



**HAL**  
open science

## Biclustering Based on FCA and Partition Pattern Structures for Recommendation Systems

Nyoman Juniarta, Victor Codocedo, Miguel Couceiro, Amedeo Napoli

► **To cite this version:**

Nyoman Juniarta, Victor Codocedo, Miguel Couceiro, Amedeo Napoli. Biclustering Based on FCA and Partition Pattern Structures for Recommendation Systems. NFMCP 2018 - 7th International Workshop on New Frontiers in Mining Complex Patterns, Sep 2018, Dublin, Ireland. hal-01889384

**HAL Id: hal-01889384**

**<https://inria.hal.science/hal-01889384>**

Submitted on 6 Oct 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Biclustering Based on FCA and Partition Pattern Structures for Recommendation Systems

Nyoman Juniarta<sup>1</sup>, Victor Codocedo<sup>2</sup>, Miguel Couceiro<sup>1</sup>, and Amedeo Napoli<sup>1</sup>

<sup>1</sup> Université de Lorraine, CNRS, Inria, LORIA, F-54000 Nancy, France  
{nyoman.juniarta, miguel.couceiro, amedeo.napoli}@loria.fr

<sup>2</sup> Universidad Técnica Federico Santa María, Santiago, Chile  
victor.codocedo@inria.cl

**Abstract.** This paper focuses on item recommendation for visitors in a museum within the framework of European Project CrossCult about cultural heritage. We present a theoretical research work about recommendation using biclustering. Our approach is based on biclustering using FCA and partition pattern structures. First, we recall a previous method of recommendation based on constant-column biclusters. Then, we propose an alternative approach that incorporates an order information and that uses coherent-evolution-on-columns biclusters. This alternative approach shares some common features with sequential pattern mining. Finally, given a dataset of visitor trajectories, we indicate how these approaches can be used to build a collaborative recommendation strategy.

**Keywords:** biclustering, FCA, pattern structures, recommendation

## 1 Introduction

CrossCult (<http://www.crosscult.eu>) is a European project whose idea is to support the emergence of a European cultural heritage by allowing visitors in different cultural sites (e.g. museum, historic city, archaeological site) to improve the quality of their visit by using adapted computer-based devices and to consider the visit at a European level. Such improvement can be accomplished by studying, among others, the possibility to build a dynamic recommendation system. This system should be able to produce a relevant suggestion on which part of a cultural site may be interesting for a specific visitor.

Here, our objective is to study a dynamic recommendation system for visitors in a museum. Given a new visitor  $V_n$ , the task is to suggest a museum item that may be interesting for him/her. Based on how a suggestion is made to a new visitor  $V_n$ , a recommendation system can be classified into one of the three following categories [1]:

- *Content-based recommendations*: The system makes a suggestion based only on the previous visited items of  $V_n$ . For example, if  $V_n$  visited mostly the items from prehistoric era, then the system recommends another item from that era.

- *Collaborative recommendations*: The system looks for previous users who have similar interest to  $V_n$ , and makes a suggestion based on their visited items. For example, if many of  $V_n$ 's similar users have visited item  $I$ , then the system recommends this item.
- *Hybrid approaches*: The combination of content-based and collaborative approaches.

Our method belongs to the second category (collaborative recommendation). First we group all previous users based on their visit trajectories using biclustering. When  $V_n$  arrives, we try to find a  $G_s$ , i.e. a group of visitors who shares a similar interest to  $V_n$ . Then, based on the behavior of the visitors in  $G_s$ , we can suggest one item that may be interesting for  $V_n$ .

In this paper we will recall an approach in [7] that uses partition pattern structures to obtain biclusters with constant (or similar) values on the columns. Then we will propose an alternative approach that relies on this approach to mine another type of biclusters: those with coherent evolution on the columns (CEC biclusters). This bicluster type is useful when we are dealing with a dataset of trajectories where each trajectory corresponds to an ordered list of items. Furthermore, the mining of CEC biclusters can be related to sequential pattern mining, which we will explore in this paper.

This paper is organized as follows. First, we mention some related works about recommendation in Section 2. Then the basic background on biclustering is given in Section 3. Section 4 explains how to perform biclustering using partition pattern structures. The application of biclustering to recommendation systems will be presented in Section 5. Finally, we conclude our paper and outline some future works in Section 6.

## 2 Related work

In this section, we will describe related work about recommendation systems, biclustering, and Formal Concept Analysis (FCA).

FCA has been studied in collaborative movie recommendations for a user by looking at the ratings given by other users. In [5], FCA is used to generate a lattice from a binary matrix (with users as rows and movies as columns) as the formal context. This matrix is derived from a rating dataset which is binarized, such that the matrix contains only the information whether a user has rated a movie. The lattice is then drawn to select some neighbors – i.e. users who have rated the same movies as the new user – regardless of the rating values. In this way, the exhaustive search of neighbors can be avoided. The neighbors' ratings can be then studied to recommend movies rated by the neighbors but not yet rated by the new user.

