

**MODUL PRAKTIKUM**  
**MATA KULIAH: DATA MINING**  
**MODUL VI**  
**[Data Graf I]**



Program Studi Sains Data  
Fakultas Sains  
**INSTITUT TEKNOLOGI SUMATERA**  
Tahun Ajaran Ganjil 2024/2025.

## 1. Tujuan Praktikum

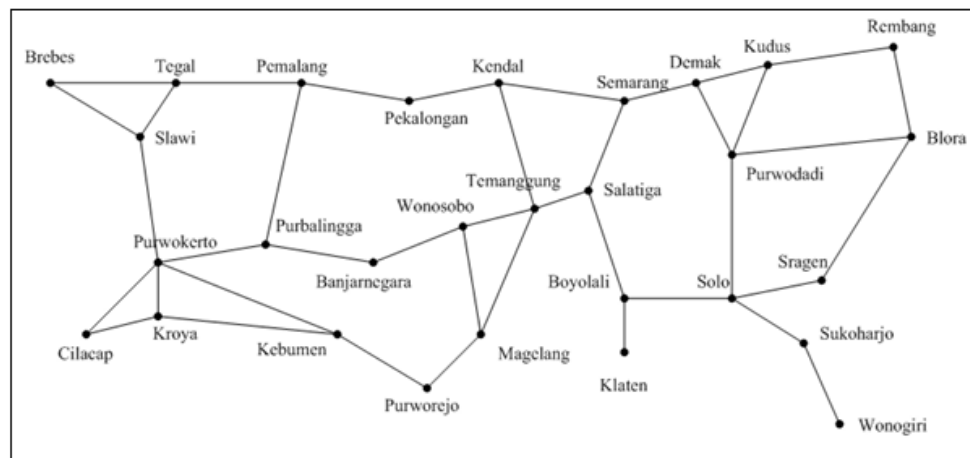
- Memahami dan menerapkan konsep-konsep dasar graf.
- Mengembangkan keterampilan analisis dan pemecahan masalah yang melibatkan struktur data graf.

## 2. Tinjauan Pustaka

- a. Konsep Graf  
i. Graf

Graf digunakan untuk merepresentasikan objek-objek diskrit dan hubungan antara objek-objek tersebut. Representasi visual dari graf adalah dengan menyatakan objek sebagai noktah, bulatan, atau titik, sedangkan hubungan antara objek dinyatakan dengan garis.

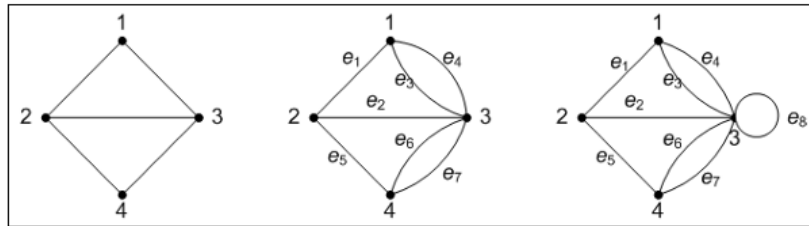
Definisi Graf adalah representasi visual atau struktural dari data di mana entitas (node atau simpul) dan hubungan antar entitas (edges atau tepi) dimodelkan untuk mencerminkan bagaimana titik-titik data saling terhubung. Graf data sangat berguna untuk memahami hubungan dan pola antar data  $G = (V, E)$  yaitu Graf terdiri dari dua hal yaitu V (vertices/nodes), Himpunan dari simpul atau node dalam graf dan E (edges). Dalam graf tak berarah, sisi adalah pasangan tidak terurut dari simpul. Salah satu contoh graf adalah sebagai berikut:



Gambar 1 contoh graf(Peta Jaringan Jalan Raya)

Diatas sudah di jelaskan bahwa graf dapat ditulis dengan menggunakan notasi  $G = (V, E)$ , dengan  $V$  = himpunan tidak-kosong dari simpul-simpul (vertices) =  $\{v_1, v_2, \dots, v_n\}$ . Sedangkan untuk nilai edges dapat diartikan  $E$  = himpunan sisi (edges) yang menghubungkan sepasang simpul =  $\{e_1, e_2, \dots, e_n\}$ .

