

INTRODUCTION TO KNOWLEDGE GRAPHS

Yogesh Haribhau Kulkarni



Outline

- ① FUNDAMENTALS OF ENTERPRISE KNOWLEDGE GRAPH
- ② INTRODUCTION TO SEMANTICS: TAXONOMY, ONTOLOGY, KNOWLEDGE GRAPHS
- ③ GRAPH THEORY FUNDAMENTALS
- ④ POPULAR GRAPH ALGORITHMS AND APPLICATIONS
- ⑤ REFERENCES

Introduction

Need Answers

- ▶ **Why** do you build a Knowledge Graph (KG)?
- ▶ **What** does it look like?
- ▶ **How** to create it and use it?

Intuition

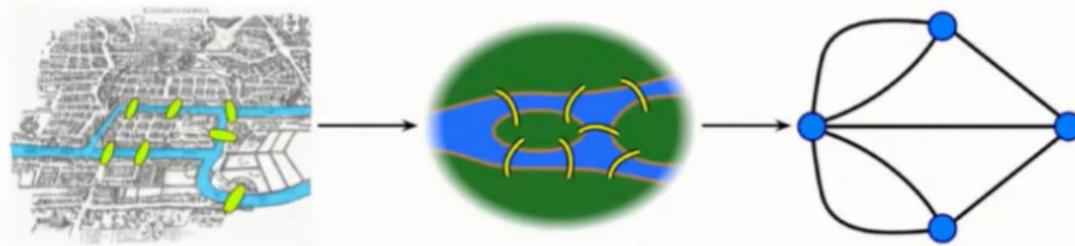
- ▶ Graph: Nodes connected by Edges
- ▶ Knowledge Graph: Entities connected by Relations.

A Graph is

... a set of discrete objects, each of which has some set of relationships with the other objects

Euler: Can we take a walk to all 4 islands, without crossing any of the bridge twice?

Abstraction (Does size of islands matter?):

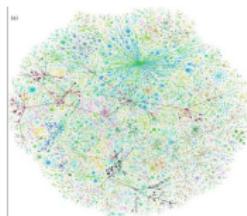


Seven Bridges of Königsberg problem. Leonhard Euler, 1735

Solution: No way!! What's the rule?

(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

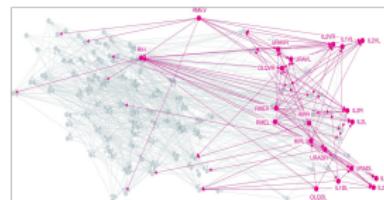
Graph-structured Data Are Ubiquitous



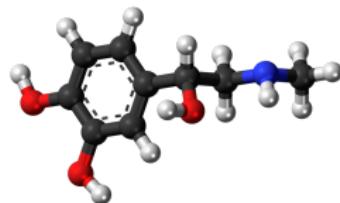
Internet



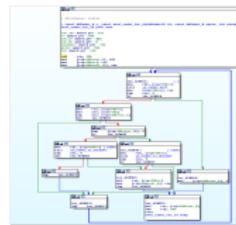
Social networks



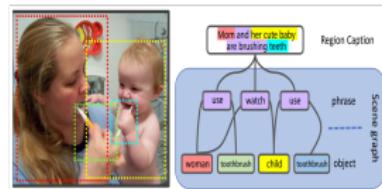
Networks of neurons



Biomedical graphs



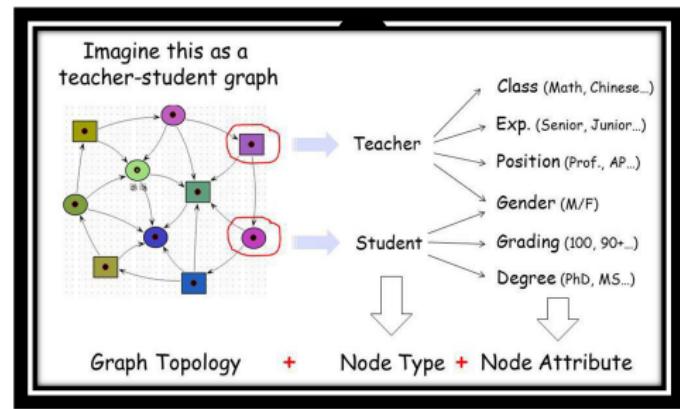
Program graphs



Scene graphs

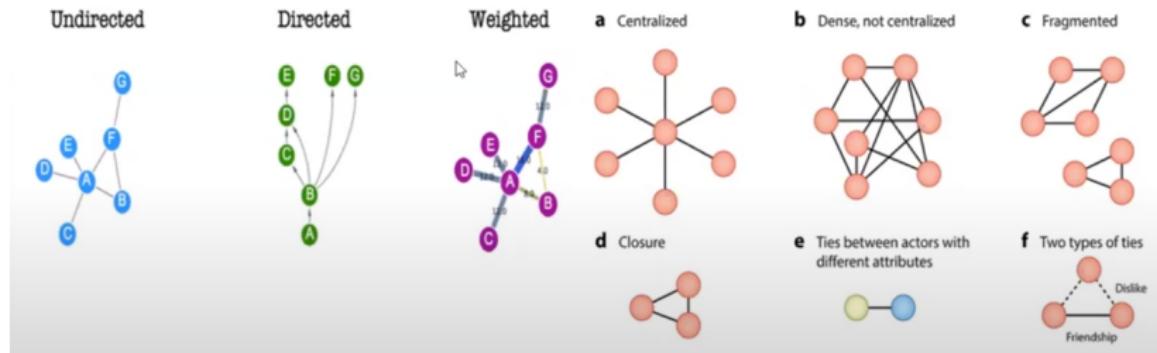
Graphs: A Universal Language

Graphs are a general language for describing and modeling complex systems



Graph!

Types of Graphs

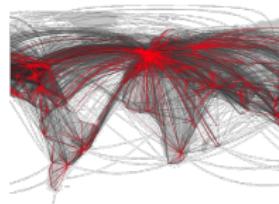


(Ref: Understanding Graph Data Science — DataHour by Vipin K.)

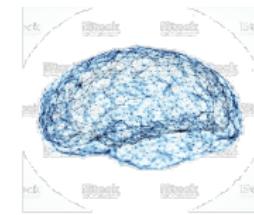
Data as Graphs - Explicit



Social Graphs



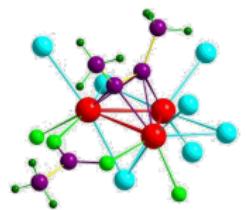
Transportation Graphs



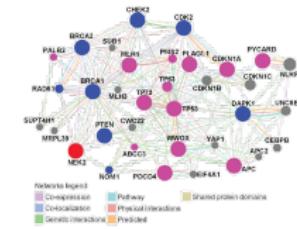
Brain Graphs



Web Graphs

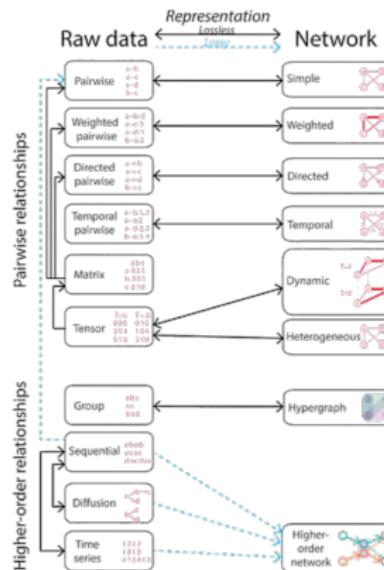
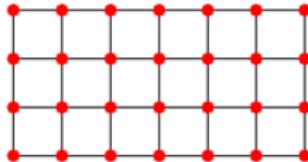


Molecular Graphs



Gene Graphs

Data as Graphs - Implicit



Jian Xu. Representing Big Data as Networks. PhD Dissertation,
University of Notre Dame

Graph Applications

Across an organization, every department can benefit from graphs to answer questions like who or what is important, what should I do next, and what's unusual about this?

Finance

- Fraud Detection
- Pricing Analysis
- Budgeting
- Forecasting

Marketing

- Customer 360
- Influencer Strategy
- Campaign Optimization
- Product Recommendations

Ops

- Product Development
- Pipeline Acceleration
- Supply Chain Optimization
- Infrastructure Planning

IT

- Network Monitoring
- Cybersecurity
- DevOps

HR

- Training
- Upskilling & Retention
- Promotions

(Ref: 5 Graph Data Science Basics Everyone Should Know - neo4j)

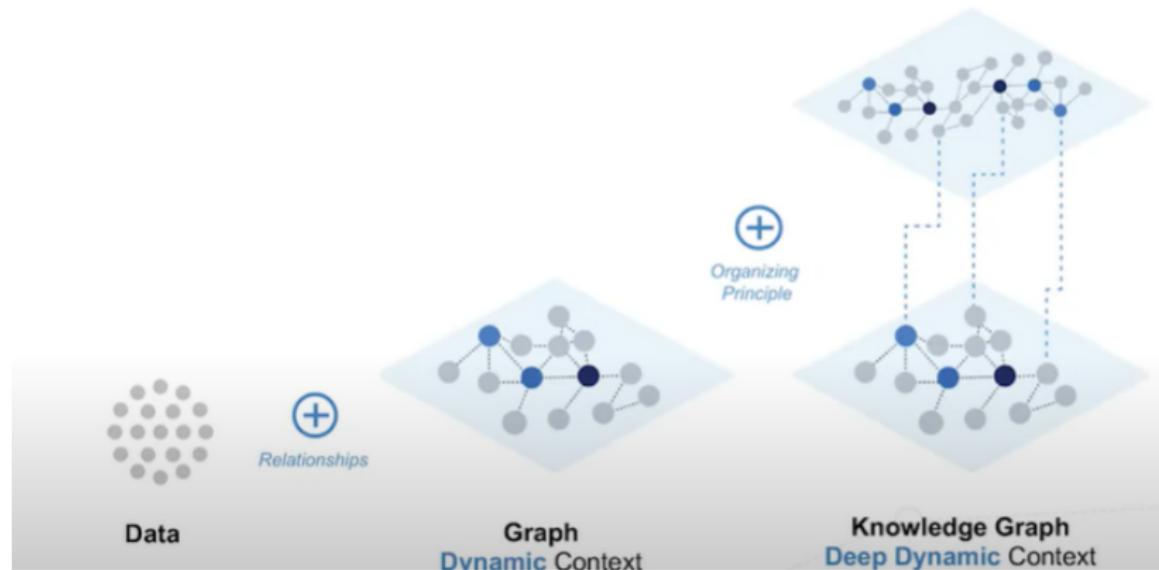
A Knowledge Graph is

...an interconnected dataset enriched with meaning so we can reason about the underlying data and use it confidently for complex decision making

- A neo4j definition

A progression

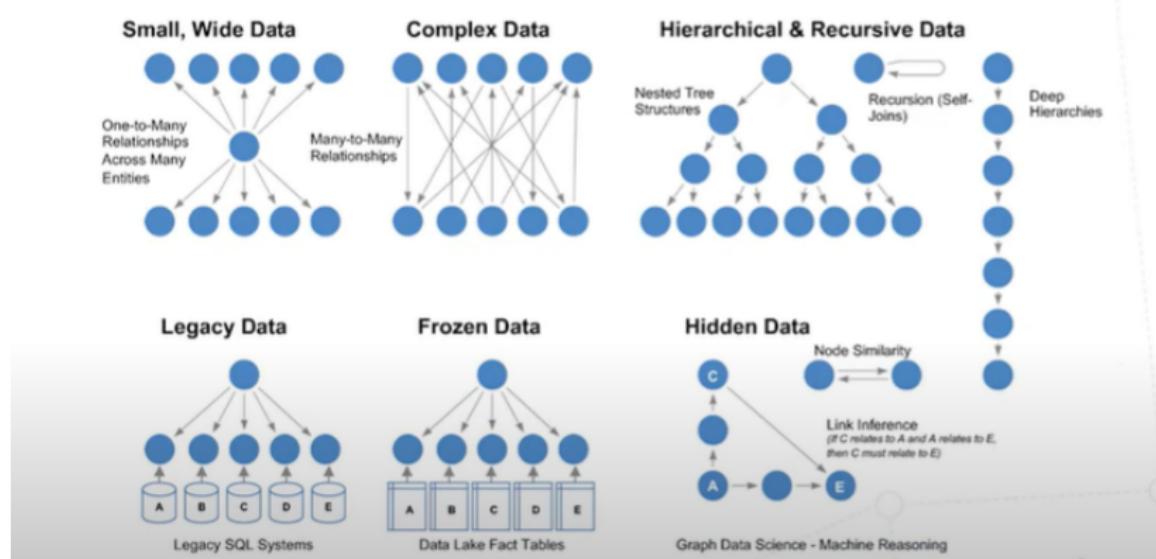
of more enrichment . . .



(Ref: A Universe of Knowledge Graphs - neo4j)

Data

Variety of type/organization of data



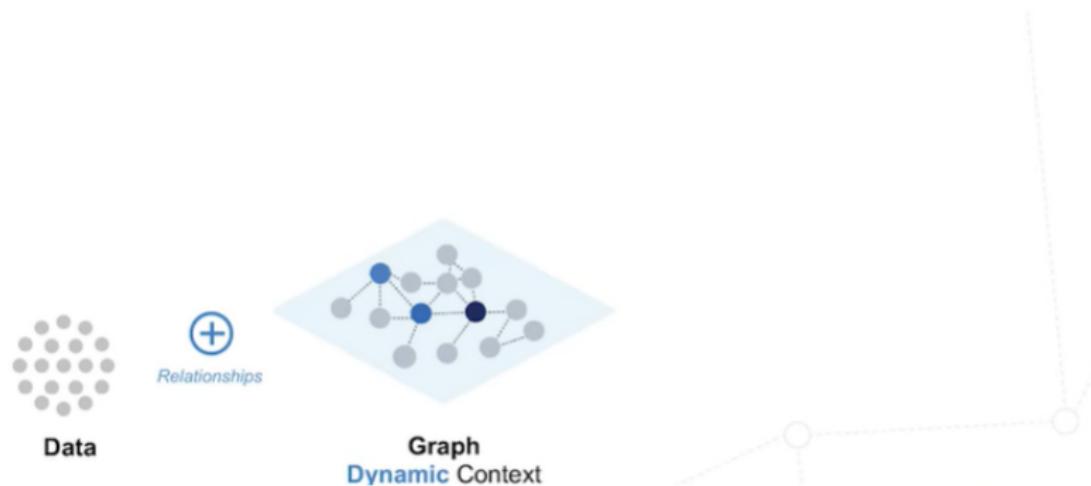
(Ref: A Universe of Knowledge Graphs - neo4j)

"No more a Big Data but need Small and Wide data" for more context in Machine Learning.



