

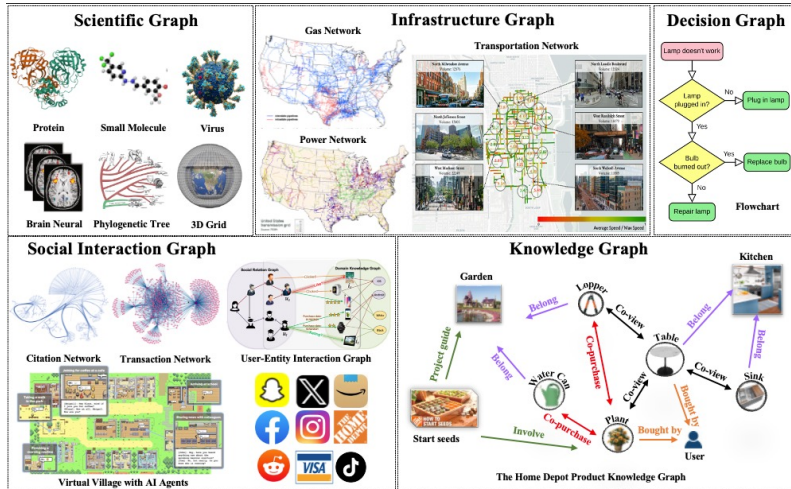


# Mining and Learning on Graphs

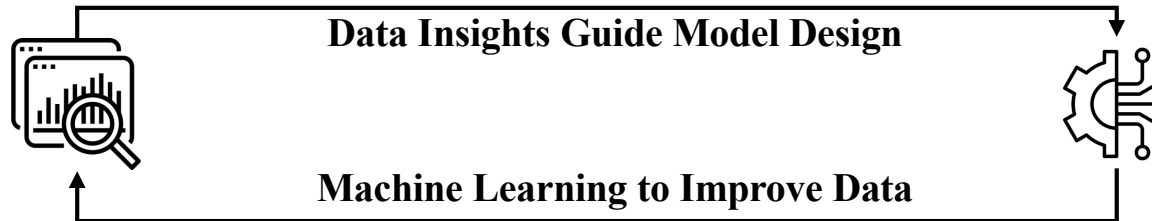
## Linear Algebra and Graph Theory

Yu Wang, Ph.D.  
Assistant Professor  
Computer and Information Science  
University of Oregon  
CS 410/510 - Fall 2024

# Data Mining & Machine Learning on Graphs



## Random Forests



### Data mining

Analyze data

Derive patterns and relationships

Solve real-world problems

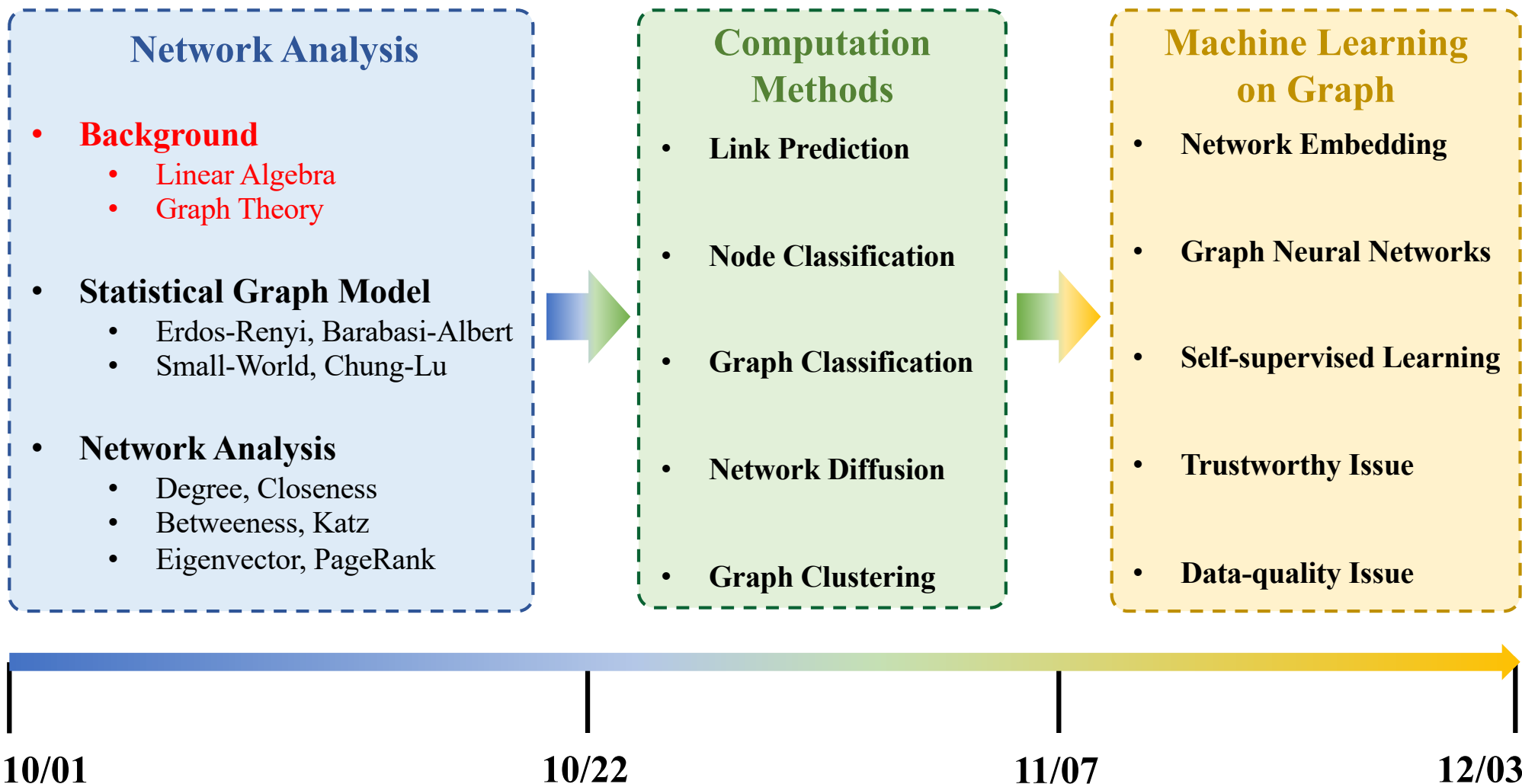
### Machine Learning

Design Model

Allow Computer to Learn and Improve

Without being explicit programmed

# Data Mining & Machine Learning on Graphs





## Vector

$$\mathbf{v} = [1 \quad 2 \quad 5]$$

$$\mathbf{u} = \begin{bmatrix} 1 \\ 2 \\ 5 \end{bmatrix}$$

**Please note that we will use  
this one by default**

## Matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \left. \vphantom{\begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix}} \right\} \begin{array}{l} 4 \text{ rows} \\ \\ \\ 3 \text{ columns} \end{array}$$

$$\mathbf{v} \in \mathbb{R}^{1 \times 3}$$

$$\mathbf{u} \in \mathbb{R}^{3 \times 1}$$

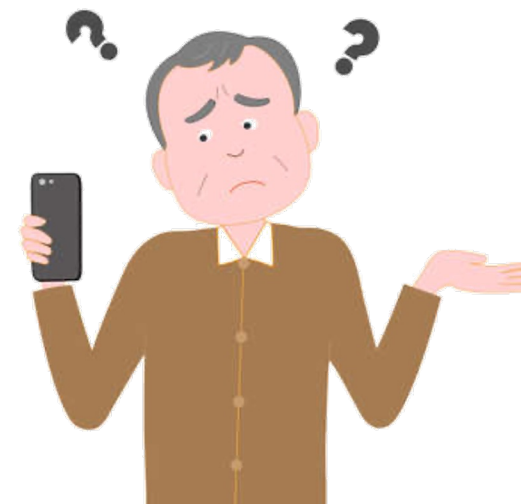
$$\mathbf{A} \in \mathbb{R}^{4 \times 3}$$