Pattern structures [9,15] are a generalization of FCA, where the objects have more complex descriptions (e.g. sequence, graph, etc.). FCA was also extended into Triadic Concept Analysis, and it was shown in [15] that triadic concepts are in 1-1-correspondence with maximal biclusters of similar values.

Partition pattern structures are an instance of the pattern structure framework. They were used in a collaborative movie recommendation [8] by identifying similar-column biclusters within the rating matrix. Such a bicluster corresponds to a set of users with similar rating behavior (hence similar interest) across a set of movies. Therefore, to recommend a movie to a new user, there is a search for biclusters whose users have similar interest to him/her. Using the real MovieLens data, a study was also conducted based on Boolean matrix factorization [2].

Moreover, recommendation systems based on FCA and/or biclustering have been applied to other real world problems such as detection of future advertising terms for a company [12], educational orientation of Russian school graduates [13], and idea recommendation at a crowdsourcing project of Witology company [11]. Other than recommendation systems, biclustering has also been applied in other fields, for example in the study of miRNA-gene target interaction data and miRNA functions and mechanisms [19]. A unified taxonomy of biclustering methods was proposed in [14].

### 3 Biclustering

In this section, we will recall the basic background and discuss illustrative examples of the different types of biclusters as described in [18].

We consider a dataset composed of a set of objects, each of which has values over a set of attributes. This dataset can be represented as a numerical matrix, where each cell  $ij$  indicates the value of object  $i$  w.r.t. attribute  $j$ .

One may be interested in finding which subset of objects possesses the same values w.r.t. a subset of attributes. Regarding the matrix representation, this is equivalent to the problem of finding a submatrix that has a constant value over all of its elements (example in Table 1). This task is called biclustering with constant values, which is a simultaneous clustering of the rows and columns of a matrix.

**Table 1.** A bicluster with constant value (shaded)

1	1	4	3	5
1	1	2	5	1
3	3	4	2	1

Other than constant values, the bicluster approach also focused on finding other types of submatrices, as shown in Table 2. A bicluster with constant columns (rows) is a submatrix where each column (row) has the same value, as illustrated in Table 2a (Table 2b, resp.).

In a bicluster with additive coherent values, the value of each cell  $ij$  follows the equation  $\gamma + \alpha_i + \beta_j$ , where  $\gamma$  is a constant,  $\alpha_i$  is a constant value for row  $i$ , and  $\beta_j$  is a constant value for column  $j$ . For example, if  $\gamma = 1$ ,  $(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = (3, 2, 4, 6)$ , and  $(\beta_1, \beta_2, \beta_3, \beta_4) = (0, -2, 1, -1)$ , then we can obtain the bicluster

**Table 2.** Examples of some types of biclusters. (a) Constant columns, (b) constant rows, (c) additive coherent values, (d) multiplicative coherent values, (e) coherent evolution on the columns, and (f) coherent evolution on the rows.

4	2	5	3	1	1	1	1	4	2	5	3	4	2	5	3	1	2	4	3	1	2	4	3
4	2	5	3	2	2	2	2	3	1	4	2	2	1	2.5	1.5	3	5	7	6	0	1	1	2
4	2	5	3	4	4	4	4	5	3	6	4	8	4	10	6	2	3	8	4	5	4	6	4
4	2	5	3	3	3	3	3	7	5	8	6	6	3	7.5	4.5	4	5	9	8	6	5	7	5
	(a)				(b)				(c)				(d)				(e)				(f)		

in Table 2c. Similarly, we can obtain a bicluster with multiplicative coherent values as shown in Table 2d using a constant for each row and each column. The main difference is that, instead of adding, we multiply them.

Another interesting type is the CEC bicluster, also known as order-preserving submatrix [4]. In this type of bicluster, each row induces the same linear order across all columns. For example, in the bicluster in Table 2e, each row follows  $column1 \leq column2 \leq column4 \leq column3$ . Moreover, a bicluster with coherent evolution on the rows can be defined similarly, as shown in Table 2f.

Those different types of biclusters are useful when we are interested in identifying a group of people who behave similarly according to a set of attributes. This group identification is necessary in the task of collaborative recommendation, because in the process of making a suggestion to a person, we first identify the people who are similar to him/her.

### 3.1 FCA and Pattern Structures

(FCA) is a mathematical framework based on lattice theory and used for classification, data analysis, and knowledge discovery [10]. From a formal context, FCA detects all formal concepts, and arranges them in a concept lattice.

