



Introduction to Graph Machine Learning

Rizal Fathony
rizal@fathony.com

Rizal Fathony



Lead Data Scientist

Grab | Trust, Identity, and Safety

2021-Present | Indonesia



Visiting Lecturer

Universitas Gadjah Mada | DIKE

2022-Present | Yogyakarta, Indonesia



Post-Doctoral Researcher

Bosch | Bosch Center of Artificial Intelligence

2020-2021 | Pittsburgh, Pennsylvania, USA



Post-Doctoral Research Fellow

Carnegie Mellon University | Computer Science Department

2019-2021 | Pittsburgh, Pennsylvania, USA

Education



Ph.D. in Computer Science

University of Illinois at Chicago

2014-2019 | Chicago, Illinois, USA



Master's in Computer Science

University of Illinois at Chicago

2012-2014 | Chicago, Illinois, USA



Bachelor's in Statistical Computing

Sekolah Tinggi Ilmu Statistik

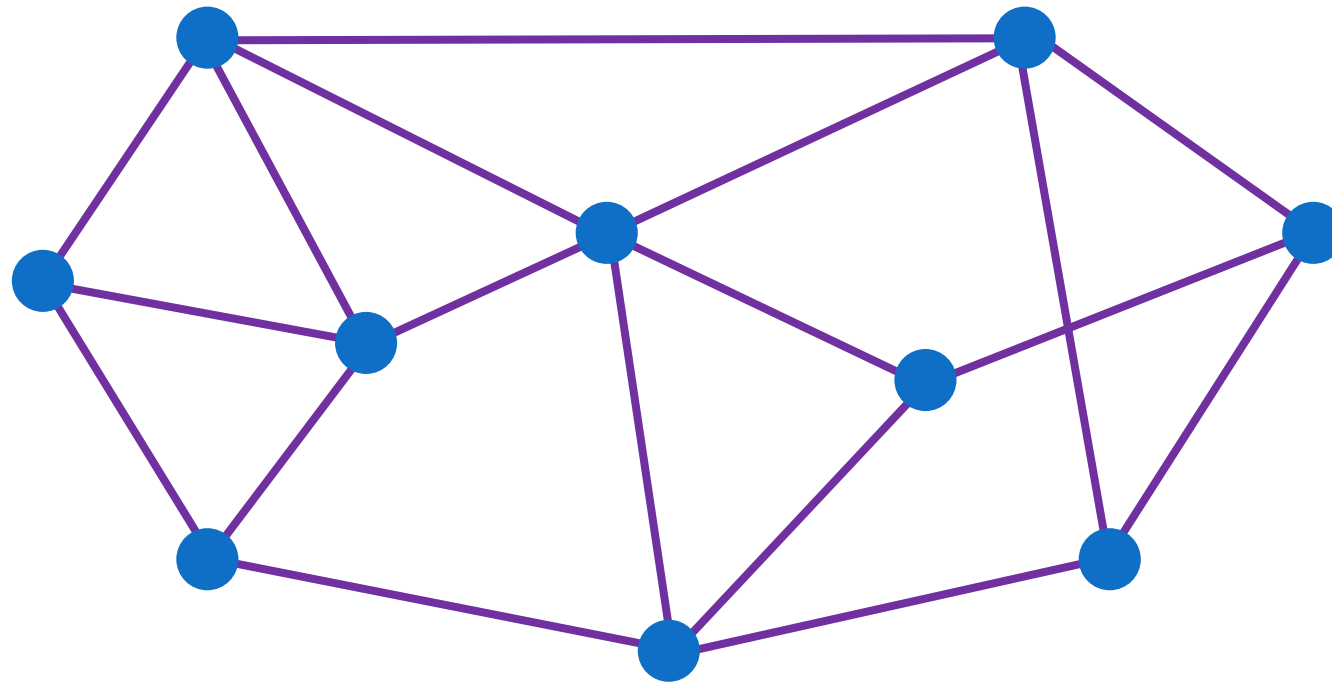
2003-2007 | Jakarta, Indonesia

What is a Graph?

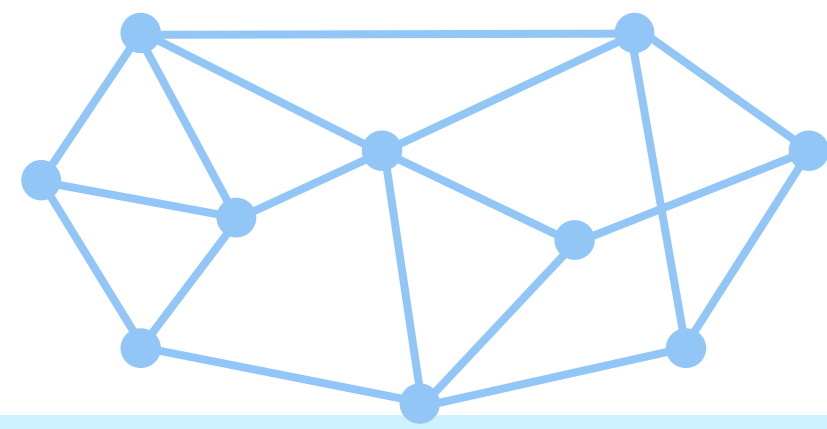
Mathematical structure for pairwise relations

Nodes

Edges

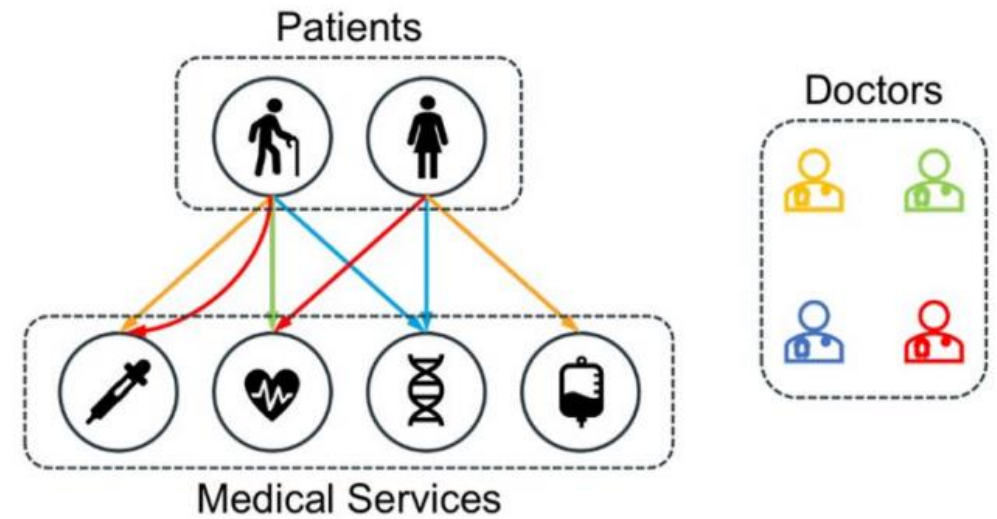
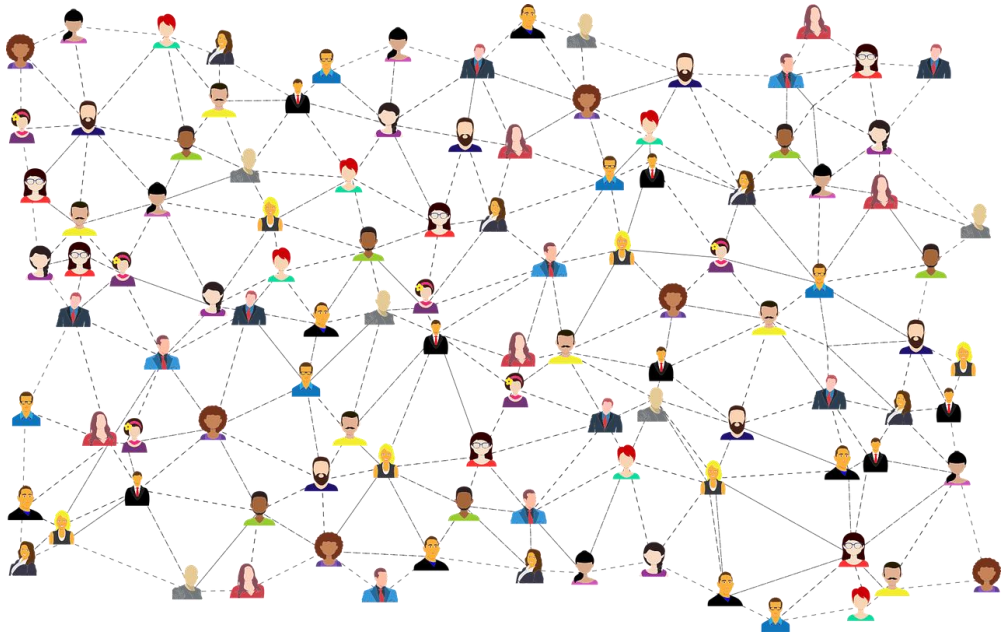


Why Graphs?



