

Jure Leskovec slides - Stanford

<http://cs224w.stanford.edu>

## Machine Learning with Graphs

01-INTRO

Book = Will Hamilton - Graph Representation Learning

Google Colab

PyTorch = Facebook  
(-)

KERAS/Tensorflow = Google

Thesis = U. Montreal

$\begin{matrix} (+) \\ \text{KERAS} \end{matrix} = \text{NN}$   
Tensorflow general

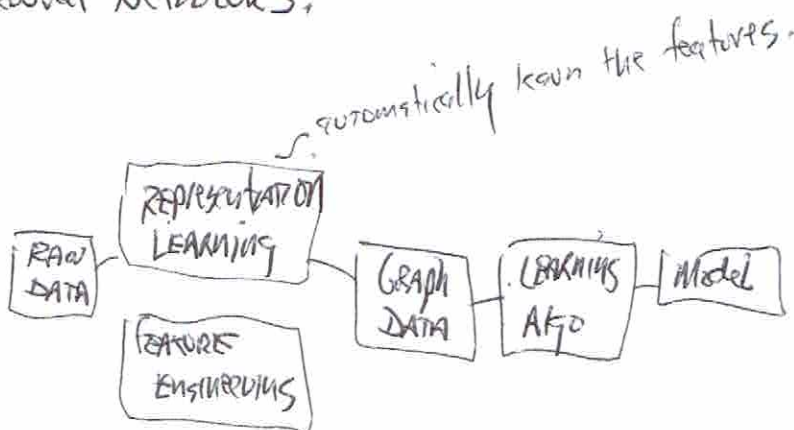
PyG = PyTorch Geometric

LIBRARY for Graph Neural Networks.

Graph Gym.

SNAP, NetworkX

Traditional NN



map nodes to  $d$ -dimensional embeddings s.t. similar nodes are embedded close together

$F: V \rightarrow \mathbb{R}^d$  Representation  
FEATURE representation = embedding.

Various Topics:

Traditional methods — Graphlets, Graph kernels.

Node Embedding Methods — DeepWalk, Node2Vec.

Graph NN — GCN, GraphSAGE, GAT, Theory & GNNs

Knowledge Graphs = TransE, BetaE

Deep Generative Models = TransE, BetaE Graph RNN

Applications

## TASKS (Graph)

Graph-level prediction  
Graph generation.

node, edge, subgraph.

## TASKS (ML)

Node classification  $\rightarrow$  Predict property of node  $\rightarrow$  online banners/items  
Link Prediction  $\rightarrow$  Link exist/missing  $\rightarrow$  knowledge graph completion.  
Graph classification  $\rightarrow$   
Clustering  
Generation  
Evolution -

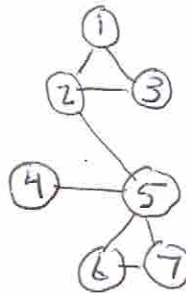
## Projected Bipartite Graph

### Heterogeneous Graphs

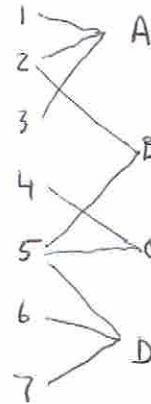
Bipartite  
Desner

"Projected" Bipartite Graph.

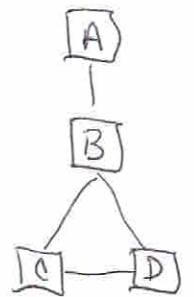
Projection U



~~Projection~~  
U V



Projection V



Adjacency Matrix

Adjacency List

Node/Edge Attributes...

Weights

Multigraph

Connected

Block-Diagonal form (Non-connected, "Components")

Strongly/Weakly Connected,

"In-Component" Our component  
SCC "Strongly Connected Component"

02 - Traditional ML

Node, Link, Graph Prediction -

Traditional ML Pipeline:

Node Features  $\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^D$

Graph Features  $\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^D$

Link Features  $\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^D$

Given a new < Node, Link, Graph > Predict ...

Features are key ...

TRADITIONAL

Hand-Designed Features

(New)

(?)



Node Classification

Structure & Position of node in network.

degree  
centrality  
clustering coefficient  
graphlet =  $\Delta$   $\square$

Katz Index = # of walks of length  $l$

Graph Kernel Method ... design kernels instead of feature vectors.

