# Analysis of the Customers' Choice Networks: An Application on Amazon Books and CDs Data

Vladimir Batagelj,[1] Nataša Kejžar,[2] and Simona Korenjak-Černe[3]

**Abstract**

Customer's choice implies some kind of relations among products. Customers' choices of products induce a network among them. Analyses of such networks can offer interesting information for marketing. In the paper some network analysis approaches are proposed to analyze such data. Two large networks obtained in 2004 from Amazon Internet bookstore and CD-store are used to illustrate these approaches. All analyses were done with program Pajek.

## 1 Amazon networks construction

Amazon.com opened its virtual doors in July 1995 with the mission to transform book buying using the Internet into the fastest and easiest way. The Company's principal corporate offices are located in Seattle, Washington. It is one of the leading online shopping sites. It offers huge selection of products, including books, CDs, videos, DVDs, toys and games, electronics, kitchenware, computers etc.

In the Amazon network $\mathcal{N} = (\mathcal{V}, \mathcal{A})$ the **vertices** $\mathcal{V}$ are books/CDs; while the **arcs** $\mathcal{A}$ are determined for each product on the basis of the list of products (books/CDs) in its description under the title: *Customers who bought this book/CD also bought*. The vertex representing a described product is linked with an arc to every product listed in the list. Figure 1 presents an example of the construction of the neighborhood of the book *The Da Vinci Code*.

Using relatively simple program written in Python we 'harvested' the books network from June 16 till June 27, 2004; and the CDs network from July 7 till July 23, 2004. We harvested only the portion of each network reachable from the selected starting book: *Introducing Social Networks* by Michel Forse and Alain Degenne (0761956042) / starting CD: *The Pros and Cons of Hitchhiking* by Roger Waters (B0000025ZF).

---

[1]University of Ljubljana, Faculty of Mathematics and Physics, Department of Mathematics; Vladimir.Batagelj@fmf.uni-lj.si

[2]University of Ljubljana, Faculty of Social Sciences, Department of Informatics and Methodology; Natasa.Kejzar@fdv.uni-lj.si

[3]University of Ljubljana, Faculty of Economics, Department of Statistics; Simona.Cerne@ef.uni-lj.si
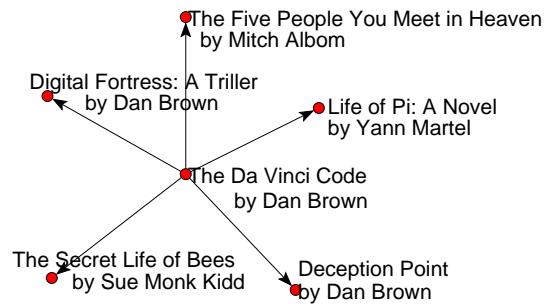
**Figure 1:** Dan Brown: *The Da Vinci Code.*

# 2   Description of obtained networks

The books network has 216737 vertices (books) and 982296 arcs. The CDs network has 79244 vertices (CDs) and 526271 arcs. By construction both networks have limited out-degree and are weakly connected. 178281 books have the maximum out-degree 5; and 55373 CDs have the maximum out-degree 8. Figure 2 shows the distributions of books/CDs by their in-degrees. The book with largest in-degree 553 is *Dan Brown: The Da Vinci Code*. The CDs with largest in-degree are *The Shins: Chutes Too Narrow* (706), and *Norah Jones: Feels Like Home* (675).

## 2.1   Strong components

The books network has 1787 nontrivial strong components, the largest of size 198808. The CDs network has 237 nontrivial strong components. The number of strong compo-nents strongly decreases with their size as can be noticed in Table 1. There are 130 strong components of size 2, 18 strong components of size 3 etc. But there is only 1 component of size 39, 1 component of size 84, 1 component of size 207, and 1 component of size 73928, and these are the only components with the sizes larger than 30.
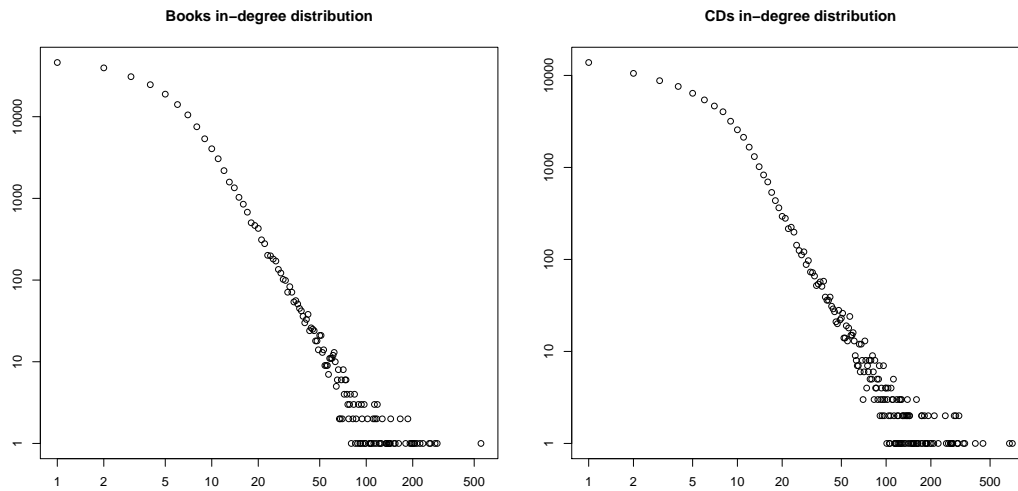
**Figure 2:** In-degree distributions.