Graph is a general method
for describing and modeling
complex systems

Many Data are Graphs



Bipartite Graph

Social Networks

Nodes: Person/Account

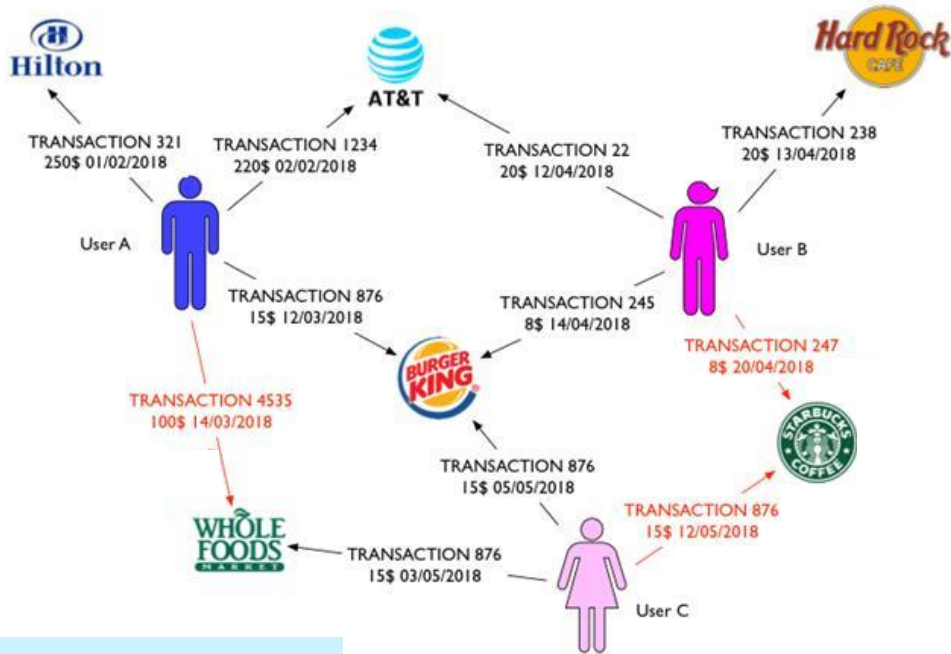
Edges: Friendship/Follows

Health Records

Nodes: Patient, Medical Service

Edges: Treatment by Doctors

Many Data are Graphs

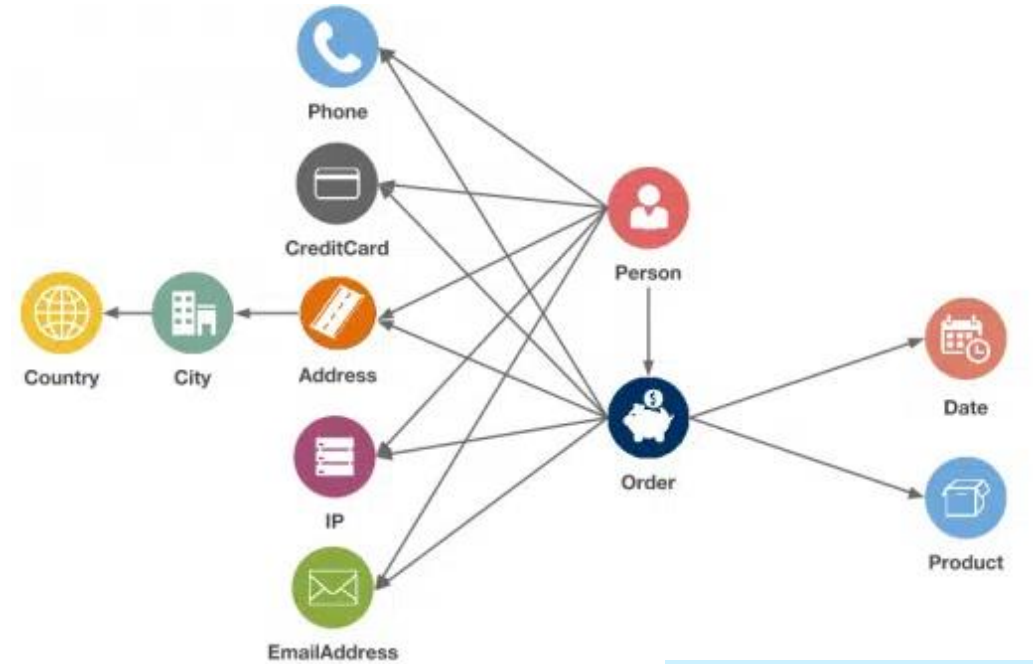


Directed Graph

Financial Transactions

Nodes: Customer, Merchant

Edges: Transaction/Payment



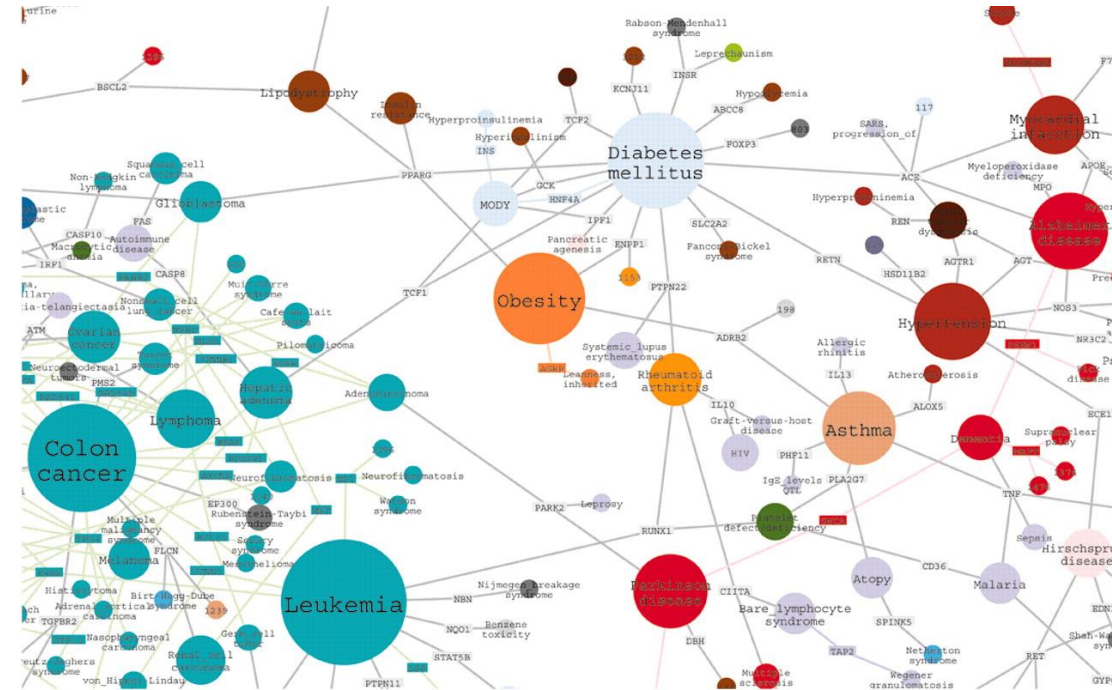
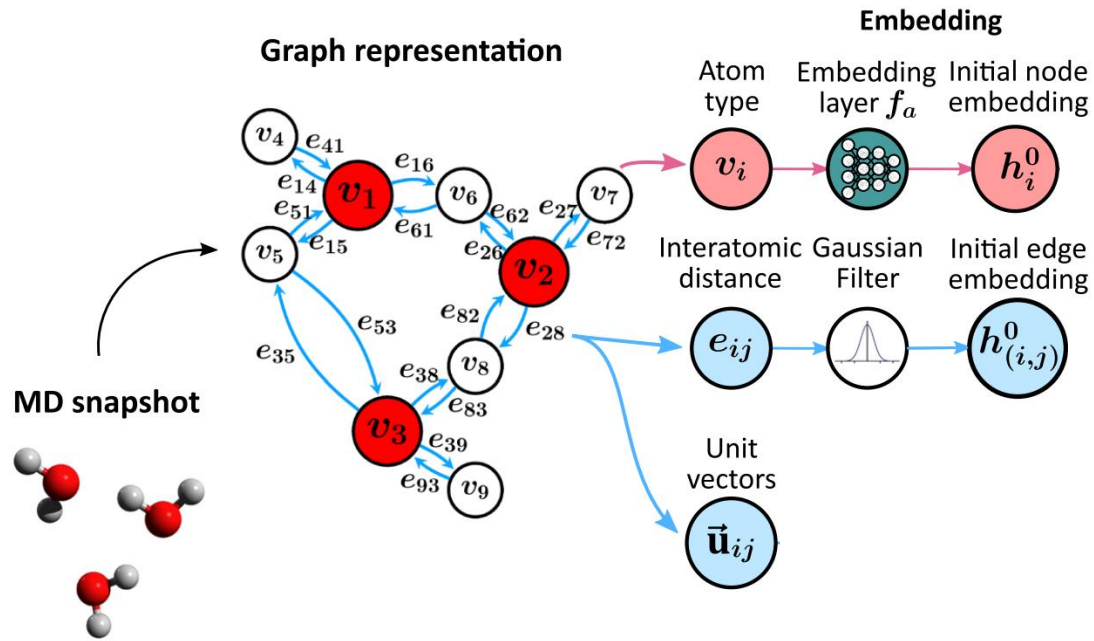
Heterogenous Graph

E-Commerce Data

Nodes: Person, Product, Credit Cards, ...

Edges: Has Phone, Has Address, Orders, ...