Simpul pada graf dapat dinomori dengan huruf seperti v, w, ..., dengan bilangan asli 1, 2, 3, ... atau gabungan keduanya. Sedangkan sisi yang menghubungkan simpul  $v_i$  dengan  $v_j$  dinyatakan dengan pasangan  $(v_i, v_j)$  atau dengan lambang  $e_1, e_2, \dots$ . Dengan kata lain, jika  $e$  adalah sisi yang menghubungkan simpul  $v_i$  dengan  $v_j$ , maka  $e$  dapat ditulis sebagai  $e = (v_i, v_j)$ . Gambar dibawah ini menunjukkan beberapa contoh graf.



Gambar 2 contoh dari graf

b. *Representasi Graf*

graf harus direpresentasikan di dalam memori. Terdapat beberapa representasi yang mungkin untuk graf. Di sini hanya diberikan tiga macam representasi yang sering digunakan pada graf sederhana, yaitu matriks ketetanggaan, matrik bersisian, dan senarai ketetanggaan.

i. *Matriks Ketetanggaan (Adjacency Matrix)*

graf dapat menggunakan matriks, Salah satunya adalah matriks ketetanggaan. Misalkan  $G = (V, E)$  adalah sebuah graf sederhana dimana  $|V| = n, n > 1$ . Misalkan simpul dari  $G$  adalah  $v_1, v_2, \dots, v_n$ . Maka, matriks ketetanggaan  $A$  dari  $G$  adalah  $n \times n$  matriks dimana:

$$A(i, j) = \begin{cases} 1 & \text{if } v_i \text{ is adjacent to } v_j \\ 0 & \text{otherwise} \end{cases}$$

Matriks ketetanggaan dari sebuah graf mengacu pada keterurutan dari simpul. Sehingga, ada sebanyak  $n!$  keterurutan yang berbeda yang terbentuk dari  $n$  simpul.

Jika grafiknya berarah, maka matriks ketetanggaan  $A$  tidak simetris, karena  $(v_i, v_j) \in E$  jelas tidak menyiratkan bahwa  $(v_j, v_i) \in E$ .

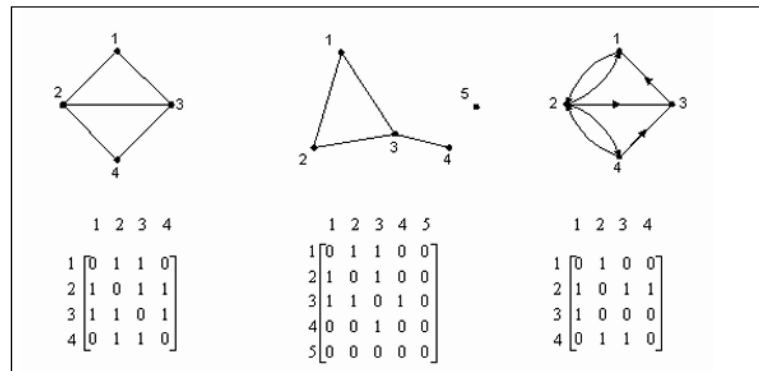
Jika grafiknya terbobot, maka kita memperoleh matriks ketetanggaan terbobot  $n \times n, A$ , yang didefinisikan sebagai:

$$A(i, j) = \begin{cases} w_{ij} & \text{if } v_i \text{ is adjacent to } v_j \\ 0 & \text{otherwise} \end{cases}$$

di mana  $w_{ij}$  adalah bobot pada sisi  $(v_i, v_j) \in E$ . Matriks ketetanggaan yang diberi bobot selalu dapat diubah menjadi matriks biner, jika diinginkan, dengan menggunakan beberapa ambang batas (*threshold*)  $\tau$  pada bobot sisi.

$$A(i, j) = \begin{cases} 1 & \text{if } w_{ij} \geq \tau \\ 0 & \text{otherwise} \end{cases}$$

