

# Decomposing Co-occurrence Matrices into Interpretable Components as Formal Concepts

Anonymous ACL submission

## Abstract

This study addresses the interpretability of word representations through an investigation of a count-based co-occurrence matrix. Employing the mathematical methodology of Formal Concept Analysis, we reveal an underlying structure that is amenable to human interpretation. Furthermore, we unveil the emergence of hierarchical and geometrical structures within word vectors as consequences of word usage. Our experiments on the PPMI matrix demonstrate that the formal concepts we identified align with interpretable categories, as shown in the *category completion task*.

## 1 Introduction

Word vector representations are central to natural language processing, as they capture semantic and syntactic features (Lenci, 2018). Their significance has amplified in recent times, as they are used as input for Transformer-based language models (Vaswani et al., 2017), where static embeddings are contextualized. Their effectiveness has been explained by the distributional hypothesis (Harris, 1954) linking similar semantics and similar distribution (Jurafsky and Martin, 2009). However, the interpretability of their dimensions remains an active research topic (Şenel et al., 2018). Levy and Goldberg (2014a) found neural word embeddings to be uninterpretable while acknowledging that sparse vectors capture some latent topics. Geva et al. (2022) pioneered efforts to interpret dynamic embeddings in GPT-2 (Radford et al., 2019) by projection into the vocabulary space, though a systematic approach to interpret dimensions of embeddings remains an open issue.

Many preceding studies have investigated the semantic properties of word embeddings and revealed that word vectors in a vector space capture relational meanings. The most well-known example is the parallelogram formed in the vector space

by the embeddings of words in analogical relations (e.g. *king:queen::man:woman*) (Mikolov et al., 2013c). Other semantic relationships also exhibit geometrical counterparts, such as semantic composition with vector addition (Mikolov et al., 2013b; Mitchell and Lapata, 2008), hypernymy captured by linear projection (Fu et al., 2014), and polysemy as a linear combination of vectors (Arora et al., 2018). Regarding the theoretical analysis of embeddings, Levy and Goldberg (2014b) suggested that word2vec (Mikolov et al., 2013a) is equivalent to the factorization of a word co-occurrence matrix. Arora et al. (2016) proposed a generative model in which PMI-based word embeddings exhibit linear structures. These related studies collectively hint that the latent structure in the co-occurrence matrix reflects linguistic regularities and is inherently embedded within vector representations. Therefore, understanding the word co-occurrence matrix represents a cornerstone in elucidating the interpretability of word representations.

In this study, we directly address the mathematical structure of a word co-occurrence matrix to uncover underlying linguistic patterns and to interpret the dimensions of word embeddings. We claim that a *formal concept*, as mathematically defined in the matrix, corresponds to interpretable categories. We substantiate our claim through the *category completion task*. Specifically, we used Formal Concept Analysis (FCA), a field of applied mathematics (Ganter and Wille, 2012), to formally characterize the internal structure of a matrix. We define a group of words as interpretable if it can be descriptively labeled. Furthermore, we demonstrate that a hierarchical structure of formal concepts emerges as a geometric formation in the vector space, which explains why relational meanings are captured by word embeddings.

Our contributions are threefold. First, we propose two methods that apply FCA to real-valued data: binarization by varying thresholds and fuzziness

fication of FCA. Second, we empirically show that the formal concepts in the co-occurrence matrix coincide with interpretable categories. Third, we present a novel algorithm to detect formal concepts, which is capable of disambiguating polysemous words. To our knowledge, this is the first study to apply FCA to a word-word co-occurrence matrix. Our study offers a new approach to uncover latent linguistic structures in co-occurrence matrices.

## 2 Formal concept analysis of word co-occurrence matrix

### 2.1 Basics of FCA

FCA is related to order theory and abstract algebra. It mathematizes *concepts* and *conceptual hierarchy* (Ganter and Wille, 2012). A concept comprises a pair of its extents (objects) and its intents (attributes). Concepts can form a hierarchy known as a *lattice*. FCA has been empirically applied for data mining and ontology (Poelmans et al., 2013), especially in bioinformatics (Roscoe et al., 2022).

A **formal context**  $\mathbb{K} := (G, M, I)$  consists of two sets  $G, M$  and a binary relation  $I \subseteq G \times M$ . The elements of  $G$  and  $M$  are called **objects** and **attributes**, respectively. For  $g \in G$  and  $m \in M$ , a relation  $(g, m) \in I$  means that the object  $g$  has the attribute  $m$ . We define two derivation operators;  $\uparrow : 2^G \rightarrow 2^M$  maps a subset of objects to a subset of attributes, and its reverse  $\downarrow : 2^M \rightarrow 2^G$  maps attributes to objects. For  $A \subseteq G, B \subseteq M$ ,

$$A\uparrow := \{m \in M \mid (g, m) \in I \ (\forall g \in A)\} \quad (1)$$

$$B\downarrow := \{g \in G \mid (g, m) \in I \ (\forall m \in B)\} \quad (2)$$

$A\uparrow \subseteq M$  is the set of attributes common to all objects in  $A$ , whereas  $B\downarrow \subseteq G$  is the set of objects that possess all the attributes in  $B$ . It can be shown that  $A \subseteq B\downarrow \Leftrightarrow B \subseteq A\uparrow$ , which is a structure-preserving (order-reversing) correspondence between ordered sets known as a Galois connection (Davey and Priestley, 2002).

A **formal concept** of the context  $(G, M, I)$  is defined as a pair  $(A, B) \in 2^G \times 2^M$  where both  $A\uparrow = B$  and  $B\downarrow = A$  hold.  $A$  and  $B$  are considered the extent and intent, respectively, of the formal concept  $(A, B)$ . The compositions of two derivation operators  $\uparrow\downarrow : 2^G \rightarrow 2^G$  and  $\downarrow\uparrow : 2^M \rightarrow 2^M$  are closure operators (Davey and Priestley, 2002), with a formal concept defined as the fixed point of these operations. If a formal context is represented as a binary matrix, it corresponds to a maximal rectangular (submatrix) with

all ones in its entries when the rows and columns are appropriately reordered.

A formal concept can also be equated with a maximal **biclique**, i.e., a complete subgraph of a bipartite graph (Chiaselotti et al., 2015). All elements of  $A$  and  $B$  are completely connected within that subgraph.

### 2.2 Rational and benefit of using FCA

A word co-occurrence matrix, used as input data to learn word embeddings, is constructed by counting the frequency of a target-context word pair that co-occurs in the neighborhood. By regarding target words as objects and context words as attributes, we can express this co-occurrence as a binary relation. Thus, we can treat a co-occurrence matrix as a formal context.

FCA is effective in analyzing co-occurrence matrices for three reasons. First, it can characterize a local structure within the matrix. Second, formal concepts can capture relations between more than three words, which cannot be represented by individual pairwise relationships, yielding a richer analysis of the structure. Third, we can define (partial) order relation between formal concepts. A semantic relationship such as hypernymy can be formalized by such an order relation. We further demonstrate the function of FCA in Section 3.

To apply the crisp (binary) FCA to a real-valued co-occurrence matrix, we tested two approaches. First, we simply binarized the matrix values by thresholds, with a varying threshold method deployed to flexibly locate formal concepts (Section 4). Second, we extended the crisp FCA to an FCA built on fuzzy logic (Section 5).

## 3 Demonstration using synthetic data

### 3.1 Artificial toy corpus

We examined how FCA handles a word co-occurrence matrix using a toy corpus. We demonstrated that formal concepts capture semantic categories emerging from word usage in the corpus and introduced a **concept lattice** of FCA to illustrate the hierarchical structure of concepts.

The demonstration contains 1) a corpus of 24 synthetic sentences with 17 words (Appendix A), 2) a co-occurrence matrix obtained from the corpus, and 3) word vectors acquired from the matrix (Fig. 1). The corpus is designed to replicate a geometric formation of the analogy relation. Specifically, we targeted eight words—*king*, *queen*, *man*,

woman, and their plurals—so that their vectors formed a parallelepiped. The sentences were expressed analogously: E.g., ‘king (queen) live in palace’, whereas ‘man (woman) live in house’. The co-occurrence matrix  $X \in \{0, 1\}^{17 \times 17}$  is bi-

	palace	house	tie	dress	alone	together
king	1	0	1	0	1	0
man	0	1	0	0	1	0
queen	1	0	0	1	0	0
woman	0	1	0	1	0	0
kings	1	0	1	0	1	0
men	0	1	0	0	1	0
queens	1	0	0	1	0	0
women	0	1	0	1	0	0

Figure 1: Binary co-occurrence (sub)matrix: Each entry is 1 if shaded and 0 otherwise. Each row is a word vector. Three submatrices with shade patterns indicate different formal concepts  $f, e, v$ .

nary, where  $X_{ij} = 1$  if two words co-occur in a sentence and  $X_{ij} = 0$  otherwise. Each row of this matrix represents a word vector. Projected on the 3-dimensional space, the eight word vectors form a parallelepiped (Fig. 2).

