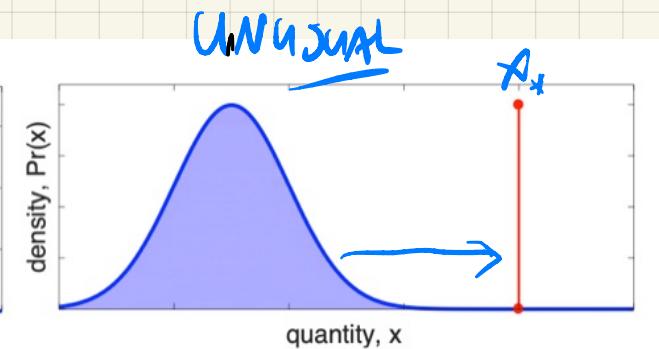
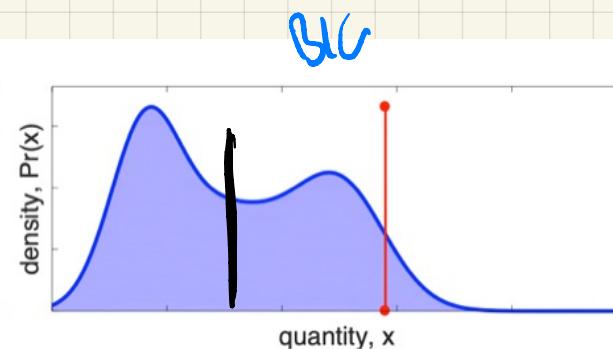
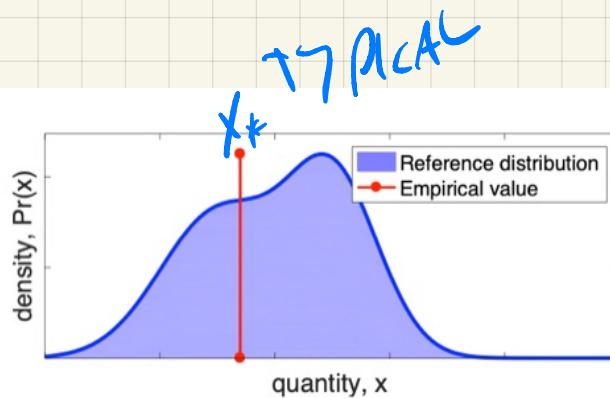


Lecture 3: Random graph models

How do we interpret the value of $\langle k \rangle$, $\langle l \rangle$? Is it **big** or **small**?
Is it **typical** or **unusual**?

reference distribution $Pr(x)$

empirical x_*



For networks \Rightarrow a random graph model defines the reference $Pr(x)$

\Rightarrow many different models, with different assumptions

3 types of random graphs (for now...)

1) Erdős - Rényi (ER) random graph

* edge density + iid

2) Configuration Model (and the Chung-Lu model)

* degree structure (otherwise random)

3) Modular random graphs

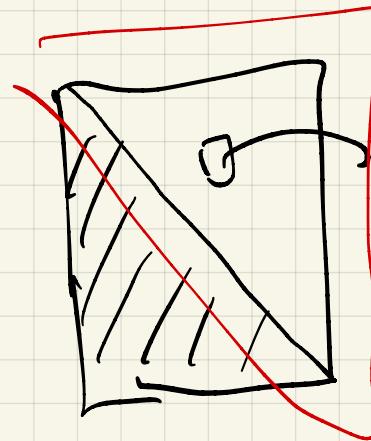
* groups of nodes

The Erdős-Rényi random graph

$$2^{\binom{n}{2}} = O(2^{n^2})$$

$G(n, p)$
 # nodes prb. $(i, j) \in E$

• defines an ensemble of graphs $\Pr(G|p)$



$$\text{flip coin } \Pr(H) = p$$

unweighted

$$\rightarrow \forall i > j \quad A_{ij} = A_{ji} = \underbrace{\quad}_{\text{undirected}}$$

$$\begin{cases} 1 & \text{with probability } p \\ 0 & \text{otherwise} \end{cases}$$

expected degree

$$p = \frac{c}{n-1} \quad \begin{matrix} \checkmark \\ \nearrow \end{matrix} \quad \begin{matrix} \checkmark \\ \searrow \end{matrix}$$

c # possible adjacencies

* $p = 1$ \rightarrow complete graph on n nodes

* $p = 0$ \rightarrow empty graph on n nodes

Properties of ER graphs

- simple graph
- connectivity is "homogeneous" \Rightarrow edge "dense" + iid
- deg. dist. $\Pr(k) \approx \text{Poisson}(c)$ $c = p(n-1)$
- diameter and MHD $\rightarrow O(\log n)$ \Rightarrow "small-world"-like nets
- clustering $\Delta \rightarrow O(1/n)$
- largest connected component (LCC) $\approx O(n)$ when $c > 1$

Generating a $G(n,p)$ network

Initial empty graph $G = (V, E)$ $|V| = n$ $E = \emptyset$

$O(n^2)$ for each possible (i,j) , draw $r \sim \text{Uniform}(0,1)$
if $r \leq p$, add (i,j) to E

can do it in \mathbb{Z}
 $O(n+m)$ time
if you're clever

Mean degree

$$\langle k \rangle = c = \sum_{j=1}^{n-1} p = p \left(\sum_{j=1}^{n-1} 1 \right) = \underline{p(n-1)}$$

$$\langle m \rangle = \sum_{i=1}^n \sum_{j>i} p = p \left(\frac{n}{2} \right) = p \left(\frac{n(n-1)}{2} \right)$$

