

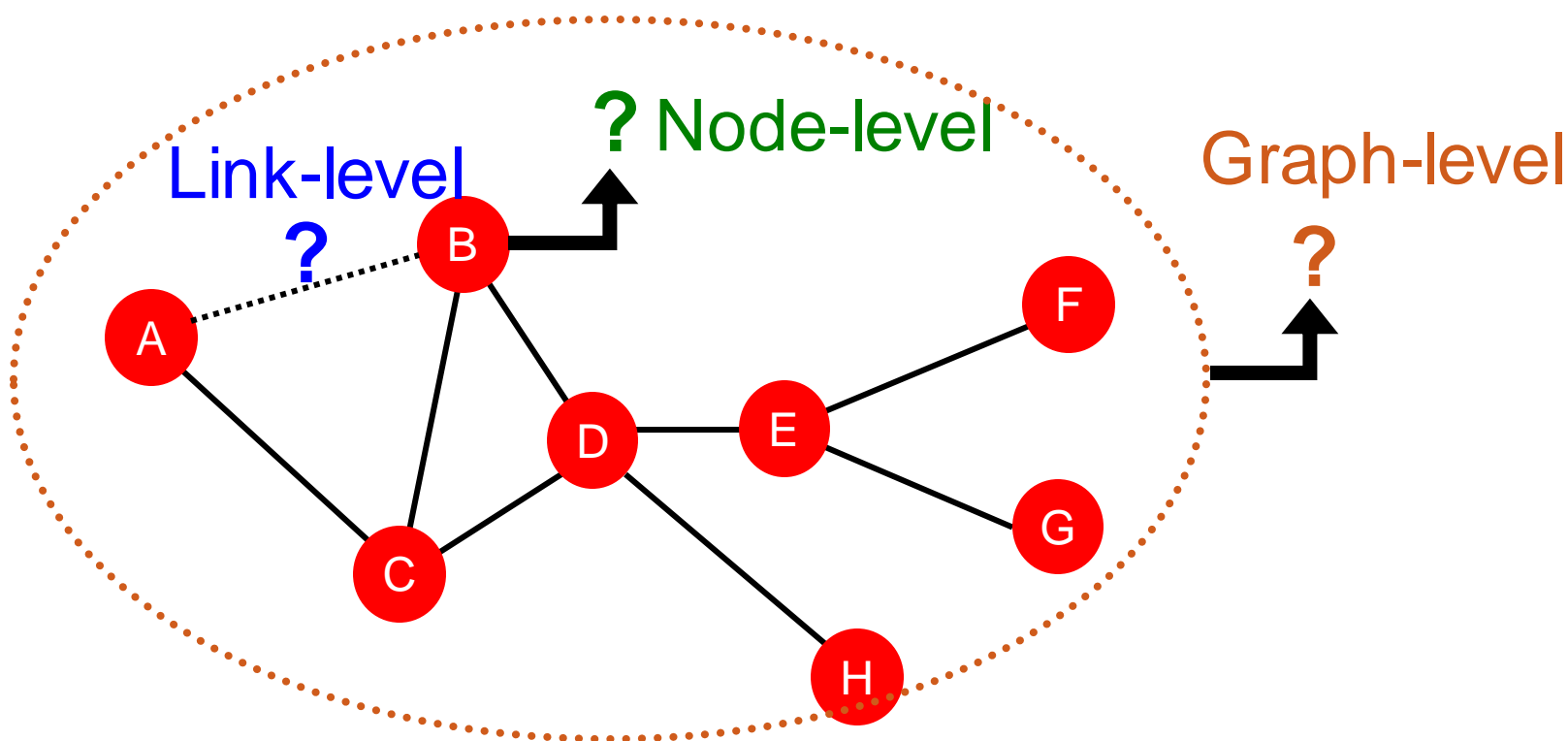
# Stanford CS224W: Traditional Methods for Machine Learning in Graphs

CS224W: Machine Learning with Graphs  
Jure Leskovec, Stanford University  
<http://cs224w.stanford.edu>