Graphs

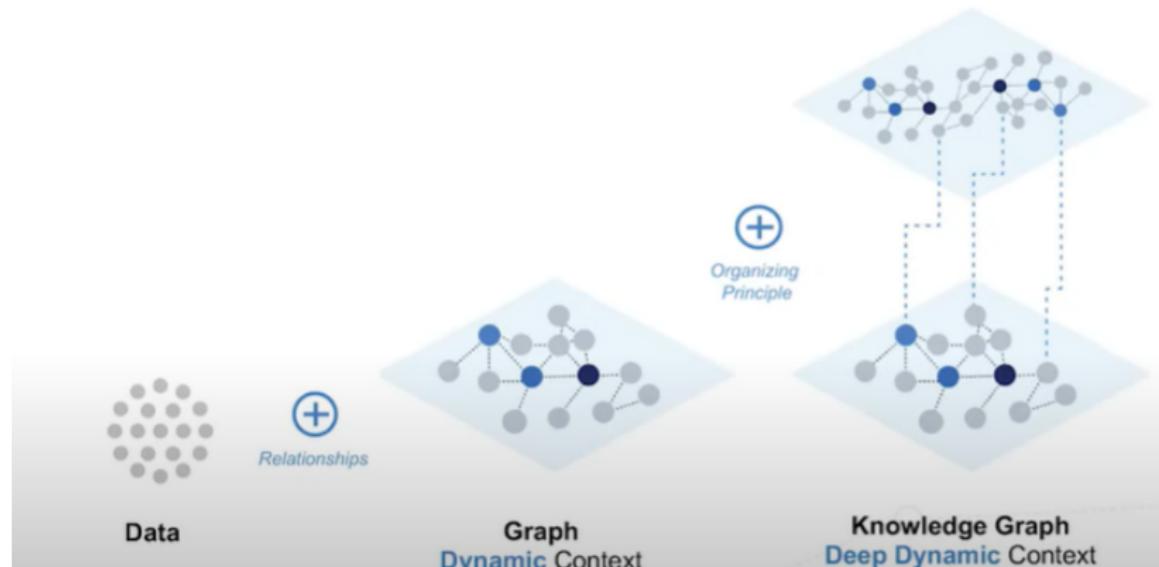
add Relationships to Data ...



(Ref: A Universe of Knowledge Graphs - neo4j)

Knowledge Graphs

add Organizing principles to Graphs . . .



(Ref: A Universe of Knowledge Graphs - neo4j)

e.g. external knowledge, ontologies injection. or 'Semantics'

Semantics

Semantics = Context = Domain Knowledge.

Controlled Vocabularies

Synonym Rings

Taxonomy

Thesaurus

Ontology

C O M P L E X I T Y

Entity Resolution
& Analysis

Classification
& Tagging

Artificial
Intelligence



(Ref: A Universe of Knowledge Graphs - neo4j)

So, a Knowledge Graph is

Knowledge in Graph form.



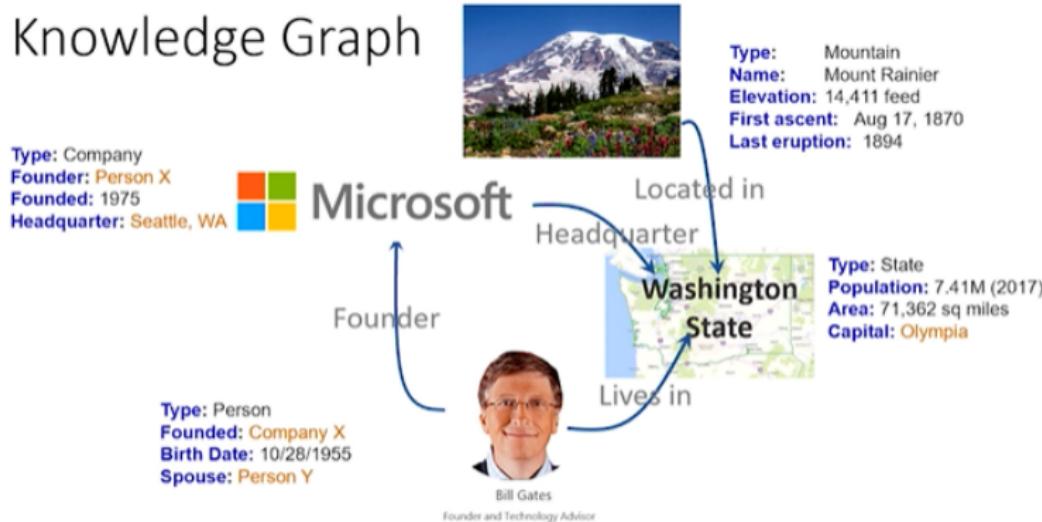
Embedded knowledge is:

- ▶ Microsoft is headquartered in Washington State.
- ▶ Bill Gates is a (co)Founder of Microsoft
- ▶ etc.
- ▶ Entities: Microsoft, Washington State, etc
- ▶ Relationships: founder, head, etc.

(Ref: DAT278x - From Graph and Knowledge Graph - EdX course)

A Knowledge Graph is

Knowledge Graph



- ▶ Entities can be locations, companies, persons, etc. Nouns
- ▶ Relationships can actions, has-a, etc. Verbs.

(Ref: DAT278x - From Graph and Knowledge Graph - EdX course)

Knowledge Graph Datasets



Google
Knowledge Graph

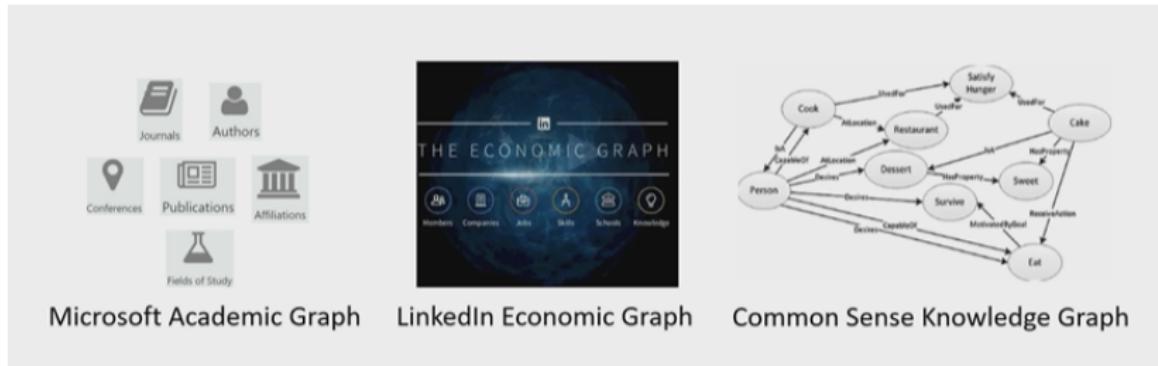
Satori Microsoft
Knowledge Graph



- ▶ Google, Microsoft KGs are private and used for search, question answering
- ▶ DBPedia, Wikidata KGs are public

(Ref: DAT278x - From Graph and Knowledge Graph - EdX course)

Domain specific Knowledge Graphs



- ▶ Microsoft Academic Graph has 170 million papers and more than 200 million authors
- ▶ LinkedIn Economic Graph has 500 million members and 20 million companies.

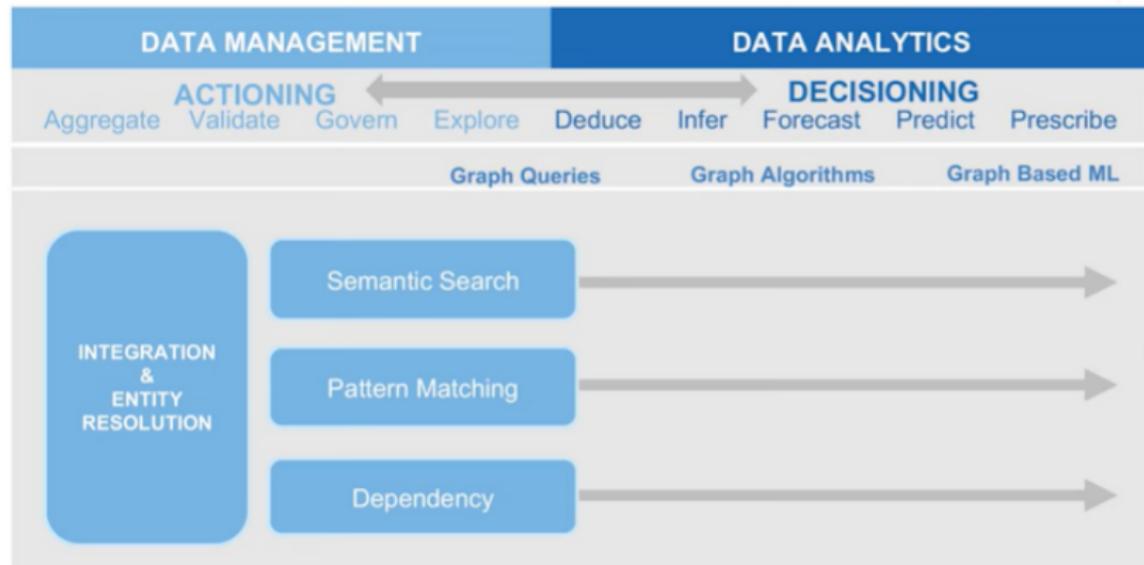
(Ref: DAT278x - From Graph and Knowledge Graph - EdX course)

Why knowledge graph is important?

- ▶ Help organize information
- ▶ Tackle information overload
- ▶ Intuitive explanation, visualization
- ▶ Supports easy querying and business decisions.
- ▶ Key component in many AI applications like Chatbot

(Ref: DAT278x - From Graph and Knowledge Graph - EdX course)

Applications

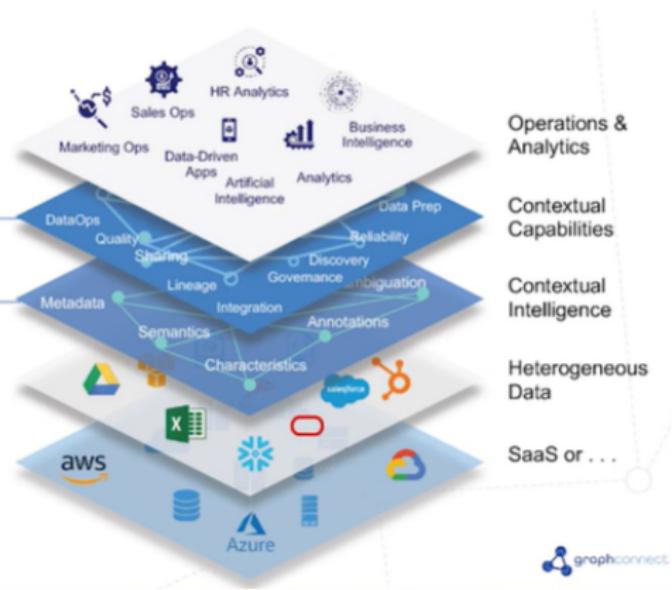
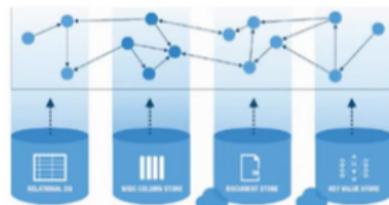


(Ref: A Universe of Knowledge Graphs - neo4j)

Applications

Bridge data Silos

Cross-Silo Connections



© 2022 Neo4j, Inc. All rights reserved.

(Ref: A Universe of Knowledge Graphs - neo4j)

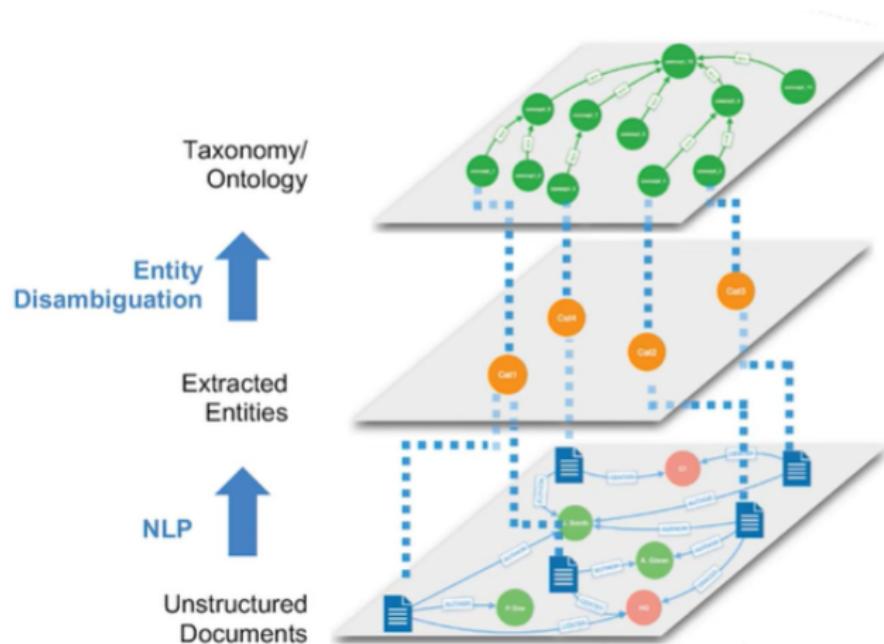
Applications

- ▶ Entity recommendation (Google right hand box): Searching via associated keywords
- ▶ Semantic Search and recommendations using click logs
- ▶ Personal Assistant like Alexa, Siri, using information associated with you.

(Ref: DAT278x - From Graph and Knowledge Graph - EdX course)



Semantic Search

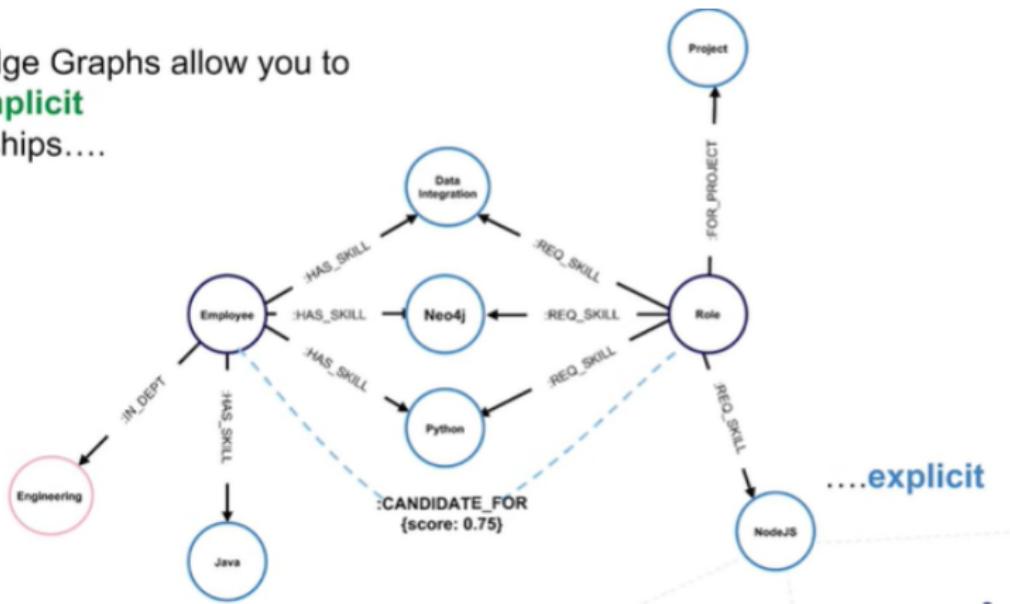


(Ref: A Universe of Knowledge Graphs - neo4j)

Implicit Relationships

Candidate to Job matching via skills

Knowledge Graphs allow you to make **implicit** relationships....



