

ĐẠI HỌC QUỐC GIA TP. HỒ CHÍ MINH  
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN



ĐỒ ÁN MÔN HỌC  
MẠNG XÃ HỘI

PHÂN TÍCH DỮ LIỆU CHUỖI CUNG ỨNG CHO  
CÁC LOẠI SẢN PHẨM THÔNG QUA CÁC THUỐC ĐO VỀ  
TRỌNG TÂM VÀ KẾT QUẢ PHÂN CỤM

SINH VIÊN THỰC HIỆN  
ĐẶNG HOÀNG QUÂN – 18520339

GIẢNG VIÊN HƯỚNG DẪN  
NGUYỄN THỊ KIM PHỤNG

TP. HỒ CHÍ MINH – 11/2021

# MỤC LỤC

<b>DANH MỤC HÌNH ẢNH .....</b>	<b>1</b>
<b>I. GIỚI THIỆU .....</b>	<b>1</b>
<b>II. XÁC ĐỊNH BÀI TOÁN.....</b>	<b>1</b>
<b>III. DỮ LIỆU.....</b>	<b>1</b>
1. Giới thiệu nguồn dữ liệu.....	1
2. Xử lý và phân tích dữ liệu .....	7
2.1. Đọc dữ liệu từ file vào DataFrame .....	7
2.2. Làm sạch dữ liệu.....	7
2.3. Chuyển đổi DataFrame thành đồ thị.....	8
<b>IV. XẾP HẠNG CÁC LOẠI SẢN PHẨM.....</b>	<b>10</b>
1. Định nghĩa các hàm hiển thị kết quả độ đo .....	10
2. Degree Centrality .....	11
2.1. Code chạy độ đo.....	11
2.2. Kết quả độ đo.....	11
2.3. Nhận xét .....	12
3. Betweenness Centrality .....	13
3.1. Code chạy độ đo.....	13
3.2. Kết quả độ đo.....	13
3.3. Nhận xét .....	13
4. Closeness Centrality.....	14
4.1. Code chạy độ đo.....	14
4.2. Kết quả độ đo.....	15
4.3. Nhận xét .....	15
5. Eigenvector Centrality.....	16
5.1. Code chạy độ đo.....	16
5.2. Kết quả độ đo.....	16
5.3. Nhận xét .....	17
6. PageRank .....	18
6.1. Code chạy độ đo.....	18

6.2. Kết quả độ đo.....	18
6.3. Nhận xét.....	19
<b>V. THUẬT TOÁN PHÁT HIỆN CỘNG ĐỒNG.....</b>	<b>20</b>
1. Định nghĩa các hàm hiển thị kết quả phân cụm.....	20
2. Thuật toán Louvain .....	21
2.1. Code chạy thuật toán .....	21
2.2. Đồ thị phân cụm.....	21
2.3. Ý nghĩa các cụm.....	22
2.4. Nhận xét .....	23
3. Thuật toán K-Means.....	23
3.1. Code chạy thuật toán .....	23
3.2. Đồ thị phân cụm.....	25
3.3. Ý nghĩa các cụm.....	25
3.4. Nhận xét .....	26
4. Gaussian Mixture Model .....	27
4.1. Code chạy thuật toán .....	27
4.2. Đồ thị phân cụm.....	27
4.3. Ý nghĩa các cụm.....	27
4.4. Nhận xét .....	28
5. Trực quan hóa các điểm chung trong các cụm .....	29
<b>VI. TÀI LIỆU THAM KHẢO.....</b>	<b>30</b>

# DANH MỤC HÌNH ẢNH

Hình 1. Đọc dữ liệu từ file vào DataFrame.....	7
Hình 2. Kiểm tra các giá trị bị khuyết .....	7
Hình 3. Loại bỏ các giá trị trùng lặp.....	7
Hình 4. Chuyển đổi DataFrame thành đồ thị.....	8
Hình 5. Đồ thị 2 phía.....	8
Hình 6. Code hiển thị đồ thị 2 phía .....	9
Hình 7. Đồ thị 1 phía.....	9
Hình 8. Code hiển thị đồ thị 1 phía .....	10
Hình 9. Code hiển thị Cliques.....	10
Hình 10. Hàm in kết quả độ đo .....	10
Hình 11. Hàm trực quan hóa kết quả độ đo.....	11
Hình 12. Code chạy Degree Centrality .....	11
Hình 13. Kết quả Degree Centrality.....	11
Hình 14. Kết quả Degree Centrality trên Gephi .....	12
Hình 15. Code chạy Betweenness Centrality .....	13
Hình 16. Kết quả Betweenness Centrality .....	13
Hình 17. Kết quả Betweenness Centrality trên Gephi.....	14
Hình 18. Code chạy Closeness Centrality.....	14
Hình 19. Kết quả Closeness Centrality .....	15
Hình 20. Kết quả Closeness Centrality trên Gephi .....	16
Hình 21. Code chạy Eigenvector Centrality.....	16
Hình 22. Kết quả Eigenvector Centrality.....	17
Hình 23. Kết quả Eigenvector Centrality trên Gephi.....	18
Hình 24. Code chạy PageRank .....	18
Hình 25. Kết quả PageRank.....	19
Hình 26. Kết quả PageRank trên Gephi.....	20
Hình 27. Hàm lấy đặc điểm chung của 1 cụm.....	20
Hình 28. Hàm in ra các cộng đồng cùng điểm chung của các node trong chúng ..	21

Hình 29. Hàm trực quan hóa các cộng đồng.....	21
Hình 30. Code chạy thuật toán Louvain .....	21
Hình 31. Đồ thị phân cụm sử dụng Louvain .....	22
Hình 32. Kết quả phân cụm sử dụng Louvain.....	22
Hình 33. Chuyển đổi đồ thị thành ma trận kề làm đầu vào cho K-Means .....	23
Hình 34. Code triển khai Elbow Method cho thuật toán K-Means .....	24
Hình 35. Chọn ra số cụm k cho thuật toán K-Means bằng Elbow Method.....	24
Hình 36. Chạy K-Means với đầu vào là ma trận kề cùng số cụm k đã chọn.....	25
Hình 37. Đồ thị phân cụm sử dụng K-Means.....	25
Hình 38. Kết quả phân cụm sử dụng K-Means .....	25
Hình 39. Code chạy mô hình Gaussian Mixture.....	27
Hình 40. Đồ thị phân cụm sử dụng GMM .....	27
Hình 41. Kết quả phân cụm sử dụng GMM.....	27
Hình 42. Các khu vực đã cung cấp các sản phẩm của cụm 0 khi dùng Louvain...	29
Hình 43. Các khu vực đã cung cấp các sản phẩm của cụm 1 khi dùng Louvain...	29
Hình 44. Các khu vực đã cung cấp các sản phẩm của cụm 2 khi dùng Louvain...	30

## I. GIỚI THIỆU

Ngày nay, nhiều công ty phải đối mặt với một thách thức dường như đầy mâu thuẫn: làm thế nào để giảm chi phí vận hành đồng thời tăng mức độ dịch vụ khách hàng. Thiết kế mạng lưới chuỗi cung ứng phù hợp cung cấp giải pháp cho cả hai vấn đề trên. Mặc dù có nhiều yếu tố cần xem xét khi thiết kế mạng lưới chuỗi cung ứng, nhưng quá trình này không quá phức tạp với các đối tác phù hợp.

Môn học mạng xã hội sẽ giúp phân tích và xem xét chuỗi cung ứng của công ty bằng cách sử dụng các thước đo về trọng tâm và phân cụm liên quan đến các khía cạnh nổi bật (dựa trên đường đi ngắn nhất, phổ, khoảng cách, ...) để từ đó đưa ra các chiến lược vận hành phù hợp nhất.

## II. XÁC ĐỊNH BÀI TOÁN

- Input: Tập dữ liệu ban đầu trên nguồn dữ liệu Kaggle được qua tiền xử lý dữ liệu.
- Output: Đưa ra độ đo, cộng đồng phục vụ cho việc phân tích mạng xã hội **DataCo SMART SUPPLY CHAIN FOR BIG DATA ANALYSIS**.

## III. DỮ LIỆU

### 1. Giới thiệu nguồn dữ liệu

- Link dataset: <https://www.kaggle.com/shashwatwork/dataco-smart-supply-chain-for-big-data-analysis>.
- Dữ liệu **DataCo SMART SUPPLY CHAIN FOR BIG DATA ANALYSIS** được cung cấp trên nền tảng Kaggle chứa dữ liệu của chuỗi cung ứng được sử dụng bởi công ty DataCo Global.
- Dataset bao gồm tập hợp các sản phẩm đã bán của công ty, chi tiết tài chính (lãi, lỗ, tổng doanh thu, v.v.), chi tiết giao hàng và chi tiết khách hàng như doanh số, nhân khẩu học và chi tiết giao dịch.
- Dữ liệu có kích thước 91 MB bao gồm 180519 dòng với 54 thuộc tính:

Thuộc tính	Type	Mô tả	Các giá trị
Type	char	Loại giao dịch thực hiện	<ul style="list-style-type: none"> <li>• CASH</li> <li>• DEBIT</li> <li>• PAYMENT</li> <li>• TRANSFER</li> </ul>
Days for shipping (real)	int	Số ngày giao hàng thực tế	
Days for shipment (scheduled)	int	Số ngày giao hàng theo dự kiến	
Benefit per order	float	Thu nhập cho mỗi đơn hàng được đặt	= Order Profit Per Order
Sales per customer	float	Tổng doanh số theo khách hàng	= Order Item Total
Delivery Status	char	Trạng thái giao hàng của đơn hàng	<ul style="list-style-type: none"> <li>• Advance shipping</li> <li>• Late delivery</li> <li>• Shipping canceled</li> <li>• Shipping on time</li> </ul>
Late_delivery_risk	int	Biến phân loại cho biết gửi muộn hay không	<ul style="list-style-type: none"> <li>• 1 – Gửi muộn</li> <li>• 0 – Không gửi muộn</li> </ul>
Category Id	int	Mã danh mục sản phẩm	= Product Category Id
Category Name	char	Tên danh mục sản phẩm	
Customer City	char	Thành phố nơi khách hàng thực hiện mua hàng	
Customer Country	char	Đất nước nơi khách hàng thực hiện mua hàng	
Customer Email	char	Email của khách hàng	XXXXXXXXXX

Customer Fname	char	Tên khách hàng	
Customer Id	int	ID khách hàng	
Customer Lname	char	Họ khách hàng	
Customer Password	char	Mật khẩu khách hàng	XXXXXXXXXX
Customer Segment	char	Phân khúc khách hàng	<ul style="list-style-type: none"> <li>• Consumer</li> <li>• Corporate</li> <li>• Home Office</li> </ul>
Customer State	char	Tiểu bang của cửa hàng đã đăng ký giao dịch mua	
Customer Street	char	Đường của cửa hàng đã đăng ký giao dịch mua	
Customer Zipcode	float	Mã bưu điện khách hàng	
Department Id	int	Mã bộ phận của cửa hàng	
Department Name	char	Tên bộ phận của cửa hàng	
Latitude	float	Vĩ độ tương ứng với vị trí của cửa hàng	
Longitude	float	Kinh độ tương ứng với vị trí của cửa hàng	
Market	char	Thị trường của nơi được giao hàng	<ul style="list-style-type: none"> <li>• Africa</li> <li>• Europe</li> <li>• LATAM</li> <li>• Pacific Asia</li> <li>• USCA</li> </ul>



Order City	char	Thành phố của đơn hàng được đặt	
Order Country	char	Quốc gia của đơn hàng được đặt	
Order Customer Id	int	Mã đặt hàng của khách	= Customer Id
order date (DateOrders)	datetime	Ngày đặt hàng	
Order Id	int	Mã đơn hàng	
Order Item Cardprod Id	int	Mã sản phẩm được tạo thông qua đầu đọc RFID	= Product Card Id
Order Item Discount	float	Giá trị chiết khấu của mặt hàng trong đơn hàng	= Sales - Order Item Total
Order Item Discount Rate	float	Phần trăm chiết khấu của mặt hàng trong đơn hàng	
Order Item Id	int	Mã mặt hàng được đặt trong đơn hàng	
Order Item Product Price	float	Giá của sản phẩm khi không giảm giá	= Product Price
Order Item Profit Ratio	float	Tỷ lệ lợi nhuận của mặt hàng trong đơn hàng	
Order Item Quantity	int	Số lượng sản phẩm c mỗi mặt hàng trong đơn hàng	
Sales	float	Doanh số bán hàng	= Order Item Product Price * Order Item Quantity
Order Item Total	float	Tổng số tiền cho mặt hàng trong đơn hàng	= Sales per customer