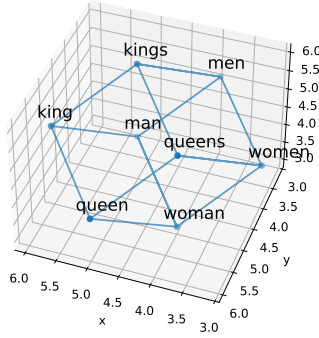


Figure 2: A parallelepiped emerges when eight word vectors (rows) are projected onto 3-dimensional space.

### 3.2 Detecting formal concepts

We now apply FCA to the matrix  $X$ . Although formal concepts can be determined by applying the closure operator  $\uparrow\downarrow$ , a simplified method is to find a rectangular in the matrix. For example, the submatrix of rows  $i \in \{1, 3, 4, 7\}$  and columns  $j \in \{1\}$  represents a formal concept, as all its entries are 1s and no other rectangular matrix contains it. This concept represents a pair of the extent  $\{king, queen, kings, queens\}$  and the intent  $\{palace\}$ , interpreted as ‘royal.’

There are a total of 28 formal concepts in this matrix (see Appendix B for the list and notation). They are classified into five types, including two trivial ones wherein one element is empty. Examples of the three non-trivial types include the following:

$$f_1 := (\{king, man, kings, men\}, \{tie\}) \quad (3)$$

$$e_1 := (\{king, man\}, \{tie, alone\}) \quad (4)$$

$$v_1 := (\{king\}, \{tie, palace, alone\}) \quad (5)$$

To see hierarchical relations between formal concepts, we first define the order relation. Let  $\mathfrak{B}(G, M, I)$  be the set of all concepts of  $(G, M, I)$ . Given  $(A_1, B_1), (A_2, B_2) \in \mathfrak{B}(G, M, I)$ ,

$$(A_1, B_1) \leq (A_2, B_2) \stackrel{\text{def}}{\iff} A_1 \subseteq A_2 \iff B_1 \supseteq B_2 \quad (6)$$

Thus, if the extent  $A_1$  is contained by the extent  $A_2$ , then the formal concept  $(A_1, B_1)$  is less than or equal to  $(A_2, B_2)$ . Owing to the Galois connection,  $A_1 \subseteq A_2$  holds if and only if  $B_1 \supseteq B_2$ . Then,  $\langle \mathfrak{B}(G, M, I) : \leq \rangle$  is a complete lattice known as a **concept lattice**, a nonempty ordered set where a join and a meet exist for all elements and subsets. Fig. 3 visualizes all ordered relations between the formal concepts identified in the matrix  $X$ . We

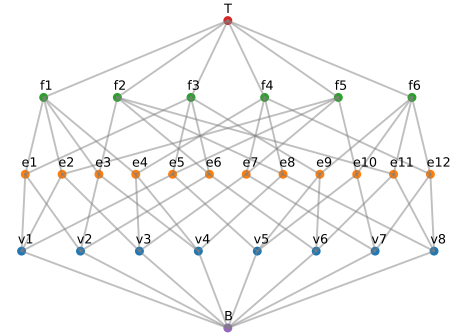


Figure 3: Concept lattice. Each node represents a formal concept. They correspond to geometric simplices of the parallelepiped: 8 vertices, 12 edges, 6 faces.

observe that the lattice of formal concepts (Fig. 3) corresponds to the parallelepiped (Fig. 2). This suggests that geometric relations between word vectors reflect the hierarchical structure latent in the word co-occurrence matrix.

### 3.3 Three implications of FCA

First, FCA allows us to easily interpret the identified formal concepts. For example,  $f_1$  should be labeled as **masculine** from its extent  $\{king, kings\}$ ,

*man, men*}, whereas  $f_6$ , with the extent  $\{queen, queens, woman, women\}$ , must be labeled as **fem-inine**. The other  $f$ -type concepts can be labeled as **royal**, **common**, **singular** and **plural**. Thus, formal concepts coincide with semantic categories.

Second,  $v_1$  (*king*) can be seen as the intersection of three others— $f_1, f_3, f_5$ —analogous to a vertex included in three faces. Semantically, *king* is something royal, masculine, and singular. This relation can be algebraically formulated as  $v_1 = f_1 \wedge f_3 \wedge f_5$  where  $\wedge$  is a *meet* operation.

Third, pairs of opposing faces in the parallelepiped form complementary concepts such as **masculine** vs. **feminine**. Mathematically, we can construct a *formal concept algebra* by defining additional operations as axioms (Wille, 2004). Using this algebra, the formal concept of **masculine** can be demonstrated to complement that of **fem-inine**;  $\neg f_1 = f_6$  where  $\neg$  is a negation. The observation that *king*  $v_1 = f_1 \wedge f_3 \wedge f_5$  and *queen*  $v_5 = f_6 \wedge f_3 \wedge f_5$  share  $f_3$  and  $f_5$  explains the phenomenon that both synonyms and antonyms appear close to each other in the vector space (Turney and Pantel, 2010).

In summary, the co-occurrence matrix exhibits the geometrical and algebraic structures formed by interpretable formal concepts.

## 4 Experiment 1: FCA by binarization

We now demonstrate that formal concepts can be defined on actual word co-occurrence data and correspond to both semantic and syntactic categories.

### 4.1 Algorithm to identify formal concepts

We designed a novel algorithm to locate formal concepts through the conversion of two derivation operators (Eq. 1 and 2). The corresponding pseudo-algorithm is shown in Algorithm 1. Given a co-occurrence matrix  $X$  and set of target words  $S$  as a seed, the algorithm returns a formal concept  $(S^{\uparrow\downarrow}, S^{\uparrow})$ , which is a pair of two subsets of the vocabulary. Here,  $S^{\uparrow\downarrow}$  is the closed set of  $S$ .

The first derivation operator  $\uparrow$  must identify context words that co-occur with all target words in  $S$ . In other words, a context word is selected when it has all entry values exceeding the threshold  $t$  for the target words in  $S$ . Equivalently, any entry value that the seed words have with the context word should not be less than  $t$ , meaning that their minimum must be greater than or equal to  $t$ . As indicated in Line 3, the algorithm finds

---

### Algorithm 1 Varying Threshold Method

---

**Input:**  $X \in \mathbb{R}^{N \times N}$ ,  $S := \{w_i\}_{i \in I_S}$ ,  $k \in \mathbb{N}$

**Output:**  $FC := (S^{\uparrow\downarrow}, S^{\uparrow})$ ,  $t \in \mathbb{R}$

```

1: function FINDFORMALCONCEPT( $S, k$ )
2:   for  $j \leftarrow 1$  to  $N$  do
3:      $m_j \leftarrow \min_{i \in I_S} X_{ij}$ 
4:   end for
5:   Sort  $[m_j]$  in descending order  $\leftarrow [m_{p(j)}]$ 
6:    $J_{S^{\uparrow}} \leftarrow \{p(j)\}_{j \leq k}$ 
7:    $S^{\uparrow} \leftarrow \{w_j\}_{j \in J_{S^{\uparrow}}}$ 
8:    $t \leftarrow m_{p(k)}$ 
9:    $I_{S^{\uparrow\downarrow}} \leftarrow \emptyset$ 
10:  for  $i \leftarrow 1$  to  $N$  do
11:     $\mu_i \leftarrow \min_{j \in J_{S^{\uparrow}}} X_{ij}$ 
12:    if  $\mu_i \geq t$  then
13:       $I_{S^{\uparrow\downarrow}} \leftarrow I_{S^{\uparrow\downarrow}} \cup \{i\}$ 
14:    end if
15:  end for
16:   $S^{\uparrow\downarrow} \leftarrow \{w_i\}_{i \in I_{S^{\uparrow\downarrow}}}$ 
17:  return  $(S^{\uparrow\downarrow}, S^{\uparrow})$ ,  $t$ 
18: end function

```

---

the minimum value that the seed words (in rows  $\forall i \in I_S$ ) have against a certain context word (in a column  $j \in \{1, \dots, N\}$ ), sorts them in descending order (Line 5), and selects the first  $k$  context words (columns)  $S^{\uparrow}$  (Line 6). The threshold is automatically determined as the  $k$ th largest minimum value (Line 8). Next, an inverse operation executes. Given  $S^{\uparrow}$ , the algorithm finds a minimum value over the context words  $S^{\uparrow}$  ( $J_{S^{\uparrow}}$  in the column index) against a target word in a row  $i$  (Line 11) and selects the target words (rows  $I_{S^{\uparrow\downarrow}}$ ) with minimum values exceeding the threshold (Line 13), which form  $S^{\uparrow\downarrow}$ . For a discussion on mathematical properties of identified submatrices, see Appendix C.

### 4.2 Category completion test

The experiment was conducted to verify that the formal concepts identified from the co-occurrence matrix coincide with interpretable categories.

**Test set** We adopted two existing test sets from Lindh-Knuutila and Honkela (2015) containing semantic categories: the Battig set (Bullinaria and Levy, 2012), comprising 53 categories with 10 words for each, and BLESS (Baroni and Lenci, 2011), containing 17 categories with 5-17 words for each. We also compiled two additional sets: Series and Syntactic. The categories tested are listed in Appendix D.



