

Question 1

Consider a small dataset of four points $D = (0, 1), (1, 0), (3, 2), (5, 4)$. Compute the equation of the linear regression, the total explained variance, the total variance, and the coefficient of determination. Please detail your steps.

To get the linear regression equation, we must first determine the slope and y-intercept of the line that best fits the data. To do so, we can utilize the following formulas:

$$\text{Slope } (m) = \frac{N \sum(xy) - \sum x \sum y}{N \sum(x^2) - (\sum x)^2} \quad (1)$$

$$\text{Y-intercept } (b) = \frac{\sum y - m \sum x}{N} \quad (2)$$

Substituting these values:

$$N = 4 \quad (3)$$

$$\sum x = 9 \quad (4)$$

$$\sum y = 7 \quad (5)$$

$$\sum(xy) = 26 \quad (6)$$

$$\sum(x^2) = 35 \quad (7)$$

We get Slope $(m) = 0.69$ and Y-intercept $(b) = 0.1975$

This gives the Linear Regression equation:

$$y = 0.69 * x + 0.1975 \quad (8)$$

To calculate the total explained variance, we can use:

$$\text{Explained Variance} = \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (9)$$

y_i is the actual value of the i-th data point, and \hat{y}_i is the predicted value of the i-th data point. For our dataset, we get:

$$\text{Explained Variance} = (1 - 0.1975)^2 + (0 - 0.8875)^2 + (2 - 2.2675)^2 + (4 - 3.6475)^2 \quad (10)$$

$$\text{Explained Variance} = 1.627475 \quad (11)$$

To calculate the total variance, we can use:

$$Total\ Variance = \sum_{i=1}^N (y_i - \bar{y})^2 \quad (12)$$

\bar{y} is the mean. that is 1.75

$$Total\ Variance = (1 - 1.75)^2 + (0 - 1.75)^2 + (2 - 1.75)^2 + (4 - 1.75)^2 \quad (13)$$

$$Total\ Variance = 8.75 \quad (14)$$

Coefficient of determination or R^2 , is given by:

$$R^2 = 1 - \frac{Explained\ Variance}{Total\ Variance} \quad (15)$$

For our dataset, we get:

$$R^2 = 1 - \frac{1.627475}{8.75} = 0.814 \quad (16)$$

Question 2

Please study the GHS Urban Centre Database from https://ghsl.jrc.ec.europa.eu/ghs_stat_ucdb2015mt_r2019a.p and write a 200-word summary of what the data represents. What type of questions can we answer with this data?

The GHS Urban Centre Database (GHS-UCDB) is a collection of information about cities all over the world. This data is based on different things like how many people live there and how much of the area is built up. The data is gathered from lots of sources and put together with GHSL data. The database tells us about lots of things in the cities like how big they are, where they are, what the weather is like, how is the economy, and how much pollution there is. The answer to these questions is provided using standard metrics that are mentioned in the data dictionary. The cities are measured using a 1x1 km grid, and they're classified as urban centers based on how many people live there and how much of the area is built up. The GHSL gives the input data, and the "degree of urbanization" (DEGURBA) is used to set the rules. Overall, the GHS-UCDB gives us lots of information about cities all over the world in a form that is conducive to analysis.

1. Data can assist service companies in determining the optimal boundaries for their service areas. For instance, delivery service aggregators like DoorDash heavily rely on urban areas, and by using the data, they can identify the extent of the urban centers and set their service boundaries accordingly.

2. Data can be used to find out correlations between interesting features like: a) Night time light emission and GDP b) Greenness and Flooding c) How Geography affects built up area, Economy and pollution.

3. Based on the multiple features available we can try to predict which urban centers are likely to experience or have the potential for rapid growth in the future based on current data and previous year data.

Question 3

From https://ghsl.jrc.ec.europa.eu/ghs-stat-ucdb2015mt_r2019a.php, download the associated csv file and perform a scaling analysis between the total population in 2015 and the urban center extension for the following nations: USA, Canada, Poland, and Russia. In your analysis, please exclude any city with less than 100 km² in area. Does the analysis match your expectations? Please explain your findings.