**Table 1:** Distribution of strong components by their size in CDs network.

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| [1] | 3512 | 130 | 18 | 7 | 2 | 1 | 1 | 3 | 5 | 5 |
| [11] | 4 | 12 | 5 | 5 | 9 | 6 | 0 | 3 | 2 | 3 |
| [21] | 1 | 3 | 2 | 2 | 2 | 1 | 1 | 0 | 0 | 0 |
| [39] | 1 | | | | | | | | | |
| [84] | 1 | | | | | | | | | |
| [207] | 1 | | | | | | | | | |
| [73928] | 1 | | | | | | | | | |

## 2.2 Symmetrical subnetworks

To obtain undirected networks from the collected directed books/CDs networks two approaches were used:

- **network skeleton**: transform each arc into an edge and delete multiple edges;

- **symmetrical subnetwork**: replace pairs of opposite arcs with edges and delete the remaining arcs. Vertices with degree 0 are removed.

In the symmetrical subnetwork two vertices are linked by an edge if and only if there exist arcs in both directions. These vertices are worthy to be considered in detail to find the topic-purchase flow. Therefore we decided to analyze symmetrical subnetwork as an undirected network.

The symmetrical subnetwork on books has 186113 vertices, 218563 edges and 18967 components – the largest of size 59289. The second largest component is much smaller with 294 books and the third one includes 257 books.

The symmetrical subnetwork on CDs has 69708 vertices, 124834 edges and 3012 components – the largest of size 51692. The second largest component includes only 102 CDs and the third one 97.

# 3   Analysis

The goal of our analysis was to **identify the important parts** of both Amazon networks and to uncover their characteristics. To determine the most popular books and CDs the **hubs and authorities procedure** (Kleinberg, 1998) was used, and to **determine the main topics** in the networks the **islands approach** (Batagelj and Zaveršnik, 2004) was applied.

## 3.1   Hubs and authorities

A vertex is a **good hub** if it points to many good authorities, and it is a **good authority** if it is pointed by many good hubs.

To formalize this idea each vertex $v$ in the network gets two weights $x_v$ and $y_v$. The corresponding vectors $x$ and $y$ are related by equations $x = W^T y$ and $y = W x$, where $W = [w_{uv}]$ is the weight matrix of the network. It can be proved that $x$ and $y$ are the principal eigenvectors of matrices $W^T W$ and $W W^T$ (Kleinberg, 1998).

In the Amazon books/CD network out-degree is truncated, therefore hubs cannot reach values as large as authorities. Figure 3 shows hubs (white) and authorities (gray) around the main authority *The Da Vinci Code* in the books network. Black vertices are both good hubs and good authorities. Figure 4 shows five hubs and five authorities with the largest weights around the main authority *The Shins: Chutes Too Narrow* in the CDs network. The sizes of the vertices correspond to their weights.
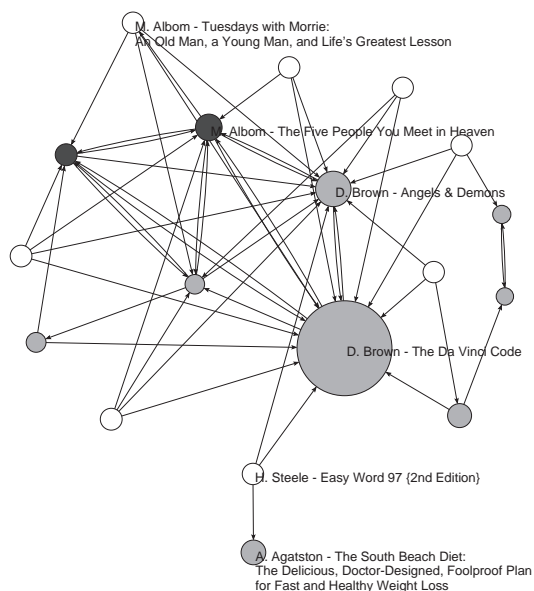


**Figure 3:** Hubs and authorities around *The Da Vinci Code* in books network.
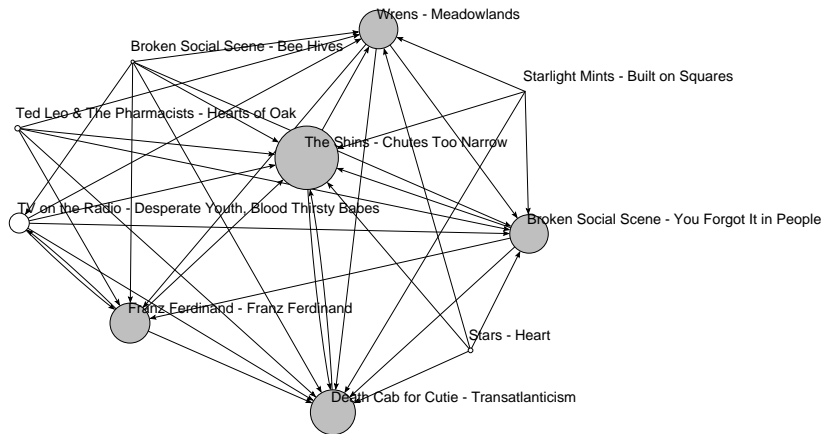
**Figure 4:** Hubs and authorities around *The Shins: Chutes Too Narrow* in CDs network.

