# What is Graph Mining?

Plan for today:
- About me and the course ✓
- Go over syllabus ✓
- Go over website ─
- Talk about project
- Talk about paper presentation
- Lecture
  - Define graphs
  - Real-world graph properties
  - Graph Mining and applications
  - Graph computation and processing
  - NetworkX and connected components

Project details

- Groups of 1-4

- Take a look at data repos

- Prior project ideas:

  * Predicting chess match outcomes using a competition network

  * Vertex centrality to identify population centers on a road network

  * Stock market / coin prediction

  * Recommender systems using

\* Recommender systems using SNAP data

- Grad students: use your ongoing project

---

Paper Presentations

- One paper selected per student (from website)

- 20-30 minute presentation

---

Graph Definitions

Graph → tuple of

$U(G)$ = vertex set of graph $G$
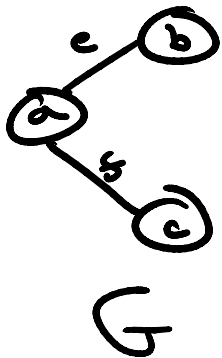$E(G)$ = edge set
$W_v(G)$ = vertex weights
$W_E(G)$ = edge weights
$D_v(G)$ = vertex metadata
$D_E(G)$ = edge metadata

$U_V(G)$ ...

$O_E(G) = $ edge metadata

$$G = \{V, E, W_V, W_E, O_V, O_E\}$$

$G = \{V, E\} \rightarrow$ at a minimum
we need vertices and
edges to detine graph



$$V(G) = \{a, b, c\}$$

$$E(G) = \{e = (a,b), \; f = (a,c)\}$$

For our GM context

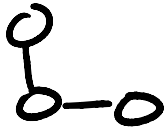vertices $\Rightarrow$ represent discrete
objects, entities, etc.

edges $\Rightarrow$ represent connections
or interactions between
those objects, entities, etc.

weights $\Rightarrow$ strength of connections
for edges or importance
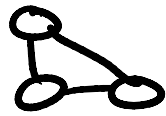or some other measure
for vertices

or some other measure
for vertices

metadata => generic catch-all
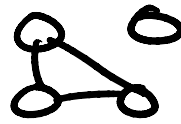for other vertex/edge
labels or properties
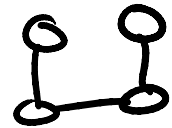
Temporal graphs => graphs that
evolve over time

$G_{t=0}$   $G_{t=1}$   $G_{t=2}$   $G_{t=3}$

---

## Real World Graphs

↳ social networks, info networks
biological networks, etc.

These graphs all have similar properties

Sparsity: $|E(G)| <<< |V(G)|^2$

Degree skew: # of low degree
vertices >>> # of
high degree vertices

$d(v)$ = # of edges
attached to v

$d(v) = $ # of edges attached to $v$

vertices ... ... ... high degree vertices

# of verts w/ degree

degree

Hubs: high degree or otherwise "important" vertices

Irregular: information, social, etc. networks are often not physically constrained

Small-world: average shortest path lengths are small relative to $|V(G)|$

"6 degrees of Kevin Bacon"

Note ✉ : Properties are typical but not every real-world graph has them

# Graph Processing

## "Vertex-centric" processing

- Each $v \in V(G)$ has same state $S(v)$
- We iteratively update $S(v)$
  $\rightarrow$ often using the states of $v$'s neighbors
- Many/most graph algorithms can be implemented in a vertex-centric way

## Algorithmically:

Input: $G = \{V, E\}$

For all $v \in V(G)$:

$\quad S(v) = \text{initialize}()$

For some # of iterations:

$\quad$ update AlgoState()

$\quad$ For all $v \in V(G)$:   $\leftarrow$ neighborhood of $v$

$\quad\quad$ For all $u \in N(v)$:

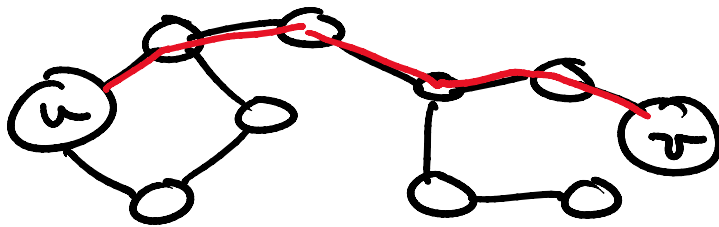$\quad\quad\quad S(v) = \text{updateState}(S(v), S(u))$

For all $v \in V(G)$:
  $S(v) = $ finalize()

---

Implement connectivity decomposition
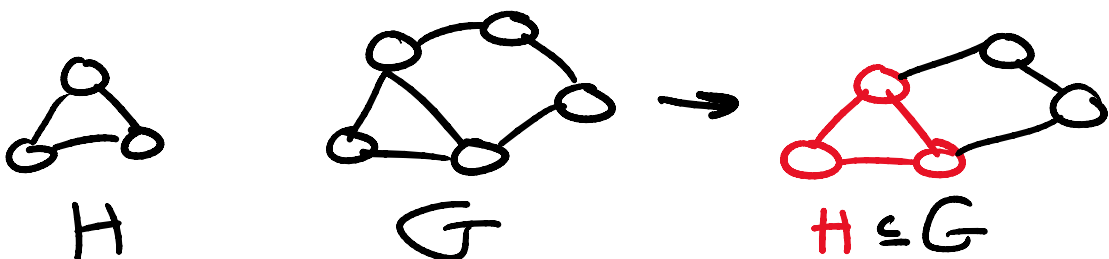  ↓
    $u, v$ are connected if there
    exists a path along edges
    from $v$ to $u$



  $G$ is connected if there exists
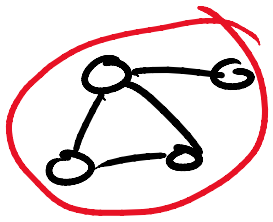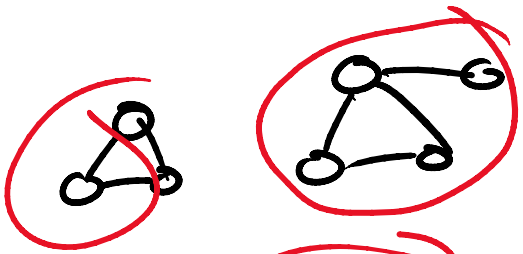  such a path for all $u, v$ pairs

  Connectivity decomposition: find all
  maximal connected subgraphs
  within same $G$



H          G          H ⊆ G

H          G          $H \subseteq G$



H is a subgraph of G

3 connected maximal
subgraphs

G

Easy way to solve using our
framework:

Each vertex gets a unique label


do
   all vertices assumes maximum
      label from its 1-hop neighborhood
while same vertex updates a label

at the end → each unique
label corresponds to a
unique connected component