(Ref: A Universe of Knowledge Graphs - neo4j)

Hidden Insights



6,100 employees
1 million calls
500,000 claims

(Ref: A Universe of Knowledge Graphs - neo4j)

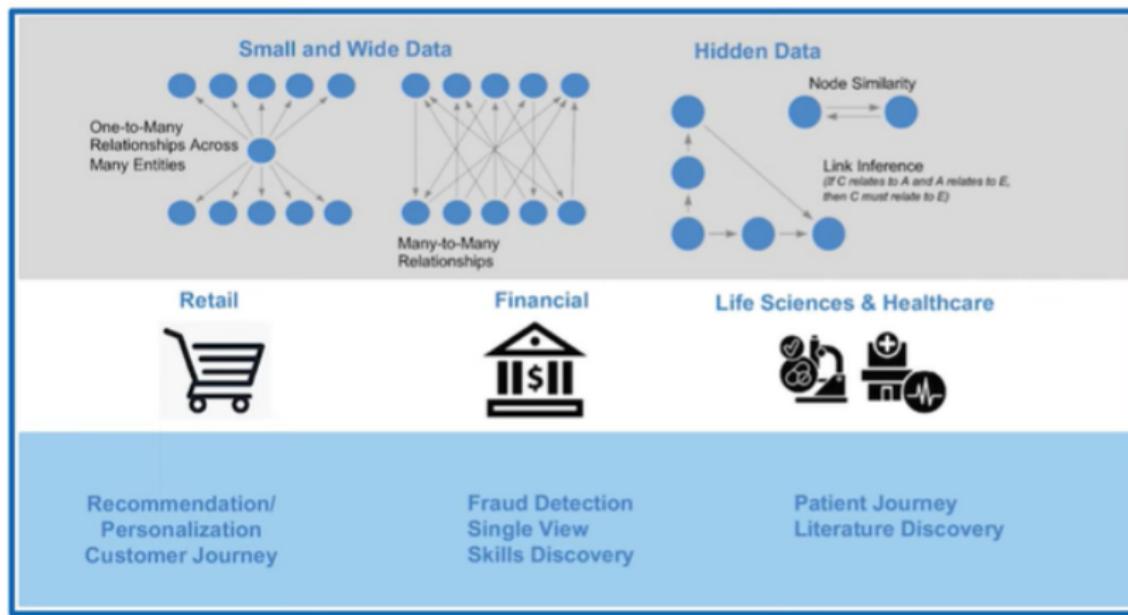
Pattern Matching = Hidden Value

Identify
Unknown Transformative Patterns



Uncover New Answers

Usages



(Ref: A Universe of Knowledge Graphs - neo4j)

Criticality



30,000 components
4,000 suppliers
1,800 parts
19 hours assembly time / car
1,450 cars / day

Small Dependency = Big Impact

Single Component Manufacturer
Shuts Down



Entire Supply Chain Affected

(Ref: A Universe of Knowledge Graphs - neo4j)

Semantics

What is Semantics?

- ▶ Semantics is the study of meaning in language and symbols.
- ▶ In the context of Knowledge Graphs, semantics provides context and understanding to data.
- ▶ It helps bridge the gap between human understanding and machine processing.

SEMANTIC ENCODING

Semantic encoding refers to the process of converting sensory input into meaningful representations in memory.

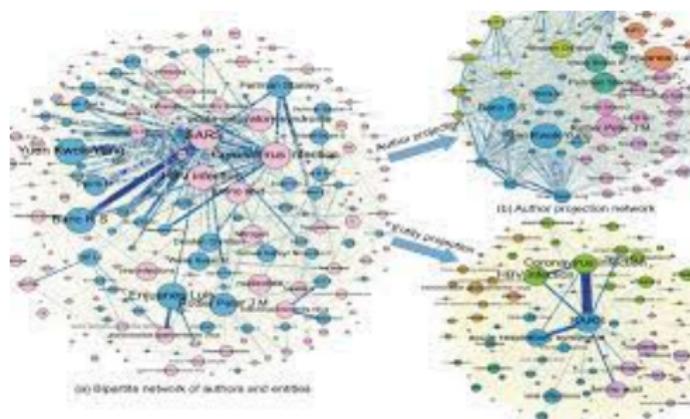
DEFINITION	EXAMPLE
<p>According to Bedford (2020),</p> <p>"...semantic encoding is the processing and encoding of sensory input that has a particular meaning, or that can be applied to a context" (p. 122).</p>	<p>A semantic encoding example is when you meet a new person and to remember their name, you tell yourself that this is the third person named "Ann" you've met this year – this simple story acts as an anecdote you can use during memory recall.</p>

HELPFULPROFESSOR.COM

(Ref: Semantic Encoding: 10 Examples and Definition (2023))

The Power of Knowledge Graphs

- ▶ Knowledge Graphs are powerful data structures representing information as nodes and edges.
 - ▶ They capture complex relationships and connections between entities.
 - ▶ Semantics enrich Knowledge Graphs by attributing meaning to these entities and relationships.
 - ▶ This deepens data understanding and enables advanced analysis and decision-making.



(Ref: PubMed Knowledge Graph – AI Health Lab)

Context Matters

- ▶ In Knowledge Graphs, semantics add context to data.
- ▶ Consider a person named "Suresh" - semantics help distinguish if it's "Suresh Sharma" or "Suresh Varma."
- ▶ By understanding context, insurers gain insights into customer preferences, risks, and more.
- ▶ Semantics enable personalized customer experiences and targeted insurance products.

Semantic Triplets

- ▶ Knowledge Graphs use semantic triplets: Subject - Predicate - Object.
- ▶ Subject: The entity about which information is stored.
- ▶ Predicate: The relationship between the subject and the object.
- ▶ Object: The value or entity related to the subject.
- ▶ For example, (Ramesh, hasOccupation, Software Engineer) is a semantic triplet.

DATA is stored in Triples, expressed as:

Subject	:	Predicate	:	Object
John Smith		livesIn		London
London		isIn		England

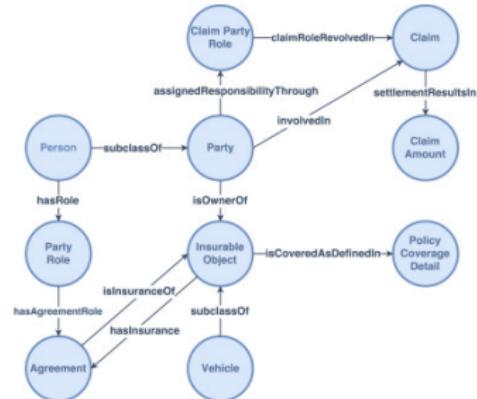
QUERY with SPARQL, gives us a simple look up... and more! Find people who live in (a place that's in) England



(Ref: Semantic Triple - SEO North)

Unleashing the Potential

- ▶ By leveraging semantics, insurers unlock hidden insights within their data.
- ▶ Semantics facilitate data integration, enriching the Knowledge Graph with diverse sources.
- ▶ Complex queries become manageable, leading to faster, data-driven decisions.
- ▶ The insurance industry can provide innovative products tailored to specific customer needs.



(Ref: A standards-based ontology and support for Big

Data Analytics in the insurance)

Real-world Applications

- ▶ Customer Profiling: Semantics enable comprehensive customer 360 views for personalized interactions.
- ▶ Risk Assessment: Understand risk factors by connecting relevant data points.
- ▶ Fraud Detection: Uncover hidden patterns and detect fraudulent activities.
- ▶ Claims Management: Improve efficiency with insights from interconnected data.

Challenges & Considerations

- ▶ Data Quality: Garbage in, garbage out - accurate data is crucial for meaningful semantics.
- ▶ Ontology Design: Thoughtful ontology creation ensures effective representation of concepts and relationships.
- ▶ Scalability: As data grows, the Knowledge Graph must handle increased complexity efficiently.
- ▶ Privacy and Security: Safeguarding sensitive data is paramount in the insurance industry.

Embracing the Future

- ▶ The adoption of Knowledge Graphs with semantics is transforming the insurance industry.
- ▶ Insurers embracing this technology gain a competitive advantage and drive innovation.
- ▶ Stay ahead of the curve by investing in semantic technologies and data-driven approaches.
- ▶ Embrace the potential of semantics to deliver better customer experiences and operational efficiency.

Knowledge Graphs in Action

- ▶ Let's take a look at a practical application of Knowledge Graph Semantics in insurance.
- ▶ Explore how semantics enhance data understanding and drive meaningful insights.
- ▶ Witness firsthand how insurers are benefiting from this transformative technology.



Taxonomy

Taxonomy in Semantics

- ▶ Taxonomy is a hierarchical classification system that organizes entities based on their characteristics and relationships.
- ▶ In Knowledge Graphs, taxonomy helps define the structure and categorization of information, improving data consistency and accuracy.
- ▶ By using taxonomy, insurers can create a clear and logical framework for their knowledge base, simplifying data integration and analysis.

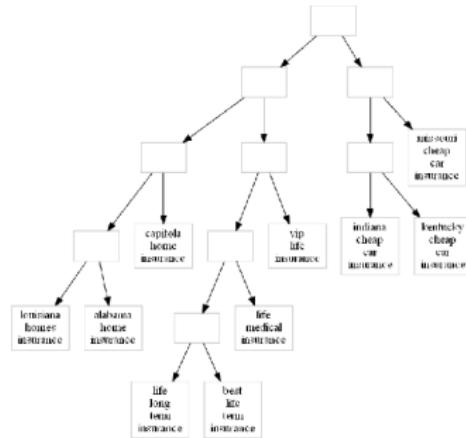


Figure 1: Example of a binary-branch taxonomy.

(Ref: Automatic taxonomy construction from
keywords)

Advantages of Taxonomy in Insurance

- ▶ Seamless Data Integration: Taxonomy ensures data from various sources can be integrated and understood uniformly.
- ▶ Faster Decision Making: Well-defined taxonomy enables quicker access to relevant information, leading to faster and more accurate decisions.
- ▶ Improved Customer Experience: Enhanced organization and retrieval of data enable a better understanding of customer needs, leading to personalized services.

Building a Taxonomy

- ▶ Identify Key Entities: Start by identifying crucial entities in your insurance domain, such as policies, claims, beneficiaries, etc.
- ▶ Define Relationships: Establish meaningful relationships between entities, like 'owns,' 'insured by,' 'dependent on,' etc.
- ▶ Hierarchical Structure: Organize entities into a hierarchical structure to create a clear taxonomy.
- ▶ Flexibility and Scalability: Keep the taxonomy flexible to adapt to evolving needs and ensure it can scale with growing data.

Example of Taxonomy

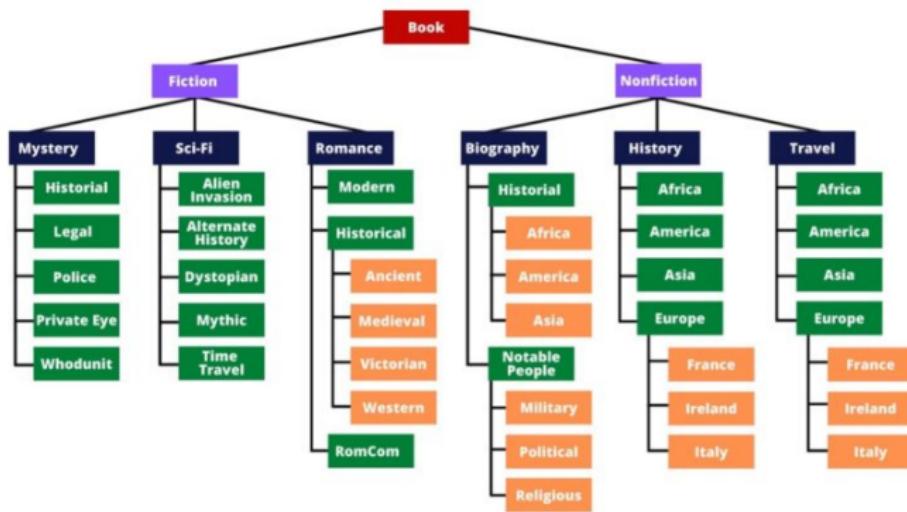


FIGURE: Example Taxonomy for Books

There are no absolute right and wrong with taxonomies, just degrees of appropriateness. The most important question to ask when creating a taxonomy is, “does this hierarchical grouping meet my needs?”

Knowledge Graph and Taxonomy Integration

- ▶ A Knowledge Graph can leverage the defined taxonomy to enhance its structure and relationships.
- ▶ Taxonomy acts as a backbone for organizing information within the Knowledge Graph, making it more intuitive to navigate.
- ▶ Queries and Insights: Taxonomy-driven Knowledge Graphs enable more precise queries and provide deeper insights into relationships between entities.

Use Cases in the Insurance Industry

- ▶ Claims Processing: Streamlining claims processing through efficient taxonomy-based data retrieval.
- ▶ Customer Service: Enhancing customer service by understanding customer profiles and offering tailored insurance products.
- ▶ Fraud Detection: Improving fraud detection by identifying suspicious patterns and connections between entities.
- ▶ Risk Assessment: Facilitating risk assessment by analyzing correlations between policyholders and claims history.

Conclusion

- ▶ Taxonomy in Semantics plays a pivotal role in shaping powerful Knowledge Graphs for the insurance industry.
- ▶ It provides a structured and standardized approach to organizing information, enabling better decision-making and customer experiences.
- ▶ Embracing taxonomy-driven Knowledge Graphs empowers insurers to stay competitive and agile in an ever-evolving industry landscape.

Ontology

Overview

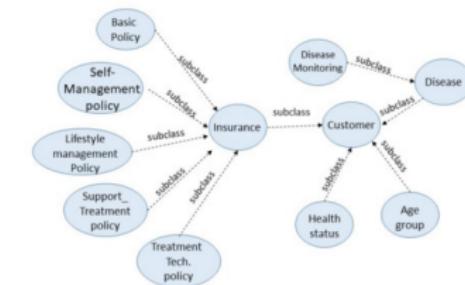
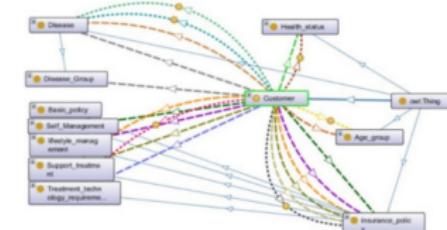
- ▶ Ontology: onto (existence, real) + logia (study) : most essential existence representation
- ▶ According to Wikipedia, an ontology “encompasses a representation, formal naming, and definition of the categories, properties, and relations between the concepts, data, and entities that substantiate one, many, or all domains of discourse.”
- ▶ In other words, ontologies allow us to organize the jargon of a subject area into a controlled vocabulary, thereby decreasing complexity and confusion.
- ▶ Without ontologies, you have no frame of reference, and understanding is lost.
- ▶ For example, an ontology will allow one to associate the Book taxonomy with the Customer taxonomy via relationships.
- ▶ An ontology is more challenging to create than a taxonomy because it needs to capture the interrelationships between business objects/concepts by encapsulating the language and terminology of the business area you are modeling.
- ▶ Ontologies take taxonomy a step-further, by providing added layers to that relationship that and take it outside of just the product domain to other domains such as customer data and digital assets.

