#### GRAPH+AI SUMMIT organized by Tiger Graph

# Graph Data Science Algorithms for Analytics and Machine Learning

Dr. Victor Lee, VP of ML/AI, TigerGraph



## Outline

- The Power of Graph Analytics and Algorithms
- In-Database Graph Data Science
- Three ways to leverage Graph Algorithms for ML
  - Unsupervised Learning
  - Feature Enrichment from Graph Features
  - In-Database Learning



## **Smarter Analytics and AI**

Good business decisions require data with context, recency, and relationship

- Context it is situational?
- Recency is it fresh?
- Relationship how do different facts relate to one another?

Good AI requires training data with the right content and right structure

Graphs provide the relationships, rich content, and flexible structure for better analytics and decision making



# Real-World Better Outcomes from Graph+AI

Healthcare:

Real-time recommendations



- 1.3TB graph brain
- Real-time care recommendations
- Improving healthcare, lowering cost

Industrial Supply Chain: Analytics for decisions



- Analytics: weeks  $\rightarrow$  minutes
- Reveal opportunities, optimize tactical & strategic decisions
- Saving \$100M+/yr

#### Financial Services: Real-time fraud detection



- Integrates multiple tools
- "Magical" real-time visual results for investigators
- Scalable for growth



## Power of Graph



#### **Richer, Smarter Data with Relationship**

- Connections-as-data
- Connects datasets, breaks down silos

#### **Deeper, Smarter Questions**

- Look for semantic patterns of relationship
- Search farther than other DBs

#### **More Computational Options**

- Graph algorithms
- Graph-enhanced machine learning

#### **Explainable Results**

- Semantic data model, queries, and answers
- Visual exploration and results

# How is Graph Model different than other Data Models ?



Graph Model can represent all other Data Models Naturally !

Author: Dan McCreary: https://www.slideshare.net/Dataversity/nosql-now-nosql-architecture-patterns-23589170

## Analytical Graph Database Characteristics

Feature	Design Difference	Benefit
Real-Time Deep-Link	<ul> <li>Native Graph design</li> <li>C++ engine for high performance</li> <li>Storage Architecture</li> </ul>	<ul> <li>Uncovers hard-to-find patterns</li> <li>Operational, real-time</li> <li>HTAP: Transactions+Analytics</li> </ul>
	<ul> <li>Distributed DB architecture</li> <li>Massively parallel processing</li> <li>Compressed storage reduces footprint and messaging</li> </ul>	<ul> <li>Integrates all your data</li> <li>Automatic partitioning</li> <li>Elastic scaling of resource usage</li> </ul>
In-Database Analytics & Marining	<ul> <li>GSQL: High-level yet Turing-complete language</li> <li>User-extensible graph algorithm library, runs in-DB</li> <li>ACID (OLTP) &amp; Accumulators (OLAP)</li> </ul>	<ul> <li>Avoids transferring data</li> <li>Richer graph context</li> <li>Graph-based feature extraction for supervised machine learning</li> <li>In-DB machine learning training</li> </ul>



Three ways to leverage Graph Algorithms for ML

- 1. Unsupervised Learning
- 2. Feature Enrichment from Graph Features
- 3. In-Database Learning



# (1) Graph Algorithms for Unsupervised Learning

Types of Graph Algorithms

- Path Finding
- Ranking and Centrality
  - PageRank, HITS
  - Closeness, Betweenness
- Similarity
  - Cosine, Jaccard
- Clustering / Community Detection
  - Louvain modularity
- Frequent Pattern Discovery

**BOLD** indicates more complex tasks, which produced a value for each vertex and can be considered **unsupervised learning** 





# Graph Algorithm Use Cases

Algorithm	Key Use Cases
Shortest Path (X to Y)	<ul> <li>Most efficient route/process from X to Y, physical or conceptual</li> <li>Likelihood of two parties X &amp; Y being aware of one another</li> </ul>
PageRank of X	<ul> <li>Predicting X's influence on others</li> <li>Likelihood of an arbitrary person/entity being aware of Entity X</li> </ul>
Closeness/ Betweenness	<ul> <li>Most efficient location for a hub</li> </ul>
Community Detection	Identifying social groups or groups of mutual awareness/influence
Similarity between (X,Y), Classification	<ul> <li>Recommendations and Substitutions (If you liked X, you'll like Y)</li> <li>Prediction (If something was true for X, it may be true for Y)</li> <li>Classification (X and Y belong to the same category)</li> </ul>



# TigerGraph In-Database Graph Data Science Library

Signals our commitment to serving the needs of data scientists

- More algorithms (15 released this week)
  - Graph Embedding
  - Link prediction
  - Similarity
  - Centrality
- More than just algorithms

#### In-Database

- No export needed
- Live, updatable data
- Scaleable, Ultra-fast engine
- GSQL query language

#### For Data Scientists

- easier and faster to run
- include ML, such as graph embeddings
- will integrate with feature & model management
- will integrate with Graph+ML Workbench



## Finding the Most Influential Health Care Providers in a Community

• Who is the most influential provider in each region for a particular medical condition?