Order Profit Per Order	float	Lợi nhuận cho mỗi mặt hàng trong đơn hàng	= Benefit per order
Order Region	char	Khu vực nơi đơn đặt hàng được tiến hành giao hay khu vực sẽ cung cấp loại sản phẩm	<ul style="list-style-type: none"> <li>• Southeast Asia</li> <li>• South Asia</li> <li>• Oceania</li> <li>• Eastern Asia</li> <li>• West Asia</li> <li>• West of USA</li> <li>• US Center</li> <li>• West Africa</li> <li>• Central Africa</li> <li>• North Africa</li> <li>• Western Europe</li> <li>• Northern</li> <li>• Caribbean</li> <li>• South America</li> <li>• East Africa</li> <li>• Southern Europe</li> <li>• East of USA</li> <li>• Canada</li> <li>• Southern Africa</li> <li>• Central Asia</li> <li>• Europe</li> <li>• Central America</li> <li>• Eastern Europe</li> <li>• South of USA</li> </ul>
Order State	char	Tiểu bang nơi đơn hàng được tiến hành giao	

Order Status	char	Trạng thái đơn hàng	<ul style="list-style-type: none"> <li>• CANCELED</li> <li>• CLOSED</li> <li>• COMPLETE</li> <li>• ON_HOLD</li> <li>• PAYMENT_REVIEW</li> <li>• PENDING</li> <li>• PENDING_PAYMENT</li> <li>• PROCESSING</li> <li>• SUSPECTED_FRAUD</li> </ul>
Order Zipcode	float	Mã bưu điện đơn hàng	
Product Card Id	int	Mã sản phẩm	= Order Item Cardprod Id
Product Category Id	int	Mã danh mục sản phẩm	= Category Id
Product Description	float	Mô tả sản phẩm	
Product Image	char	Liên kết đến hình ảnh của sản phẩm	
Product Name	char	Tên sản phẩm	
Product Price	float	Giá sản phẩm	= Order Item Product Price
Product Status	int	Trạng thái sản phẩm	<ul style="list-style-type: none"> <li>• 1 – Không có sẵn</li> <li>• 0 – Có sẵn</li> </ul>
shipping date (DateOrders)	datetime	Ngày và thời gian chính xác của lô hàng	
Shipping Mode	char	Chế độ vận chuyển	<ul style="list-style-type: none"> <li>• First Class</li> <li>• Same Day</li> <li>• Second Class</li> <li>• Standard Class</li> </ul>

## 2. Xử lý và phân tích dữ liệu

### 2.1. Đọc dữ liệu từ file vào DataFrame

```
1 import matplotlib.pyplot as plt
2 plt.style.use('seaborn')

1 import pandas as pd
2 PRIMARY = 'Category Name'
3 SECONDARY = 'Order Region'
4
5 df = pd.read_csv(
6     'DataCoSupplyChainDataset.csv',
7     usecols = [PRIMARY, SECONDARY],
8     encoding = 'unicode_escape'
9 ).apply(lambda col: col.str.strip())
10 df.head()
```

	Category Name	Order Region
0	Sporting Goods	Southeast Asia
1	Sporting Goods	South Asia
2	Sporting Goods	South Asia
3	Sporting Goods	Oceania
4	Sporting Goods	Oceania

Hình 1. Đọc dữ liệu từ file vào DataFrame

### 2.2. Làm sạch dữ liệu

- Kiểm tra các giá trị bị khuyết ➔ không phát hiện giá trị nào nên không cần loại bỏ:

```
1 df.isnull().sum().sort_values(ascending=False)
```

Category Name	0
Order Region	0

Hình 2. Kiểm tra các giá trị bị khuyết

- Loại bỏ các giá trị trùng lặp ➔ Kết quả cuối cùng nhận được là 1 bộ dữ liệu gồm 691 dòng và 2 cột:

```
1 df.drop_duplicates(inplace=True)
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 691 entries, 0 to 162002
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Category Name    691 non-null    object
1   Order Region     691 non-null    object
dtypes: object(2)
memory usage: 16.2+ KB
```

Hình 3. Loại bỏ các giá trị trùng lặp

## 2.3. Chuyển đổi DataFrame thành đồ thị

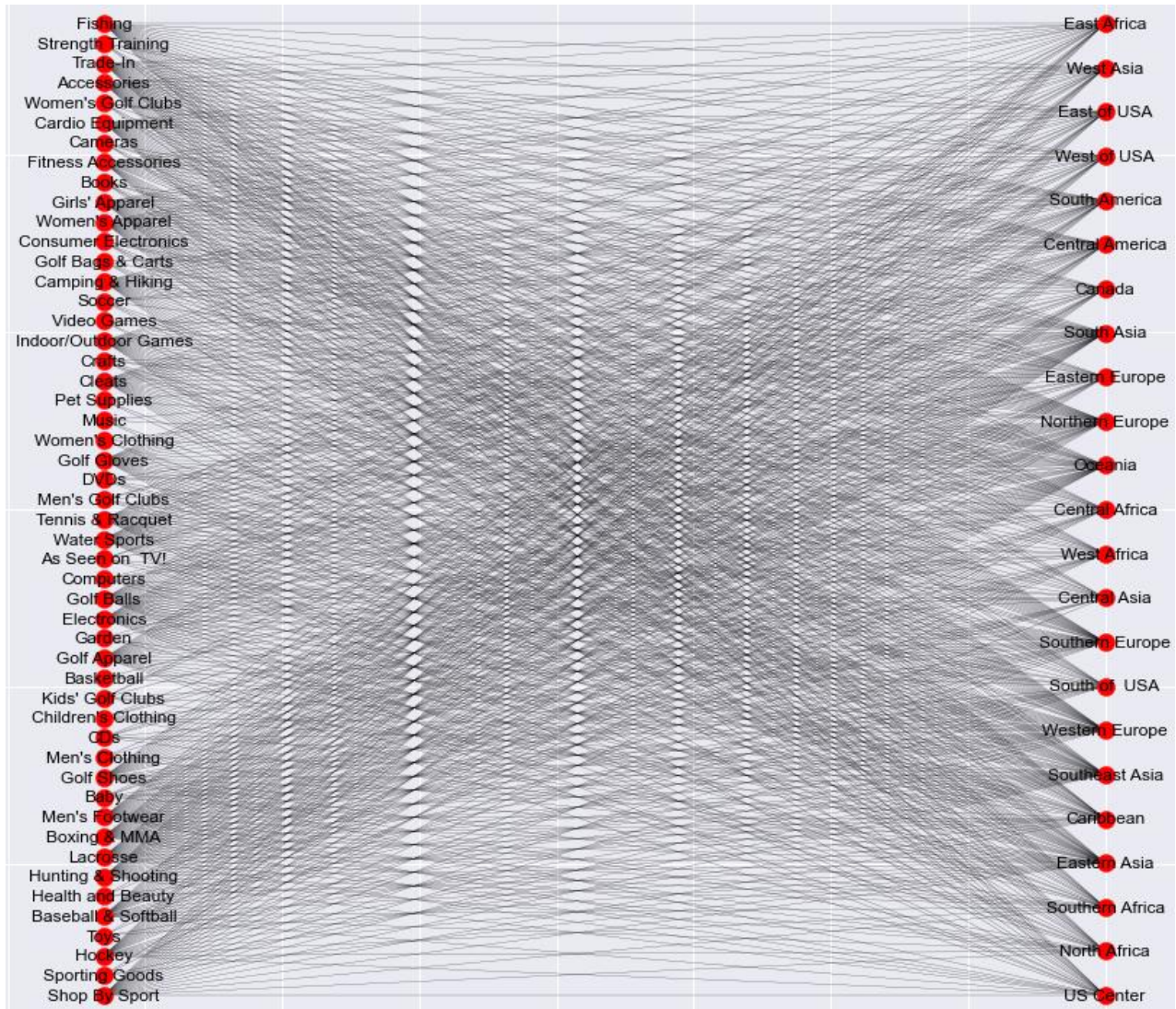
```
1 primary_col, secondary_col = df[PRIMARY], df[SECONDARY]
2 print(f'Number of {PRIMARY}:', primary_col.nunique())
3 print(f'Number of {SECONDARY}:', secondary_col.nunique())
4 print('Number of edges:', len(df))
```

Number of Category Name: 50  
Number of Order Region: 23  
Number of edges: 691

Hình 4. Chuyển đổi DataFrame thành đồ thị

### 2.2.1. Đồ thị 2 phía

- Node: danh mục các loại sản phẩm (Category Name) và khu vực nơi mà đơn đặt hàng được tiến hành giao hay có thể nói cách khác là khu vực sẽ cung cấp loại sản phẩm (Order Region).
- Edge: mối quan hệ hình thành khi 1 loại sản phẩm được cung cấp từ 1 khu vực.



Hình 5. Đồ thị 2 phía



⇒ Nhìn vào đồ thị có thể thấy cùng 1 loại sản phẩm có thể được cung cấp từ nhiều khu vực và 1 khu vực cũng có thể cung cấp nhiều loại sản phẩm.

- Code hiển thị đồ thị 2 phía:

```
1 import networkx as nx
2 from networkx.algorithms import bipartite
3
4 B = nx.Graph()
5 for index, row in df.iterrows():
6     B.add_edge(row[PRIMARY], row[SECONDARY], weight=1)
7
8 B.add_nodes_from(primary_col, bipartite=0)
9 B.add_nodes_from(secondary_col, bipartite=1)

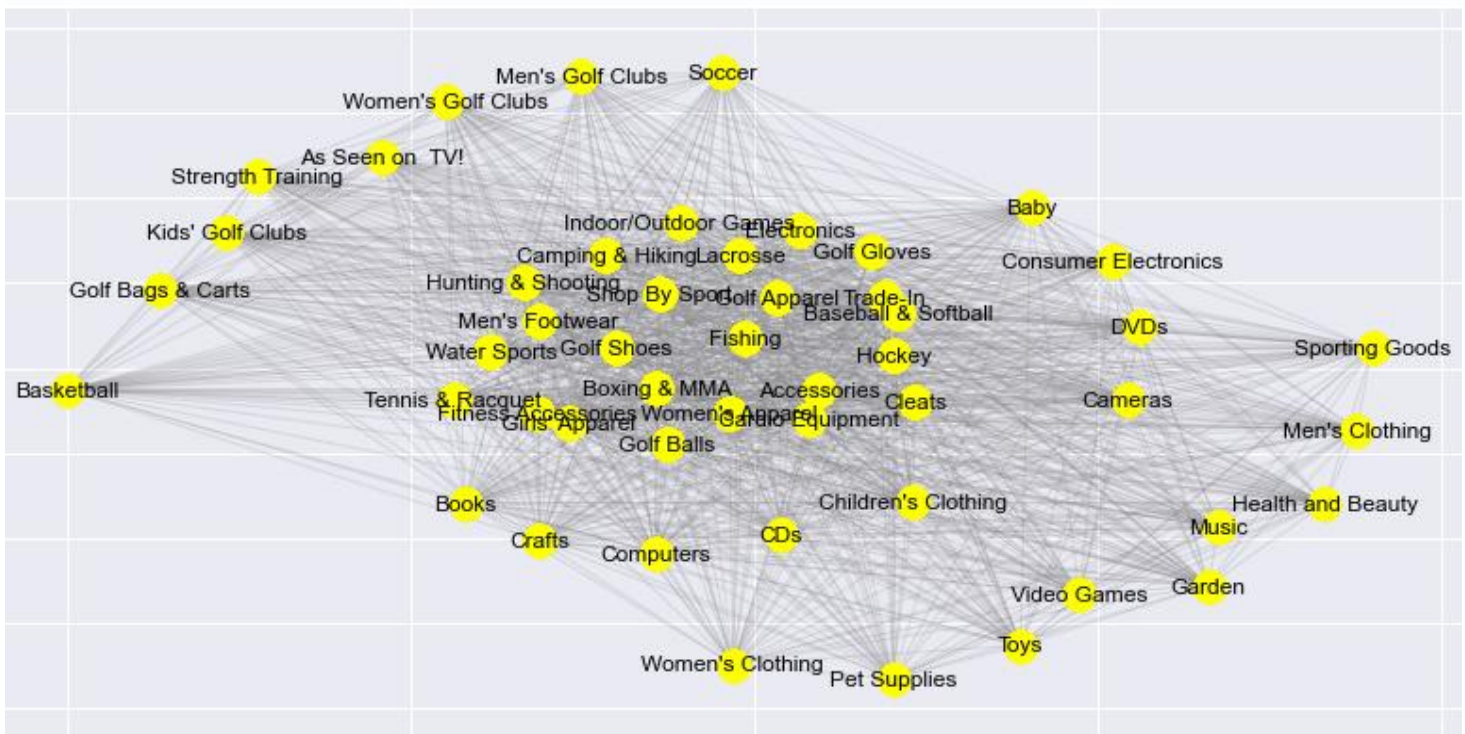
```

```
1 import matplotlib.pyplot as plt
2 plt.figure(figsize=(15, 15))
3 pos = nx.drawing.layout.bipartite_layout(B, primary_col)
4 nx.draw_networkx(B, pos=pos, node_size=150, width=0.2, node_color='red')
```

Hình 6. Code hiển thị đồ thị 2 phía

### 2.2.2. Đồ thị 1 phía

- Node: danh mục các loại sản phẩm (Category Name).
- Edge: 2 loại sản phẩm cùng được cung cấp từ 1 khu vực sẽ được nối với nhau tạo thành cạnh, ý nghĩa nói lên với cùng 1 khu vực có thể cung cấp nhiều loại sản phẩm.
- Weight: trọng số là số khu vực giống nhau đã cung cấp 2 loại sản phẩm.