## 3.2 Islands

Islands are connected parts of a network containing locally the most important vertices/ lines with respect to a given property/weight. Formally:

For a given network $\mathcal{N} = (\mathcal{V}, \mathcal{L}, p)$ with a vertex property $p$ a **regular vertex island** is a cluster $\mathcal{C}$ of vertices for which the corresponding induced subgraph is connected, and the property values of vertices in the cluster $\mathcal{C}$ are larger than the values of vertices in the cluster's neighborhood $N(\mathcal{C})$:

$$\max_{u \in N(\mathcal{C})} p(u) < \min_{v \in \mathcal{C}} p(v)$$

The set of vertices is a **local vertex peak**, if it is a regular vertex island and all of its vertices have the same value. Vertex island with a single local vertex peak is called a **simple vertex island**.

For a given network $\mathcal{N} = (\mathcal{V}, \mathcal{L}, w)$ with a line weight $w$ a **regular line island** is determined by a cluster $\mathcal{C}$ of vertices in which exists a spanning tree $\mathcal{T}$ on $\mathcal{C}$ such that the weights of lines with exactly one endpoint in the cluster are smaller than the weights of lines of the tree:

$$\max_{e(u:v) \in \mathcal{L}:(u \in \mathcal{C} \wedge v \notin \mathcal{C})} w(e) < \min_{e \in \mathcal{L}(\mathcal{T})} w(e)$$

where $e(u : v)$ denotes a line $e$ linking vertex $u$ with vertex $v$.

Similarly as simple vertex island we define simple line island as an island with only one peak.

We analyzed Amazon networks using several vertex properties and line weights. Here we present some results obtained with **clustering coefficient** (Watts and Strogatz, 1998) as a vertex property on CDs network and **triangular count** as a line weight on books network, although all presented approaches can be used on both networks.

### 3.2.1 Vertex property – clustering coefficient

The clustering coefficient measures the local density of a network in given vertex. It is defined as a proportion of links between the vertices within the neighborhood of the vertex

by the number of all possible links in its neighborhood. Let $\deg(v)$ denote the degree of vertex $v$, $|\mathcal{L}(G_1(v))|$ the number of lines among vertices ($\mathcal{L} = \mathcal{A}$ for arcs and $\mathcal{L} = \mathcal{E}$ for edges) in 1-neighborhood of vertex $v$, and $\Delta$ the maximum vertex degree in a network.

**Clustering coefficient** $CC_1(v)$ of vertex $v$ is defined as follows:

$$
\begin{aligned}
\text{for a } \textit{directed} \text{ network:} \quad CC_1(v) &= \frac{|\mathcal{A}(G_1(v))|}{\deg(v) \cdot (\deg(v) - 1)} \\
\text{for an } \textit{undirected} \text{ network:} \quad CC_1(v) &= \frac{2|\mathcal{E}(G_1(v))|}{\deg(v) \cdot (\deg(v) - 1)}
\end{aligned}
$$

The problem with the clustering coefficient $CC_1$ is that it has high values on vertices of small degree. To neutralize this effect we propose to use in data analytic tasks the **corrected clustering coefficient**:

$$
CC_1'(v) = \frac{\deg(v)}{\Delta} CC_1(v)
$$

### Directed CDs network

In the directed network on CDs are $4415$ simple vertex islands for $p(v) = CC_1'(v)$. The distribution of vertex islands by their size (number of vertices) is presented in the Table 2.

**Table 2:** Distribution of vertex islands by their size in the directed CDs network.

| | | | | | | | | | | |
|------|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| [1]  | 1625 | 576 | 356 | 273 | 221 | 154 | 176 | 152 | 158 | 144 |
| [11] | 128  | 105 | 87  | 72  | 38  | 34  | 21  | 19  | 12  | 19  |
| [21] | 13   | 8   | 5   | 2   | 5   | 1   | 2   | 2   | 0   | 0   |
| [31] | 2    | 1   | 2   | 1   | 0   | 0   | 1   | 0   | 0   | 0   |

The island of the maximal size contains 37 CDs. The analysis of CDs within each island shows that they are either from the same author or they belong to the same type of music. For example the seven largest vertex islands based on corrected clustering coefficient in the directed network on CDs, can be identified as:

| #of CDs | Description |
|---------|-------------|
| 37 | Barbra Streisand |
| 34 | Modern Scandinavian folk music |
| 33 | Soul and blues |
| 33 | Hed Kandi's house and disco music |
| 32 | Progressive rock music |
| 31 | Hip-hop |
| 31 | Julio Iglesias |

### Undirected CDs network

In the undirected network (symmetrical subnetwork) on CDs are $8302$ simple vertex islands based on the corrected clustering coefficient $p(v) = CC_1'(v)$. The distribution of vertex islands by their size is presented in the Table 3.