# Linear Algebra – Basic Operations



## 1 Basics

$$\begin{aligned}(AB)^{-1} &= B^{-1}A^{-1} & (1) \\(ABC\dots)^{-1} &= \dots C^{-1}B^{-1}A^{-1} & (2) \\(A^T)^{-1} &= (A^{-1})^T & (3) \\(A+B)^T &= A^T+B^T & (4) \\(AB)^T &= B^T A^T & (5) \\(ABC\dots)^T &= \dots C^T B^T A^T & (6) \\(A^H)^{-1} &= (A^{-1})^H & (7) \\(A+B)^H &= A^H+B^H & (8) \\(AB)^H &= B^H A^H & (9) \\(ABC\dots)^H &= \dots C^H B^H A^H & (10)\end{aligned}$$



## Matrix Codebook

<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

## The Matrix Cookbook

[ <http://matrixcookbook.com> ]

Kaare Brandt Petersen  
Michael Syskind Pedersen

VERSION: NOVEMBER 15, 2012

# Linear Algebra – Matrix Multiplication



## Matrix Multiplication

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \times \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} \longrightarrow \mathbf{C} = [ \quad ] \quad ?$$

$4 \times 3$                        $3 \times 2$

**Dimensions much match!**

**What is the dimension of C?  $(4 \times 3)(3 \times 2) \rightarrow 4 \times 2$**

# Linear Algebra – Matrix Multiplication



## Matrix Multiplication

$$\begin{matrix} \mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} & \times & \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} & \longrightarrow & \mathbf{C} = \begin{bmatrix} 20 \\ \phantom{0} \\ \phantom{0} \\ \phantom{0} \end{bmatrix} \\ 4 \times 3 & & 3 \times 2 & & \end{matrix}$$

$1 \times 1 + 2 \times 2 + 3 \times 5$

$$\begin{matrix} \mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} & \times & \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} & \longrightarrow & \mathbf{C} = \begin{bmatrix} 20 & 29 \\ \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} \end{bmatrix} \\ 4 \times 3 & & 3 \times 2 & & \end{matrix}$$

$1 \times 2 + 2 \times 3 + 3 \times 7$

# Linear Algebra – Matrix Multiplication



## Matrix Multiplication

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \times \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} \longrightarrow \mathbf{C} = \begin{bmatrix} 20 & 29 \\ 15 & \\ & \end{bmatrix}$$

$4 \times 3 \qquad \qquad \qquad 3 \times 2$

$0 \times 1 + 5 \times 2 + 1 \times 5$

*(Note: In the original image, the second row of A and the first column of B are highlighted in red. An arrow points from the calculation above to the element 15 in the second row, first column of C.)*

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \times \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} \longrightarrow \mathbf{C} = \begin{bmatrix} 20 & 29 \\ 15 & 22 \\ & \end{bmatrix}$$

$4 \times 3 \qquad \qquad \qquad 3 \times 2$

$0 \times 2 + 5 \times 3 + 1 \times 7$

*(Note: In the original image, the second row of A and the second column of B are highlighted in red. An arrow points from the calculation above to the element 22 in the second row, second column of C.)*



# Linear Algebra – Matrix Multiplication



## Matrix Multiplication

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \times \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} \longrightarrow \mathbf{C} = \begin{bmatrix} 20 & 29 \\ 15 & 22 \\ 43 & \end{bmatrix}$$

$4 \times 3 \qquad 3 \times 2$

$2 \times 1 + 3 \times 2 + 7 \times 5$

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \times \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} \longrightarrow \mathbf{C} = \begin{bmatrix} 20 & 29 \\ 15 & 22 \\ 43 & 62 \end{bmatrix}$$

$4 \times 3 \qquad 3 \times 2$

$2 \times 2 + 3 \times 3 + 7 \times 7$

# Linear Algebra – Matrix Multiplication



## Matrix Multiplication

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \times \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} \longrightarrow \mathbf{C} = \begin{bmatrix} 20 & 29 \\ 15 & 22 \\ 43 & 62 \\ 61 & 62 \end{bmatrix}$$

$4 \times 3 \qquad 3 \times 2$

$3 \times 1 + 2 \times 9 + 5 \times 8$

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix} \times \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 5 & 7 \end{bmatrix} \longrightarrow \mathbf{C} = \begin{bmatrix} 20 & 29 \\ 15 & 22 \\ 43 & 62 \\ 61 & 89 \end{bmatrix}$$

$4 \times 3 \qquad 3 \times 2$

$3 \times 2 + 3 \times 9 + 7 \times 8$

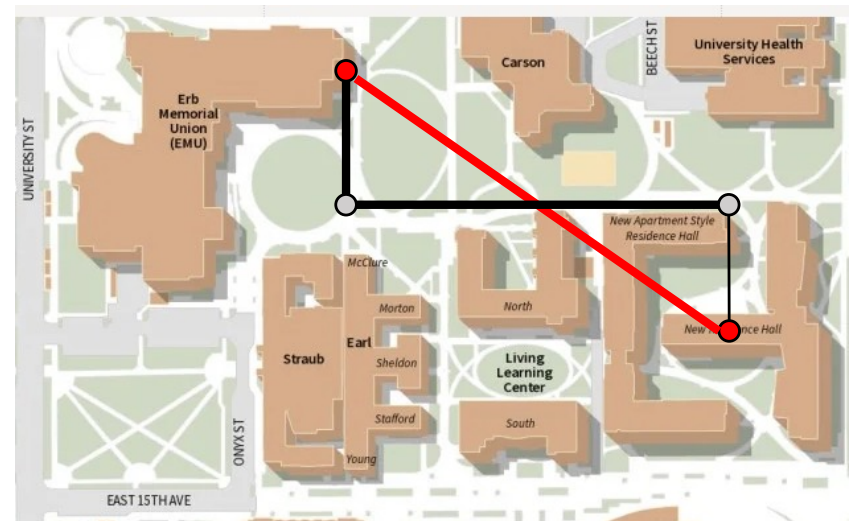


# Linear Algebra – Vector Norms

- $\|\mathbf{v}\|^p$
- Function from a vector space to a single positive real value:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$
- Length of  $\mathbf{v}$

$$\|\mathbf{v}\|^p = \left( \sum_{i=1}^d v_i^p \right)^{\frac{1}{p}}$$

- Examples:
  - Manhattan distance ( $L_1$ ):  $\|\mathbf{v}\|^1 = \left( \sum_{i=1}^d |v_i| \right)$
  - Euclidean distance ( $L_2$ ):  $\|\mathbf{v}\|^2 = \left( \sum_{i=1}^d v_i^2 \right)^{\frac{1}{2}}$
  - **How about  $L_0$ ?**



# Linear Algebra – Transpose Matrix



- $\mathbf{A}^T$  or  $\mathbf{A}'$
- Flip rows and columns

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 5 & 1 \\ 2 & 3 & 7 \\ 3 & 9 & 8 \end{bmatrix}$$

$4 \times 3$

$$\mathbf{A}^T = \begin{bmatrix} 1 & 0 & 2 & 3 \\ 2 & 5 & 3 & 9 \\ 3 & 1 & 7 & 8 \end{bmatrix}$$

$3 \times 4$

$$\mathbf{u} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$3 \times 1$

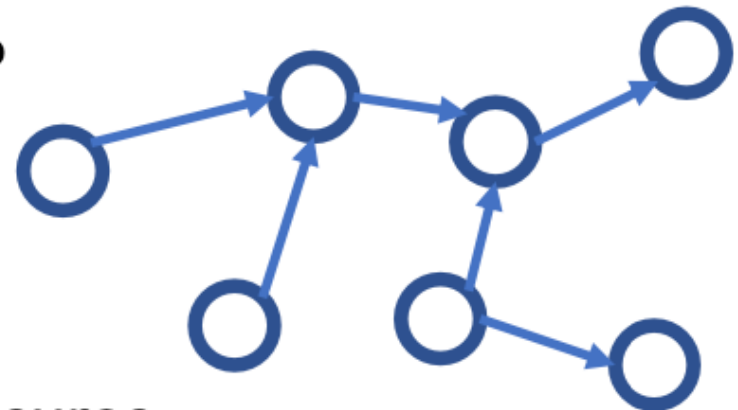
$$\mathbf{u}^T = [1 \quad 2 \quad 3]$$

$1 \times 3$



## ■ Motivation

- Given a graph (e.g., web pages that have a keyword that was searched for)
- Which node is the most important?



