

# Unsupervised Estimation of Subjective Content Descriptions

Magnus Bender\*, Tanya Braun†, Ralf Möller\* and Marcel Gehrke\*

\*University of Lübeck  
Institute of Information Systems  
Ratzeburger Allee 160, 23562 Lübeck  
{bender, moeller, gehrke}@ifis.uni-luebeck.de

†University of Münster  
Computer Science Department  
Einsteinstr. 62, 48155 Münster  
tanya.braun@uni-muenster.de

**Abstract**—An agent in pursuit of a task may work with a corpus containing text documents. One possible task of the agent is to retrieve documents of similar content and highlight relevant locations in retrieved documents. To perform information retrieval on the corpus, the agent may need additional data associated with the documents. Subjective Content Descriptions (SCDs) provide additional location-specific data for text documents. However, the agent needs SCDs referencing sentences of similar content across various documents in the corpus and most text documents are not associated with SCDs. Therefore, this paper presents UESM, an unsupervised approach to estimate SCDs for text documents, i.e., to associate any corpus with SCDs. In an evaluation, we show that the performance of UESM is on par with latent Dirichlet allocation, while UESM provides SCDs referencing sentences of similar content.

## I. INTRODUCTION

An agent in pursuit of a task, explicitly or implicitly defined, may work with a corpus of text documents as a reference library. From an agent-theoretic perspective, an agent is a rational, autonomous unit acting in a world fulfilling a defined task, e.g., providing document retrieval services given requests from users. We assume that the corpus represents the context of the task defined by the users of the retrieval service, since document retrieval is not an end in itself. Further, documents in a given corpus might be associated with additional location-specific data making the content nearby the location explicit by providing descriptions, references, or explanations. We refer to these additional location-specific data as Subjective Content Descriptions (SCDs) [1].

Associating documents with SCDs support agents in the task of information retrieval (IR). Returning to our agent providing an document retrieval service, it is valuable for the agent to have for each sentence references to similar sentences in the documents across its corpus. Such references can be modeled by SCDs being associated with the sentences. In our understanding, each SCD represents a concept or topic mentioned in the corpus. Whereby, each SCD’s concept is implicitly defined by the content of the sentences referenced by each SCD. Thus, the sentences of an SCD are similar.

Let us have a closer look why sentence similarity can help with the IR task. Given a corpus of three documents  $\{d_1, d_2, d_3\}$ , dealing with three different car

models, and each document consisting of ten sentences  $\{d_i = (s_1^{d_i}, s_2^{d_i}, \dots, s_{10}^{d_i})\}_{i=1}^3$ , the agent could be interested in the SCDs  $t_1$  and  $t_2$ . SCD  $t_1$  is referencing the sentences  $s_2^{d_2}$  and  $s_8^{d_1}$  because both sentences are about the engine’s horsepower. Thus,  $t_1$  represents the concept *engine power*. Furthermore, SCD  $t_2$  represents the concept *car manufacturer* because it references two sentences  $s_7^{d_3}$  and  $s_2^{d_1}$  about the car’s manufacturer. Then, the agent returns  $d_1$  and highlights  $s_8^{d_1}$  answering a request about sentence  $s_2^{d_2}$  because both sentences cover *engine power*. So for this request, the additional information of  $t_1$  turned out to be useful.

However, most corpora are not associated with SCDs nor contain references to sentences representing the same concepts. In a first step, we are interested to identifying similar sentences—preferably in an unsupervised manner. Then, using the identified similar sentences, we form SCDs, where each SCD represents a different concept in the corpus and references multiple locations in text documents of the corpus. Later, we enrich the SCDs in an semi-supervised fashion.

In addition to the applications of SCDs just described, SCDs can support agents by performing the following tasks:

- (i) Estimating SCDs for a single previously unseen text document using the Most Probably Suited SCD (MPS<sup>2</sup>CD) algorithm [2],
- (ii) classifying a text document as related, extended, revised, or unrelated to a corpus [2],
- (iii) moving the SCDs from one corpus to another similar corpus by adapting the SCDs’ domain [3],
- (iv) separating SCDs and actual content being interleaved in text documents [4],
- (v) enriching SCDs in a corpus already sparsely associated with SCDs [5], or
- (vi) detecting complementary documents to a corpus [6].

Whereby, all the approaches have in common that initially a corpus associated with SCDs is needed. More precisely, an SCD matrix—built on the corpus associated with SCDs—is used to model the SCDs and their locations in the documents of the corpus. Therefore, this paper presents Unsupervised Estimator for SCD Matrices (UESM), an approach to estimate an SCD matrix for any corpus in an unsupervised manner.

Thus, UESM associates any corpus of text documents with SCDs. Mainly, UESM detects similar concepts referenced in the text documents of the corpus and then forms an SCD, which groups all occurrences of the same concept.

UESM works on any corpus of text documents, in contrast to the tentative approaches using Wiktionary<sup>1</sup> or OpenIE [7], e.g., used in the evaluations of [4] or [3], respectively. Additionally, supervised learning, e.g., using Wiktionary, is not always possible. For example, Wiktionary can only be used for corpora about fields where it contains enough definitions.

Summarized, UESM allows for estimating an SCD matrix and thus SCDs for any corpus. Hence, UESM enables the above described approaches to be applied to any corpus without need for SCDs on the corpus in beforehand, as UESM provides the required SCD matrix.

Moreover, each estimated SCD represents a topic of the corpus, which is why an SCD matrix can be interpreted as a topic model of the corpus. A well-known topic modelling technique is Latent Dirichlet Allocation (LDA) [8]. In contrast to LDA, UESM associates each sentence in the corpus with an SCD, while LDA associates each single word with a topic and each text document with a topic distribution. An SCD consists of multiple referenced sentences in the corpus and the sentences' overall word distribution, while LDA's topics consist of a distribution of words associated with each topic. Hence, associating SCDs with sentences instead of words or text documents is the important difference.