Understanding Ontology

- ▶ Ontology is the backbone of Knowledge Graphs, providing a formal, shared understanding of the domain's concepts and relationships.
- ▶ It defines the vocabulary and rules for representing and reasoning about the data, ensuring consistency and accuracy.
- ▶ In the insurance industry, ontology can capture complex relationships between insurance products, coverage options, and risk factors.

Advantages of Ontology in Insurance

- ▶ Enhanced Data Integration: Ontology allows for seamless integration of diverse data sources, leading to a more comprehensive view of customers and policies.
- ▶ Improved Search and Navigation: Ontology-driven Knowledge Graphs enable intuitive search and navigation, making it easier to find relevant information.
- ▶ Better Decision Making: With a formal understanding of concepts and relationships, insurers can make more informed and data-driven decisions.



(Ref: A-SHIP: Ontology-Based Adaptive Sustainable Healthcare Insurance Policy)

Key Components of an Ontology

1. Classes:

- ▶ Classes represent concepts or categories of objects in the domain.
- ▶ In insurance, classes could include "Policy," "Customer," "Claim," etc.

2. Properties:

- ▶ Properties define relationships between classes or properties and classes.
- ▶ Examples include "hasPolicy," "isInsuredBy," "belongsTo," etc.

3. Individuals:

- ▶ Individuals are instances of classes, representing specific objects or entities in the domain.
- ▶ For instance, a specific insurance policy, customer, or claim can be considered individuals.

4. Axioms and Constraints:

- ▶ Axioms and constraints specify rules and logical relationships between classes and properties.
- ▶ They enable reasoning and inferencing in the Knowledge Graph.

5. Hierarchy and Inheritance:

- ▶ Ontologies can have hierarchies, with classes inheriting characteristics from parent classes.
- ▶ For example, a "Health Insurance Policy" inherits properties from the "Insurance Policy" class.

Building an Ontology

- ▶ Identify Key Concepts: Start by identifying the fundamental concepts in the insurance domain, such as policies, claims, beneficiaries, etc.
- ▶ Define Relationships: Establish relationships like 'has policy,' 'covered by,' 'owns,' etc., to represent connections between concepts.
- ▶ Formalize Vocabulary: Create a formal vocabulary that standardizes the representation of concepts and their attributes.
- ▶ Iterative Process: Ontology development is iterative and may evolve as domain knowledge deepens and business requirements change.

Example of Ontology

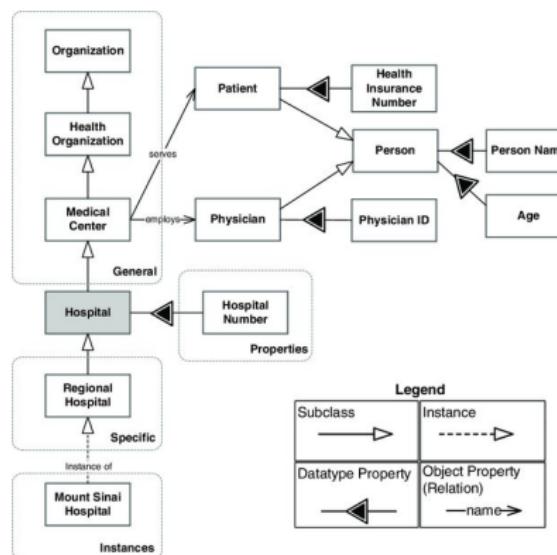


FIGURE: Example Ontology

(Ref: A semantic approach to approximate service retrieval)

Knowledge Graph and Ontology Integration

- ▶ Ontology forms the foundation of a Knowledge Graph, providing structure and meaning to the data.
- ▶ It ensures that data is organized in a way that aligns with the domain's concepts and relationships.
- ▶ Integration allows for powerful querying, reasoning, and inferencing within the Knowledge Graph.

Use Cases in the Insurance Industry

- ▶ Customer Profiling: Creating comprehensive customer profiles based on their policies, claims, and interactions.
- ▶ Risk Assessment: Analyzing data from various sources to assess risks associated with insured entities.
- ▶ Personalized Recommendations: Offering tailored insurance products and coverage options based on individual customer needs.
- ▶ Fraud Detection: Identifying potential fraudulent activities by detecting anomalies in data patterns.

Conclusion

- ▶ Ontology in Semantics is a fundamental component of Knowledge Graphs that empowers the insurance industry with data-driven insights.
- ▶ It enables a formal understanding of complex relationships between insurance concepts, leading to improved decision-making and customer experiences.
- ▶ Embracing ontology-driven Knowledge Graphs positions insurers to thrive in an increasingly competitive and data-centric landscape.

Key Differences between Taxonomy, Ontology, and Knowledge Graphs

Taxonomy

- ▶ Focuses on hierarchical classification of entities based on their shared characteristics.
- ▶ Typically employs a tree-like structure with parent-child relationships.
- ▶ Provides a simple and intuitive way to categorize and organize data.
- ▶ Primarily used for browsing and basic information retrieval.

Ontology

- ▶ Defines a formal, shared vocabulary that represents concepts and their relationships within a specific domain.
- ▶ Captures the meaning and semantics of the data, enabling reasoning and inferencing.
- ▶ Utilized for more advanced data integration, reasoning, and complex information retrieval.
- ▶ Acts as the backbone of a Knowledge Graph, providing a foundation for structuring and understanding data.

Knowledge Graphs

- ▶ Represent interconnected data as entities and their relationships, forming a graph structure.
- ▶ Utilizes ontology to define the schema and provide a formal representation of domain knowledge.
- ▶ Enables sophisticated queries, data analysis, and inference through graph-based algorithms.
- ▶ Offers a holistic view of data, empowering data-driven insights and decision-making in the insurance industry.

Semantics Standards: RDF, RDFS, OWL

Understanding RDF, RDFS, and OWL

- ▶ RDF (Resource Description Framework), RDFS (RDF Schema), and OWL (Web Ontology Language) are fundamental semantic web standards.
- ▶ These standards play a crucial role in knowledge representation and reasoning within Knowledge Graphs.
- ▶ In the insurance industry, these standards enable the development of rich and structured data models to drive better insights and decision-making.

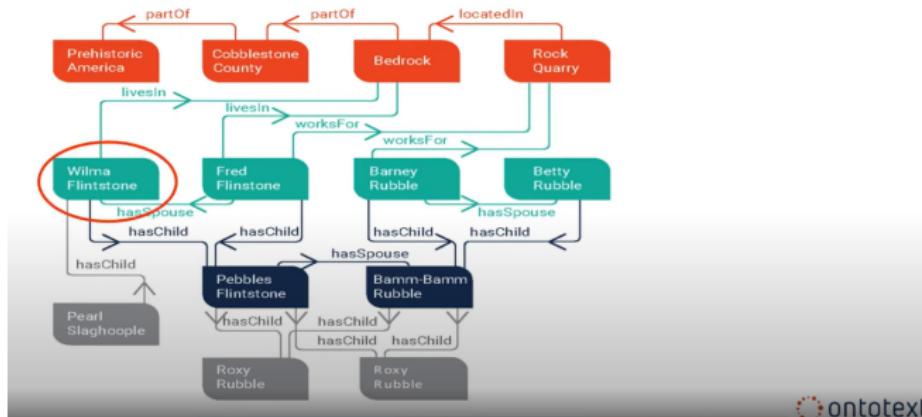
Key Features and Use Cases

- ▶ **RDF (Resource Description Framework):**
 - ▶ Defines a simple model to describe resources and their relationships using subject-predicate-object triples.
 - ▶ Use Cases: Capturing and linking data about policies, claims, customers, and other entities in the insurance domain.
- ▶ **RDFS (RDF Schema):**
 - ▶ Adds schema, resources, class, datatype, etc
 - ▶ Provides a basic vocabulary for creating hierarchies and defining classes and properties.
 - ▶ Use Cases: Specifying insurance-related concepts, properties, and their relationships in a structured manner.
- ▶ **OWL (Web Ontology Language):**
 - ▶ Offers more expressive power for defining classes, properties, and relationships, enabling advanced reasoning and inferencing.
 - ▶ Use Cases: Building complex insurance domain ontologies to support intelligent querying, fraud detection, and risk assessment.

Key Features and Applications of RDF

- ▶ **Subject-Predicate-Object Triples:**
 - ▶ RDF represents information as triples, consisting of subject-predicate-object, forming a relationship statement.
 - ▶ Example: "Policy A (subject) hasCoverage (predicate) Auto (object)."
- ▶ **Resource Identification:**
 - ▶ Each resource in RDF is assigned a unique URI, enabling easy identification and linking of related data.
 - ▶ Example: A URI can uniquely represent a policy, a customer, or a claim in the insurance domain.
- ▶ **Linked Data:**
 - ▶ RDF facilitates the creation of Linked Data by connecting resources across the web, forming a vast knowledge network.
 - ▶ Use Cases: Cross-referencing insurance policies from different providers or analyzing industry trends by linking relevant data sources.
- ▶ **Use RDF when**
 - ▶ You need a simple and straightforward representation of data.
 - ▶ Your knowledge graph primarily involves describing relationships between resources in a triple format.
 - ▶ You want to enable easy data exchange and integration on the web.

RDF



Key Features and Applications of RDFS

► Defining Classes:

- RDFS allows us to define classes, representing groups of resources with shared characteristics.
- Example: In the insurance domain, we can define classes like "Policy," "Customer," and "Claim."

► Defining Properties:

- RDFS enables us to define properties to describe relationships between resources and classes.
- Example: We can define properties like "hasCoverage," "belongsTo," and "hasClaim" to link policies, customers, and claims.

► Hierarchy and Inheritance:

- RDFS supports class hierarchies and inheritance, allowing for more nuanced representations of data.
- Use Cases: Representing the hierarchical relationship between different types of insurance policies or customer categories.

► Use RDFS when

- You need a bit more structure in your knowledge graph with the ability to define classes, properties, and hierarchies.
- Basic inferencing capabilities (e.g., subsumption) are sufficient for your application.

RDFS

```
br:hasSpouse a rdf:Property;  
rdfs:domain br:Human ;  
rdfs:range br:Human .
```

```
br:Fred br:hasSpouse br:Wilma .
```

```
br:Human a rdf:Class;  
rdfs:subClassOf br:Mammal .
```

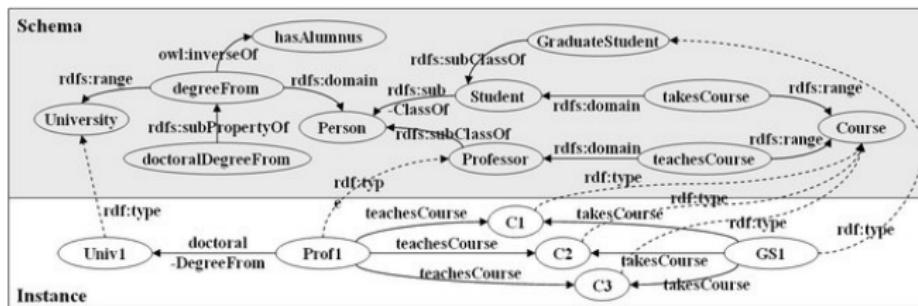
```
⇒ br:Fred a br:Human .  
⇒ br:Wilma a br:Human .
```

hsSpouse relationship is restricted to humans so we can infer that Fred and Wilma are humans and they are mammals.

Key Features and Applications of OWL

- ▶ **Richer Vocabulary:**
 - ▶ OWL provides a wide range of constructs for defining classes, properties, and relationships, enabling more precise modeling.
 - ▶ Example: Representing intricate insurance policies with multiple coverage options and exclusions.
- ▶ **Advanced Reasoning:**
 - ▶ OWL supports automated reasoning and inferencing, allowing for the deduction of implicit information from explicit data.
 - ▶ Use Cases: Detecting potential inconsistencies in insurance policies or automatically inferring customer preferences based on their interactions.
- ▶ **Real-World Complexities:**
 - ▶ OWL enables the representation of real-world complexities in the insurance domain, such as temporal aspects, uncertainty, and probabilistic relationships.
 - ▶ Use Cases: Modeling the dynamic nature of insurance claims processing or incorporating risk probabilities into underwriting decisions.
- ▶ **Use OWL when**
 - ▶ Your knowledge graph requires advanced modeling capabilities and complex relationships between entities.
 - ▶ You need robust reasoning and inferencing to draw logical conclusions from your data.
 - ▶ You want to build ontologies with formal semantics to ensure consistency and reasoning.

OWL



(Ref: An Efficient and Scalable Management of Ontology)

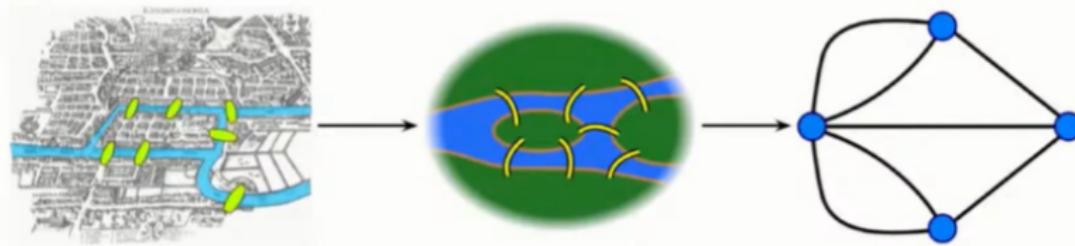
Introduction to Graphs

A graph is

... a set of discrete objects, each of which has some set of relationships with the other objects

Euler: Can we take a walk to all 4 islands, without crossing any of the bridge twice?

Abstraction (Does size of islands matter?):

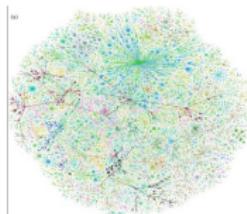


Seven Bridges of Königsberg problem. Leonhard Euler, 1735

Solution: No way!! What's the rule?