**Procedure** For each category, we systematically furnished the algorithm with all possible word pairs as seeds derived from the category’s word set. Next, we identified the optimal seed that yields the most extensive set of accurately classified words. We then assessed how effectively the algorithm retrieves the correct words from the optimal seed for the given category (**Precision, Recall**). Because the word sets are not necessarily exhaustive, we also regarded those missed words as correct, based on our human judgement (**Extended precision**)<sup>1</sup>.

**Baseline** We used a similarity-based approach as a baseline. Specifically, we applied the k-nearest neighbor algorithm with cosine similarity. To ensure a fair comparison, we utilized the identical optimal seeds derived by the FCA method and found the nearest vectors to their mean vectors.

**Data** The co-occurrence matrix was constructed from the English Wikipedia dump (20171001)<sup>2</sup> (2.9B tokens), counted with a window of 10. We adopted PPMI (positive point-wise mutual information) as it yields the best results in the semantic task (Bullinaria and Levy, 2012). To keep the matrix size manageable, we limited the vocabulary to the 10K most frequent words.

### 4.3 Results

**Qualitative results** Table 1 presents output samples produced by the algorithm. When given *{large, huge}* as a seed, the algorithm returned *{large, huge, enormous, vast}* as the extent and *{sums, amounts, quantities}* as the intent, which constitutes a formal concept. All PPMI values within this concept exceeded 3.95. This formal concept can be labeled as "largeness" or Adjective of size, which implies that it is indeed interpretable. Interestingly, another formal concept consisting of *{large, small}* arises from the different seed instead. Similar results held for other seeds.

**Quantitative results** Table 2 shows that 61.5–84.3% of the identified extent words matched the category labels in the test sets (**Extended precision**). Furthermore, 56.3–76.8% of the words in the test sets were retrieved by the algorithm (**Recall**). Semantic categories in Battig, BLESS, and Series were more effectively captured by formal concepts than syntactic categories. We also observed that

homogeneous categories (e.g., Country) frequently formed formal concepts. With the exception of the Extended Precision metric for the Syntactic test set, our proposed method consistently achieved higher scores compared to the baseline.

The use of optimal seeds in the evaluation is justified because the objective is to measure the extent to which a mathematically identified formal concept best matches categories provided in the test set. Other non-optimal seeds return different formal concepts, which indicate the heterogeneity of human-made categories in the test set. See Appendix F for performance spread and a further discussion on the roles of seed words.

### 4.4 Analysis

The results suggest that formal concepts overlap with interpretable categories, which are defined as a set of words that human can descriptively label. Furthermore, the FCA method exhibited a notable advantage over the cosine similarity-based approach in concept retrieval. This is because the latter broadly identifies related words, whereas the former delves into specifying the underlying context. For example, given the seed words *{church, chapel}*, FCA additionally retrieves *{cathedral, shrine}*, emphasizing the context of "religious buildings." In contrast, the cosine method returns *{cathedral, catholic}* as output, failing to extract the feature of "buildings."

This advantage of FCA stems from its ability to locate mathematical structures within the matrix. Higher PPMI values discriminate the submatrix of a formal concept from its neighbors, forming a local plateau-like structure that is not necessarily captured by the cosine similarity (see a mathematical discussion in Appendix C). This insight offers a use case for the proposed algorithm.

**Disambiguating polysemy** A target word can participate in multiple formal concepts. By inputting seed words with different associations, we found that polysemous words such as *tie* and *spring* have multiple formal concepts, as shown in Table 3. We observed that separate formal concepts (e.g., clothing, match, fasten) may contain the same word (e.g., *tie*) in their extents. Three separate plateaus may share the same row as visualized in Fig. 4.

Arora et al. (2018) discovered that the embeddings of polysemous words can be decomposed as linear combinations of sense vectors. Our finding suggests that these vectors reflect separate formal

<sup>1</sup>The annotation was done by one of the authors, who is non-native but has educational experience in the U.S.

<sup>2</sup>CC BY-SA 3.0; <https://dumps.wikimedia.org/legal.html>

Seed	Formal Concept (upper:extents; lower:intents)	Th.	Category
<i>large, huge</i>	<i>large, huge, enormous, vast</i>	3.95	Adjectives of size
	<i>sums, amounts, quantities</i>		
<i>large, small</i>	<i>large, small</i>	3.47	Adjectives of scale
	<i>amounts, quantities, intestine</i>		
<i>church, temple, mosque</i>	<i>chapel, church, mosque, synagogue, temple</i>	2.85	Religious buildings
	<i>worship, jpg, ruined</i>		
<i>quicker, bigger, warmer</i>	<i>bigger, brighter, colder, cooler, heavier, hotter, louder,...</i>	2.45	Comparatives
	<i>than, considerably, deeper</i>		

Table 1: Examples of formal concepts identified from a binarized PPMI matrix. Given seed words, the algorithm returns an extent-intent pair representing a formal concept. The parameter  $k$  was set to 3. Th. means threshold.

Testset	Mtd	Pr	Ext.P	Re	LKH
Battig	FCA	<b>51.0</b>	<b>81.7</b>	<b>64.4</b>	(37.0)
	BL	36.9	67.8	50.5	-
BLESS	FCA	<b>57.8</b>	<b>84.3</b>	<b>67.0</b>	(64.7)
	BL	50.5	74.5	57.5	-
Series	FCA	<b>62.8</b>	<b>82.7</b>	<b>76.8</b>	-
	BL	53.5	75.6	67.5	-
Syntactic	FCA	<b>57.1</b>	61.5	<b>56.3</b>	-
	BL	53.6	<b>61.8</b>	54.6	-

Table 2: Average precision (**Pr**), extended precision (**Ext.P**), and recall (**Re**) over the categories ( $k = 3$ ), expressed as percentages. LKH lists % of the categories identified by Lindh-Knuutila and Honkela (2015). BL=baseline

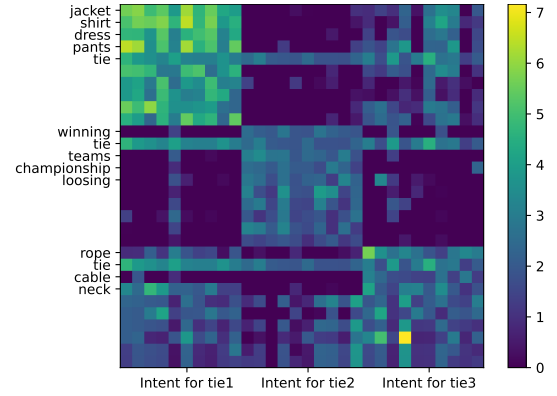


Figure 4: PPMI submatrix of three formal concepts containing the same polysemous word *tie*. For ease of visibility, the row for *tie* is presented multiple times.

concepts, and that the embeddings inherit the inner structure of the co-occurrence matrix.

## 5 Experiment 2: Applying Fuzzy FCA

### 5.1 Fuzzification of FCA

Our second application of FCA to a real-valued matrix involves the fuzzification of the crisp FCA by incorporating fuzzy set theory (Belohlavek, 2007). A fuzzy set formalizes an ambiguous set, such as "a set of tall people," by assigning a degree of membership to each element. In Appendix E, we give the definition of a fuzzy formal concept and show that it is equivalent to a rank-one submatrix under our proposed specification. Thus, the problem of finding fuzzy formal concepts can be regarded as that of identifying nonnegative rank-one submatrices in a PPMI matrix.

Because it is NP-hard to exactly decompose a matrix into nonnegative factors (Vavasis, 2010), we obtained an approximation by deploying nonnegative matrix factorization (NMF; Lee and Seung, 1999), as its  $L_1$  regularization is considered

effective in making them sparse. We controlled the sparseness so that the decomposed submatrices became disjoint. NMF decomposes  $X \in R_+^{m \times n}$  into two matrices  $W \in R_+^{m \times r}$  and  $H \in R_+^{n \times r}$  so that  $X = WH^T = \sum_{k=1}^r w_k h_k^T$ , where  $w_k$  and  $h_k$  are the  $k$ th columns of  $W$  and  $H$ , respectively. The outer product  $w_k h_k^T$  is of rank one and preferably sparse, thereby approximating a fuzzy formal concept. The loss function is  $\mathcal{L}_\alpha(W, H) = \frac{1}{2} \|X - WH^T\|_F^2 + \alpha(n\|W\|_1 + m\|H\|_1)$ . We recursively applied NMF<sup>3</sup> over three rounds—first to the PPMI matrix as in Section 4, then twice to the positive residual matrices resulting from decomposition—factorizing into  $r = 300$  components each round. Parameters for the  $L_1$  norms were set to  $\alpha = 5, 3, 1 \times 10^{-4}$  for each round.

### 5.2 Results

We manually labeled 900 rank-one submatrices by reviewing the words corresponding to the largest entries in  $w_k$  and  $h_k$  (see Appendix G.1 for details). We then classified the submatrices among

<sup>3</sup>NMF from Scikit-learn library: BSD license.