Hình 7. Đồ thị 1 phía

- Code hiển thị đồ thị 1 phía:

```

1 G = bipartite.weighted_projected_graph(B, primary_col)
2 node_labels = dict(zip(primary_col, primary_col))
3 nx.to_pandas_edgelist(G).to_csv('Edges table.csv', index=False)

1 plt.figure(figsize=(15, 7))
2 pos = nx.spring_layout(G)
3
4 nx.draw_networkx_nodes(G, pos, nodelist=primary_col, node_color='yellow')
5 nx.draw_networkx_edges(G, pos, edge_color='grey', alpha=0.2)
6 nx.draw_networkx_labels(G, pos, labels=node_labels)
7 plt.show()

```

Hình 8. Code hiển thị đồ thị 1 phía

### 2.2.3. Cliques

- Clique được định nghĩa là 1 đồ thị con hoàn chỉnh **cực đại** của 1 đồ thị trong đó mỗi node được kết nối với tất cả các node khác. Từ **cực đại** có nghĩa là nếu ta thêm một node khác vào Clique, nó sẽ không còn là 1 Clique.
- Code hiển thị Cliques:

```

1 cliques = list(nx.find_cliques(G))
2 print('Number of cliques:', len(cliques))
3 for index, clique in enumerate(cliques):
4     print(f'\nClique {index}:', clique)

```

Number of cliques: 2

Clique 0: ['Golf Shoes', 'Electronics', 'CDs', 'Lacrosse', 'Women's Apparel', 'Fishing', 'Baseball & Softball', 'Water Sports', 'Baby', 'Boxing & MMA', 'Indoor/Outdoor Games', 'DVDs', 'Cleats', 'Men's Footwear', 'Books', 'Tennis & Racquet', 'Fitness Accessories', 'Hunting & Shooting', 'Trade-In', 'Golf Balls', 'Crafts', 'Shop By Sport', 'Hockey', 'Computers', 'Cameras', 'Golf Gloves', 'Golf Apparel', 'Accessories', 'Girls' Apparel', 'Children's Clothing', 'Consumer Electronics', 'Camping & Hiking', 'Cardio Equipment', 'Men's Golf Clubs', 'Basketball', 'As Seen on TV!', 'Golf Bags & Carts', 'Women's Golf Clubs', 'Kids' Golf Clubs', 'Soccer', 'Strength Training']

Clique 1: ['Golf Shoes', 'Electronics', 'CDs', 'Lacrosse', 'Women's Apparel', 'Fishing', 'Baseball & Softball', 'Water Sports', 'Baby', 'Boxing & MMA', 'Indoor/Outdoor Games', 'DVDs', 'Cleats', 'Men's Footwear', 'Books', 'Tennis & Racquet', 'Fitness Accessories', 'Hunting & Shooting', 'Trade-In', 'Golf Balls', 'Crafts', 'Shop By Sport', 'Hockey', 'Computers', 'Cameras', 'Golf Gloves', 'Golf Apparel', 'Accessories', 'Girls' Apparel', 'Children's Clothing', 'Consumer Electronics', 'Camping & Hiking', 'Cardio Equipment', 'Toys', 'Video Games', 'Music', 'Health and Beauty', 'Garden', 'Pet Supplies', 'Sporting Goods', 'Women's Clothing', 'Men's Clothing']

Hình 9. Code hiển thị Cliques

⇒ Dễ thấy 2 Clique rút ra được chỉ khác nhau ở các node ngoài rìa (màu đỏ) có thể nhìn thấy rõ ràng ở đồ thị 1 phía.

## IV. XẾP HẠNG CÁC LOẠI SẢN PHẨM

### 1. Định nghĩa các hàm hiển thị kết quả độ đo

- Hàm in kết quả độ đo đã sắp xếp giảm dần:

```

1 def print_centralty(centrality, name):
2     result = pd.DataFrame(centrality.items(), columns=['Category Name', name])
3     result.sort_values(name, ascending=False, inplace=True)
4     print(result.to_records(index=False).tolist())
5     return result

```

Hình 10. Hàm in kết quả độ đo

- Hàm trực quan hóa kết quả độ đo:

```

1 def plot_centrality(centrality):
2     plt.figure(figsize=(15, 7))
3     node_color = [centrality[i] for i in centrality.keys()]
4
5     cmap = plt.cm.ScalarMappable(
6         cmap = 'Wistia',
7         norm = plt.Normalize(vmin=min(node_color), vmax=max(node_color))
8     )
9     cmap.set_array([])
10    plt.colorbar(cmap)
11
12    pos = nx.spring_layout(G)
13    nx.draw_networkx_nodes(G, pos, node_color=node_color, cmap='Wistia')
14    nx.draw_networkx_edges(G, pos, edge_color='grey', alpha=0.2)
15    nx.draw_networkx_labels(G, pos)
16    plt.show()

```

Hình 11. Hàm trực quan hóa kết quả độ đo

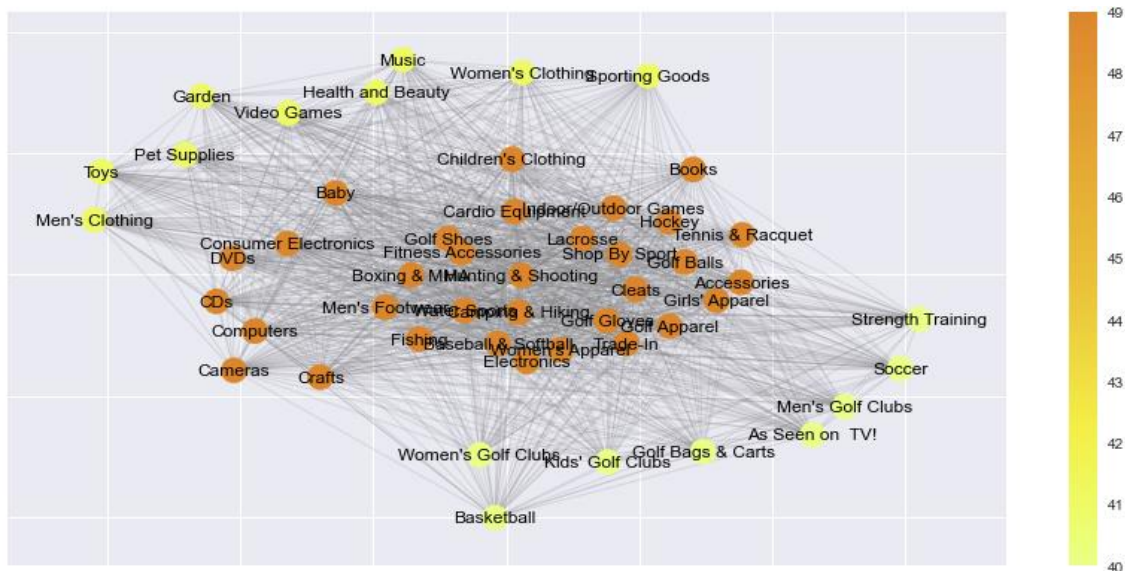
## 2. Degree Centrality

### 2.1. Code chạy độ đo

```
1 degree = dict(nx.degree(G))
2 plot_centrality(degree)
3 print_centrality(degree, 'Degree Centrality').head()
```

Hình 12. Code chạy Degree Centrality

## 2.2. Kết quả độ đo



[('Indoor/Outdoor Games', 49), ('Golf Balls', 49), ('Golf Shoes', 49), ('Crafts', 49), ('Golf Gloves', 49), ('Tennis & Racquet', 49), ('Fitness Accessories', 49), ('Cleats', 49), ('Children's Clothing', 49), ('Golf Apparel', 49), ('Lacrosse', 49), ('Baby', 49), ('Fishing', 49), ('Books', 49), ('DVDs', 49), ('CDs', 49), ('Accessories', 49), ('Girls' Apparel', 49), ('Hockey', 49), ('Shop By Sport', 49), ('Computers', 49), ('Cameras', 49), ('Consumer Electronics', 49), ('Camping & Hiking', 49), ('Men's Footwear', 49), ('Baseball & Softball', 49), ('Hunting & Shooting', 49), ('Water Sports', 49), ('Trade-In', 49), ('Cardio Equipment', 49), ('Boxing & MMA', 49), ('Electronics', 49), ('Women's Apparel', 49), ('Music', 41), ('Health and Beauty', 41), ('Toys', 41), ('Video Games', 41), ('Sporting Goods', 41), ('Pet Supplies', 41), ('Garden', 41), ('Men's Clothing', 41), ('Women's Clothing', 41), ('Strength Training', 40), ('Golf Bags & Carts', 40), ('Soccer', 40), ('Women's Golf Clubs', 40), ('Men's Golf Clubs', 40), ('Basketball', 40), ('Kids' Golf Clubs', 40), ('As Seen on TV!', 40)]

*Hình 13. Kết quả Degree Centrality*



### 2.3. Nhận xét

- Degree của 1 node thể hiện số lượng node mà 1 node nhất định kết nối. Loại sản phẩm có Degree càng cao thì càng có nhiều liên kết, nghĩa là các khu vực mà cung cấp loại sản phẩm này cũng sẽ cung cấp các loại sản phẩm khác.
- Những loại sản phẩm có Degree Centrality cao thể hiện các loại sản phẩm này được cung cấp bởi nhiều khu vực khác nhau, vì vậy nó mới kết nối với nhiều loại sản phẩm khác nhau.
- Như trong hình 13, ta có thể thấy các node ở phía ngoài rìa đồ thị sẽ có Degree Centrality thấp hơn. Bên cạnh đó, kết quả độ đo thực hiện trên Python và Gephi không có sự khác biệt.