(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

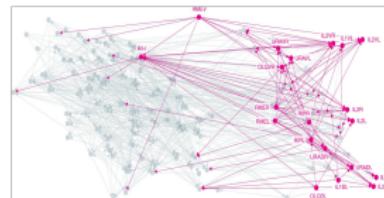
Graph-structured Data Are Ubiquitous



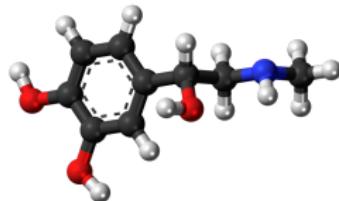
Internet



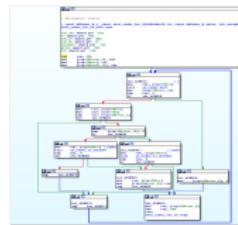
Social networks



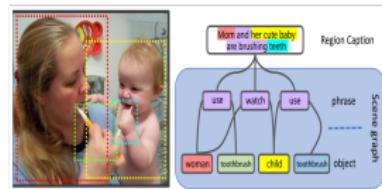
Networks of neurons



Biomedical graphs

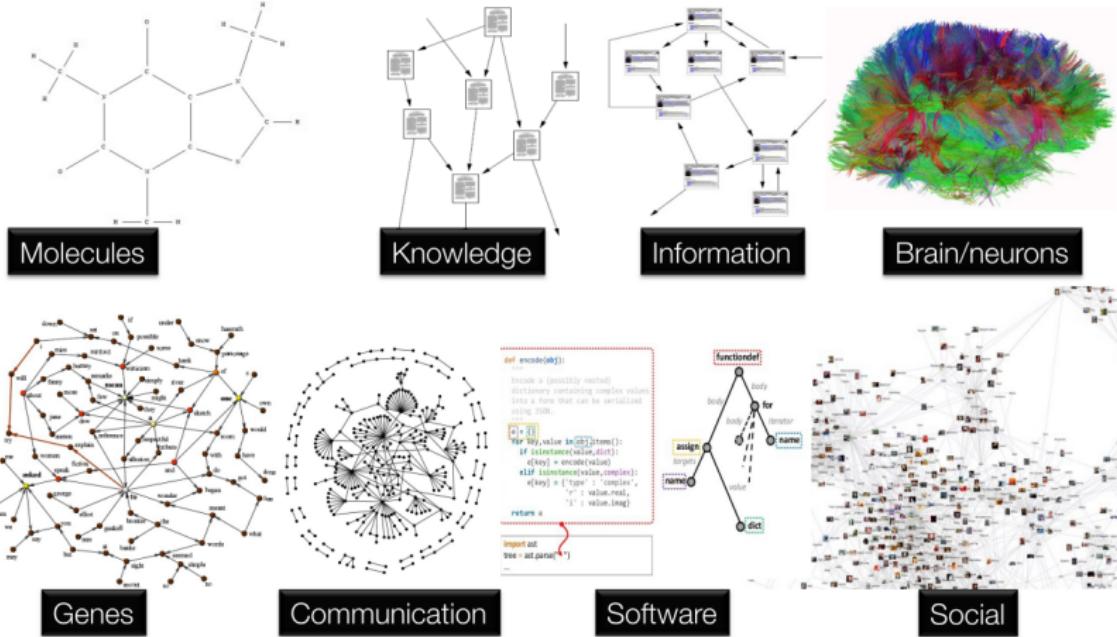


Program graphs



Scene graphs

Networks around us!

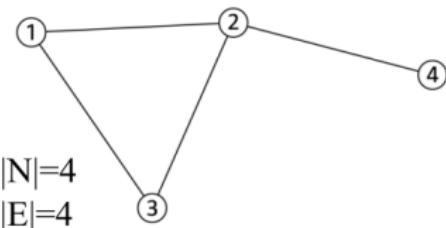
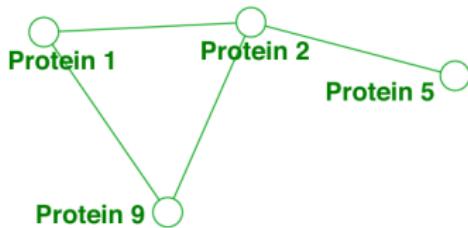
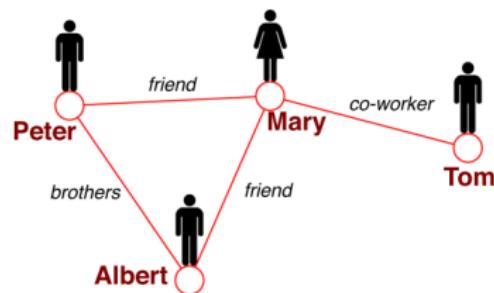
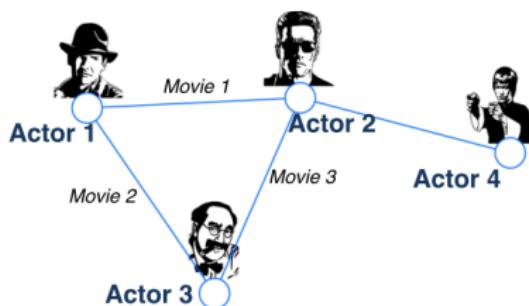


Jure Leskovec (@jure), Stanford University

2

YHK

Graphs: Common Language



Why do graphs matter?

Across an organization, every department can benefit from graphs to answer questions like who or what is important, what should I do next, and what's unusual about this?

Finance

- Fraud Detection
- Pricing Analysis
- Budgeting
- Forecasting

Marketing

- Customer 360
- Influencer Strategy
- Campaign Optimization
- Product Recommendations

Ops

- Product Development
- Pipeline Acceleration
- Supply Chain Optimization
- Infrastructure Planning

IT

- Network Monitoring
- Cybersecurity
- DevOps

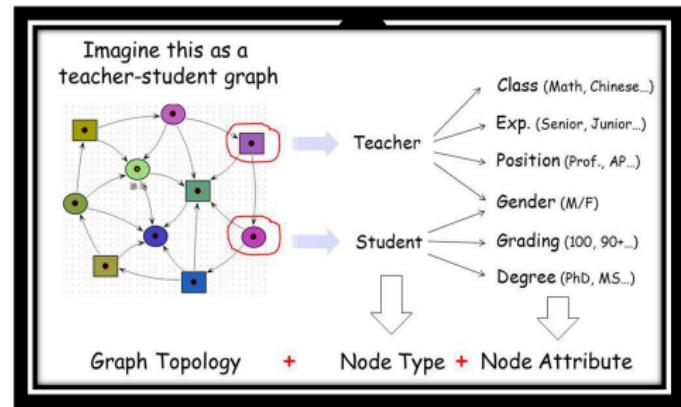
HR

- Training
- Upskilling & Retention
- Promotions

(Ref: 5 Graph Data Science Basics Everyone Should Know - neo4j)

Graphs: A Universal Language

Graphs are a general language for describing and modeling complex systems



Graph!

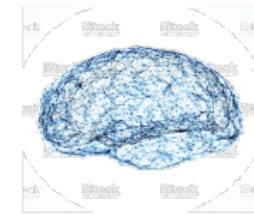
Data as Graphs - Explicit



Social Graphs



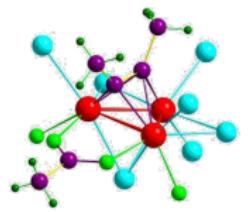
Transportation Graphs



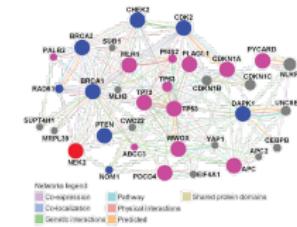
Brain Graphs



Web Graphs

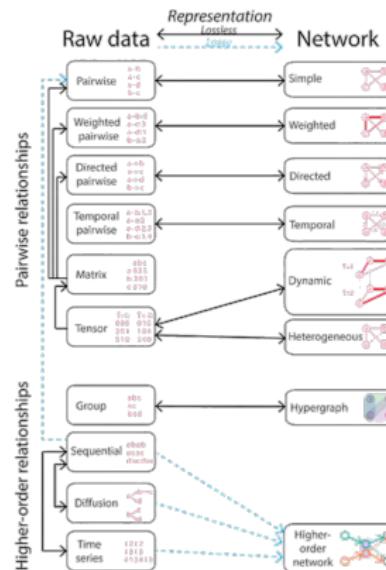
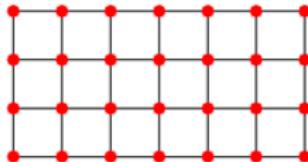


Molecular Graphs



Gene Graphs

Data as Graphs - Implicit



Jian Xu. Representing Big Data as Networks. PhD Dissertation,
University of Notre Dame



Homogeneous vs Multi-relational vs Heterogeneous Graphs

Graph types	Homogeneous	Multi-relational	Heterogeneous
# of node types	1	1	> 1
# of edge types	1	> 1	≥ 1

64

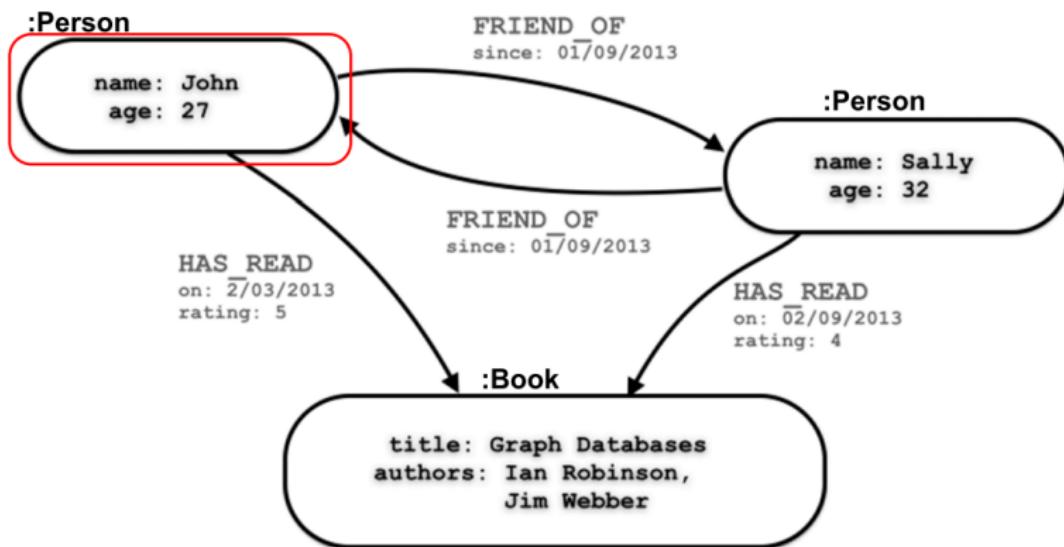


Graph Components

- ▶ Node (Vertex): A must data element for constructing a graph
- ▶ Relationship (Edge) : Link between two nodes, can have direction and type.
- ▶ Label: Node category/type such as PERSON, ORG, etc. One node can have many types.
- ▶ Properties: Attributes or fields in Nodes or Edges, eg. A node can have Label PERSON and Property such as "name: Jane"

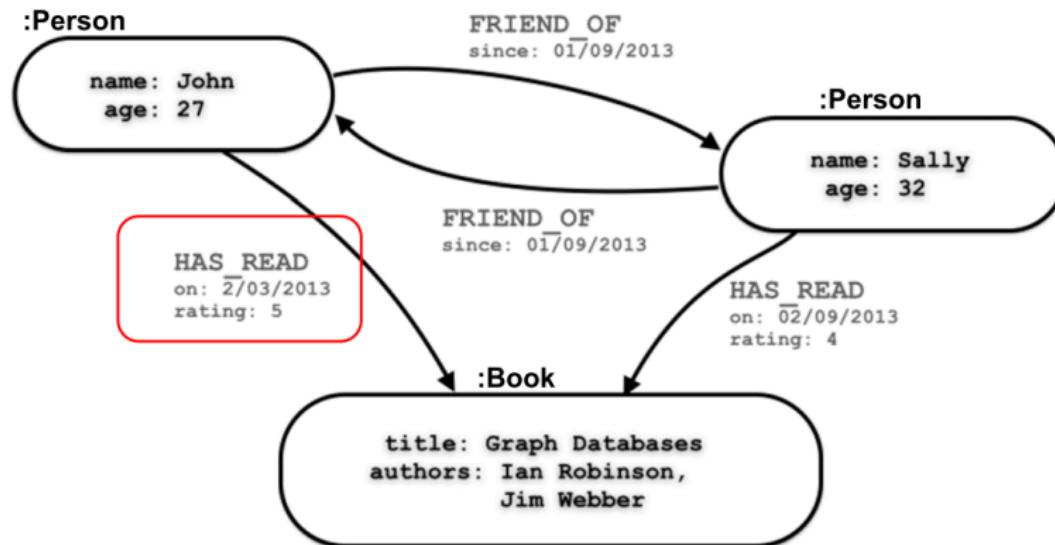
(Ref: Introduction to Neo4j – a hands-on crash course - neo4j)

Nodes



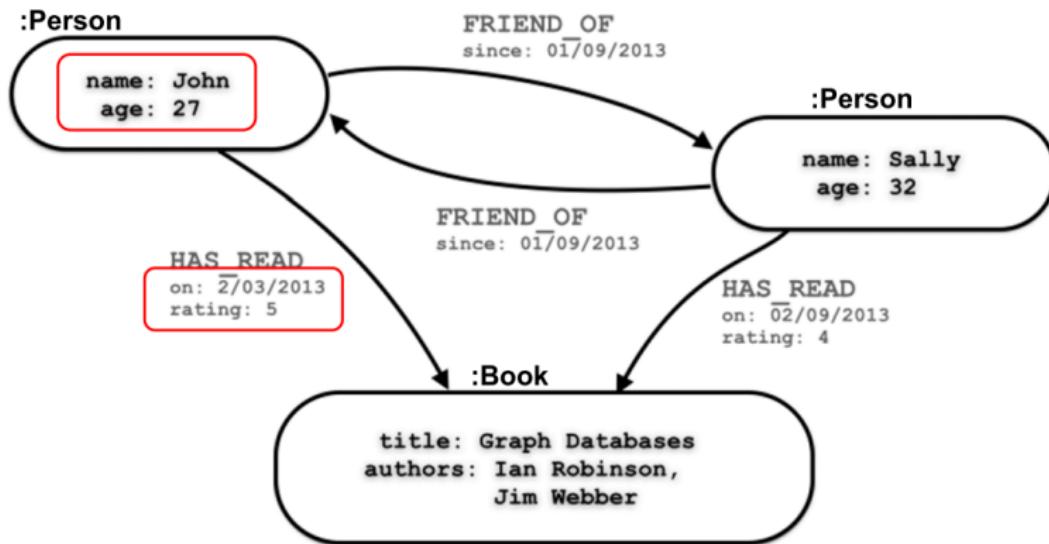
(Ref: CIS 6930 - Advanced Databases - Neo4j)

Relationships



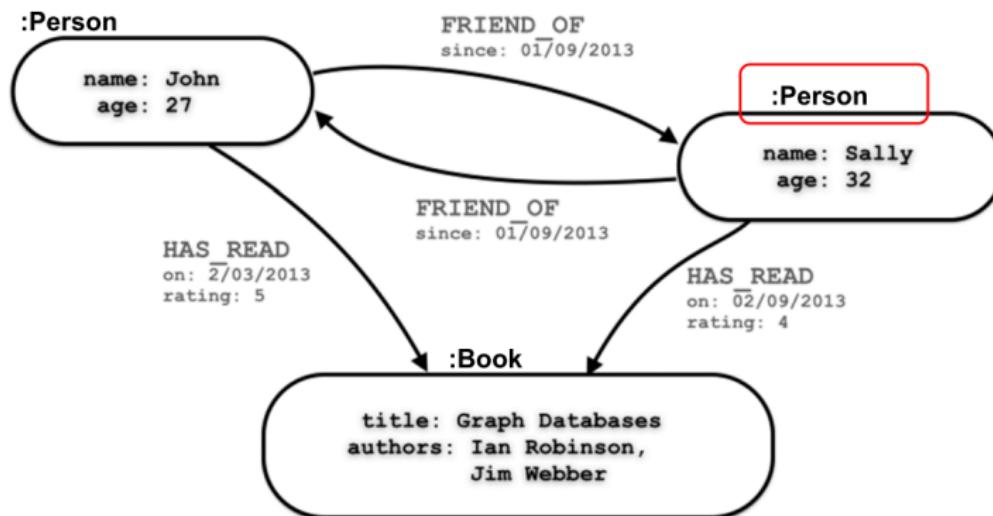
(Ref: CIS 6930 - Advanced Databases - Neo4j)

Properties



(Ref: CIS 6930 - Advanced Databases - Neo4j)

Labels



(Ref: CIS 6930 - Advanced Databases - Neo4j)

Graph-Edge Types

- ▶ Undirected graph: edges are bidirectional, e.g. Michale is brother of John, necessarily means that John is brother of Michale.
- ▶ Directed graph: edges have one direction, e.g. A likes B does not necessarily mean that B likes A.
- ▶ Weighted graph: edges have weights, e.g. connection from city A to city B via road R1 will have weight, say, 8, due to high traffic, but via road R2, may have weight 2, due to lower traffic.