**Definition 1.** A formal context is a triple  $(G, M, I)$ , where  $G$  is a set of objects,  $M$  is a set of attributes, and  $I$  is a binary relation between  $G$  and  $M$ , i.e.  $I \subseteq G \times M$ .

If an object  $g$  has an attribute  $m$ , then  $(g, m) \in I$ . An example of a formal context is shown in Table 3. This table shows whether a visitor ( $V_1$ - $V_4$ ) visits an item (102, 302, 402, or 704).

The Galois connection for a formal context  $(G, M, I)$  is defined as follows:

**Definition 2.** For a subset of objects  $A \subseteq G$ ,  $A'$  is the set of attributes that are possessed by all objects in  $A$ , i.e.:

$$A' = \{m \in M \mid \forall g \in A, (g, m) \in I\}, \quad A \subseteq G$$

Dually, for a subset of attributes  $B \subseteq M$ ,  $B'$  is the set of objects that have all attributes in  $B$ , i.e.:

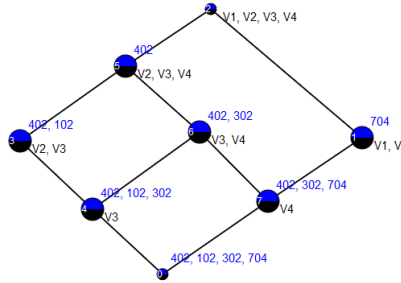
$$B' = \{g \in G \mid \forall m \in B, (g, m) \in I\}, \quad B \subseteq M$$

**Definition 3.** A formal concept is a pair  $(A, B)$ , where  $A \subseteq G$  and  $B \subseteq M$ , and such that  $A' = B$  and  $B' = A$ .

**Table 3.** A formal context for four objects, with four items: 102, 302, 402, and 704 as an example.

	102	302	402	704
$V_1$				×
$V_2$	×		×	
$V_3$	×	×	×	
$V_4$		×	×	×

A formal concept  $(A, B)$  is a *subconcept* of  $(C, D)$  – denoted by  $(A, B) \leq (C, D)$  – if  $A \subseteq C$  (or equivalently  $D \subseteq B$ ). A concept lattice can be formed using the  $\leq$  relation which defines the order among concepts. For the context in Table 3, the formal concepts and their corresponding lattice are shown in Fig. 1.



**Fig. 1.** Concept lattice for the formal context in Table 3

FCA is restricted to specific datasets where each attribute is binary (e.g. has only yes/no value). For more complex values (e.g. numbers, strings, trees, graphs...), FCA is then generalized into pattern structures [9].

**Definition 4.** A pattern structure is a triple  $(G, (D, \sqcap), \delta)$ , where  $G$  is a set of objects,  $(D, \sqcap)$  is a complete meet-semilattice of descriptions, and  $\delta : G \rightarrow D$  maps an object to a description.

The operator  $\sqcap$  is a similarity operation that returns the common elements between any two descriptions. A description can be a set, a sequence, or other complex structure. In the case of set as a description,  $\sqcap$  corresponds to set intersection ( $\cap$ ), i.e.  $\{a, b, c\} \sqcap \{a, b, d\} = \{a, b\}$ , and  $\sqsubseteq$  corresponds to subset inclusion ( $\subseteq$ ). In the case of sequence as a description,  $\sqcap$  is a set of common closed subsequences (SCCS) [6]. Similarly,  $\sqsubseteq$  corresponds to subsequence inclusion ( $\preceq$ ).

**Definition 5.** The Galois connection for a pattern structure  $(G, (D, \sqcap), \delta)$  is defined as:

$$A^\circ = \sqcap_{g \in A} \delta(g), \quad A \subseteq G$$

$$d^\circ = \{g \in G \mid d \sqsubseteq \delta(g)\}, \quad d \in D$$

Finally, a pattern concept is similar to a standard formal concept:

**Definition 6.** A pattern concept is a pair  $(A, d)$ ,  $A \subseteq G$  and  $d \in D$ , where  $A^\circ = d$  and  $d^\circ = A$ .

## 4 Biclustering Using Partition Pattern Structures

Biclustering has many common elements with FCA. In FCA, from a binary matrix we try to find a maximal submatrix whose elements are 1. In other words, the objective is to identify maximal constant-value biclusters (but only for biclusters whose values are 1). Hence, a formal concept can be considered as a bicluster of objects and attributes. Furthermore, formal concepts are arranged in a concept lattice, that can describe the hierarchical relation among all biclusters.