**Table 3:** Distribution of vertex islands by the size in the undirected CDs network.

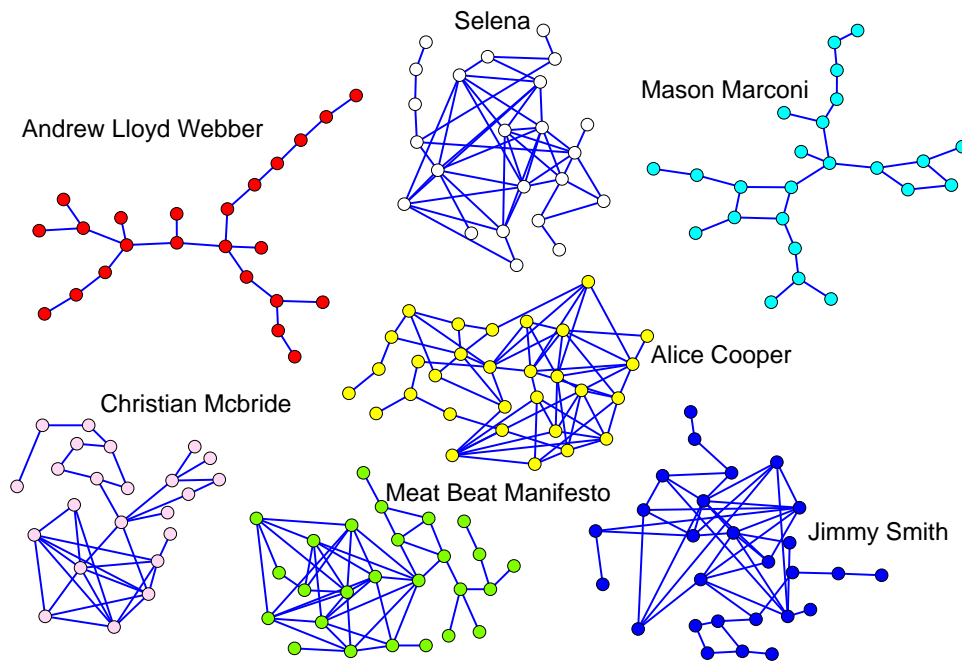| [1] | 1229 | 2158 | 1237 | 769 | 520 | 441 | 380 | 360 | 332 | 296 |
|-----|------|------|------|-----|-----|-----|-----|-----|-----|-----|
| [11] | 184 | 141 | 76 | 61 | 35 | 24 | 12 | 15 | 9 | 4 |
| [21] | 5 | 7 | 3 | 1 | 0 | 1 | 1 | 0 | 0 | 1 |



**Figure 5:** Seven largest vertex islands in the undirected network on CDs.

Also in these islands the CDs inside each island are either from the same author or they belong to the same type of music. The island with the maximal size contains 30 CDs. The seven largest vertex islands based on the corrected clustering coefficient in the undirected network on CDs, shown in the Figure 5, can be identified as:

| #of CDs | Author(s) or type of music |
|---------|----------------------------|
| 30 | Alice Cooper |
| 27 | Jimmy Smith, Eddie Harries, Lou Rawls |
| 26 | Meat Beat Manifesto, L.S.G. |
| 24 | Erotic Fantasy (Mason Marconi, unknown artists) |
| 23 | Musicals (Andrew Lloyd Webber and others) |
| 23 | Jazz (Christian Mcbride, Nicholas Payton, Mark Whitfield) |
| 23 | Selena |

The vertex island that includes CDs with songs from well known **musicals** such as *Sunset Boulevard* and *Starlight Express* from Andrew Lloyd Webber, *Cabaret* and

*Chicago* from John Kander, and *My One and Only* from George and Ira Gershwin is presented in detail separately in the Figure 6.

Andrew Lloyd Webber, Alan Ayckbourn: By Jeeves (1996, London)
Andrew Lloyd Webber, Jim Steinman: Whistle Down The Wind (1998, London)
Andrew Lloyd Webber: Sunset Boulevard (1994, Los Angeles)
Andrew Lloyd Webber: Sunset Boulevard (1993, London)
Andrew Lloyd Webber: Aspects Of Love (1989, London)
Andrew Lloyd Webber: Dance (1982, London)
Leonard Bernstein (Composer), et al: Candide (1956, Broadway)
Andrew Lloyd Webber, Richard Stilgoe: Starlight Express (1984, London)
Andrew Lloyd Webber, Richard Stilgoe: The New Starlight Express (1992, London)
Jerry Bock, et al: She Loves Me (1963, Broadway)
Andrew Lloyd Webber: The Songs (Broadway)
Jerry Bock (Composer), et al: Fiorello! (1959, Broadway)
Andrew Lloyd Webber: The Beautiful Game (2000, London)
Edith Adams (Composer), et al: Li'l Abner (1956, Broadway)
Richard Adler, Jerry Ross: The Pajama Game (1954, Broadway)
Jerry Herman: Mabel (1974, Broadway)
Cy Coleman (Composer), et al: Barnum (1980, Broadway)
Noel Gay, et al: Me And My Girl (1986, Broadway)
George Gershwin, Ira Gershwin: My One And Only (1983, Broadway)
Harry Warren, Al Dubin: 42nd Street (1980, Broadway)
John Kander: Chicago - A Musical Vaudeville (1975, Broadway)
John Kander: Cabaret (Broadway)
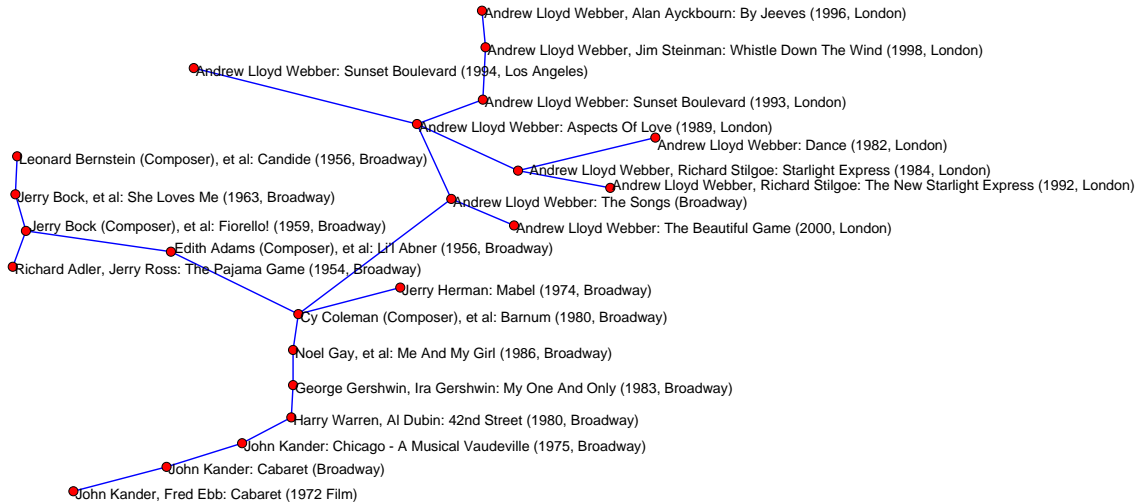John Kander, Fred Ebb: Cabaret (1972 Film)