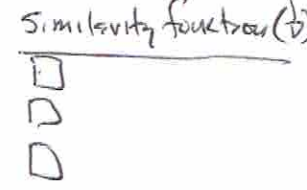
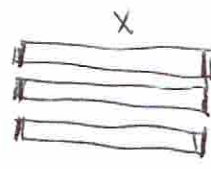
- Remember the  $i$ th training example
- Instance-based Learning.

$$\begin{pmatrix} \bar{x}_n \\ y_n \end{pmatrix} \rightarrow w_n$$

"Similarity function" = "kernel"

given input  $x^*$

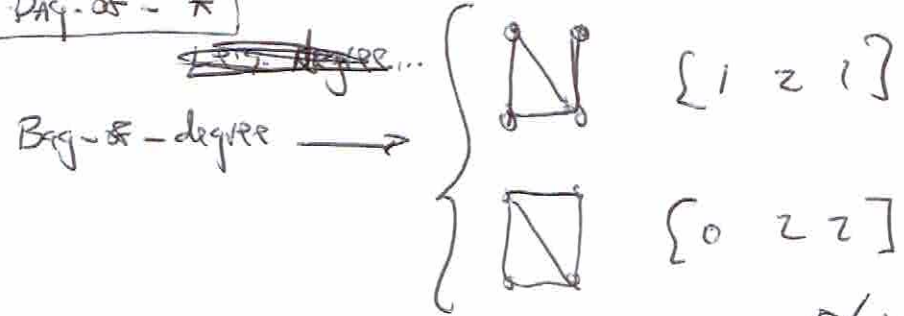
Trainings



Replaces S.V.M.

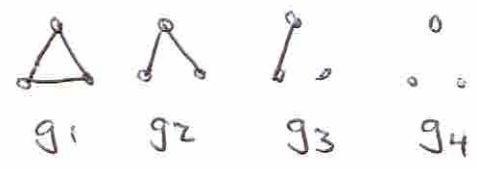
OR - Traditional - ML

Bag-of-words for our graph kernel



"Graphlet kernel" and Weisfeiler-Lehman - ~~Kernel is based on~~ colour assignment  
Computationally efficient

Count graphlets ...

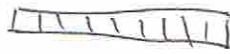


Feature vector is  
Count of graphlets.

Embeddings: Why?

Map nodes  $\rightarrow$  Embeddings

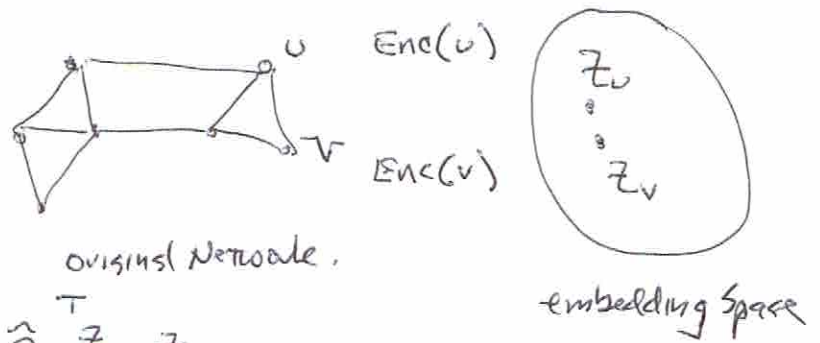
Similarity & embeddings indicates similarity



embedding  $\rightarrow$  Tasks:

- Node classification
- Link prediction
- Graph classification
- Clustering
- ...

$\rightarrow$  Goal is to map nodes to embedding space s.t. similarity (dot product) approximates similarity in graph.



$$\text{Similarity}(u, v) \approx z_u^T z_v$$

$$(1 \times N)(N \times 1) \Rightarrow 1$$

## Learning Embeddings

Encoder: node  $\rightarrow$  Embedding

Similarity function = similarity in original Network

Decoder:

Embeddings  $\rightarrow$  Similarity score.  $\text{DEC}(z_u^T z_v)$

Optimize parameters & encoder s.t.

$$\text{Similarity}(u, v) \approx \underbrace{z_u^T}_{\text{orig. NW.}} \underbrace{z_v}_{\text{embedding space.}}$$



Shallow encoding - each node assigned a unique vector

Random walk embeddings via  
Stochastic Gradient Descent...  
node2vec...

} see Google and Ferrara 2017 survey



"Embedding entire graphs" / Subgraphs

Run on subgraph then average  
Virtual node to represent subgraph  
Anonymous walk

Hierarchical Embeddings (clusters)

PageRank

Spider Trap, Teleport.  
Node proximity (consumer, item)  
Matrix factorization.

04 - PageRank

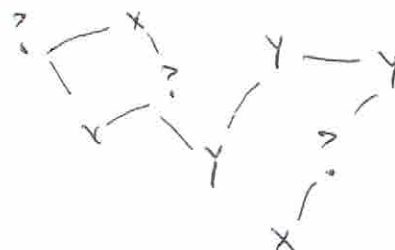
Message Passing & Node Classification

05 - Message Passing

Given a network with labels on some nodes how do we assign labels on other nodes?

"Semi-supervised"

X = trusted  
Y = fraudster  
? = unknown / predict.



Method 1 = node embeddings.

Method 2 (today) = message passing

- Relational classification
- Iterative classification
- Correct & Smooth

Correlations in network

Guilt-by-Association

Homophily

Individual

↓  
Social

Tendency to  
Bond. Birds  
& a feather

Influence

Social

↓  
Individual

Social connections  
can influence  
individual. I try  
~~what~~ what my  
friends like.