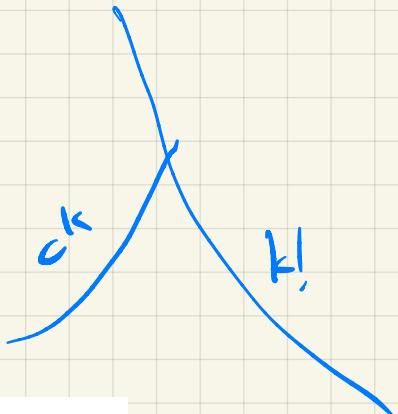
$$\langle b \rangle = \frac{\sum}{n} \left(p \left(\frac{n(n-1)}{2} \right) \right) = \underline{p(n-1)}$$

recall
 $\langle k \rangle = \frac{2n}{n}$

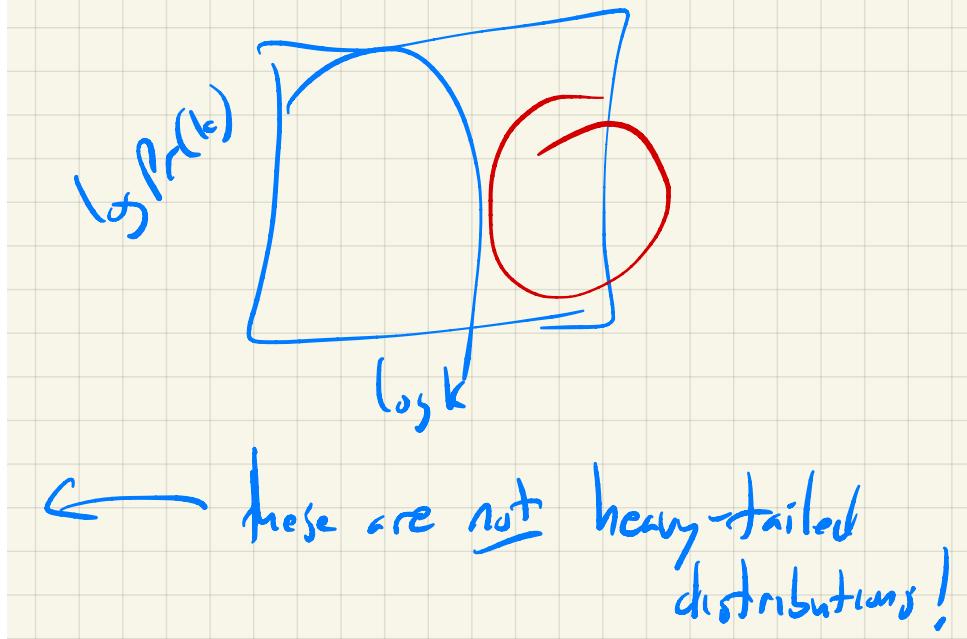
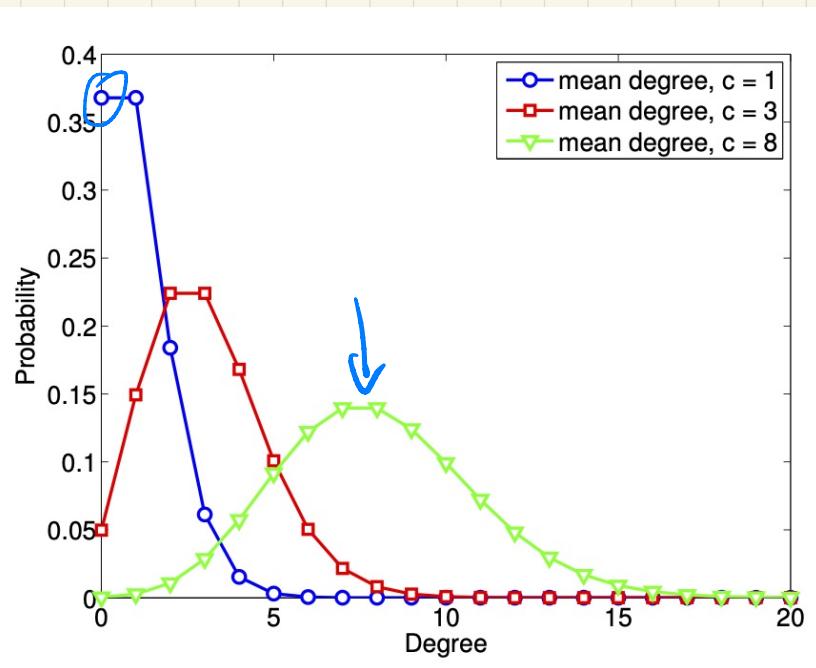
degree distribution

$$\Pr(k) = \binom{n-1}{k} p^k (1-p)^{n-1-k}$$

$$= \frac{c^k}{k!} e^{-c}$$



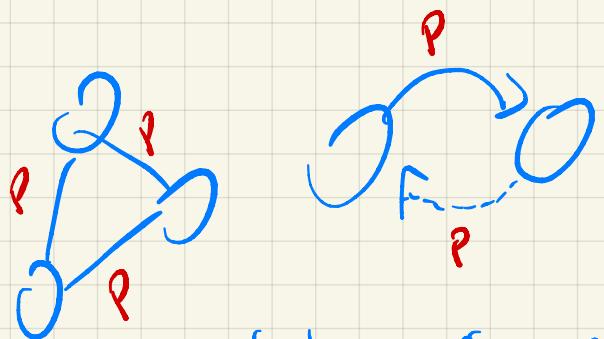
Most nets are sparse
 $\langle k \rangle = O(1)$
 $= p(n-1)$
 $p = \frac{c}{n-1}$