Using the code as displayed below:

```
1 import pandas as pd
2 import numpy as np
3
4
5 data = pd.read_csv('GHS_STAT_UCDB2015MT_GLOBE_R2019A_V1_2.csv', encoding='ISO-8859-1',
6                   low_memory=False)
7 data1 = data[['CTR_MN_NM', 'UC_NM_MN', 'P15', 'BBX_LATMN', 'BBX_LATMX', 'BBX_LONMN', 'BBX_LONMX',
8               'AREA']]
9 countries = ['United States', 'Canada', 'Poland', 'Russia']
10 data1 = data1[data1['CTR_MN_NM'].isin(countries)]
11 data1 = data1[data1['AREA']>=100]
12
13 data1 = data1[['P15', 'AREA', 'CTR_MN_NM']].dropna()
14
15 import numpy as np
16 from sklearn.linear_model import LinearRegression
17 import matplotlib.pyplot as plt
18
19 for country in countries:
20     data_country = data1[data1['CTR_MN_NM'] == country]
21     data_log = data_country[['P15', 'AREA']].dropna().apply(np.log)
22     X = data_log['P15'].values.reshape(-1, 1)
23     y = data_log['AREA']
24     model = LinearRegression().fit(X, y)
25     print("\n\n", country)
26     print('Intercept:', model.intercept_)
27     print('Coefficient:', model.coef_[0])
28     print('R^2:', model.score(X, y))
29     if model.coef_[0] == 1:
30         print('Linear relationship')
31     elif model.coef_[0] < 1:
32         print('Sublinear relationship')
33     else:
34         print('Superlinear relationship')
35     plt.figure(figsize=(10, 10))
36     plt.scatter(X, y, alpha=0.75, color='darkblue')
37     plt.plot(X, model.predict(X), color='red')
38     plt.xlabel('log(Population in 2015)', fontsize=14)
39     plt.ylabel('log(AREA)', fontsize=14)
40     plt.title('Relationship between log(Population in 2015) and log(AREA) in ' + country,
41              fontsize=16)
42     plt.tick_params(axis='both', labelsize=12)
```

```

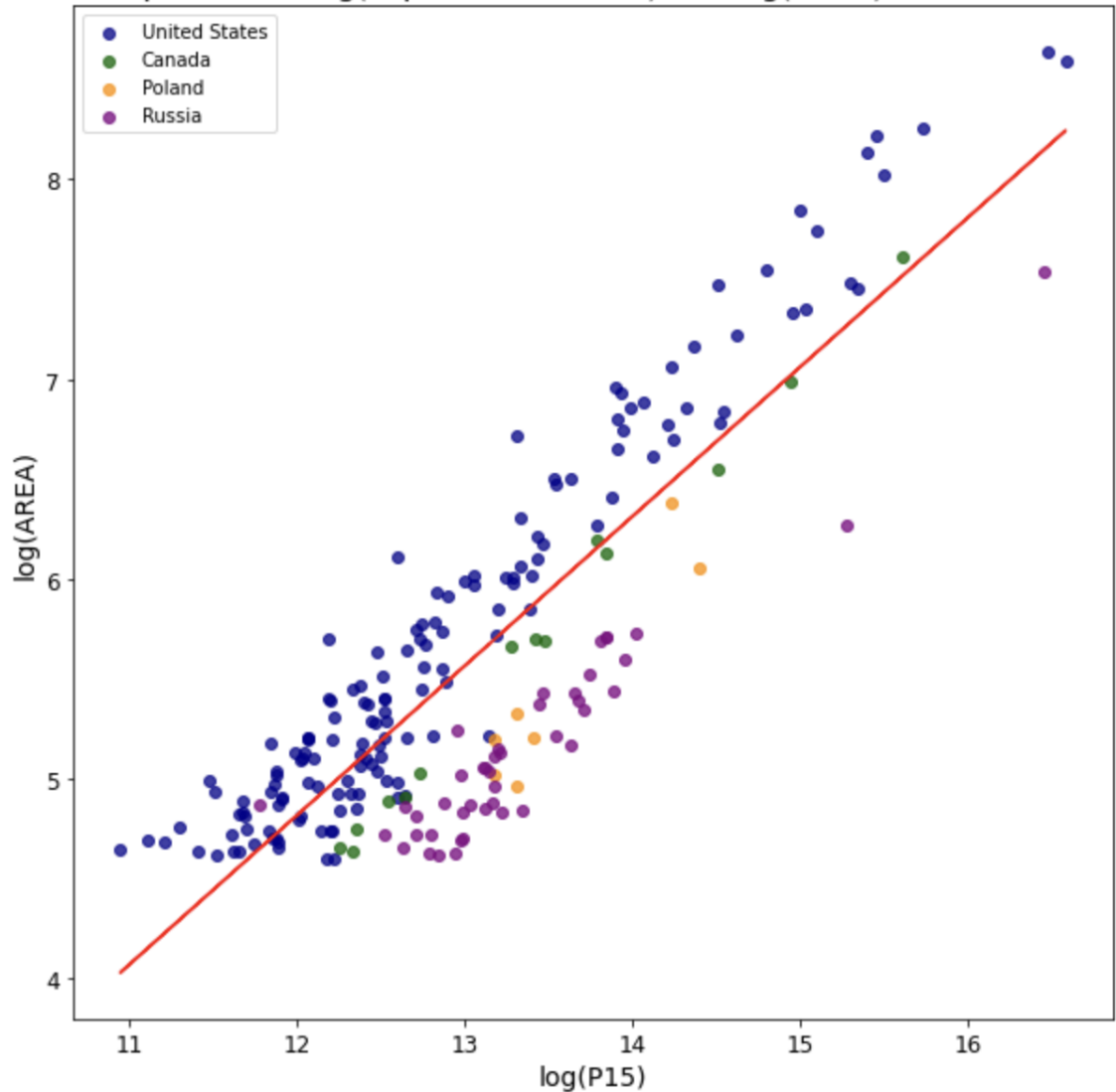
40     plt.show()
41
42
43 data_log = data1[['P15', 'AREA', 'CTR_MN_NM']].dropna()
44 data_log['P15'] = np.log(data_log['P15'])
45 data_log['AREA'] = np.log(data_log['AREA'])
46
47
48 X = data_log['P15'].values.reshape(-1, 1)
49 y = data_log['AREA']
50 model = LinearRegression().fit(X, y)
51 print("\n\nAll 4 Countries")
52 print('Intercept:', model.intercept_)
53 print('Coefficient:', model.coef_[0])
54 print('R^2:', model.score(X, y))
55 if model.coef_[0] == 1:
56     print('Linear relationship')
57 elif model.coef_[0] < 1:
58     print('Sublinear relationship')
59 else:
60     print('Superlinear relationship')
61
62 countries_colors = {'United States': 'darkblue', 'Canada': 'darkgreen', 'Poland': '
    darkorange', 'Russia': 'purple'}
63 data_log['color'] = data_log['CTR_MN_NM'].map(countries_colors)
64 fig, ax = plt.subplots(figsize=(10, 10))
65 for country, color in countries_colors.items():
66     data_country = data_log[data_log['CTR_MN_NM'] == country]
67     plt.scatter(data_country['P15'], data_country['AREA'], alpha=0.75, color=color, label
        =country)
68 plt.plot(X, model.predict(X), color='red')
69 plt.xlabel('log(P15)', fontsize=14)
70 plt.ylabel('log(AREA)', fontsize=14)
71 plt.title('Relationship between log(Population in 2015) and log(AREA) in different
    countries', fontsize=16)
72 plt.tick_params(axis='both', labelsize=12)
73 plt.legend()
74 plt.show()