The remainder of this paper is structured as follows: First, we recap the basics of SCDs, especially of the SCD matrix. Second, we describe the problem of estimating SCDs in an unsupervised manner and our solution UESM. We introduce three methods for UESM, namely a greedy, a K-Means [9], and a DBSCAN [10] based method. Third, we compare UESM with its three methods in an evaluation against the well-know LDA. Finally, we look at related work and conclude afterwards.

## II. PRELIMINARIES

This section specifies notations and recaps the basics of SCDs.

### A. Notations

First, we formalize our setting of a corpus.

- A word  $w_i$  is a basic unit of discrete data from a vocabulary  $\mathcal{V} = \{w_1, \dots, w_L\}$ ,  $L \in \mathbb{N}$ .
- A sentence  $s$  is defined as a sequence of words  $s = (w_1, \dots, w_N)$ ,  $N \in \mathbb{N}$ , where each word  $w_i \in s$  is an element of vocabulary  $\mathcal{V}$ . Commonly, a sentence is terminated by punctuation symbols like “.”, “!”, or “?”.
- A document  $d$  is defined as a sequence of sentences  $d = (s_1^d, \dots, s_M^d)$ ,  $M \in \mathbb{N}$ .
- A corpus  $\mathcal{D}$  represents a set of documents  $\{d_1, \dots, d_D\}$ ,  $D \in \mathbb{N}$ .

<sup>1</sup><https://www.wiktionary.org/>

- An SCD  $t$  is a tuple of the SCD's additional data  $\mathcal{C}$  and the referenced sentences  $\{s_1, \dots, s_S\}$ ,  $S \in \mathbb{N}$ . Thus, each SCD references sentences in documents of  $\mathcal{D}$ , while in the opposite direction a sentence is associated with an SCD.
- A sentence associated with an SCD is called SCD window, inspired by a tumbling window moving over the words of a document. Generally, an SCD window might not be equal to a sentence and may be a subsequence of a sentence or the concatenated subsequences of two sentences, too. Even though, in this paper, an SCD window always equals a sentence.
- For a corpus  $\mathcal{D}$  there exists a set  $g$  called SCD set containing  $K$  associated SCDs

$$g(\mathcal{D}) = \left\{ t_j = \left( \mathcal{C}_j, \bigcup_{d \in \mathcal{D}} \{s_1^d, \dots, s_S^d\} \right) \right\}_{j=1}^K.$$

Given a document  $d \in \mathcal{D}$ , the term  $g(d)$  refers to the set of SCDs associated with sentences from document  $d$ .

- Each word  $w_i \in s^d$  is associated with an influence value  $I(w_i, s^d)$  representing the relevance of  $w_i$  in the sentence  $s^d$ . For example, the closer  $w_i$  is positioned to the object of the sentence  $s^d$ , the higher its corresponding influence value  $I(w_i, s^d)$ . The influence value is chosen according to the task and might be distributed binomial, linear, or constant.

### B. Subjective Content Descriptions

SCDs provide additional location-specific data for documents [1]. The data provided by SCDs may be of various types, like additional definitions or links to knowledge graphs. However, in this paper we do not focus on the additional data, instead we focus on how to determine which sentences belong to one SCD.

Kuhr et al. use an SCD-word distribution represented by a matrix when working with SCDs [1]. The SCD-word distribution matrix, in short SCD matrix, can be interpreted as a generative model. A generative model for SCDs is characterized by the assumption that the SCDs generate the words of the documents. We assume that each SCD shows a specific distribution of words of the referenced sentences in the documents.

Before we describe UESM, we outline the details of SCD matrices and an algorithm training an SCD matrix  $\delta(\mathcal{D})$  for a corpus  $\mathcal{D}$  given the SCD set  $g(\mathcal{D})$  in a supervised manner.

The SCD matrix  $\delta(\mathcal{D})$  models the distributions of words for all SCDs  $g(\mathcal{D})$  of a corpus  $\mathcal{D}$  and is structured as follows:

$$\delta(\mathcal{D}) = \begin{matrix} & w_1 & w_2 & w_3 & \cdots & w_L \\ \begin{matrix} t_1 \\ t_2 \\ \vdots \\ t_K \end{matrix} & \begin{pmatrix} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,L} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,L} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{M,1} & v_{K,2} & v_{K,3} & \cdots & v_{K,L} \end{pmatrix} \end{matrix}$$

---

**Algorithm 1** Supervised estimation of SCD matrices  $\delta(\mathcal{D})$ 

---

```
1: function BUILDMATRIX( $\mathcal{D}$ ,  $g(\mathcal{D})$ )
2:   Input: Corpus  $\mathcal{D}$ ; Set of SCDs  $g(\mathcal{D})$ 
3:   Output: SCD-word distribution matrix  $\delta(\mathcal{D})$ 
4:   Initialize an  $K \times L$  matrix  $\delta(\mathcal{D})$  with zeros
5:   for each document  $d \in \mathcal{D}$  do
6:     for each SCD  $t = (\mathcal{C}, \{s_1^d, \dots, s_S^d\}) \in g(d)$  do
7:       for  $j = 1, \dots, S$  do  $\triangleright$  Iterate over sentences
8:         for each word  $w_i \in s_j^d$  do
9:            $\delta(\mathcal{D})[t][w_i] += I(w_i, s_j^d)$ 
10:  return  $\delta(\mathcal{D})$ 
```

---

The SCD matrix consists of  $K$  rows, one for each SCD in  $g(\mathcal{D})$ , and each row contains the word probability distribution for the SCD. Therefore, the SCD matrix has  $L$  columns, one for each word in the vocabulary of the corresponding corpus.