Berikut contoh graf dengan matriks ketegangan/ *Adjacency Matrix*



Gambar 3 graf dengan Adjacency Matrix

Matriks ketetanggaan merupakan graf sederhana yang simetris, yaitu  $a_{ij} = a_{ji}$ . Hal ini disebabkan oleh kedua-duanya adalah 1 ketika  $v_i$  dan  $v_j$  mempunyai sisi, dan adalah 0 jika tidak ada sisi diantaranya, dinamakan juga matriks nol-satu, sedangkan untuk graf berarah, matriks ketetanggaannya belum tentu simetri (akan berupa simetri jika berupa graf berarah lengkap). Selain itu, graf sederhana tidak mempunyai gelang, sehingga diagonal utamanya selalu 0 karena  $a_{ii}, i = 1, 2, 3, \dots, n$  adalah 0. Sayangnya, matriks ketetanggaan nol-satu tidak dapat digunakan untuk merepresentasikan graf yang mempunyai sisi ganda. Untuk mengatasinya, maka elemen  $a_{ij}$  pada matriks ketetanggaan sama dengan jumlah yang berasosiasi dengan  $(v_i, v_j)$ . Matriks ketetanggaannya bukan lagi matriks nol-satu. Untuk graf semu, gelang pada simpul  $v_i$  dinyatakan dengan nilai satu pada posisi  $(i, i)$  di matriks ketetanggaannya. Keuntungan representasi dengan matriks adalah, elemen matriksnya dapat diakses langsung melalui indeks. Selain itu, kita dapat mengetahui secara langsung apakah simpul  $i$  dan simpul  $j$  bertetangga.

## ii. Graf dari Data Matrix

Banyak himpunan data yang tidak berbentuk grafik dapat diubah menjadi satu. Misalkan  $D$  adalah himpunan data yang terdiri dari  $n$  titik

$x_i \in \mathbb{R}^d$  dalam ruang berdimensi  $d$ . Kita dapat mendefinisikan grafik berbobot  $G = (V, E)$ , di mana terdapat simpul untuk setiap titik di  $D$ , dan terdapat sisi di antara setiap pasangan titik, dengan bobot:

$$w_{ij} = \text{sim}(\mathbf{x}_i, \mathbf{x}_j) \quad \dots\dots(a)$$

di mana  $\text{sim}(\mathbf{x}_i, \mathbf{x}_j)$  menyatakan kesamaan antara titik  $\mathbf{x}_i$  dan  $\mathbf{x}_j$ . Misalnya, kesamaan dapat didefinisikan sebagai berbanding terbalik dengan jarak Euclidean antara titik-titik melalui transformasi.

$$w_{ij} = \text{sim}(\mathbf{x}_i, \mathbf{x}_j) = \exp \left\{ -\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right\} \quad \dots\dots(b)$$

di mana  $\sigma$  adalah parameter sebaran (setara dengan deviasi standar dalam fungsi kepadatan normal). Transformasi ini membatasi fungsi kesamaan  $\text{sim}()$  agar berada dalam rentang  $[0, 1]$ . Kita kemudian dapat memilih ambang batas  $\tau$  yang sesuai dan mengubah matriks ketetanggaan tertimbang menjadi matriks biner melalui Persamaan (a).

c. Atribut Topologi graf

i. Degree(derajat)

Kita telah mendefinisikan derajat sebuah simpul melalui jumlah tetangganya. Definisi yang lebih umum yang berlaku bahkan ketika grafik diberi bobot adalah sebagai berikut:

$$d_i = \sum_j \mathbf{A}(i, j)$$

Derajat jelas merupakan atribut lokal dari setiap node. Salah satu atribut global yang paling sederhana adalah derajat rata-rata:

$$\mu_d = \frac{\sum_i d_i}{n}$$

Definisi sebelumnya dapat dengan mudah digeneralisasi untuk grafik berarah (berbobot). Misalnya, kita dapat memperoleh derajat dan derajat keluar dengan mengambil penjumlahan dari sisi masuk dan keluar, sebagai berikut.

$$id(v_i) = \sum_j \mathbf{A}(j, i)$$

$$od(v_i) = \sum_j \mathbf{A}(i, j)$$

Rata-rata indeks dan rata-rata derajat keluar dapat diperoleh dengan cara yang sama.

ii. Average Path Length (Panjang Jalur Rata-rata)

Panjang jalur rata-rata, juga disebut panjang jalur karakteristik, dari grafik terhubung diberikan sebagai berikut:

$$\mu_L = \frac{\sum_i \sum_{j>i} d(v_i, v_j)}{\binom{n}{2}} = \frac{2}{n(n-1)} \sum_i \sum_{j>i} d(v_i, v_j)$$

di mana  $n$  adalah jumlah simpul dalam grafik, dan  $d(v_i, v_j)$  adalah jarak antara  $v_i$  dan  $v_j$ . Untuk grafik berarah, rata-rata adalah semua pasangan simpul yang diurutkan:

$$\mu_L = \frac{1}{n(n-1)} \sum_i \sum_j d(v_i, v_j)$$

Untuk grafik terputus, rata-rata hanya diambil dari pasangan simpul yang terhubung.