Many Data are Graphs



Molecular Modeling

Nodes: Atom

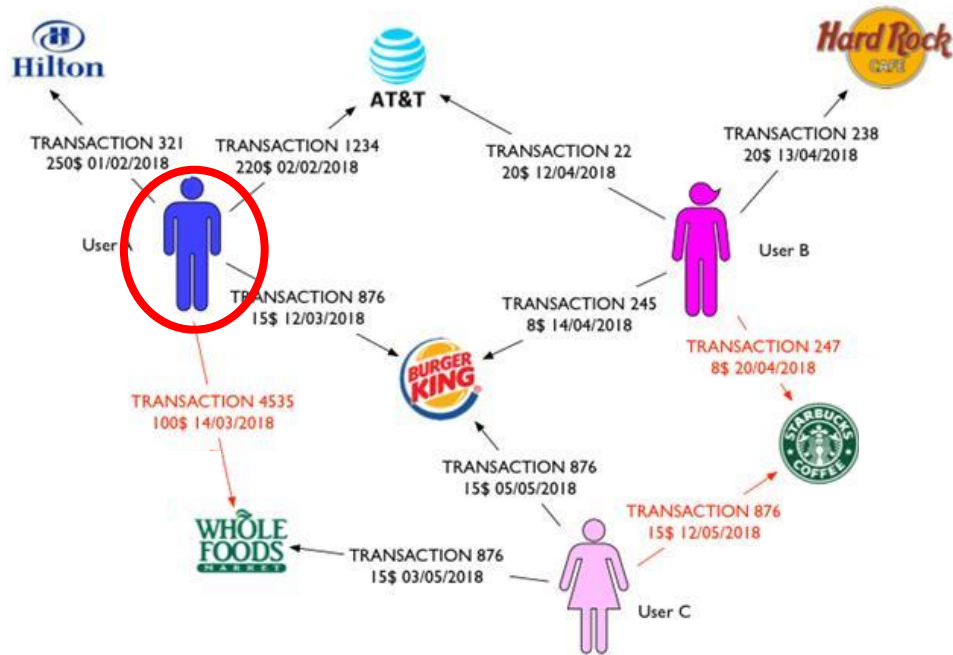
Edges: Chemical Bond

Human Disease Network

Nodes: Disease

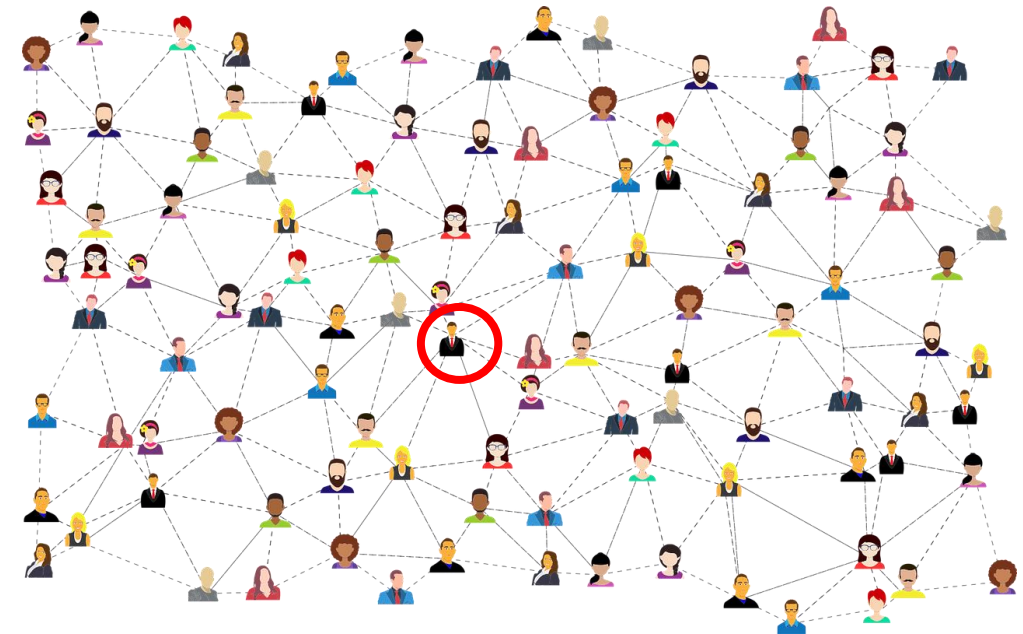
Edges: Genetic Link

Task Example: Node Prediction



Input: Transactional Graph

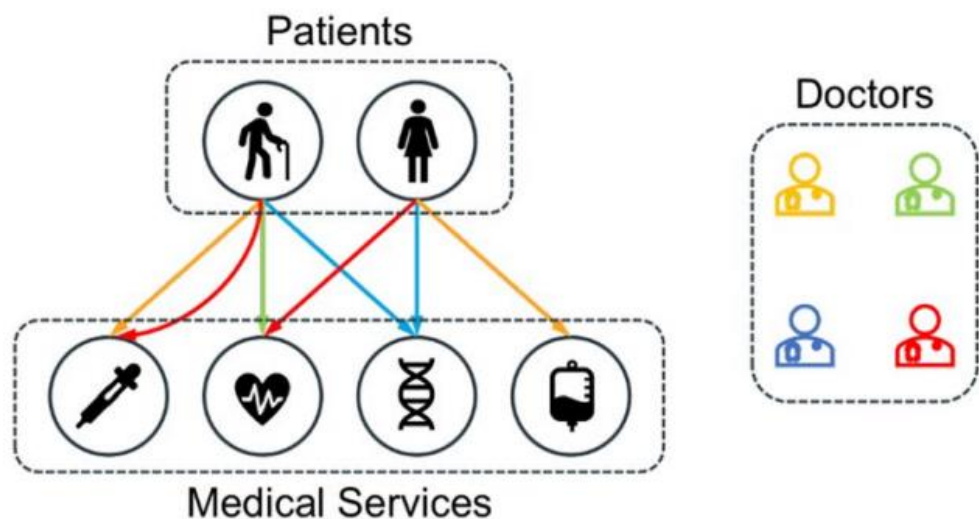
Task: Find user that use stolen credit card in the transactions



Input: Social Network

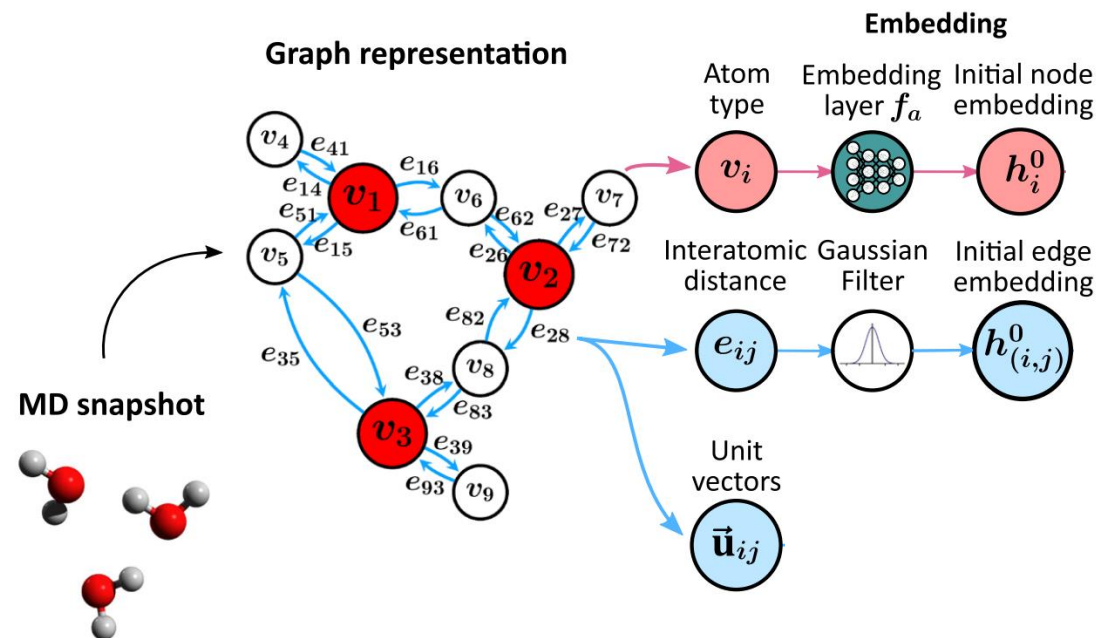
Task: Identify fake user with influence power

Task Example: Edge Prediction



Input: Health Records Graph

Task: Predict if a patient need to see a doctor for medical treatment



Input: Molecular Graph

Task: Predict how strong the chemical bonds' force for a given molecule

What's Next?