05. message

Probabilistic Relational Classifier = Propagate node labels across network.

Iterative Classification - classify based on attributes of node as well as labels of neighbor nodes

(does not use node attributes - only labels)

Correct & Smooth

Given partially labeled graph

- ① Train base classifier predictor
- ② Use base predictor to predict node labels
- ③ Post-process using graph structure to get final node values

06 - Graph Neural Networks

Insuition = map nodes to d-dimensional embeddings s.t. similar nodes are "close"

Encoder

Encode  $(u, v)$   
Similarity function  
Similarity  $(u, v)$

$f(\text{input graph}) =$

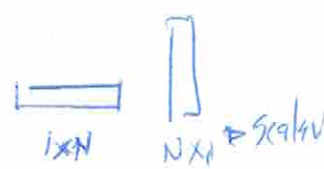
$z \rightarrow$  node embeddings



Similarity  $(u, v) \approx z_v^T z_u$

similarity in original network

(dot product of embeddings)



"Shallow" Encoder

"Encoder"





66-GRAPH-NN

8.  
9

# Deep Graph Encoders.



## "Aggregate Neighbors"

Each node defines a computation graph based on neighbors.

## "Deep Encoder"

- 1) Instance vs class (1st order vs 2nd order)
- 2) Forward/Backward (Loss & Diagnostic)
- 3)

66 How to Train a GNN... (Slide 66)

GNN vs CNN

GNN vs Transformer

## BASICS of Deep Learning

15. input LABEL  
x y  
vector  
sequence  
matrix  
graph

## Optimization Problem

$$\min_{\theta} L(y, f(x))$$

Loss function  
Set of parameters we optimize

Example Loss Function = Cross Entropy

y = categorical encoding  
[0 0 1 0 0] "onehot"

$$CE(y, f(x)) = - \sum_{i=1}^C (y_i \log f(x)_i)$$

Actual Predicted

the lower the loss the closer the prediction is to one-hot.

How to Optimize...

$$\nabla_{\theta} L = \frac{\partial L}{\partial \theta_n}$$



30

## Deep Learning for Graphs

- Local Network Neighborhoods
- Stacking multiple Layers

\* Graph  $G = (V, A, X \in \mathbb{R}^{m \times |V|})$

$\underbrace{\hspace{10em}}_{\text{Adjacency Matrix}} \quad \underbrace{\hspace{10em}}_{\text{node features}}$

$v$  is a node in  $V$

$N(v)$  is neighbors of  $v$

### Naive Approach:

Join Adjacency matrix + features

$\Rightarrow O(|V|)$  parameters.

	A	B	C	D	features	
					$f_1$	$f_2$
A	0	1	1	0	1	0
B	1	1	0	1	0	1
C	1	1	1	0	1	1
D	1	0	0	0	1	1

Convolution = generalize and apply across all.

(But graph has no "Locality" or "Sliding window")

\* Graph is "permutation invariant"  $\Rightarrow$  no pre-set ordering.

"Order Perm"  $O\text{-Perm} = A_1 X_1$      $O\text{-Perm} = A_2 X_2$

A permutation invariant function is needed.