Motifs, reciprocity, clustering coefficient

$$\text{reciprocity } r = \frac{\# \text{ of recip. links}}{\# \text{ of links}} = \frac{(n^2-n) p^2}{(n^2-n) p} = p = O(k_n)$$

$$= \frac{c}{n-1}$$



$$\text{clust. coeff } C = \frac{\#\Delta}{\#\Lambda} = \frac{\binom{n}{3} p^3}{\binom{n}{2} p^2} = p = O(k_n)$$

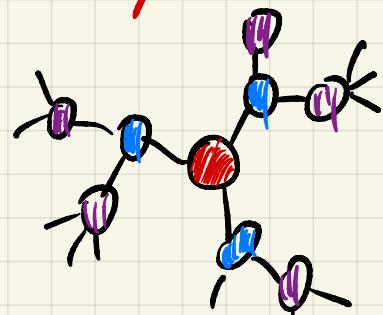
Diameter and MGD

$$(c-1)^{d_{\max}} = n$$

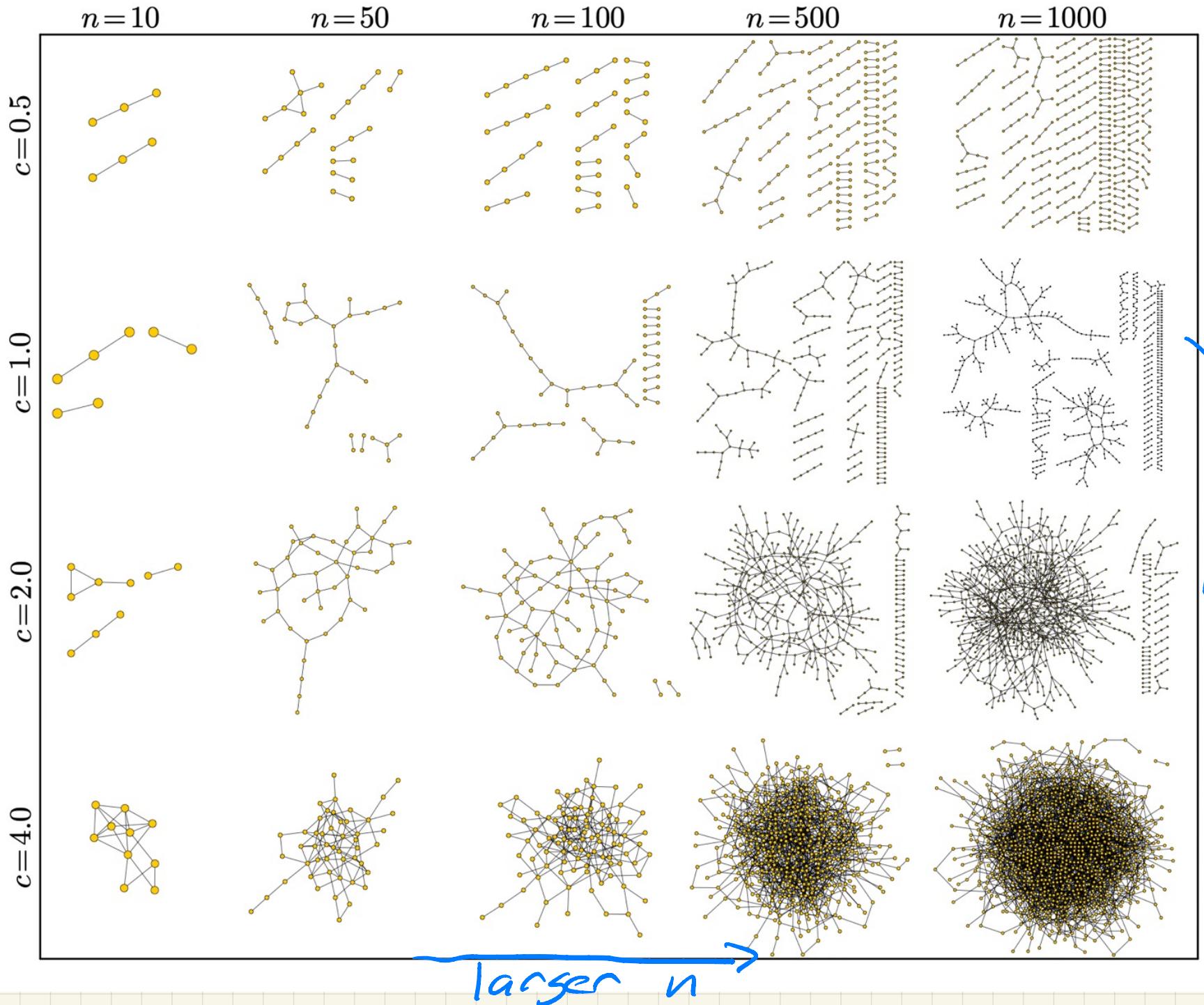
↓

$$d_{\max} = O(\log n)$$

ER graphs are
locally tree-like

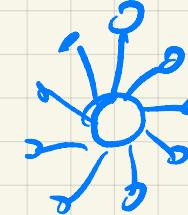
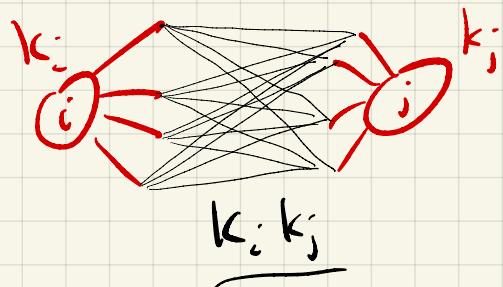


What do ER $[G(n,p)]$ graphs look like?



The Configuration Model

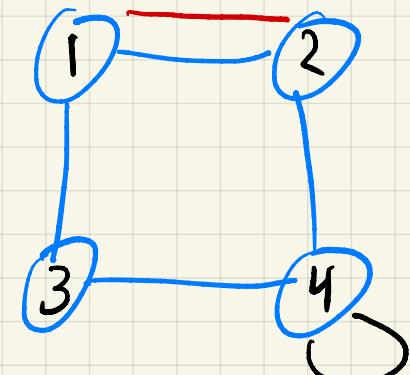
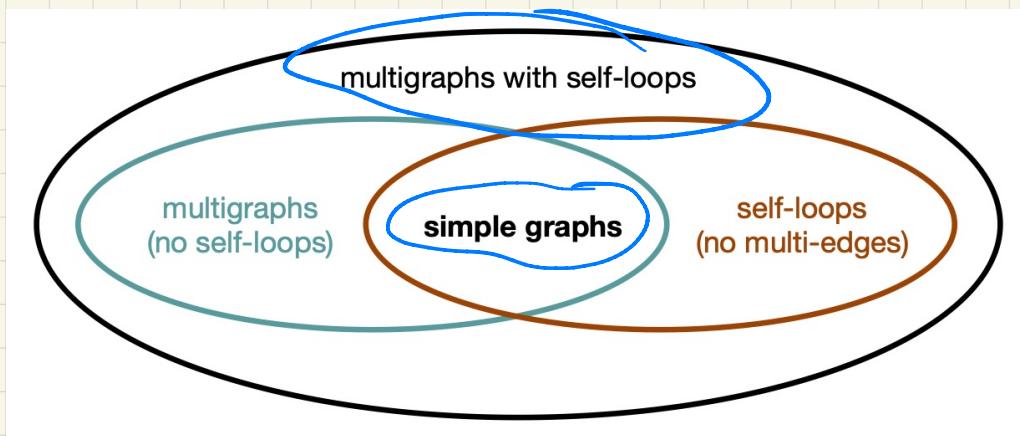
assume degree seq $\vec{k} = \{k_1, k_2, k_3, \dots, k_n\}$
take from data



$$\forall_{i > j} A_{ij} = A_{ji} = \begin{cases} 1 & \text{with probability } \propto k_i k_j \\ 0 & \text{otherwise} \end{cases}$$

(\propto means "is proportional to")

contrast this with ER model
of $p = \frac{2n}{(n-1)}$ (a constant value)



Properties of the configuration model

$G(n, \bar{k})$

- four flavors: simple or multigraph / w/ or w/o self loops
- connectivity specified by \bar{k}
- diameter and $\langle \text{cc} \rangle \sim O(\lg n)$
- motif freq $\Delta D \sim O(1/n)$
- LCC $\sim \Theta(n)$ if G not too sparse

(4)

locally tree-like

depends on mean $\langle k \rangle$
and variance $\langle k^2 \rangle$

generating a config. model

2 facts make it tricky:

1) not every \bar{k} is graphical

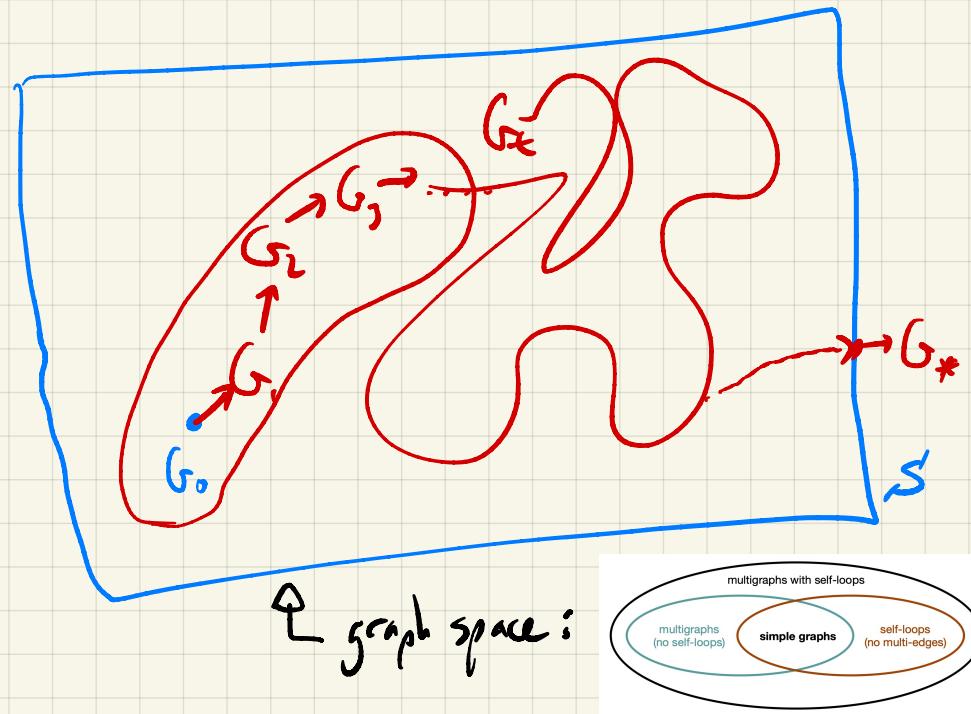
2) no single procedure for $G \sim \text{Pr}(G)$

Solution: Markov chain Monte Carlo (MCMC)

two parts:

1) generate $\xrightarrow{\sim} G_0 \in \mathcal{S}$
graph space

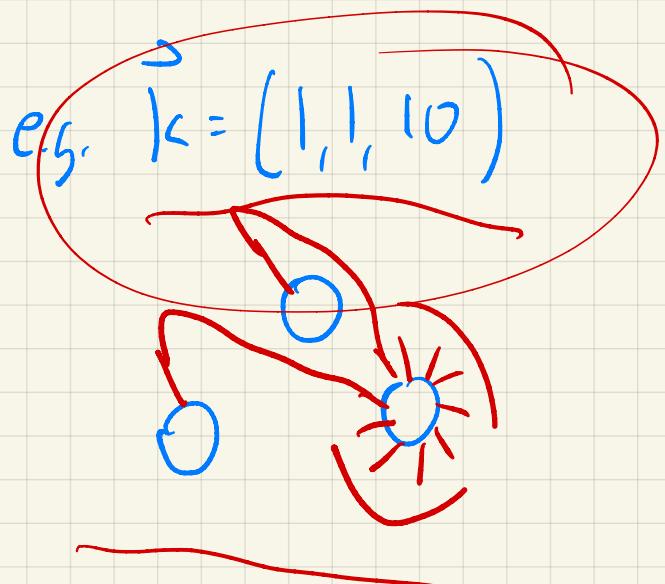
2) $g: (G_t) \rightarrow G_{t+1} \in \mathcal{S}$



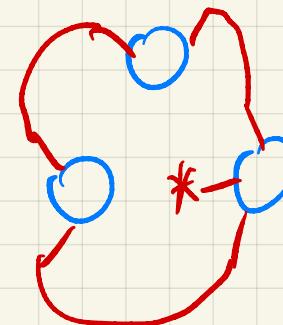
$$\Pr(G | \tilde{v})$$

Is it graphical?

- Can we make a graph with $\vec{k} = \{k_1, k_2, k_3, \dots, k_n\}$?
↳ if so, then every "stub" can be paired up correctly



or $\vec{k} = (2, 2, 3)$



$\cancel{\text{if } k_i = \text{odd}}$
 $\cancel{\text{not graphical}}$

• Harrel-Hakimi:

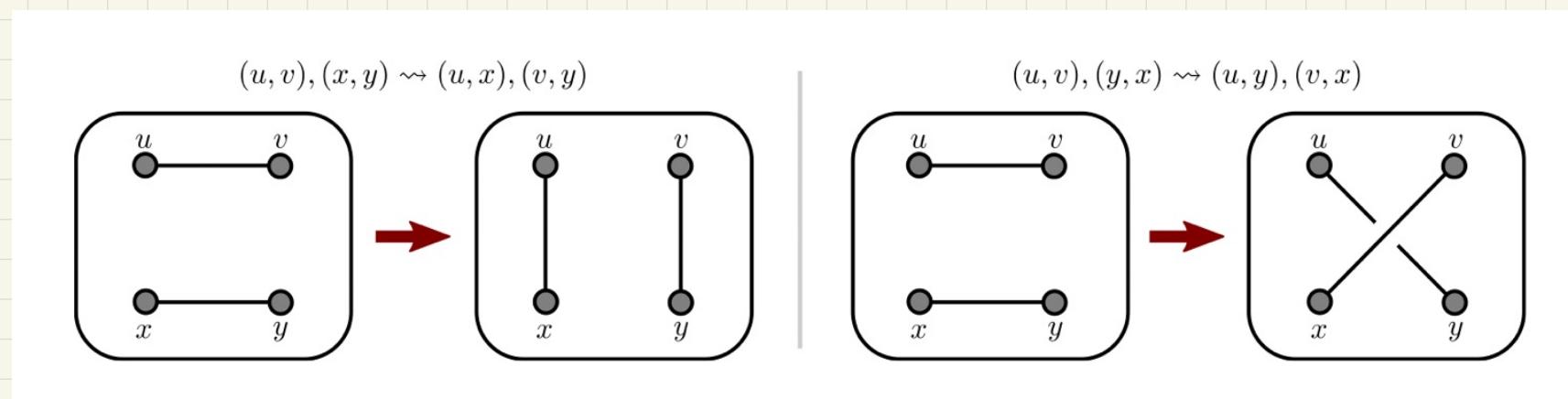
$\text{sort } \vec{k}$
 $k_{n-1}, k_n)$

makes a synthetic
seed for MCMC

G_0

g : the double edge swap: given $G_t \in S$, apply $g(G_t) \rightarrow G_{t+1} \in S$