**Figure 6:** Vertex island: Musicals.

Not only the size of the island, but also the weight of vertices inside the island is worth considering, because it contains more detailed information about the importance of the vertex. In the undirected CDs network 9 vertices have the highest weight (equal to 1) and all of them are in the same vertex island presented in the Figure 7. They form a clique of order 9. The island can be described as 'Fred Astaire and Ella Fitzgerald island'. Six CDs have the second largest weight. All of them are in the same vertex island with 4 other CDs that are related with the Beatles. This island is presented in the Figure 8.
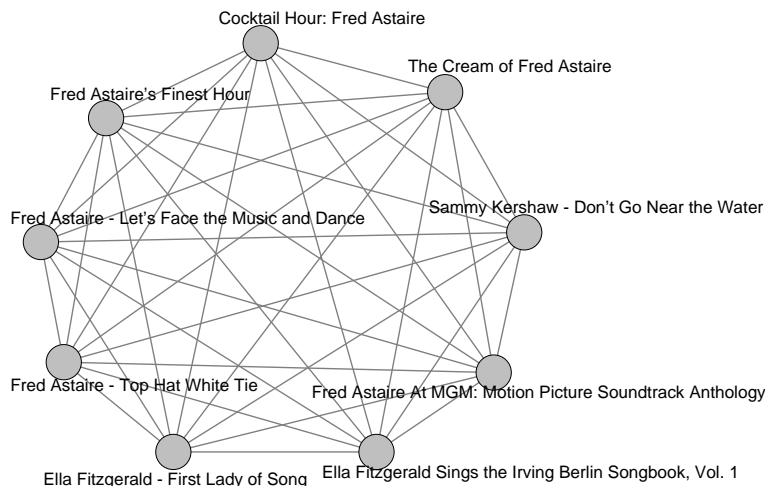
Cocktail Hour: Fred Astaire
The Cream of Fred Astaire
Fred Astaire's Finest Hour
Sammy Kershaw - Don't Go Near the Water
Fred Astaire - Let's Face the Music and Dance
Fred Astaire - Top Hat White Tie
Fred Astaire At MGM: Motion Picture Soundtrack Anthology
Ella Fitzgerald - First Lady of Song
Ella Fitzgerald Sings the Irving Berlin Songbook, Vol. 1

**Figure 7:** Vertex island with the largest weights.
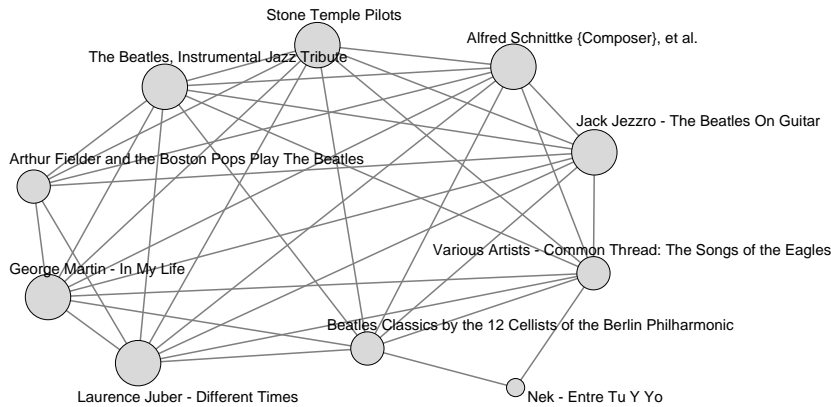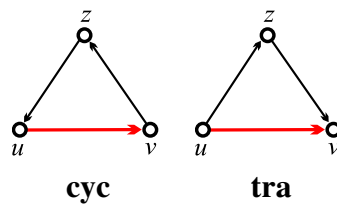
**Figure 8:** Vertex island with the second largest weights.

### 3.2.2 Line weight – triangle count

**Triangular weight** $w_3(e)$ of a line $e$ counts the number of triangles to which the line belongs. If the line $e$ belongs to a $k$-clique then $w_3(e) \geq k - 2$. Therefore the use of triangle count enables efficient detection of dense parts of networks. This approach can be used on:

- *undirected networks*, to get edge weights

- *directed networks*, to get arc weights

In the directed network two basic types of triangles can be observed: cyclic and transitive.



#### Books network as directed network

Figure 9 presents the distributions of the size of line islands in the directed networks of books and CDs considering the number of cyclic triangles as arcs' weights.

Examining the books included in the islands, the same or similar topics can be noticed inside the same island. Figure 10 shows line (arc) islands with at least 25 books and the main topics of the books inside them. Two of these islands are presented separately: the island including novels written by the **same author** Catherine Cookson in Figure 11, and the island including books about the **same topic** precious stones in Figure 12.