The supervised estimation of an SCD matrix is described in Algorithm 1. Given a corpus  $\mathcal{D}$ , the algorithm iterates over each document  $d$  in the corpus and the document's SCDs. For each associated SCD  $t$ , the referenced sentences  $s_1^d, \dots, s_S^d$  are used to update the SCD matrix. Thereby, the row of the matrix representing SCD  $t$  gets incremented for each word in each sentence by each word's influence value.

Finally, the SCD matrix needs to be normalized row-wise to meet the requirements of a probability distribution. However, we skip the normalization because multiple calculations on small decimal values on a computer reduce the accuracy. Later, we use the cosine similarity with the rows of the matrix and the cosine similarity does a normalization by definition. Thus, by skipping the normalization, we save computational resources and get slightly more accurate results.

Next, we present UESM, which estimates an SCD matrix  $\delta(\mathcal{D})$  without needing the SCD set  $g(\mathcal{D})$ .

### III. UNSUPERVISED ESTIMATION OF SCD MATRICES

This section introduces UESM, the Unsupervised Estimator for SCD Matrices. Algorithm 2 outlines UESM. The SCD matrix represents in its rows each SCD found in the corpus. Whereby, each row contains the word distribution of the sentences associated with the row's SCD. UESM is also a topic estimation algorithm because each SCD represents a concept in the corpus and the SCD references the sentences dealing about this concept.

In the beginning, UESM only gets a corpus of text documents for which the SCD matrix has to be estimated. Commonly, a sentence is associated with an SCD and each SCD references one or multiple sentences. UESM initially starts by associating each sentence to one unique SCD. The SCD's word distribution of each SCD then only contains the words of the referenced sentence. Lines 10 - 14 of Algorithm 2 show how to create this initial SCD matrix, which consists of a row for each sentence in the document's corpus. The word distributions are calculated using the influence value the same way as in Algorithm 1.

---

**Algorithm 2** Unsupervised Estimator for SCD Matrices  $\delta(\mathcal{D})$ 

---

```
1: function UESM( $\mathcal{D}$ ,  $m$ ,  $[\theta, ]$  [ $K, ]$  [ $\varepsilon, ms$ ])
2:   Input: Corpus  $\mathcal{D}$ ; Method with hyperparameters, i.e.,
3:      $m =$  Greedy and threshold  $\theta$ ,
4:      $m =$  K-Means and number of SCDs  $K$ , or
5:      $m =$  DBSCAN, distance  $\varepsilon$ , and threshold  $ms$ 
6:   Output: SCD-word distribution matrix  $\delta(\mathcal{D})$ 
7:   Initialize an  $(\sum_{d \in \mathcal{D}} M^d) \times L$  matrix  $\delta(\mathcal{D})$  with zeros
8:    $l \leftarrow 0$ 
9:    $\triangleright$  Build initial SCD matrix
10:  for each document  $d \in \mathcal{D}$  do
11:    for each sentence  $s^d \in d$  do
12:      for each word  $w_i \in s^d$  do
13:         $\delta(\mathcal{D})[l][w_i] += I(w_i, s^d)$ 
14:       $l \leftarrow l + 1$ 
15:     $\triangleright$  Use method  $m$  to merge rows
16:  if  $m =$  Greedy then
17:    repeat  $\triangleright$  Detect similar rows and merge
18:       $(r_i, r_j) \leftarrow$  MOSTSIMILARROWS( $\delta(\mathcal{D})$ )
19:       $\delta(\mathcal{D})[r_i] \leftarrow \delta(\mathcal{D})[r_i] + \delta(\mathcal{D})[r_j]$   $\triangleright$  Sum rows
20:       $\delta(\mathcal{D})[r_j] \leftarrow Nil$   $\triangleright$  Delete row
21:    until SIMILARITY( $r_i, r_j$ )  $< \theta$ 
22:  else
23:     $\triangleright$  Create clusters of similar rows
24:  if  $m =$  K-Means then
25:     $clusters \leftarrow$  KMEANS( $\delta(\mathcal{D})$ ,  $K$ )
26:  else
27:     $clusters \leftarrow$  DBSCAN( $\delta(\mathcal{D})$ ,  $\varepsilon$ ,  $ms$ )
28:  for each cluster  $c \in clusters$  do
29:     $\triangleright$  Create sum of all cluster's rows in first row
30:     $r_i \leftarrow$  FIRSTROW( $c$ )
31:     $\delta(\mathcal{D})[r_i] \leftarrow \sum_{r_j \in c} \delta(\mathcal{D})[r_j]$ 
32:    for each row  $r_j \in c$  do
33:      if  $r_i \neq r_j$  then  $\triangleright$  Delete all non-first rows
34:         $\delta(\mathcal{D})[r_j] \leftarrow Nil$ 
35:  return  $\delta(\mathcal{D})$ 
```

---

Next, the sentences representing the same concept have to be found and grouped into one SCD. Thereby, we differentiate between three methods to detect similar rows in the initial SCD matrix. Lines 16 - 34 of Algorithm 2 show the three methods and how the rows are merged. The main idea of merging two rows is to sum up the quantities of each word in both distributions of words and deleting the second row from the matrix.

To identify similar sentences, UESM has three different methods. The first is a greedy approach followed by two well-known clustering techniques, K-Means and DBSCAN.

a) *Greedy by Similarity:* The first method greedily selects the next two rows to merge. It calculates the cosine similarity between all rows, containing the word distributions, in the matrix and always merges the two most similar rows. This is repeated until the similarity between the two most similar rows is below the threshold  $\theta$  (Algorithm 2 lines 17

- 21). Thus, with a lower threshold less SCDs with more referenced sentences each will be estimated and a higher threshold leads to more SCDs with less referenced sentences.