```

We get the following outputs as listed. Let's see the final output with all 4 countries followed by each individual country.

```
All 4 Countries
Intercept: -4.143808264706207
Coefficient: 0.7468159086791123
R^2: 0.7586324868929072
Sublinear relationship
```

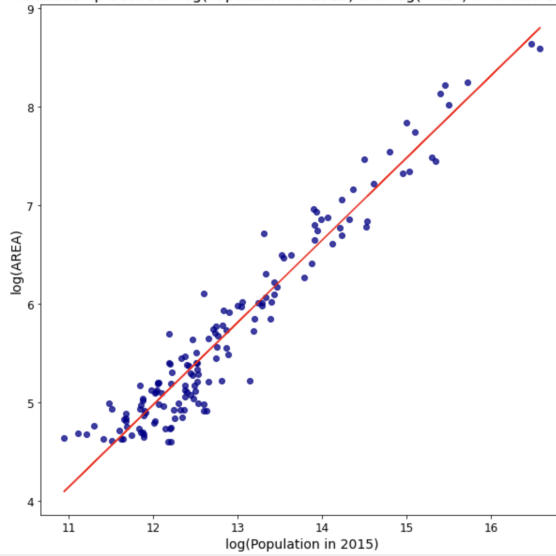
Relationship between $\log(\text{Population in 2015})$ and $\log(\text{AREA})$ in different countries



Here We can notice a sub-linear relationship as backed by theory. Following this, we have interesting results in other countries where the USA and CANADA are similar with the scaling coefficient in the range of 0.8 - 0.9 while Russia has a much lower scaling coefficient of 0.66. Indicating that cities in Russia have a much lesser increase in area as the population increases. Meanwhile, in Poland we have 0.988 indicating an almost linear increase in the area of the city along with an increase in population. The graphs for individual countries in the following page.

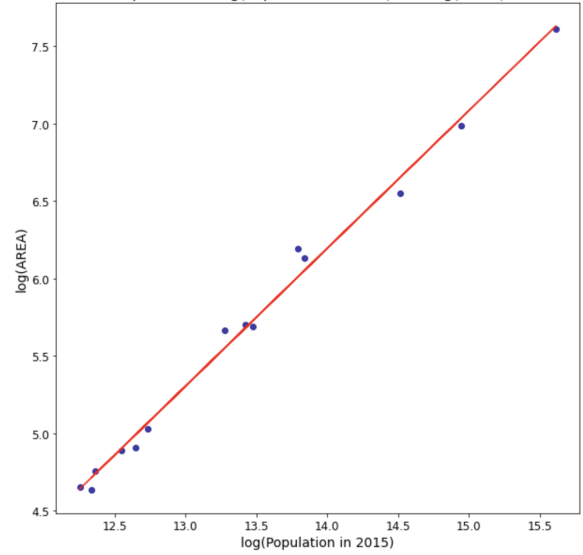
United States
Intercept: -5.036636236695314
Coefficient: 0.8342868837859487
R²: 0.9313119228493343
Sublinear relationship

Relationship between log(Population in 2015) and log(AREA) in United States



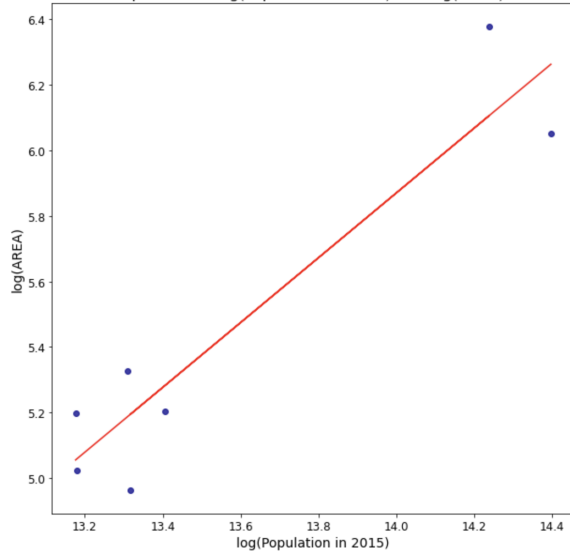
Canada
Intercept: -6.25676946786586
Coefficient: 0.889359630243167
R²: 0.9923562716968196
Sublinear relationship

Relationship between log(Population in 2015) and log(AREA) in Canada



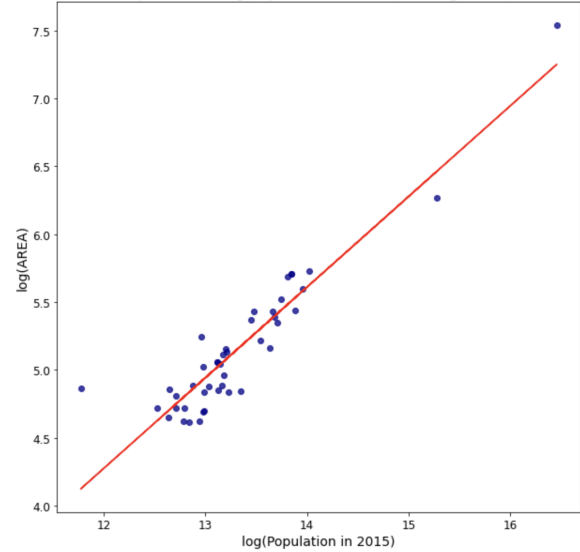
Poland
Intercept: -7.974310871947761
Coefficient: 0.988887124217453
R²: 0.8765464831894277
Sublinear relationship

Relationship between log(Population in 2015) and log(AREA) in Poland



Russia
Intercept: -3.7237621335524933
Coefficient: 0.6667213924716087
R²: 0.8631859543474508
Sublinear relationship

Relationship between log(Population in 2015) and log(AREA) in Russia



Question 4

Consider the amorphous settlement model in Part 1 of Chapter 3. Assuming that the distance traveled by a person per day is on average 3km and that the physical distance used for defining an interaction is about 1 m, what is the expected total number of interactions per capita in Brooklyn (pull data from Wikipedia on population and area)? Does the number make sense to you based on your daily experience?