Word (sense)	Seed	Extent of Formal concept
tie1 (clothing)	<i>tie, pants, shirt</i>	<i>collar, jacket, pants, shirt, tie, wears</i>
tie2 (match)	<i>tie, teams, winning</i>	<i>championship, playoffs, teams, tie, winning</i>
tie3 (fasten)	<i>tie, cable, rope</i>	<i>cable, loose, neck, rope, tie</i>
spring1 (season)	<i>spring, autumn, month</i>	<i>autumn, cold, coldest, cooler, dry, month, rainfall,...</i>
spring2 (metal)	<i>spring, wheel, suspension</i>	<i>fitted, mounted, rear, spring, suspension, wheel, wheels</i>
spring3 (water)	<i>spring, creek, river</i>	<i>basin, brook, creek, reservoir, river, spring, stream</i>

Table 3: The extent of multiple formal concepts comprises polysemous words. The proposed algorithm is able to disambiguate these contexts in response to the seeds associated with them. The parameter  $k$  was set to 5 except for the case of spring1 ( $k = 10$ ).

four categories to assess how well the labels describe the words in each formal concept<sup>4</sup> (Table 4). Out of 900 acquired formal concepts, 95.7% were

Class	R1	R2	R3	LKH
Descriptive	182	75	73	27
Partial	56	63	48	72
Meaningful	56	150	158	2
Nonsense	6	12	21	11
Total	300	300	300	112

Table 4: Decomposed rank-one submatrices in four classes for each round, indicating how the submatrices coincide with labeled categories. Definitions are provided in Appendix G.2 and the numbers under LKH are cited from Lindh-Knuutila and Honkela (2015).

labeled descriptively or partially descriptively, or at least consisted of meaningfully related words.

### 5.3 Analysis

We found that Fuzzy FCA reveals the same formal concepts as the crisp FCA. For example, all categories listed in Table 2 also appear as rank-one submatrices. Of the 108 formal concepts identified in Experiment 1, 89 formal concepts (82.4%) were included by those found by the fuzzy method (Supplemental statistics in Appendix H). In fact, Fuzzy FCA detected more eligible words (e.g. *im-mense, massive* for Largeness, *shrine* for Religious Buildings). This observation demonstrates the robustness of FCA, as well as the correlation between the two methods.

Another interesting finding is that two types of rank-one submatrices were discovered: a *clique* type with identical rows and columns, and a *biclique* type with different rows and columns. An example of the latter is (*{explain, describe, discuss, ...}*, *{beliefs, concepts, ideas, ...}*), which represents a verb phrase for an act of communication.

<sup>4</sup>The same as the footnote 1.

## 6 Discussion

### 6.1 Why do formal concepts correspond to interpretable categories ?

As noted in Section 2, a formal concept is equivalent to a biclique, which means that the words in it are densely connected. A group of words forms a dense community if the words are repeatedly used together. Furthermore, if the same latent state always emits the same set of words, then those words are repeatedly counted as co-occurrence, thereby forming a formal concept. The random walk model of Arora et al. (2016) captures the same mechanism to generate linearly structured embeddings.

However, a latent state is not necessarily limited to a topic, i.e., a state based on thematic proximity. As revealed in Section 4, formal concepts may reflect functional proximity, e.g. the comparative. Furthermore, we observed phrasal proximity, as in a verb phrase. Thus, a broad range of semantic and syntactic patterns of word usage may be captured as a formal concept.

These results open questions about the potential relationship between formal concepts and human cognition (e.g., Wu et al., 2022 showed that fMRI patterns contain information to solve analogical reasoning), which may be the subject of future studies of semantic cognition. Our approach may provide a quantitative method to address these questions.

### 6.2 How can formal concepts be fully captured ?

We designed two methods that apply FCA to a real-valued matrix to detect interpretable formal concepts, although we do not yet have a theory to assess and relate the two methods.

In general, the challenge of FCA in applied studies is scalability stemming from computational complexity, which must be addressed when increasing the size of a co-occurrence matrix. Another

challenge is posed by heterogeneous data from large corpora. Specifically, we observed that interpretable formal concepts are detected at different threshold levels (Section 4) and by layered factorization (Section 5). The latent structures at different scales indicate that multiple formal contexts co-exist in the matrix as if they were superposed, and they were probably generated separately. Thus, the rank-one submatrices may be disjoint, superposed, or overlapping. To extract such latent structures more precisely, the algorithm must depend upon the modeling of generative processes, which is also a topic for a future study.

### 6.3 What do embeddings represent after all?

Recall that formal concepts defined as rank-one submatrices appear as components of matrix factorization  $X = WH^T$  (Section 5). While a column of  $W$  corresponds to a fuzzy set that constitutes each formal concept, a row of  $W$  is used as a word embedding. Thus, a value in each dimension of the embedding can be seen as the "coordinate" of the corresponding formal concept. The other matrix  $H$  is considered to encode attributes. The embeddings, acquired by matrix factorization or implicit factorization (Levy and Goldberg, 2014b), must inherit the structures of formal concepts, as the factor matrices can be mutually transformed.

### 6.4 Are FCA methods beneficial in practice?

The FCA method exhibited an advantage over the cosine similarity-based approach in the category completion task (Section 4). Although both methods can capture a similarity between words, a fundamental distinction lies in the subspace where these similarities are assessed. While cosine similarity utilizes the entire vector space and treats vectors as static entities, the FCA method dynamically narrows the subspace based on a given set of words. It identifies subvectors with significantly high occurrences, a task that cosine-based methods cannot perform. This merit makes FCA beneficial in various tasks such as polysemy disambiguation and concept completion (e.g., Shani et al., 2023).

Another potential benefit of FCA is its use as an analytical tool for contextualized embeddings. Dar et al., 2023 reveals that the inner representations of GPT-2 can be interpreted by projecting vectors into the vocabulary space. They report actual pairs of words processed in layers of GPT-2, some of which seem to be similar to the formal concepts identified in Experiment 2. Thus, there is a possibility that

the contextualization process can be interpreted as certain algebraic operations on formal concepts, though this is still speculative.

## 7 Related studies

Several studies demonstrated that sparse embeddings are interpretable. Murphy et al. (2012) and Biggs et al. (2008) applied nonnegative matrix factorization with a sparsity constraint to word-document co-occurrence data and discovered topics. Other studies (Faruqui et al., 2015; Park et al., 2017; Jang and Myaeng, 2017) investigated word embeddings to restore interpretability by using sparsity. We mathematically formalized the latent structure in the word co-occurrence matrix, which prior studies might have empirically detected.

FCA has been applied in linguistics (Priss, 2005), primarily for ontology. Cimiano et al. (2005) applied FCA for the automatic acquisition of taxonomies from a corpus. Moraes and Lima (2012) built a semantic structure by setting the S-V-C tuples of the annotated corpora as a formal context. Berend et al. (2018) used FCA by binarizing sparse word embedding for hypernymy discovery. In contrast to these studies, we deployed FCA to explore the structure of the matrix itself, which revealed the underlying structure of word-word co-occurrence matrices.

Gastaldi (2021) delved into the underlying mechanism of word embeddings from a linguistic-philosophical perspective and pointed out simultaneous codetermination or *bi-duality* between terms and contexts as a significant feature of language, which we believe to have successfully formalized via FCA. Our mathematical approach to interpreting co-occurrence data may shed light on the structure of language, as Bradley et al. (in press) frames language in category theory.

## 8 Summary

This study establishes a mathematical characterization of semantic relations represented as geometrical formations in a vector space, employing FCA to investigate a word co-occurrence matrix. Our experiments demonstrate that identified formal concepts align with interpretable categories. Using synthetic data, we also illustrated the emergence of hierarchical structures from word usage. Subsequent challenges include theoretical sophistication in applying FCA, exploring generative modeling, and delving into cognitive inquiries.



## 9 Limitations and risks

Our study is inherently exploratory, with the aim of communicating critical insights in a timely manner before exhaustively diving into a comprehensive analysis. Consequently, a more thorough investigation and nuanced analysis are deferred to future work, acknowledging that the current study serves as a preliminary exploration that lays the foundation for deeper scrutiny. One direction is to identify the entire formal concepts and construct a concept lattice to identify all the possible semantic relations within a vocabulary. Another direction is to investigate more thoroughly the relationship between formal concepts and linguistic concepts to quantify how far the distributional hypothesis holds.

Another limitation of this work stems from the reliance on a singular dataset for our analysis. Although our findings reveal compelling patterns within the chosen dataset, generalizability across diverse data sets remains an unexplored avenue. We anticipate similar trends in other data sets, but a comprehensive cross-validation across various sources is pending. Future research efforts should extend our methodology to encompass a wider spectrum of data sets, ensuring the robustness and applicability of our observed trends across different contexts.

The study constitutes a fundamental analysis aimed at identifying mathematical properties within linguistic statistical data, thus enhancing interpretability. Notably, no significant material risks were identified throughout the investigation and will not be seen due to the nature of the analytical approach employed.

## Acknowledgements

## References

Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2016. [A latent variable model approach to PMI-based word embeddings](#). *Transactions of the Association for Computational Linguistics*, 4:385–399.

Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2018. [Linear algebraic structure of word senses, with applications to polysemy](#). *Transactions of the Association for Computational Linguistics*, 6:483–495.

Marco Baroni and Alessandro Lenci. 2011. [How we BLESSed distributional semantic evaluation](#). In *Proceedings of the GEMS 2011 Workshop on GEometrical Models of Natural Language Semantics*, pages

1–10, Edinburgh, UK. Association for Computational Linguistics.

William F Battig and William E Montague. 1969. [Category norms of verbal items in 56 categories a replication and extension of the connecticut category norms](#). *Journal of experimental Psychology*, 80(3p2):1.

Radim Belohlavek. 2007. A note on variable threshold concept lattices: threshold-based operators are reducible to classical concept-forming operators. *Information Sciences*, 177(15):3186–3191.

Radim Belohlavek and Vilem Vychodil. 2012. [Formal concept analysis and linguistic hedges](#). *International Journal of General Systems*, 41(5):503–532.

Gábor Berend, Márton Makrai, and Péter Földiák. 2018. [300-sparsans at SemEval-2018 task 9: Hypernymy as interaction of sparse attributes](#). In *Proceedings of the 12th International Workshop on Semantic Evaluation*, pages 928–934, New Orleans, Louisiana. Association for Computational Linguistics.

Michael Biggs, Ali Ghodsi, and Stephen Vavasis. 2008. [Nonnegative matrix factorization via rank-one down-date](#). In *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, page 64–71, New York, NY, USA. Association for Computing Machinery.

Tai-Danae Bradley, Juan Luis Gastaldi, and John Terilla. in press. [The structure of meaning in language: parallel narratives in linear algebra and category theory](#). *Notices of the American Mathematical Society*, February 2024.

John A Bullinaria and Joseph P Levy. 2012. [Extracting semantic representations from word co-occurrence statistics: stop-lists, stemming, and SVD](#). *Behavior research methods*, 44:890–907.

Giampiero Chiaselotti, Davide Ciucci, and Tommaso Gentile. 2015. [Simple undirected graphs as formal contexts](#). *Formal Concept Analysis: 13th International Conference, ICFCA 2015, Nerja, Spain, June 23–26, 2015, Proceedings 13*, pages 287–302.

Philipp Cimiano, Andreas Hotho, and Steffen Staab. 2005. [Learning concept hierarchies from text corpora using formal concept analysis](#). *Journal of artificial intelligence research*, 24:305–339.

Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. 2023. [Analyzing transformers in embedding space](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16124–16170, Toronto, Canada. Association for Computational Linguistics.

Brian A Davey and Hilary A Priestley. 2002. [Introduction to lattices and order](#). Cambridge university press.

710	Manaal Faruqui, Yulia Tsvetkov, Dani Yogatama, Chris	Omer Levy and Yoav Goldberg. 2014b. <a href="#">Neural word</a>	765
711	Dyer, and Noah A. Smith. 2015. <a href="#">Sparse overcom-</a>	<a href="#">embedding as implicit matrix factorization</a> . <i>Ad-</i>	766
712	<a href="#">plete word vector representations</a> . In <i>Proceedings</i>	<i>advances in Neural Information Processing Systems</i> ,	767
713	<i>of the 53rd Annual Meeting of the Association for</i>	27.	768
714	<i>Computational Linguistics and the 7th International</i>		
715	<i>Joint Conference on Natural Language Processing</i>	Tiina Lindh-Knuutila and Timo Honkela. 2015. <a href="#">Ex-</a>	769
716	<i>(Volume 1: Long Papers)</i> , pages 1491–1500, Beijing,	<a href="#">ploratory analysis of semantic categories: comparing</a>	770
717	China. Association for Computational Linguistics.	<a href="#">data-driven and human similarity judgments</a> . <i>Com-</i>	771
		<i>putational Cognitive Science</i> , 1:1–25.	772
718	Ruiji Fu, Jiang Guo, Bing Qin, Wanxiang Che, Haifeng		
719	Wang, and Ting Liu. 2014. <a href="#">Learning semantic hier-</a>	Tomas Mikolov, Kai Chen, Gregory S. Corrado, and	773
720	<a href="#">archies via word embeddings</a> . In <i>Proceedings of the</i>	Jeffrey Dean. 2013a. <a href="#">Efficient estimation of word</a>	774
721	<i>52nd Annual Meeting of the Association for Compu-</i>	<a href="#">representations in vector space</a> . <i>International Con-</i>	775
722	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	<i>ference on Learning Representations</i> .	776
723	1199–1209, Baltimore, Maryland. Association for		
724	Computational Linguistics.	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Cor-	777
		rado, and Jeff Dean. 2013b. <a href="#">Distributed representa-</a>	778
725	Bernhard Ganter and Rudolf Wille. 2012. <a href="#">Formal con-</a>	<a href="#">tions of words and phrases and their compositionality</a> .	779
726	<a href="#">cept analysis: mathematical foundations</a> . Springer	In <i>Advances in Neural Information Processing Sys-</i>	780
727	Science & Business Media.	<i>tems</i> , volume 26. Curran Associates, Inc.	781
728	Juan Luis Gastaldi. 2021. <a href="#">Why can computers under-</a>	Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig.	782
729	<a href="#">stand natural language? the structuralist image of</a>	2013c. <a href="#">Linguistic regularities in continuous space</a>	783
730	<a href="#">language behind word embeddings</a> . <i>Philosophy &amp;</i>	<a href="#">word representations</a> . In <i>Proceedings of the 2013</i>	784
731	<i>Technology</i> , 34(1):149–214.	<i>Conference of the North American Chapter of the</i>	785
		<i>Association for Computational Linguistics: Human</i>	786
732	Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Gold-	<i>Language Technologies</i> , pages 746–751, Atlanta,	787
733	berg. 2022. <a href="#">Transformer feed-forward layers build</a>	Georgia. Association for Computational Linguistics.	788
734	<a href="#">predictions by promoting concepts in the vocabulary</a>		
735	<a href="#">space</a> . In <i>Proceedings of the 2022 Conference on</i>	Jeff Mitchell and Mirella Lapata. 2008. <a href="#">Vector-based</a>	789
736	<i>Empirical Methods in Natural Language Process-</i>	<a href="#">models of semantic composition</a> . In <i>Proceedings</i>	790
737	<i>ing</i> , pages 30–45, Abu Dhabi, United Arab Emirates.	<i>of ACL-08: HLT</i> , pages 236–244, Columbus, Ohio.	791
738	Association for Computational Linguistics.	Association for Computational Linguistics.	792
739	Zellig S Harris. 1954. <a href="#">Distributional structure</a> . <i>Word</i> ,	Sílvia Moraes and Vera Lima. 2012. <a href="#">Combining formal</a>	793
740	10(2-3):146–162.	<a href="#">concept analysis and semantic information for build-</a>	794
		<a href="#">ing ontological structures from texts : an exploratory</a>	795
741	Tatsunori B. Hashimoto, David Alvarez-Melis, and	<a href="#">study</a> . In <i>Proceedings of the Eighth International</i>	796
742	Tommi S. Jaakkola. 2016. <a href="#">Word embeddings as met-</a>	<i>Conference on Language Resources and Evaluation</i>	797
743	<a href="#">ric recovery in semantic spaces</a> . <i>Transactions of the</i>	<i>(LREC’12)</i> , pages 3653–3660, Istanbul, Turkey. Eu-	798
744	<i>Association for Computational Linguistics</i> , 4:273–	ropean Language Resources Association (ELRA).	799
745	286.		
746	Kyoung-Rok Jang and Sung-Hyon Myaeng. 2017. <a href="#">Elu-</a>	Brian Murphy, Partha Talukdar, and Tom Mitchell. 2012.	800
747	<a href="#">cidating conceptual properties from word embed-</a>	<a href="#">Learning effective and interpretable semantic mod-</a>	801
748	<a href="#">dings</a> . In <i>Proceedings of the 1st Workshop on Sense,</i>	<a href="#">els using non-negative sparse embedding</a> . In <i>Pro-</i>	802
749	<i>Concept and Entity Representations and their Appli-</i>	<i>ceedings of COLING 2012: Technical Papers</i> , pages	803
750	<i>cations</i> , pages 91–95, Valencia, Spain. Association	1933–1950, Mumbai.	804
751	for Computational Linguistics.		
752	Daniel Jurafsky and James H. Martin. 2009. <a href="#">Speech and</a>	Sungjoon Park, JinYeong Bak, and Alice Oh. 2017.	805
753	<a href="#">Language Processing (2nd Edition)</a> . Prentice-Hall,	<a href="#">Rotated word vector representations and their inter-</a>	806
754	Inc., USA.	<a href="#">pretability</a> . In <i>Proceedings of the 2017 Conference</i>	807
		<i>on Empirical Methods in Natural Language Process-</i>	808
755	Daniel D Lee and H Sebastian Seung. 1999. <a href="#">Learning</a>	<i>ing</i> , pages 401–411, Copenhagen, Denmark. Associ-	809
756	<a href="#">the parts of objects by non-negative matrix factoriza-</a>	ation for Computational Linguistics.	810
757	<a href="#">tion</a> . <i>Nature</i> , 401(6755):788–791.		
758	Alessandro Lenci. 2018. <a href="#">Distributional models of word</a>	Jonas Poelmans, Dmitry I. Ignatov, Sergei O. Kuznetsov,	811
759	<a href="#">meaning</a> . <i>Annual review of Linguistics</i> , 4:151–171.	and Guido Dedene. 2013. <a href="#">Formal concept analysis in</a>	812
		<a href="#">knowledge processing: A survey on applications</a> . <i>Ex-</i>	813
		<i>pert Systems with Applications</i> , 40(16):6538–6560.	814
760	Omer Levy and Yoav Goldberg. 2014a. <a href="#">Dependency-</a>	Uta Priss. 2005. <a href="#">Linguistic applications of formal con-</a>	815
761	<a href="#">based word embeddings</a> . In <i>Proceedings of the 52nd</i>	<a href="#">cept analysis</a> . In <i>Formal Concept Analysis: Founda-</i>	816
762	<i>Annual Meeting of the Association for Computational</i>	<i>tions and Applications</i> , page 149–160, Berlin, Hei-	817
763	<i>Linguistics (Volume 2: Short Papers)</i> , pages 302–	delberg. Springer-Verlag.	818
764	308.		