- (select) choose uniformly at random
- (swap) choose output 1 or output 2 with equal prob
- (accept) if $G_{t+1} \in S'$



degree preserving: k_u, k_v, k_x, k_y don't change

$G_0, G_1, G_2, \dots, G_t$ \equiv $G_0, g(G_0), g(g(G_0)), \dots, \underbrace{g^t(G_0)}$

the chain

How long of a chain? (mixing times)

want independent samples:

$$G_t \text{ and } G_{t+n_0} \text{ sampling gap}$$

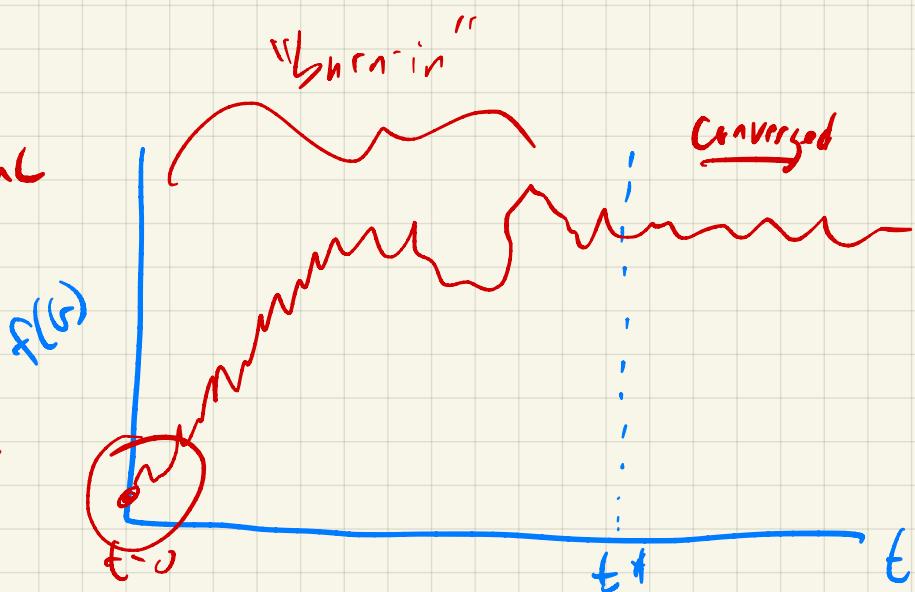
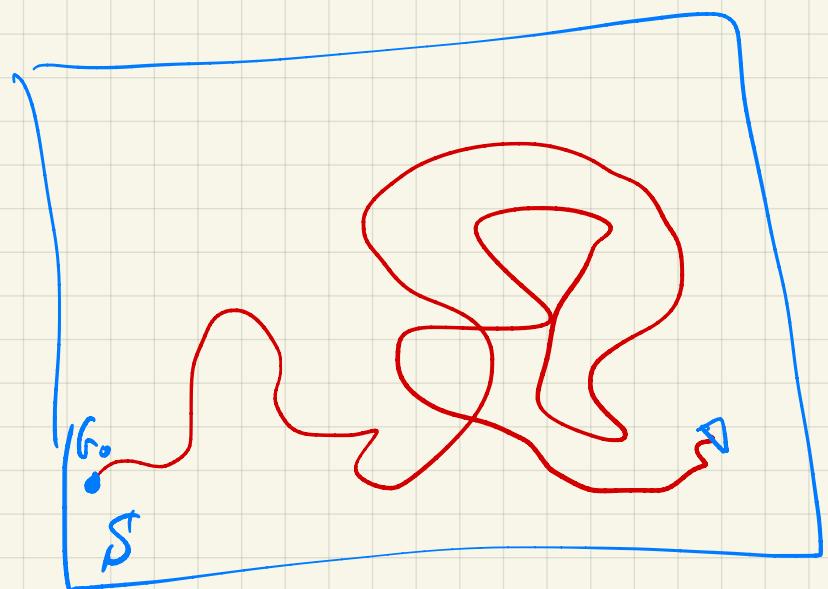
$$n_0 = \frac{M}{3} + 300$$

want $\Pr(b)$ to be uniform over S

detect convergence

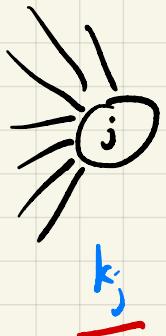
test or two samples from mcmc

summary statistic



Chung-Lu Model (simple graphs or directed graphs)

- like configuration model, but only $E[\vec{k}]$ rather than exactly \vec{k}



$$P_{ij} = k_i \left(\frac{k_j}{2n - k_i} \right) = \frac{k_i \cdot k_j}{(2n - k_i)}$$

if $\max_i k_i < 2n - k_i$



$$\forall_{i > j} A_{ij} = A_{ji} = \begin{cases} 1 & \text{with prob. } p_{ij} \\ 0 & \text{otherwise} \end{cases}$$

generating Chung-Lu:

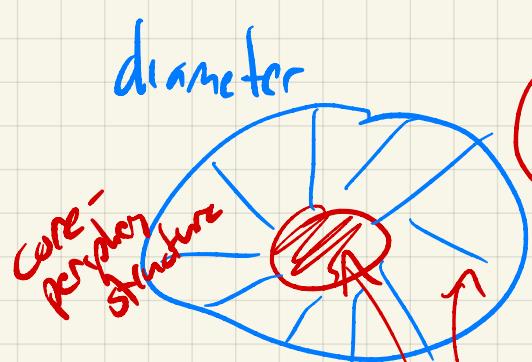
for $i = 1$ to n
 for $j = i+1$ to n
 $A_{ij} = A_{ji} = 1$

Clustering coefficient

$$C = \frac{1}{n} \left[\frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle^3} \right]^2 = O(1/n) \text{ when } \langle k^2 \rangle \text{ is finite}$$

Giant component [$LCC = \Theta(n)$]

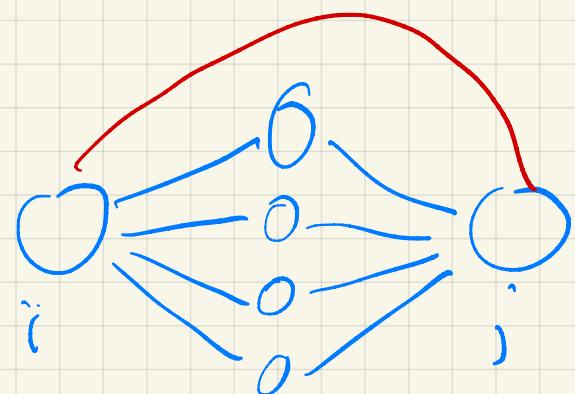
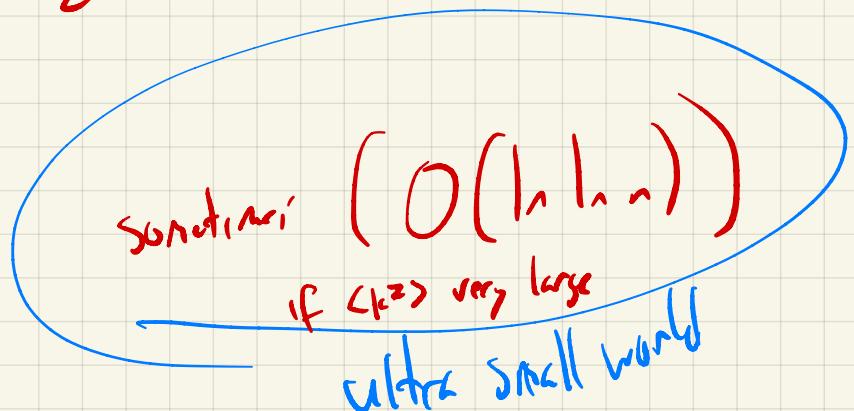
$$\langle k^2 \rangle - 2\langle k \rangle > 0$$



of common neighbors

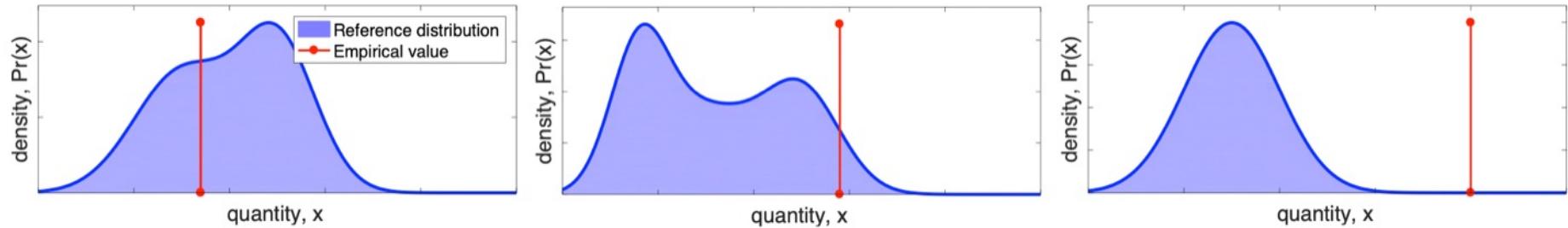
$$n_{ij} = P_{ij} \left(\frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle} \right)$$

high degree nodes here
low degree nodes here



Random graphs as null Models

recall: random graphs produce $\Pr(G)$ as a reference distribution



Null Model:

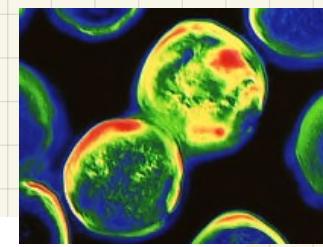
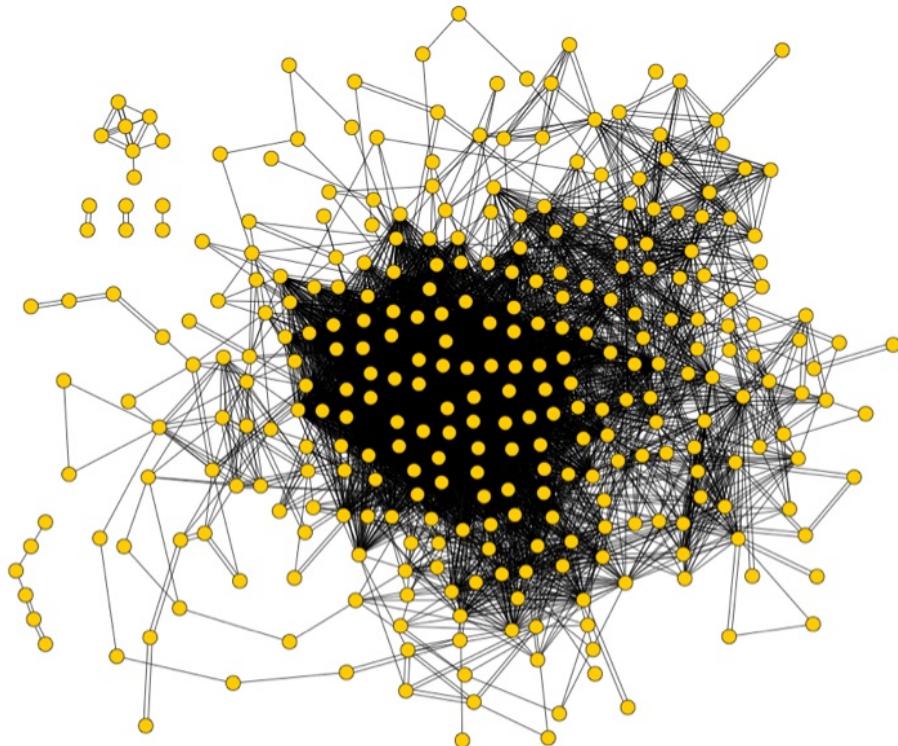
How much of an observed pattern is explained by {edge density or degrees} alone, under randomness?

need $\Pr(G | \theta)$
est. from our data

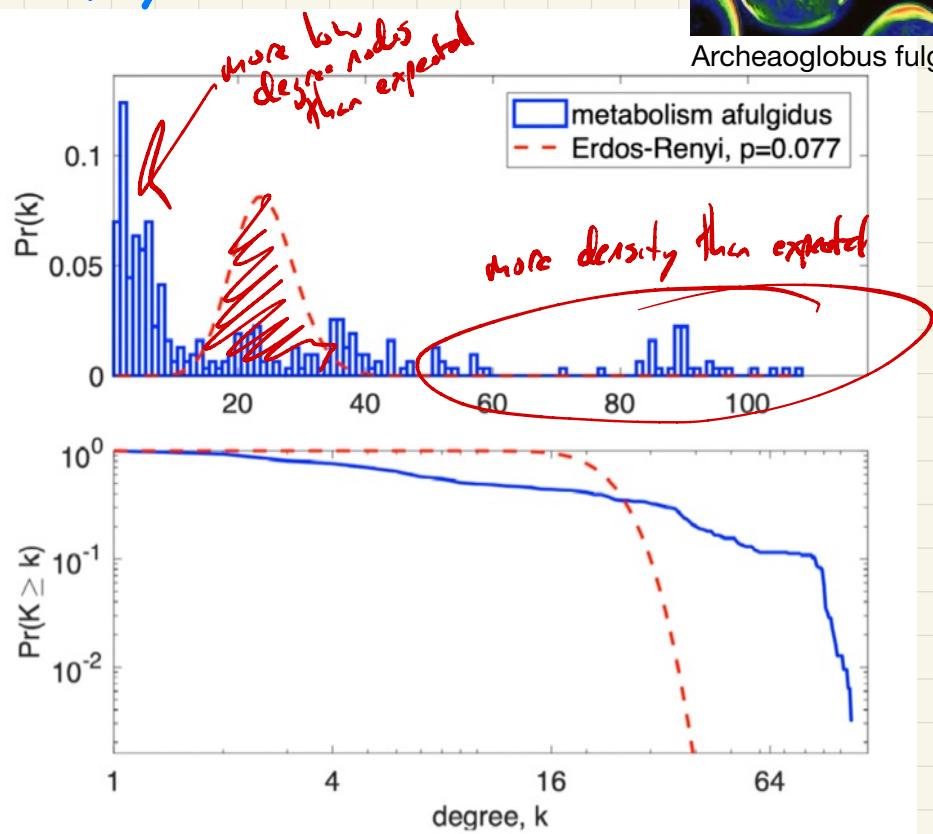
$$\hat{P} = \frac{\langle k \rangle}{n-1}$$
$$k = \{k_1, k_2, \dots, k_n\}$$

$$n = 315$$

$m = 3793$ (undirected edges)



Archeoglobus fulgidus



$$\langle k \rangle = 24.08$$

LCC: 296 nodes (94%)

$$\text{ER model: } \hat{\rho} = \frac{\langle k \rangle}{n-1} = 0.077$$

$$k_{\max} = 108$$

$$\Pr(K_{\max} = 108 | \langle k \rangle, \hat{\rho}) \\ = 4.3937 \times 10^{-36}$$