According to the amorphous model we have

$$k = \frac{N - 1}{A} * a_0 * l_0 \quad (17)$$

$$a_0 = 1m \quad l_0 = 3km$$

From <https://en.wikipedia.org/wiki/Brooklyn> we have $N = 2736074$ and $A = 183.4km^2$ considering only land area.

$$k = \frac{273074 - 1}{183.4} * \frac{1}{1000} * 3 \quad (18)$$

$$k = 4.46 \quad (19)$$

From my daily experience, I think this number is more or less correct on average. There are days when I have interactions with 0 people and days when I have interactions with 10 people but more or less can balance out.

Question 5

Consider again the amorphous settlement model in Part 1 of Chapter 3. Assume that the cost of transport is independent of any distance traveled in an imaginary city in which public transport is free for anyone and the only means of transport. How would the area of the city change as a function of the population? Please explain your results.

If the cost of transport is independent of any distance traveled, it means that the cost of transport per capita, C_t , is constant, regardless of the radius of the city. In other words, C_t does not depend on R , the radius of the city.

Given that

$$N * C_t = Y \quad (20)$$

according to the amorphous model, where Y is the total income of the population, and C_t is constant, we have

$$N \propto Y \quad (21)$$

and

$$Y \propto \frac{N^2}{A} \quad (22)$$

This gives

$$A \propto N \quad (23)$$

Therefore, the area of the city is proportional to the population raised to the power of 1, which means that A is directly proportional to N . This is different from the previous case where the cost of transport per capita depended on the radius of the city and $A \propto N^{(\frac{2}{3})}$.

Question 6

Consider the classical (compartmental) SIR model discussed in class for the following parameters $\beta = 0.3$, $\gamma = 0.1$, and $N = 10,000$. Simulate the model and plot the three variables from $t = 0$ to 200 for the case in which there is only one infected subject at the beginning of the infection and nobody is immune to the disease.

Using these system of equations, we can simulate the SIR model.

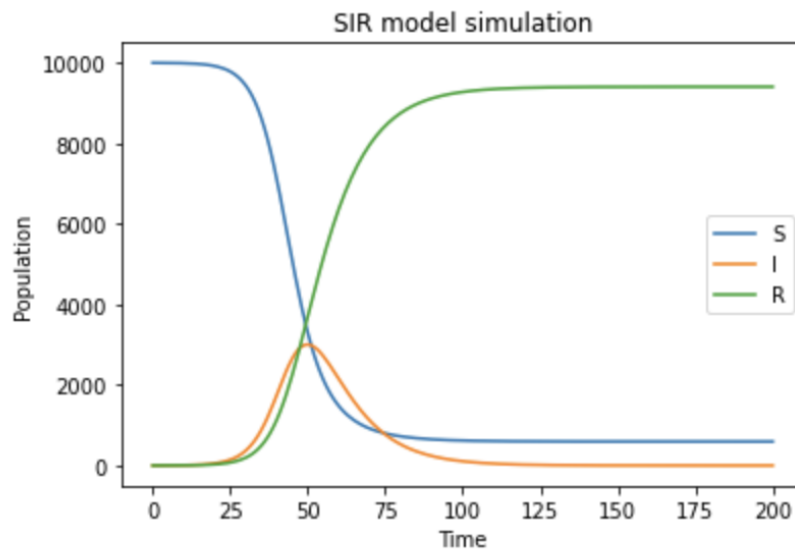
$$\begin{aligned}\frac{dS}{dt} &= -\frac{\beta}{N}IS \\ \frac{dI}{dt} &= \frac{\beta}{N}IS - \gamma I \\ \frac{dR}{dt} &= \gamma I\end{aligned}$$

```
1 import numpy as np
2 from scipy.integrate import odeint
3 import matplotlib.pyplot as plt
4
5 #Initial
6 S0 = 9999
7 I0 = 1
8 R0 = 0
9
10
11 beta = 0.3
12 gamma = 0.1
13 N = 10000
14
15
16 t = np.linspace(0, 200, 201)
17
18
19 def sir(y, t, beta, gamma, N):
20     S, I, R = y
21     dSdt = (-beta/N) * I * S
22     dIdt = (beta/N) * I * S - gamma * I
23     dRdt = gamma * I
24     return [dSdt, dIdt, dRdt]
25
26
27 sol = odeint(sir, [S0, I0, R0], t, args=(beta, gamma, N))
28
29
30 plt.plot(t, sol[:, 0], label='S')
31 plt.plot(t, sol[:, 1], label='I')
32 plt.plot(t, sol[:, 2], label='R')
33 plt.legend()
34 plt.xlabel('Time')
35 plt.ylabel('Population')
36 plt.title('SIR model simulation')
37 plt.show()
```

```

38
39
40 print("At t=0:")
41 print("S = ", sol[0, 0])
42 print("I = ", sol[0, 1])
43 print("R = ", sol[0, 2])
44 print("")
45
46 print("At t=200:")
47 print("S = ", sol[-1, 0])
48 print("I = ", sol[-1, 1])
49 print("R = ", sol[-1, 2])
50 print("")
51
52
53 Imax = np.max(sol[:, 1])
54 print("Maximum value of I = ", Imax)

```



```

At t=0:
S = 9999.0
I = 1.0
R = 0.0

At t=200:
S = 595.1358608053775
I = 0.02864570401681712
R = 9404.835493490615

Maximum value of I = 3004.8159818370027

```

Question 7

For the same SIR considered above. Compute the reproductive number and the effective reproductive number. Would you expect the disease to be epidemic? If so, does this match the numerical results? What do you think is a disease with a similar behavior? Please compute analytically the peak number of infections and compare with simulations. Also, please compute the steady-state number of susceptible individuals by solving the associated implicit equation and compare with numerical results.

$$R_o = \frac{\beta}{\gamma} = \frac{0.3}{0.1} = 3 \quad (24)$$

$$R_e = \frac{S_o}{N} * \frac{\beta}{\gamma} = \frac{9999}{10000} * \frac{0.3}{0.1} = 2.9997 \quad (25)$$

Since R_e is very close to but slightly below 3, we would expect the disease to be epidemic. This matches the numerical results obtained by simulating the SIR model.