Consider the matrix given by Table 4, where we are interested in finding constant-column biclusters. We recognize that the values of  $\mathbf{m}_1$  “break” the objects into two sets:  $\{\mathbf{g}_1, \mathbf{g}_2\}$  and  $\{\mathbf{g}_3, \mathbf{g}_4\}$ . The same “break” is also obtained from the values of  $\mathbf{m}_4$ . In particular, we can see that the pair  $(\{\mathbf{g}_1, \mathbf{g}_2\}, \{\mathbf{m}_1, \mathbf{m}_4\})$  corresponds to a constant-column bicluster. Therefore, it is possible to mine this type of bicluster using this “breaking” – or “partitioning” – technique. Moreover, this technique can be performed using partition pattern structures – an extension of FCA.

In this section, first we will recall the constant-column biclustering approach using partition pattern structures [7]. We will then propose an extension of this approach to perform CEC biclustering.

### 4.1 Biclustering with Constant Columns

A partition  $\mathbf{d} = \{p_i\}$  of a set  $\mathbf{G}$  is a collection of  $p_i \subseteq \mathbf{G}$  such that:

$$\bigcup_{p_i \in \mathbf{d}} p_i = \mathbf{G} \quad \text{and} \quad p_i \cap p_j = \emptyset \quad \text{whenever} \quad i \neq j. \quad (1)$$

Notice that when calculating the initial partitions, missing values can produce overlapping partitions (i.e.  $p_i \cap p_j \neq \emptyset$ ). Consider the dataset given by Table 4 that has  $\mathbf{G} = \{\mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3, \mathbf{g}_4\}$  as the set of objects and  $\mathbf{M} = \{\mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3, \mathbf{m}_4, \mathbf{m}_5\}$  as the set of attributes. Here we can define a partition mapping  $\delta : \mathbf{M} \rightarrow \mathbf{D}$ . The partition is based on the fact that the values of the attribute are equal for all objects in a subset. For example,  $\delta(\mathbf{m}_1) = \{\{\mathbf{g}_1, \mathbf{g}_2\}, \{\mathbf{g}_3, \mathbf{g}_4\}\}$  because  $\mathbf{G}$  is partitioned as such regarding the value of  $\mathbf{m}_1$ , whereas  $\delta(\mathbf{m}_4) = \{\{\mathbf{g}_1, \mathbf{g}_2\}, \{\mathbf{g}_1, \mathbf{g}_3, \mathbf{g}_4\}\}$ . This partition overlaps on  $\mathbf{g}_1$  since this object has a missing value on  $\mathbf{m}_4$ . Intuitively,  $\mathbf{g}_1$  can be grouped with either  $\{\mathbf{g}_2\}$  or  $\{\mathbf{g}_3, \mathbf{g}_4\}$  w.r.t.  $\mathbf{m}_4$ .

The space  $\mathbf{D}$  of all partitions over  $\mathbf{G}$  is a complete lattice, where the meet and join of two partitions  $\mathbf{d}_1 = \{p_i\}$  and  $\mathbf{d}_2 = \{p_j\}$  are defined as:

**Table 4.** A dataset with 4 objects and 5 attributes

	$m_1$	$m_2$	$m_3$	$m_4$	$m_5$
$g_1$	1	5	3	?	7
$g_2$	1	1	4	2	7
$g_3$	2	5	4	5	3
$g_4$	2	5	4	5	7

$$d_1 \sqcap d_2 = \left( \bigcup_{i,j} p_i \cap p_j \right)^+ \quad (2)$$

$$d_1 \sqcup d_2 = \left( \bigcup_{p_i \cap p_j \neq \emptyset} p_i \cup p_j \right)^+ \quad (3)$$

where  $(\cdot)^+$  is a closure that preserves only the maximal components in  $d$ . For example,  $\delta(m_1) \sqcap \delta(m_4) = \{\{g_1, g_2\}, \{g_1\}, \{g_3, g_4\}\}^+ = \{\{g_1, g_2\}, \{g_3, g_4\}\}$ , and  $\delta(m_1) \sqcup \delta(m_4) = \{\{g_1, g_2\}, \{g_1, g_2, g_3, g_4\}, \{g_1, g_3, g_4\}\}^+ = \{\{g_1, g_2, g_3, g_4\}\}$ .

The order between any two partitions is given by the subsumption relation:

$$d_1 \sqsubseteq d_2 \iff d_1 \sqcap d_2 = d_1 \quad (4)$$

Given a set of attributes  $M$ , a partition space  $D$ , and a mapping  $\delta$ , a partition pattern structures for constant-column biclustering is determined by the triple  $(M, D, \delta)$ . A pair  $(A, d)$  is then called a partition pattern concept (pp-concept) iff  $A^\square = d$  and  $d^\square = A$ , where:

$$A^\square = \bigcap_{m \in A} \delta(m) \quad A \subseteq M \quad (5)$$

$$d^\square = \{m \in M \mid d \sqsubseteq \delta(m)\} \quad d \in D \quad (6)$$

For any partition component  $p \in d$ , each pair  $(p, A)$  corresponds to a constant-column bicluster. For example, from the concept  $(\{m_1, m_4\}, \{\{g_1, g_2\}, \{g_3, g_4\}\})$ , two biclusters can be obtained:  $(\{g_1, g_2\}, \{m_1, m_4\})$  and  $(\{g_3, g_4\}, \{m_1, m_4\})$ .