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. [Language models are unsupervised multitask learners](#). *OpenAI blog*, 1(8):9.

Sarah Roscoe, Minal Khatri, Adam Voshall, Surinder Batra, Sukhwinder Kaur, and Jitender Deogun. 2022. [Formal concept analysis applications in bioinformatics](#). *ACM Comput. Surv.*, 55(8).

Lütfi Kerem Şenel, Ihsan Utlu, Veysel Yücesoy, Aykut Koc, and Tolga Cukur. 2018. [Semantic structure and interpretability of word embeddings](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 26(10):1769–1779.

Chen Shani, Jilles Vreeken, and Dafna Shahaf. 2023. [Towards concept-aware large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13158–13170, Singapore. Association for Computational Linguistics.

Peter D Turney and Patrick Pantel. 2010. [From frequency to meaning: Vector space models of semantics](#). *Journal of artificial intelligence research*, 37:141–188.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). *Advances in neural information processing systems*, 30.

Stephen A. Vavasis. 2010. [On the complexity of non-negative matrix factorization](#). *SIAM Journal on Optimization*, 20(3):1364–1377.

Rudolf Wille. 2004. [Preconcept algebras and generalized double boolean algebras](#). In *Concept Lattices: Second International Conference on Formal Concept Analysis, ICFCA 2004, Sydney, Australia, February 23-26, 2004. Proceedings 2*, pages 1–13. Springer.

Meng-Huan Wu, Andrew J. Anderson, Robert A. Jacobs, and Rajeev D. S. Raizada. 2022. [Analogy-Related Information Can Be Accessed by Simple Addition and Subtraction of fMRI Activation Patterns, Without Participants Performing any Analogy Task](#). *Neurobiology of Language*, 3(1):1–17.

## A Toy corpus

The corpus contains 24 synthetic sentences shown in Table 5. The target words—*king*, *queen*, *man*, *woman* and their plurals—are subjects of the sentences. Each of the eight words appears with three verbs—*live-in*, *wear*, *eat*—once for each. The remaining six words—*palace*, *house*, *tie*, *dress*, *alone*, *together*—discriminate the subject words so that they are in the analogical relations of three dimensions.

<i>king live-in palace</i>	<i>kings live-in palace</i>
<i>queen live-in palace</i>	<i>queens live-in palace</i>
<i>man live-in house</i>	<i>men live-in house</i>
<i>woman live-in house</i>	<i>women live-in house</i>
<i>king wear tie</i>	<i>kings wear tie</i>
<i>queen wear dress</i>	<i>queens wear dress</i>
<i>man wear tie</i>	<i>men wear tie</i>
<i>woman wear dress</i>	<i>women wear dress</i>
<i>king eat alone</i>	<i>kings eat together</i>
<i>queen eat alone</i>	<i>queens eat together</i>
<i>man eat alone</i>	<i>men eat together</i>
<i>woman eat alone</i>	<i>women eat together</i>

Table 5: 24 sentences in the toy corpus

## B List of formal concepts

There are 28 formal concepts in the co-occurrence matrix derived from the toy corpus.

Suppose that the set of objects (target words) and the set of attributes (context words) be  $G, M$  respectively, defined as:

$$\begin{aligned} G &= \{king, man, queen, queens, \\ &\quad kings, men, queens, women, \} \\ M &= \{tie, dress, \\ &\quad palace, house, \\ &\quad alone, together\} \end{aligned}$$

Then, all the formal concepts are identified as below:

$$\begin{aligned} T &= (G, \emptyset) \\ f_1 &= (\{king, man, kings, men\}, \{tie\}) \\ f_2 &= (\{man, woman, men, women\}, \{house\}) \\ f_3 &= (\{king, queen, man, woman\}, \{alone\}) \\ f_4 &= (\{kings, queens, men, women\}, \{together\}) \\ f_5 &= (\{king, queen, kings, queens\}, \{palace\}) \\ f_6 &= (\{queen, woman, queens, women\}, \{dress\}) \\ e_1 &= (\{king, man\}, \{tie, alone\}) \\ e_2 &= (\{king, kings\}, \{tie, palace\}) \\ e_3 &= (\{man, men\}, \{tie, house\}) \\ e_4 &= (\{kings, men\}, \{tie, together\}) \\ e_5 &= (\{king, queen\}, \{palace, alone\}) \\ e_6 &= (\{man, woman\}, \{house, alone\}) \\ e_7 &= (\{kings, queens\}, \{palace, together\}) \\ e_8 &= (\{men, women\}, \{house, together\}) \\ e_9 &= (\{queen, woman\}, \{dress, alone\}) \\ e_{10} &= (\{queen, queens\}, \{palace, dress\}) \\ e_{11} &= (\{woman, women\}, \{house, dress\}) \\ e_{12} &= (\{queens, women\}, \{dress, together\}) \\ v_1 &= (\{king\}, \{tie, palace, alone\}) \\ v_2 &= (\{man\}, \{tie, house, alone\}) \\ v_3 &= (\{kings\}, \{tie, palace, together\}) \\ v_4 &= (\{men\}, \{tie, house, together\}) \\ v_5 &= (\{queen\}, \{dress, palace, alone\}) \\ v_6 &= (\{woman\}, \{dress, house, alone\}) \\ v_7 &= (\{queens\}, \{dress, palace, together\}) \\ v_8 &= (\{women\}, \{dress, house, together\}) \\ B &= (\emptyset, M) \end{aligned}$$

## C Mathematical property of an identified submatrix

Let  $I_{S\downarrow}$  and  $J_{S\uparrow}$  be subsets of rows and columns corresponding to  $S\downarrow$  and  $S\uparrow$  that define an identified formal concept, respectively, and let  $t$  be the threshold determined in the binarization method for a matrix  $X$ . Then, Algorithm 1 ensures that a submatrix  $(X_{ij})_{i \in I_{S\downarrow}, j \in J_{S\uparrow}}$  satisfies:

$$X_{ij} \geq t \quad (i \in I_{S\downarrow}, j \in J_{S\uparrow}) \quad (7)$$

$$X_{ij} < t \quad (\forall j \notin J_{S\uparrow}, \exists i \in I_{S\downarrow}) \quad (8)$$

$$X_{ij} < t \quad (\forall i \notin I_{S\downarrow}, \exists j \in J_{S\uparrow}) \quad (9)$$

Note that the submatrix of  $I_{S\downarrow} \times J_{S\uparrow}$  is discriminated from its neighbouring area. Its inner region has higher values than  $t$  (Eq. 7), whereas each of its exterior rows and columns horizontally (Eq. 8) and vertically (Eq. 9) adjacent to the submatrix contains at least one cell below the threshold. In other words, the higher entry values discriminate the submatrix of a formal concept from its neighbors, forming a local plateau-like structure that is not necessarily captured by the cosine similarity.

## D Category completion test

We used the four test sets for the category completion test: Battig, BLESS, Series and Syntactic.

Battig test (Bullinaria and Levy, 2012), originated from Battig and Montague (1969), contains 53 categories with 10 words for each category, of which we used 44 categories in the experiments, since the others have less than two words in our vocabulary of the co-occurrence matrix.

BLESS (Baroni and Lenci, 2011) contains 17 categories with 5-17, of which we used 12 categories for the same reason.

Both of Series and Syntactic are developed by the authors to supplement Battig and BLESS, which contain only common nouns. Series is hinted by Hashimoto et al. (2016) that proposed the series completion task (*penny:nickel:dime:?*) for word embeddings. Syntactic is motivated by our early finding that comparative adjectives such as *quicker*, *faster*, ... emerge as a salient formal concept with a high threshold in the binary FCA experiment. In both test sets, each category consists of 4 to 5 words, which are manually selected by one of the authors. In the development process, we partly use AI assistance<sup>5</sup> to generate a list of candidates for

<sup>5</sup><https://chat.openai.com/>



a category and its word set, by prompting with an example "Direction: *north, east, south, west*".