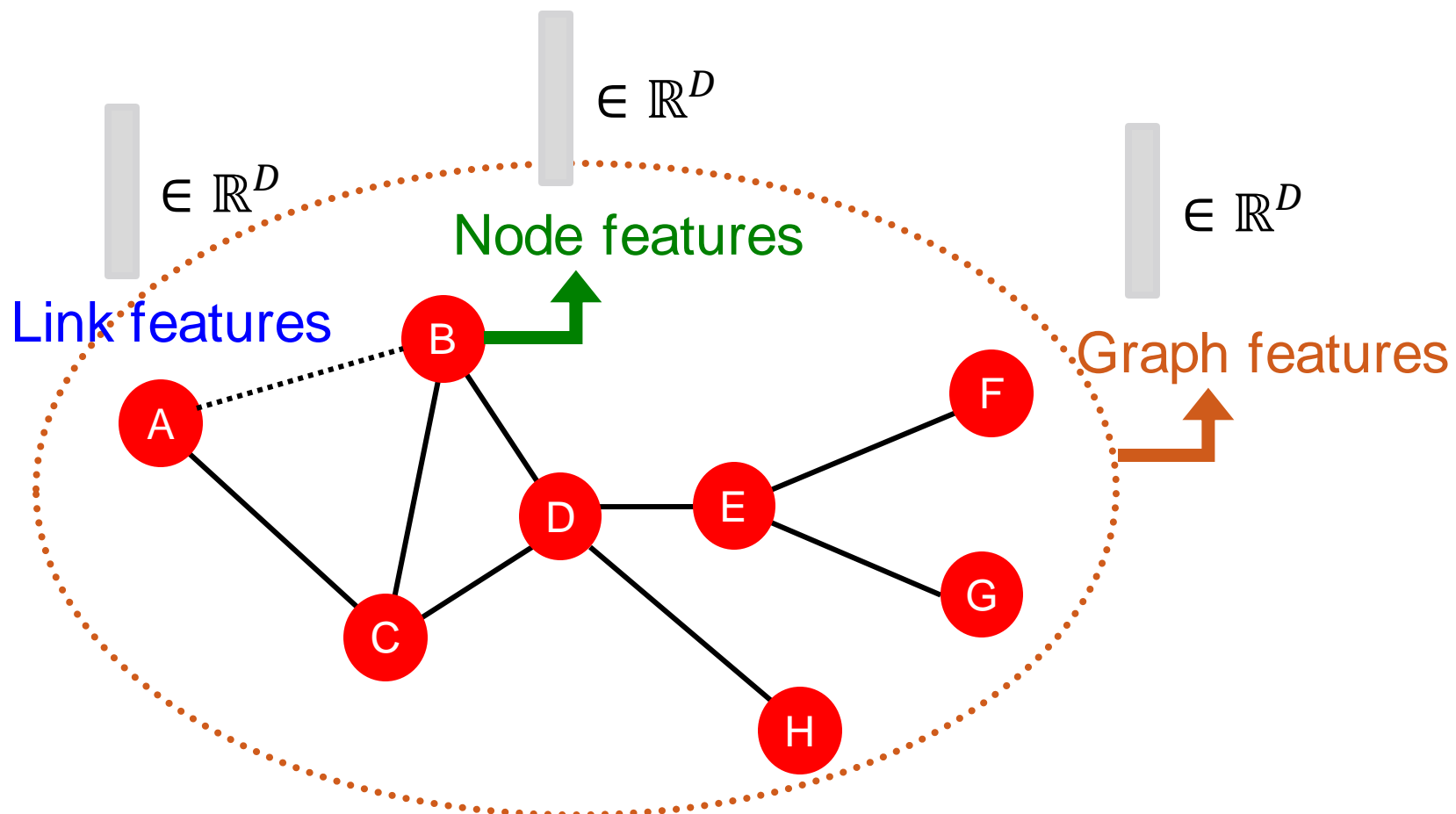
# Machine Learning Tasks: Review

- Node-level prediction
- Link-level prediction
- Graph-level prediction



# Traditional ML Pipeline

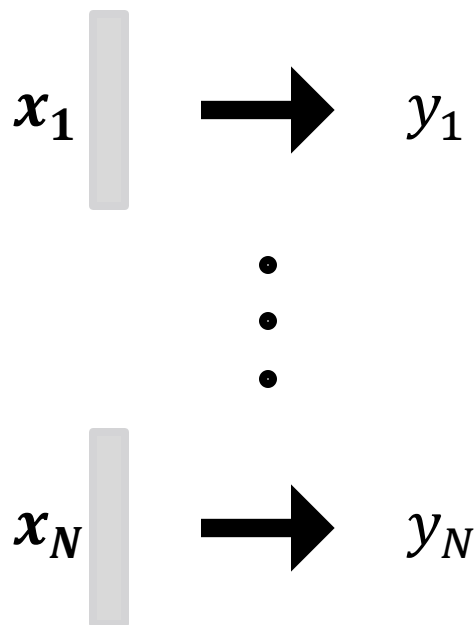
- Design features for nodes/links/graphs
- Obtain features for all training data



# Traditional ML Pipeline

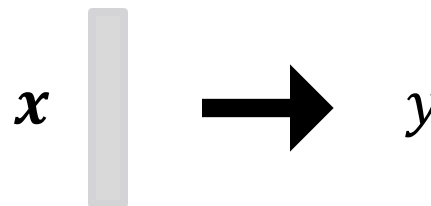
## ■ Train an ML model:

- Random forest
- SVM
- Neural network, etc.



## ■ Apply the model:

- Given a new node/link/graph, obtain its features and make a prediction



# This Lecture: Feature Design

- Using effective features over graphs is the key to achieving good model performance.
- Traditional ML pipeline uses hand-designed features.