## 4.2 Biclustering with Coherent Evolution on the Columns

In a dataset of movie ratings, bicluster with constant columns is useful to identify a set of users with the same taste regarding a set of movies. Another interesting problem arises, e.g. when the dataset contains watching order. In that case, we may be interested in finding a set of users who watch a set of movies in the same order. This problem corresponds to **CEC** biclustering, where the objective is to find a set of rows which has coherent evolution over a set of columns, as



**Table 5.** A dataset with 5 objects and 5 attributes

	m <sub>1</sub>	m <sub>2</sub>	m <sub>3</sub>	m <sub>4</sub>	m <sub>5</sub>
g <sub>1</sub>	1	2	3	4	5
g <sub>2</sub>	4	2	1	?	3
g <sub>3</sub>	2	3	4	1	1
g <sub>4</sub>	5	4	2	3	1
g <sub>5</sub>	2	1	5	4	3

**Table 6.** Some examples of partitions over Table 5

Pair	Partition
P <sub>1,2</sub>	{{g <sub>1</sub> , g <sub>3</sub> }, {g <sub>2</sub> , g <sub>4</sub> , g <sub>5</sub> }}
P <sub>1,3</sub>	{{g <sub>1</sub> , g <sub>3</sub> , g <sub>5</sub> }, {g <sub>2</sub> , g <sub>4</sub> }}
P <sub>1,4</sub>	{{g <sub>1</sub> , g <sub>2</sub> , g <sub>5</sub> }, {g <sub>2</sub> , g <sub>3</sub> , g <sub>4</sub> }}
P <sub>2,3</sub>	{{g <sub>1</sub> , g <sub>3</sub> , g <sub>5</sub> }, {g <sub>2</sub> , g <sub>4</sub> }}
P <sub>2,5</sub>	{{g <sub>1</sub> , g <sub>2</sub> , g <sub>5</sub> }, {g <sub>3</sub> , g <sub>4</sub> }}

previously described in Section 3. In the current section, we will explain the possible application of partition pattern structures to discover CEC biclusters.

Consider the dataset given by Table 5, with the set of attributes  $G = \{g_1, g_2, g_3, g_4, g_5\}$ . First, we have to list each pair of attributes and the partition according to the pair’s evolution. For the pair  $p_{1,2} = (m_1, m_2)$ , the partition is  $\{\{g_1, g_3\}, \{g_2, g_4, g_5\}\}$  because in  $g_1$  and  $g_3$ ,  $m_1$  is less than  $m_2$ , whereas in  $g_2$ ,  $g_4$ , and  $g_5$ ,  $m_1$  is greater. As in Subsection 4.1, missing values generate an overlapping partition. For instance, the pair  $p_{1,4}$  gives rise to the partition is  $\{\{g_1, g_2, g_5\}, \{g_2, g_3, g_4\}\}$ . Furthermore, two columns with the same value (e.g.  $g_3$  in  $m_4$  and  $m_5$ ) can also produce an overlapping partition because, by our definition of CEC bicluster in Section 3, they satisfy  $m_4 \leq m_5$  and  $m_5 \leq m_4$ . Therefore, the partition for  $p_{4,5}$  is  $\{\{g_1, g_2, g_3\}, \{g_2, g_3, g_4, g_5\}\}$ . Some pairs and their partitions are listed in Table 6.

Since a CEC partition is defined by at least two attributes, the partition mapping becomes  $\gamma : P \rightarrow D$ . For instance,  $\gamma(p_{1,2}) = \{\{g_1, g_3\}, \{g_2, g_4, g_5\}\}$ .

As in Subsection 4.1, given a set of attribute pairs  $P$ , a partition space  $D$ , and the mapping function  $\gamma$ , a partition pattern structures for coherent-evolution biclustering is determined by the triple  $(P, (D, \sqsubseteq), \gamma)$ . A pp-concept is a pair  $(B, d)$  such that  $B^\square = d$  and  $d^\square = B$ , where:

$$B^\square = \bigcap_{p \in B} \gamma(p) \quad B \subseteq P \quad (7)$$

$$d^\square = \{p \in P \mid d \sqsubseteq \gamma(p)\} \quad d \in D \quad (8)$$

Here, the extent of a pp-concept is a set of attribute pairs. We can obtain a CEC bicluster in a pp-concept if there is a clique among the attributes in the pairs. For example, consider the pp-concept  $\text{ppc}_1$  with extent  $\{p_{1,2}, p_{1,3}, p_{2,3}\}$  and intent  $\{\{g_1, g_3\}, \{g_5\}, \{g_2, g_4\}\}$ . Its extent forms a clique among  $m_1$ ,  $m_2$ , and  $m_3$ , since all pairings of any two of those attributes are included.