### Permutation Equivariance

"Graph Representation"  $\Rightarrow$

"Node Representation"  $\Rightarrow$

function maps graph  $G = (A, X)$  to  $\mathbb{R}^d$

function maps node to  $\mathbb{R}^{m \times d}$

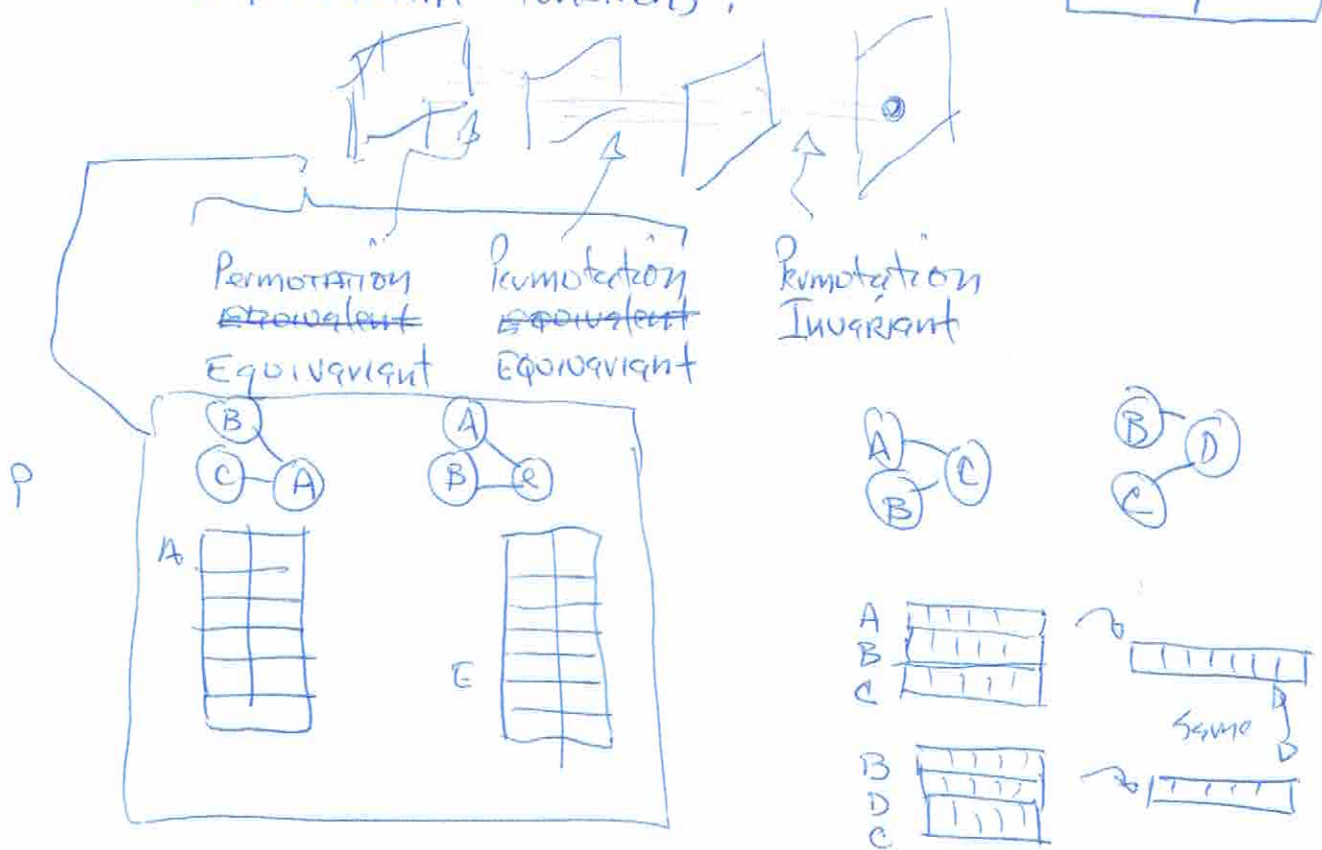
embedding  $d$

nodes  $m$


42

A Graph Neural Network consists of multiple Permutation Equivalent / Invariant functions.

06 Graph NN



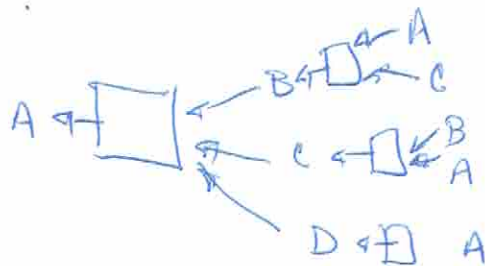
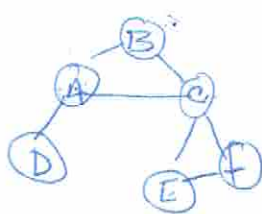
51

## GRAPH CONVOLUTIONAL NETWORK

~~IDEA: Computational Neighborhood~~

- IDEA 1: Nodes Neighborhood defines a "computational Graph"
- 2: Generate node embeddings based on Local network neighborhood.
- 3: Nodes Aggregate information from their neighbors

54

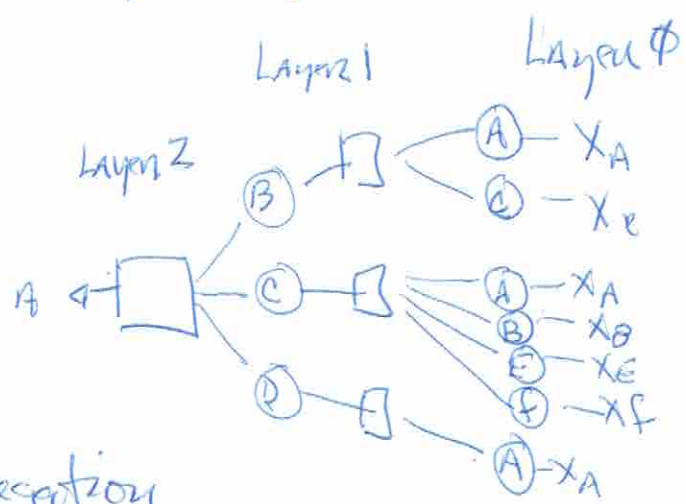


56 Deep Model: Many Layers

Nodes have embeddings at each layer.