- Later we will discuss centrality measures
  - E.g., HITS and PageRank
- That are based on SVD and Eigendecomposition



- $\mathbf{A}_{[n \times m]} = \mathbf{U}_{[n \times r]} \mathbf{\Lambda}_{[r \times r]} (\mathbf{V}_{[m \times r]})^T$
- **A**:  $n \times m$  matrix (e.g.,  $n$  documents,  $m$  terms)
- **U**:  $n \times r$  matrix ( $n$  documents,  $r$  concepts)
- **$\Lambda$** :  $r \times r$  diagonal matrix ( $r$  is rank of the matrix)
  - Can be seen as strength of the topic
- **V**:  $m \times r$  matrix ( $m$  terms,  $r$  concepts)

# Linear Algebra – Singular Value Decomposition (SVD)



$$\begin{array}{c}
 \uparrow \\
 \text{CS} \\
 \downarrow \\
 \uparrow \\
 \text{MD} \\
 \downarrow
 \end{array}
 \begin{array}{c}
 \text{retrieval} \\
 \text{data} \quad \text{inf.} \downarrow \quad \text{brain} \quad \text{lung} \\
 \begin{bmatrix}
 1 & 1 & 1 & 0 & 0 \\
 2 & 2 & 2 & 0 & 0 \\
 1 & 1 & 1 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 0 & 0 & 2 & 2 \\
 0 & 0 & 0 & 3 & 3 \\
 0 & 0 & 0 & 1 & 1
 \end{bmatrix}
 =
 \begin{bmatrix}
 0.18 & 0 \\
 0.36 & 0 \\
 0.18 & 0 \\
 0.90 & 0 \\
 0 & 0.53 \\
 0 & 0.80 \\
 0 & 0.27
 \end{bmatrix}
 \times
 \begin{bmatrix}
 9.64 & 0 \\
 0 & 5.29
 \end{bmatrix}
 \times
 \begin{bmatrix}
 0.58 & 0.58 & 0.58 & 0 & 0 \\
 0 & 0 & 0 & 0.71 & 0.71
 \end{bmatrix}
 \end{array}$$

# Linear Algebra – Singular Value Decomposition (SVD)



$$\begin{array}{c}
 \uparrow \\
 \text{CS} \\
 \downarrow \\
 \uparrow \\
 \text{MD} \\
 \downarrow
 \end{array}
 \begin{array}{c}
 \text{data} \\
 \text{inf.} \\
 \text{retrieval} \\
 \text{brain} \\
 \text{lung}
 \end{array}
 \begin{bmatrix}
 1 & 1 & 1 & 0 & 0 \\
 2 & 2 & 2 & 0 & 0 \\
 1 & 1 & 1 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 0 & 0 & 2 & 2 \\
 0 & 0 & 0 & 3 & 3 \\
 0 & 0 & 0 & 1 & 1
 \end{bmatrix}
 =
 \begin{array}{c}
 \text{CS-concept} \\
 \text{MD-concept}
 \end{array}
 \begin{bmatrix}
 0.18 & 0 \\
 0.36 & 0 \\
 0.18 & 0 \\
 0.90 & 0 \\
 0 & 0.53 \\
 0 & 0.80 \\
 0 & 0.27
 \end{bmatrix}
 \times
 \begin{bmatrix}
 9.64 & 0 \\
 0 & 5.29
 \end{bmatrix}
 \times
 \begin{bmatrix}
 0.58 & 0.58 & 0.58 & 0 & 0 \\
 0 & 0 & 0 & 0.71 & 0.71
 \end{bmatrix}$$



# Linear Algebra – Singular Value Decomposition (SVD)



term-to-concept  
similarity matrix

retrieval  
inf. ↓ brain lung

↑ CS  
↓  
↑ MD

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

CS-concept

# Linear Algebra – Singular Value Decomposition (SVD)



retrieval  
inf. ↓ brain lung

‘strength’ of CS-concept

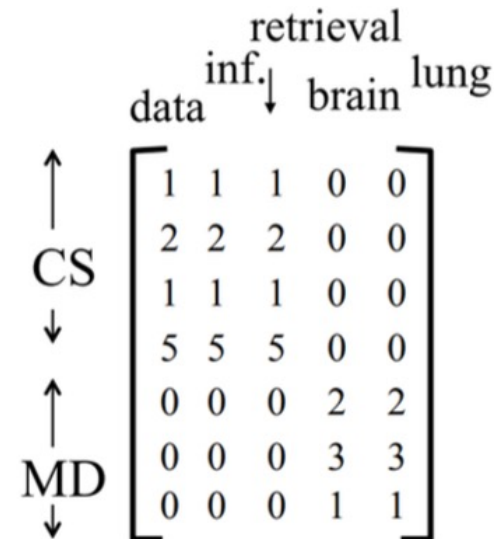
↓

$$\begin{array}{c} \uparrow \\ \text{CS} \\ \downarrow \\ \uparrow \\ \text{MD} \\ \downarrow \end{array}
 \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}
 =
 \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}
 \times
 \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}
 \times
 \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$



# Linear Algebra – Singular Value Decomposition (SVD)

```
[>>> import numpy as np
[>>> a = np.array([[0.18, 0], [0.36, 0], [0.18, 0], [0.90, 0], [0, 0.53], [0, 0.80], [0, 0.27]])
[>>> b = np.array([[0.58, 0.58, 0.58, 0, 0], [0, 0, 0, 0.71, 0.71]])
[>>> a
array([[0.18, 0. ],
       [0.36, 0. ],
       [0.18, 0. ],
       [0.9 , 0. ],
       [0. , 0.53],
       [0. , 0.8 ],
       [0. , 0.27]])
[>>> b
array([[0.58, 0.58, 0.58, 0. , 0. ],
       [0. , 0. , 0. , 0.71, 0.71]])
[>>> np.matmul(a, b)
array([[0.1044, 0.1044, 0.1044, 0. , 0. ],
       [0.2088, 0.2088, 0.2088, 0. , 0. ],
       [0.1044, 0.1044, 0.1044, 0. , 0. ],
       [0.522 , 0.522 , 0.522 , 0. , 0. ],
       [0. , 0. , 0. , 0.3763, 0.3763],
       [0. , 0. , 0. , 0.568 , 0.568 ],
       [0. , 0. , 0. , 0.1917, 0.1917]])
```





- **'documents', 'terms', 'concepts'**
- **U**: document-to-concept similarity matrix
- **V**: term-to-concept similarity matrix
- **$\Lambda$** : the diagonal elements are the 'strength' of each concept
  
- **A**:  $n \times m$  matrix (e.g.,  $n$  documents,  $m$  terms)
- What is  **$A^T A$** ?
  - Term-to-term ( $m \times m$ ) similarity matrix
- What is  **$AA^T$** ?
  - ...