- In this lecture, we overview the traditional features for:
  - Node-level prediction
  - Link-level prediction
  - Graph-level prediction
- For simplicity, we focus on undirected graphs.

# Machine Learning in Graphs

**Goal:** Make predictions for a set of objects

**Design choices:**

- **Features:**  $d$ -dimensional vectors
- **Objects:** Nodes, edges, sets of nodes, entire graphs
- **Objective function:**
  - What task are we aiming to solve?

# Machine Learning in Graphs

## Example: Node-level prediction

- Given:  $G = (V, E)$
- Learn a function:  $f : V \rightarrow \mathbb{R}$

How do we learn the function?

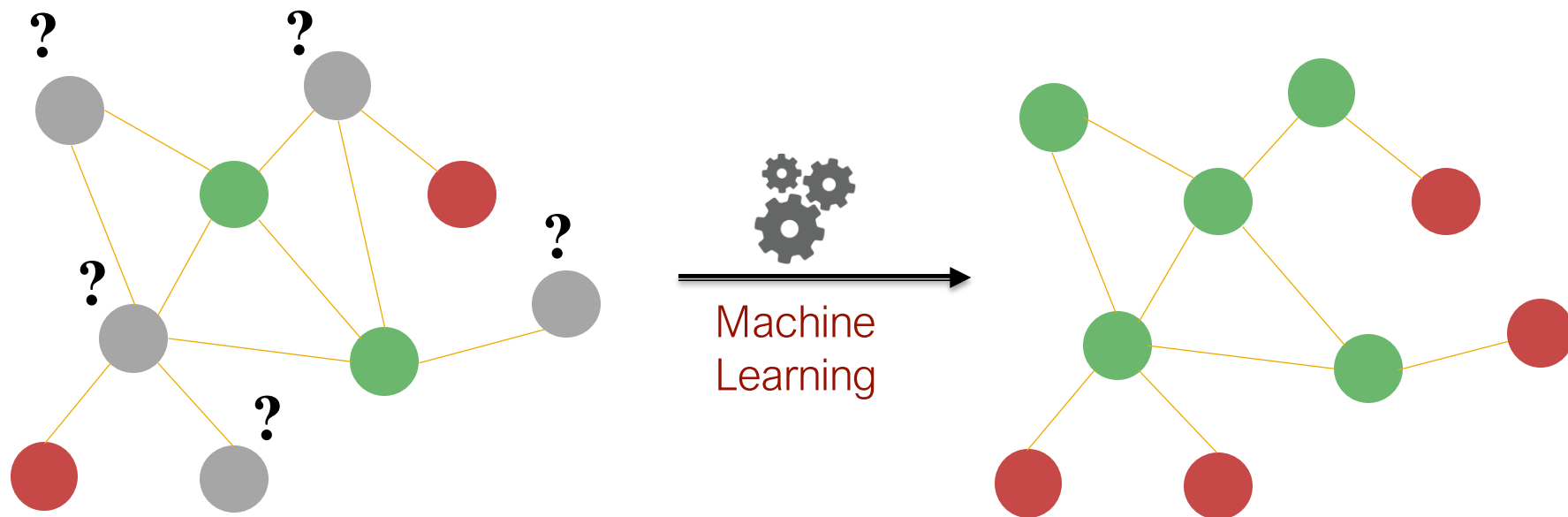
# Stanford CS224W: Node-Level Tasks and Features