iii. Eccentricity

Eksentrisitas suatu simpul  $v_i$  adalah jarak maksimum dari  $v_i$  ke simpul lain dalam grafik.

$$e(v_i) = \max_j \{d(v_i, v_j)\}$$

Jika grafiknya terputus, eksentrisitas hanya dihitung pada pasangan titik sudut dengan jarak terhingga, yaitu hanya empat titik sudut yang dihubungkan oleh suatu lintasan.

iv. Radius and Diameter

Jari-jari grafik terhubung, dilambangkan  $r(G)$ , adalah eksentrisitas minimum dari setiap simpul dalam grafik:

$$r(G) = \min_i \{e(v_i)\} = \min_i \left\{ \max_j \{d(v_i, v_j)\} \right\}$$

Diameter, dilambangkan  $d(G)$ , adalah eksentrisitas maksimum dari setiap titik sudut pada grafik.

$$d(G) = \max_i \{e(v_i)\} = \max_{i,j} \{d(v_i, v_j)\}$$

Untuk grafik terputus, diameter adalah e c sentrisitas maksimum atas semua komponen grafik yang terhubung. Diameter grafik  $G$  sensitif terhadap outlier.

### 3. Prosedur Praktikum

#### a. Konsep graf

Untuk memahami graf kita dapat meninjau studi kasus dibawah ini:

Graf memiliki  $|V| = 8$  simpul dan  $|E| = 11$  sisi. Karena  $(v_1, v_5) \in E$ , kita mengatakan bahwa  $v_1$  dan  $v_5$  adalah tetangga. Derajat dari  $v_1$  adalah  $d(v_1) = d_1 = 4$ . Urutan derajat dari graf ini adalah:

$$(4, 4, 4, 3, 2, 2, 2, 1)$$

dan distribusi frekuensi derajatnya adalah:

$$(N_0, N_1, N_2, N_3, N_4) = (0, 1, 3, 1, 3)$$

Kita memiliki  $N_0 = 0$  karena tidak ada simpul yang terisolasi, dan  $N_4 = 3$  karena ada tiga simpul, yaitu  $v_1, v_4$ , dan  $v_5$ , yang memiliki derajat  $k = 4$ ; angka lainnya diperoleh dengan cara yang sama.

Distribusi derajat diberikan sebagai:

$$(f(0), f(1), f(2), f(3), f(4)) = (0, 0.125, 0.375, 0.125, 0.375)$$

Urutan simpul  $(v_3, v_1, v_2, v_5, v_1, v_2, v_6)$  adalah jalur dengan panjang 6 antara  $v_3$  dan  $v_6$ . Kita dapat melihat bahwa simpul  $v_1$  dan  $v_2$  telah dikunjungi lebih dari satu kali.

Source Code:

```
1 import networkx as nx
2 import matplotlib.pyplot as plt
3 from collections import Counter
4 # Membuat graf
5 G = nx.Graph()
6 # Menambahkan simpul dan sisi (edges)
7 edges = [(1, 2), (1, 5), (1, 3), (1, 4), (2, 5), (2, 6),
8          (3, 4), (3, 6), (4, 5), (5, 6), (6, 7)]
9 G.add_edges_from(edges)
10 # Menghitung urutan derajat (degree sequence)
11 degree_sequence = sorted([d for n, d in G.degree()], reverse=True)
12 print("Urutan derajat:", degree_sequence)
13 # Menghitung distribusi frekuensi derajat
14 degree_count = Counter(degree_sequence)
```

```

16 N_k = [degree_count.get(i, 0) for i in range(max_degree + 1)]
17 print("Distribusi frekuensi derajat (N_k):", N_k)
18 # Menghitung distribusi derajat (f(k))
19 n = len(G.nodes())
20 f_k = [N_k[i] / n for i in range(max_degree + 1)]
21 print("Distribusi derajat (f(k)):", f_k)
22 # Menggambar graf
23 nx.draw(G, with_labels=True, node_color='lightblue', node_size=2000,
24         font_size=12, font_color='black', edge_color='gray')
25 plt.title("Graf dari studi kasus")
26 plt.show()

```

### Graphs from Data Matrix

Database yang akan digunakan adalah data iris dapat dilihat dari gambar dibawah ini :