Id	Degree		
Accessories	49	Indoor/Outdoor Games	49
Baby	49	Lacrosse	49
Baseball & Softball	49	Men's Footwear	49
Books	49	Shop By Sport	49
Boxing & MMA	49	Tennis & Racquet	49
Cameras	49	Trade-In	49
Camping & Hiking	49	Water Sports	49
Cardio Equipment	49	Women's Apparel	49
CDs	49	Garden	41
Children's Clothing	49	Health and Beauty	41
Cleats	49	Men's Clothing	41
Computers	49	Music	41
Consumer Electronics	49	Pet Supplies	41
Crafts	49	Sporting Goods	41
DVDs	49	Toys	41
Electronics	49	Video Games	41
Fishing	49	Women's Clothing	41
Fitness Accessories	49	As Seen on TV!	40
Girls' Apparel	49	Basketball	40
Golf Apparel	49	Golf Bags & Carts	40
Golf Balls	49	Kids' Golf Clubs	40
Golf Gloves	49	Men's Golf Clubs	40
Golf Shoes	49	Soccer	40
Hockey	49	Strength Training	40
Hunting & Shooting	49	Women's Golf Clubs	40

Hình 14. Kết quả Degree Centrality trên Gephi

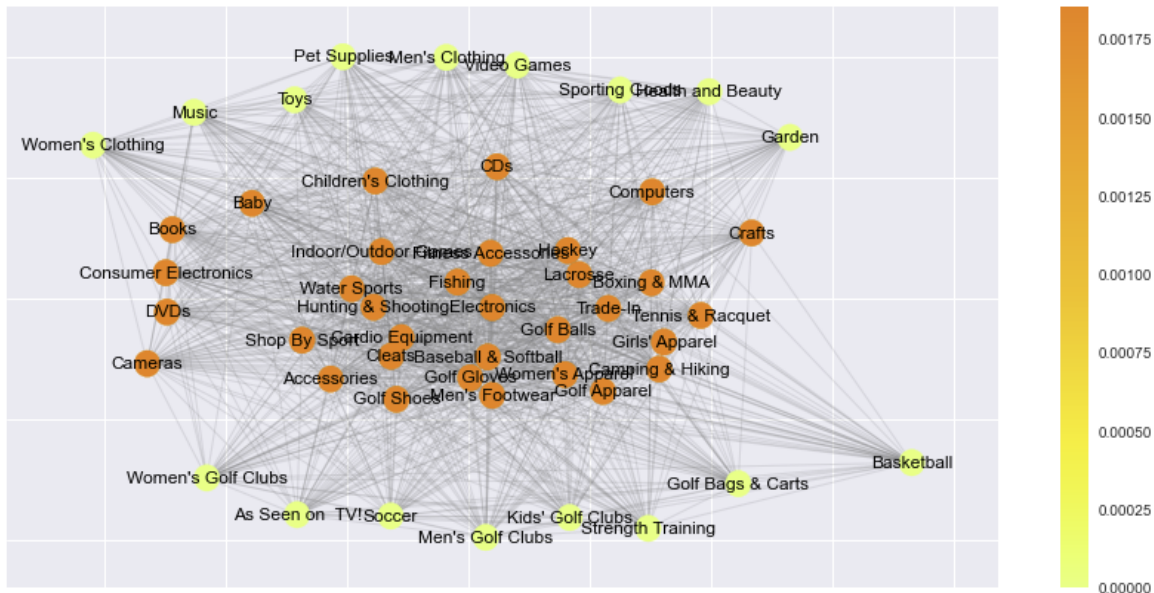
### 3. Betweenness Centrality

#### 3.1. Code chạy độ đo

```
1 betweenness = nx.betweenness centrality(G)
2 plot_centrality(betweenness)
3 print_centrality(betweenness, 'Betweenness Centrality').head()
```

Hình 15. Code chạy Betweenness Centrality

#### 3.2. Kết quả độ đo



```
[('Indoor/Outdoor Games', 0.0018552875695732828), ('Girls' Apparel', 0.0018552875695732828), ('Golf Apparel', 0.0018552875695732828), ('Tennis & Racquet', 0.0018552875695732828), ('Fitness Accessories', 0.0018552875695732828), ('Cleats', 0.0018552875695732828), ('Golf Balls', 0.0018552875695732828), ('Golf Shoes', 0.0018552875695732828), ('Children's Clothing', 0.0018552875695732828), ('Lacrosse', 0.0018552875695732828), ('Baby', 0.0018552875695732828), ('Fishing', 0.0018552875695732828), ('Books', 0.0018552875695732828), ('DVDs', 0.0018552875695732828), ('CDs', 0.0018552875695732828), ('Hockey', 0.0018552875695732828), ('Golf Gloves', 0.0018552875695732828), ('Crafts', 0.0018552875695732828), ('Accessories', 0.0018552875695732828), ('Hunting & Shooting', 0.0018552875695732828), ('Electronics', 0.0018552875695732828), ('Shop By Sport', 0.0018552875695732828), ('Women's Apparel', 0.0018552875695732828), ('Computers', 0.0018552875695732828), ('Cameras', 0.0018552875695732828), ('Consumer Electronics', 0.0018552875695732828), ('Camping & Hiking', 0.0018552875695732828), ('Men's Footwear', 0.0018552875695732828), ('Baseball & Softball', 0.0018552875695732828), ('Water Sports', 0.0018552875695732828), ('Trade-In', 0.0018552875695732828), ('Cardio Equipment', 0.0018552875695732828), ('Boxing & MMA', 0.0018552875695732828), ('Golf Bags & Carts', 0.0), ('Music', 0.0), ('Women's Golf Clubs', 0.0), ('Toys', 0.0), ('Men's Golf Clubs', 0.0), ('Video Games', 0.0), ('Sporting Goods', 0.0), ('Health and Beauty', 0.0), ('Pet Supplies', 0.0), ('Garden', 0.0), ('Strength Training', 0.0), ('Men's Clothing', 0.0), ('Women's Clothing', 0.0), ('Soccer', 0.0), ('Basketball', 0.0), ('Kids' Golf Clubs', 0.0), ('As Seen on TV!', 0.0)]
```

Hình 16. Kết quả Betweenness Centrality

#### 3.3. Nhận xét

- Betweenness Centrality là thước đo độ “ở giữa” trong đồ thị dựa trên ý tưởng về đường đi ngắn nhất. Betweenness Centrality của 1 node là 1 phần nhỏ của các đường đi ngắn nhất đi qua node đó.
- Các node có độ “ở giữa” cao hoạt động như 1 cầu nối giữa các phần biệt lập khác nhau của mạng. Các loại sản phẩm có Betweenness Centrality cao thể hiện nó có liên kết với

nhều loại sản phẩm khác, vì vậy các loại sản phẩm này được các khu vực cung cấp nhiều hơn các loại sản phẩm khác.

- Như trong hình 16, ta có thể thấy các node màu cam có Betweenness Centrality cao nhất, nghĩa là nó quan trọng trong việc kết nối giữa các node, các node muốn kết nối với nhau phải thông qua node đó. Bên cạnh đó, kết quả độ đo thực hiện trên Python và Gephi không có sự khác biệt.

Id	Betweenness Centrality		
Accessories	0.001855	Indoor/Outdoor Games	0.001855
Baby	0.001855	Lacrosse	0.001855
Baseball & Softball	0.001855	Men's Footwear	0.001855
Books	0.001855	Shop By Sport	0.001855
Boxing & MMA	0.001855	Tennis & Racquet	0.001855
Cameras	0.001855	Trade-In	0.001855
Camping & Hiking	0.001855	Water Sports	0.001855
Cardio Equipment	0.001855	Women's Apparel	0.001855
CDs	0.001855	Garden	0.0
Children's Clothing	0.001855	Health and Beauty	0.0
Cleats	0.001855	Men's Clothing	0.0
Computers	0.001855	Music	0.0
Consumer Electronics	0.001855	Pet Supplies	0.0
Crafts	0.001855	Sporting Goods	0.0
DVDs	0.001855	Toys	0.0
Electronics	0.001855	Video Games	0.0
Fishing	0.001855	Women's Clothing	0.0
Fitness Accessories	0.001855	As Seen on TV!	0.0
Girls' Apparel	0.001855	Basketball	0.0
Golf Apparel	0.001855	Golf Bags & Carts	0.0
Golf Balls	0.001855	Kids' Golf Clubs	0.0
Golf Gloves	0.001855	Men's Golf Clubs	0.0
Golf Shoes	0.001855	Soccer	0.0
Hockey	0.001855	Strength Training	0.0
Hunting & Shooting	0.001855	Women's Golf Clubs	0.0

Hình 17. Kết quả Betweenness Centrality trên Gephi

## 4. Closeness Centrality

### 4.1. Code chạy độ đo

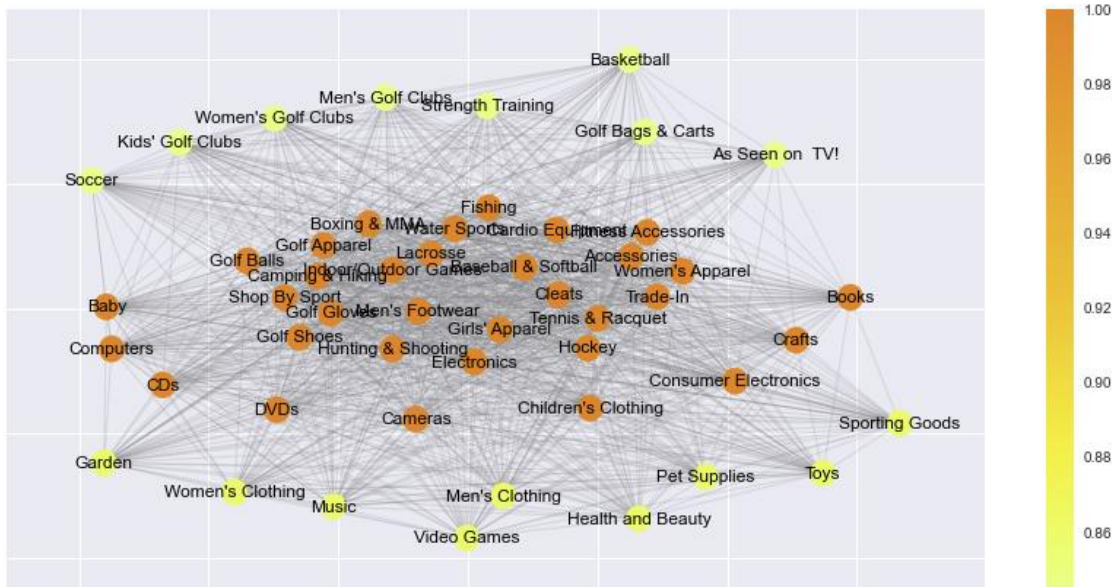
```

1 closeness = nx.closeness_centrality(G)
2 plot_centrality(closeness)
3 print_centrality(closeness, 'Closeness Centrality').head()

```

Hình 18. Code chạy Closeness Centrality

## 4.2. Kết quả độ đo



```
[('Indoor/Outdoor Games', 1.0), ('Golf Balls', 1.0), ('Golf Shoes', 1.0), ('Crafts', 1.0), ('Golf Gloves', 1.0), ('Tennis & Racquet', 1.0), ('Fitness Accessories', 1.0), ('Cleats', 1.0), ('Children's Clothing', 1.0), ('Golf Apparel', 1.0), ('Lacrosse', 1.0), ('Baby', 1.0), ('Fishing', 1.0), ('Books', 1.0), ('DVDs', 1.0), ('CDs', 1.0), ('Accessories', 1.0), ('Girls' Apparel', 1.0), ('Hockey', 1.0), ('Shop By Sport', 1.0), ('Computers', 1.0), ('Cameras', 1.0), ('Consumer Electronics', 1.0), ('Camping & Hiking', 1.0), ('Men's Footwear', 1.0), ('Baseball & Softball', 1.0), ('Hunting & Shooting', 1.0), ('Water Sports', 1.0), ('Trade-In', 1.0), ('Cardio Equipment', 1.0), ('Boxing & MMA', 1.0), ('Electronics', 1.0), ('Women's Apparel', 1.0), ('Music', 0.8596491228070176), ('Health and Beauty', 0.8596491228070176), ('Toys', 0.8596491228070176), ('Video Games', 0.8596491228070176), ('Sporting Goods', 0.8596491228070176), ('Pet Supplies', 0.8596491228070176), ('Garden', 0.8596491228070176), ('Men's Clothing', 0.8596491228070176), ('Women's Clothing', 0.8596491228070176), ('Strength Training', 0.8448275862068966), ('Golf Bags & Carts', 0.8448275862068966), ('Soccer', 0.8448275862068966), ('Women's Golf Clubs', 0.8448275862068966), ('Men's Golf Clubs', 0.8448275862068966), ('Basketball', 0.8448275862068966), ('Kids' Golf Clubs', 0.8448275862068966), ('As Seen on TV!', 0.8448275862068966)]
```

Hình 19. Kết quả Closeness Centrality

## 4.3. Nhận xét

- Closeness Centrality là 1 cách phát hiện các node có thể phân phối luồng đi 1 cách hiệu quả qua mạng, node quan trọng là node ở gần với các node khác.
- Kết quả của Closeness Centrality thể hiện top “độ gần” của các node đến tất cả các node trong mạng, đồng nghĩa với việc top các loại sản phẩm sẽ cần được cung cấp nhiều nhất. Bên cạnh đó, kết quả độ đo thực hiện trên Python và Gephi không có sự khác biệt.