#### Books network as undirected network

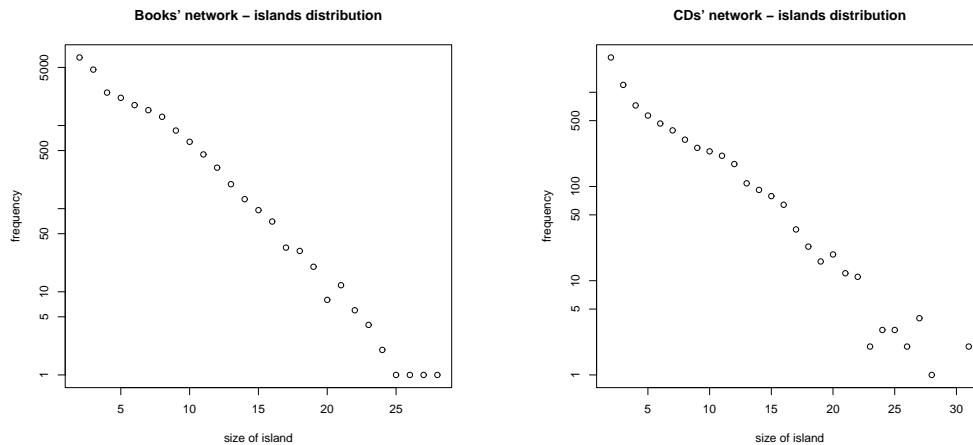Figure 13 shows line (edge) islands with at least 35 books in the symmetric subnetwork

**Figure 9:** Simple line (arc) islands distribution by their size. Note the log-linear scale of the graphs, which shows the sharp drop in frequency of larger line islands.
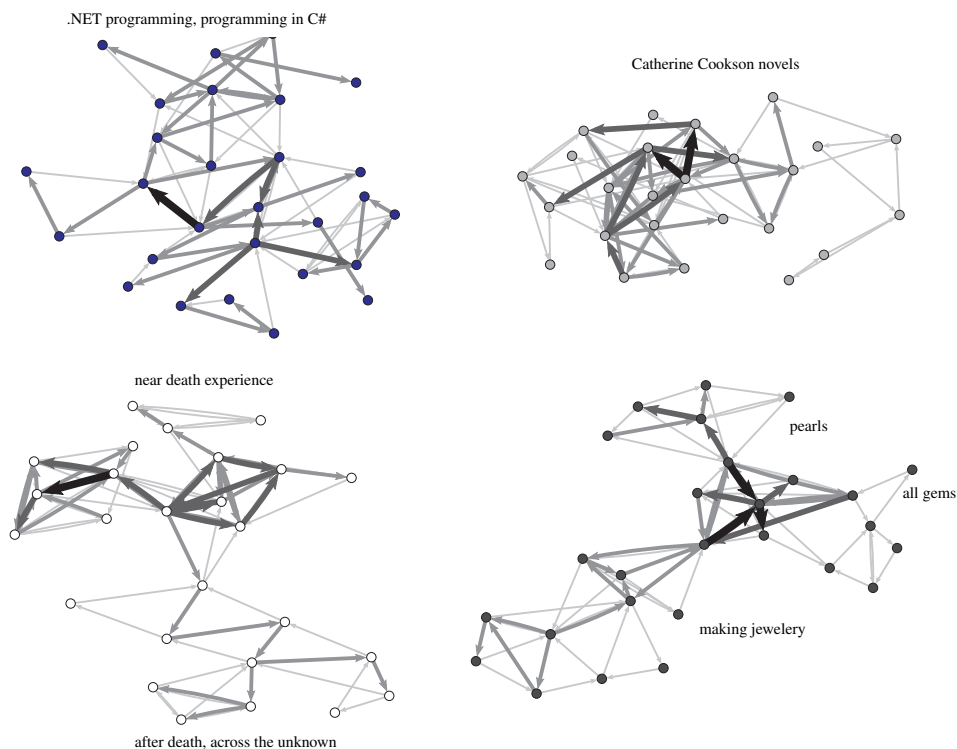


**Figure 10:** Arc islands with at least 25 vertices.

of the books network. Chaining is observed. We conjecture that the chaining is the 'backbone' of the topic, because two vertices in this network are connected only if there exist arcs both ways between them. These vertices therefore have to be really important for topic-purchase flow. Two of these islands are presented separately: the island including

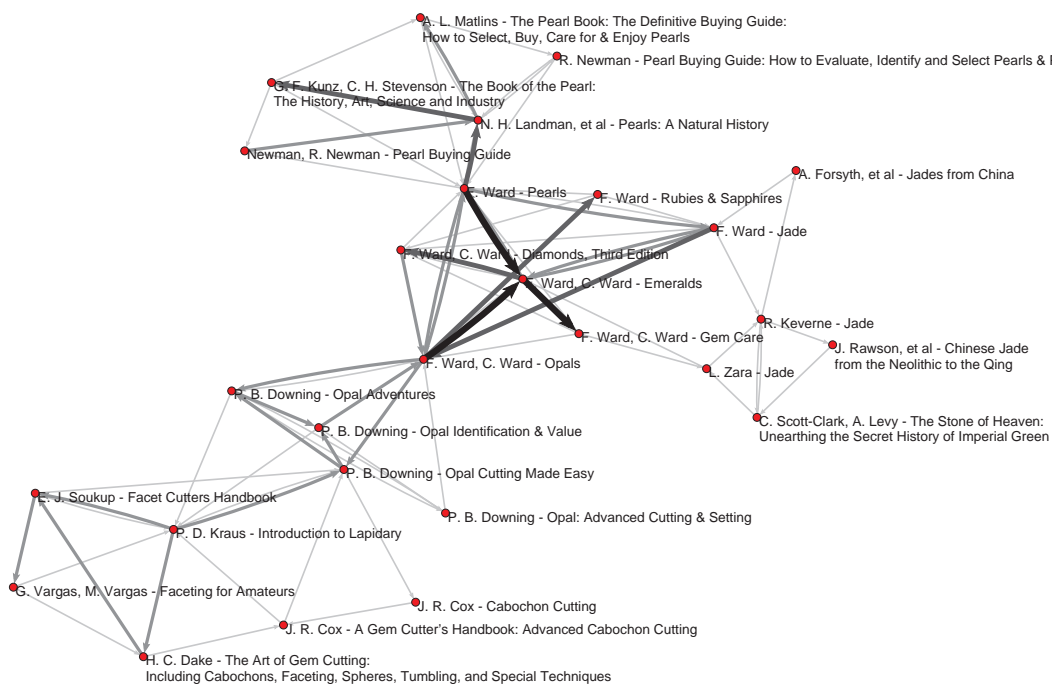**Figure 11:** Island of Catherine Cookson novels.