The calculation of the cosine similarity between all rows is realized as a matrix multiplication:

$$S_{\delta(\mathcal{D})} = \frac{\delta(\mathcal{D}) \cdot \delta(\mathcal{D})^T}{\|\delta(\mathcal{D})\|_2 \cdot \|\delta(\mathcal{D})\|_2^T}$$

The numerator represents the dot product between each row of the matrix to each other and the denominator contains the lengths of each row to normalize the matrix's rows, as  $\|v\|_2$  represents a vector of the euclidean norm of each row in  $v$  and the symbol  $\cdot$  the matrix multiplication. Numerator and denominator are matrices of size  $K \times K$  each, which are then divided element-wise to form the cosine similarity matrix  $S_{\delta(\mathcal{D})}$ . After doing so,  $S_{\delta(\mathcal{D})}$  contains the cosine similarity between each pair of rows in the matrix  $\delta(\mathcal{D})$ . The two most similar rows in  $\delta(\mathcal{D})$  can now be identified by searching for the highest value in  $S_{\delta(\mathcal{D})}$ , of course without searching the diagonal. Row and column index of the highest value in  $S_{\delta(\mathcal{D})}$  represent the most similar rows in  $\delta(\mathcal{D})$ .

Matrix multiplications on huge matrices can be computationally expensive. In case of the SCD matrix, it is a sparse matrix and sparse matrix multiplication is reasonably fast. Additionally, the euclidean norms of the rows can be cached and updated partially for the changed rows, only.

*b) K-Means:* One well-know clustering techniques is K-Means [9]. We will not get into the details how K-Means works, but focus on how to apply K-Means. K-Means is initialized with  $K$  centroids whereof each centroid represents a cluster. Each point is assigned the nearest centroid in terms of the euclidean distance using a vector representation of the point. Iteratively the clusters are optimized by aligning each centroid in the center of all the points contained in the centroid's cluster.

We run K-Means on the rows of the SCD matrix to detect clusters of similar rows in the initial SCD matrix. Each row represents a point and the word distribution is the vector representation of this point. After K-Means is finished, we merge the rows of the matrix included in the same cluster (Algorithm 2 lines 28 - 34). Hence, the number of clusters is equal to the number of SCDs in the end. As hyperparameter the number of SCDs to estimate  $K$  is specified. Additionally, a factor to multiply with the initial number of sentences in the corpus or a technique to estimate a *good* number of cluster for K-Means on the corpus may be used [11].

*c) DBSCAN:* Another well-know clustering technique is DBSCAN [10]. In contrast to K-Means, DBSCAN is able to detect concave structures in data and works density based. DBSCAN clusters two points together if both are in a neighborhood, the distance making up a neighborhood is defined by the hyperparameter  $\epsilon$ . A cluster then grows by adding all points in the neighborhood to the same cluster. Additionally, there is a minimum samples threshold  $ms$  which defines the minimum number of points needed to form a cluster.

We run DBSCAN on the cosine similarity matrix  $S_{\delta(\mathcal{D})}$  and again merge the rows of the matrix included in the same cluster (Algorithm 2 lines 28 - 34).

Comparing the three methods, when using K-Means the number of SCDs to estimate  $K$  has to be specified in beforehand. The greedy method and DBSCAN determine the number of SCDs on their own. Though, the greedy method needs a similarity threshold  $\theta$  and DBSCAN  $\epsilon$  and the minimum samples threshold  $ms$ .

We use DBSCAN and K-Means because each method represents a clustering method following a different approach, i.e., density based and distance based clustering. We can not predict which clustering method works better for a given corpus. As typical for greedy methods, we expect the greedy method working well for higher thresholds and more SCDs to estimate, while for smaller thresholds and a small number of SCDs, the greedy method will miss the global optimum.

We have now proposed UESM including how to determine similar sentences based on three different methods greedy, K-Means, and DBSCAN. Next, we describe and discuss the workflow, dataset, and implementation used in our evaluation along with the results comparing UESM against LDA.

#### IV. DISCUSSION AND EVALUATION

After we have introduced UESM with its three methods, we present an evaluation. First, we describe the used corpus and evaluation metrics. Finally, we present the results of the evaluation and demonstrate the performance of UESM in comparison to LDA.

##### A. Dataset

In this evaluation we use the Bürgerliches Gesetzbuch (BGB)<sup>2</sup>, the civil code of Germany, in German language as corpus. However, OpenIE can not be used on this German language text and thus it is a example where we need UESM. The BGB is freely available and can be downloaded as XML file. Therefore, it is easily parsable and processable. As the corpus is a law text it consists of correct language, i.e., punctuation and spelling follow the orthographic rules. Thus, less preprocessing and no data cleaning is needed. Furthermore, the words used in text documents have a clear meaning and mostly the same words are used instead of using synonyms.

The entire corpus consists of 2462 law paragraphs and overall 8020 sentences which are used as SCD windows. Each law paragraph contains between 1 and 45 sentences with an average of 3.3 sentences. The vocabulary consist of 5294 words, where each sentence is between 1 and 51 with an average of 10.9 words long.

##### B. Metrics

Topic models are trained unsupervised using statistical methods, thus, the topics gained by LDA are statistically optimized but may not match human judgement of *good*

<sup>2</sup><https://www.gesetze-im-internet.de/bgb/>, Englisch translation [https://www.gesetze-im-internet.de/englisch\\_bgb/](https://www.gesetze-im-internet.de/englisch_bgb/)

topics. In general, automatically evaluating the quality of a model from a human point of view is a difficult task. A common measure to evaluate the interpretability of topics regarding human judgement is coherence. Röder et al. [12] compare and evaluate multiple coherence measures against human judgement as gold standard. The authors gain the best results using the  $C_V$  measure. However, due to negative correlations and problems reproducing the  $C_V$  values in their paper, Röder does not recommend to use the  $C_V$  coherence any more<sup>3</sup>. Therefore, in our evaluation we use the UMass coherence calculated using Gensim’s coherence model.