Id	Closeness Centrality		
Cleats	1.0	Baseball & Softball	1.0
Children's Clothing	1.0	Camping & Hiking	1.0
Fitness Accessories	1.0	Cardio Equipment	1.0
Hockey	1.0	Shop By Sport	1.0
Tennis & Racquet	1.0	Golf Gloves	1.0
Consumer Electronics	1.0	Crafts	1.0
Golf Apparel	1.0	Electronics	1.0
Hunting & Shooting	1.0	Fishing	1.0
Golf Balls	1.0	Sporting Goods	0.859649
Books	1.0	Women's Clothing	0.859649
Water Sports	1.0	Video Games	0.859649
CDs	1.0	Toys	0.859649
Golf Shoes	1.0	Music	0.859649
Indoor/Outdoor Games	1.0	Men's Clothing	0.859649
DVDs	1.0	Health and Beauty	0.859649
Trade-In	1.0	Garden	0.859649
Baby	1.0	Pet Supplies	0.859649
Men's Footwear	1.0	As Seen on TV!	0.844828
Accessories	1.0	Basketball	0.844828
Cameras	1.0	Men's Golf Clubs	0.844828
Girls' Apparel	1.0	Golf Bags & Carts	0.844828
Computers	1.0	Soccer	0.844828
Women's Apparel	1.0	Strength Training	0.844828
Lacrosse	1.0	Women's Golf Clubs	0.844828
Boxing & MMA	1.0	Kids' Golf Clubs	0.844828

*Hình 20. Kết quả Closeness Centrality trên Gephi*

## 5. Eigenvector Centrality

### 5.1. Code chạy độ đo

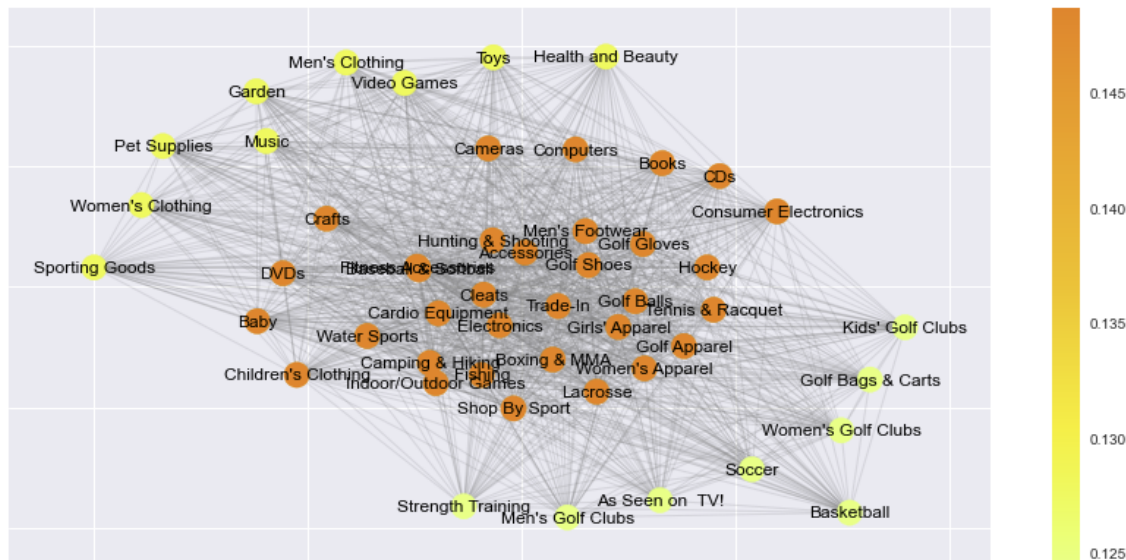
```

1 eigenvector = nx.eigenvector_centrality(G)
2 plot_centrality(eigenvector)
3 print_centrality(eigenvector, 'Eigenvector Centrality').head()

```

*Hình 21. Code chạy Eigenvector Centrality*

### 5.2. Kết quả độ đo



```
[('Indoor/Outdoor Games', 0.14867064721396608), ('Children's Clothing', 0.14867064721396608), ('Crafts', 0.14867064721396608), ('Golf Gloves', 0.14867064721396608), ('Tennis & Racquet', 0.14867064721396608), ('Fitness Accessories', 0.14867064721396608), ('Cleats', 0.14867064721396608), ('Golf Balls', 0.14867064721396608), ('Lacrosse', 0.14867064721396608), ('Girls' Apparel', 0.14867064721396608), ('Baby', 0.14867064721396608), ('Fishing', 0.14867064721396608), ('Books', 0.14867064721396608), ('DVDs', 0.14867064721396608), ('CDs', 0.14867064721396608), ('Hockey', 0.14867064721396608), ('Golf Shoes', 0.14867064721396608), ('Accessories', 0.14867064721396608), ('Golf Apparel', 0.14867064721396608), ('Water Sports', 0.14867064721396608), ('Shop By Sport', 0.14867064721396608), ('Women's Apparel', 0.14867064721396608), ('Electronics', 0.14867064721396608), ('Boxing & MMA', 0.14867064721396608), ('Cardio Equipment', 0.14867064721396608), ('Trade-In', 0.14867064721396608), ('Hunting & Shooting', 0.14867064721396608), ('Baseball & Softball', 0.14867064721396608), ('Men's Footwear', 0.14867064721396608), ('Camping & Hiking', 0.14867064721396608), ('Consumer Electronics', 0.14867064721396608), ('Cameras', 0.14867064721396608), ('Computers', 0.14867064721396608), ('Video Games', 0.12767929833125463), ('Music', 0.12767929833125463), ('Garden', 0.12767929833125463), ('Toys', 0.12767929833125463), ('Pet Supplies', 0.12767929833125463), ('Health and Beauty', 0.12767929833125463), ('Sporting Goods', 0.1276792983312546), ('Men's Clothing', 0.1276792983312546), ('Women's Clothing', 0.1276792983312546), ('Strength Training', 0.12444087284685672), ('Golf Bags & Carts', 0.12444087284685672), ('Soccer', 0.12444087284685672), ('Women's Golf Clubs', 0.12444087284685672), ('Men's Golf Clubs', 0.12444087284685672), ('Basketball', 0.12444087284685672), ('Kids' Golf Clubs', 0.12444087284685672), ('As Seen on TV!', 0.12444087284685672)]
```

Hình 22. Kết quả Eigenvector Centrality

### 5.3. Nhận xét

- Eigenvector Centrality đo lường tầm quan trọng của 1 node trong mạng đồng thời xem xét tầm quan trọng của các node lân cận. Các kết nối đến các node trung tâm của Eigenvector có điểm số cao đóng góp nhiều hơn vào điểm số chung so với các kết nối ngang bằng với các nút có điểm số thấp. Nói cách khác, 1 node có nhiều kết nối có thể có điểm Eigenvector thấp nếu tất cả các kết nối của nó đều có các node có điểm thấp.
- Loại sản phẩm có điểm Eigenvector cao khi nó được kết nối với nhiều node có điểm cao khác, đồng nghĩa với việc các loại sản phẩm có điểm Eigenvector càng cao sẽ càng được cung cấp bởi nhiều khu vực, vì các loại sản phẩm kết nối với loại sản phẩm đó cũng được cung cấp bởi nhiều khu vực.
- Như trong hình 22, ta có thể thấy các node màu cam có ảnh hưởng nhất mạng. Bên cạnh đó, kết quả độ đo thực hiện trên Python cũng có sự khác biệt so với Gephi.

Id	Eigenvector Centrality		
Hockey	1.0	Shop By Sport	1.0
Golf Apparel	1.0	Crafts	1.0
Books	1.0	Electronics	1.0
Water Sports	1.0	Consumer Electronics	1.0
Golf Shoes	1.0	Men's Footwear	1.0
Indoor/Outdoor Games	1.0	Cameras	1.0
Golf Gloves	1.0	Baseball & Softball	1.0
Fishing	1.0	Camping & Hiking	1.0
Cleats	1.0	Sporting Goods	0.858678
Children's Clothing	1.0	Women's Clothing	0.858678
Fitness Accessories	1.0	Video Games	0.858678
Tennis & Racquet	1.0	Toys	0.858678
Hunting & Shooting	1.0	Music	0.858678
Golf Balls	1.0	Men's Clothing	0.858678
CDs	1.0	Health and Beauty	0.858678
DVDs	1.0	Garden	0.858678
Trade-In	1.0	Pet Supplies	0.858678
Baby	1.0	As Seen on TV!	0.836945
Accessories	1.0	Basketball	0.836945
Girls' Apparel	1.0	Men's Golf Clubs	0.836945
Computers	1.0	Golf Bags & Carts	0.836945
Women's Apparel	1.0	Soccer	0.836945
Lacrosse	1.0	Strength Training	0.836945
Boxing & MMA	1.0	Women's Golf Clubs	0.836945
Cardio Equipment	1.0	Kids' Golf Clubs	0.836945

*Hình 23. Kết quả Eigenvector Centrality trên Gephi*

## 6. PageRank

### 6.1. Code chạy độ đo

```

1 pagerank = nx.pagerank(G)
2 plot_centrality(pagerank)
3 print_centrality(pagerank, 'PageRank').head()

```

*Hình 24. Code chạy PageRank*

### 6.2. Kết quả độ đo





Id	PageRank	Consumer Electronics	0.01461
Cleats	0.029094	Books	0.01461
Tennis & Racquet	0.029094	CDs	0.01461
Hunting & Shooting	0.029094	DVDs	0.01461
Golf Balls	0.029094	Baby	0.01461
Water Sports	0.029094	Computers	0.01461
Indoor/Outdoor Games	0.029094	Crafts	0.01461
Trade-In	0.029094	Cameras	0.012947
Men's Footwear	0.029094	As Seen on TV!	0.011339
Accessories	0.029094	Men's Golf Clubs	0.011339
Girls' Apparel	0.029094	Golf Bags & Carts	0.011339
Women's Apparel	0.029094	Soccer	0.011339
Lacrosse	0.029094	Strength Training	0.011339
Boxing & MMA	0.029094	Women's Golf Clubs	0.011339
Cardio Equipment	0.029094	Kids' Golf Clubs	0.011339
Shop By Sport	0.029094	Sporting Goods	0.009728
Golf Gloves	0.029094	Women's Clothing	0.009728
Electronics	0.029094	Video Games	0.009728
Fishing	0.029094	Toys	0.009728
Golf Apparel	0.029094	Music	0.009728
Camping & Hiking	0.029094	Men's Clothing	0.009728
Hockey	0.028388	Health and Beauty	0.009728
Golf Shoes	0.028388	Garden	0.009728
Baseball & Softball	0.028388	Pet Supplies	0.009728
Fitness Accessories	0.028388	Basketball	0.00782
Children's Clothing	0.01461		

Hình 26. Kết quả PageRank trên Gephi

## V. THUẬT TOÁN PHÁT HIỆN CỘNG ĐỒNG

### 1. Định nghĩa các hàm hiển thị kết quả phân cụm

- Hàm lấy đặc điểm chung của các node trong 1 cụm:

```

1 df_groupby_secondary = df.groupby([SECONDARY])[PRIMARY].apply(lambda x: sorted(set(x)))
2 df_groupby_secondary = df_groupby_secondary.reset_index()
3 # print(df_groupby_secondary.iloc[14][1])
4 df_groupby_secondary.head()

```

	Order	Region	Category Name
0	Canada	[Accessories, Baseball & Softball, Boxing & MM...	
1	Caribbean	[Accessories, As Seen on TV!, Baseball & Soft...	
2	Central Africa	[Accessories, Baseball & Softball, Boxing & MM...	
3	Central America	[Accessories, As Seen on TV!, Baseball & Soft...	
4	Central Asia	[Accessories, Boxing & MMA, Camping & Hiking, ...	

```

1 def get_cluster_common_values(cluster):
2     commons = []
3     for index, row in df_groupby_secondary.iterrows():
4         if set(cluster).issubset(row[PRIMARY]):
5             commons.append(row[SECONDARY])
6     return commons

```

Hình 27. Hàm lấy đặc điểm chung của 1 cụm

- Hàm in ra các cộng đồng đã phát hiện được cùng điểm chung của các node trong chúng:

```

1 def print_communities(node_groups):
2     print('Number of communities:', len(node_groups))
3     for index, cluster in enumerate(node_groups):
4         cluster = sorted(cluster)
5         common_values = sorted(get_cluster_common_values(cluster))
6
7         print(f'\nCluster {index}:')
8         print(f"- {len(cluster)} Nodes: {' '.join(cluster)}")
9         print(f"- {len(common_values)} Common values: {' '.join(common_values)}")

