

2. a) $\int_0^{\infty} \frac{\alpha x_0^\alpha}{x^{\alpha+1}} dx = \alpha x_0^\alpha \int_{x_0}^{\infty} x^{-\alpha-1} dx = \alpha x_0^\alpha \left(\frac{1}{-\alpha} x^{-\alpha} \right) \Big|_{x_0}^{\infty}$
 $= \alpha x_0^\alpha \left(\frac{1}{-\alpha} - \frac{1}{x_0^\alpha} \right)$
 $\bar{F}(X \geq x) = \int_x^{\infty} \frac{\alpha x_0^\alpha}{x^{\alpha+1}} dx$
 $= \alpha x_0^\alpha \left(\frac{1}{-\alpha} x^{-\alpha} \right) \Big|_x^{\infty}$
 $= \frac{\alpha x_0^\alpha}{-\alpha} \left(-\frac{1}{x^\alpha} \right) = -x_0^\alpha \left(-\frac{1}{x^\alpha} \right)$
 $= \frac{\alpha x_0^\alpha}{-\alpha} \left(\frac{1}{\infty^\alpha} - \frac{1}{x^\alpha} \right) = -x_0^\alpha \left(-\frac{1}{x^\alpha} \right) = \left(\frac{x_0}{x} \right)^\alpha$
 $\bar{F}(X \geq x) = \left(\frac{x_0}{x} \right)^\alpha$ = 1 \Rightarrow valid PDF

c) $\ln(\alpha) = \log \left(\prod_{i=1}^n \frac{\alpha x_0^\alpha}{x_i^{\alpha+1}} \right) = \sum_{i=1}^n \log \left(\frac{\alpha x_0^\alpha}{x_i^{\alpha+1}} \right) = \sum_{i=1}^n (\log(\alpha x_0^\alpha) - \log(x_i^{\alpha+1}))$
 $\ln(\alpha) = \sum_{i=1}^n (\log \alpha + \alpha \log x_0 - (\alpha+1) \log x_i)$
 $\frac{d \ln(\alpha)}{d \alpha} = \sum_{i=1}^n \left(\frac{1}{\alpha} + \log x_0 - \log x_i \right)$
 $0 = \frac{n}{\alpha} + n \log x_0 - \sum_{i=1}^n \log x_i$
 $\frac{n}{\alpha} = \sum_{i=1}^n \log x_i - n \log x_0 \Rightarrow \hat{\alpha} = \frac{n}{\sum_{i=1}^n \log x_i - n \log x_0}$

3. a) $cl(G) = \frac{3 \cdot \# \text{ triangles}}{\# \text{ connected triples}}$; $d_i = 2k \forall i$

triangles per node, clockwise: choose 2 of the node's clockwise nbs within k hops
 avoid double-counting $\Rightarrow \binom{k}{2} = \frac{k!}{2!(k-2)!} = \frac{k(k-1)}{2}$

connected triplets per node: $\binom{2k}{2} = \frac{2k(2k-1)}{2}$

$\Rightarrow cl(G) = \frac{3nk(k-1)}{2nk(2k-1)} = \frac{3(k-1)}{2(2k-1)} = \frac{3k-3}{4k-2}$

b) • # triangles;
- still $\frac{n k(k-1)}{2}$ original

- new:

• can be made by connecting nodes between $k+1$ and $2k$ hops

away; prob that they're connected is $\frac{1}{2} n 2k \cdot p / \frac{n(n-1)}{2} \approx \frac{2kp}{n}$

• # triples;
- still $\frac{2nk(2k-1)}{2}$ original

- each new edge creates triples w/all of the original $2k$ nodes @ either of its ends ($= (2k)^2$ triples)

• n possible, w/prob. p

$$\Rightarrow (2k)^2 np$$

→ however, as $n \rightarrow \infty$ this term goes to 0, so it's negligible.

$$\Rightarrow \# \Delta \approx \frac{nk(k-1)}{2}$$

- new edges also create pairs w/other new edges:

- l new edges $\rightarrow \binom{l}{2} = \frac{l(l-1)}{2}$ new triples

- expected # of new edges to a node is $\frac{2|E|p}{n} = \frac{2(\frac{2nk}{2}p)}{n} = \frac{2kp}{n}$