Furthermore, the number of referenced SCD windows per SCD is relevant. For example, having 1000 SCD windows and 100 SCDs, each SCD should have a similar number around 10 referenced SCD windows. It would be bad, if 99 SCDs reference 1 window each and the 1 remaining SCD references the remaining 901 windows. Therefore, we evaluate the number of referenced windows per SCD. Besides showing all numbers of referenced windows, we also show the numbers only for SCDs with two or more referenced windows, i.e., we interpret SCDs with only one referenced window as an irrelevant SCD to omit.

For LDA an evaluation of referenced documents per topic is not necessary, as the training ensures a similar number of referenced topics per document.

### C. Workflow and Implementation

UESM is implemented using Python and runs inside a Docker container. The implementation uses the libraries Gensim<sup>4</sup>, NumPy<sup>5</sup>, and NLTK<sup>6</sup>. The evaluation of UESM follows this workflow:

- (i) Extract the law paragraphs from the BGB’s XML file and divide each paragraph into its sentences, which are then used as initial SCD windows.
- (ii) Lowercase all characters, tokenize the sentences into words, stem the words, and eliminate stop words from a wordlist containing 232 German words. These four tasks are called preprocessing tasks. Preprocessing a text of a document transforms the text in a more digestible form for machine learning algorithms and increases their performance [13].
- (iii) Form an initial SCD matrix where each row contains the word probability distribution for one sentence of the corpus.
- (iv) Apply UESM with one of the three methods greedy, K-Means, or DBSCAN to detect similar rows in the SCD matrix. Afterwards, merge the similar rows or the rows in the same cluster by summing the distributions’ values. We run each method with different hyperparameters influencing the number of SCDs estimated. To be able

to show the results of all methods in one figure, we represent the results by the number of SCDs estimated. We show this number of SCDs by the reduction of the number of windows in percent, i.e., if an initial SCD matrix of 8020 rows is reduced to 802 rows, the matrix is reduced by 90 %. For example, in this case  $K$  would have been 802 for the method K-Means.

- (v) Calculate the UMass coherence using Gensim for the newly estimated SCD matrix on the corpus. Hereby, for each SCD the word probability distribution is used to determine the 20 most probable words of the referenced SCD windows. For each SCD these 20 words are interpreted as the SCD’s topic.

For comparison, we train two topic models by LDA using Gensim and the hyperparameters  $\alpha = 0.01$  and  $\beta = 0.05$ . Small  $\alpha$  and  $\beta$  lead the model to assign each document a single topic with a high probability, this matches the idea of associating an SCD window with one SCD. We train models with different numbers of topics and represent the number of topics by the reduction of the number of documents given to the model in percent, analogously to the reduction described in (iv) previously.

*a) LDA Windows:* This topic model is trained on the 8020 sentences as documents. Therefore, the model’s document topic distributions allow to determine the topic of each sentence and thus the model’s topics are comparable to the SCDs referencing multiple sentences in the corpus. However, LDA is not designed to be trained with very short documents like single sentences.

*b) LDA Documents:* This topic model is trained on the 2462 law paragraphs as documents and applies LDA in its typical fashion with medium sized documents. However, using this model’s document topic distributions it is not possible to determine the topic of each sentence, as each of the model’s documents contain more than one sentence.

Again, we calculate the UMass coherence for each topic model directly using Gensim’s functionality.

### D. Results

In this section, we present the results gained using UESM and the previously described workflow.

In Figure 1, the coherences of the three methods using UESM and both topic models are shown. The UMass scores calculated by Gensim are negative, higher values are better. On the left side, the reduction of the number of windows is small, thus many SCDs are created. Going to the right, the number of SCDs decreases, e.g., the rightmost triangle of greedy similarity represents 834 SCDs gained from initially 8020 windows.

The lines of DBSCAN, greedy similarity, K-Means, and LDA Documents are all close together, while LDA Windows shows poor results far below all other lines. This demonstrates that LDA Windows is not capable of estimating SCDs in an unsupervised manner, because the windows used as documents are too small. LDA Documents however demonstrates the UMass score a good topic model reaches on the BGB and

<sup>3</sup>"The usage of the  $C_V$  coherence is not recommended anymore!", stated on <https://github.com/dice-group/Palmetto/wiki/How-Palmetto-can-be-used>, last accessed 24. September 2022

<sup>4</sup><https://radimrehurek.com/gensim/>

<sup>5</sup><https://numpy.org/>

<sup>6</sup><https://www.nltk.org/>

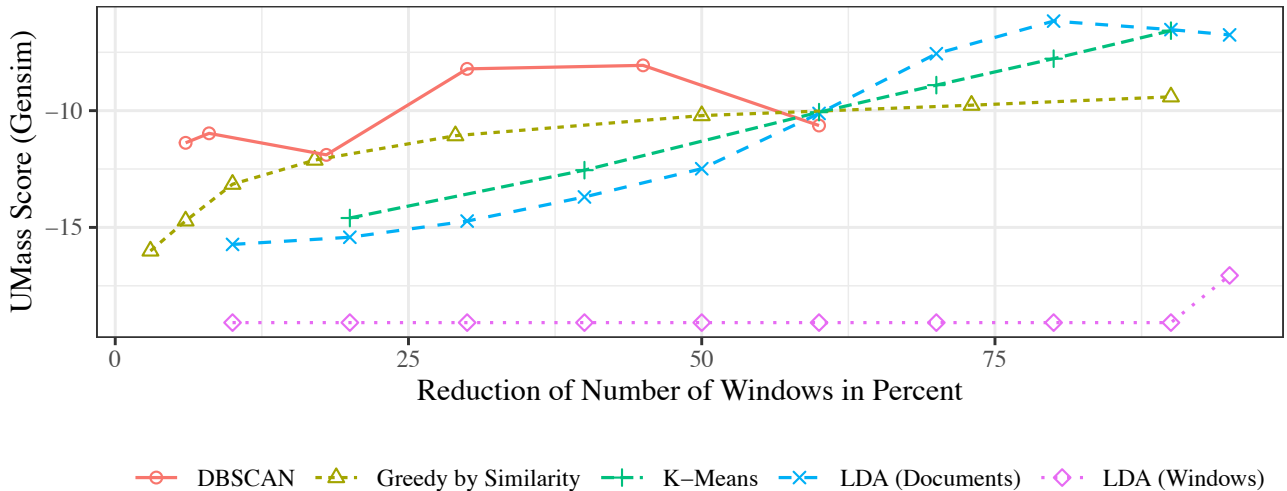


Figure 1. UMass coherence of the three methods using UESM and the coherences of both topic models trained using LDA for comparison.