A disease with similar behavior to the SIR model with $R_0 = 3$ and $R_e = 2.997$ could be COVID-19 as inferred from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7657547/>.

Computing I_{max} analytically:

We have

$$\frac{I_{max}}{N} = 1 - \frac{1}{R_o} * (1 + \ln(R_o)) \quad (26)$$

$$R_o = \frac{\beta}{\gamma} = \frac{0.3}{0.1} = 3 \quad (27)$$

$$I_{max} = 3004.6 \quad (28)$$

This number matches our graph

computing the steady-state number of susceptible individuals using the equation:

$$\log\left(\frac{S_{\infty}}{N}\right) = R_0 \cdot \left(\frac{S_{\infty}}{N} - 1\right) \quad (29)$$

```
1 import math
2 from scipy.optimize import root_scalar
3
4 def solve_for_S_inf(R0, N):
5     def equation(S_inf):
6         return math.log(S_inf/N) - R0 * ((S_inf-N)/N)
7
8     sol = root_scalar(equation, bracket=[1, N/10])
9     S_inf = sol.root
10
11     return S_inf
12
13 R0 = 3
14 N = 10000
15 S_inf = solve_for_S_inf(R0, N)
16 print("S_inf =", S_inf)
```

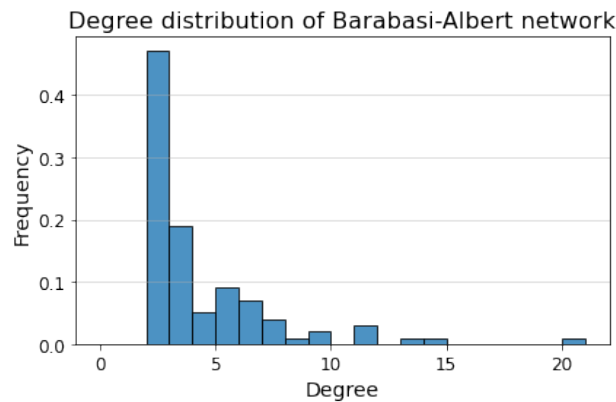
This gives a value close to the $S_{t=200}$ value we got from the simulation

$$S_{inf} = 595.20$$

Question 8

Create any scale-free network of your choice with at least 100 nodes. Plot the degree distribution and compute the mean degree $\langle k \rangle$ and the second moment of the degree $\langle k^2 \rangle$

```
1 import networkx as nx
2 import matplotlib.pyplot as plt
3
4 ba = nx.barabasi_albert_graph(100, 2)
5
6 degrees = [ba.degree(node) for node in ba.nodes()]
7 plt.hist(degrees, bins=range(max(degrees)+2), density=True, alpha=0.8, edgecolor='black')
8 plt.xlabel('Degree', fontsize=14)
9 plt.ylabel('Frequency', fontsize=14)
10 plt.title('Degree distribution of Barabasi-Albert network', fontsize=16)
11 plt.tick_params(axis='both', labelsize=12)
12 plt.grid(axis='y', alpha=0.5)
13 plt.tight_layout()
14 plt.show()
```



```
1 import numpy as np
2
3 mean_degree = np.mean(degrees)
4 print(f"Mean degree: {mean_degree:.2f}")
5
6 second_moment = np.mean(np.power(degrees, 2))
7 print(f"Second moment of degree: {second_moment:.2f}")
```

Giving an output of

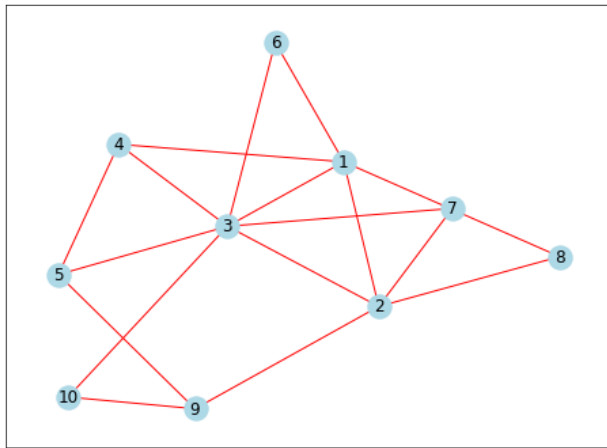
Mean degree: 3.92 Second moment of degree: 24.66

Question 9

Consider the network depicted in Fig. 1, compute the searchability of the path $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$.

Given the graph, we can use the following formula to calculate the searchability of the path. Where x excludes the first and last node in the path.

```
1 #Plotting the network
2 import networkx as nx
3 import matplotlib.pyplot as plt
4
5 G = nx.Graph()
6 G.add_nodes_from([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
7 G.add_edges_from([(1, 2), (1, 3), (1, 4), (1, 6), (1, 8), (2, 3), (2, 7), (2, 8), (2, 9),
8                   (3, 4), (3, 5), (3, 6), (3, 7), (3, 10), (4, 5), (5, 9), (9, 10)])
9
10 fig, ax = plt.subplots(figsize=(8, 6))
11 pos = nx.spring_layout(G)
12 nx.draw_networkx(G, pos=pos, with_labels=True, node_color='lightblue', edge_color='red')
13 plt.show()
```