Introduction



ML Algorithms



Node Embedding

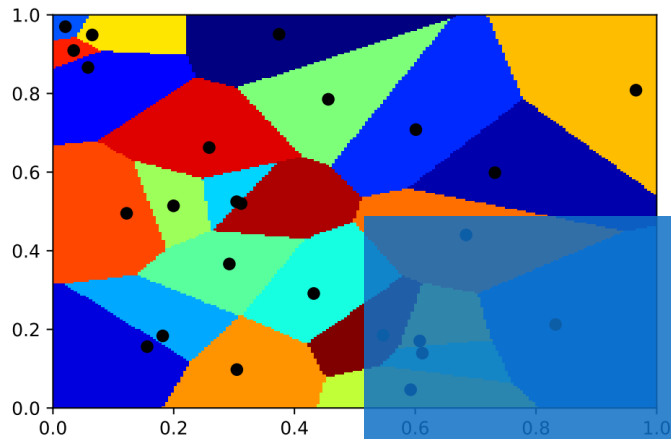


Graph Neural Networks

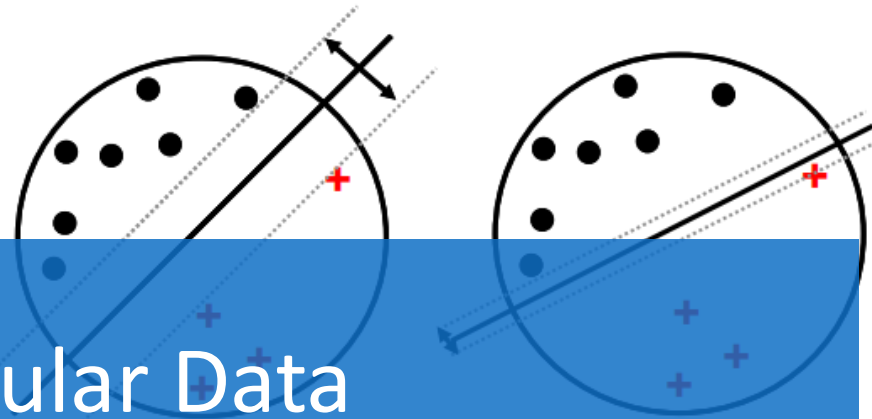
Machine Learning Algorithms

From Classical to Graph ML

Classical Machine Learning Algorithms

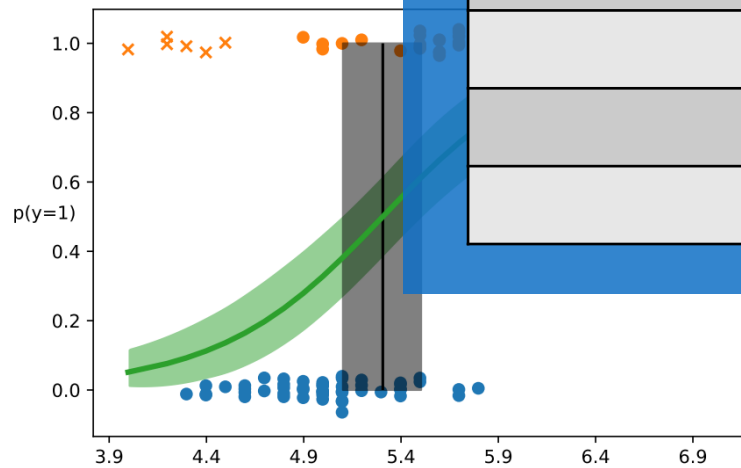


K-Nearest Neighbors

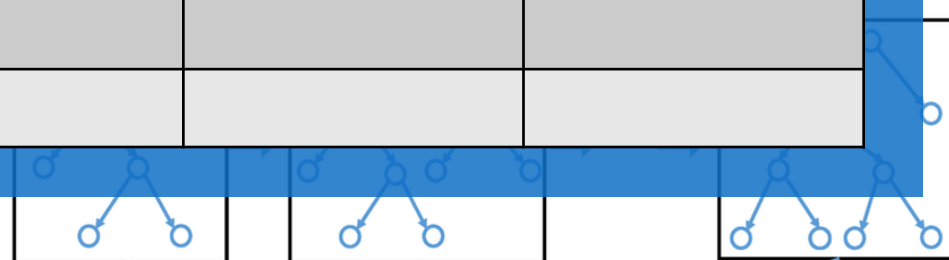


Tabular Data

Feature 1	Feature 2	Feature 3	Feature 4



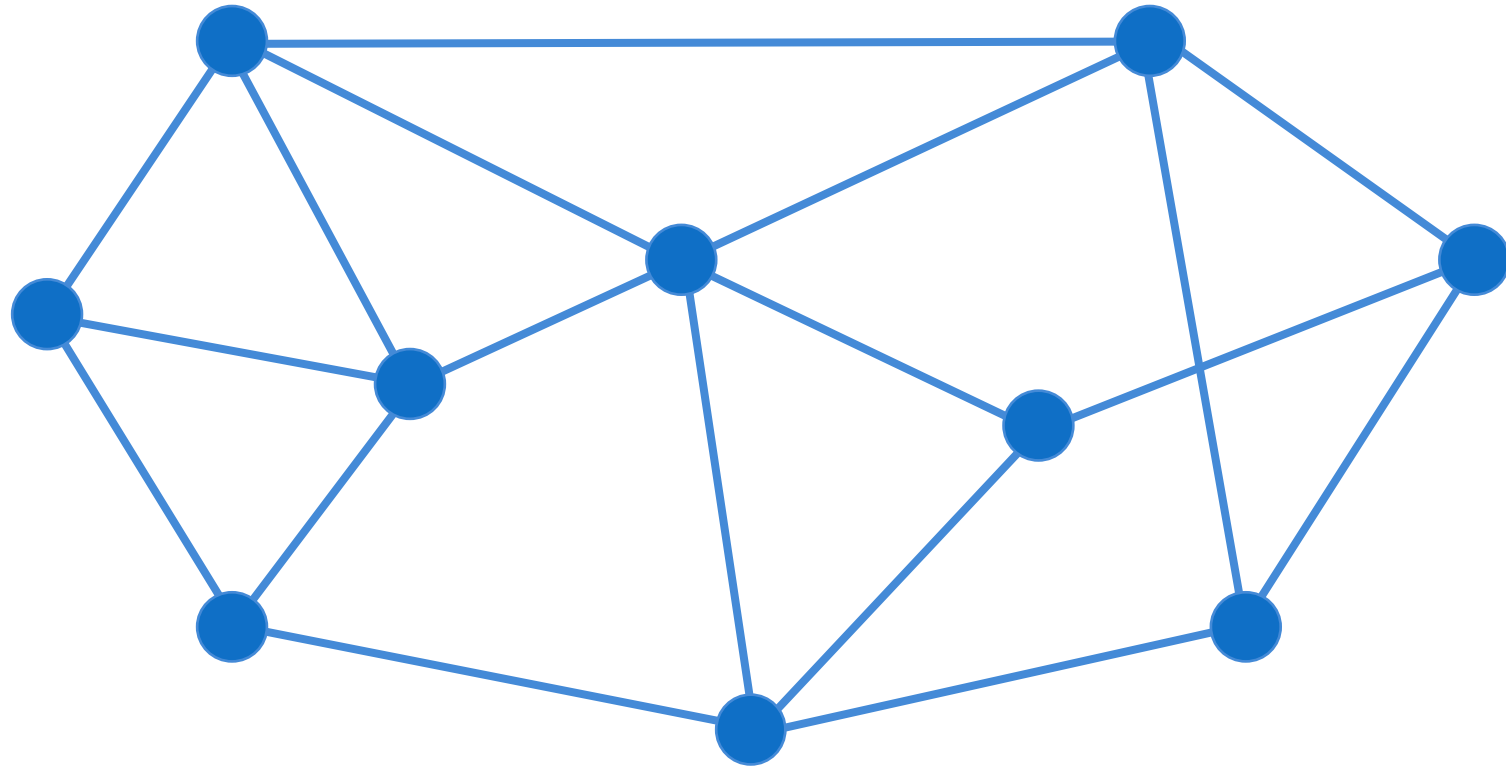
Logistic Regression



$$\hat{y} = \sum_{k=1}^n f_k(x)$$

(Boosted) Decision Tree

Classical ML for Graph Data?



Graph

Classical ML for Graph Data?

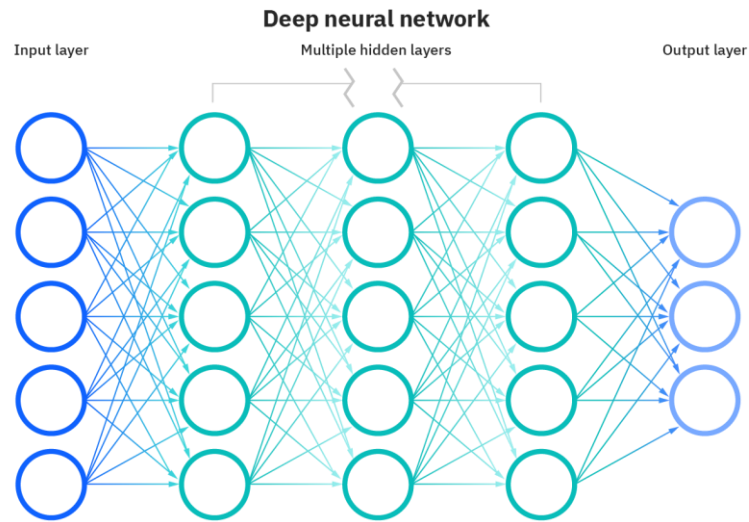
The diagram consists of a 10x10 grid of cells. The top row is solid blue. The remaining 9 rows are light blue with a white grid pattern. To the left of the grid are eight blue circles. A central blue box with white text 'Information Loss' is overlaid on the grid, covering approximately the middle four rows and the middle four columns.