Summary: More Kinds of Graphs

- Weighted Graph



- Labeled Graph

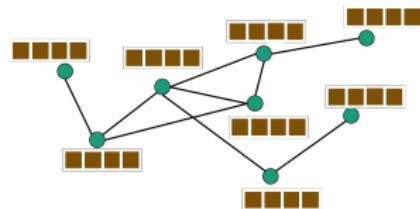


- Property Graph



(Ref: CIntroduction to Graph Databases - Max De Marzi)

Graphs and features



$$\mathcal{V} = \{v_1, \dots, v_N\}$$

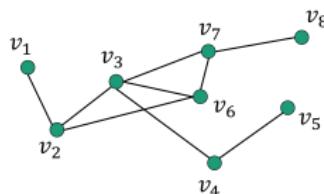
$$\mathcal{E} = \{e_1, \dots, e_M\}$$

$$\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$$

Graph Signal: $f : \mathcal{V} \rightarrow \mathbb{R}^{N \times d}$

$$\mathcal{V} \longrightarrow \begin{bmatrix} f(1) \\ f(2) \\ f(3) \\ f(4) \\ f(5) \\ f(6) \\ f(7) \\ f(8) \end{bmatrix}$$

Matrix Representations of Graphs



Adjacency Matrix: $A[i,j] = 1$ if v_i is adjacent to v_j
 $A[i,j] = 0$, otherwise

Degree Matrix: $\mathbf{D} = \text{diag}(\text{degree}(v_1), \dots, \text{degree}(v_N))$

Degree Matrix	Adjacency Matrix	Laplacian Matrix
$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 4 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$	$- \quad \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}$	$= \begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 3 & -1 & 0 & 0 & -1 & 0 & 0 \\ 0 & -1 & 4 & -1 & 0 & -1 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & -1 & 0 & -1 & 4 & -1 & 0 \\ 0 & 0 & -1 & 0 & 0 & -1 & 3 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix}$

\mathbf{D}

\mathbf{A}

\mathbf{L}

Spectral graph theory. American Mathematical Soc.; 1997.

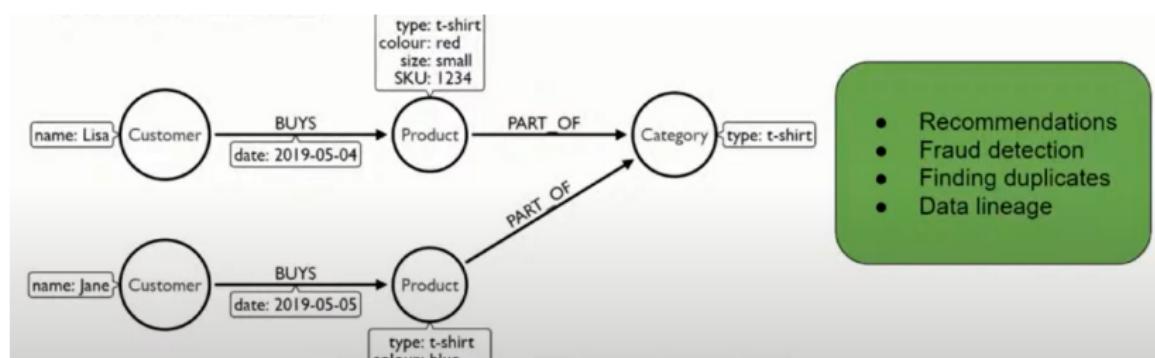
12

Why Graphs? Why Now?

- ▶ Universal language for describing complex data: Networks/graphs from science, nature, and technology are more similar than one would expect
- ▶ Shared vocabulary between fields: Computer Science, Social science, Physics, Biology, Economics
- ▶ Data availability (+ computational challenges): Social/Internet, text, logic, program, bio, health, and medical
- ▶ Impact: Social networking, Social media, Drug design, Event detection, Natural language processing, Computer vision, and Logic reasoning

Identifying good Graph scenarios - 1/4

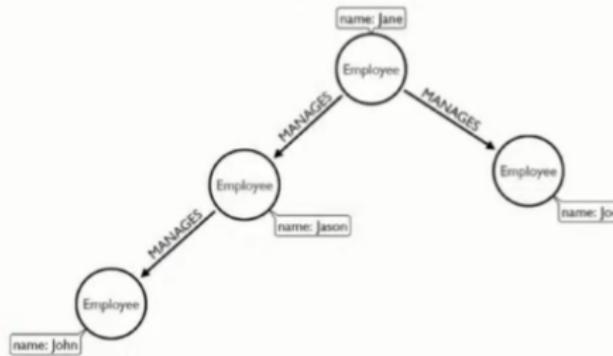
Does the problem involve understanding relationships between entities?
Behavioral analysis:



(Ref: Introduction to Neo4j – a hands-on crash course - neo4j)

Identifying good Graph scenarios - 2/4

Does the problem involve a lot of self-referencing to the same type of entity?
Org chart of employees:

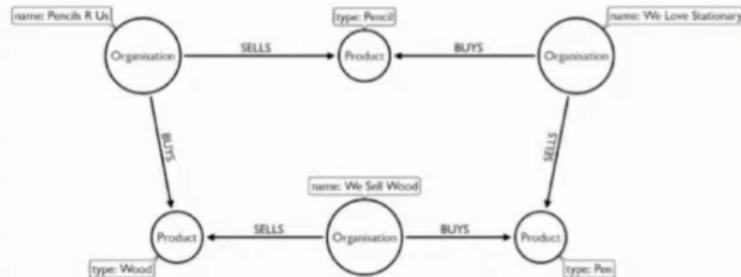


- Organisational hierarchies
- Access management
- Social influencers
- Friends of friends

(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

Identifying good Graph scenarios - 3/4

Does the problem explore relationships of varying and unknown depth?
Changes in manufacturing process:

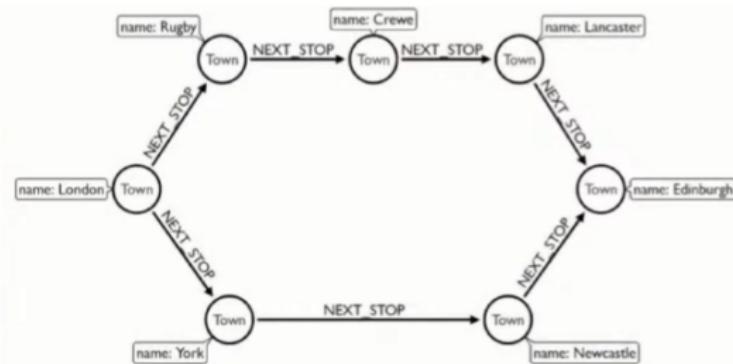


(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

- Supply chain visibility
- Bill of Materials
- Network management

Identifying good Graph scenarios - 4/4

Does the problem involve discovering lots of different routes or paths?
Optimum logistics:



- Logistics and routing
- Infrastructure management
- Dependency tracing

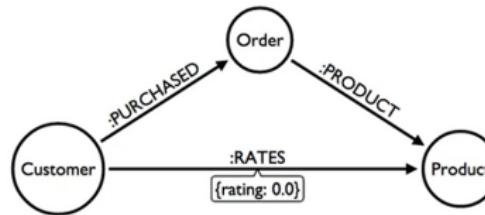
(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

Common Use cases: E-Commerce Recommendations

Easy in graph databases: those who bought A also bought B.

cypher

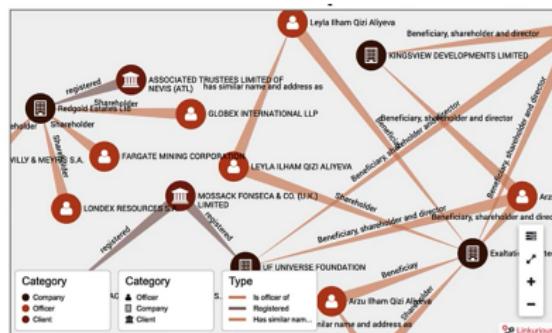
```
MATCH (c:Category)-[:HAS_CHILD|HAS_PRODUCT*1..3]->(p:Product)  
RETURN p.id, p.title, collect(c.name) AS categories
```



(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

Common Use cases: Investigative Journalism

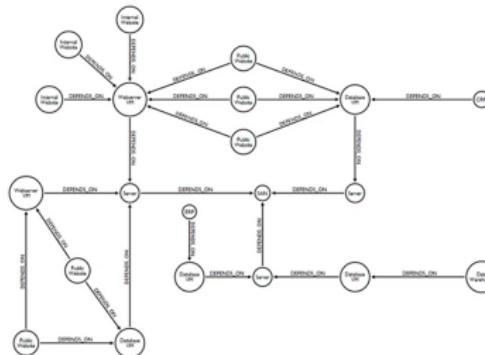
Panama papers: Identify corruption based on relationships between people/companies/financial-institutions.



- What families with the name that contains the string 'aliye' are Officers of Companies?
- How is the family with the name that contains the string 'aliye' related to Companies?
- How are Officers related to each other?
- What are the connections between multiple companies and a family?

(Ref: Introduction to Neo4j – a hands-on crash course - neo4j)

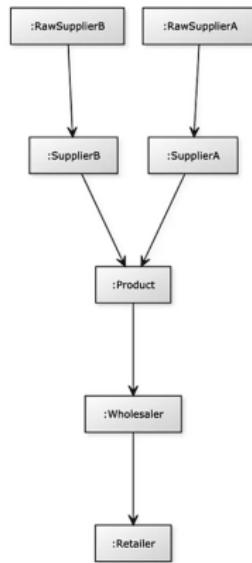
Common Use cases: Network Dependencies



- What are the direct dependencies of public websites?
 - What are the direct dependencies of internal websites?
 - What is the most depended-upon component?
 - Find the dependency chain for a business critical component.
 - What is the impact of removing a server?

(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

Common Use cases: Supply Chain



- Who is the best wholesaler for each retailer based upon distance?
- Which raw supplier will give a particular retailer the freshest products?
- Which retailer provides locally grown products?
- How can we rate each supply chain?

(Ref: Introduction to Neo4j - a hands-on crash course - neo4j)

Graph Algorithms

Centrality

What is Centrality?

- ▶ In a knowledge graph, centrality measures the relative significance of individual nodes.
- ▶ Centrality algorithms help identify key nodes that act as pivotal points of influence in the network.
- ▶ It aids in understanding which nodes are critical for information flow and decision-making processes.
- ▶ Some types of centrality: Degree centrality, Closeness centrality, and Betweenness centrality.
- ▶ More algorithmic types are: Pagerank, Eigenvector

Graph Centralities

► Degree Centrality

- Measures the number of direct connections a node has.
- A node with a higher degree centrality is more connected to other nodes.

► Closeness Centrality

- Measures how quickly a node can reach other nodes in the graph.
- A node with higher closeness centrality is closer to all other nodes.

► Betweenness Centrality

- Measures how often a node appears on the shortest path between other nodes.
- A node with higher betweenness centrality acts as a bridge between different parts of the graph.

The Power of PageRank

- ▶ PageRank, developed by Google's Larry Page and Sergey Brin, revolutionized web search rankings.
- ▶ Originally designed for web pages, PageRank extends to any interconnected network, including insurance-related data.
- ▶ PageRank assigns importance scores to nodes based on incoming links from other important nodes.
- ▶ Nodes with high PageRank are considered influential and contribute significantly to information dissemination.
- ▶ In insurance, PageRank can identify influential policies, clients, or agents within the network.

PageRank Calculation

1. Start with assigning an initial score to all nodes.
2. Iteratively update each node's score based on incoming links and their corresponding scores.
3. Continue the iterations until convergence (scores stabilize).
4. The final scores represent the PageRank values of the nodes.

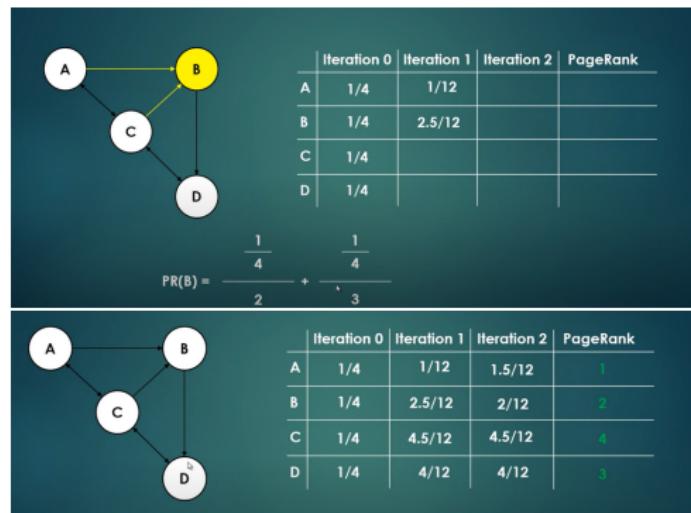
$$\text{score}(v) = \frac{(1-\text{damping_factor})}{N} + \text{damping_factor} \times \sum_{u \in G} \frac{\text{score}(u)}{\text{out_degree}(u)}$$

In this formula, `damping_factor` represents the probability that a user continues to another page rather than following a link, and `out_degree(u)` is the number of outgoing links from node u .