**Figure 12:** Island of precious stones.

books for children in Figure 14, and the island including books for students of literature in Figure 15.

Since the number of vertices in the island is not the only characteristic that has to be considered in the analysis, we inspected also the weights on lines. They range from 0 to 4, which implies that this network is relatively sparsely connected (there are at most 4 triangles above each line, and they represent only 0.5 % of all the edges). There are 73 islands (of 2 or more vertices) with at least one line weight 4, they are mostly of size 2 or 3 (68 of them) and 5 islands with 6 vertices, which are complete graphs. Four of

them consist of books of the same author(s): (1) Will Durant and Ariel Durant (history books), (2) Peter J. D'Adamo and Catherine Whitney (books about food and blood type), (3) Immanuel Velikovsky (history and legend books) and (4) Rachel Rubin Wolf (a book series). The last one of the islands consists of books of the Silva mind control method, where only 2 books are written by Silva himself.
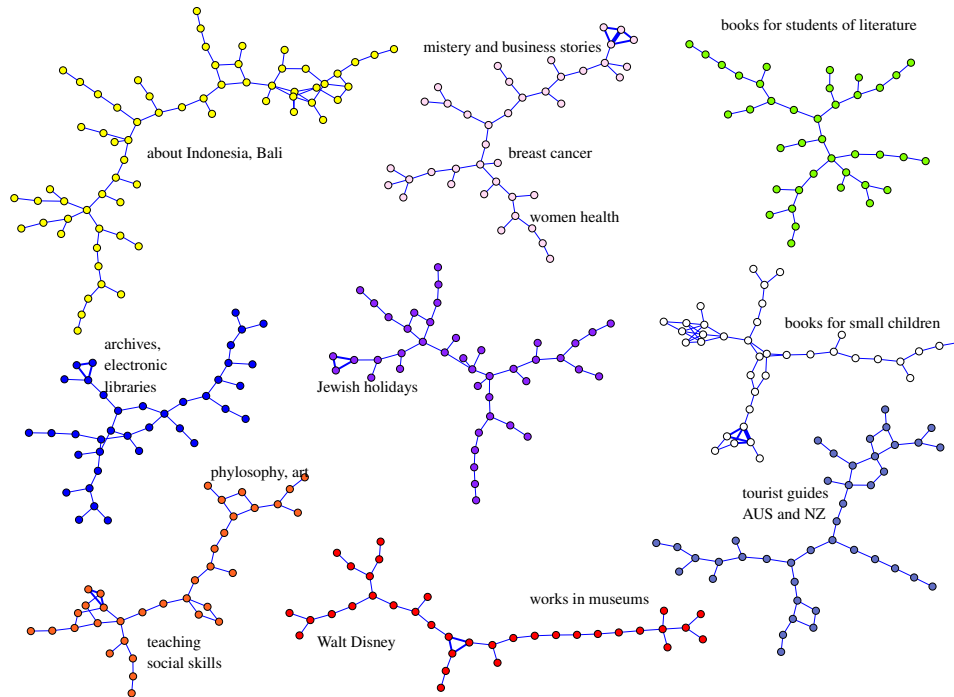


**Figure 13:** Edge islands with more than 35 vertices.

# 4 Discussion

Like well known data mining techniques (Berry and Linoff, 2004) used to extract information from large amounts of data, approaches from network analysis can show many interesting connections among products. In the paper some of these approaches have been shown for customers' choice networks constructed from the Amazon Internet bookstore, based on the list of products (books/CDs) under the title: *Customers who bought this book/CD also bought*.