●									
●									
●									
●									
●									
●									
●									
●									

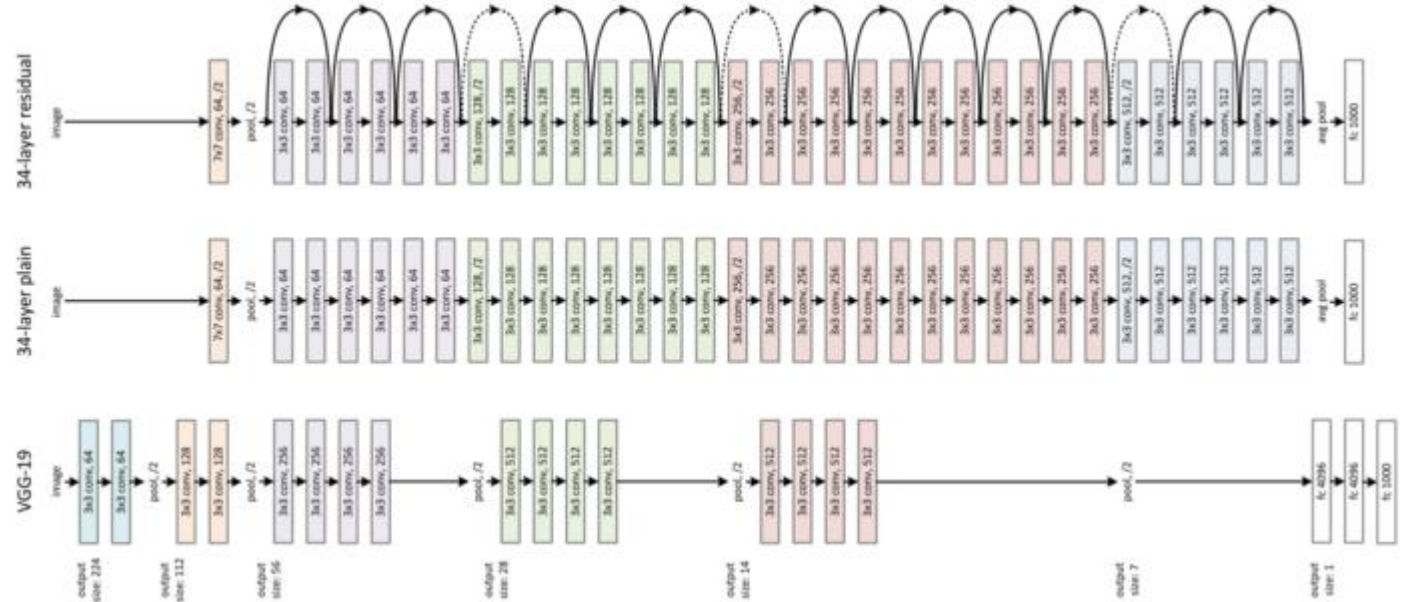
Tabular

Neural Networks and Deep Learning

Multiple layers of learning



Multi Layer Perceptron (MLP)

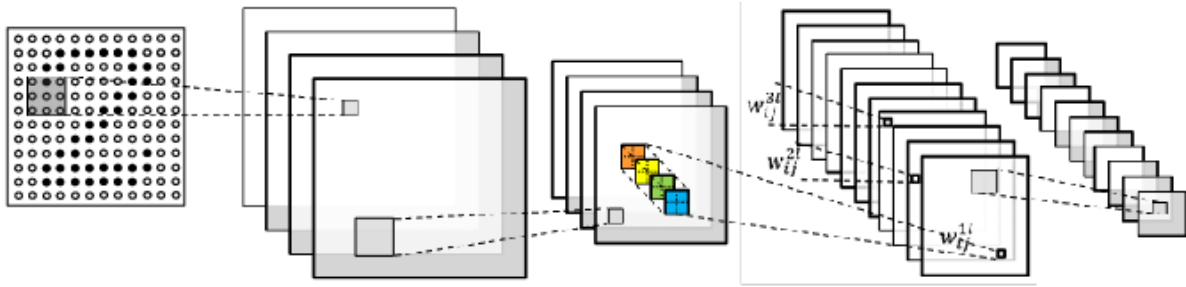


Residual Networks (ResNet)

Capable to learn from “raw” data

Grids and Sequences

Grids

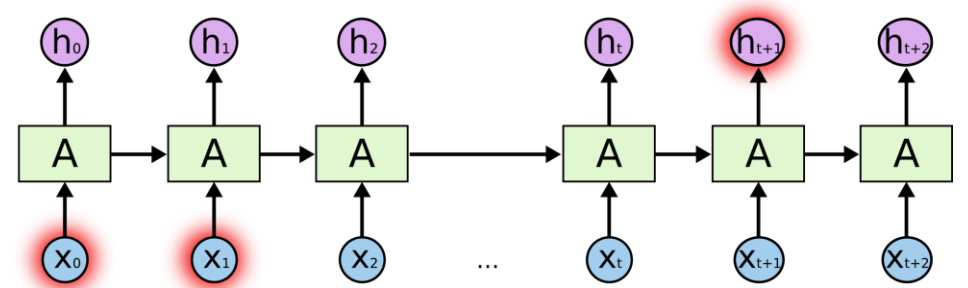


Convolutional Neural Networks (CNN)

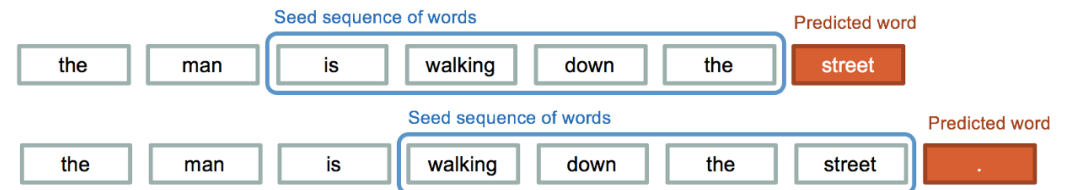


Images

Sequences



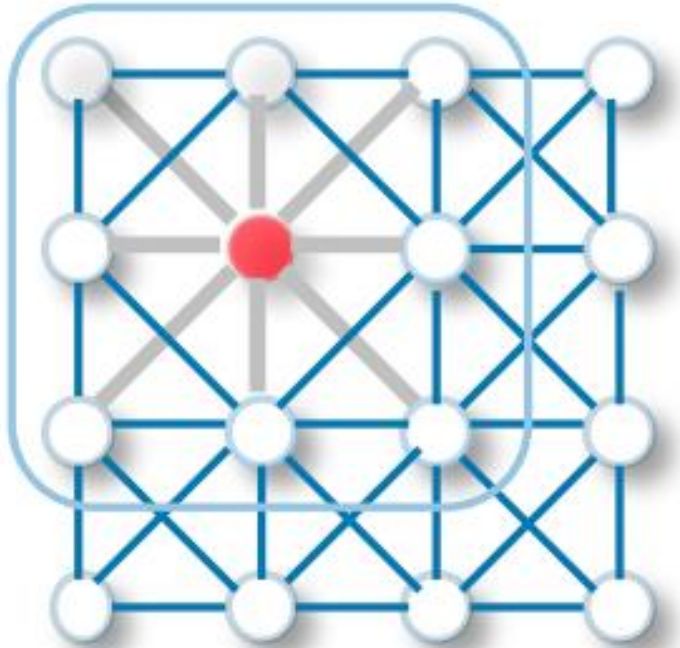
Recurrent Networks (RNN)