Examples of a category in each test set are shown below (Table 6)

Test set	Category	Word set
Battig	Metal	<i>gold, iron, lead, steel,...</i>
BLESS	Fruit	<i>apple, banana, pear,...</i>
Series	Direction	<i>north, east, south, west</i>
Syntactic	Verb (go)	<i>go, goes, went, gone</i>

Table 6: Examples of test sets

We used only the categories that contain more than or equal to three words in our vocabulary, which are listed in Table 7.

## E Fuzzification of FCA

Formally, a fuzzy set  $A$  is a function  $A : X \rightarrow L$  where  $X$  is a ground set and  $L = [0, 1]$ , which assigns the value to each member of  $X$ . A subsumption relation  $A \subseteq B$  holds if and only if  $A(x) \leq B(x)$  for all  $x \in X$ . In Fuzzy FCA, a formal concept is  $\mathbb{K} := (G, M, I, L)$ . We consider two fuzzy sets  $A \in L^G, B \in L^M$  as objects and attributes and a fuzzy relation  $I \in L^{G \times M}$ . Mathematically,  $L$  can be generalized to a *residuated lattice* that includes  $[0, 1]$  as its special case. Similar to the crisp setting, two fuzzy derivation operators  $\uparrow : L^G \rightarrow L^M$  and  $\downarrow : L^M \rightarrow L^G$  are defined as follows: For all  $m \in M$  and  $g \in G$ ,

$$A \uparrow (m) := \bigwedge_{g \in G} (A(g) \rightarrow I(g, m)) \in L \quad (10)$$

$$B \downarrow (g) := \bigwedge_{m \in M} (B(m) \rightarrow I(g, m)) \in L \quad (11)$$

Note that  $A \uparrow \in L^M, B \downarrow \in L^G$  and  $(\rightarrow) : L \times L \rightarrow L$ , which is a binary operation defined on  $L$ . In plain English, the degree to which an object  $g$  belongs to the fuzzy set  $A$  should imply the level of co-occurrence between  $g$  and an attribute  $m$ , which retrospectively should determine the degree to which the attribute  $m$  belongs to another fuzzy set  $A \uparrow$ . Then, fuzzy formal concepts are defined as a pair of fuzzy sets  $(A, B)$  where  $A \uparrow = B$  and  $B \downarrow = A$  hold as in the crisp FCA.

We need to specify operations such as  $(\rightarrow)$  to numerically compute them. Three specifications, named as Lukasiewicz, Gödel and Goguen, have already been proposed (Belohlavek and Vychodil, 2012), but instead we propose our own specification tailored to the analysis of a word co-occurrence

matrix.

$$a \rightarrow b := \begin{cases} b/a & \text{if } a > 0 \\ \top & \text{if } a = 0 \end{cases} \quad (12)$$

where  $\top$  is the greatest element in  $L$ . This specification is a slight modification of the one proposed by Goguen. The meet  $\wedge$  is numerically calculated as a minimum.

Our specification is equivalent to defining  $(A, A \uparrow)$  and  $(B \downarrow, B)$  as a rank-one submatrix. Recall that the fuzzy set  $A \in L^G$  assigns a value  $x \in L$  to the element  $g$ . Similarly, the fuzzy set  $A \uparrow \in L^M$  assigns a value  $y \in L$  to the element  $m$ . Thus, the specification in Eq. (12) ensures that  $y = xy/x$  for  $x > 0$ . This means that nonnegative entries in both fuzzy sets  $A, A \uparrow$  constitute a rank-one submatrix.

<b>Battig</b>	<b>BLESS</b>	<b>Series</b>	<b>Syntactic</b>
Disease	Ground mammal	Emotion	Demonstrative adverb
Metal	Furniture	Season	Comparative adjective
Carpenter's tool	Tool	Sea	Preposition
Crime	Container	Great lakes	Verb conjugation
Substance for flavoring food	Fruit	Direction	Manner adverb
Elective Office	Vehicle	Art form	Adverb of frequency
Toy	Appliance	Part of a tree	Personal pronoun
Weapon	Weapon	Book part	Linking verb
Member of clergy	Musical instrument	Continent	Demonstrative determiner
Four-footed animal	Building	Movie genre	Coordinating conjunction
Nonalcoholic beverages	Clothing	Number	Adjective of taste
Building for religious services	Bird	US president	Possessive pronoun
Precious stone		Stage of life	Frequency adverb
Part of human body		Planet	Quantitative determiner
Fruit		Weekday	Subordinating conjunction
Sport		Music genre	Action verb
Part of a building		Natural disaster	Modal auxiliary
Male's first name		Decathlon	Total pronoun
Relative		Family	Adjective of size
Human dwelling		Ocean	Interrogative pronoun
Insect		Adverb of time	Article
Type of fuel		Month	Totality adverb
Music instrument		Communication act	Verb conjugation
Furniture		Match	
Ship		Religion	
Kind of money		Time of day	
Color		Writing	
Kind of cloth		Style of architecture	
Unit of distance		Midwest U.S. state	
Type of music			
City			
Country			
Reading material			
Military title			
Natural earth formation			
Unit of time			
Part of speech			
Kitchen utensil			
Vehicle			
Science			
Weather phenomenon			
Occupation or profession			
Bird			

Table 7: Used categories of the test sets

## F Role of seed words and performance spread

In Algorithm 1 (Section 4), seed words are required as input to the algorithm since a formal concept is detected as a closed set containing those seed words. A closed set is the fixed point of a closure operator. In this algorithm, any set of words can be a seed. If seed words are randomly chosen, then the algorithm will find and return a formal concept in an unsupervised way. Different sets of words return the same formal concept if all of the used words belong to the same closed set, and derive different formal concepts otherwise due to the mathematical property of a closure operator and closure system.

In the category completion test in Experiment 1, all possible pairs from the word list of the same category are used as seeds (2 words) to derive a formal concept. Therefore, different formal concepts can be identified by a chosen pair of words from the test set.

Table 8 shows the statistics of the distribution of recalls calculated for all combinations of two words, reflecting the degree to which the FCA retrieved elements from the test set based on seed choices. To verify, cross-reference the numbers in the Max column with the corresponding recalls reported in Table 2 where the optimal seed pairs were selected. In cases where a word set within a specific category in the test set comprises 10 words, there exist  $_{10}C_2 = 45$  possible pairs. We computed the minimum, maximum, and median values for each category, subsequently averaging them between categories for each test set and the entire dataset.

These statistics suggest that employing a "right seed" (optimal pair) results in the formal concept covering 56.3% (Syntactic) to 76.8% (Series). On the contrary, the use of a different seed may yield a distinct formal concept. This divergence can be attributed to the non-cohesive nature of word groups within the test set.

For instance, the "Occupation or profession" category of Battig comprises words such as *doctor*, *lawyer*, *teacher*, *engineer*, *professor*, *carpenter*, *salesman*, *nurse*, *psychologist* (with one word omitted due to limited vocabularies in the matrix). Notably, the FCA found that the maximum formal concept within this category is only four words: *lawyer*, *nurse*, *psychologist*, *teacher*, which seem to represent the "profession" part of the category.

Test set	Min	Max	Median
Battig	33.8	64.4	37.3
BLESS	36.5	67.0	43.9
Series	49.2	76.8	53.6
Syntactic	46.1	56.3	48.6

Table 8: Spread % of Recall over different choice of seeds

## G Decomposition by NMF

### G.1 Decomposed submatrices by NMF

We applied NMF recursively in three rounds. In the first round, we decomposed the PPMI matrix as in  $X_0 \approx W_1 H_1^T$  into 300 components ( $\alpha = 0.0005$ ). In the second round, we applied NMF to the positive residual matrix after the first decomposition:  $X_1 := \max(X_0 - W_1 H_1^T, 0)$  as decomposed as in  $X_1 \approx W_2 H_2^T$  ( $\alpha = 0.0003$ ). In the third round, the residual matrix  $X_2 := \max(X_1 - W_2 H_2^T, 0)$  was decomposed into  $X_2 \approx W_3 H_3^T$  ( $\alpha = 0.0001$ ). Note that each component (rank-one matrix)  $w_k h_k^T$  was forced to be sparse by  $L_1$  regularization. Thus, their nonnegative rows and columns make a non-negative rank-one submatrix, which we regard as a fuzzy formal concept.

The components derived in the first round were indexed from 1 to 300. Similarly, those in the second round were indexed from 301 to 600, and the ones from the third round were indexed from 601 to 900. We ordered each component according to the Frobenius norm within each round. Therefore, the smaller ID number implies that the submatrix has a greater norm in each round.

Samples of the components are presented in Table 9. The class was evaluated by one of the authors according to the definition given in Appendix G.2. The author also labeled a category from the words that comprise the submatrix  $w_k h_k^T$ . More specifically, for each vector  $w_k$  and  $h_k$ , we picked 20 words that correspond to the largest elements in the vectors, respectively. In Table 9, the only four top words are presented for both  $w_k$  as extents and  $h_k$  as intents. For ease of visibility, categories were labeled with more general expressions, although they could be labeled with more focused category names.