⇒ Use PageRank to rank each provider based on the relative importance of their referrals

• Who is influenced by these leaders (e.g. other doctors, chiropractors, physical therapists, facilities)?

⇒ Use Community Detection to find the groups surrounding Influencers



Graph with Patients, Providers, and Service Claims



# Finding Similar Cases to deliver better healthcare

က်ိဳ Mem		Jouri Tiger Grap	ney			$\checkmark$	×		A		Q	
Member Name: Age: 78 Doris Smith Gender: Female Age: 78 DOB: 04/17/41			Phone Number: (650) 888-9090 Email: dsmith41@gma	Home Address: 3 Main St. Redwood City, CA 94065				Find Similar Members				
EVENTS			May	Ju	ne	Ju	ıly	Aug	ust Septe	ember Od	ctober N	ovember
Errollment     Findlment     Prescher Claims     Prescher Claims     Weiness Check     Dertal Claim     Heathcare Advisor Visits     Heathvorae Advisor Visits     Labs     Admissions     Program Outreach     Outbound Call     Inbound Call	Enro	llment	P									
	Wellne	ss Check		(			<b>(</b>		$\triangle$		2	
	Prescib	er Claims			Ì							
	Testing & Cla	Procedure aims				À					A.B.	
TIMELINE Last 1 Day Last 7 Days Last 30 Days Last 90 Days Last 1 Year Y Custom T 0 1 totoors @ Custom	Pharma	cy Claims					3				8	
	Healthca Vi	re Advisor sits			C.A.		, A	L	Δ	⚠	S A	€
	Inbou	ind Call		9		9				9		



- Find similar members with a click of a button in real-time
- Deliver care path recommendations for similar members



# Finding Similar Cases to deliver better healthcare

က်ိဳ Mem	ber	Jouri Tiger Grap	ney	/		1	×		A	4			α.	•
Member Name: Age: 78 Doris Smith DoB: 04/17/41			Phone Number: (650) 888-9090 Email: dsmith41@gma	Home Address: 3 Main St. Redwood City, CA 94065				Find Similar Members				ers		
EVENTS			May	Jur	ne	Ju	ıly	Aug	gust Si	ptember	Oct	ober	Nover	nber
Enrollment     Pharmacy Claim     Presciber Claims     Melanan Charle	Enro	llment	P											
Verlense Check     Dental Claim     Testing & Procedure Claims     Heathcare Advisor Visits     Behavioral Claim     Labs     Admissions     Program Outreach     Outbound Call     Inbound Call	Wellne	ss Check		6	a)		1		♪		(	2		
	Prescib	er Claims			]									
	Testing & Cla	Procedure aims				Ø						6	<b>C</b>	
TIMELINE Last 1 Day Last 7 Days Last 30 Days Last 90 Days Last 1 Year C Custom To C m transmit	Pharma	cy Claims					3					8		
	Healthca Vi	re Advisor sits			¢д	(	\$.Ê	۷	ß	⚠	(	ŝ	Ğ.	•
	Inbou	ind Call		۷		۲					Ŷ		9	



- Find similar members with a click of a button in real-time
- Deliver care path recommendations for similar members



# **Entity Resolution** using Similarity Scores



Several scoring systems: Jaccard, Cosine similarity, Kolmogorov distance, etc.



## Entity Resolution using Similarity Scores





# (2) Graph Feature Extraction - Better Training Data



I GRAPHAISUMMIT.COM I #GRAPHAISUMMIT

#### Challenge

Find and report fraudsters among billions of calls per week.

#### Solution

- **Build graph**: Real-time operational graph with 600M phone nodes & 15B call detail records.
- **Get features and labels**: Domain experts write GSQL queries to extract 118 features/phone. Some past calls are labeled for 3 types of unwanted calls.
- **Train**: Feed machine learning with training data for fraud detection with 118 features/phone for 30M calls.
- **Deploy**: For each incoming call, extract the current 118 features (subsecond) and apply model for real-time answer.

#### Results

- If unwanted call is predict, display alert on recipient's phone
- Process 2000+ calls/sec
- Improved customer satisfaction

17

## Graph Feature Example





## Graph Features Support Explainable AI



visualization, exploration and features

feature computation

GRAPH+AI SUMMIT

TigerGraph

| GRAPHAISUMMIT.COM | #GRAPHAISUMMIT

# Feature Extraction using Graph Embeddings

, Graph Embedding transforms graph structure into a compact set of vertex vectors.