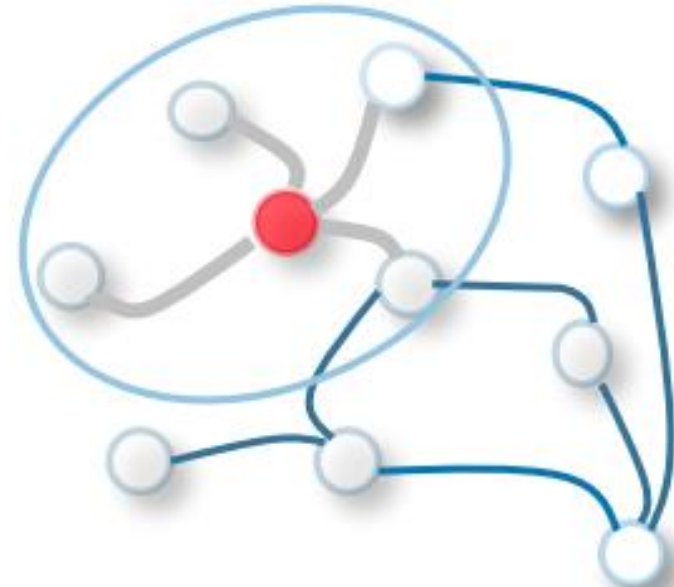
Text

16

Grid and Sequence as Graph



Grid Computation Flow



Graph Computation Flow

Node Embedding

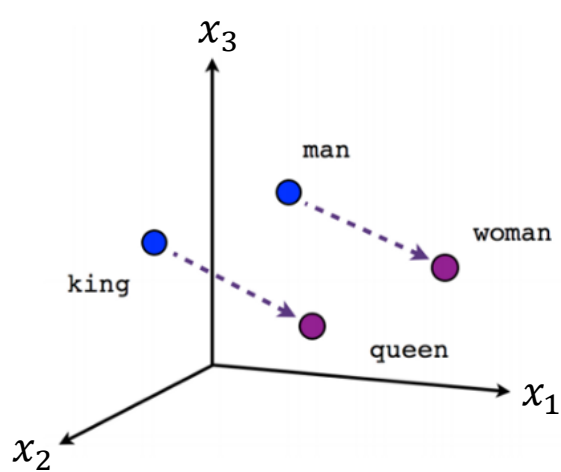
Nodes + neighbors \rightarrow numbers

Inspiration from word embedding: word2vec

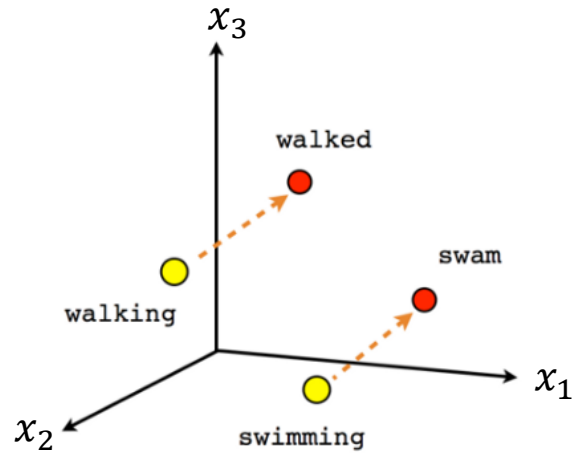
Map words to numerical features

similar word \rightarrow similar values

preserve word associations



Male-Female



Verb tense

king - man + woman \approx queen

walked - walking \approx swam - swimming

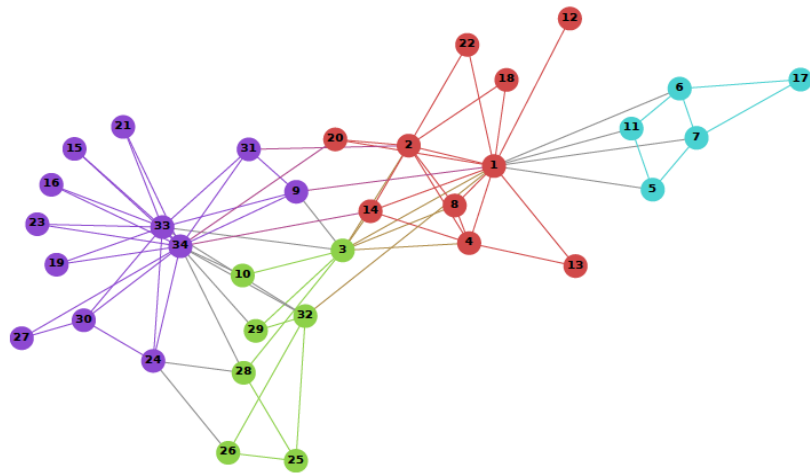
word2vec training process:
predict the neighboring words

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. \rightarrow	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. \rightarrow	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. \rightarrow	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. \rightarrow	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

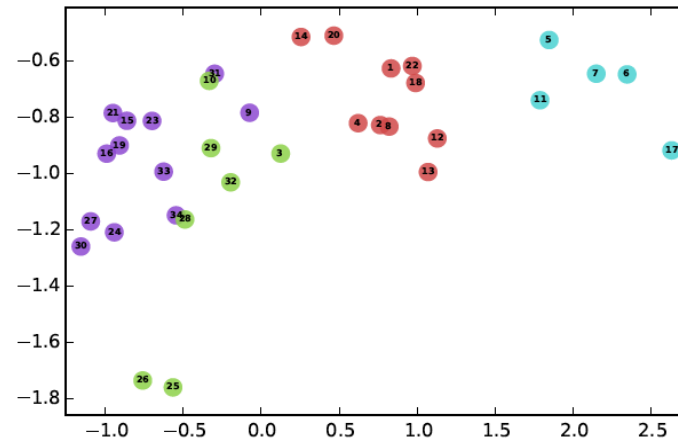
Node embedding algorithm

Map nodes in a graph to **numerical features** (embedding)

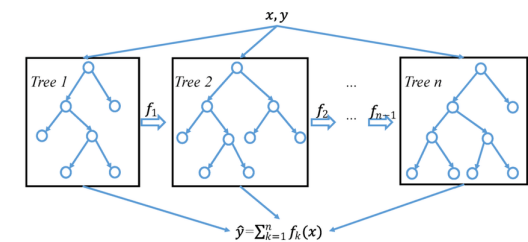
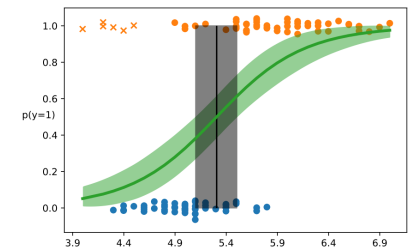
similar nodes → **similar embeddings**



Graph



Embedding



Classifier
(ML models)

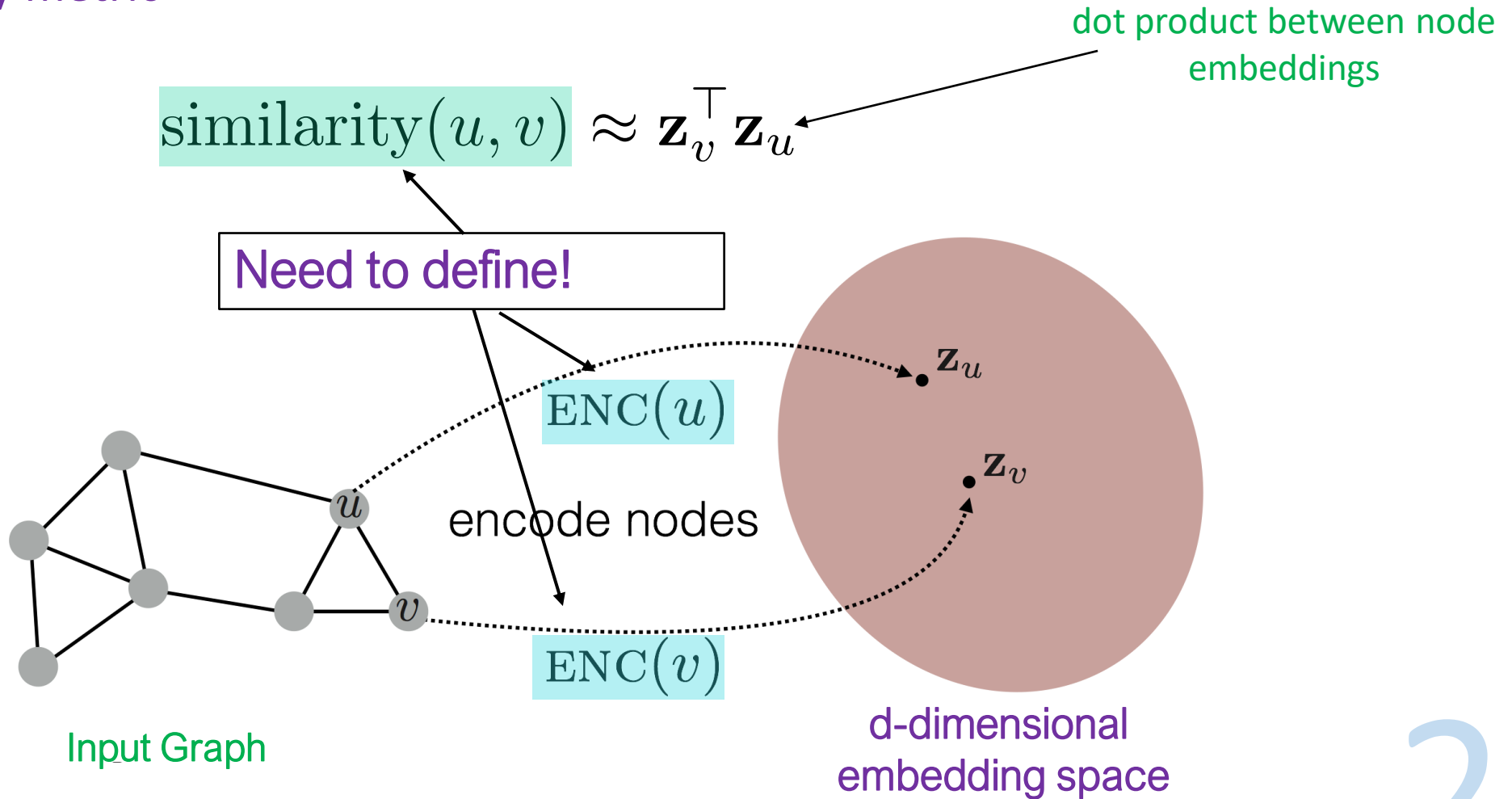
20

Node embedding components

1. Encoder

2. Similarity metric

3. Optimization Algorithm



Simple node embedding (example)

1. Encoder

“Shallow encoder”

Each node is assigned a unique embedding vector (i.e., we directly optimize the embedding of each node)

used by most node embedding models

$$\text{ENC}(v) = \mathbf{z}_v$$

node in the input graph

d-dimensional embedding

2. Similarity metric

Key ingredient that differentiate node embedding methods

- Adjacency based similarity

if two nodes are connected via an edge = similar

Adjacency matrix

$$\mathbf{A}_{u,v} = 1, \text{ if } u \text{ and } v \text{ are connected via an edge}$$
$$\mathbf{A}_{u,v} = 0, \text{ otherwise}$$

3. Optimization Algorithm

Loss function:

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \| \mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v} \|^2$$