Layer 0 embedding of node is its feature set  $X_v$

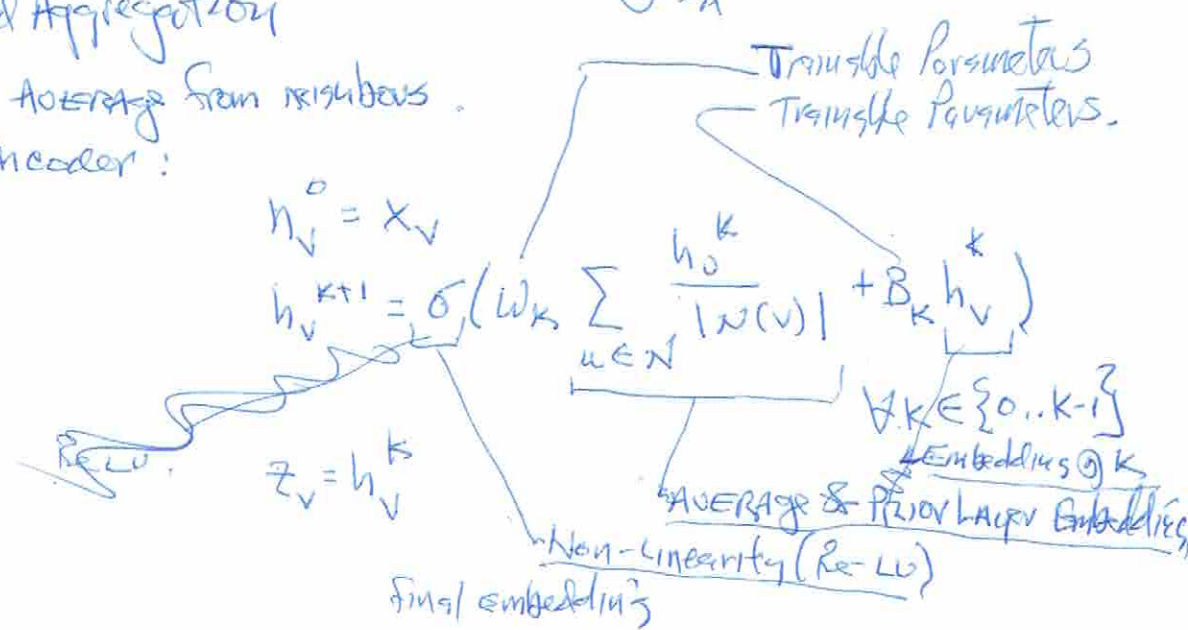
Layer-k embedding gets information from nodes that are k-hops away.



Neighborhood Aggregation

BASIC: Average from neighbors.

DEEP Encoder:



59

63

## 06 Graph NN

How To TRAIN?

64

Need A Loss Function...

Matrix Formulation...

66

~~Supervised~~ How to train =

Supervised...  $\mathcal{L}$  could be L2 (numeric) or Cross-Entropy (categorical)

Unsupervised... use the graph structure.

74

Inductive capability...

76

GNNs subsume CNNs and Transformers

CNN can be seen as a special case of GNN where neighbor (structure) is fixed.  
CNN is not permutation independent.

Transformer can be seen as a GNN that runs on a fully connected "word" graph.



## GT GRAPH NN - Z

Project: GRAPH NN using PyG.

PyG.org - Library for graph learning pipelines

SLACK channel  $\rightarrow$  Q: Conferences are ICML NeurIPS  
ICLR KDD WWW

Blog Posts: Should contain

① Step-by-step explanation of Graph ML Techniques  
Assume reader familiar w/ ML, deep learning, PyTorch.  
Not familiar w/ graph ML, PyG.

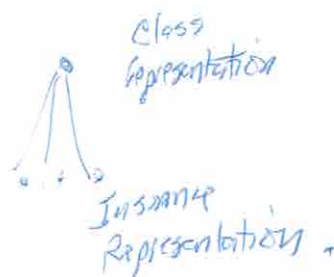
② Visualization - more better.  
Code snippets of PyG/PyTorch.

\*\*\* ③ Link to Google Colab.

④ Finding Graph ML models  $\leftarrow$

open Graph Benchmark -  
LEADER BOARD - OGB -  
ogb.stanford.edu/docs/leader\_nodeprop

\*\*\*  $\rightarrow$  OGB is like KAGGLE for GRAPH NN, and run by Stanford.



First comes in as an instance

Several come  $\rightarrow$  establish a class

Instances get spawned into "context". Appropriate disappropiate.

Each instance is associated with "features" to the class, and to the evidence



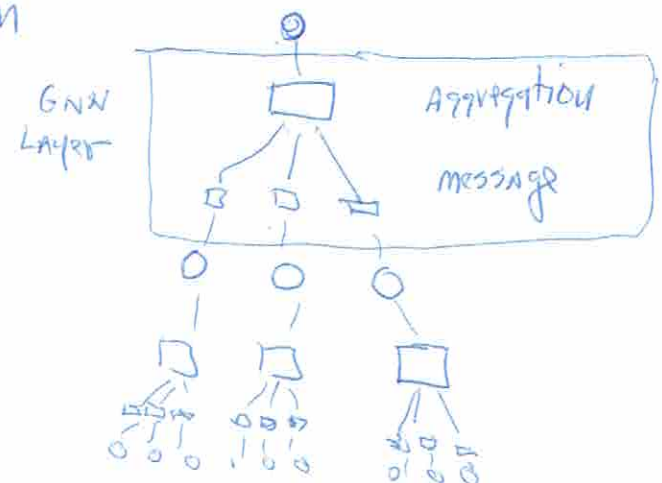


Recap: Deep Graph Encoder  
~~Graph Neural Network~~

Node's neighborhood defines a "computation graph".  
Aggregate from neighbors.