- Captures the essence of a vertex's "nature" as a set of latent features
- Enables graph data to run efficiently on non-graph neural networks



Works for numerous cases

- Recommendation (similarity)
- Fraud detection (classification)

 $\begin{array}{c} \text{Compact} \rightarrow \\ \textbf{scalability} \\ \text{Set of vectors} \rightarrow \\ \textbf{compatibility,} \\ \textbf{reduced complexity} \end{array}$ 



## (3) In-Database Machine Learning

2. Enrich with Graph Features







Prediction



- Simplified pipeline
- No need to export/import
- Fresh data, up to date
- DB needs to be scalable, with parallel processing



- Entity resolution
- Recommendation
- Fraud detection

• ...



| GRAPHAISUMMIT.COM | #GRAPHAISUMMIT

## In Database ML for Movie Recommendation



MARVEL'S THE AVENGERS PG13, 2 hr.22 min. Action & Adventure , Science Fiction & Fantasy Directed By: Joss Whedon In Theaters: May 4, 2012 Wide On DVD: Sep 25, 2012 Walt Disney Pictures





movie features



All Critics Top C	Critics All Audience	
users	ratings NEXT →	Low-Rank Approximation Machine Learning v3
Danny D	★★★★★ 5d ago How many movies did it take to come up with this mundane plot ?	
Benjamin C	<ul> <li>Goals:</li> <li>Predict users' ratings for movies, based on previous ratings</li> </ul>	TigerGraph GraphGurus
Martyn K	<ul> <li>Recommend movies to users based on rating prediction</li> </ul>	EPISODE 28 An In-Database Machine Learning Solution For Real-Time Recommendations

## User-Rating-Movie Graph



#### MovieLens dataset

https://grouplens.org/datasets/movielens/

- 100K ratings and 40K tags that 1K users gave to 17K movies
- Ratings are from 0 to 5 stars

#### **Recommendation** Approaches

- Collaborative filtering
- Content based method
- K-nearest neighbors
- Latent factor (model-based)
- Hybrid method

•••



## Movie Rating Prediction (Latent Factor Model)

\_ \_ \_ \_ \_ \_ \_ ,

		$\theta^{(1)} = [5, 0]  \theta$	$^{(2)} = [5, 0]$	$\theta^{(3)} = [0, 5]$	$\theta^{(4)} = [0, 5]$	romance
,	Movie	Alice	Bob	Carol	Dave	action
$x^{(1)} = [0.9, 0]$	Love at last	5 <mark>4.5</mark>	5	0	0	
$x^{(2)} = [1, 0.1]$	Romance forever	5 <mark>5.0</mark>	_	-	0	
$x^{(3)} = [0.9, 0]$	Cute puppies of love	_ <mark>4.5</mark>	4	0	-	Each movie has a latent
$x^{(4)} = [0.1, 1]$	Toy story	_ 0.5	_	-	5	factor vector: $\boldsymbol{\theta}^{(j)}$
$x^{(5)} = [0.1, 1]$	Sword vs. karate	0 <mark>0.5</mark>	0	5	-	Each user has a latent factor vector: <b>x</b> <sup>(i)</sup>
$x^{(6)} = [0, 0.9]$	Nonstop car chases	0 <mark>0.0</mark>	0	5	4	Predict the user j's rating to
						movie i by: ( <b>θ</b> <sup>(j)</sup> ) <sup>T</sup> <b>x</b> <sup>(i)</sup>



# **Graph Machine Learning Techniques**

#### Neural Networks for Graphs

Including the graph's connection information in the ML training produces more accurate models and/or reduces the need for explicit features.



Fig. 2: Network Embedding v.s. Graph Neural Networks.

https://medium.com/@terngoodod/a-comprehensive-survey-on-graph-neural-networks-part-1-types-of-graph-neural-network-1dd93b823c70



# Graph Neural Networks





- Combines the added insight from connected data with the modeling power of neural networks
- Uses the graph during training; in-database training is the ideal.

# Summary

- Natural Data Model Graph is how we think
- Richer Data connections between entities, graph-based features
- Graphs have always had a **natural role in machine learning**:
  - Unsupervised learning through **graph algorithms**, frequent pattern mining
  - Graph features provide richer training data
  - Graph embeddings transform graph data into easier-to-use vectors
  - Learning through **neural networks** and deep learning
- Graph data models are uniquely qualified to provide **explanatory AI**.
- Native Graphs with Massively Parallel Processing like TigerGraph enable large scale feature extraction and in-graph analytics