CS224W: Machine Learning with Graphs  
Jure Leskovec, Stanford University  
<http://cs224w.stanford.edu>





# Node-Level Tasks



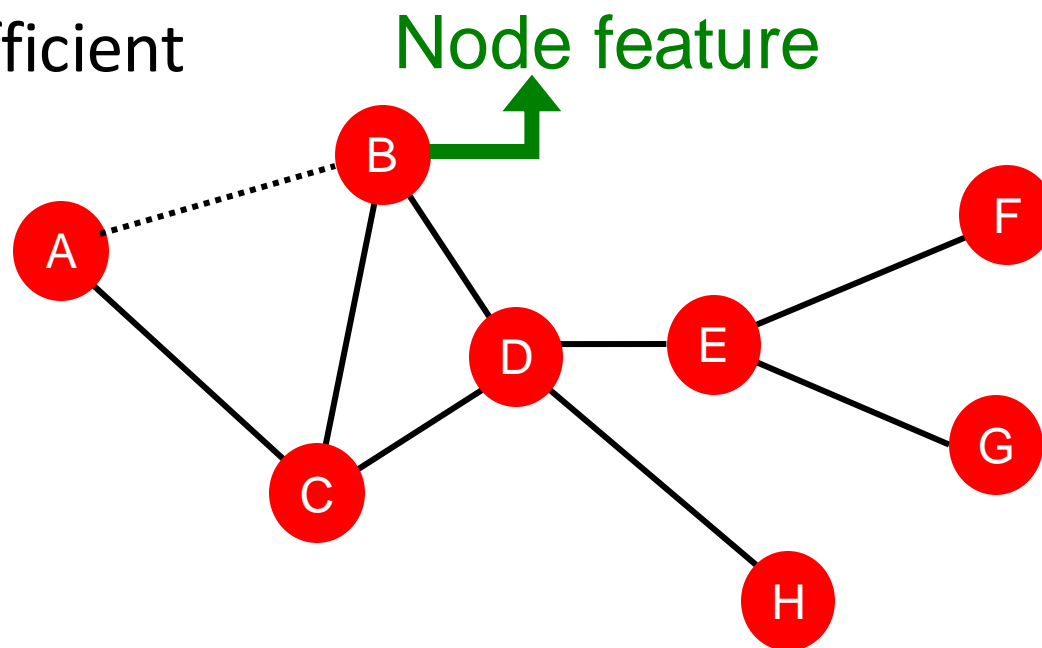
Node classification

ML needs features.

