# Clustering Introductory Computer Science Exercises Using Topic Modeling Methods

Laura O. Moraes and Carlos Eduardo Pedreira, *Senior Member, IEEE*

arXiv:2104.10748v1 [cs.LG] 21 Apr 2021

*Abstract*—This is the accepted version, to read the final version published in 2021 in the IEEE Transactions on Learning Technologies (IEEE TLT), please go to: http://doi.org/10.1109/TLT.2021.3056907. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Manually determining concepts present in a group of questions is a challenging and time-consuming process. However, the process is an essential step while modeling a virtual learning environment since a mapping between concepts and questions using mastery level assessment and recommendation engines are required. We investigated unsupervised semantic models (known as topic modeling techniques) to assist computer science teachers in this task and propose a method to transform Computer Science 1 teacher-provided code solutions into representative text documents, including the code structure information. By applying non-negative matrix factorization and latent Dirichlet allocation techniques, we extract the underlying relationship between questions and validate the results using an external dataset. We consider the interpretability of the learned concepts using 14 university professors' data, and the results confirm six semantically coherent clusters using the current dataset. Moreover, the six topics comprise the main concepts present in the test dataset, achieving 0.75 in the normalized pointwise mutual information metric. The metric correlates with human ratings, making the proposed method useful and providing semantics for large amounts of unannotated code.

*Index Terms*—Topic modeling, clustering, educational data mining, computer science education.

## I. INTRODUCTION

MEASURING student's knowledge requires an understanding of which educational concepts are needed to answer each question. Recently, open online courses and intelligent tutoring systems are widely adopted learning environments. Their popularity increases the demand for tools to map questions to concepts correctly since students' mastery level assessment and next steps recommendation depend on these mappings.

The authors are with the Systems and Computing Engineering Department (COPPE-PESC), Universidade Federal do Rio de Janeiro (UFRJ), Rio de Janeiro, 21941-914, Brazil (e-mail: lmoraes@cos.ufrj.br; pedreira56@gmail.com).

However, manually identifying the concepts required to answer questions can be time-consuming and difficult, increasing the need for tools to assist teachers in the tasks. Desmarais [1] suggested that even partial automation of the process can be highly desirable. Besides decreasing the manual labeling required from the experts, the process automation also results in a more objective and replicable mapping. Applying supervised machine learning-based solutions is not entirely appropriate because it requires a considerable amount of labeled data. Also, a question can relate to multiple concepts, increasing the complexity of the labeling task. The traditional unsupervised learning methods, such as K-means and hierarchical clustering, are also not suitable for this task because it is hard to determine each cluster's features.

This paper proposes unsupervised semantic methods, known as topic modeling techniques [2]–[5], as more interpretable methods for experts, to be applied in introductory computer science problems. Specifically, we propose modeling code snippets as text documents and use topic modeling techniques to extract and improve the semantic relationships between them, providing technology to support concept identification experts.

Our key research questions are summarized as follows: 1) how can semantic relationships be extracted and structured from code? 2) how can humans read, interpret, and use the extracted relationships?

The main contributions of this paper include:

(i) A tokenization structure to transform raw code snippets into a document-term matrix.
(ii) A code-clustering method to optimize positively correlated metrics for human-interpretability.
(iii) Experts validation, illustrating how the proposed method can support questions' exercise labeling using each topic's terms.

The proposed code-clustering pipeline builds a document-term matrix with a code tokenizer by comparing various methods, including clustering algorithms. We compare non-negative matrix factorization (NMF) [2], latent Dirichlet allocation (LDA) [3], and K-means [6], as a baseline. The models were evaluated with the UMass and UCI coherence metrics [7]–[10] using the top-5 and top-10 terms. According to the metrics scores, we selected the two best-ranked models using Fagin's algorithm [11]. The LDA-based clustering approach provides the most interpretable results from these models. Fourteen professors manually contextualize the LDA-based clustering in the Computer Science 1 (CS1) domain, demonstrating how the proposed method could be used to facilitate the clusters' interpretability.