$$H_x = \log(K_s) + \sum_{j \in \hat{x}} \log(K_j - 1)$$

Based on this we have

$$H_x = \log(5) + \log(5 - 1) + \log(7 - 1)$$

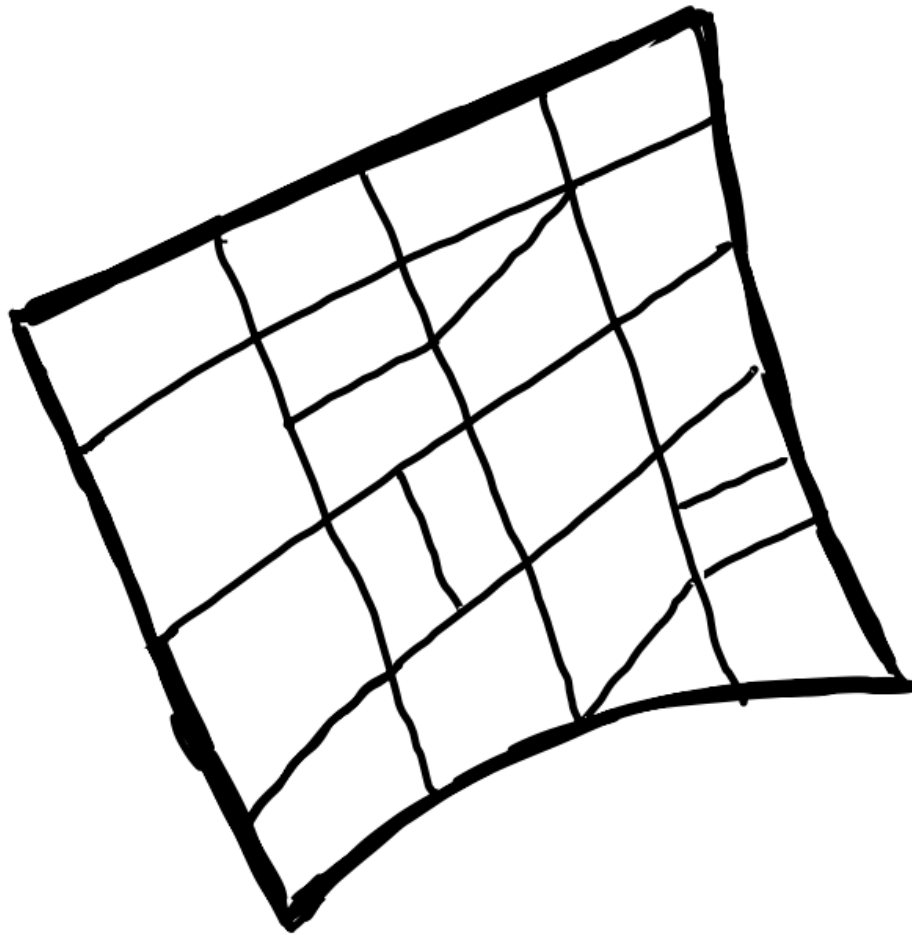
So,

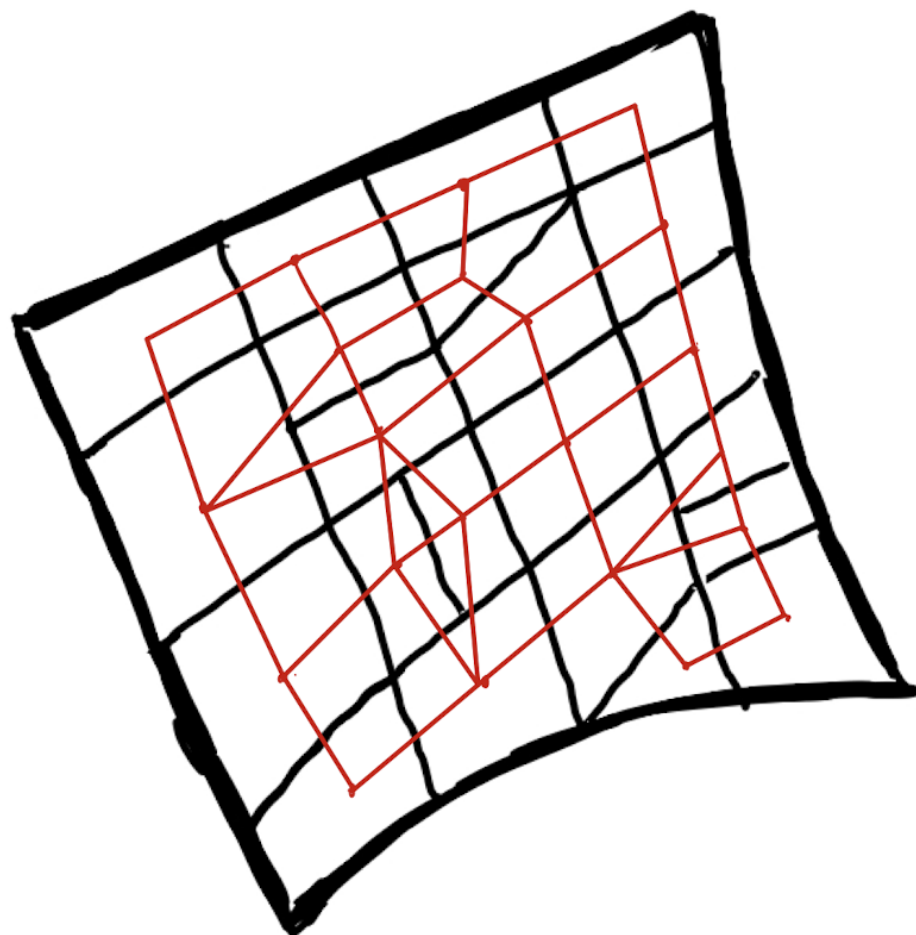
