

Recall: Graph Mining is...

Link Prediction

Centrality

Clustering/CO

Vertex classification

Now: Subgraph mining

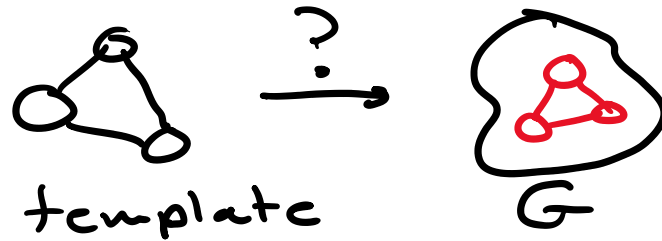
you can tell it's relevant
to graph mining

Many different problems fit
under this umbrella



- Triangle counting (clustering/CC, triadic closure, etc)
- Template matching
aka subgraph isomorphism

aka subgraph isomorphism



also: subgraph counting/enumeration
↳ How many? ↳ where?

- Motif/antimotif detection

what subgraphs show up more/less frequently than expected

aka H3p2

- subgraph-based comparative analytics

compare two or more graph via constituent subgraphs
(usually normalized)

- graph alignment (next class)

approximate (sub)graph isomorphism

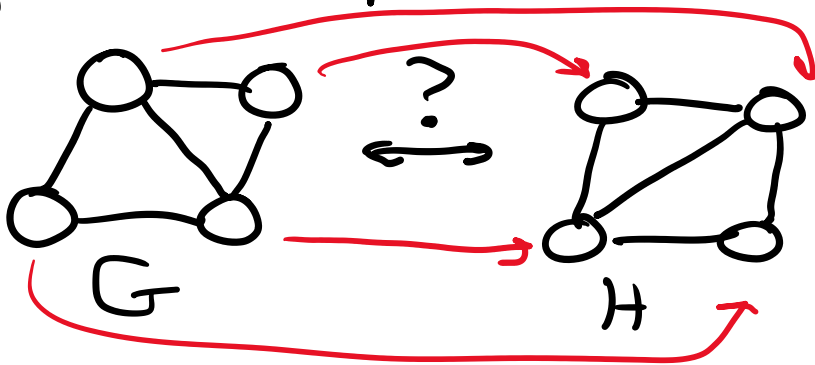
we align networks based on some cost function





Subgraph isomorphism

Graph isomorphism



Are G and H equivalent?

↳ Tougher problem than you might think

(complexity unknown)

Usually: algorithms turn G and H into strings and compare them

Useful in (bio)chemistry

What we care more about:

Subgraph isomorphism

subgraph

Subgraph isomorphism

Good

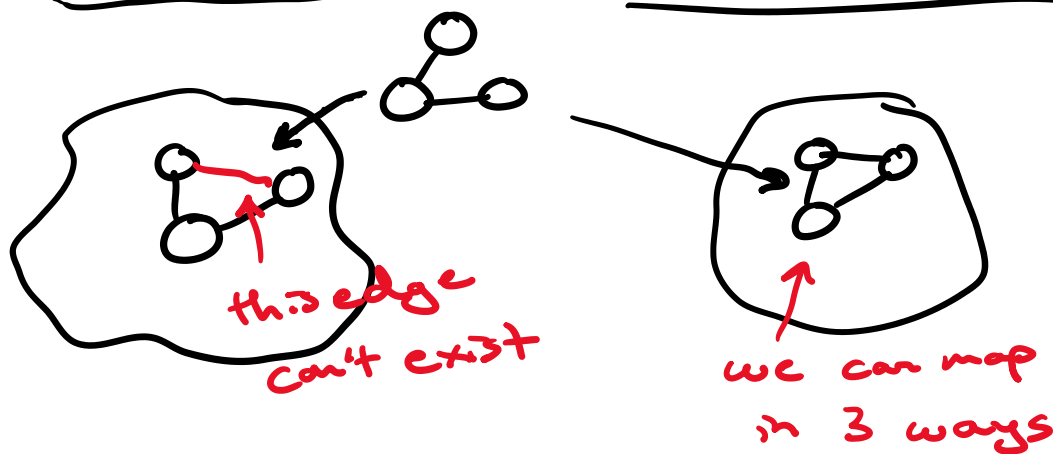
ton of applications
broadly useful
computationally feasible
(in certain cases)

Bad

NP-complete $O(n^3)$
we need to $|V(G)|$
solve for all subgraph-
based analytics
However: can use
approximating counts



Induced vs. non-induced subgraphs



Triangle Counting

Triangle are particularly useful
→ why? very useful for clustering
analytics, prediction, etc.