\Rightarrow exp. new triples is $n \frac{2kp(2kp-1)}{2} \approx (2k)^2 p^2 n$

$$\Rightarrow cl(G) = \frac{\frac{1}{2}nk(k-1) \cdot 3}{\frac{1}{2}2nk(2k-1) + (2k)^2 np + \underbrace{\frac{1}{2}n(2k^2)p^2}_{\text{ignore}}} = \frac{\frac{3}{2}nk(k-1)}{nk(2k-1) + (2k)^2 np} = \frac{\frac{3}{2}nk(k-1)}{2nk^2 - nk + 4k^2 np}$$

$$= \frac{3k(k-1)}{2(2k^2 - k + 4k^2 p)}$$

$$= \frac{3k-3}{2(2k-1+4kp)}$$

$$= \frac{3k-3}{4k-2+8kp}$$

4. a) $\overset{i > m}{k_i(i) = M}$ (each node connects to m others @ birth; node i is born @ time i)

b) $\boxed{\frac{dk_i(t)}{dt} \approx \frac{M}{t}}$ (a ratio of $\frac{M}{t}$ nodes have new connections added from time $t \rightarrow t+1$)

$$k_i(t) = \int \frac{M}{t} dt = M \log t + c \rightarrow k_i(i) = m = m \log i + c$$

$$\Rightarrow \underline{c = M - m \log i}$$

c) $k_i(t) = m + m \log\left(\frac{t}{i}\right) \geq d \Rightarrow k_i(t) = m \log t + m - m \log i$

$$m \log\left(\frac{t}{i}\right) \geq d - m$$

$$\frac{t}{i} \geq e^{\frac{d-m}{m}}$$

$$i \leq \frac{t}{e^{\frac{d-m}{m}}} \rightarrow \text{t nodes by time t, so } \boxed{\bar{F}(d) = e^{\frac{(m-d)}{d}}}$$

$$= m + m \log\left(\frac{t}{i}\right) \quad \square$$

$$p(d) = -\frac{d \bar{F}(d)}{d d} = -\frac{d}{d d} e^{\frac{m-d}{d}} = -\frac{d}{d d} e^{\frac{m}{d} - 1} = (e^{\frac{m}{d} - 1}) \cdot -m d^{-2}$$

$$p(d) = \boxed{\frac{1}{m d^2} \cdot e^{\frac{m-d}{d}}}$$

```
In [ ]: import networkx as nx
import numpy as np
import matplotlib.pyplot as plt
from collections import Counter
```

1. ER phase transitions

a) ER graph component sizes

```
In [ ]: def analyze_er_graph_components(n, p_range, graphs_per_p=20):
    means, sds = [], []
    for p in p_range:
        print('Analyzing ER graphs for p={}'.format(p))
        comp_sizes = []
        for i in range(graphs_per_p):
            G = nx.generators.fast_gnp_random_graph(n, p, seed=i)
            comp_sizes.append(max([len(comp) for comp in nx.connected_components(G)]))
        means.append(np.mean(comp_sizes))
        sds.append(np.std(comp_sizes))

    return np.array(means), np.array(sds)
```

```
In [ ]: n = 5000
p_range = np.logspace(-5, -2, num=20)
means, sds = analyze_er_graph_components(n, p_range)
```

```
Analyzing ER graphs for p=1e-05
Analyzing ER graphs for p=1.438449888287663e-05
Analyzing ER graphs for p=2.06913808111479e-05
Analyzing ER graphs for p=2.9763514416313192e-05
Analyzing ER graphs for p=4.281332398719396e-05
Analyzing ER graphs for p=6.158482110660267e-05
Analyzing ER graphs for p=8.858667904100833e-05
Analyzing ER graphs for p=0.00012742749857031334
Analyzing ER graphs for p=0.00018329807108324357
Analyzing ER graphs for p=0.00026366508987303583
Analyzing ER graphs for p=0.000379269019073225
Analyzing ER graphs for p=0.0005455594781168515
```