PageRank Pseudo-code

```
1. Initialize the score for each node in the graph:  
2   For each node v in G:  
3     score(v) = 1 / N  
4  
5. Set damping factor (usually 0.85):  
6   damping.factor = 0.85  
7  
8. Set maximum number of iterations and tolerance:  
9   max_iterations = 1000  
10  tolerance = 0.0001  
11  
12. Repeat the following steps until convergence or max_iterations:  
13   For iteration = 1 to max_iterations:  
14     5. Create a copy of the current scores:  
15       For each node v in G:  
16         old_score(v) = score(v)  
17     6. Update the score for each node:  
18       For each node v in G:  
19         score(v) = (1 - damping_factor) / N + damping_factor * sum(score(u) / out_degree(u)) for all nodes u pointing to v  
20     7. Check for convergence:  
21       For each node v in G:  
22         If abs(score(v) - old_score(v)) < tolerance:  
23           Converged for node v  
24     8. If scores have converged for all nodes, break out of the loop  
25  
26. Normalize the final scores:  
27   total_score = sum(score(v)) for all nodes v  
28   For each node v in G:  
29     score(v) = score(v) / total_score  
30  
31. Return the PageRank scores for each node.
```

PageRank



(Ref: PageRank Algorithm - Example Global Software Support)

Find nodes pointing to current node, to summation of those nodes previous page rank, divided by outgoing edges from that node Start with initialization for each node, $1/n$, n is total nodes. For A, lets consider C for now which is pointing to A. For C, ratio is $1/4$ divided by 3 (C's outgoing links), so final is $1/12$

PageRank code

```
1 def pagerank_centrality(graph, damping_factor=0.85, max_iterations=100, tolerance=1e-6):
2     # Initialize PageRank scores for all nodes
3     num_nodes = len(graph)
4     pagerank = {node: 1 / num_nodes for node in graph}
5
6     # Perform iterative PageRank calculation
7     for _ in range(max_iterations):
8         new_pagerank = {}
9         for node in graph:
10             incoming_nodes = [incoming_node for incoming_node in graph if node in graph[incoming_node]]
11             incoming_pagerank = sum(pagerank[incoming_node] / len(graph[incoming_node]) for incoming_node in incoming_nodes)
12             new_pagerank[node] = (1 - damping_factor) / num_nodes + damping_factor * incoming_pagerank
13
14         # Check for convergence
15         if all(abs(new_pagerank[node] - pagerank[node]) < tolerance for node in graph):
16             break
17
18     pagerank = new_pagerank
19
20 return pagerank
```

PageRank code

```
1 # Example usage:  
2 # Create a sample graph representing insurance data  
3 graph = {  
4     1: [2, 3],  
5     2: [3],  
6     3: [4],  
7     4: [3]  
8 }  
9  
10 # Compute the PageRank centrality scores  
11 page_rank_centrality = pagerank.centrality(graph)  
12  
13 # Display the PageRank centrality scores for each node  
14 print("PageRank Centrality Scores:")  
15 for node, score in page_rank_centrality.items():  
16     print(f"Node {node}: {score:.4f}")
```

Eigenvector Centrality

- ▶ Degree centrality means the one who has more neighbors is more influential. But that's not the case in real life. So, a better measure is Eigenvector centrality as it incorporate influence of neighbors also.
- ▶ Eigenvector centrality is another vital centrality measure used in various fields, including finance and insurance.
- ▶ The centrality score of a node depends on the centrality of its neighboring nodes.
- ▶ Nodes connected to other highly central nodes will have higher eigenvector centrality scores.
- ▶ Eigenvector centrality can reveal nodes that have indirect influence through their connections.
- ▶ In insurance, this can help identify key agents with a vast network of influential clients.

Eigenvector Centrality Calculation

1. Start with initializing centrality scores for all nodes.
2. Update node scores iteratively based on their neighbors' scores.
3. NORMALIZE so that max becomes 1.
4. Continue the iterations until convergence.
5. The final scores represent the Eigenvector centrality values of the nodes.

$$\text{score}(v) = \sum_{u \in G} \frac{\text{score}(u)}{\text{degree}(u)}$$

PageRank Pseudo-code

```
1. Initialize the score for each node in the graph:  
    For each node v in G:  
        score(v) = 1  
  
2. Set maximum number of iterations and tolerance:  
    max_iterations = 1000  
    tolerance = 0.0001  
  
3. Repeat the following steps until convergence or max_iterations:  
    For iteration = 1 to max_iterations:  
  
4. Create a copy of the current scores:  
    For each node v in G:  
        old_score(v) = score(v)  
5. Update the score for each node:  
    For each node v in G:  
        score(v) = sum(score(u) / degree(u)) for all nodes u adjacent to v  
6. Calculate the normalization factor (used to prevent score explosion):  
    normalization_factor = max(score)  
7. Normalize the scores:  
    For each node v in G:  
        score(v) = score(v) / normalization_factor  
8. Check for convergence:  
    For each node v in G:  
        If abs(score(v) - old_score(v)) < tolerance:  
            Converged for node v  
9. If scores have converged for all nodes, break out of the loop  
  
10. Return the Eigenvector Centrality scores for each node.
```

Eigenvector Centrality

Calculate: new power \leftarrow sum of power of my neighbors

	0	1	2	3	4	norm
1	1	1	3	6	17	0.32
2	1	3	6	17	38	0.73
3	1	2	6	13	37	0.71
4	1	3	8	19	52	1.00
5	1	2	6	15	39	0.75
6	1	3	7	20	47	0.90



$$\mathbf{v} = \begin{bmatrix} 0.31 \\ 0.79 \\ 0.69 \\ 1.00 \\ 0.78 \\ 0.97 \end{bmatrix}$$

And the update rule is given by: $\mathbf{x}_t = \mathbf{A}^\top \mathbf{x}_{t-1}$

(Ref: NetSci 04-2 Eigenvector Centrality - Andrew Beveridge)

Eigenvector Centrality Code

```
1 import numpy as np
2
3 def eigenvector_centrality(graph, max_iterations=100, tolerance=1e-6):
4     # Initialize Eigenvector Centrality scores for all nodes
5     num_nodes = len(graph)
6     eigenvector_centrality = {node: 1 for node in graph}
7     # Perform iterative Eigenvector Centrality calculation
8     for _ in range(max_iterations):
9         new_eigenvector_centrality = {}
10        for node in graph:
11            incoming_nodes = [incoming_node for incoming_node in graph if node in graph[incoming_node]]
12            incoming_centrality_sum = sum(eigenvector_centrality[incoming_node] for incoming_node in incoming_nodes)
13            new_eigenvector_centrality[node] = incoming_centrality_sum
14        # Normalize the centrality scores
15        norm = np.linalg.norm(list(new_eigenvector_centrality.values()))
16        new_eigenvector_centrality = {node: score / norm for node, score in new_eigenvector_centrality.items()}
17        # Check for convergence
18        if all(abs(new_eigenvector_centrality[node] - eigenvector_centrality[node]) < tolerance for node in graph):
19            break
20        eigenvector_centrality = new_eigenvector_centrality
21    return eigenvector_centrality
```

Eigenvector Centrality Code

```
1 # Example usage:  
2 # Create a sample graph representing insurance data  
3 graph = {  
4     1: [2, 3],  
5     2: [3],  
6     3: [4],  
7     4: [3]  
8 }  
9  
10 # Compute the Eigenvector Centrality scores  
11 eigenvector_centrality_scores = eigenvector_centrality(graph)  
12  
13 # Display the Eigenvector Centrality scores for each node  
14 print("Eigenvector Centrality Scores:")  
15 for node, score in eigenvector_centrality_scores.items():  
16     print(f"Node {node}: {score:.4f}")
```

Real-world Insurance Applications

- ▶ Centrality algorithms find numerous applications in the insurance industry:
 - ▶ Identifying key policyholders for targeted marketing campaigns.
 - ▶ Assessing agents' influence and performance within the network.
 - ▶ Detecting potential fraud by analyzing suspicious connections.
 - ▶ Enhancing risk assessment by understanding complex interdependencies.
 - ▶ Improving customer service through personalized recommendations.
- ▶ Centrality provides actionable insights for better decision-making and operational efficiency.

Data Challenges and Ethical Considerations

- ▶ Insurance companies handle vast amounts of sensitive data, requiring secure handling.
- ▶ Ensuring data privacy and compliance with regulations is critical when applying centrality algorithms.
- ▶ High-quality data and accurate network representations are essential for meaningful results.
- ▶ Ethical considerations include avoiding bias and discrimination in centrality-based decisions.
- ▶ Striking the right balance between data-driven insights and responsible use of centrality is paramount.

Embracing the Power of Centrality

- ▶ Leveraging centrality algorithms can provide a competitive advantage in the insurance industry.
- ▶ Identifying key influencers and critical nodes helps optimize operations and mitigate risks.
- ▶ Centrality complements traditional analytics, leading to better data-driven strategies.
- ▶ Continuous refinement of centrality models enhances decision-making capabilities.
- ▶ Stay ahead of the curve by embracing graph centrality in your insurance workflows.

Similarity

Understanding Similarity

- ▶ In the insurance domain, data is often represented as a knowledge graph.
- ▶ Graph Similarity measures quantify the likeness between nodes or subgraphs within this graph.
- ▶ Similarity algorithms enable us to find nodes or entities that share common characteristics or behaviors.
- ▶ It plays a crucial role in recommendation systems, fraud detection, and customer segmentation.
- ▶ Let's explore one powerful subtype of Graph Similarity: K-Nearest Neighbors (K-NN).

Introducing K-Nearest Neighbors (K-NN)

- ▶ K-NN is a popular algorithm used in various machine learning tasks, including similarity analysis.
- ▶ It is a non-parametric and lazy learning method, making it suitable for dynamic and changing insurance data.
- ▶ K-NN finds K nearest neighbors to a given node based on a similarity metric, such as Euclidean distance or Jaccard index.
- ▶ The algorithm can be adapted to various use cases, such as finding similar clients for personalized recommendations.
- ▶ Next, we'll dive into the mechanics of K-NN and see how it works in the insurance context.

How K-NN Works

1. Start with defining a similarity metric to measure the distance between nodes or entities.
2. Select the value of K (number of neighbors) based on the problem's requirements and data size.
3. For each node in the graph, find its K nearest neighbors based on the similarity metric.
4. Based on the K neighbors, the node's attributes or behavior can be predicted or analyzed.
5. K-NN is intuitive, easy to implement, and doesn't require a training phase.

KNN Pseudo-code

```
1 Input: Graph G with N nodes, target node v, and parameter K
Output: List of K-nearest neighbors of the target node v
3
KNN Algorithm for Graphs:
5 1. Compute the similarity score between the target node v and all other nodes in the graph:
   For each node u in G:
      similarity(u, v) = [Some similarity metric between u and v]
7
9 2. Sort the nodes based on their similarity scores in descending order:
   SortedNodes = Sort(G.nodes, key=lambda u: similarity(u, v), reverse=True)
11
13 3. Select the top K nodes with the highest similarity scores as the K-nearest neighbors:
   KNearstNeighbors = SortedNodes[:K]
15 4. Return the list of K-nearest neighbors KNearstNeighbors.
```

K-NN code

```
1 import numpy as np
2 from collections import defaultdict
3
4 def euclidean_distance(a, b):
5     return np.sqrt(np.sum((a - b) ** 2))
6
7 def k_nearest_neighbors(data, labels, new_data_point, k=3):
8     distances = defaultdict(list)
9
10    for i, data_point in enumerate(data):
11        distance = euclidean_distance(data_point, new_data_point)
12        distances[distance].append(i)
13
14    k_nearest_indices = sorted(distances.keys())[:k]
15    k_nearest_labels = [labels[i] for index in k_nearest_indices for i in distances[index]]
16
17    return max(set(k_nearest_labels), key=k_nearest_labels.count)
```

K-NN code

```
1 # Example usage:  
# Sample insurance data with attributes (age, income)  
3 data = np.array([[25, 50000], [30, 60000], [40, 70000], [45, 80000]])  
labels = np.array(['low risk', 'medium risk', 'high risk', 'low risk'])  
5  
# New data point to classify  
7 new_data_point = np.array([35, 75000])  
9  
# Classify the new data point using K-Nearest Neighbors  
predicted_label = k_nearest_neighbors(data, labels, new_data_point)  
11  
print(f"Predicted label for new data point: {predicted_label}")
```

Real-world Applications in Insurance

- ▶ K-NN has several practical applications in the insurance industry:
 - ▶ Customer Segmentation: Grouping policyholders with similar needs for targeted marketing.
 - ▶ Fraud Detection: Identifying anomalous behavior based on the behavior of similar clients.
 - ▶ Risk Assessment: Predicting risk profiles based on similarities with historical claims data.
 - ▶ Customer Lifetime Value: Estimating the long-term value of clients based on similar customer behaviors.
 - ▶ Cross-selling: Recommending relevant insurance products to clients with similar preferences.
- ▶ K-NN empowers insurers to offer personalized services and make data-driven decisions.

Challenges and Considerations

- ▶ While K-NN offers valuable insights, there are certain challenges and considerations to keep in mind:
 - ▶ Choice of Similarity Metric: Selecting an appropriate similarity measure is crucial for accurate results.
 - ▶ Curse of Dimensionality: High-dimensional data can impact the algorithm's performance.
 - ▶ Scalability: As the dataset grows, computation time may increase, affecting real-time applications.
 - ▶ Data Quality: Inaccurate or incomplete data can lead to misleading similarity results.
 - ▶ Bias: Bias in the data can affect the fairness of recommendations and decisions.
- ▶ Addressing these challenges ensures the responsible and effective use of K-NN in insurance applications.

Integrating Graph Similarity in Workflows

- ▶ Incorporating Graph Similarity and K-NN into your insurance workflows can yield significant benefits:
 - ▶ Enhanced Customer Experience: Personalized services and tailored product recommendations.
 - ▶ Improved Risk Management: Early fraud detection and accurate risk assessment.
 - ▶ Optimal Resource Allocation: Targeted marketing efforts and efficient cross-selling.
 - ▶ Better Decision-making: Data-driven insights for strategic planning and policy optimization.
- ▶ By leveraging these algorithms, insurers can stay competitive and agile in a dynamic market landscape.

Pathfinding

Understanding Graph Pathfinding

- ▶ In the insurance industry, knowledge graphs model complex relationships between entities like policies, claims, and customers.
- ▶ Pathfinding algorithms help us navigate this interconnected web and find the most efficient routes between nodes.
- ▶ The fundamental problem is finding the shortest or most optimal path between two nodes in the graph.
- ▶ These algorithms have significant implications for claim processing, cost optimization, and even customer satisfaction.
- ▶ Now, let's delve deeper into the magic of Shortest Path algorithms.

Introducing Shortest Path Algorithms

- ▶ Shortest Path algorithms are a class of graph algorithms that focus on finding the most efficient route between two nodes.
- ▶ These algorithms are vital in scenarios like:
 - ▶ Determining the shortest route for claims processing, reducing turnaround time.
 - ▶ Identifying the quickest path for field agents to reach clients in emergencies.
 - ▶ Minimizing operational costs by optimizing logistics and resource allocation.
 - ▶ Evaluating risk exposure by analyzing paths through interconnected policies.
- ▶ Next, we'll explore two popular Shortest Path algorithms: Dijkstra's algorithm and A* algorithm.

Dijkstra's Algorithm

- ▶ Dijkstra's algorithm, proposed by Edsger W. Dijkstra in 1956, is a widely used shortest path algorithm.
- ▶ It guarantees finding the shortest path in non-negative weighted graphs.
- ▶ Dijkstra's algorithm is employed in various insurance-related applications, such as optimizing claim processing workflows.
- ▶ This algorithm is efficient for finding the shortest path from a single source to all other nodes in the graph.
- ▶ Let's see how Dijkstra's algorithm works in action.

How Dijkstra's Algorithm Works

1. Start by initializing distances from the source node to all other nodes as infinity.
2. Set the distance from the source node to itself as zero.
3. Repeatedly select the node with the smallest distance from the source (not yet visited).
4. Update the distances to its neighbors if a shorter path is found.
5. Continue until all nodes are visited and their shortest distances are determined.