# Node-Level Features: Overview

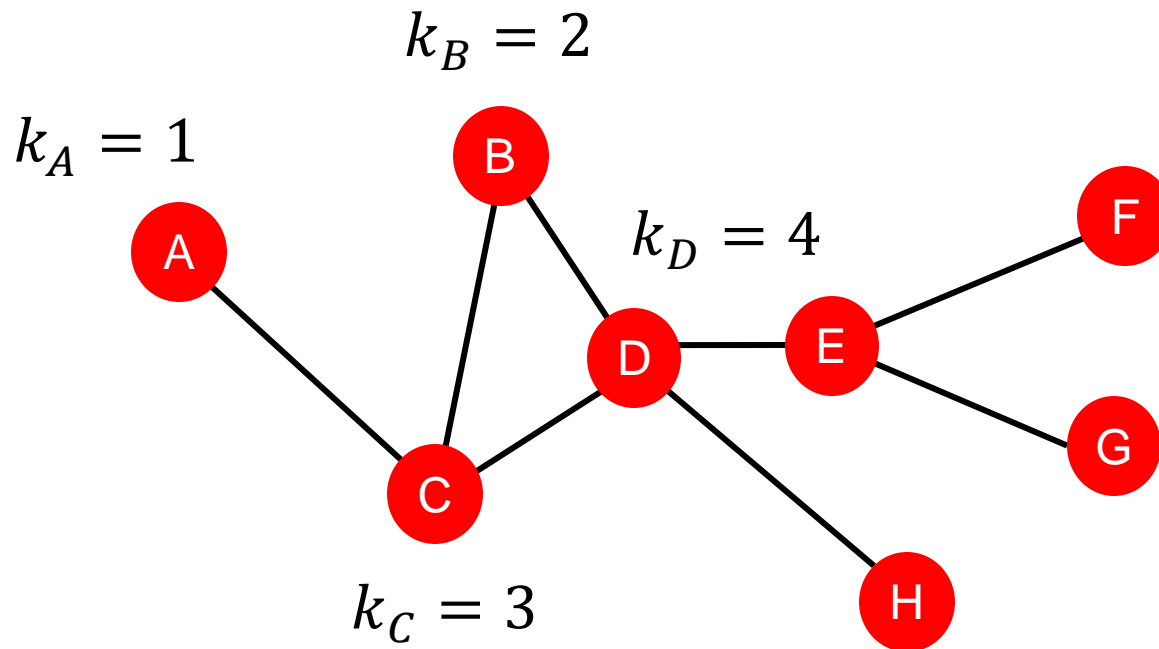
**Goal:** Characterize the structure and position of a node in the network:

- Node degree
- Node centrality
- Clustering coefficient
- Graphlets



# Node Features: Node Degree

- The degree  $k_v$  of node  $v$  is the number of edges (neighboring nodes) the node has.
- Treats all neighboring nodes equally.



# Node Features: Node Centrality

- Node degree counts the neighboring nodes without capturing their importance.
- Node centrality  $c_v$  takes the node importance in a graph into account
- **Different ways to model importance:**
  - Eigenvector centrality
  - Betweenness centrality
  - Closeness centrality
  - and many others...

# Node-Level Feature: Summary

- We have introduced different ways to obtain node features.
- They can be categorized as:
  - Importance-based features:
    - Node degree
    - Different node centrality measures
  - Structure-based features:
    - Node degree
    - Clustering coefficient
    - Graphlet count vector

# Node-Level Feature: Summary

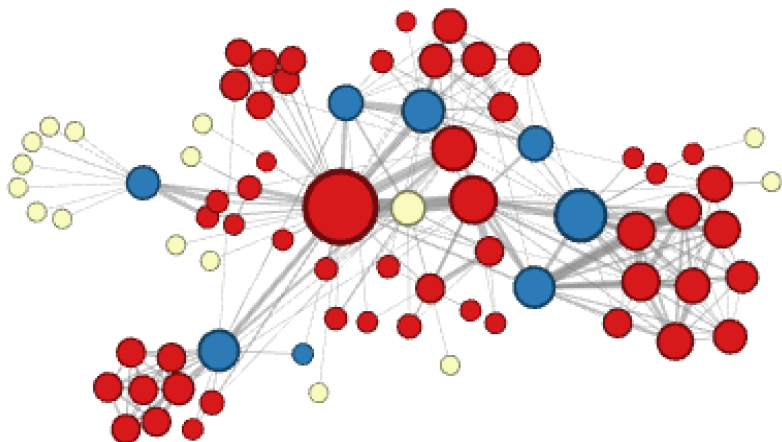
- **Importance-based features:** capture the importance of a node in a graph
  - Node degree:
    - Simply counts the number of neighboring nodes
  - Node centrality:
    - Models **importance of neighboring nodes** in a graph
    - Different modeling choices: eigenvector centrality, betweenness centrality, closeness centrality
- Useful for predicting influential nodes in a graph
  - **Example:** predicting celebrity users in a social network

# Node-Level Feature: Summary

- **Structure-based features:** Capture topological properties of local neighborhood around a node.
  - **Node degree:**
    - Counts the number of neighboring nodes
  - **Clustering coefficient:**
    - Measures how connected neighboring nodes are
  - **Graphlet degree vector:**
    - Counts the occurrences of different graphlets
- **Useful for predicting a particular role a node plays in a graph:**
  - **Example:** Predicting protein functionality in a protein-protein interaction network.