```

Hình 28. Hàm in ra các cộng đồng cùng điểm chung của các node trong chúng

- Hàm trực quan hóa các cộng đồng:

```

1 import matplotlib.cm as cm
2 def plot_communities(nodes, labels):
3     cmap = cm.get_cmap('autumn', max(labels) + 1)
4     pos = nx.spring_layout(G)
5
6     plt.figure(figsize=(15, 7))
7     nx.draw_networkx_nodes(G, pos, nodes, cmap=cmap, node_color=labels)
8     nx.draw_networkx_edges(G, pos, edge_color='grey', alpha=0.2)
9     nx.draw_networkx_labels(G, pos)
10    plt.show()

```

Hình 29. Hàm trực quan hóa các cộng đồng

## 2. Thuật toán Louvain

### 2.1. Code chạy thuật toán

```

1 import community.community_louvain as community_louvain
2 partition = community_louvain.best_partition(G)
3 louvain_node_groups = [[] for _ in set(partition.values())]
4
5 for node, cluster in sorted(partition.items()):
6     louvain_node_groups[cluster].append(node)

```

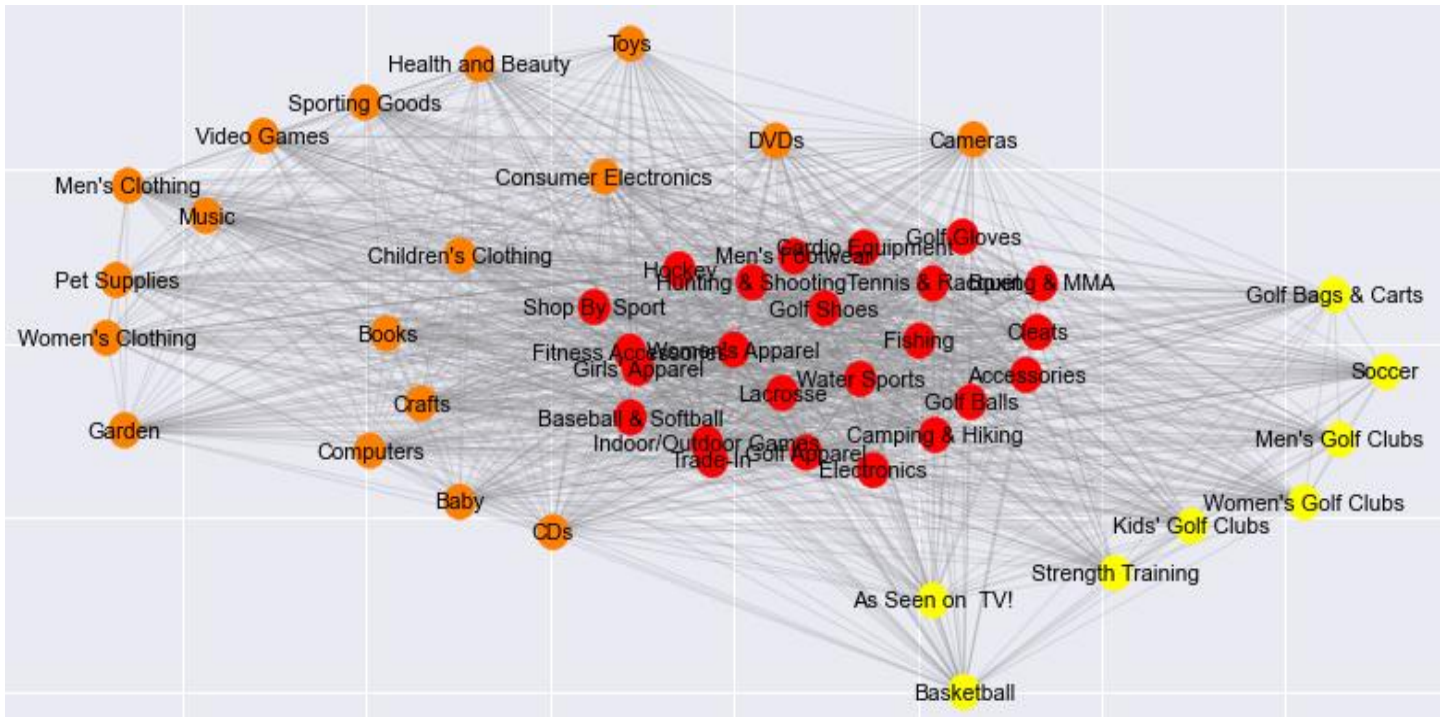
```

1 print_communities(louvain_node_groups)
2 plot_communities(partition.keys(), list(partition.values()))

```

Hình 30. Code chạy thuật toán Louvain

### 2.2. Đồ thị phân cụm



Hình 31. Đồ thị phân cụm sử dụng Louvain

⇒ Dễ thấy thuật toán rút ra được 3 cụm

### 2.3. Ý nghĩa các cụm

Number of communities: 3

Cluster 0:

- 24 Nodes: Accessories, Baseball & Softball, Boxing & MMA, Camping & Hiking, Cardio Equipment, Cleats, Electronics, Fishing, Fitness Accessories, Girls' Apparel, Golf Apparel, Golf Balls, Golf Gloves, Golf Shoes, Hockey, Hunting & Shooting, Indoor/Outdoor Games, Lacrosse, Men's Footwear, Shop By Sport, Tennis & Racquet, Trade-In, Water Sports, Women's Apparel  
- 22 Common values: Canada, Caribbean, Central Africa, Central America, East Africa, East of USA, Eastern Asia, Eastern Europe, North Africa, Northern Europe, Oceania, South America, South Asia, South of USA, Southeast Asia, Southern Africa, Southern Europe, US Center, West Africa, West Asia, West of USA, Western Europe

Cluster 1:

- 18 Nodes: Baby, Books, CDs, Cameras, Children's Clothing, Computers, Consumer Electronics, Crafts, DVDs, Garden, Health and Beauty, Men's Clothing, Music, Pet Supplies, Sporting Goods, Toys, Video Games, Women's Clothing  
- 3 Common values: Oceania, South Asia, Southeast Asia

Cluster 2:

- 8 Nodes: As Seen on TV!, Basketball, Golf Bags & Carts, Kids' Golf Clubs, Men's Golf Clubs, Soccer, Strength Training, Women's Golf Clubs  
- 3 Common values: Northern Europe, Southern Europe, Western Europe

Hình 32. Kết quả phân cụm sử dụng Louvain

#### ❖ Cụm thứ 0:

- **Gồm 24 Node:** Accessories, Baseball & Softball, Boxing & MMA, Camping & Hiking, Cardio Equipment, Cleats, Electronics, Fishing, Fitness Accessories, Girls' Apparel, Golf Apparel, Golf Balls, Golf Gloves, Golf Shoes, Hockey, Hunting &

Shooting, Indoor/Outdoor Games, Lacrosse, Men's Footwear, Shop By Sport, Tennis & Racquet, Trade-In, Water Sports, Women's Apparel.

- Ý nghĩa: những loại sản phẩm được cung cấp ở cả **22 khu vực** gồm Canada, Caribbean, Central Africa, Central America, East Africa, East of USA, Eastern Asia, Eastern Europe, North Africa, Northern Europe, Oceania, South America, South Asia, South of USA, Southeast Asia, Southern Africa, Southern Europe, US Center, West Africa, West Asia, West of USA, Western Europe.

❖ **Cụm thứ 1:**

- Gồm **18 Node**: Baby, Books, CDs, Cameras, Children's Clothing, Computers, Consumer Electronics, Crafts, DVDs, Garden, Health and Beauty, Men's Clothing, Music, Pet Supplies, Sporting Goods, Toys, Video Games, Women's Clothing.
- Ý nghĩa: những loại sản phẩm được cung cấp ở cả **3 khu vực** gồm Oceania, South Asia, Southeast Asia.

❖ **Cụm thứ 2:**

- Gồm **8 Node**: As Seen on TV!, Basketball, Golf Bags & Carts, Kids' Golf Clubs, Men's Golf Clubs, Soccer, Strength Training, Women's Golf Clubs.
- Ý nghĩa: những loại sản phẩm được cung cấp ở cả **3 khu vực** gồm Northern Europe, Southern Europe, Western Europe.

## 2.4. Nhận xét

- Nếu nhìn bằng mắt thường ta cũng dễ dàng nhận ra đồ thị có thể có 3 cụm. Louvain đã phát hiện được 3 cụm đó khá rõ rệt với nhau. Như vậy chứng tỏ thuật toán Louvain đã phân cụm khá tốt.

## 3. Thuật toán K-Means

### 3.1. Code chạy thuật toán

#### 3.1.1. Chuyển đổi đồ thị thành ma trận kề

```
1 from sklearn.cluster import KMeans
2 from scipy.spatial.distance import cdist
3 adj_matrix = nx.to_numpy_array(G)
```

*Hình 33. Chuyển đổi đồ thị thành ma trận kề làm đầu vào cho K-Means*

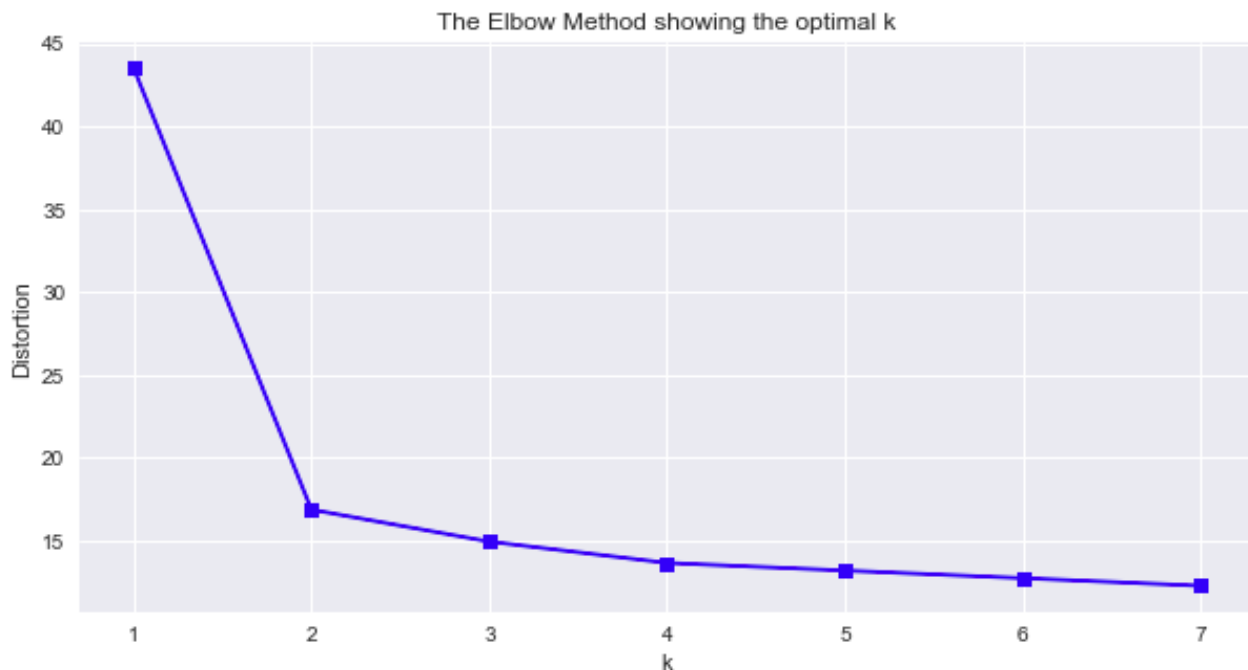
### 3.1.2. Chọn ra số cụm bằng Elbow Method

```
1 import numpy as np
2 distortions = []
3 K = range(1, 8)
4
5 for k in K:
6     kmean_model = KMeans(n_clusters=k)
7     kmean_model.fit(adj_matrix)
8     dist = sum(np.min(cdist(adj_matrix, kmean_model.cluster_centers_, 'euclidean'), axis=1))
9     distortions.append(dist / adj_matrix.shape[0])
```

Hình 34. Code triển khai Elbow Method cho thuật toán K-Means

- Tính khoảng cách từng phần tử trong ma trận kề (50 x 50) với các tâm cụm (k x 50).
- Ma trận (50 x k) tương ứng với 50 loại sản phẩm, với hàng là các khoảng cách giữa mỗi phần tử trong ma trận kề với các tâm cụm.
- Lấy tổng các khoảng cách của mỗi phần tử trong ma trận kề với tâm cụm mà tại đó có khoảng cách giữa 2 bên là nhỏ nhất / số lượng phần tử trong ma trận kề.
- Cuối cùng, tính tổng biến thiên khoảng cách nhỏ nhất của trong cụm.

```
1 plt.figure(figsize=(10, 5))
2 plt.plot(K, distortions, 'bs-')
3 plt.xlabel('k')
4 plt.ylabel('Distortion')
5 plt.title('The Elbow Method showing the optimal k')
6 plt.show()
```