If a pp-concept  $(B, d)$  contains a set of attributes  $A$  that forms a clique, then each pair  $(p, A)$ , for any partition component  $p \in d$ , corresponds to a CEC bicluster. For example, from  $\text{ppc}_1$ , we can obtain bicluster  $(\{g_1, g_3\}, \{m_1, m_2, m_3\})$ .

### 4.3 Comparison with Sequential Pattern Mining

A sequence is an ordered list  $\langle s_1 s_2 \dots s_m \rangle$ , where  $s_i$  is an itemset  $\{i_1, \dots, i_n\}$ . A sequence  $s = \langle s_1 s_2 \dots s_m \rangle$  is a subsequence of  $s' = \langle s'_1 s'_2 \dots s'_n \rangle$ , denoted by

$s \preceq s'$ , if there exist indices  $1 \leq i_1 < i_2 < \dots < i_m \leq n$  such that  $s_j \subseteq s'_{i_j}$  for all  $j = 1 \dots m$  and  $m \leq n$ . For example, the sequence  $\langle \{a\}\{d\} \rangle$  is a subsequence of  $\langle \{a, b\}\{a, c, d\} \rangle$ , while sequence  $\langle \{c\}\{d\} \rangle$  is not.

Notice that the problem of retrieving CEC biclusters can be thought of as a particular type of sequential pattern mining where each itemset is composed by only one item, i.e. the sequences are an ordered list of items. Mining sequential patterns means retrieving frequent subsequences (i.e., subsequences that are present in more than  $n$  sequences) and for which there exist many efficient algorithms [21,20].

The CEC biclustering differs from sequential pattern mining when we allow overlaps in the partitions. Consider Table 5 as a sequential dataset. Each number in row  $x$  column  $y$  corresponds to the itemset when item  $y$  appears in the sequence  $x$ . For example, the sequence of object  $\mathbf{g}_3$  is  $\langle \{\mathbf{m}_4, \mathbf{m}_5\}\{\mathbf{m}_1\}\{\mathbf{m}_2\}\{\mathbf{m}_3\} \rangle$ .

Let us consider items  $\mathbf{m}_1$ ,  $\mathbf{m}_4$ , and  $\mathbf{m}_5$ . According to sequential pattern mining,  $\mathbf{g}_3$  is different from  $\mathbf{g}_4$ , because  $\mathbf{m}_4$  and  $\mathbf{m}_5$  appear in the same itemset in  $\mathbf{g}_3$ . On the other hand, according to CEC biclustering with overlaps,  $\mathbf{g}_3$  is similar to  $\mathbf{g}_4$ , because in both objects  $\mathbf{m}_5 \leq \mathbf{m}_4 \leq \mathbf{m}_1$ .

## 5 Recommendation

In the context of CrossCult, we are working on a visitor dataset that comprises several trajectories in a museum. Within this project, our main objective is to build a dynamic recommendation system for new visitors. This system should be able to suggest a museum item to visitors based on their trajectories and by looking at the trajectories of previous visitors. Also, it should be able to update the suggestion as they move inside the museum.

### 5.1 Matrix as order of interest

For each item in the museum, we can measure (e.g. by rating, duration of visit, etc.) their level of interestingness from a set of visitors. An example is shown in Table 7, where the number in cell  $xy$  is the ranking of item  $y$  according to visitor  $x$ . Here we have 3 visitors ( $\mathbf{v}_1$ ,  $\mathbf{v}_2$ , and  $\mathbf{v}_3$ ) in the database and 1 target visitor ( $\mathbf{v}_a$ ). Among the existing visitors, only  $\mathbf{v}_1$  has complete values over all five items. According to this visitor  $\mathbf{i}_1$  is the best, followed by  $\mathbf{i}_2$ ,  $\mathbf{i}_3$ ,  $\mathbf{i}_4$ , and the worst  $\mathbf{i}_5$ . For  $\mathbf{v}_2$  ( $\mathbf{v}_3$ ), the order of interest of  $\mathbf{i}_2$  ( $\mathbf{i}_4$  resp.) is not known.

The target visitor ( $\mathbf{v}_a$ ) has visited only three items, with the same order of preference in  $\mathbf{i}_1$  and  $\mathbf{i}_2$ . This visitor will be included in the bicluster  $(\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_a\}, \{\mathbf{i}_1, \mathbf{i}_2, \mathbf{i}_3\})$ . The other two members of this bicluster do not agree on the order of interest of  $\mathbf{i}_4$  and  $\mathbf{i}_5$ . Hence, we should suggest  $\mathbf{i}_5$  to him/her, since one visitor ( $\mathbf{v}_2$ ) similar to him/her ranked it first.