Table 1.1. Extract from the Iris dataset

	Sepal length	Sepal width	Petal length	Petal width	Class
	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
$x_1$	5.9	3.0	4.2	1.5	Iris-versicolor
$x_2$	6.9	3.1	4.9	1.5	Iris-versicolor
$x_3$	6.6	2.9	4.6	1.3	Iris-versicolor
$x_4$	4.6	3.2	1.4	0.2	Iris-setosa
$x_5$	6.0	2.2	4.0	1.0	Iris-versicolor
$x_6$	4.7	3.2	1.3	0.2	Iris-setosa
$x_7$	6.5	3.0	5.8	2.2	Iris-virginica
$x_8$	5.8	2.7	5.1	1.9	Iris-virginica
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_{149}$	7.7	3.8	6.7	2.2	Iris-virginica
$x_{150}$	5.1	3.4	1.5	0.2	Iris-setosa

**Studi kasus 2** : dengan menggunakan database iris buatlah Kemiripan (similarity) berpasangan antara pasangan titik yang berbeda dengan  $\sigma = \frac{1}{\sqrt{2}}$  (kita tidak memperbolehkan adanya lingkaran, untuk menjaga kesederhanaan grafik). Rata-rata kemiripan antara titik adalah 0.197, dengan standar deviasi 0.290. Matriks adjacency biner menggunakan ambang batas  $\tau=0.777$ , yang menghasilkan sisi antara titik-titik yang memiliki kemiripan lebih tinggi dari dua standar deviasi dari rata-rata. Grafik yang dihasilkan memiliki 150 node dan 753 sisi.

#### Source code python:

kita memanggil dataset iris menggunakan Sklearn dan mengambil data sebanyak 150 data dengan random.



```

1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import load_iris
4 import networkx as nx
5 import matplotlib.pyplot as plt
6 from sklearn.utils import shuffle
7 # Load and shuffle iris data
8 data = load_iris()
9 iris_features, iris_target = shuffle(data.data, data.target, random_state=10000)
10 # Use all 150 samples
11 iris_features = iris_features[:150]

```

Definisikan perumusan dari graf:

```

13 # Define sigma and calculate similarity matrix using the given formula
14 sigma = 1 / np.sqrt(2)
15 similarity_matrix = np.exp(-np.square(np.linalg.norm(iris_features[:, np.newaxis]
16                                                         - iris_features, axis=2)) / (2 * sigma**2))
17 # Calculate mean and standard deviation of similarities
18 mean_similarity = np.mean(similarity_matrix)#mean
19 std_similarity = np.std(similarity_matrix) #standard deviation of similarities
20 # Set threshold to mean + 2 * std deviation
21 threshold = mean_similarity + 2 * std_similarity

22 # Create a graph using NetworkX
23 graph = nx.Graph()
24 # Add nodes to the graph
25 for i in range(len(iris_features)):
26     graph.add_node(i, label=data.target_names[iris_target[i]])
27 # Add edges based on similarity threshold
28 for i in range(len(iris_features)):
29     for j in range(i + 1, len(iris_features)):
30         if similarity_matrix[i, j] >= threshold:
31             graph.add_edge(i, j, weight=similarity_matrix[i, j])
32 # Draw the graph with different symbols for each species
33 pos = nx.spring_layout(graph, seed=150) # Position nodes using spring layout
34 plt.figure(figsize=(8, 5))

```

```

35 # Define markers and colors for each species
36 markers = {'setosa': 's', 'versicolor': 'o', 'virginica': '^'}
37 colors = {'setosa': 'green', 'versicolor': 'blue', 'virginica': 'coral'}
38 node_shapes = {species: [] for species in markers.keys()}
39 # Assign nodes to their respective shapes and colors based on species
40 for i in range(len(iris_features)):
41     species = data.target_names[iris_target[i]]
42     node_shapes[species].append(i)
43 # Draw nodes with different shapes and transparent colors for each species
44 for species, nodes in node_shapes.items():
45     nx.draw_networkx_nodes(graph, pos, nodelist=nodes, node_shape=markers[species],
46                           node_color=colors[species], alpha=0.77, node_size= 5, label=species)

47 # Draw edges and labels
48 nx.draw_networkx_edges(graph, pos, edge_color='grey', alpha=0.3)
49 nx.draw_networkx_labels(graph, pos, labels={i: '' for i in range(len(iris_features))}, font_size=10)
50 plt.legend(scatterpoints=1)
51 plt.title("Iris Similarity Graph")
52 plt.show()

```

b. Topological Attributes

dalam graf ada beberapa atribut yang digunakan seperti mencari radius dan diameter connected graph serta derajat dari graf.