# Discussion

## Different ways to label nodes of the network:



Node features defined so far would allow to distinguish nodes in the above example



However, the features defines so far would not allow for distinguishing the above node labelling



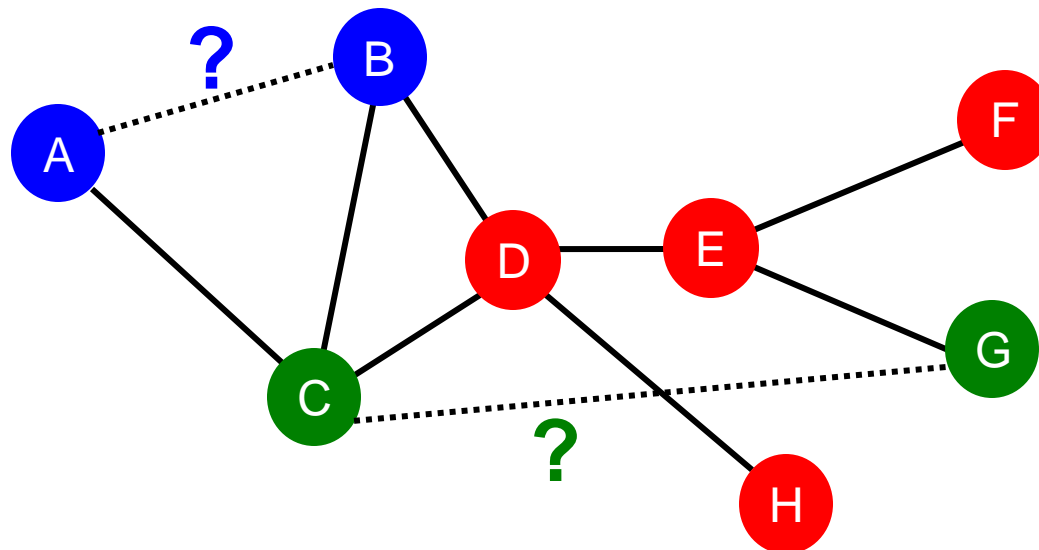
# Stanford CS224W: Link Prediction Task and Features

CS224W: Machine Learning with Graphs  
Jure Leskovec, Stanford University  
<http://cs224w.stanford.edu>



# Link-Level Prediction Task: Recap

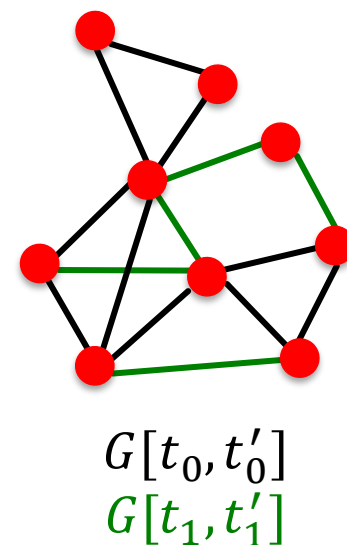
- The task is to predict **new links** based on the existing links.
- At test time, node pairs (with no existing links) are ranked, and top  $K$  node pairs are predicted.
- The key is to design features for a **pair of nodes**.