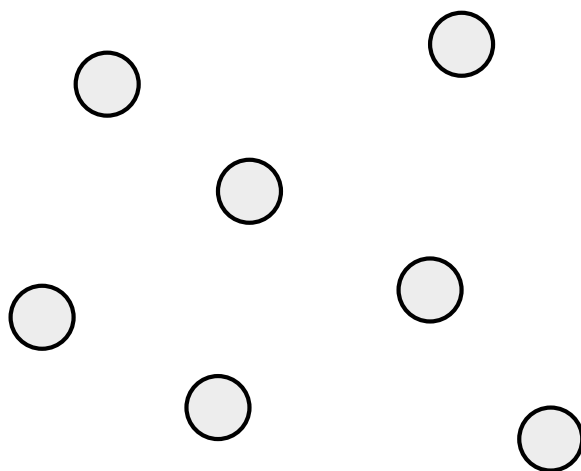
- Text domain:
  - Latent Semantic Indexing (LSI)
    - Analyze the relationship between a set of documents and their terms contained by generating a set of concepts
- Dimensionality reduction
  - In general, for a data matrix  $\mathbf{X}$
  - Or can be used for an adjacency matrix  $\mathbf{A}$

# Graph and Network Theory- Why?

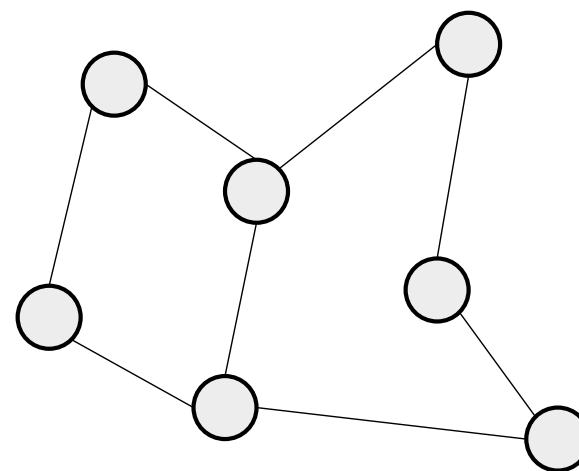


**They provide a general language for describing highly complex systems in a unified way**

**Traditional Data View**



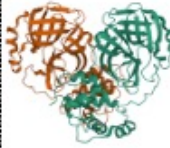
**Network Data View**




# Data Mining & Machine Learning on Graphs



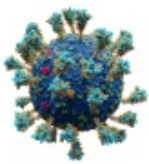
### Scientific Graph




Protein



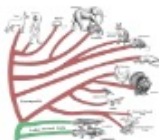
Small Molecule




Virus



Brain Neural




Phylogenetic Tree




3D Grid

### Infrastructure Graph

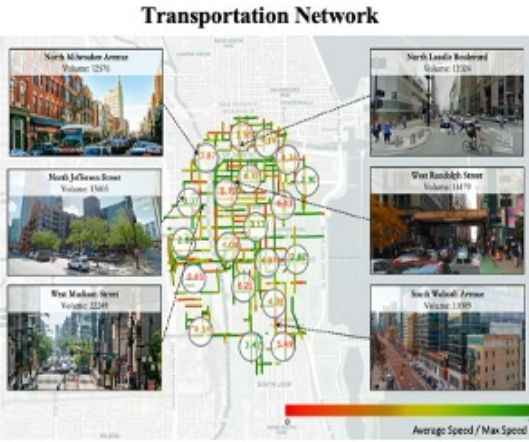
#### Gas Network



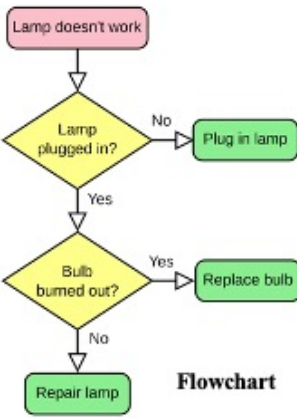
#### Power Network



#### Transportation Network




### Decision Graph




**Flowchart**


### Social Interaction Graph




Citation Network



Transaction Network

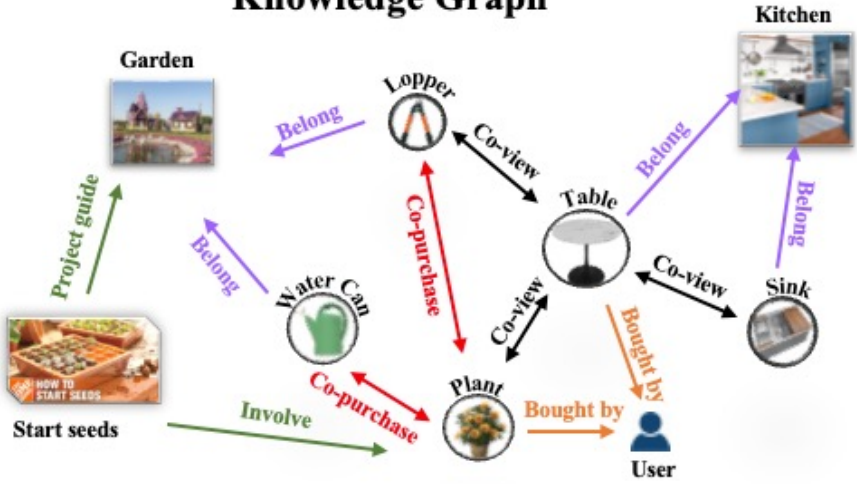


User-Entity Interaction Graph



Virtual Village with AI Agents

### Knowledge Graph



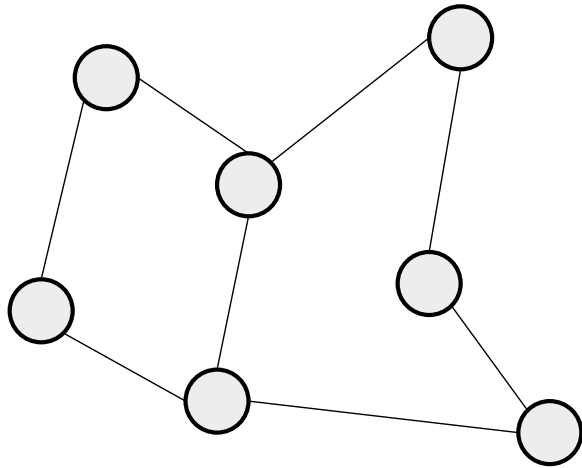
**The Home Depot Product Knowledge Graph**



# Graph and Network Theory - Basic



## Network/Graph-Structured Data



Vertices/Nodes:  $\mathcal{V}$  with feature  $\mathbf{X} \in \mathbb{R}^{|\mathcal{V}| \times d}$

Edges/Links:  $\mathcal{E}$  with feature  $\mathbf{E} \in \mathbb{R}^{|\mathcal{E}| \times d'}$

Overall system:  $G = (\mathcal{V}, \mathcal{E}, \mathbf{X}, \mathbf{E})$

**Network**

Real-world system



**Graph**

Model representation of a network in mathematics



# Graph and Network Theory - Examples



- **Connect Conference/Journal Papers with each other where nodes are papers, and the links represent a citation from one to another.....**
  - Citation Network (e.g., DBLP)
- **If we connect people based on their dating relations where the nodes are people, and the links are their relations...**
  - Dating Network (e.g., Tinder)
- **If we connect all the words in the dictionary where the nodes are the words and the links connect words having semantic relations between them ...**
  - Word Network (e.g., WordNet)

**Can you name some examples?**



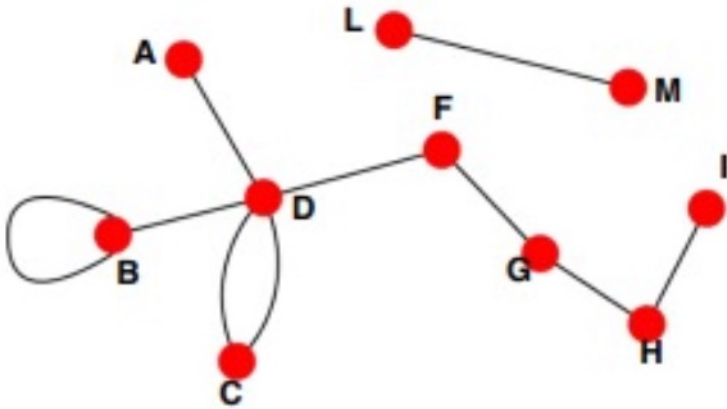
## Given a specific scenario, how should we construct the network?

