



Week 2

🕒 Created @February 10, 2025 3:29 PM

1. Some interesting facts :

- Take all the words in English language as nodes, and draw edge from one word to another if it is a synonym. The resultant graph is connected. Interesting part here is, we can even travel from node 'beautiful' to 'ugly'.
- Take all websites/webpages as nodes and draw an edge from one node to another, if there is a hyperlink between them. This resultant graph is called "Web graph". This was used by Google to rank the pages/websites for their search engine.
- Take each person as node, and draw edge between them if they are friends on social media. This graph can be used to predict who is the next probable friend of a person.

2. Network Datasets Format :

- CSV
- GML
- PajekNet
- GraphML
- GEXF

2.1 : CSV

Comma Separated Values

Extension : .csv, .txt

Types : Edgelist, Adjacency list

Edgelist :

CSV Format- Edgelist	
0	344
0	345
0	346
0	347
1	48
1	53
1	54
1	73
1	88
1	92

these edges you can also do that, for example, you want to add weights to these edges like this you

CSV Format- Weighted Edgelist		
1	2	0.3
3	4	0.1
5	7	0.2
4	7	0.8

can add third value to every row. So for example the edge which is going to go from one to two

- Format : Sourcenode Targetnode
- It is simple and contains list of edges from source node to target node.
- It can also have weights associated with each edge .
- It cannot have any non-numeric weight, add any attribute, change any labels.

Adjacency List :

CSV Format- Adjlist	
1	2 5 7
2	4 6
3	1 4 7
4	6 2 3
5	7 2
6	1
7	2 8 4
8	9
9	1 7 4 8

values in every row are the nodes which are adjacent to this source node. So this first

- Format : Soucenode Neighbours
- It is simple and contains sourcenode and its neighbours in each row.
- It also cannot add any attribute , change any labels.

2.2 : GML

```

GML Format : Graph Modeling Language

graph
[
  node
  [
    id A
  ]
  node
  [
    id B
  ]
  node
  [
    id C
  ]
  edge
  [
    source B
    target A
  ]
  edge
  [
    source C
    target A
  ]
]

```

```

GML Format with Labels

graph
[
  node
  [
    id A
    label "Node A"
  ]
  node
  [
    id B
    label "Node B"
  ]
  node
  [
    id C
    label "Node C"
  ]
  edge
  [
    source B
    target A
    label "Edge B to A"
  ]
  edge
  [
    source C
    target A
    label "Edge C to A"
  ]
]

```

- Graph Modelling Language
- Most commonly used format for datasets.
- Each node , edge is defined with keyword and [] .
- Labels and attributes can also be added here.

2.3 : PajekNet

```

Pajek Net Format : Uses .NET extension

*Vertices 82670
1 "entity" 4244 107
2 "thing" 4244 238
3 "anything" 4244 4292
4 "something" 4247 107
5 "nothing" 4248 1
6 "whole" 4248 54

```

```

Pajek Net Format with attributes

*arcs
4244 107 5

```

- Extension : .net , .paj
- Format : *Vertices no.of vertices , Then each row contains vertex number and label
 * arcs/ * Edges , Then each row contains source node and target node and attribute

2.4 : GraphML

GraphML Format: Uses .graphml extension

```
<?xml version="1.0" encoding="UTF-8"?>
<graphml xmlns="http://graphml.graphdrawing.org/xmlns"
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xsi:schemaLocation="http://graphml.graphdrawing.org/xmlns
http://graphml.graphdrawing.org/xmlns/1.0/graphml.xsd">
<graph id="G" edgedefault="undirected">
<node id="n0"/>
<node id="n1"/>
<edge id="e1" source="n0" target="n1"/>
</graph>
</graphml>
```

- ML denotes XML, it uses various tags like in html.
- <graph id="g" edgedefault = "undirected">
- <node id="n0">
- <edge id="e1" soucenode="n0" targetnode="n1">

2.5 : GEXF

GEXF Format : Graph Exchange XML Format

```
<?xml version="1.0" encoding="UTF-8"?>
<gexf xmlns="http://www.gexf.net/1.2draft" version="1.2">
<meta lastmodifieddate="2009-03-20">
<creator>Gexf.net</creator>
<description>A hello world! file</description>
</meta>
<graph mode="static" defaultedgetype="directed">
<nodes>
<node id="0" label="Hello" />
<node id="1" label="Word" />
</nodes>
<edges>
<edge id="0" source="0" target="1" />
</edges>
</graph>
</gexf>
```

- Graph Exchange XML Format
- Similar to GraphML but more cleaner than it.

3. Analyzing Network Datasets :

```

import networkx as nx
import matplotlib.pyplot as plt

G = nx.read_edgelist('dataset path') # reads edgelist dataset and return a graph
print(nx.info(G)) # show basic details of the graph
print(nx.number_of_nodes(G)) # return no. of nodes
print(nx.number_of_edges(G)) # return no. of edges
print(nx.is_directed(G)) # return True or False

G = nx.read_pajek('dataset path') # to read a pajek dataset
G = nx.read_gml('dataset path') # to read gml dataset
G = nx.read_graphml('dataset path') # to read graphml dataset
G = nx.read_gexf('dataset path') # to read gexf dataset

nx.draw(G) # to draw the graph
nx.draw_circular(G) # nodes will be arranged circularly
plt.show() # to see the output

```

Degree :

The degree of the node is the total number of edges it has.

For undirected graph - degree

For directed graph - indegree, outdegree

Density :

Number of edges / Total possible edges

Density value of complete graph is 1.0

Total possible edges = $n(n-1) / 2$ [n - number of nodes]

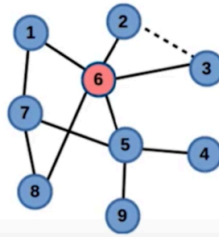
Clustering Coefficient :

Actual edges between neighbours / Total possible edges between neighbours

It is calculated for each node.

Clustering coefficient of a graph is average of all nodes.

Clustering Coefficient contd.



$$\text{Clustering Coefficient}(6) = \frac{1}{10}$$

Here , consider Node 6 , its neighbours are 1,2,3,5,8

Actual edges between neighbours = 1 (from 2 to 3)

Total possible edges between neighbours = 10 [$n*(n-1)/2$]

so, $1/10$

Diameter :

Take shortest path from each node to each other node. Diameter is the longest path in this.

In other words. it is the maximum of all shortest paths from one node to other.

```
import networkx as nx
import matplotlib.pyplot as plt

G = nx.read_gml('dataset path')
nx.degree(G)
# returns dictionary with node : degree

nx.density(G)
# returns density of graph

nx.clustering(G)
# returns dictionary with node : clustering coefficient
nx.average_clustering(G)
# returns average clustering coefficient
```

```
nx.diameter(G)  
#returns the diameter of graph
```

4. Emergence of connectedness :

If there is a graph with “n” nodes and if we add “n log n” edges uniformly at random the probability that a node will be isolated is $1/n^2$.

So, For 100 nodes after $100 \log 100$ edges, probability that a node will be isolated = $1/10,000$

That is, probability the graph is not connected = $1/10,000$

Number of edges required to make a graph connected = $n \log n$

More specifically, it is order of $n \log n$