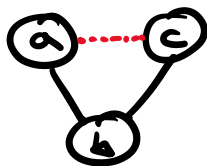
analytics, prediction, etc,

Also: most real, interesting, large networks exhibit some form of clustering

Can be easier to count (enumerate) than larger templates

$O(n^w)$ ← exponent of fast matrix multiplication complexity

↳ essentially computing closure of all wedges



Sparse graphs: $m \ll n^2$

We have $O(m^{3/2})$

↑
operations on edge list itself

Basic Algorithm:

Count = 0

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

For all $u, w \in N(v)$:

if $(u, w) \in E(G)$:

if $(u, w) \in E(G)$:

count $+1$

count $\neq 6$

(or 3)

Most general-purpose subgraph counting algorithms use some variant of the above

Template matching

↳ general term for "subgraph isomorphism"

We consider some template T and search for isomorphic subgraphs in some G

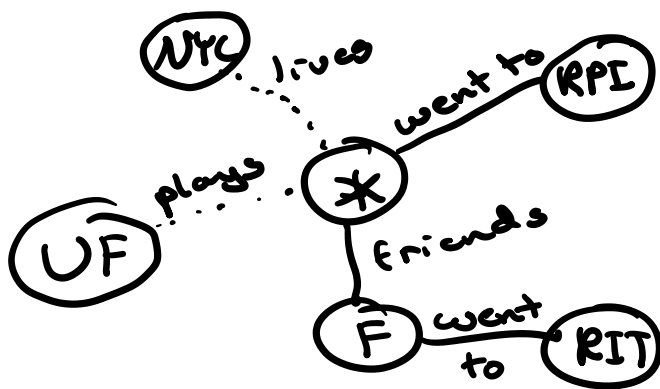
Note: we can also consider labels on our template

↳ Luckily, it makes the problem a lot easier

Template search applications

Facebook Graph Search

↳ "Find all people who went to RPI and are friends with someone who went to RIT and currently live in NYC OR currently play ultimate frisbee"



=> greatly restricts search space

stalkers, int. agency hackers, etc

Unfortunately: immediately abused and removed

Also widely applicable:

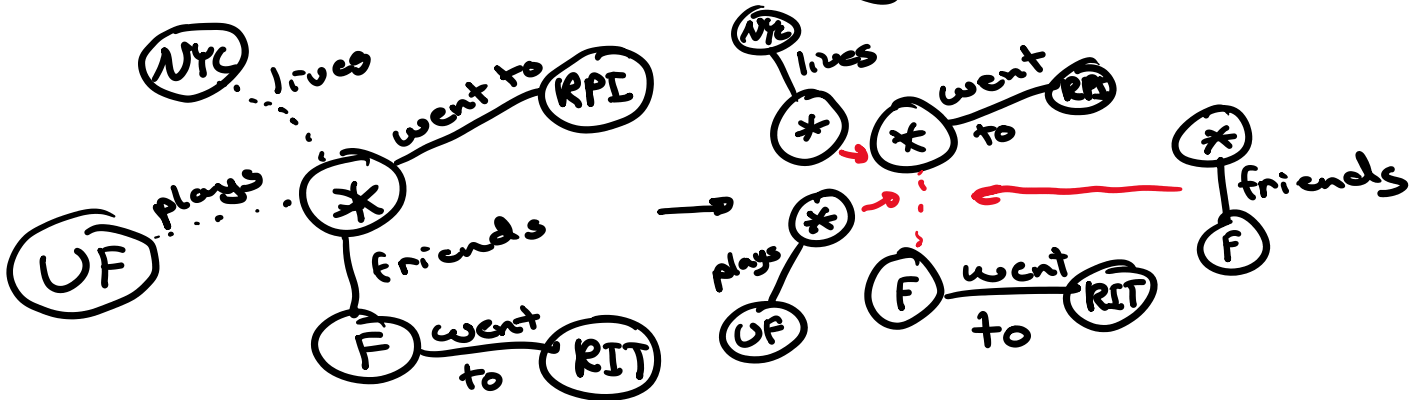
- Social networks, financial networks
- Information networks,
 ETC.
 ↑
 flow of money

Algorithms for template search:

Naive brute force $\rightarrow O(n^2)$

Dynamic Programming

\rightarrow Better, especially for labeled data



However: worst-case complexity still the same

Approximate Algorithms:

we can actually improve our bound

But: we either aren't getting an exact count OR we aren't enumerating all subgraphs

How: sampling, color-coding

\nearrow
only search a portion of the network

\uparrow
we color G randomly with at least k colors, and

partition of the
network

with at least k colors,
and only search for
"colorful" embedding

↑
all vertices have
a different color

One more thing: Parallelization



applies to all the above

=> generally, parallelizing where
in the graph we're searching