D. Koop, CSCI 490/680, Spring 2020



#### [R. Angles and C. Gutierrez, 2017]



Northern Illinois University





## Basic Labeled Model (Gram)

- Directed graph with nodes and edges labeled by some vocabulary • Gram is a directed labeled multigraph
- - Each node is labeled with a symbol called a type
  - Each edge has assigned a label representing a relation between types







## Hypergraph Model (Groovy)

- nodes
- dependencies (directed), object-ID and (multiple) structural inheritance



D. Koop, CSCI 490/680, Spring 2020

#### Notion of edge is extended to hyperedge, which relates an arbitrary set of

# • Hypergraphs allow the definition of complex objects (undirected), functional





### Hypernode Model

- hypernodes), allowing **nesting** of graphs
- Encapsulates information



#### D. Koop, CSCI 490/680, Spring 2020

# Hypernode is a directed graph whose nodes can themselves be graphs (or

Northern Illinois University





## Semistructured (Tree) Model: (OEM Graph)

- "Self-describing" data like JSON and XML
- OEM uses pointers to data in the tree



#### D. Koop, CSCI 490/680, Spring 2020

Northern Illinois University





## RDF (Triple) Model

- Schema and instance are mixed together
- SPAQL to query
- Semantic web



D. Koop, CSCI 490/680, Spring 2020

# Interconnect resources in an extensible way using graph-like structure for data









## Property Graph Model (Cypher in neo4j)

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







## Types of Graph Queries

- Adjacency queries (neighbors or neighborhoods)
- Pattern matching queries (related to graph mining)
  - Graph patterns with structural extension or restrictions
  - Complex graph patterns
  - Semantic matching
  - Inexact matching
  - Approximate matching
- Reachability queries (connectivity)

#### D. Koop, CSCI 490/680, Spring 2020





Northern Illinois University





## Types of Graph Queries (continued)

- Analytical queries
  - Summarization queries
  - Complex analytical queries (PageRank, characteristic path length,

# connected components, community detection, clustering coefficient)









### Graph Query Languages



D. Koop, CSCI 490/680, Spring 2020

#### [R. Angles and C. Gutierrez, 2017]



Northern Illinois University









## Sypher

- Implemented by neo4j system
- Expresses reachability queries via path expressions -p = (a) - [:knows\*] -> (b): nodes from a to b following knows edges • START x=node:person(name="John")
- MATCH  $(x) [:friend] \rightarrow (y)$ RETURN y.name

#### D. Koop, CSCI 490/680, Spring 2020





Northern Illinois University



## SPARQL (RDF)

- Uses SELECT-FROM-WHERE pattern like SQL
- SELECT ?N FROM <http://example.org/data.rdf> WHERE { ?X rdf:type voc:Person . ?X voc:name ?N }

#### D. Koop, CSCI 490/680, Spring 2020





Northern Illinois University



## Comparing Graph Database Systems: Features

#### Data Storage

Graph	Main	External	Backend	Indexes
Database	memory	memory	Storage	
AllegroGraph	•	•		•
DEX		•		•
Filament	•		•	
G-Store		•		
HyperGraphDB	•	•	•	•
InfiniteGraph		•		•
Neo4j	•	•		•
Sones				•
vertexDB		•	•	

D. Koop, CSCI 490/680, Spring 2020

#### **Operations/Manipulation**

	Data	Data	Query	API	GU
Graph	Definition	Manipulat.	Language		
Database	Language	Language			
AllegroGraph	•	•	•	•	
DEX					
Filament					
G-Store	•		•		
HyperGraphDB					
InfiniteGraph					
Neo4j					
Sones	•				
vertexDB					











### Comparing Graph Database Systems: Representation

#### Graph Data Structures

	Graphs				Nodes		Edges		
Graph Database	Simple graphs	Hypergraphs	Nested graphs	Attributed graphs	Node labeled	Node attribution	Directed	Edge labeled	Edge attribution
AllegroGraph	•				•		•	•	
DEX					•	•	•	•	•
Filament	•				•		•	•	
G-Store	•				•		•	•	
HyperGraphDB		•			•		•	•	
InfiniteGraph					•	•	•	•	•
Neo4j							•	•	
Sones		•							
vertexDB	•						•	•	

#### D. Koop, CSCI 490/680, Spring 2020

#### Entites & Relations

	S	Schem	a	Instance					
Graph Database	Node types	Property types	Relation types	Object nodes	Value nodes	Complex nodes	Object relations	Simple relations	Complex relations
AllegroGraph									
DEX	•		•		•		•	•	
Filament					•			•	
G-Store					•				
HyperGraphDB	•		•		•			•	•
InfiniteGraph	•		•		•		•		
Neo4j				•	•			•	
Sones					•			•	•
vertexDB									









### Comparing Graph Database Systems: Queries

#### Query Support

	Туре			Use			
Graph Database	Query Lang.	API	Graphical Q. L.	Retrieval	Reasoning	Analysis	
AllegroGraph	0		•		•		
DEX				•			
Filament		•		•			
G-Store				•			
HyperGraphDB				•			
InfiniteGraph				•			
Neo4j	0	•		•			
Sones			•	•			
vertexDB							

#### D. Koop, CSCI 490/680, Spring 2020

#### Types of Queries

	Adjacency		Reachability				
Graph Database	Node/edge adjacency	k-neighborhood	Fixed-length paths	Regular simple paths	Shortest path	Pattern matching	Summarization
Allegro	•						
DEX				●	•		
Filament			•				
G-Store					•		
HyperGraph							
Infinite			•	•	●		
Neo4j					●		
Sones							
vertexDB							

[R. Angles, 2012]





### Reminder: Discussions

- Please post at least once this week on the discussion board in Blackboard with a question, answer, or discussion point about one of the lectures
- Feedback is useful—I hope the lectures are clear, but I am pretty sure there are still places where I can clarify things better
- You may also post questions about the assignment there if you believe they are relevant to all students

#### D. Koop, CSCI 490/680, Spring 2020