Dijkstra's Algorithm Pseudo-code

```
1 Input: Graph G with N nodes, source node s
2 Output: Shortest path distances from the source node s to all other nodes
3
4 Dijkstra Algorithm:
5 1. Initialize the distance scores for all nodes in the graph:
6     For each node v in G:
7         distance(v) = infinity
8
9 2. Set the distance score of the source node s to zero:
10    distance(s) = 0
11
12 3. Create an empty priority queue Q to store nodes and their distance scores.
13
14 4. Add the source node s to the priority queue with a distance of zero.
15
16 5. While the priority queue Q is not empty:
17     a. Dequeue the node u with the minimum distance score from Q.
18     b. For each neighbor v of node u:
19         i. Calculate the tentative distance from the source node s to node v:
20             tentative.distance = distance(u) + weight(u, v) // weight(u, v) is the weight of the edge from u to v.
21         ii. If tentative.distance is less than the current distance(v), update distance(v) to the tentative.distance.
22         iii. Enqueue node v into Q with the updated distance(v).
23
24 6. Return the shortest path distances distance(v) for all nodes v in G.
```

Dijkstra's Algorithm code

```
import math
2
def dijkstra(graph, start_node):
    distances = {node: math.inf for node in graph}
    distances[start_node] = 0
6
    visited = set()
8    while len(visited) < len(graph):
        # Find the node with the minimum distance
        min_node = None
10       for node in graph:
            if node not in visited and (min_node is None or distances[node] < distances[min_node]):
                min_node = node
14
            # Calculate distances to neighbors from the current node
16       for neighbor, weight in graph[min_node].items():
            total_distance = distances[min_node] + weight
18           if total_distance < distances[neighbor]:
                distances[neighbor] = total_distance
20
            visited.add(min_node)
22
    return distances
```

Dijkstra's Algorithm code

```
1 # Example usage:  
2 # Create a sample graph representing insurance data with weighted edges  
3 graph = {  
4     'A': {'B': 5, 'C': 3},  
5     'B': {'C': 2, 'D': 4},  
6     'C': {'D': 6, 'E': 8},  
7     'D': {'E': 3},  
8     'E': {}  
9 }  
10  
11 # Find the shortest path distances from node 'A'  
12 shortest_distances = dijkstra(graph, 'A')  
13  
14 # Display the shortest path distances from node 'A' to other nodes  
15 print("Shortest Path Distances from Node 'A':")  
16 for node, distance in shortest_distances.items():  
17     print(f"To node '{node}': {distance}")
```

Usage

1. In the insurance domain, this can be used for route optimization, claims processing, or finding the most efficient way to reach policyholders in case of emergencies.
2. Claims Processing: Suppose an insurance company receives a claim from a policyholder located in a remote area. Using Dijkstra's algorithm, the insurer can find the shortest path to reach the insured's location, minimizing response time and ensuring prompt assistance.

A* Algorithm

- ▶ The A* algorithm is a popular extension of Dijkstra's algorithm and is widely used in pathfinding problems.
- ▶ It employs a heuristic function to guide the search and improve efficiency.
- ▶ A* is particularly useful in scenarios where we need to find the shortest path from a source to a destination node.
- ▶ In insurance, the A* algorithm can optimize navigation for agents visiting clients in different locations.
- ▶ Let's uncover the mechanics of the A* algorithm.

How A* Algorithm Works

1. Similar to Dijkstra's algorithm, initialize distances from the source node to all others as infinity.
2. Set the distance from the source node to itself as zero.
3. In addition to the distance, compute a heuristic value (estimated distance) from each node to the destination node.
4. Repeatedly select the node with the smallest sum of distance and heuristic value (not yet visited).
5. Update the distances to its neighbors if a shorter path is found, considering the heuristic value.
6. Continue until the destination node is reached, and the shortest path is determined.

A star Algorithm Pseudo-code

```
Input: Graph G with N nodes, source node s, target node t
2 Output: Shortest path from source node s to target node t

4 A* Algorithm:
1. Initialize the G-Score and F-Score for all nodes in the graph:
6   For each node v in G:
     G-Score(v) = infinity // The actual cost of the cheapest path from s to v
8     F-Score(v) = infinity // G-Score(v) + heuristic_estimate(v, t)

10 2. Set the G-Score of the source node s to zero:
    G-Score(s) = 0
12 3. Set the F-Score of the source node s using the heuristic estimate:
14   F-Score(s) = heuristic_estimate(s, t)

16 4. Create an open set to store nodes to be evaluated. Initially, it contains only the source node s.

18 5. While the open set is not empty:
a. Dequeue the node u with the lowest F-Score from the open set.
20 b. If u is the target node t, reconstruct the path from s to t and return it.
c. For each neighbor v of node u:
22   i. Calculate the tentative G-Score for node v:
      tentative_G_Score = G-Score(u) + weight(u, v) // weight(u, v) is the weight of the edge from u to v.
24   ii. If tentative_G_Score is less than G-Score(v), update G-Score(v) to tentative_G_Score and set F-Score(v) using the
      heuristic estimate:
      G-Score(v) = tentative_G_Score
26     F-Score(v) = G-Score(v) + heuristic_estimate(v, t)
28   iii. If v is not already in the open set, enqueue it.

6. No path from s to t is found, return failure.
```

A* Algorithm code

```
1 import math
3
3 def euclidean_distance(a, b):
4     return math.sqrt((a[0] - b[0])**2 + (a[1] - b[1])**2)
5
6 def reconstruct_path(graph, start, goal):
7     current = goal
8     path = [current]
9     while current != start:
10         current = min(graph[current], key=lambda node: g_score[node] + euclidean_distance(node, goal))
11         path.append(current)
12     return list(reversed(path))
13
14 def a_star(graph, start, goal):
15     open_set = set([start])
16     closed_set = set()
17
18     g_score = {node: math.inf for node in graph}
19     g_score[start] = 0
20
21     f_score = {node: math.inf for node in graph}
22     f_score[start] = euclidean_distance(start, goal)
23
```

A* Algorithm code

```
1
2
3     while open_set:
4         current = min(open_set, key=lambda node: f_score[node])
5         if current == goal:
6             return reconstruct_path(graph, start, goal)
7         open_set.remove(current)
8         closed_set.add(current)
9
10        for neighbor, cost in graph[current].items():
11            if neighbor in closed_set:
12                continue
13            tentative_g_score = g_score[current] + cost
14            if neighbor not in open_set:
15                open_set.add(neighbor)
16            elif tentative_g_score >= g_score[neighbor]:
17                continue
18            g_score[neighbor] = tentative_g_score
19            f_score[neighbor] = g_score[neighbor] + euclidean_distance(neighbor, goal)
20
21    return None
```

A* Algorithm code

```
# Example usage:  
2 # Create a sample graph representing insurance data with weighted edges (Euclidean distances)  
graph = {  
    4   (0, 0): {(1, 1): 1.41, (0, 1): 1.0},  
    (1, 1): {(2, 2): 1.41, (1, 0): 1.0},  
6   (0, 1): {(0, 2): 1.0},  
    (2, 2): {(2, 3): 1.0},  
8   (0, 2): {(1, 2): 1.0},  
    (1, 2): {(2, 2): 1.0, (1, 3): 1.0},  
10  (2, 3): {}  
12 }  
14  
# Define start and goal nodes  
start_node = (0, 0)  
goal_node = (2, 3)  
16  
# Find the shortest path using A* algorithm  
18 shortest_path = a_star(graph, start_node, goal_node)  
20  
# Display the shortest path  
print("Shortest Path from", start_node, "to", goal_node, ":", shortest_path)
```

Real-world Insurance Applications

- ▶ Shortest Path algorithms offer numerous practical applications in the insurance industry:
 - ▶ Optimizing claims processing by finding the most efficient route through interconnected policies and claims data.
 - ▶ Streamlining field agent operations for faster response times and improved customer service.
 - ▶ Analyzing risk exposure and dependencies between different policies and insurance products.
 - ▶ Minimizing operational costs by optimizing resource allocation and logistics.
- ▶ These algorithms empower insurers to make data-driven decisions, enhancing efficiency and customer satisfaction.

Challenges and Considerations

- ▶ While Shortest Path algorithms offer invaluable insights, there are challenges and considerations to address:
 - ▶ Data Accuracy: Reliable data is essential to ensure the accuracy of the shortest path.
 - ▶ Dynamic Environments: Insurance data is ever-changing; algorithms should adapt accordingly.
 - ▶ Scalability: As the graph size grows, computational complexity can impact real-time applications.
 - ▶ Heuristic Selection: Choosing appropriate heuristics affects the algorithm's efficiency and accuracy.
 - ▶ Ethical Use: Fair treatment and transparency are paramount when optimizing paths for claims processing or resource allocation.
- ▶ Responsible implementation ensures that Shortest Path algorithms deliver tangible benefits to the insurance industry.

Community Detection

Understanding Graph Community Detection

- ▶ In the insurance industry, knowledge graphs represent complex relationships among policies, claims, and customers.
- ▶ Community detection algorithms help identify cohesive groups or communities within these interconnected networks.
- ▶ These algorithms reveal clusters of related entities, enabling targeted analysis and decision-making.
- ▶ Community detection has applications in customer segmentation, fraud detection, and network optimization.
- ▶ Now, let's explore one popular subtype of community detection: Label Propagation.

Introducing Label Propagation Algorithm

- ▶ Label Propagation is a widely used algorithm for community detection in graphs.
- ▶ It is an iterative algorithm that propagates labels (community assignments) based on local connectivity.
- ▶ Label Propagation is particularly useful for detecting communities in large graphs with unknown community structures.
- ▶ The algorithm identifies communities based on the assumption that nodes with similar attributes are likely to belong to the same community.
- ▶ Now, let's delve into the mechanics of the Label Propagation algorithm.

How Label Propagation Works

1. Start by assigning a unique label (community identifier) to each node in the graph.
2. In each iteration, nodes update their labels based on the labels of their neighbors.
3. The update rule considers the labels of neighboring nodes, giving more weight to frequently occurring labels.
4. The process continues iteratively until the labels stabilize or reach a predefined convergence condition.
5. Nodes with the same labels after convergence belong to the same community.

Community Detection by Label Propagation Pseudo-code

```
1 Input: Graph G with N nodes
2 Output: Community assignments for each node in G
3
4 Label Propagation Algorithm:
5 1. Initialize a unique label for each node in the graph:
6     For each node v in G:
7         label(v) = v
8
9 2. Repeat the following steps until convergence or a maximum number of iterations:
10    For iteration = 1 to max_iterations:
11
12    3. Create an empty dictionary for counting label occurrences in the neighborhood:
13        neighbor_labels_count = {}
14
15    4. For each node v in G:
16        5. Get the neighboring labels of node v:
17            neighbor_labels = [label(u) for all nodes u adjacent to v]
18
19        6. Count the occurrences of each label in the neighbor_labels list:
20            For each label l in neighbor_labels:
21                neighbor_labels_count[l] += 1
22
23        7. Find the most frequent label(s) in the neighbor_labels_count:
24            most_frequent_labels = labels with the highest counts in neighbor_labels_count
25
26        8. If there is a tie (multiple labels with the same highest count), randomly choose one label from most_frequent_labels.
27
28        9. Assign the most frequent label to node v:
29            label(v) = randomly chosen label from most_frequent_labels
30
31 10. Return the community assignments label(v) for each node v in G.
```

Label Propagation code

```
1 def label_propagation(graph):
2     labels = {node: node for node in graph}
3
4     # Perform iterative label propagation
5     max_iterations = 100
6     for _ in range(max_iterations):
7         modified = False
8         nodes = list(graph.keys())
9         random.shuffle(nodes) # Randomly shuffle the order of nodes
10        for node in nodes:
11            neighbor_labels = [labels[neighbor] for neighbor in graph[node]]
12            most_common_label = max(set(neighbor_labels), key=neighbor_labels.count)
13            if labels[node] != most_common_label:
14                labels[node] = most_common_label
15                modified = True
16
17        if not modified:
18            break
19
20    return labels
```

Label Propagation code

```
# Example usage:  
2 # Create a sample undirected graph representing insurance data  
graph = {  
4     1: {2, 3},  
5     2: {1, 3},  
6     3: {1, 2, 4},  
7     4: {3},  
8     5: {6},  
9     6: {5}  
10 }  
11  
12 # Perform community detection using Label Propagation  
community_labels = label_propagation(graph)  
14  
13 # Display the community labels for each node  
16 print("Community Labels:")  
for node, label in community_labels.items():  
    print(f"Node {node}: Community {label}")
```

Real-world Applications in Insurance

- ▶ Community detection using Label Propagation offers numerous practical applications in the insurance industry:
 - ▶ Customer Segmentation: Identifying groups of policyholders with similar needs and preferences.
 - ▶ Fraud Detection: Uncovering clusters of suspicious activities or fraud rings.
 - ▶ Network Analysis: Analyzing relationships between agents, clients, and policies.
 - ▶ Product Recommendations: Targeting specific insurance products to relevant customer communities.
 - ▶ Claims Analysis: Understanding patterns and correlations among claims to optimize processing.
- ▶ These applications enable insurers to tailor services, detect anomalies, and optimize operations.

Benefits and Considerations

- ▶ Label Propagation for community detection provides several benefits but also requires careful consideration:
 - ▶ Scalability: The algorithm performs well on large graphs, making it suitable for insurance networks.
 - ▶ Unsupervised Learning: Label Propagation does not require predefined community labels or training data.
 - ▶ Real-time Analysis: The algorithm is iterative and can adapt to dynamic insurance data.
 - ▶ Interpretability: Detected communities can be interpreted to gain insights into the underlying structure.
- ▶ However, parameter tuning and handling disconnected nodes should be considered for accurate results.

Conclusions

- ▶ **Centrality Algorithms**
 - ▶ PageRank assigns scores to nodes based on their incoming links and the importance of the linking nodes.
 - ▶ Eigenvector Centrality considers both the number of connections and the centrality of the nodes' neighbors.
- ▶ **Similarity Algorithms**
 - ▶ K-Nearest Neighbors is a versatile algorithm for similarity-based tasks.
 - ▶ It classifies or predicts based on the similarity of a node to its neighbors.
 - ▶ KNN is useful in collaborative filtering, recommendation systems, and clustering applications.
- ▶ **Pathfinding Algorithms**
 - ▶ Shortest path algorithms, like Dijkstra's and A*, find the most efficient routes between nodes.
 - ▶ Dijkstra's algorithm works on non-negative edge weights, while A* incorporates heuristics for faster convergence.
 - ▶ These algorithms are essential in navigation, network routing, and logistics optimization.
- ▶ **Community Detection Algorithms**
 - ▶ Label Propagation is a simple but effective method for detecting communities in graphs.
 - ▶ It iteratively propagates labels based on the labels of neighboring nodes.
 - ▶ Label Propagation is widely used for social network analysis, identifying clusters, and understanding network structures.

References

- ▶ Neo4j website: <https://neo4j.com>
- ▶ "Graph Algorithms" book by Mark Needham and Amy E. Hodler
- ▶ "Graph Databases" book by Ian Robinson, Jim Webber, and Emil Eifrem
- ▶ Neo4j documentation and developer guides
- ▶ "an introduction to neo4j (graph database tutorial for beginners)" - Chris Hay
- ▶ CIS 6930 - Advanced Databases - Neo4j
- ▶ "Learning Neo4j" - Wabri/LearningNeo4j - Github
- ▶ Neo4J Certification Sample Questions -
https://wiki.glitchdata.com/index.php/Neo4J_Certification
- ▶ Neo4j Certified Professional: Exam Practice Tests - By Cristian Scutaru
- ▶ "Graph Data Science with Neo4j Graph Algorithms - Will Lyon" Youtube
- ▶ "Mark Needham - Intro to Graph Data Science with Neo4j" Youtube
- ▶ "THE POWER OF "GRAPH DATA SCIENCE" using Neo4j" Youtube
- ▶ "Using Neo4j Graph Data Science in Python to Improve Machine Learning Models" - Tomaz Bratanic
- ▶ "Graph Data Science for Supply Chains – Part 1" - Zach Blumenfeld