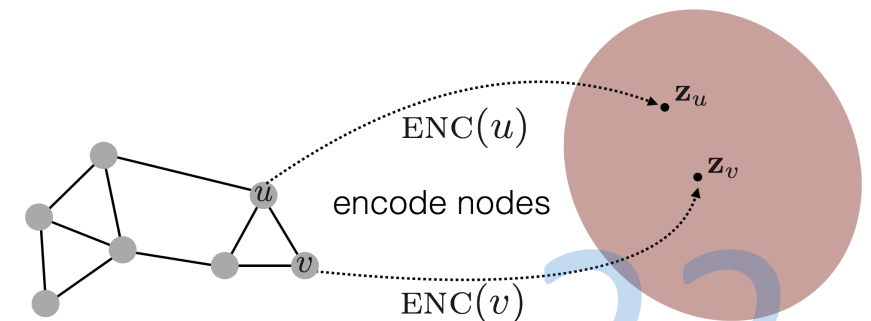
loss (what we want to minimize)

sum over all node pairs

embedding similarity

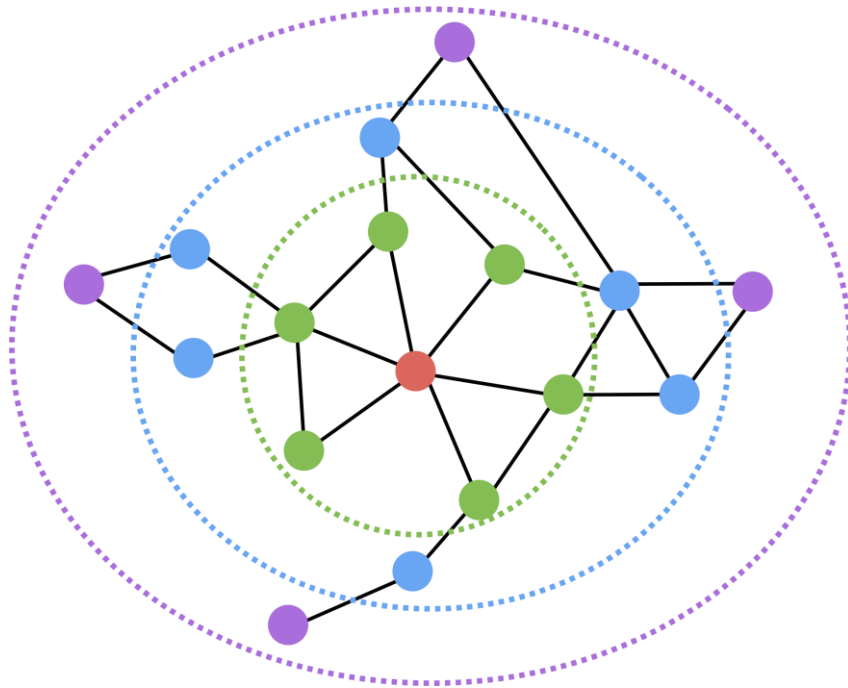
(weighted) adjacency matrix for the graph

Optimize with SGD!



K-Hop Similarity

Neighboring nodes achievable in k-hop should have similar embedding



- **Red:** Target node
- **Green:** 1-hop neighbors
 - \mathbf{A} (i.e., adjacency matrix)
- **Blue:** 2-hop neighbors
 - \mathbf{A}^2
- **Purple:** 3-hop neighbors
 - \mathbf{A}^3

Objective:

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \|\mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v}^k\|^2$$

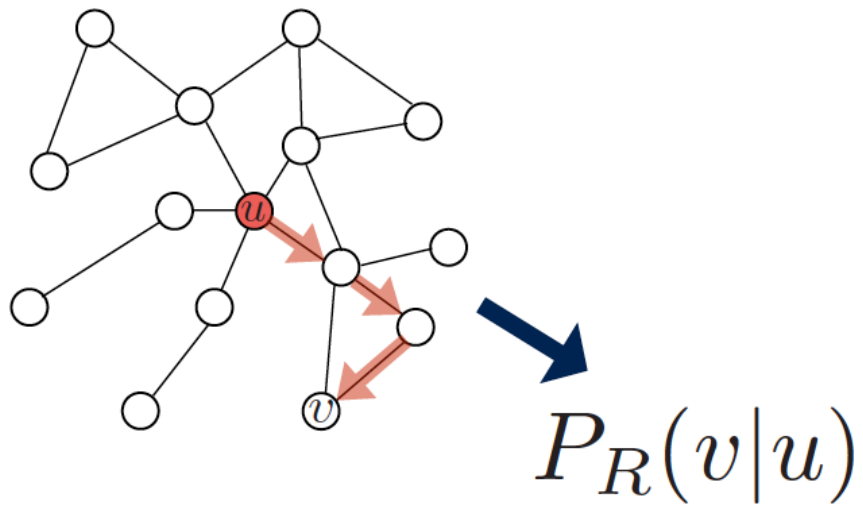
adjacency matrix's k-th power = k-hop

Theorem: Raising an adjacency matrix \mathbf{A} of simple graph G to the n -th power gives the number of n -length walks between two vertices v_i, v_j of G in the resulting matrix.

Random Walk Similarity

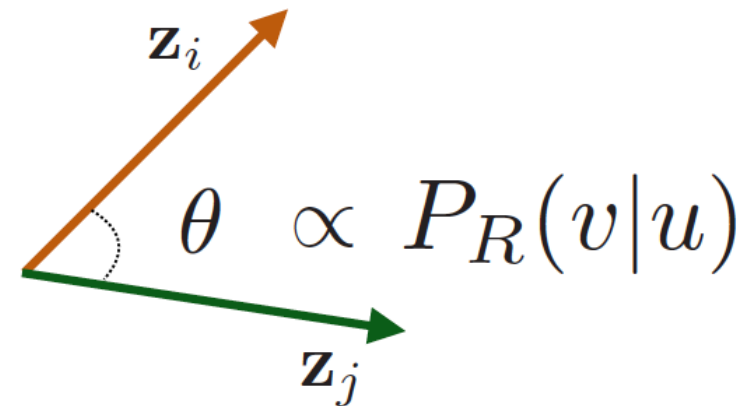
Random walk: start from node u , repeatedly jump (walk) to a neighboring node

$P_R(v|u)$: probability of visiting node v from random walks starting from node u



A random walk from u to v

Embedding similarity should approximate $P_R(v|u)$



Cosine similarity

Random Walk Optimization

Objective

maximize likelihood of
random walk co-occurrences

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log(P(v|\mathbf{z}_u))$$

For each node u collect $N_R(u)$, the set of nodes visited on random walks starting from u .

Parameterize

use softmax

$$P(v|\mathbf{z}_u) = \frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)}$$

Optimization: minimize \mathcal{L}

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log \left(\frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$

sum over all nodes u sum over nodes v seen on random walks starting from u predicted probability of u and v co-occurring on random walk

Nested sum over nodes gives $O(|V|^2)$ complexity!!

Approximate the normalization constant

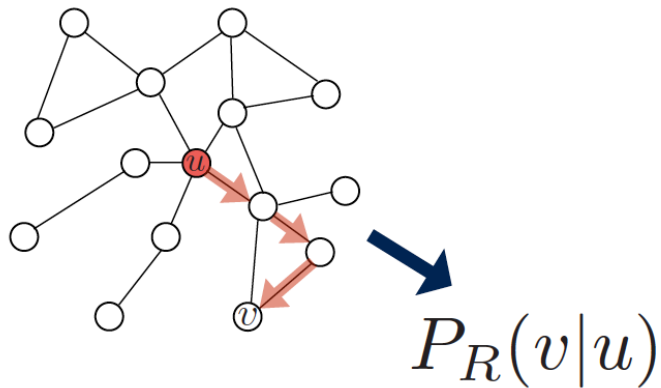
Negative sampling

- Most pairs of nodes are not connected (negative sample)
- Instead of normalizing w.r.t. all nodes, just normalize against k random negative samples.

Why random walk?

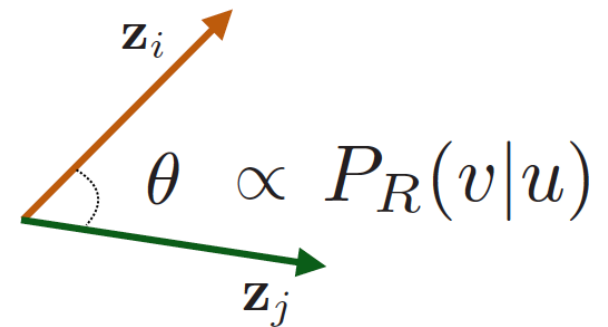
Efficiency

Do not need to consider all node pairs when training; only need to consider pairs that co-occur on random walks.



Expressiveness

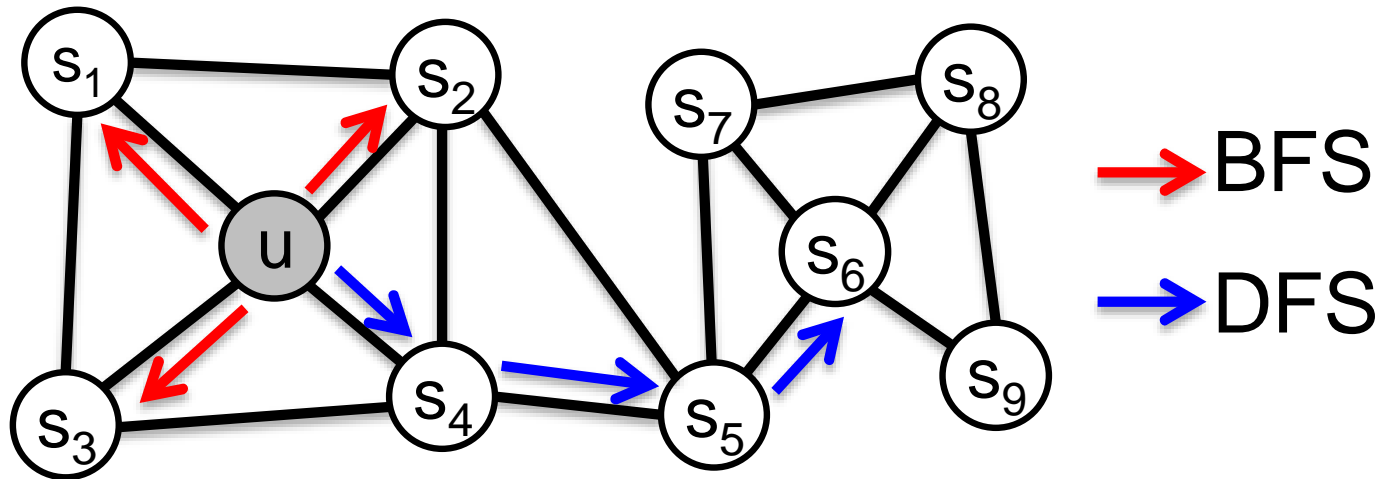
Flexible stochastic definition of node similarity that incorporates both local and higher-order neighborhood information.