### 5.2 Matrix as order of visit

Consider the dataset given by Table 8 about 4 visitors in a museum with 7 items. As explained in Subsection 4.3, this table can be regarded as a sequential

**Table 7.** Order of interest of 5 items, observed from certain visitors

	$i_1$	$i_2$	$i_3$	$i_4$	$i_5$
$v_1$	1	2	3	4	5
$v_2$	3	?	1	2	4
$v_3$	2	4	3	?	1
$v_a$	1	1	2	?	?

**Table 8.** Order of visit of 7 items, observed from certain visitors

	$i_1$	$i_2$	$i_3$	$i_4$	$i_5$	$i_6$	$i_7$
$v_1$	1	2	3	4	5	6	7
$v_2$	2	4	5	3	7	1	6
$v_3$	4	2	1	5	6	3	7
$v_4$	7	3	1	4	2	6	5
$v_a$	?	?	1	2	?	?	?

dataset. The numbers in row  $x$  indicate the path of visitor  $v_x$ . For example, the path of  $v_2$  is  $i_6 \rightarrow i_1 \rightarrow i_4 \rightarrow i_2 \rightarrow i_3 \rightarrow i_7 \rightarrow i_5$ .

Now we have a new visitor  $v_a$  who recently arrives to the museum. He/She visits  $i_3$ , followed by  $i_4$ . Our task is to recommend a new item to her, by studying the CEC biclusters over the first four visitors. Some of those biclusters are listed in Table 9. From B7, we can see that all four visitors visit  $i_7$  after  $i_2$ . In B3,  $v_1$  and  $v_3$  follow the same order w.r.t.  $\{i_3, i_4, i_5, i_7\}$ :  $i_3 \rightarrow i_4 \rightarrow i_5 \rightarrow i_7$ . These two visitors also agree in the order of items  $\{i_1, i_4, i_5, i_7\}$ , as seen in B4.

**Table 9.** Some CEC biclusters in Table 8

#	Visitors	Items (in order)
B1	$v_1$	$i_1, i_2, i_3, i_4, i_5, i_6, i_7$
B2	$v_2$	$i_6, i_1, i_4, i_2, i_3, i_7, i_5$
B3	$v_1, v_3$	$i_3, i_4, i_5, i_7$
B4	$v_1, v_3$	$i_1, i_4, i_5, i_7$
B5	$v_1, v_4$	$i_3, i_5, i_6$
B6	$v_2, v_3$	$i_6, i_1, i_4, i_5$
B7	$v_1, v_2, v_3, v_4$	$i_2, i_7$

Those CEC biclusters can be studied to give a recommendation to  $v_a$  by focusing on those visitors that are similar to him/her. Thus, we can propose a recommendation strategy that follows sequential patterns in the dataset. The idea behind is the following: if many visitors have the path  $i_a \rightarrow i_b \rightarrow i_c$ , then we should recommend item  $i_c$  to a visitor who has done  $i_a \rightarrow i_b$ .

Since  $v_a$  has path  $i_3 \rightarrow i_4$ , we focus on the CEC biclusters that have those two items, i.e. B1, B2, and B3. One of those biclusters (B2) has a different ordering ( $i_3$  after  $i_4$ ), and thus we filter it out. Then, in B3 for example, the path is  $i_3 \rightarrow i_4 \rightarrow i_5 \rightarrow i_7$ . Therefore, we can recommend  $i_5$  to  $v_a$ .

### 5.3 Application to Hecht Museum

In the framework of the CrossCult project, we are working on a specific dataset about the trajectories of 254 visitors in Hecht Museum in Haifa, Israel [17]. In

this dataset, there are 52 items over 8 rooms (A–G). Therefore, the visitor–item matrix for biclustering will have 254 rows and 52 columns.

A visitor’s trajectory corresponds to a list of visits, which is composed by three elements: “start time”, “end time”, and “item name”. From the trajectories, we can build either matrix of order of interest or matrix of order of visit. Then, the application of recommendation can be tested over this dataset.

## 6 Conclusion

In this work, we have explored an approach to build collaborative recommendation strategy for visitors in a museum. This strategy takes into account the order of interest or the order of visit for each visitor, and we showed how to use CEC biclustering to obtain a set of similar visitors. We also presented a technique for mining CEC biclusters based on FCA using partition pattern structures. This recommendation strategy can be applied to any dataset where the order of items is relevant.

As future work, we intend to explore recommendations based on the order of visit. The problem to be solved is how to model visitors who visit a single item multiple times (for example,  $i_1$ ,  $i_2$ , and back to  $i_1$ ).