# Link Prediction as a Task

## Two formulations of the link prediction task:

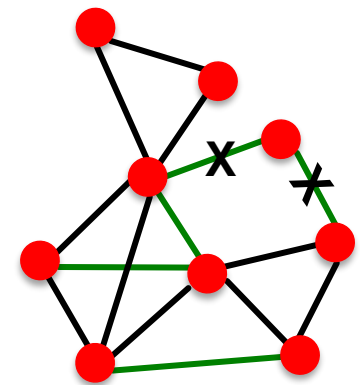
- **1) Links missing at random:**
  - Remove a random set of links and then aim to predict them
- **2) Links over time:**
  - Given  $G[t_0, t'_0]$  a graph defined by edges up to time  $t'_0$ , **output a ranked list  $L$**  of edges (not in  $G[t_0, t'_0]$ ) that are predicted to appear in time  $G[t_1, t'_1]$
  - **Evaluation:**
    - $n = |E_{new}|$ : # new edges that appear during the test period  $[t_1, t'_1]$
    - Take top  $n$  elements of  $L$  and count correct edges



# Link Prediction via Proximity

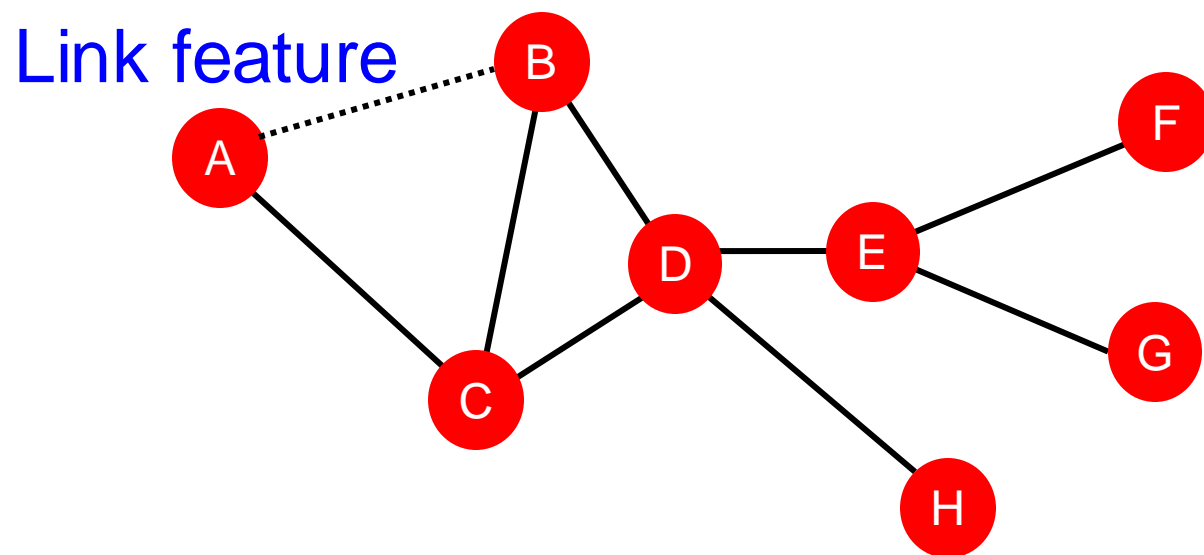
## ■ Methodology:

- For each pair of nodes  $(x,y)$  compute score  $c(x,y)$ 
  - For example,  $c(x,y)$  could be the # of common neighbors of  $x$  and  $y$
- Sort pairs  $(x,y)$  by the decreasing score  $c(x,y)$
- **Predict top  $n$  pairs as new links**
- **See which of these links actually appear in  $G[t_1, t'_1]$**



# Link-Level Features: Overview

- Distance-based feature
- Local neighborhood overlap
- Global neighborhood overlap



# Link-Level Features: Summary

- **Distance-based features:**

- Uses the shortest path length between two nodes but does not capture how neighborhood overlaps.

- **Local neighborhood overlap:**

- Captures how many neighboring nodes are shared by two nodes.
- Becomes zero when no neighbor nodes are shared.

- **Global neighborhood overlap:**

- Uses global graph structure to score two nodes.
- Katz index counts #walks of all lengths between two nodes.