- **How to decide what are the nodes and edges?**
  - Sometimes it is unambiguous and unique
    - User-user interaction, customer-item interaction
  - Other times it is left up to the application needs
    - Route-Net Example
- **In either way, constructing the network/graph is very importance for downstream tasks**
  - If we connect two users based on whether they have the same first name instead of based on whether they have following relations on a Twitter dataset...



## Undirected

- **Links:** undirected (symmetrical, reciprocal)

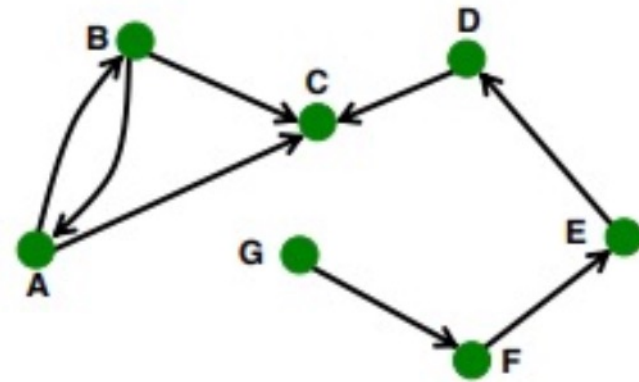


- **Examples:**

- Collaborations
- Friendship on Facebook

## Directed

- **Links:** directed (arcs)



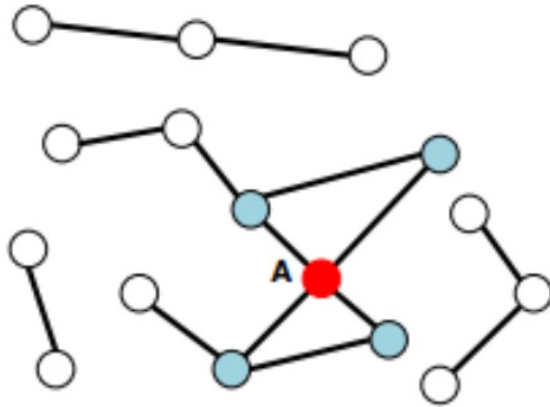
- **Examples:**

- Phone calls
- Following on Twitter

# Graph and Network Theory – Node Degree



Undirected

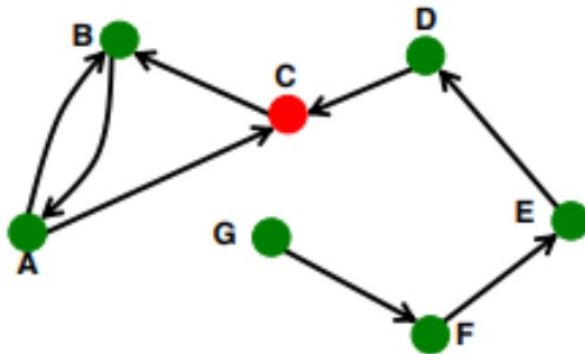


**Node degree,  $k_i$ :** the number of edges adjacent to node  $i$

$$k_A = 4$$

**Avg. degree:**  $\bar{k} = \langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i = \frac{2E}{N}$

Directed



In directed networks we define an **in-degree** and **out-degree**.

The (total) degree of a node is the sum of in- and out-degrees.

$$k_C^{in} = 2 \quad k_C^{out} = 1 \quad k_C = 3$$

$$\bar{k} = \frac{E}{N}$$

$$\overline{k^{in}} = \overline{k^{out}}$$

**Source:** Node with  $k^{in} = 0$

**Sink:** Node with  $k^{out} = 0$

# Graph and Network Theory – Examples

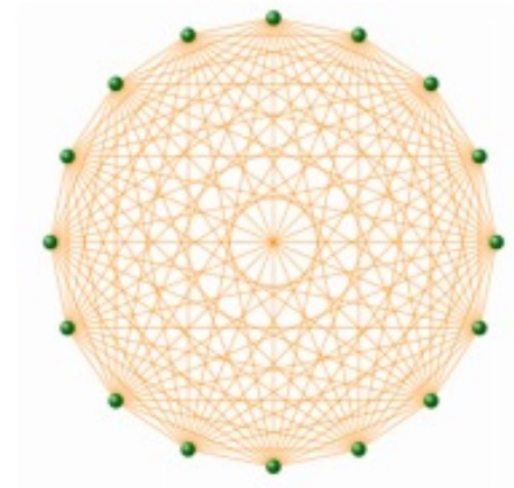


NETWORK	NODES	LINKS	DIRECTED UNDIRECTED	N	L	$\langle k \rangle$
Internet	Routers	Internet connections	Undirected	192,244	609,066	6.33
WWW	Webpages	Links	Directed	325,729	1,497,134	4.60
Power Grid	Power plants, transformers	Cables	Undirected	4,941	6,594	2.67
Mobile Phone Calls	Subscribers	Calls	Directed	36,595	91,826	2.51
Email	Email addresses	Emails	Directed	57,194	103,731	1.81
Science Collaboration	Scientists	Co-authorship	Undirected	23,133	93,439	8.08
Actor Network	Actors	Co-acting	Undirected	702,388	29,397,908	83.71
Citation Network	Paper	Citations	Directed	449,673	4,689,479	10.43
E. Coli Metabolism	Metabolites	Chemical reactions	Directed	1,039	5,802	5.58
Protein Interactions	Proteins	Binding interactions	Undirected	2,018	2,930	2.90



The **maximum number of edges** in an undirected graph on  $N$  nodes is

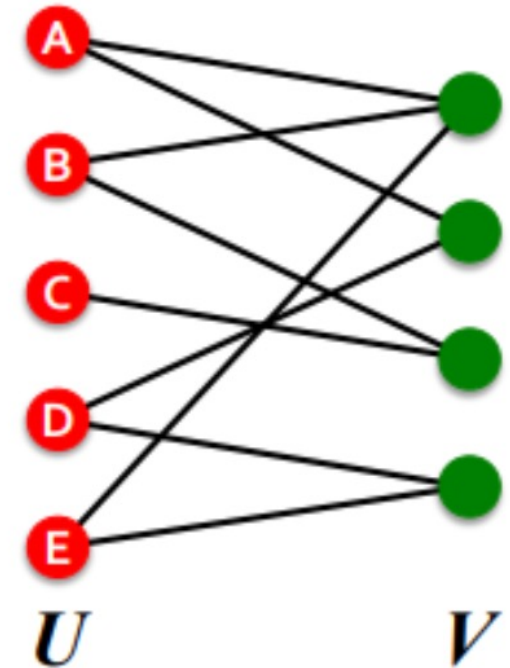
$$E_{\max} = \binom{N}{2} = \frac{N(N-1)}{2}$$



An undirected graph with the number of edges  $E = E_{\max}$  is called a **complete graph**, and its average degree is  $N-1$