UESM using K-Means reaches similarly good values. UESM works well with greedy similarity and less reduction of windows, but K-Means becomes better for more reduction. DBSCAN is quite unstable and the amount of reduction is difficult to configure using the hyperparameters  $\epsilon$  and number of minimum samples  $m_s$ . Although, the coherences of DBSCAN are good, we later see in Figure 2 why DBSCAN is not a good choice.

To summarize, using UESM with K-Means yields coherences on par with LDA. However, LDA is not able to estimate SCDs what UESM does.

In Figure 2 for each of the three methods two plots are shown. In the upper row, for each percentage of reduction the numbers of referenced windows are shown by boxplots on a logarithmic scale. The lower row shows the same, but SCDs referencing only one window are omitted. We focus on the lower row: For K-Means and greedy similarity most SCDs reference less than 10 windows, which is a good number of references. However, there are also many outliers referencing more windows. For K-Means the largest number of references is 952 and 1741 with greedy similarity. An SCD referencing 1741 windows references 21 % of the corpus and it is hard to imagine that 21 % of the corpus share the same concept. Again, this demonstrates that greedy similarity does not work well with a high reduction of the number of windows.

Using DBSCAN there are more SCDs referencing a large number of windows, which also implies that there are many SCDs referencing only one window. Also, the largest number of references is 3832 for DBSCAN, which means that a single SCD references 48 % of the corpus. An SCD referencing nearly half of the corpus can not be good. Furthermore, only this single SCD can reference 48 % of the sentences, while no other SCD can reference the same sentences.

Summarized, K-Means shows an overall very good distribution of referenced windows per SCD and greedy similarity is

good, too. Though, DBSCAN generates an SCD referencing nearly half of the corpus. Thus, DBSCAN should not be used.

## V. RELATED WORK

Before we conclude, we take a look at related work. Adding data to corpora of text documents has been investigated for a long time [14]. Often the data associated with a corpus is denoted as an annotation.

In the beginning of natural language annotation, most annotations had to be added manually to the corpora. Even today, crowdsourcing can be used to manually annotate text documents [15]. Furthermore, semi-automatic and automatic annotation systems were developed, too, e.g., OpenIE [7]. Thus, unsupervised corpus annotation remains an important field of research.

In the context of SCDs, we interpret an SCD as a corpus annotation and in context of this paper, an SCD annotates multiple sentences of similar concepts all over the corpus' text documents. Topic models assign a distribution over the topics, estimated by the model itself, to each text document in the corpus and each topic is characterized by a distribution of occurring words in the topic's documents. Thus, similarly to the topics of a topic model, UESM associates text documents with SCDs representing concepts. A well-known topic modelling technique is LDA [8]. LDA is a generative model representing documents as a probability distribution over topics. Many extensions have been proposed to optimize the performance of LDA, e.g., the author-topic model [16], which extends LDA to couple each author of a document with a multinomial over words, and the dynamic topic model [17], which allows for analysing topic changes over time.

Documents assigned with a similar distribution over the topics, are assumed to be similar in terms of an topic model. However, LDA's perception of similar documents may not always match the human perception of similar documents [18].