$$H_x = 4.78$$

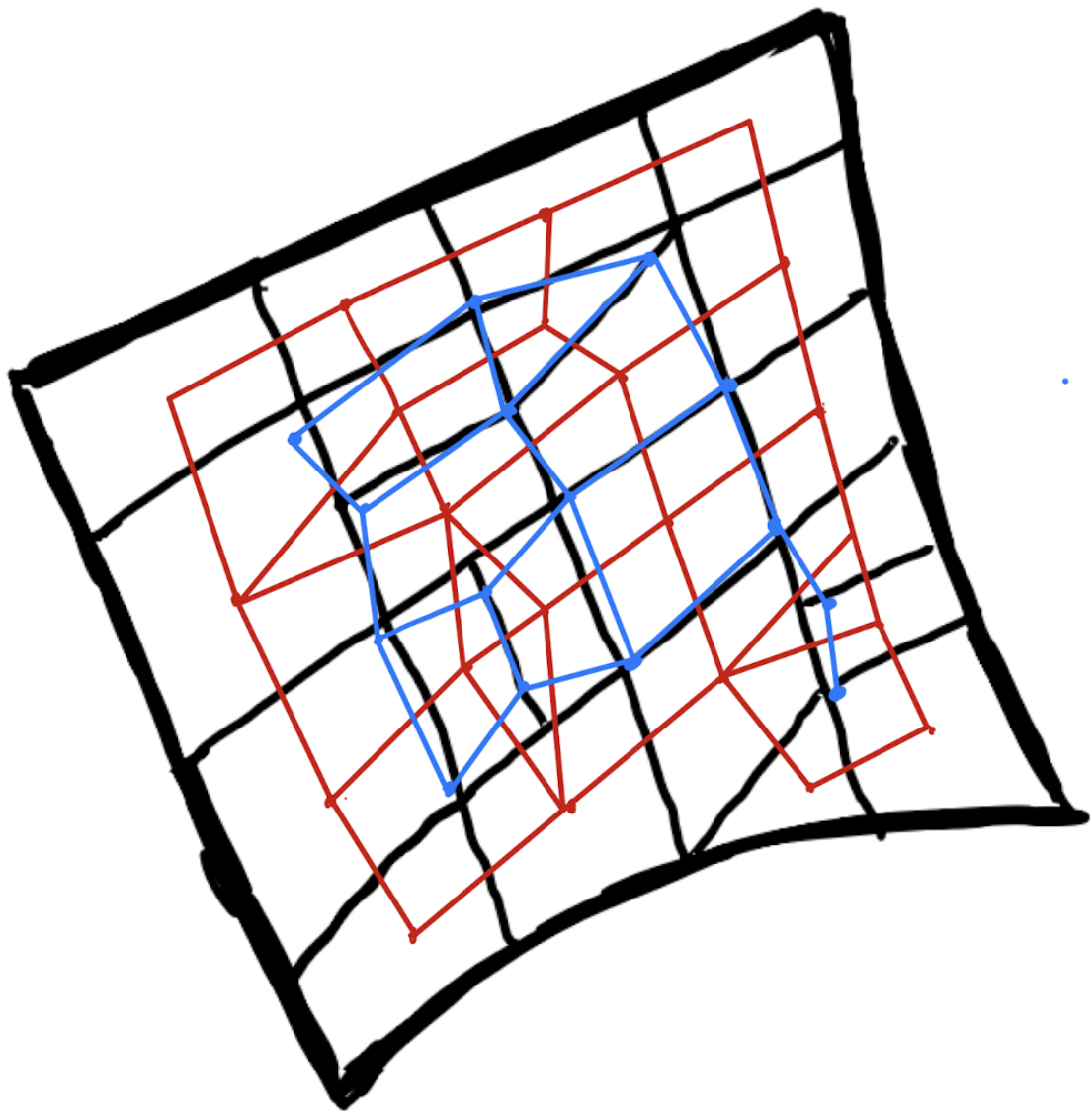
Question 10

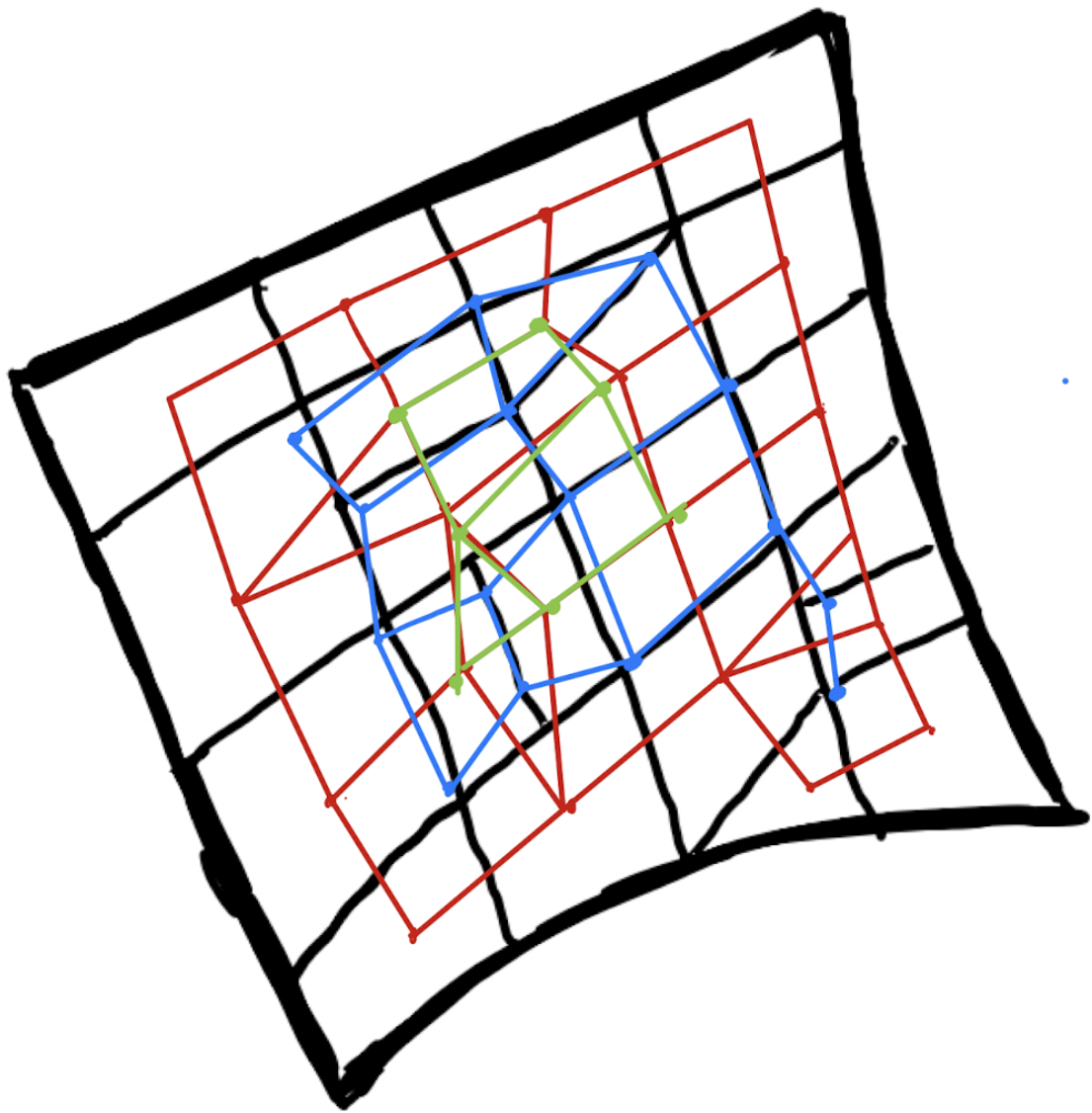
Consider the schematic of the block depicted in Fig. 2, what is the block complexity k_{max} ?