Table 10 shows a supplemental analysis of the type of relatedness between words participating in each submatrix.

ID	Class	Category	Extents (top 4 words)	Intent (top 4 words)
2	D	Geography	<i>iran, kerman, khorasan, province</i>	<i>iran, kerman, khorasan, province</i>
5	N	None	<i>pineapples, tasteful, lilongwe, unimpressive</i>	<i>dawn, windsor, batting, relegation</i>
8	D	Music	<i>chart, charts, billboard, singles</i>	<i>chart, charts, billboard, singles</i>
14	D	Sports	<i>discus, javelin, jump, hurdles</i>	<i>discus, javelin, jump, hurdles</i>
22	D	Education	<i>degree, bachelor, doctorate, laude</i>	<i>degree, bachelor, doctorate, laude</i>
35	D	Diplomacy	<i>embassy, ambassador, diplomatic, relations</i>	<i>turkmenistan, tajikistan, kyrgyzstan, uzbekistan</i>
46	D	Sports	<i>baseman, pitcher, outfielder, shortstop</i>	<i>baseman, pitcher, outfielder, shortstop</i>
89	D	Religion	<i>rabbi, yeshiva, synagogue, hebrew</i>	<i>rabbi, yeshiva, synagogue, hebrew</i>
90	D	US states	<i>idaho, montana, dakota, wyoming</i>	<i>idaho, montana, dakota, wyoming</i>
95	D	Climates	<i>cyclone, hurricane, storm, typhoon</i>	<i>cyclone, hurricane, storm, typhoon</i>
98	D	Politics	<i>polling, votes, voters, vote</i>	<i>polling, votes, voters, vote</i>
102	D	Phrases	<i>increases, decreases, decrease, increase</i>	<i>temperature, concentrations, accuracy, velocity</i>
104	D	Politics	<i>incumbent, reelection, democrat, republican</i>	<i>incumbent, reelection, democrat, republican</i>
116	P	Medical	<i>ligament, knee, ankle, injury</i>	<i>ligament, knee, ankle, injury</i>
125	P	Career	<i>postdoctoral, professor, adjunct, emeritus</i>	<i>postdoctoral, professor, adjunct, emeritus</i>
137	P	TV show	<i>starring, roommate, daughters, actress</i>	<i>jennifer, laura, jessica, nicole</i>
146	P	Legal	<i>convicted, guilty, sentenced, imprisonment</i>	<i>convicted, guilty, sentenced, imprisonment</i>
147	P	History	<i>nazi, nazis, deported, camps</i>	<i>nazi, nazis, deported, camps</i>
159	P	Geography	<i>mountain, peaks, summit, mountains</i>	<i>mountain, peaks, summit, mountains</i>
160	M	Expression	<i>acclaim, garnered, reviews critical</i>	<i>garnered, acclaim, reviews, critical</i>
165	P	Expression	<i>regain, recover, conquer, attract</i>	<i>trying, attempting, attempt, attempts</i>
181	M	Expression	<i>tasked, thereby, prevented, intention</i>	<i>securing, obtaining, capturing, creating</i>
184	M	Expression	<i>lied, intentions, poisoned, whereabouts</i>	<i>reveals, realizes, believing, realises</i>
192	D	Music	<i>punk, hop, hip, folk</i>	<i>punk, hop, hip, folk</i>
210	D	Religion	<i>christianity, catholicism, islam, beliefs</i>	<i>christianity, catholicism, islam, beliefs</i>
212	M	Religion	<i>you, think, really, know</i>	<i>you, think, really, know</i>
214	M	Syntactic	<i>various, numerous, several, these</i>	<i>genera, disciplines, locations, dialects</i>
237	D	Comparative	<i>faster, stronger, heavier, than</i>	<i>faster, stronger, heavier, than</i>
239	D	Politics	<i>obama, barack, reagan, clinton</i>	<i>obama, barack, reagan, clinton</i>
313	D	Religion	<i>quantities, amounts, sums, amassed</i>	<i>enormous, huge, immense, considerable</i>
329	P	Time	<i>spends, spend, spent, spending</i>	<i>summers, much, time, remainder</i>
341	P	Geography	<i>maui, oahu, hawaii, honolulu</i>	<i>maui, oahu, hawaii, honolulu</i>
370	D	Unit	<i>millions, billions, million, billion</i>	<i>millions, billions, million, dollars</i>
405	M	Linguistics	<i>vowel, vowels, stressed, accent</i>	<i>vowel, vowels, stressed, accent</i>
408	P	Travel	<i>immigration, nationals, emigration, citizen</i>	<i>immigration, nationals, emigration, citizen</i>
419	D	Proposal	<i>proposal, offer, invitation, plea</i>	<i>rejected, accepted, rejects, accepting</i>
431	M	Expression	<i>poorly, properly, carefully, fully</i>	<i>handled, treated, understood, trained</i>
435	D	Buildings	<i>housed, built, constructed, build</i>	<i>synagogue, mosque, mansion, convent</i>
484	D	Auxiliary	<i>did, does, doesn, didn</i>	<i>speak, exist, suffer, appear</i>
507	D	Movement	<i>down, forth, out, into</i>	<i>fell, put, falling, fallen</i>
514	D	War	<i>pistol, revolver, magnum, rifle</i>	<i>pistol, revolver, magnum, rifle</i>
517	M	Plants	<i>botanical, zoological, garden, gardens</i>	<i>botanical, zoological, garden, gardens</i>
577	P	Expression	<i>totally, completely, virtually, almost</i>	<i>totally, virtually, completely, vanished</i>
605	D	Number	<i>vii, ix, viii, xiii</i>	<i>fantasy, corps, intensity, chapter</i>
626	P	Accounting	<i>collect, collecting, exception, collected</i>	<i>taxes, debt, debts, fees</i>
645	D	Month	<i>june, july, august, september</i>	<i>premiered, consecrated, baptised, inaugurated</i>
667	D	Expression	<i>taking, take, taken, takes</i>	<i>hostage, advantage, seriously, refuge</i>
669	D	Geography	<i>gaza, palestinians, palestinian, israeli</i>	<i>strip, gaza, rockets, barrier</i>
679	D	Geography	<i>colombian, venezuelan, peruvian, chilean</i>	<i>peso, divisi, primera, aut</i>
774	P	IT	<i>java, server, windows, software</i>	<i>java, server, windows, software</i>
781	D	Expression	<i>bought, purchased, buying, buy</i>	<i>shares, stake, tickets, tracts</i>
784	M	Marketing	<i>advertising, commercials, campaigns, marketing</i>	<i>advertising, commercials, campaigns, marketing</i>
804	M	Expression	<i>about, detail, matters, topics</i>	<i>discuss, discussed, discussing, discusses</i>
855	M	Expression	<i>heavily, originally, by, recently</i>	<i>influenced, inspired, invented, borrowed</i>
864	D	Expression	<i>currently, presently, still, today</i>	<i>currently, resides, owns, produces</i>
874	P	Expression	<i>launching, pursued, launched, developed</i>	<i>ventures, venture, scheme, initiative</i>

Table 9: Samples of decomposed submatrices labeled with a category name. Classes are abbreviated; D:Descriptive, P:Partial, M:Meaningful, N:Nonsense



Proximity	R1	R2	R3
Categorical	74	64	53
Contextual	171	147	148
Combinatorial	41	59	62
Syntactic	9	18	19
None	5	12	18
Total	300	300	300

Table 10: Proximity types of word relations in each NMF-decomposed component. Categorical: words are in the same category, Contextual: words are related in a shared context, Combinatorial: words are a part of possible phrases, i.e., paradigmatic, Syntactic: words are in the same syntactic category.

## G.2 Types of qualitative classes

The set of words corresponding to the largest dimensions within each component is classified into four qualitative classes, as in the below definition (Table 11), following Lindh-Knuutila and Honkela (2015). These classes indicate how well an identified formal concept (a rank-one matrix) is interpretable as a category.

Class	Description
Descriptive	Words are related in some way, and the majority label given is as descriptive as possible of the words in the set.
Partial	Words are related in some way, and the majority label is somewhat descriptive, but a more descriptive account can be easily given.
Meaningful	Words are related, but no majority label describes the words.
Nonsense	There is no majority label, nor is there any perceived relation between the words in the set.

Table 11: Definition of qualitative classes assessing how well the labels describe the words in each formal concept. (Lindh-Knuutila and Honkela, 2015)

## H Overlap of two FCA methods

In addition, we performed an analysis of set overlap at the word level. For each of the 89 groups, we calculated the set overlap using the Jaccard index, which is defined as the number of words in the intersection divided by the number in the union. The results are presented in Table 12 as percentages.

Min	Max	Mean	Median
1.4	64.5	23.7	20.0

Table 12: Jaccard index between the corresponding formal concepts of Binarization method and Fuzzy method over 89 categories