**Studi kasus untuk eksentrisitas, jari-jari, diameter:** jika anda lakukan analisis graf menggunakan dataset **Iris** dengan beberapa langkah. Anda akan memanfaatkan jarak Euclidean untuk membangun graf, menghitung properti graf seperti eksentrisitas, jari-jari, dan diameter. Kita akan meng

Source code python :

```

1 import pandas as pd
2 import numpy as np
3 import networkx as nx
4 import matplotlib.pyplot as plt
5 from sklearn.datasets import load_iris
6 # 1. Load the Iris dataset and prepare it as a DataFrame
7 iris = load_iris()
8 data = pd.DataFrame(data=iris.data, columns=iris.feature_names)

```

```

9 # 2. Create Graph using Euclidean distance
10 def create_graph(data):
11     G = nx.Graph()
12     for i in range(len(data)):
13         for j in range(i + 1, len(data)):
14             distance = np.linalg.norm(data.iloc[i] - data.iloc[j])
15             G.add_edge(i, j, weight=distance)
16     return G
17 # Create the graph from the data
18 G = create_graph(data)
19 # 3. Compute Eccentricity
20 eccentricity = nx.eccentricity(G)
21 print("Eccentricity of each node:", eccentricity)
22 # 4. Compute Radius
23 radius = min(eccentricity.values())
24 print("Radius of the graph:", radius)
25 # 5. Compute Diameter
26 diameter = nx.diameter(G)
27 print("Diameter of the graph:", diameter)

```

**Studi kasus untuk derajat(degree):** buatlah distribusi derajat (*degree distribution*) dari data iris dengan menggunakan histogram.

**Studi kasus untuk degree distribution:**

```

1 import numpy as np
2 import networkx as nx
3 import matplotlib.pyplot as plt
4 from sklearn import datasets
5 from sklearn.preprocessing import StandardScaler
6 # 1. Load the Iris dataset
7 iris = datasets.load_iris()
8 X = iris.data
9 # 2. Standardize the data
10 scaler = StandardScaler()
11 X_scaled = scaler.fit_transform(X)

```

```

12 # 3. Compute pairwise similarities using Euclidean distance
13 def pairwise_similarity(X):
14     n_samples = X.shape[0]
15     similarity_matrix = np.zeros((n_samples, n_samples))
16     for i in range(n_samples):
17         for j in range(i+1, n_samples):
18             distance = np.linalg.norm(X[i] - X[j])
19             similarity_matrix[i, j] = np.exp(-distance) # Similarity function
20             similarity_matrix[j, i] = similarity_matrix[i, j] # Ensure symmetry
21     return similarity_matrix

22 # Compute similarity matrix
23 similarity_matrix = pairwise_similarity(X_scaled)
24 # 4. Compute mean and standard deviation of similarities
25 mean_similarity = np.mean(similarity_matrix)
26 std_similarity = np.std(similarity_matrix)
27 threshold = mean_similarity + 2 * std_similarity
28 # 5. Create binary adjacency matrix based on the threshold
29 adjacency_matrix = (similarity_matrix > threshold).astype(int)
30 # Ensure the adjacency matrix is symmetric
31 adjacency_matrix = np.maximum(adjacency_matrix, adjacency_matrix.T)
32 # 6. Build the graph from the adjacency matrix
33 G = nx.DiGraph(adjacency_matrix)
34 # 7. Calculate the degree of each node
35 degrees = [G.degree(n) for n in G.nodes()]

36 # 8. Visualize the degree distribution using a histogram
37 plt.figure(figsize=(8, 6))
38 plt.hist(degrees, bins=range(min(degrees), max(degrees) + 2),
39         align='left', rwidth=0.8, color='lightgrey', edgecolor='black')
40 plt.title("Degree Distribution of the Iris Graph")
41 plt.xlabel("Degree")
42 plt.ylabel("Frequency")
43 # Add frequency numbers on top of each bar
44 for i, count in enumerate(np.bincount(degrees)):
45     if count > 0:
46         plt.text(i, count, str(count), ha='center', va='bottom')
47 plt.show()

```