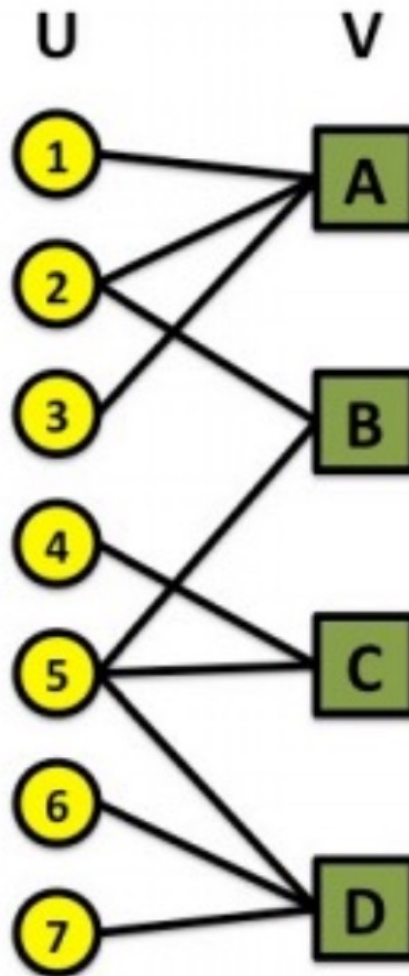
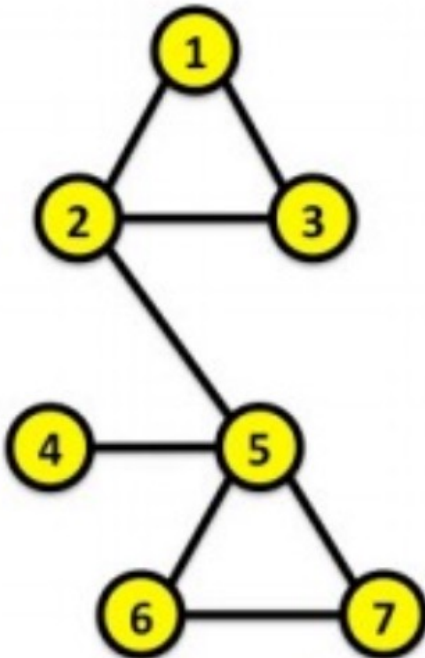
- **Bipartite graph** is a graph whose nodes can be divided into two disjoint sets  $U$  and  $V$  such that every link connects a node in  $U$  to one in  $V$ ; that is,  $U$  and  $V$  are **independent sets**
- **Examples:**
  - Authors-to-Papers (they authored)
  - Actors-to-Movies (they appeared in)
  - Users-to-Movies (they rated)
  - Recipes-to-Ingredients (they contain)
- **“Folded” networks:**
  - Author collaboration networks
  - Movie co-rating networks



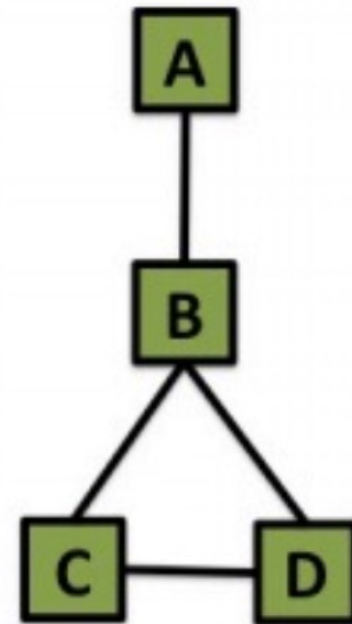
# Graph and Network Theory – Bipartite Projections



Projection U



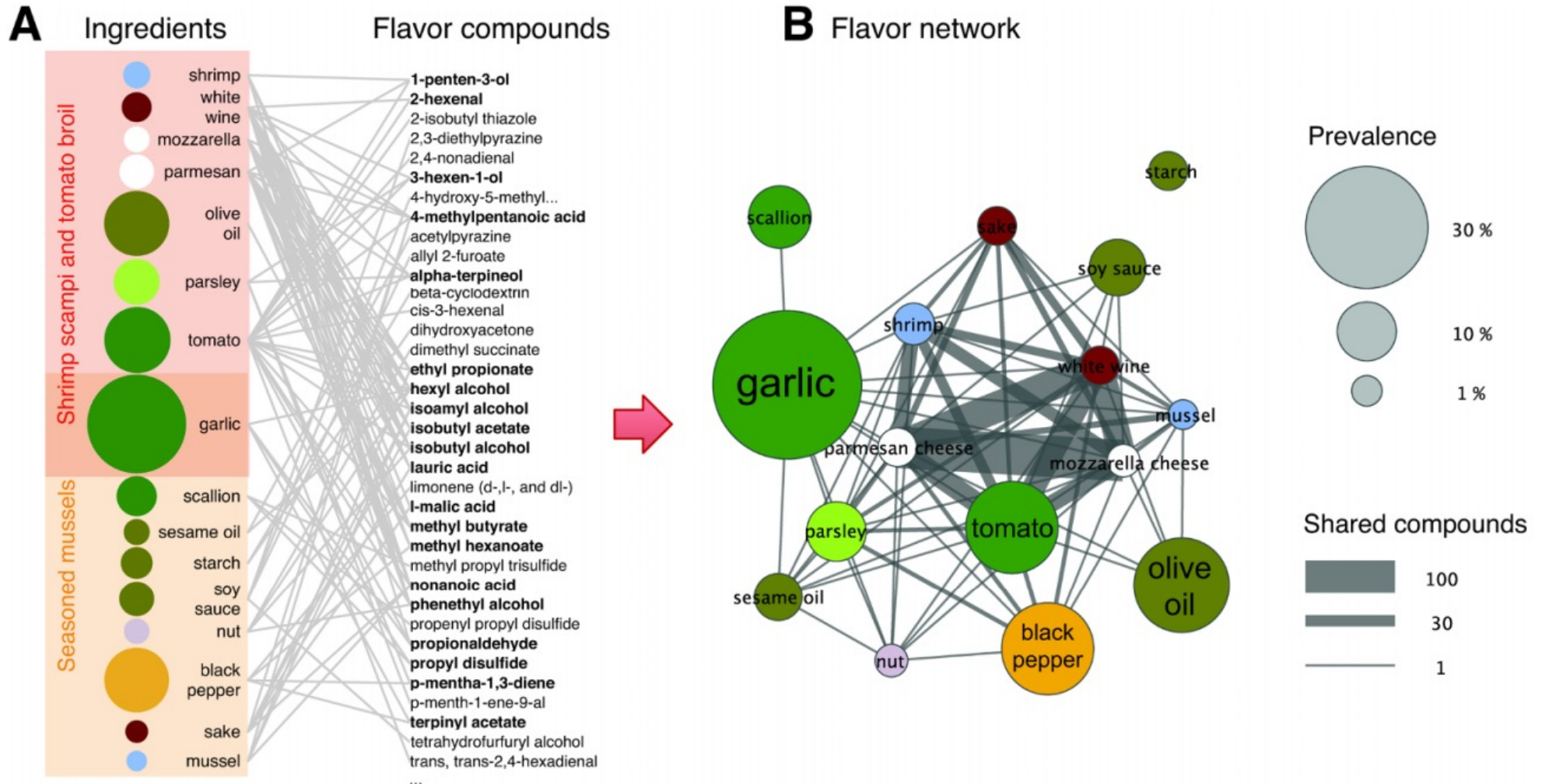
Projection V





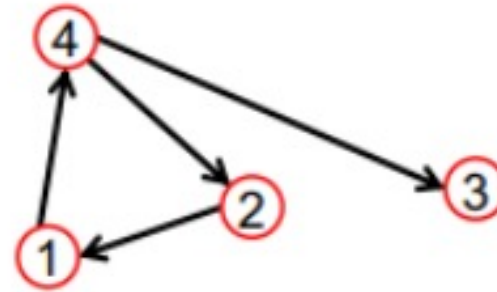
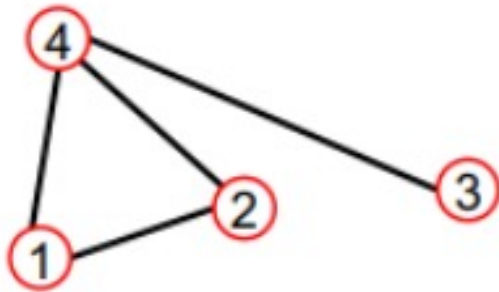