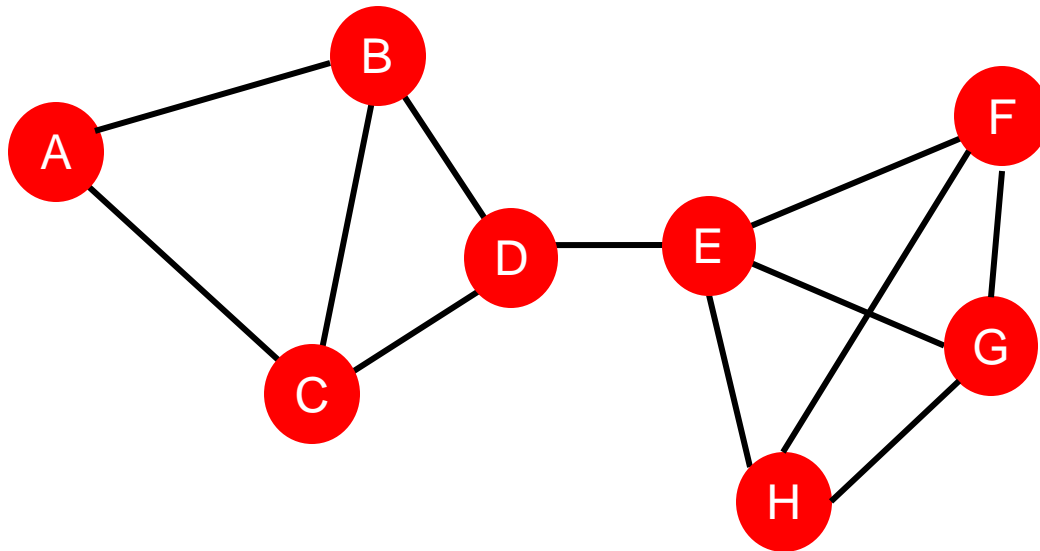
# Stanford CS224W: Graph-Level Features and Graph Kernels

CS224W: Machine Learning with Graphs  
Jure Leskovec, Stanford University  
<http://cs224w.stanford.edu>



# Graph-Level Features

- **Goal:** We want features that characterize the structure of an entire graph.
- **For example:**





# Background: Kernel Methods

- **Kernel methods** are widely-used for traditional ML for graph-level prediction.
- **Idea: Design kernels instead of feature vectors.**
- **A quick introduction to Kernels:**
  - Kernel  $K(G, G') \in \mathbb{R}$  measures similarity b/w data
  - Kernel matrix  $\mathbf{K} = (K(G, G'))_{G, G'}$  must always be positive semidefinite (i.e., has positive eigenvalues)
  - There exists a feature representation  $\phi(\cdot)$  such that  $K(G, G') = \phi(G)^T \phi(G')$
  - Once the kernel is defined, off-the-shelf ML model, such as **kernel SVM**, can be used to make predictions.

# Graph-Level Features: Overview

- **Graph Kernels:** Measure similarity between two graphs:
  - Graphlet Kernel [1]
  - Weisfeiler-Lehman Kernel [2]
  - Other kernels are also proposed in the literature (beyond the scope of this lecture)
    - Random-walk kernel
    - Shortest-path graph kernel
    - And many more...

[1] Shervashidze, Nino, et al. "Efficient graphlet kernels for large graph comparison." Artificial Intelligence and Statistics. 2009.

[2] Shervashidze, Nino, et al. "Weisfeiler-lehman graph kernels." Journal of Machine Learning Research 12.9 (2011).

# Graph-Level Features: Summary

## ■ Graphlet Kernel

- Graph is represented as **Bag-of-graphlets**
- **Computationally expensive**

## ■ Weisfeiler-Lehman Kernel

- Apply  $K$ -step color refinement algorithm to enrich node colors
  - Different colors capture different  $K$ -hop neighborhood structures
- Graph is represented as **Bag-of-colors**
- **Computationally efficient**
- Closely related to Graph Neural Networks (as we will see!)

# Today's Summary

- **Traditional ML Pipeline**

- Hand-crafted feature + ML model

- **Hand-crafted features for graph data**

- **Node-level:**

- Node degree, centrality, clustering coefficient, graphlets

- **Link-level:**

- Distance-based feature
    - local/global neighborhood overlap

- **Graph-level:**

- Graphlet kernel, WL kernel