The next section begins by reviewing the manual, supervised, and unsupervised concept identification approaches. Code-clustering techniques inside the educational data mining (EDM) and software engineering contexts are reviewed together with existing LDA proposals to handle short texts. We describe a proposed method to cluster code in Section III. The main challenges facing the CS1 context are the answers' small size and document term-matrix sparsity. Our proposed code tokenizer overcame these problems using code structure information to augment the corpus. The results are shown in Section IV, where the best two clustering schemes are analyzed based on the coherence evaluation metrics. This section also demonstrates how professors can get an overview of the required concepts from these results. Finally, Section V presents the conclusion and future work directions.

## II. RELATED WORK

The existing methods to identify concepts from a set of CS1 exercises involve manual work and input from experts [12], [13]. For example, Sheard *et al.* [12] characterized introductory programming examination questions according to their concept areas, question style, and required skills. Participants manually classified the questions and the determined topics covered alongside the necessary skill levels to solve them. Nonetheless, applying a successful approach in a different set of exercises requires a new manual labeling stage, which may not be achievable.

One strategy to overcome this issue and minimize the domain experts' workload is to apply supervised learning. Previous research in question classifications used supervised learning to classify questions according to the level of difficulty [14], Bloom's taxonomy [15], answer type [16], and subject [17]. In Godea *et al.* [16], the features are derived from the questions, using part-of-speech tags, word embeddings, interclass correlations, and manual annotation. Supraja *et al.* [15] use a grid search to analyze different combinations of weight schemes and methods to find the best set of parameters to build a supervised model to classify questions given Bloom's Taxonomy. Its main cost is the manual annotation of all labels, impractical when applying to large datasets. Unsupervised learning can group similar items without a predefined label, but it is harder to ascertain the results since there is no objective goal to analyze, and evaluating the clustering outcomes becomes a subjective task. Unsupervised learning techniques have been used to address EDM problems [18]–[22]. For example, Trivedi *et al.* [23] use spectral clustering with linear regression to predict student performance. In the questions' classification context, an unsupervised approach using K-means, as a clustering algorithm [24], was proposed to group similar learning objects (such as handouts, exercises, complementary readings, and suggested activities). Still, K-means does not provide a list of features that best characterize each cluster, making the expert infer them manually by reading a sample of each cluster's exercises.

It is possible to use the code provided as answers to the exercises as input while restricting the concept identification to the CS1 domain. Unsupervised learning using code as input is achieved by calculating their abstract syntax trees (ASTs) and clustering the most similar ones; besides, various strategies calculate similarity among trees. Huang *et al.* [25] and Nguyen *et al.* [26] use edit tree distance, Paasen *et al.* [27] apply sequential clustering algorithms, and Price *et al.* [28] and Mokbel *et al.* [29] fragment each AST into subgraphs, pair similar subgraphs and compute the distance between the subgraphs within a pair. Concepts identification in these scenarios requires experts to read samples of exercises from each cluster to infer the relevant features. It is also unclear how to generalize it to other domains where the answers are not constructed directly using a tree or a graph.

We propose a CS1 concept identification method to support experts in labeling concepts by providing the relevant features in each cluster. This method could potentially be extended to different domains since it only requires text inputs. We develop a code tokenizer that uses its code structure to augment and improve the corpus. We apply topic modeling techniques where the documents are the code provided as answers to the exercises. In the software engineering field, concepts are extracted from code using programmers' annotations [30]. Similarly, code clustering, identifying similar features/concepts in a repository, predicting bugs, analyzing/predicting source code evolution, tracing links between modules, and detecting clones are widely applied in the same field [31]. For the EDM, Azcona *et al.* [32] created a Python code submission tokenizer to setup features to generate code embeddings. Unlike the software engineering context, code snippets in CS1 are small in size and lack annotations.

Although the LDA in Blei *et al.* [3] is a common technique in topic modeling, it does not perform well in short texts (code in this context) because the traditional way of extracting terms does not provide enough textual words to characterize a specific topic [33], [34]. It is necessary to decrease the latent document-topic or word-topic spaces, making them more specific for each context. Hsiao *et al.* [35], [36] propose a topic-facet LDA model using sentence LDA (SLDA) with a facet representing a more specific topic and all words from a sentence belonging to the same facet. Zhao *et al.* [37] decrease the latent space by creating a common word distribution with denominated background words, which are the same for every topic. Steyers *et al.* [38] and Rosen-Zvi *et al.* [39] adopted a similar strategy. In their method, the generative process to create a document decreases the space by choosing an author and then choosing a topic. Li *et al.* [40] use a distribution over tags to restrict the latent topic space before inferring the documents' topic distribution.