GNN Layer = message + Aggregation

GCN, GraphSage, GAT



Message =  $m_v^{(l)} = \text{MSG}^{(l)}(h_v^{(l-1)})$

each node creates a message which will be sent to other nodes later.

$m_v^{(l)} = W^{(l)} h_v^{(l-1)}$

AGGREGATION

Node Aggregates messages from neighbors -  
Information from node itself could get lost

$$h_v^{(l)} = \text{CONCAT} \left( \text{AGG} \left( \sum_{u \in N(v)} m_u^{(l)} \right), \underbrace{m_v^{(l)}}_{\text{self}} \right)$$

Nonlinearity (  $\sigma(\cdot)$ ,  $\text{ReLU}(\cdot)$ ,  $\text{Sigmoid}(\cdot)$  )

Example:

GCN

GraphSAGE

Long Short-term Memory

Aggregation functions = mean, pool, Lstm, L2 Normalization

GAT = Graph Attention Network.

## Graph Attention Networks

42

## Attention Mechanism

## GNN Layers in Practice

Batch normalization - stabilize re-center (zero-mean)

Dropout - regularize to prevent overfitting,  
Attention / Gating

Activation - ReLU or similar Parametric ReLU



\*\*\*

## Graph Gyr

59

Over-smoothing - all the embeddings converge to the same value

Receptive Field = set of nodes that determine (contribute to)  
the embedding of a node.In a  $k$ -layer GNN each node has a  
receptive field of  $k$ -hop neighborhood.

D Avoid too many layers.

65

Expressive Power for Shallow GNNs

Within each layer - each aggregation/transition is a NN.

Skip connections

## GNN Augmentation & Pitfalls

Stacking Layers  
 Over-smoothing  
 Receptive field  
 Do not set # Layers too large.  
 Shallow GNN.

14 Skip connections

## 17 Graph Augmentation for GNNs

Idea: Raw input graph  $\neq$  computational graph.

Reason to consider input graph  $\neq$  computational graph.

- a) graph too sparse
- b) graph too dense.
- c) graph too large.

Lacks features  $\rightarrow$  feature augmentation

to sparse, dense, large

sparse  $\rightarrow$  add virtual nodes

dense  $\rightarrow$  ~~add~~ sample neighbors.

large  $\rightarrow$  sample subgraphs.

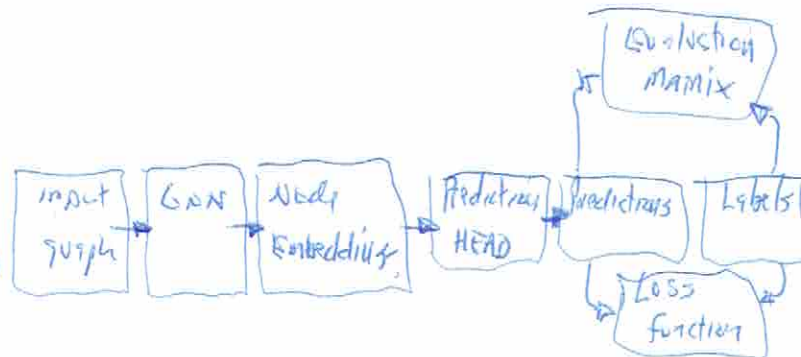
23 Feature augmentation

28 Add virtual nodes

30 Node neighborhood sampling

## Prediction With GNNs

Training Pipeline 1- Prediction HEAD



34

38

## Prediction Heads

Different tasks require different prediction heads.

Node prediction — directly

Edge prediction — pair-wise

Graph prediction — using all nodes.

Global pooling loses information:

global mean, max, sum.

Hierarchical Global Pooling

Training Pipeline (2) = Predictions, Labels

Supervised, Unsupervised.

Training Pipeline (3) = Loss Function.

54

Classification Loss vs Regression Loss.

Classification Loss  $\rightarrow$  cross-entropy =  $y_j^{(n)} \log(y_j^{(n)})$

—  $n$ th data point

—  $j$ th class

Regression Loss (mean squared error) MSE

a.k.a. L2 Loss.

$$= \sum_{j=1}^K \left( y_j^{(n)} - \bar{y}_j^{(n)} \right)^2$$

Training Pipeline (4) = Evaluation Metrics.