# Graph and Network Theory – Examples



Flavor network and the principles of food pairing

# Graph and Network Theory – Graph Adjacency Matrix



$A_{ij} = 1$  if there is a link from node  $i$  to node  $j$

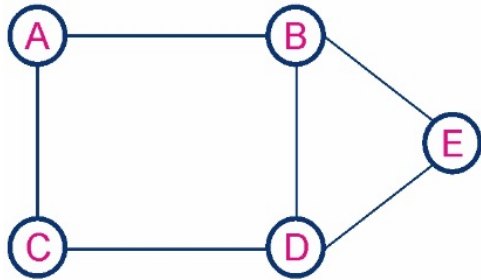
$A_{ij} = 0$  otherwise

$$A = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

$$A = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

Note that for a directed graph (right) the matrix is not symmetric.

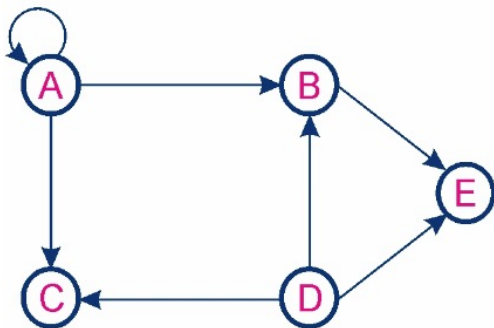
# Graph and Network Theory – Graph Adjacency Matrix



Undirected Graph

	A	B	C	D	E
A	0	1	1	0	0
B	1	0	0	1	1
C	1	0	0	1	0
D	0	1	1	0	1
E	0	1	0	1	0

$$A_{ij}$$

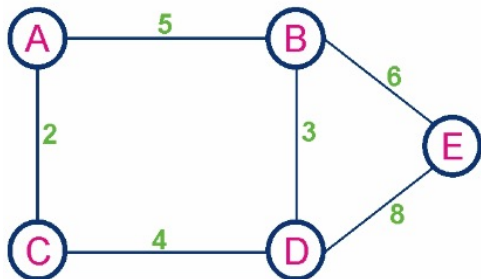


Directed Graph

	A	B	C	D	E
A	1	1	1	0	0
B	0	0	0	0	1
C	0	0	0	0	0
D	0	1	1	0	1
E	0	0	0	0	0

$$D_i$$

$$D_i^+$$

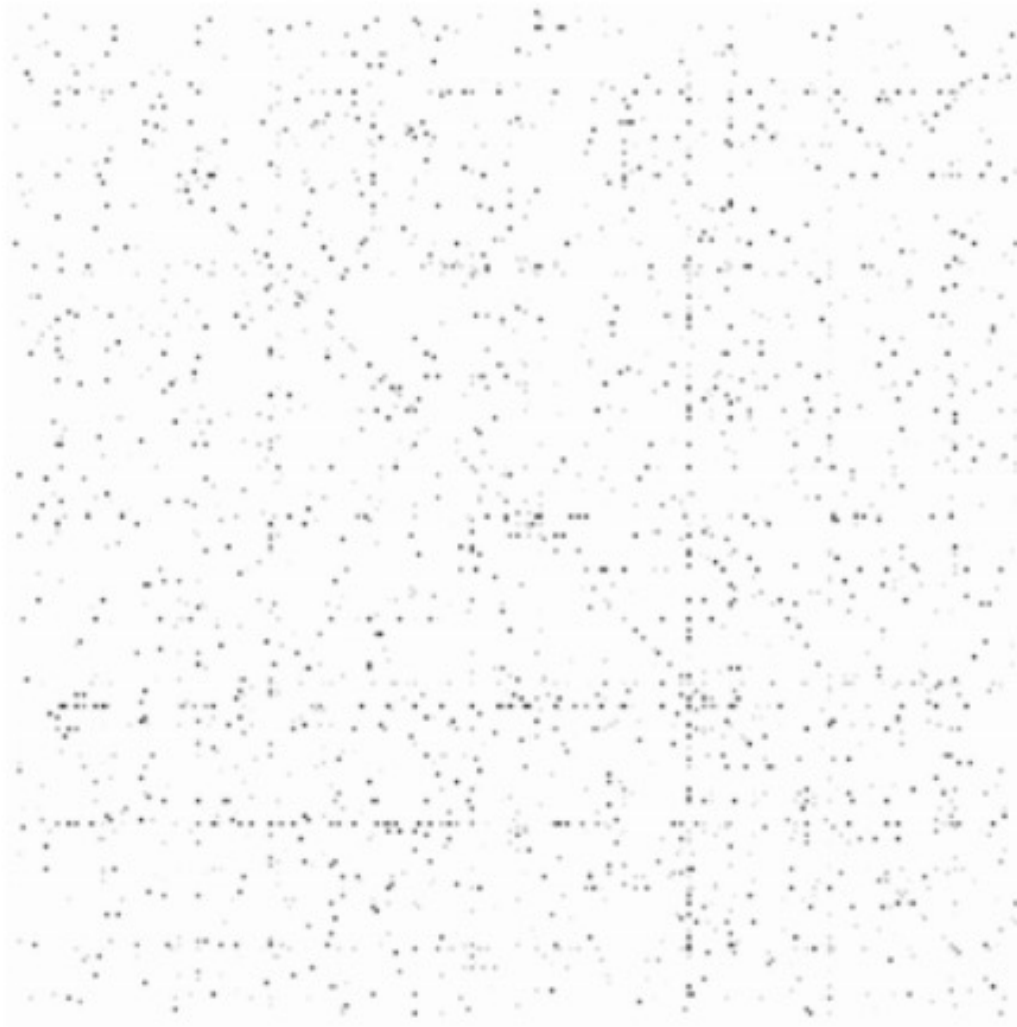


Weighted Graph

	A	B	C	D	E
A	0	5	2	0	0
B	5	0	0	3	6
C	2	0	0	4	0
D	0	3	4	0	8
E	0	6	0	8	0

$$D_i^-$$

# Graph and Network Theory – Adjacency Matrix



**Most adjacency matrices are sparse**



Most real-world networks are **sparse**

$$E \ll E_{\max} \quad (\text{or } \bar{k} \ll N-1)$$

WWW (Stanford-Berkeley):	$N=319,717$	$\langle k \rangle=9.65$
Social networks (LinkedIn):	$N=6,946,668$	$\langle k \rangle=8.87$
Communication (MSN IM):	$N=242,720,596$	$\langle k \rangle=11.1$
Coauthorships (DBLP):	$N=317,080$	$\langle k \rangle=6.62$
Internet (AS-Skitter):	$N=1,719,037$	$\langle k \rangle=14.91$
Roads (California):	$N=1,957,027$	$\langle k \rangle=2.82$
Proteins ( <i>S. Cerevisiae</i> ):	$N=1,870$	$\langle k \rangle=2.39$

(Source: Leskovec et al., *Internet Mathematics*, 2009)