Hình 35. Chọn ra số cụm k cho thuật toán K-Means bằng Elbow Method



### 3.1.3. Huấn luyện với số cụm đã chọn

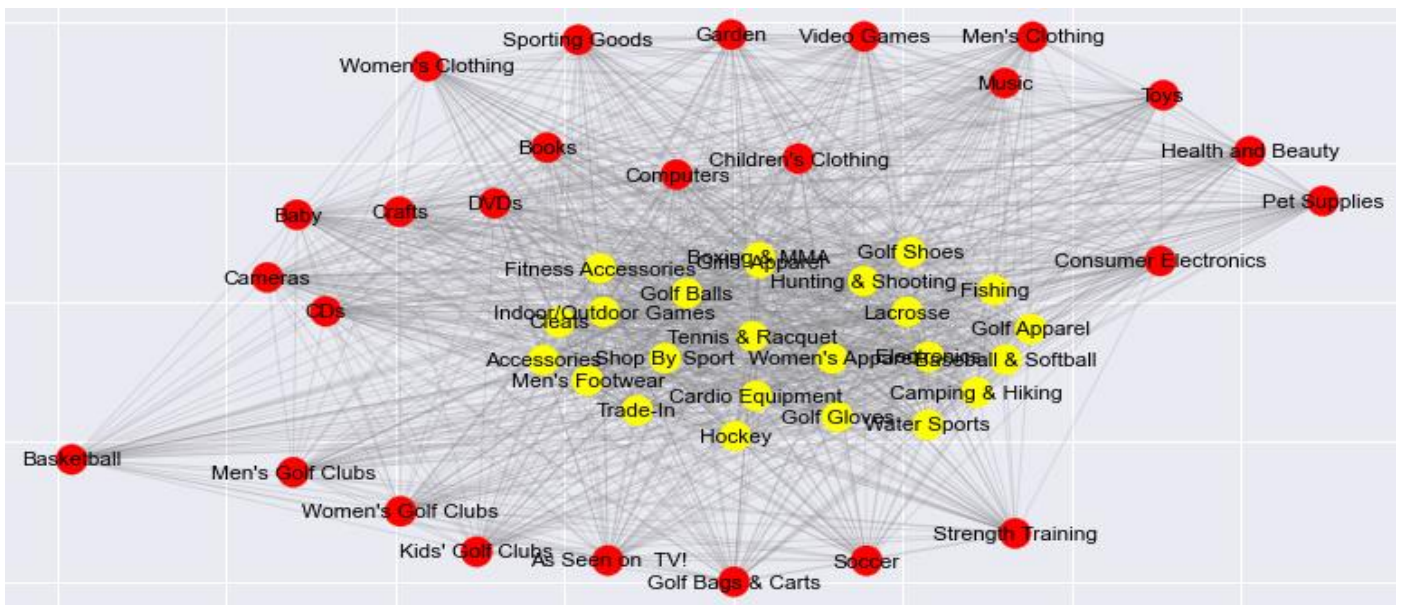
- Chọn số cụm k là 2 để huấn luyện cũng như sẽ là số lượng centroid sẽ tạo ra.
- Tiến hành huấn luyện và hiển thị kết quả phân cụm với đầu vào là ma trận kề đã có.

```
1 kmeans = KMeans(n_clusters=2)
2 kmeans.fit(adj_matrix)
3
4 kmeans_node_groups = [[] for _ in range(kmeans.n_clusters)]
5 for node, cluster in zip(G.nodes(), kmeans.labels_):
6     kmeans_node_groups[cluster].append(node)

1 print_communities(kmeans_node_groups)
2 plot_communities(G.nodes(), kmeans.labels_)
```

Hình 36. Chạy K-Means với đầu vào là ma trận kề cùng số cụm k đã chọn

### 3.2. Đồ thị phân cụm



Hình 37. Đồ thị phân cụm sử dụng K-Means

### 3.3. Ý nghĩa các cụm

Number of communities: 2

Cluster 0:

- 26 Nodes: As Seen on TV!, Baby, Basketball, Books, CDs, Cameras, Children's Clothing, Computers, Consumer Electronics, Crafts, DVDs, Garden, Golf Bags & Carts, Health and Beauty, Kids' Golf Clubs, Men's Clothing, Men's Golf Clubs, Music, Pet Supplies, Soccer, Sporting Goods, Strength Training, Toys, Video Games, Women's Clothing, Women's Golf Clubs  
- 0 Common values:

Cluster 1:

- 24 Nodes: Accessories, Baseball & Softball, Boxing & MMA, Camping & Hiking, Cardio Equipment, Cleats, Electronics, Fishing, Fitness Accessories, Girls' Apparel, Golf Apparel, Golf Balls, Golf Gloves, Golf Shoes, Hockey, Hunting & Shooting, Indoor/Outdoor Games, Lacrosse, Men's Footwear, Shop By Sport, Tennis & Racquet, Trade-In, Water Sports, Women's Apparel  
- 22 Common values: Canada, Caribbean, Central Africa, Central America, East Africa, East of USA, Eastern Asia, Eastern Europe, North Africa, Northern Europe, Oceania, South America, South Asia, South of USA, Southeast Asia, Southern Africa, Southern Europe, US Center, West Africa, West Asia, West of USA, Western Europe

Hình 38. Kết quả phân cụm sử dụng K-Means

❖ **Cụm thứ 0:**

- **Gồm 26 Node:** As Seen on TV!, Baby, Basketball, Books, CDs, Cameras, Children's Clothing, Computers, Consumer Electronics, Crafts, DVDs, Garden, Golf Bags & Carts, Health and Beauty, Kids' Golf Clubs, Men's Clothing, Men's Golf Clubs, Music, Pet Supplies, Soccer, Sporting Goods, Strength Training, Toys, Video Games, Women's Clothing, Women's Golf Clubs.
- Ý nghĩa: các node trong cụm này **không** có điểm chung gì đồng thời với tất cả các node khác. Đây có thể đơn giản chỉ là phần còn lại sau khi phát hiện được cụm 1.

❖ **Cụm thứ 1:**

- **Gồm 24 Node:** Accessories, Baseball & Softball, Boxing & MMA, Camping & Hiking, Cardio Equipment, Cleats, Electronics, Fishing, Fitness Accessories, Girls' Apparel, Golf Apparel, Golf Balls, Golf Gloves, Golf Shoes, Hockey, Hunting & Shooting, Indoor/Outdoor Games, Lacrosse, Men's Footwear, Shop By Sport, Tennis & Racquet, Trade-In, Water Sports, Women's Apparel.
- Ý nghĩa: những loại sản phẩm được cung cấp ở cả **22 khu vực** gồm Canada, Caribbean, Central Africa, Central America, East Africa, East of USA, Eastern Asia, Eastern Europe, North Africa, Northern Europe, Oceania, South America, South Asia, South of USA, Southeast Asia, Southern Africa, Southern Europe, US Center, West Africa, West Asia, West of USA, Western Europe.

### 3.4. Nhận xét

- Từ kết quả có thể thấy K-Means chỉ phân biệt được các node có phân bố dạng cụm lại với nhau tạo thành 1 cụm (màu vàng) nhưng không phân biệt được hay tìm được các đặc điểm chung của các node có phân bố dạng đường vòng (màu đỏ).
- Nguyên nhân có thể do khi sử dụng thuật toán K-Means thì mỗi điểm dữ liệu được gán 1 cách dứt khoát thuộc 1 trung tâm cụm. Ngoài ra, K-Means không tính đến covariance (hiệp phương sai) để thể hiện tính liên quan của dữ liệu.
- Nếu ta có 2 điểm cách đều nhau từ trung tâm cụm nhưng 1 điểm theo xu hướng này và điểm kia thì không, K-Means sẽ coi chúng là như nhau, vì K-Means sử dụng khoảng cách Euclide.

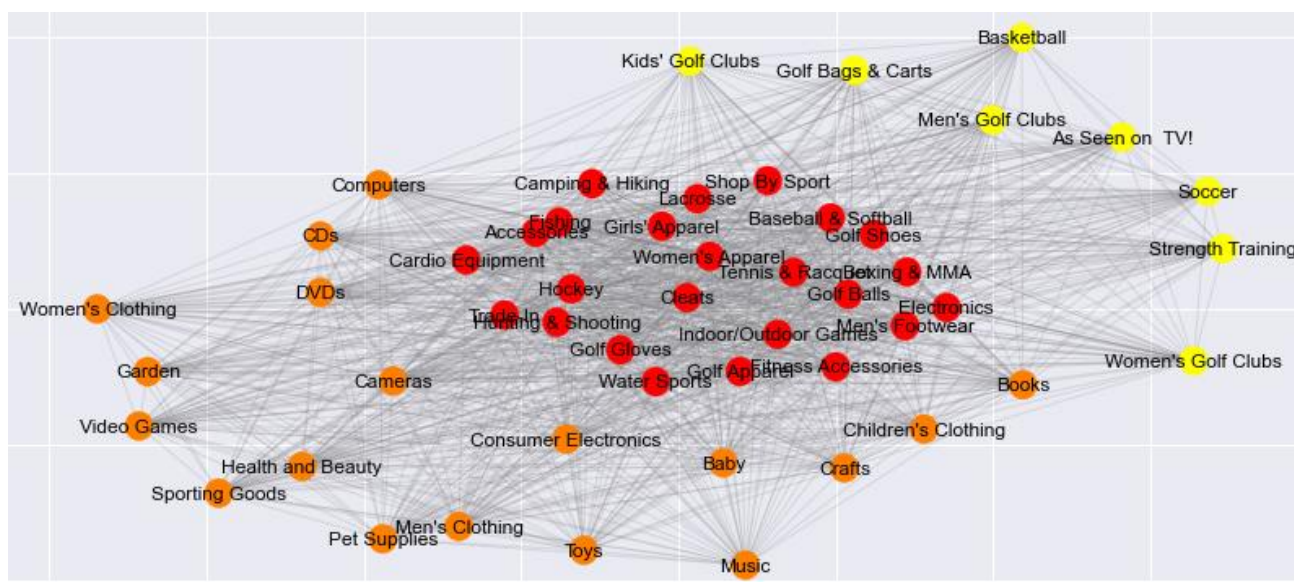
## 4. Gaussian Mixture Model

### 4.1. Code chạy thuật toán

```
1 from sklearn.mixture import GaussianMixture
2 n_clusters = 3
3 gmm = GaussianMixture(n_components=n_clusters)
4 gmm.fit(adj_matrix)
5
6 labels = gmm.predict(adj_matrix)
7 gmm_node_groups = [[] for _ in range(n_clusters)]
8 for node, cluster in zip(G.nodes(), labels):
9     gmm_node_groups[cluster].append(node)
10
11 print_communities(gmm_node_groups)
12 plot_communities(G.nodes(), labels)
```

Hình 39. Code chạy mô hình Gaussian Mixture

### 4.2. Đồ thị phân cụm



Hình 40. Đồ thị phân cụm sử dụng GMM

### 4.3. Ý nghĩa các cụm

Number of communities: 3

Cluster 0:

- 24 Nodes: Accessories, Baseball & Softball, Boxing & MMA, Camping & Hiking, Cardio Equipment, Cleats, Electronics, Fishing, Fitness Accessories, Girls' Apparel, Golf Apparel, Golf Balls, Golf Gloves, Golf Shoes, Hockey, Hunting & Shooting, Indoor/Outdoor Games, Lacrosse, Men's Footwear, Shop By Sport, Tennis & Racquet, Trade-In, Water Sports, Women's Apparel  
- 22 Common values: Canada, Caribbean, Central Africa, Central America, East Africa, East of USA, Eastern Asia, Eastern Europe, North Africa, Northern Europe, Oceania, South America, South Asia, South of USA, Southeast Asia, Southern Africa, Southern Europe, US Center, West Africa, West Asia, West of USA, Western Europe

Cluster 1:

- 18 Nodes: Baby, Books, CDs, Cameras, Children's Clothing, Computers, Consumer Electronics, Crafts, DVDs, Garden, Health and Beauty, Men's Clothing, Music, Pet Supplies, Sporting Goods, Toys, Video Games, Women's Clothing  
- 3 Common values: Oceania, South Asia, Southeast Asia

Cluster 2:

- 8 Nodes: As Seen on TV!, Basketball, Golf Bags & Carts, Kids' Golf Clubs, Men's Golf Clubs, Soccer, Strength Training, Women's Golf Clubs  
- 3 Common values: Northern Europe, Southern Europe, Western Europe

Hình 41. Kết quả phân cụm sử dụng GMM



❖ **Cụm thứ 0:**

- **Gồm 24 Node:** Accessories, Baseball & Softball, Boxing & MMA, Camping & Hiking, Cardio Equipment, Cleats, Electronics, Fishing, Fitness Accessories, Girls' Apparel, Golf Apparel, Golf Balls, Golf Gloves, Golf Shoes, Hockey, Hunting & Shooting, Indoor/Outdoor Games, Lacrosse, Men's Footwear, Shop By Sport, Tennis & Racquet, Trade-In, Water Sports, Women's Apparel.
- **Ý nghĩa:** những loại sản phẩm được cung cấp ở cả **22 khu vực** gồm Canada, Caribbean, Central Africa, Central America, East Africa, East of USA, Eastern Asia, Eastern Europe, North Africa, Northern Europe, Oceania, South America, South Asia, South of USA, Southeast Asia, Southern Africa, Southern Europe, US Center, West Africa, West Asia, West of USA, Western Europe.