## Evaluation Metrics: Regression

(use sklearn)

RMS E Root Mean-Squared-Error

$$\sqrt{\frac{\sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2}{N}}$$

MAE Mean Absolute Error

$$\frac{\sum_{i=1}^N |y^{(i)} - \hat{y}^{(i)}|}{N}$$

## Evaluation Metrics: Classification

① multiclass classification (Accuracy)  $\frac{1}{N} \left[ \text{argmax}(\hat{y}^{(i)} = y^{(i)}) \right]$

② Binary Classification

Accuracy ~~Precision~~  
Precision | Recall

If range of prediction is [0, 1] we use 0.5

Metrics Agnostic to classification threshold

ROC AOC

Recall =  $\frac{TP}{TP + FN}$

number predicted  
+ of actual

the number correctly predicted

Precision

$\frac{TP}{TP + FP}$

or those predicted +  
how many are actually +

ACTUAL

+

-

Predicted

+

-

TP

FN

FP

TN

Accuracy =  $\frac{TP + TN}{| \text{Dataset} |}$

F1-Score

$\frac{2PR}{P+R}$

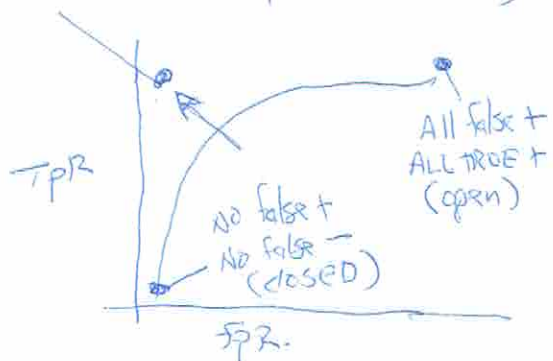
\*\*\* sklearn classification Report



**ROC CURVE**: Captures the tradeoffs

OS GNN Application

All true + no false +  
In TPR  $\neq$  FPR as classification threshold varies (for a binary classifier)



$$TPR = Recall = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

		Actual	
		+	-
Pre dict	+	TP	FP
	-	FN	TN

$$ROC = \frac{TPR}{FPR} = \frac{TP}{TP + FN} \cdot \frac{FP + TN}{FP}$$

Dataset split. — this is different for graphs...

Image Database — each (independent)

Graph — not independent.

Node Prediction

~~Transductive~~

Solution 1 = Transductive Setting.

- Only split the node labels.
- All edges (entire graph)
- Only certain nodes ~~are~~ data

Solution 2

- Break up the graph to get multiple graphs

Graph Prediction

For Graph classification only inductive setting.

Link Prediction

Hide some edges.

89 **Resources**

87 **Debugging Deep Networks**

Loss/Accuracy not converging  
overfit on subset of training data  
scrutinize & Loss visualizations

Resources  
Graph Gym  
PyG  
PyTorch  
GraphNets  
TorchVis

85 **When things don't go as planned...**  
data preprocessing  
Adam optimizer  
Activation f  
model dimensions

## 26 How Expressive Are GNNs

## Expressive Power

Characterized by neighbor Aggregation Function used

"GCN" = mean-pool - uses element-wise mean pooling over  
neighboring node features  
(Kipf, Welling)

$$\text{Mean} \left( \{x_u\}_{u \in N(v)} \right)$$

over the neighboring nodes  
mean of features

"GraphSAGE" = max-pool - uses element-wise pooling over neighboring  
node-features  
(Hamilton et al.)

(above are shown to be "not injective" - thus not maximally powerful)

## 62 Injective Multi-Set

7 Heterogeneous Graphs :  $G = (V \in R \ T)$

Nodes  $v_n \in V$

Edges  $(v_i, r, v_j) \in E \rightarrow$  directed

Node Type  $T(v_n)$

Relation Type  $r \in R$

Node has type

Edge has type

8 Relational GCN

directed w/ 1 relation

28 Knowledge Graphs

entities (nodes) w/ Relation (edges)

40 Trans E

50 Trans R

60 DisMul

69 Complex

"Completion" ?

|

## REASONING in Knowledge Graphs

Can we 'complete' the K.G?

Given a head, Relation we predict the missing tails...

e.g. Predict the tail "science fiction" for "J.K. Rowling" "genre"

- Multi-hop queries

Query2Box

- Path queries

ANSWERING Predictive queries

34

KGS USING Box Embeddings

67

TRAINING Query2Box

80

VISUALIZING

\*\*\*

"We use t-SNE to reduce the embedding space to a 2-dimensional space to visualize the query results"

Fast Neural Subgraphs

Building Blocks -  
characterize  
discriminate

- 7
- ① Node-induced subgraph  
subset of nodes, all edges induced by them.
  - ② Edge-induced subgraph -  
subset of edges and all corresponding nodes.