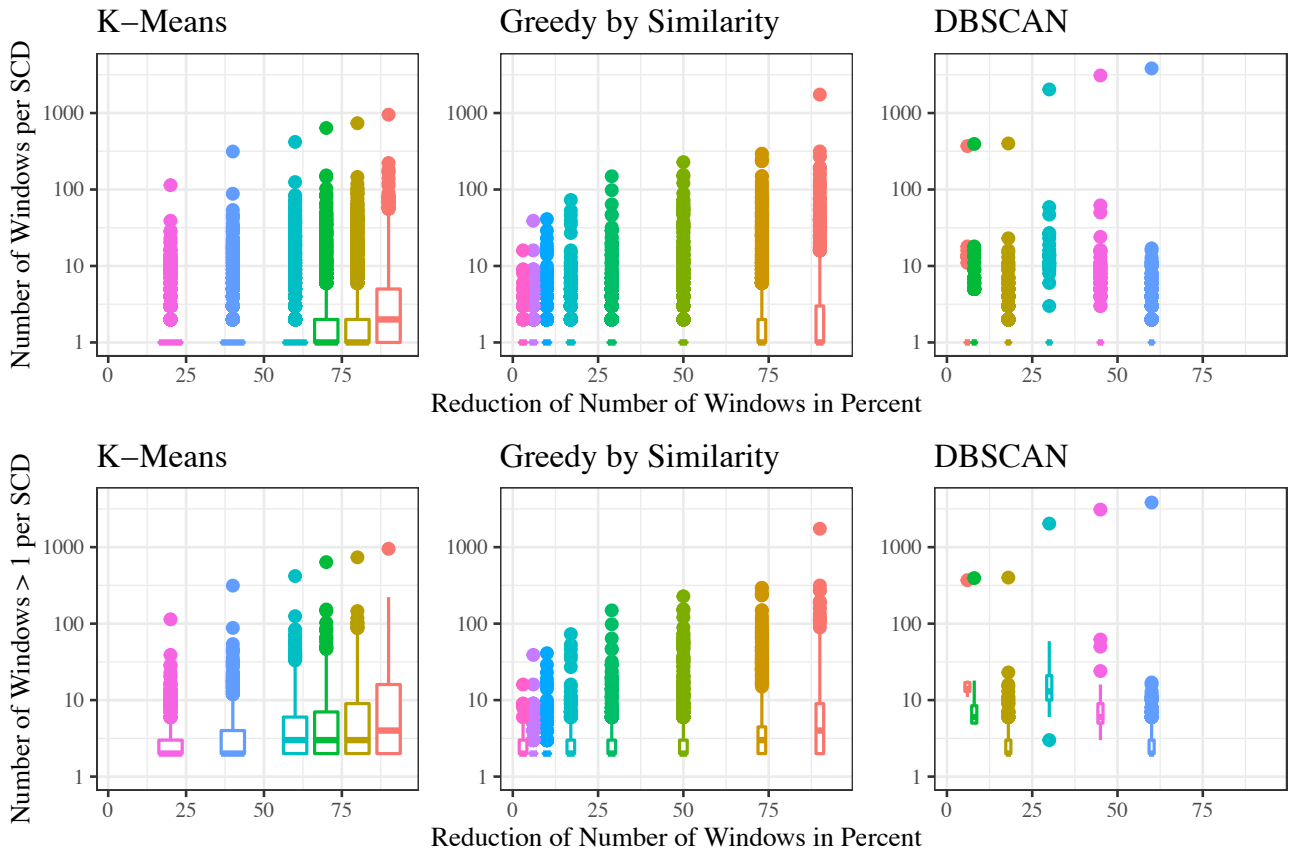


Figure 2. Number of windows referenced by one SCD for the three different methods of UESM. In the lower row, SCDs referencing only one window are omitted.

UESM uses greedy similarity, K-Means or DBSCAN to identify similar sentences. Another technique to find similar sentences in a corpus of text documents is Similar Short Passages Identifier (SiSP) [19]. SiSP first extracts features from the sentences and then creates clusters of similar sentences. The authors evaluate the clusters found by SiSP against human labeled sentences. UESM may be used with SiSP, however, SiSP was developed for the Portuguese language.

A further idea is to cluster sentences hierarchically [20]. In difference to the clustering techniques used by UESM, the authors start with a sentence in the corpus and build a hierarchical clustering from this sentence. The hierarchical clustering has a tree-like structure, i.e., after starting from the first sentence, the tree branches across multiple levels to different concepts in the corpus.

Clustering can not only be used to identify similar sentences, it can also help to annotate sentences with their sentiment [21]. The authors assume, that two sentences in the same cluster have a similar sentiment and thus they can enrich the number of labels in a sparsely labeled corpus. In cases, where short sentences do not contain enough shared words to apply the cosine similarity, a ranking of the suitable clusters for each sentence can be used [22] to increase the performance of clustering techniques.

Text summarization is another field of research, where clusters of similar sentences are used. Thereby, the idea is to remove most of the similar sentences and only keep one sentence from each cluster. SimFinder [23] clusters small pieces of text, like sentences, into tight clusters. Unlike UESM, SimFinder does not work unsupervised, as it needs feature words in beforehand.

Another approach for text summarization is to extract the word vectors from each sentence and weight each word using tf-idf [24]. Then, the weighted word vectors are clustered using the cosine similarity with K-Means [9], again, from each cluster one sentence makes up the summary [25]. This approach overlaps with UESM in using word vectors, the cosine similarity, and K-Means. However, the authors pursue only the goal of text summarization, while UESM uses three methods and represents the concepts and topics of a corpus in the estimated SCD matrix.

## VI. CONCLUSION

This paper introduces UESM with three methods, namely K-Means, greedy similarity, and DBSCAN. UESM estimates SCD matrices for corpora of text documents in an unsupervised manner. Thereby, UESM detects sentences of similar concepts or topics in a corpus and then associate the same SCD to these similar sentences.

An SCD matrix for a corpus can be interpreted as a topic model of the corpus. Hence, the well-known LDA [8] is used to evaluate the performance of UESM. We use the UMass coherence to evaluate the quality of each model and show, that especially UESM using K-Means performs as good as LDA. Whereby in particular, LDA does not estimate an SCD matrix, but especially the SCD matrix is needed by approaches introduced by Kuhr et al. [2], [3] and Bender et al. [4], [6]. UESM enables the authors' approaches to be used in an unsupervised manner, i.e., without needing SCDs for a corpus in beforehand. Generally, without a focus on SCDs, UESM provides a new and powerful technique to create a topic model for a corpus.

Overall, the evaluation shows that UESM using K-Means can keep up with LDA, while the greedy method is only slightly less powerful. Because DBSCAN associates too many sentences with the same SCD, DBSCAN is not suitable for most use-cases.

So far, we have introduced an SCD as a tuple of the SCD's additional data  $\mathcal{C}$  and the referenced sentences. In this paper, we put efforts in finding the referenced sentences but we did not estimate the data  $\mathcal{C}$ . Thus, currently we only get the word probability distribution and the referenced sentences for each SCD without  $\mathcal{C}$ . Future work will focus on estimating  $\mathcal{C}$ , which is similar to automated topic naming for topic models [26]. Additionally, we will focus on more approaches applying the SCD matrix to solve natural language processing related problems.

#### ACKNOWLEDGMENT

The research of Marcel Gehrke was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy – EXC 2176 “Understanding Written Artefacts: Material, Interaction and Transmission in Manuscript Cultures”, project no. 390893796. The research was conducted within the scope of the Centre for the Study of Manuscript Cultures (CSMC) at Universität Hamburg. The authors thank the AI Lab Lübeck for providing the hardware used in the evaluation.