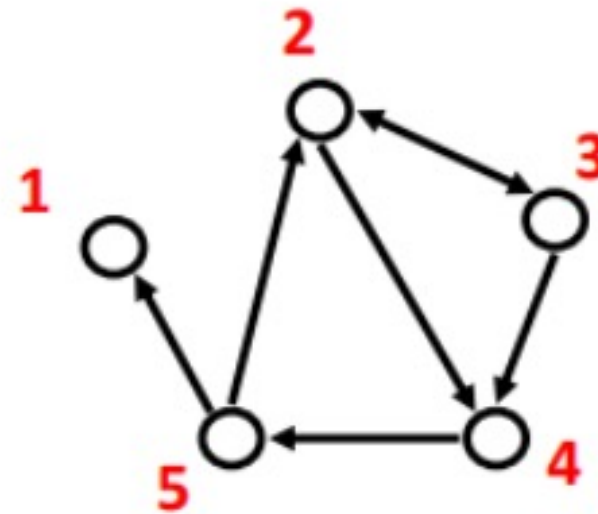
**Consequence: Adjacency matrix is filled with zeros!**

(Density of the matrix ( $E/N^2$ ): WWW= $1.51 \times 10^{-5}$ , MSN IM =  $2.27 \times 10^{-8}$ )



## ■ Represent graph as a set of edges:

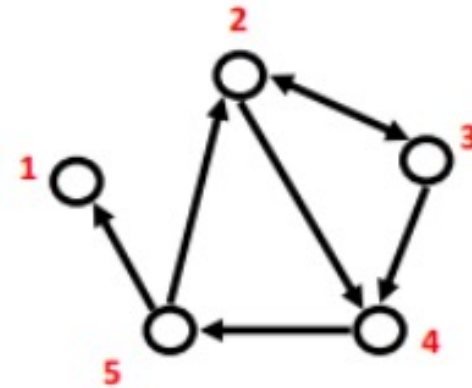
- (2, 3)
- (2, 4)
- (3, 2)
- (3, 4)
- (4, 5)
- (5, 2)
- (5, 1)





## ■ Adjacency list:

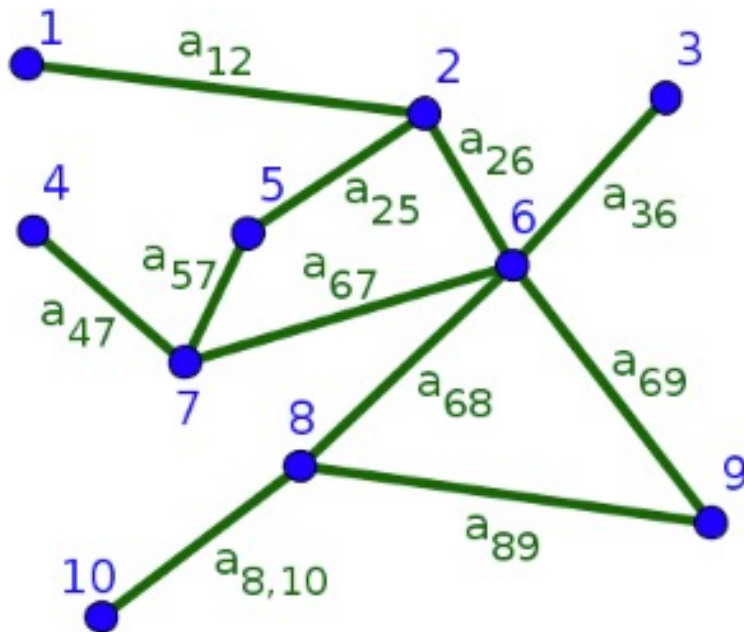
- Easier to work with if network is
  - Large
  - Sparse
- Allows us to quickly retrieve all neighbors of a given node
  - 1:
  - 2: 3, 4
  - 3: 2, 4
  - 4: 5
  - 5: 1, 2



# Graph and Network Theory – Degree Distribution



A node's degree is the number of edges or connections it has to other nodes in a network.



$$k_i = \sum_j a_{ij}$$

$$\sum_{v_i \in \mathcal{V}} k_i = 2 \#edges$$



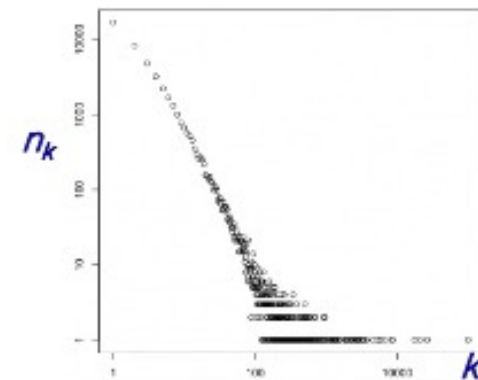
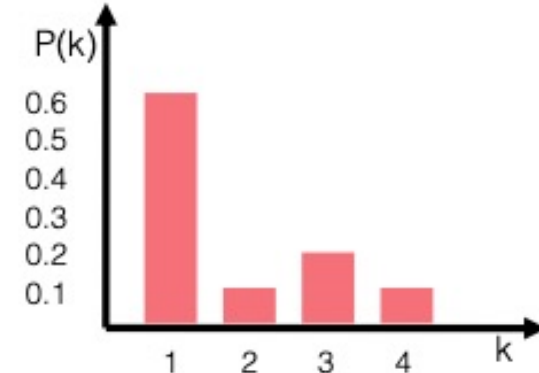
# Graph and Network Theory – Degree Distribution



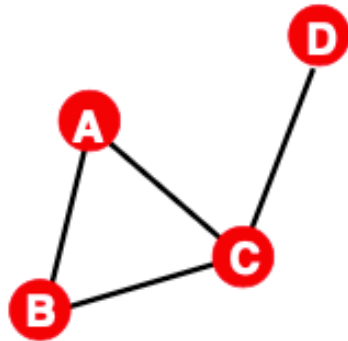
- Degree distribution  $P(k)$ : Probability that a randomly chosen node has degree  $k$

$$n_k = \# \text{ nodes with degree } k$$

- Normalized histogram:  
 $P(k) = n_k / n \rightarrow \text{plot}$



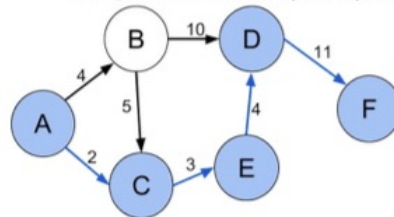
# Graph and Network Theory – Distance and Shortest Path



$$h_{B,D} = 2$$

- **Distance (shortest path, geodesic)** between a pair of nodes is defined as the number of edges along the shortest path connecting the nodes
  - If the two nodes are disconnected, the distance is usually defined as infinite

- **Definition:** find a path between two nodes in a graph, in such a way that the sum of the weights of its constituent edges is minimized
  - Many applications (e.g., road networks, community detection, communications)
- Variants
  - **Single-source** shortest path problem
  - **Single-destination** shortest path problem
  - **All-pairs** shortest path problem



Shortest path (A, C, E, D, F) between vertices A and F in the weighted directed graph

Various algorithms:

- Dijkstra
- Bellman-Ford  
(works with negative edge weights)



- **Diameter:** the maximum (shortest path) distance between any pair of nodes in a graph
- **Average path length** for a connected graph (component) or a strongly connected (component of a) directed graph
  - Many times we compute the average only over the connected pairs of nodes (that is, we ignore “infinite” length paths)

$$\bar{h} = \frac{1}{n(n-1)} \sum_{i,j \neq i} h_{ij} \quad \text{where } h_{ij} \text{ is the distance from node } i \text{ to node } j$$



- **Edge Attributes**
  - Weight (e.g., frequency of communication)
  - Ranking (best friend, second best friend)
  - Type (friend, relative, co-worker)
  - Sign: Friend vs Foe, Trust vs Distrust
  - Properties depending on the structure of the rest of the graph: number of common friends
- **Node Attributes**
  - Bag-of-words feature for documents
  - Customer profile
  - Product meta data

# Any Question?