With the hubs and authorities procedure (Kleinberg, 1998) the most important products were identified. By determining different types of islands (Batagelj and Zaveršnik, 2004) the important topics of books/CDs were detected. For the weights of vertices (books/CDs) the corrected clustering coefficient as a measure of relative local density in a given vertex was chosen. For line weights the number of triangles was chosen which also enables efficient detection of dense parts of networks. Due to large amount of results we presented only some of them to show the main possibilities of such analysis. To obtain the insight into the logic of groups formation, their characteristics have to be analyzed.
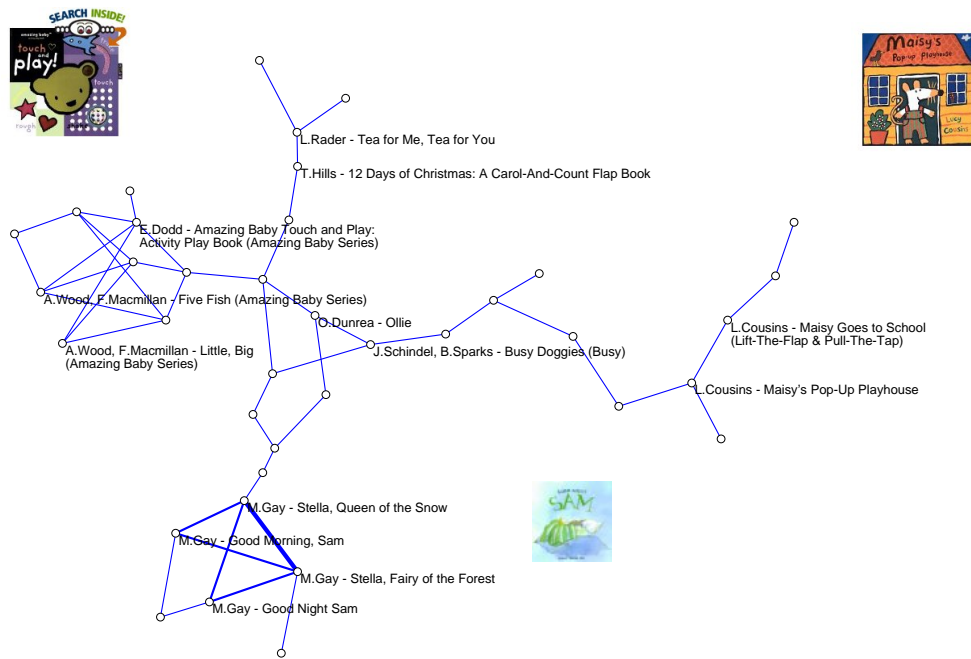
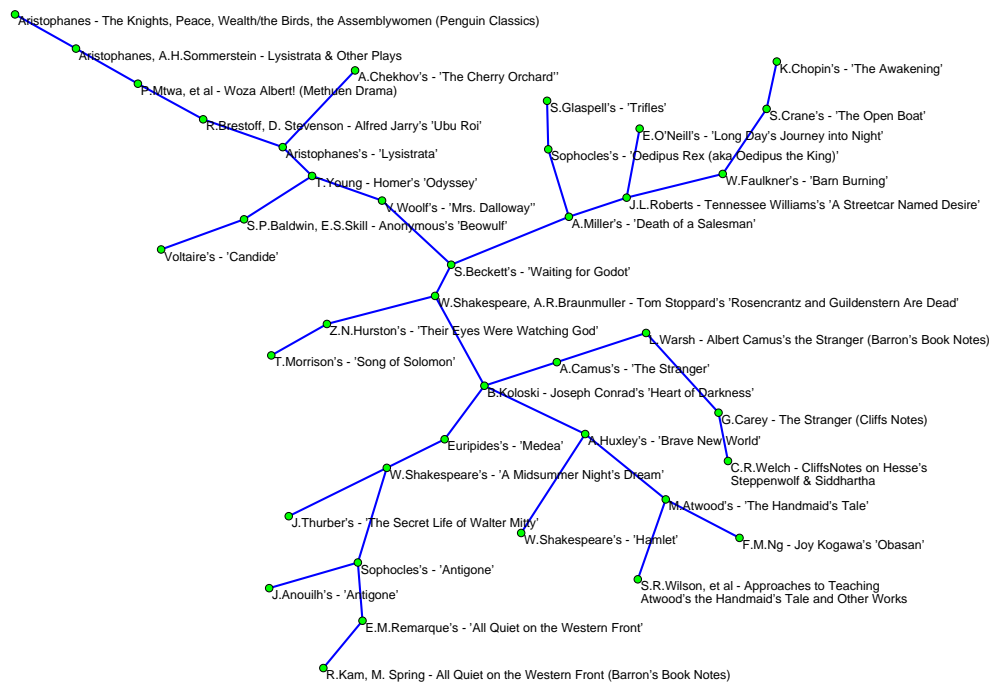**Figure 14:** Island of books for small children.



**Figure 15:** Island of books for students of literature.

Some of them can be obtained from article descriptions from Amazon; the others can be deduced from additional sources.

We believe that analysis of customers' choice networks based on customer's purchases

can offer interesting information for marketing, for example in advertising (interesting themes, most popular books, most popular authors etc.).

Another interesting theme would also be to observe how networks are changing through time.

## Acknowledgment

## References

[1] Batagelj, V. (2004): Collecting network data from the Amazon in Python. `http://vlado.fmf.uni-lj.si/pub/networks/data/econ/amazon/amazon.htm`

[2] Batagelj, V. and Zaveršnik, M. (2004): Islands – identifying themes in large networks. Presented at Sunbelt XXIV Conference, Portorož, May 2004.

[3] Berry, M. J. A. and Linoff, G. S. (2004): *Data Mining Techniques*. Second Edition. Wiley Publishing, Inc.

[4] Kleinberg, J. (1998): Authoritative sources in a hyperlinked environment, *Proc. 9th ACM-SIAM Symposium on Discrete Algorithms*.

[5] De Nooy, W., Mrvar, A., and Batagelj, V. (2005): *Exploratory Social Network Analysis with Pajek*. New York: Cambridge University Press.

[6] Zaveršnik, M. (2003): *Razčlembe omrežij (Network decompositions)*. PhD. Thesis, FMF, University of Ljubljana.

[7] Watts, D. J. and Strogatz, S. (1998): Collective dynamics of 'small-world' networks. *Nature*, **393**. 440-442.

[8] The Amazon Internet Store: `http://www.amazon.com/`

[9] The Pajek program – home page: `http://vlado.fmf.uni-lj.si/pub/networks/pajek/`

[10] The Python Programming Language: `http://www.python.org/`