Another noteworthy question is how to measure the “score” of each bicluster in order to rank recommendations for a new visitor. Ranking candidate items was studied in [3] for constant-value biclusters, and it is possible to extend this work to CEC biclusters. Moreover, further comparisons of CEC biclustering and sequential pattern mining should be investigated, in particular, regarding their complexities and their results. Biclustering-based recommender systems could also be compared to social network analysis and cluster-indexing collaborative filtering [16]. Finally, an implementation of the CEC biclustering using partition pattern structures and an empirical study on real-world data should be performed to measure its complexity and efficiency.

## References

1. Adomavicius, G., Tuzhilin, A.: Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE transactions on knowledge and data engineering* **17**(6), 734–749 (2005)
2. Akhmaturov, M., Ignatov, D.I.: Context-aware recommender system based on boolean matrix factorisation. In: *CLA*. pp. 99–110 (2015)
3. Alqadah, F., Reddy, C.K., Hu, J., Alqadah, H.F.: Biclustering neighborhood-based collaborative filtering method for top-n recommender systems. *Knowledge and Information Systems* **44**(2), 475–491 (2015)
4. Ben-Dor, A., Chor, B., Karp, R., Yakhini, Z.: Discovering local structure in gene expression data: the order-preserving submatrix problem. *Journal of computational biology* **10**(3-4), 373–384 (2003)
5. du Boucher-Ryan, P., Bridge, D.: Collaborative recommending using formal concept analysis. *Knowledge-Based Systems* **19**(5), 309–315 (2006)
6. Codocedo, V., Bosc, G., Kaytoue, M., Boulicaut, J.F., Napoli, A.:

7. Codocedo, V., Napoli, A.: Lattice-based biclustering using partition pattern structures. In: Proceedings of the Twenty-first European Conference on Artificial Intelligence. pp. 213–218. IOS Press (2014)
8. Codocedo-Henríquez, V.: Contributions à l’indexation et à la récupération d’information utilisant l’analyse formelle de concepts. Ph.D. thesis, Université de Lorraine (2015), <http://www.theses.fr/2015LORR0143>
9. Ganter, B., Kuznetsov, S.O.: Pattern structures and their projections. In: International Conference on Conceptual Structures. pp. 129–142. Springer (2001)
10. Ganter, B., Wille, R.: Formal Concept Analysis: Mathematical Foundations (1999)
11. Ignatov, D.I., Kaminskaya, A.Y., Konstantinova, N., Malyukov, A., Poelmans, J.: FCA-based recommender models and data analysis for crowdsourcing platform Witology. In: International Conference on Conceptual Structures. pp. 287–292. Springer (2014)
12. Ignatov, D.I., Kuznetsov, S.O., Poelmans, J.: Concept-Based Biclustering for Internet Advertisement. In: 2012 IEEE 12th International Conference on Data Mining Workshops. pp. 123–130. Brussels, Belgium (Jun 2009), <http://arxiv.org/abs/0906.4982>, arXiv: 0906.4982
13. Ignatov, D.I., Poelmans, J., Zaharchuk, V.: Recommender system based on algorithm of bicluster analysis RecBi. arXiv preprint arXiv:1202.2892 (2012)
14. Ignatov, D.I., Watson, B.W.: Towards a unified taxonomy of biclustering methods. arXiv preprint arXiv:1702.05376 (2017)
15. Kaytoue, M., Kuznetsov, S.O., Macko, J., Napoli, A.: Biclustering meets triadic concept analysis. *Annals of Mathematics and Artificial Intelligence* **70**(1-2), 55–79 (2014)
16. Kim, K.j., Ahn, H.: Recommender systems using cluster-indexing collaborative filtering and social data analytics. *International Journal of Production Research* **55**(17), 5037–5049 (2017)
17. Lanir, J., Kuflik, T., Dim, E., Wecker, A.J., Stock, O.: The influence of a location-aware mobile guide on museum visitors’ behavior. *Interacting with Computers* **25**(6), 443–460 (2013)
18. Madeira, S.C., Oliveira, A.L.: Biclustering algorithms for biological data analysis: a survey. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)* **1**(1), 24–45 (2004)
19. Pio, G., Ceci, M., Malerba, D., D’Elia, D.: ComiRNet: a web-based system for the analysis of miRNA-gene regulatory networks. *BMC Bioinformatics* **16**(9), S7 (2015)
20. Wang, J., Han, J.: BIDE: Efficient mining of frequent closed sequences. In: Data Engineering, 2004. Proceedings. 20th International Conference on. pp. 79–90. IEEE (2004)
21. Yan, X., Han, J., Afshar, R.: CloSpan: Mining: Closed sequential patterns in large datasets. In: Proceedings of the 2003 SIAM international conference on data mining. pp. 166–177. SIAM (2003)