node2vec: biased random walk similarity

Biased random walk to encourage:
local and **global** views

local microscopic view → breadth first search (BFS) walk
global macroscopic view → depth first search (DFS) walk



$$N_{BFS}(u) = \{s_1, s_2, s_3\} \quad \text{Local microscopic view}$$

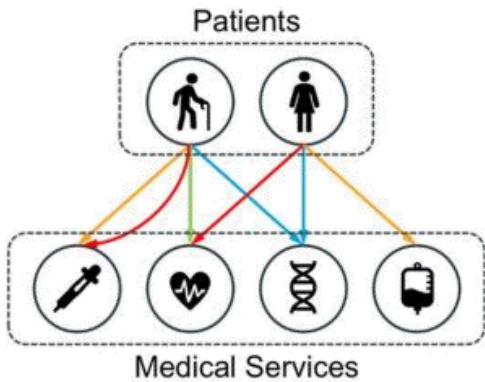
$$N_{DFS}(u) = \{s_4, s_5, s_6\} \quad \text{Global macroscopic view}$$

Application: Health Records

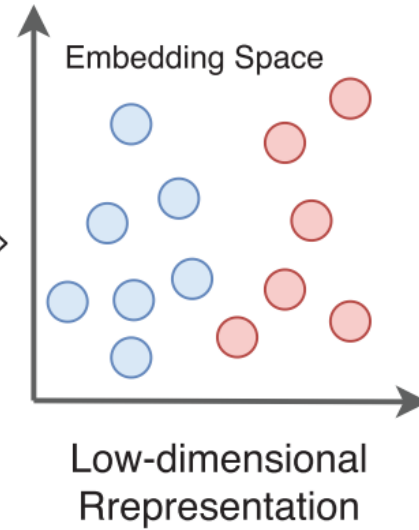
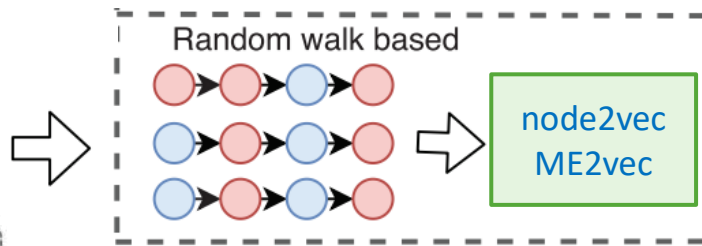
Leveraging graph-based hierarchical medical entity embedding for healthcare applications

Tong Wu et.al (2021) [Nature Scientific Reports | Advanced Analytics, IQVIA Inc]

node2vec, ME2vec

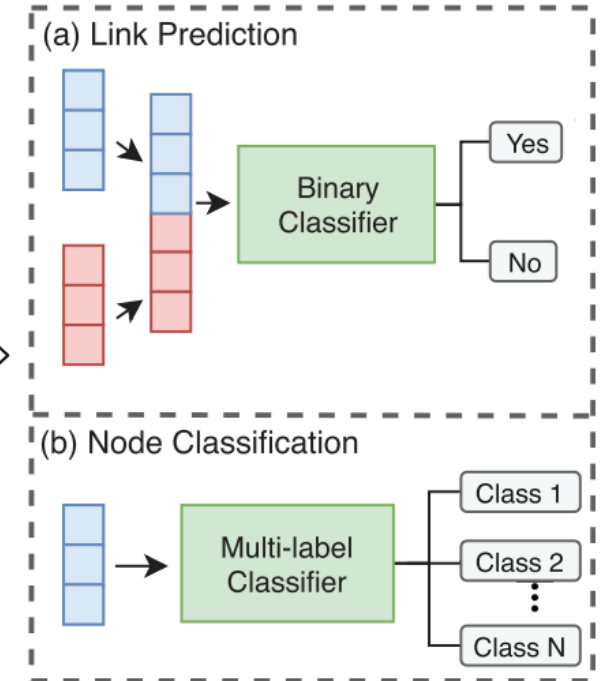


patient – medical service graph
patient – doctor graph



patient, medical service,
doctor embedding

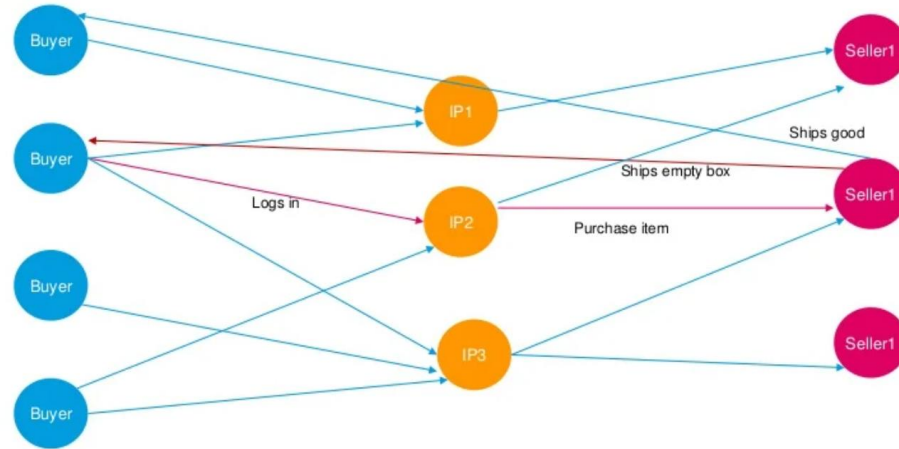
Logistic Regression



Downstream Prediction Tasks

- Prediction Tasks:
1. Predict patient diagnostic [node classification]
 2. Predict if patient need to see a doctor [link prediction]
 3. Readmission prediction [node classification]

Application: PayPal's Collusion Fraud Prevention



DATA

Training Data:

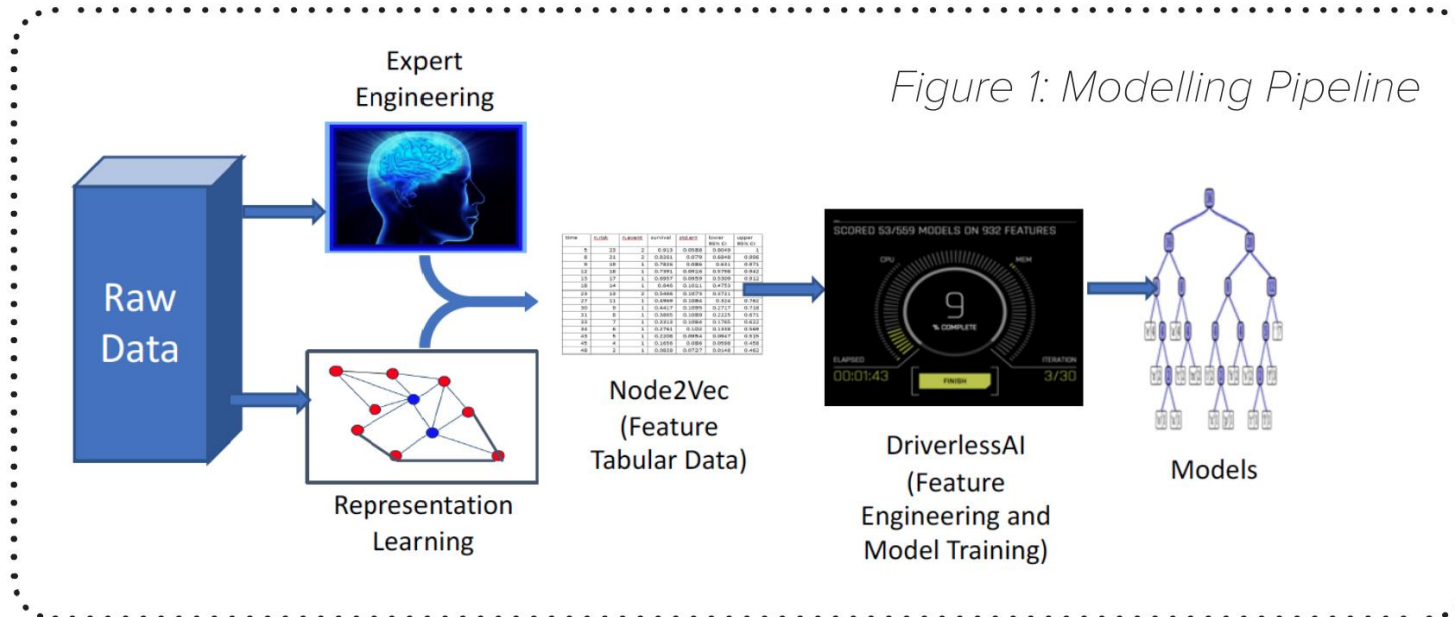
- Subset of one year's transactions.
- 1.5 billion edges, .5 million nodes.

Test Data:

- 3 months

Number of Features:

- 400-600



Thank You