"Network Building Blocks"

- 11
- "Graph Isomorphism"  
check whether two graphs are identical  
1 sub

- 15
- Network Motifs  
Recovering significant patterns -  
ie. "Induced Subgraphs"

- 18
- Subgraph Frequency  
Motif Significance - null model  $\rightarrow$  test "Significance"  
Random Graph i.e. Erdős Rany

- 33
- Neural Subgraph Representations  
Neural architecture for subgraph matching

- 60
- Finding Frequent Subgraphs  
solution = Representation Learning  
Combinatorial Explosion  $\Rightarrow$  Discrete Search Space.
- \*\*\*  
SPMiner = identify frequent motifs  
subgraph isomorphism



GNNs For Recommender Systems

Netflix Amazon Spotify Youtube Pinterest.

movies

Products

Music

Videos

Pins (images)

Bipartite graph

- NGCF - Wang - NEURAL GRAPH COLLABORATIVE FILTERING
- LightGCN - He
- PinSAGE - Ying

TRADITIONAL COLLABORATIVE FILTERING - BASED ON MATRIX FACTORIZATION,  
USES "SHALLOW ENCODERS" FOR USERS/ITEMS

Limitation  $\Rightarrow$  only first-hop relations expressed  
GNN  $\Rightarrow$  does not consider graph structure.  
addressed both.

32 NGCF

39 LightGCN

55 PinSAGE

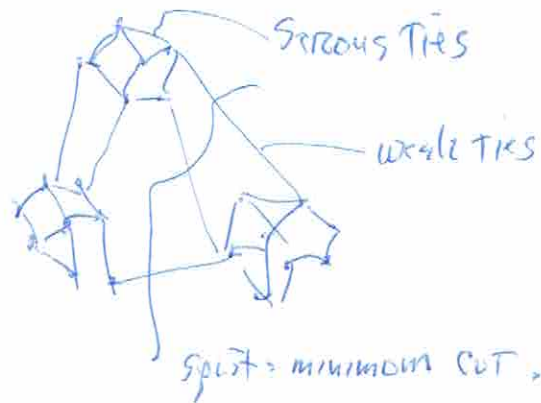
## Community Detection in Networks -

Granovetter - friendship is structural, Interpersonal

Triadic closure  
clusters/groups/modules/  
communities

Eachans KARATE Club Network

find micro-markets -



26 Null Model

30 Louvain Algorithm = GREEDY Algorithm for Community Detection  
 $O(n \log n)$  run time

Algorithm greedily maximises "modularity"

① Local changes to node-community

② communities are aggregated into super-nodes

47 Detecting Overlapping Communities

Community Affiliation Graph Model (AGM)

Node Model

Deep Generative models for Graphs

IDEA: So far we have been learning from Graphs.

Assumed Graphs are given

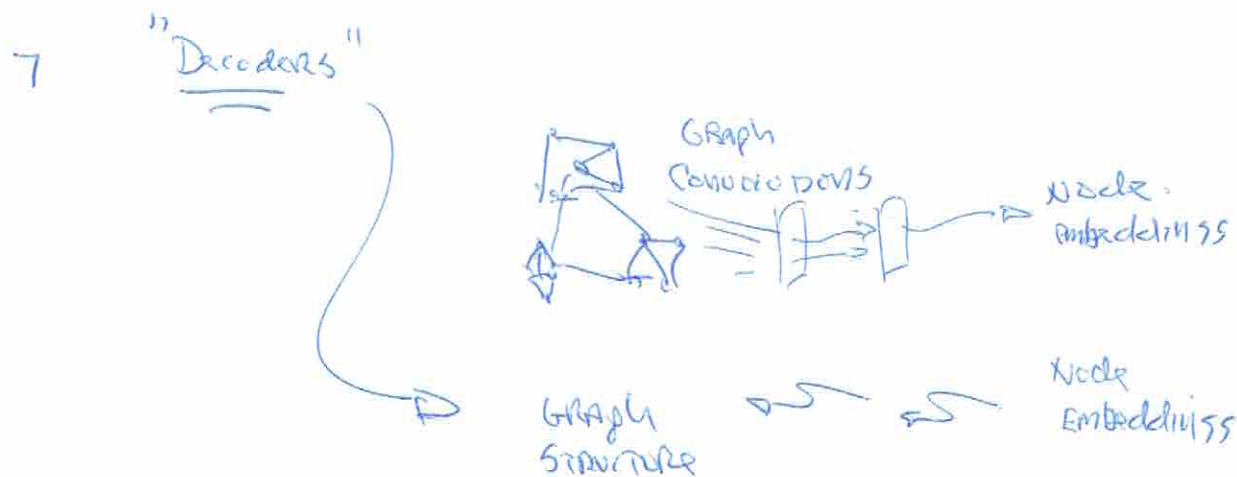
→ We want to generate realistic graphs using graph generative models.

Synthetic graph → insights  
 → Predictions  
 → Simulations  
 → Anomaly

History Properties of Real-world graphs  
 Traditional graph generative

Deep graph generative models

Learn graph formulation process from the data



9 Goal: generate graph similar to set of graphs

## 15 Graph RNN

(Graph as sequence.)

44 Scaling & Evaluating Graph Generation

58 Application: DEEP Graph Generative Model to Molecule Generation