❖ **Cụm thứ 1:**

- **Gồm 18 Node:** Baby, Books, CDs, Cameras, Children's Clothing, Computers, Consumer Electronics, Crafts, DVDs, Garden, Health and Beauty, Men's Clothing, Music, Pet Supplies, Sporting Goods, Toys, Video Games, Women's Clothing.
- **Ý nghĩa:** những loại sản phẩm được cung cấp ở cả **3 khu vực** gồm Oceania, South Asia, Southeast Asia.

❖ **Cụm thứ 2:**

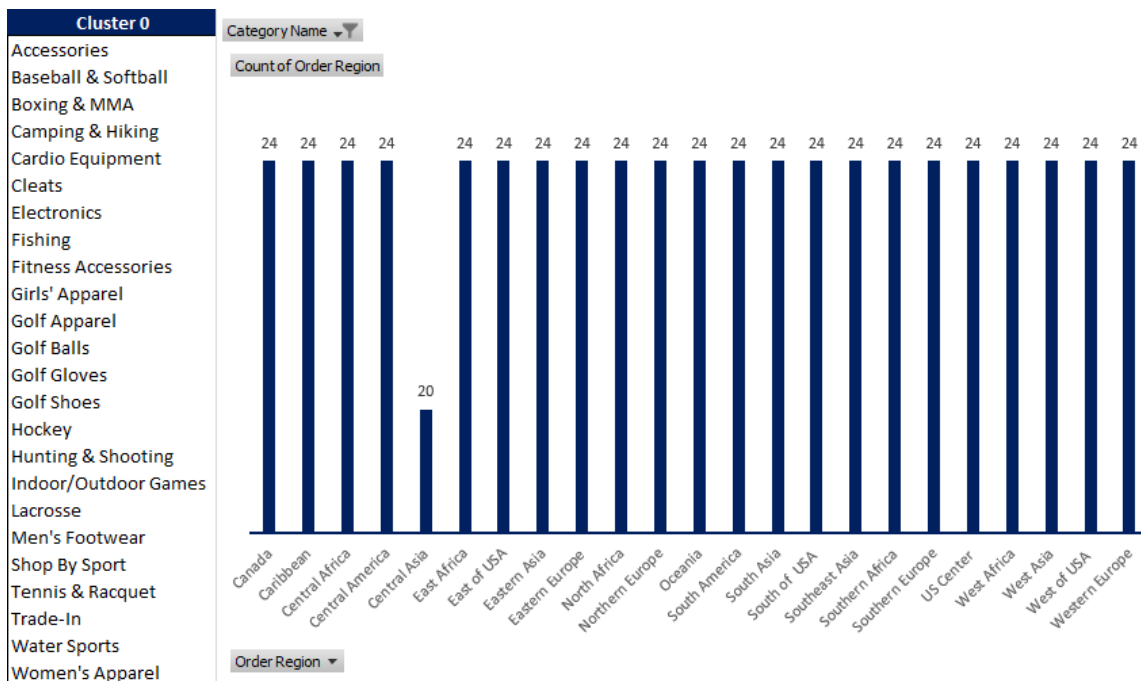
- **Gồm 8 Node:** As Seen on TV!, Basketball, Golf Bags & Carts, Kids' Golf Clubs, Men's Golf Clubs, Soccer, Strength Training, Women's Golf Clubs.
- **Ý nghĩa:** những loại sản phẩm được cung cấp ở cả **3 khu vực** gồm Northern Europe, Southern Europe, Western Europe.

#### 4.4. Nhận xét

- Gaussian Mixture Model là một mô hình xác suất giả định tất cả các điểm dữ liệu được tạo ra từ một hỗn hợp của 1 số hữu hạn của các phân phối Gaussian với các tham số không xác định và là phần mở rộng các ý tưởng đằng sau K-Means.
- Vì vậy Gaussian Mixture Model (GMM) đã khắc phục được các nhược điểm khi sử dụng thuật toán K-Means và đã cho ra kết quả phân cụm tương tự khi sử dụng thuật toán Louvain.

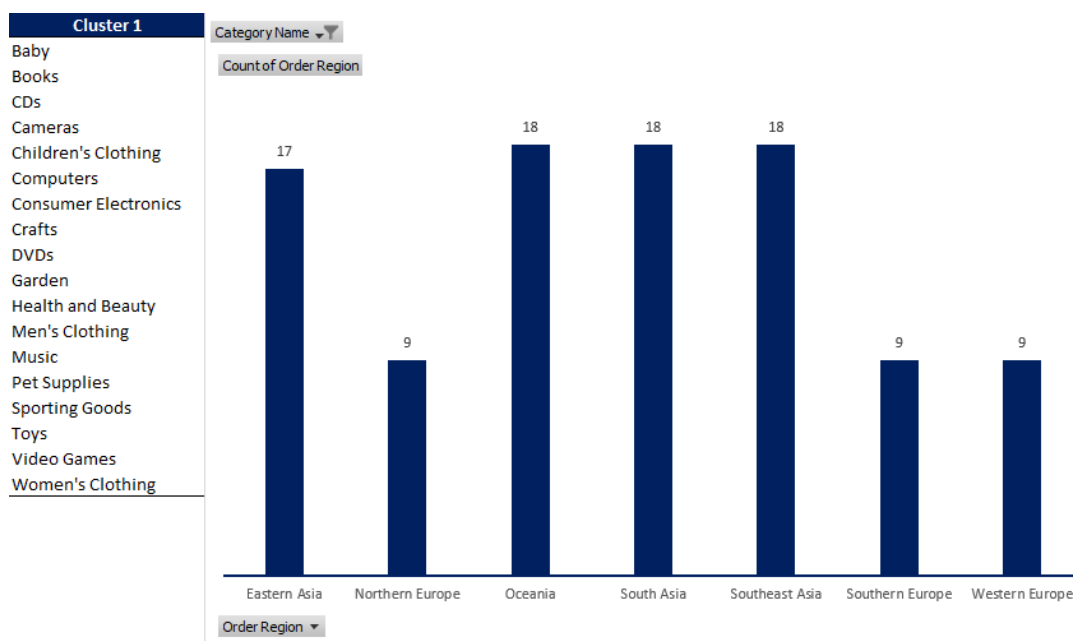
## 5. Trực quan hóa các điểm chung trong các cụm

- Do K-Means có vẻ không detect cụm tốt với dataset này, còn Louvain thì có kết quả phân cụm tương tự như khi sử dụng GMM nên ở đây ta sẽ trực quan hóa các khu vực đã cung cấp các loại sản phẩm của các cụm được detect khi sử dụng thuật toán Louvain.
- Cụm 0: các loại sản phẩm thường được cung cấp trên toàn cầu.



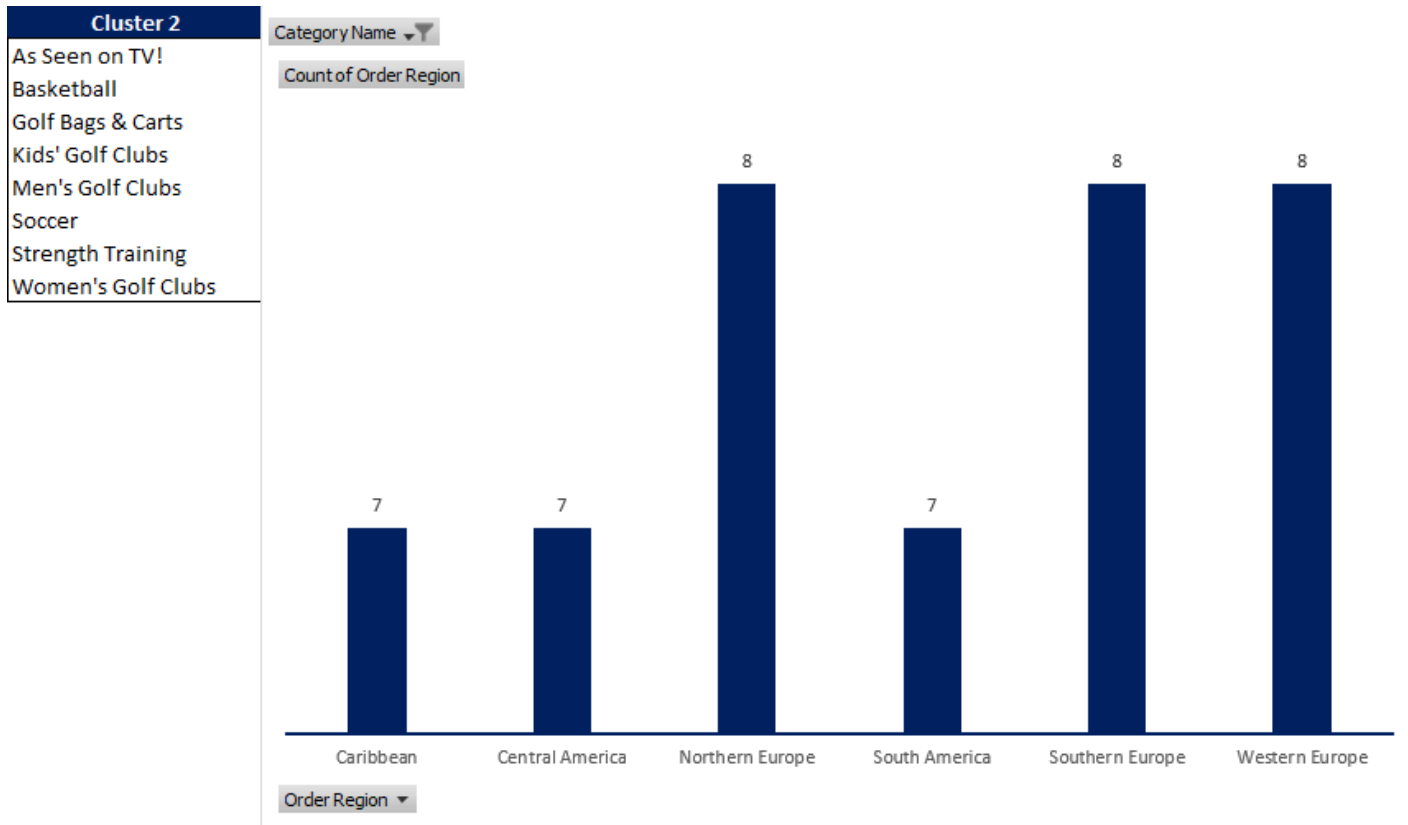
Hình 42. Các khu vực đã cung cấp các sản phẩm của cụm 0 khi dùng Louvain

- Cụm 1: các loại sản phẩm thường được cung cấp ở các khu vực Châu Á.



Hình 43. Các khu vực đã cung cấp các sản phẩm của cụm 1 khi dùng Louvain

- Cụm 2: các loại sản phẩm thường được cung cấp ở các khu vực Châu Âu.



Hình 44. Các khu vực đã cung cấp các sản phẩm của cụm 2 khi dùng Louvain

## VI. TÀI LIỆU THAM KHẢO

1. <https://youtube.com/playlist?list=PLoROMvodv4rPLKxIpqhjhPgDQy7imNkDn>
2. <https://www.coursera.org/learn/python-social-network-analysis>
3. <https://github.com/Geometrein/helsinki-city-bikes>