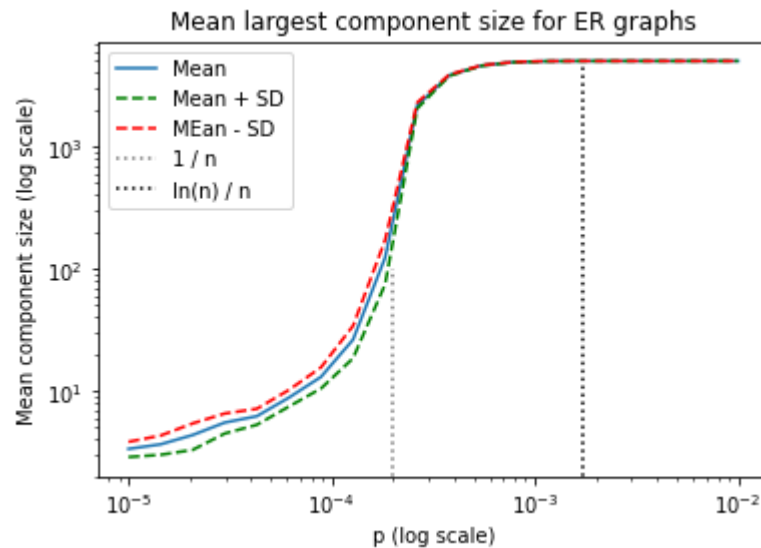
Analyzing ER graphs for $p=0.0007847599703514606$
 Analyzing ER graphs for $p=0.0011288378916846883$
 Analyzing ER graphs for $p=0.001623776739188721$
 Analyzing ER graphs for $p=0.002335721469090121$
 Analyzing ER graphs for $p=0.003359818286283781$
 Analyzing ER graphs for $p=0.004832930238571752$
 Analyzing ER graphs for $p=0.0069519279617756054$
 Analyzing ER graphs for $p=0.01$

```
In [ ]: def plot_with_sds(x, y, y_sds, n=5_000):
    plt.plot(x, y)
    plt.plot(x, y - y_sds, c='green', linestyle='dashed')
    plt.plot(x, y + y_sds, c='red', linestyle='dashed')

    # plot vertical lines at 1 / n and ln(n) / n
    plt.vlines(1 / n, 0, 100, linestyle='dotted', colors='gray')
    plt.vlines(np.log(n) / n, 0, 4900, linestyle='dotted', colors='black')

    plt.yscale('log')
    plt.xscale('log')
    plt.title('Mean largest component size for ER graphs')
    plt.xlabel('p (log scale)')
    plt.ylabel('Mean component size (log scale)')
    plt.legend(['Mean', 'Mean + SD', 'Mean - SD', '1 / n', 'ln(n) / n'])
```

```
In [ ]: plot_with_sds(p_range, means, sds)
```



b) Interesting values of p

From class, we have that when $p = \lambda / n$, $\lambda > 1$ makes the largest component size converge to n , while $\lambda < 1$ makes it converge to $\ln(n)$. Here, $\lambda = 1$ represents a point at which the largest component size is transitioning rapidly, as it has the graph's peak slope. Also, I notice that as λ shrinks less than 1, the size remains low (near $\ln(n) \approx 8$), while as λ grows above 1, the size approaches n .

Also, we have that when $p = \lambda \ln(n) / n$, $\lambda < 1$ makes the graph's probability of being connected converge to 0, while $\lambda > 1$ makes it converge to 1. Here, to the left of $\ln(n) / n$ (i.e. $\lambda > 1$), the mean component size has not yet converged to n , but to the right, the mean is equal to n and the standard deviation is 0, meaning that, as expected, all sampled graphs converge to being connected.

2. Fitting power-law exponents

b) Pareto sampling

```
In [ ]: def sample_pareto(alpha=1, x_0=1, size=100_000):
    pareto_func = np.vectorize(lambda u: x_0 / u**(1 / alpha))
    u_samples = np.random.sample(size=size)
    return pareto_func(u_samples)

def compute_pdf(samples):
```

```

samples_ints = np.round(samples, decimals=0)
samples_int_counts = Counter(samples_ints)
x = np.array(sorted(samples_int_counts.keys()))
px_x = np.array([samples_int_counts[k] for k in x])
px_x = px_x / px_x.sum()
return x, px_x

def compute_ccdf(x, px_x):
    s = 0
    fbarx_x = np.zeros(px_x.shape)
    for i in range(len(x) - 1, -1, -1):
        s += px_x[i]
        fbarx_x[i] = s

    return x, fbarx_x

def plot_observations(x, y, alpha=1):
    plt.scatter(x, y, s=5)
    plt.xscale('log')
    plt.yscale('log')
    y_expected = 1 / x ** (alpha + 1)
    plt.plot(x[y_expected >= 1e-5], y_expected[y_expected >= 1e-5], c='green')
    plt.title('Sampled and actual PDFs for Pareto distribution')
    plt.legend(['Sampled', 'Actual'])
    plt.xlabel('x (log scale)')
    plt.ylabel('p(X = x) (log scale)')