#### REFERENCES

- [1] F. Kuhr, T. Braun, M. Bender, and R. Möller, “To Extend or not to Extend? Context-specific Corpus Enrichment,” *Proceedings of AI 2019: Advances in Artificial Intelligence*, pp. 357–368, 2019.
- [2] F. Kuhr, M. Bender, T. Braun, and R. Möller, “Augmenting and automating corpus enrichment,” *Int. J. Semantic Computing*, vol. 14, no. 2, pp. 173–197, 2020.
- [3] —, “Context-specific adaptation of subjective content descriptions,” *15th IEEE International Conference on Semantic Computing, (ICSC 2021), Laguna Hills, CA, USA, January 27-29*, pp. 134–139, 2021.
- [4] M. Bender, T. Braun, M. Gehrke, F. Kuhr, R. Möller, and S. Schiff, “Identifying and translating subjective content descriptions among texts,” *Int. J. Semantic Computing*, vol. 15, no. 4, pp. 461–485, 2021.
- [5] F. Kuhr, B. Witten, and R. Möller, “Corpus-Driven Annotation Enrichment,” *13th IEEE International Conference on Semantic Computing, (ICSC 2019), Newport Beach, CA, USA, January 30 - February 1*, pp. 138–141, 2019.
- [6] M. Bender, F. Kuhr, and T. Braun, “To extend or not to extend? complementary documents,” *16th IEEE International Conference on Semantic Computing, (ICSC 2022), Virtual, January 26-28*, pp. 17–24, 2022. [Online]. Available: <https://doi.org/10.1109/ICSC52841.2022.00011>
- [7] G. Angeli, M. J. Johnson Premkumar, and C. D. Manning, “Leveraging linguistic structure for open domain information extraction,” *Proceedings of the Association of Computational Linguistics (ACL)*, pp. 344–354, 2015.
- [8] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent dirichlet allocation,” *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003. [Online]. Available: <http://jmlr.org/papers/v3/blei03a.html>
- [9] S. Lloyd, “Least squares quantization in pcm,” *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, 1982.
- [10] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, “A density-based algorithm for discovering clusters in large spatial databases with noise,” pp. 226–231, 1996.
- [11] D. Pelleg, A. W. Moore *et al.*, “X-means: Extending k-means with efficient estimation of the number of clusters,” in *Icml*, vol. 1, 2000, pp. 727–734.
- [12] M. Röder, A. Both, and A. Hinneburg, “Exploring the space of topic coherence measures,” in *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, ser. WSDM '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 399–408.
- [13] S. Vijayarani, J. Ilamathi, and S. Nithya, “Preprocessing techniques for text mining - an overview,” *International Journal of Computer Science & Communication Networks*, vol. 5, p. 7–16, 2015.
- [14] G. V. Maverick, “Computational analysis of present-day american english. henry kučera, w. nelson francis,” *International Journal of American Linguistics*, vol. 35, no. 1, pp. 71–75, 1969. [Online]. Available: <https://doi.org/10.1086/465045>
- [15] M. Sabou, K. Bontcheva, L. Derczynski, and A. Scharl, “Corpus annotation through crowdsourcing: Towards best practice guidelines,” *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC'14)*, 01 2014.
- [16] M. Rosen-Zvi, T. Griffiths, M. Steyvers, and P. Smyth, “The author-topic model for authors and documents,” *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, p. 487–494, 2004.
- [17] D. M. Blei and J. D. Lafferty, “Dynamic topic models,” ser. ICML '06. New York, NY, USA: Association for Computing Machinery, 2006, p. 113–120.
- [18] W. B. Towne, C. P. Rosé, and J. D. Herbsleb, “Measuring similarity similarly: LDA and human perception,” *ACM Trans. Intell. Syst. Technol.*, vol. 8, no. 1, pp. 7:1–7:28, 2016. [Online]. Available: <https://doi.org/10.1145/2890510>
- [19] E. R. M. Seno and M. d. G. V. Nunes, “Some experiments on clustering similar sentences of texts in portuguese,” in *Computational Processing of the Portuguese Language*, A. Teixeira, V. L. S. de Lima, L. C. de Oliveira, and P. Quaresma, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 133–142.
- [20] D. Kavayasarjana and B. C. Rao, “Hierarchical clustering for sentence extraction using cosine similarity measure,” in *Emerging ICT for Bridging the Future - Proceedings of the 49th Annual Convention of the Computer Society of India (CSI) Volume 1*, S. C. Satapathy, A. Govardhan, K. S. Raju, and J. K. Mandal, Eds. Cham: Springer International Publishing, 2015, pp. 185–191.
- [21] M. Hadano, K. Shimada, and T. Endo, “Aspect identification of sentiment sentences using a clustering algorithm,” *Procedia - Social and Behavioral Sciences*, vol. 27, pp. 22–31, 2011, computational Linguistics and Related Fields. [Online]. Available: <https://doi.org/10.1016/j.sbspro.2011.10.579>
- [22] L. Yang, X. Cai, Y. Zhang, and P. Shi, “Enhancing sentence-level clustering with ranking-based clustering framework for theme-based summarization,” *Information Sciences*, vol. 260, pp. 37–50, 2014. [Online]. Available: <https://doi.org/10.1016/j.ins.2013.11.026>
- [23] V. Hatzivassiloglou, J. L. Klavans, M. L. Holcombe, R. Barzilay, M.-Y. Kan, and K. McKeown, “SimFinder: A flexible clustering tool for summarization,” 2001.
- [24] J. E. Ramos, “Using tf-idf to determine word relevance in document queries,” 2003.
- [25] K. Shetty and J. S. Kallimani, “Automatic extractive text summarization using k-means clustering,” in *2017 International Conference on Electrical, Electronics, Communication, Computer, and Optimization Techniques (ICEECCOT)*, 2017. [Online]. Available: <https://doi.org/10.1109/ICEECCOT.2017.8284627>
- [26] A. Hindle, N. A. Ernst, M. W. Godfrey, and J. Mylopoulos, “Automated topic naming,” *Empirical Software Engineering*, vol. 18, no. 6, pp. 1125–1155, 2013.