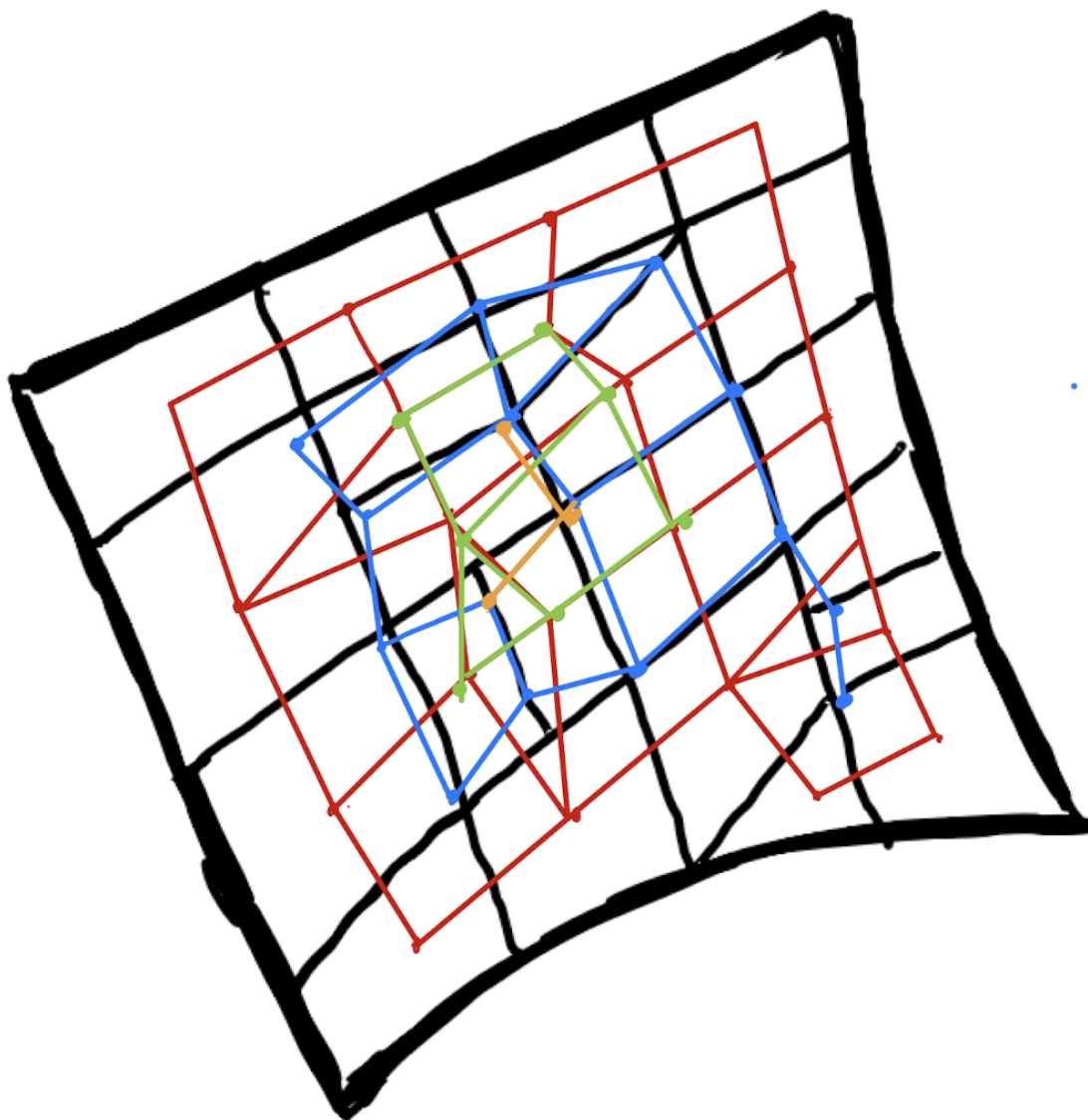
The following images show each iteration of connecting midpoints of adjacent blocks until a dot or line is obtained.











As we went through 4 iterations, the block complexity k_max is 4.