```

```

In [ ]:
pareto = sample_pareto()
x, px_x = compute_pdf(pareto)

```

```

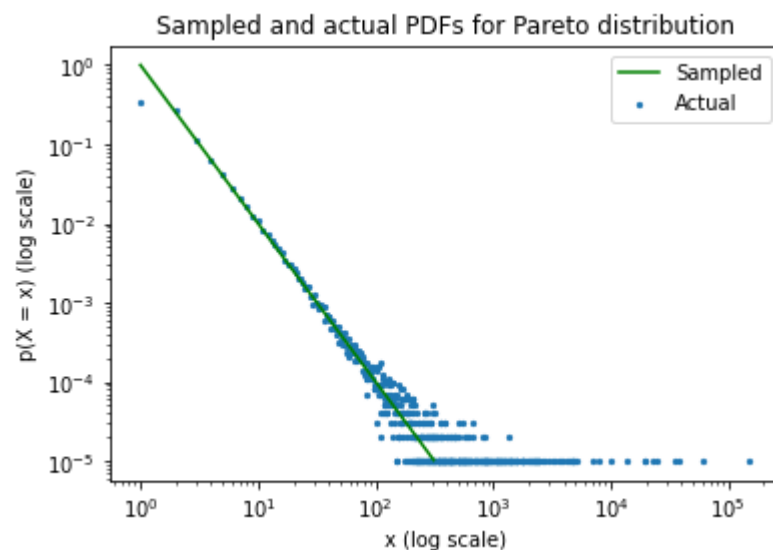
In [ ]:
x, fbarx_x = compute_ccdf(x, px_x)

```

```

In [ ]:
plot_observations(x, px_x)

```



d) Estimating alpha

In []:

```
def estimate_alpha_pdf_ls(samples):
    x, p_x = compute_pdf(samples)

    x_log = np.log(x.reshape((x.shape[0], 1)))
    y_log = np.log(p_x)

    alpha_pred = -np.linalg.lstsq(x_log, y_log, rcond=None)[0] - 1
    return alpha_pred

def estimate_alpha_ccdf_ls(samples):
    x, px_x = compute_pdf(samples)
    x, fbarx_x = compute_ccdf(x, px_x)

    x_log = np.log(x.reshape((x.shape[0], 1)))
    y_log = np.log(fbarx_x)

    # integrating the PDF gives line with slope -alpha
    alpha_pred = -np.linalg.lstsq(x_log, y_log, rcond=None)[0]

    return alpha_pred

def estimate_alpha_mle(samples, x_0=1):
    n = len(samples)
```



```

# from c), MLE is  $\alpha = n / (\sum_i (\log(x_i)) - (n * \log(x_0)))$ 
return n / (np.log(samples).sum() - (n * np.log(x_0)))

def eval_alpha_estimate(estimate_alpha, name, samples=100):
    preds = []
    for _ in range(samples):
        samples = sample_pareto()
        preds.append(estimate_alpha(samples))
    m, sd = np.mean(preds), np.std(preds)
    print('{}:'.format(name))
    print('Mean: {}, SD: {}'.format(m, sd))

```

```
In [ ]: eval_alpha_estimate(estimate_alpha_pdf_ls, 'Least-squares on PDF');
```

```

Least-squares on PDF:
Mean: 0.7535055686849037, SD: 0.01445416636769601

```

```
In [ ]: eval_alpha_estimate(estimate_alpha_ccdf_ls, 'Least-squares on CCDF');
```

```

Least-squares on CCDF:
Mean: 0.9985720588961311, SD: 0.010272736553571053

```

```
In [ ]: eval_alpha_estimate(estimate_alpha_mle, 'Log-likelihood maximization');
```

```

Log-likelihood maximization:
Mean: 1.0003310052420762, SD: 0.0032189410484160037

```

Least-squares regression on the PDF gives a very poor result (25% too small), while least-squares on the CCDF and the log-likelihood MLE both give results very close to the true value. The MLE is closer to the true value than CCDF least-squares by one order of magnitude, and has a smaller standard deviation, so I'd say that the MLE gives the best estimate.