Another approach to overcome the lack of textual words is to increase the representation, so Weng *et al.* [41] proposed an LDA to cluster subjects' twitterers (people who tweet). Their method increases document representation by aggregating all tweets from a single user into one document. Abolhassani and Ramaswamy [34] enhance short semi-structured texts by augmenting the corpus' structure, which is similar to the method we adopted in this paper. In our proposed tokenizer, we increase the vocabulary size (total terms) from 287 to 2388 and the average of terms per document from 23 to 137. We maintain a 95% sparsity, which agrees with the sparsity of

the long-text documents from Syed and Spruit [42] and Zhao *et al.* [37], i.e., successfully clustered using topic modeling techniques.

## III. METHODS

The methodology is categorized into three main tasks: 1) we generate a database by crawling different Python web tutorials, 2) run code-clustering experiments to group exercises into topics, and 3) ask experts to contextualize the clusters into CS1 concepts. Task 1 prepares the data to investigate the research questions. Task 2 explores ways of extracting semantic relationships from code [research question 1 (RQ1)] by proposing a code tokenizer and comparing various data transformation methods and topic modeling algorithms. In task 3, we ask professors to read and interpret task 2, giving support for RQ2. All the scripts used in this study analyses were achieved using Python and open-source Python libraries. The code is available at https://github.com/laura-moraes/machine-teaching.

### A. Dataset

The objective of the experiment is to find semantically related CS1 code solutions written in Python. We chose four introductory online tutorials: Practice Python [43], Python School [44], Python Programming Exercises [45], and W3Resource [46] that provide both solutions and exercise statements. Since the sources do not have label topics or follow a course curriculum with structured syllabus topics, we work in an unsupervised environment. We crawled 54 exercises for the training set. The code snippets are functions with an average of 9 lines/code.

For the test dataset, we collected solutions from another set of exercises given to us by the CS1 professors at the Federal University of Rio de Janeiro (UFRJ). There are 65 different problems with their respective solutions in the dataset. As the training set, the code snippets are functions with an average of 7 lines/code.

### B. Code-Clustering Pipeline

The code-clustering pipeline takes as input Python code snippets, which are semi-structured text documents. By using topic modeling techniques, the pipeline outputs an underlying structure within the semi-structured corpus. It contains the topics present in the code snippets and the most relevant words that characterize them. This paper is based on the assumption that code snippets with similar CS1 concepts share identical terms. Therefore, based on this assumption, the extracted topic underlying structure can be interpreted as CS1 concepts or groups of CS1 concepts present in the code snippets.

The code-clustering pipeline starts by transforming the original data to the proper format expected by the topic modeling methods. We augmented the data and constructed a matrix D (the document-term matrix) where each element $D_{ij}$ contains the weight of term $w_j$ in document $d_i$. Then, using topic modeling, we calculated the relevance of each topic $t_k$ for each document $d_i$ and the relevance of each term $w_j$

for each topic $t_k$. Finally, we applied a grid search and topic coherence to choose the best models and evaluate the external corpus results. In the topic filter and selection phase, we also processed the resulting topics by merging similar or removing topics with few documents. These results are presented in Section IV, while Fig. 1 illustrates the code-clustering pipeline. External evaluation is not depicted in this overview.

*1) Data Transformation:* In this application, the CS1 code solutions written in Python are considered documents. The document-term matrix creation process starts by splitting each code snippet into words. The first proposed tokenizer includes only split word tokens. Henceforth this tokenizer will be referred to as the *standard tokenizer*. As stated in the related work section: 1) the LDA usually does not perform well on short texts and 2) augmenting the corpus by adding the text's structure on semi-supervised documents demonstrated improved results. We propose a new tokenizer to augment the standard tokenizer with extra features and refer to it as the *augmented tokenizer*. The augmented tokenizer parses the code and makes special annotations by adding extra features if the token is a number, an array (or a list), a dictionary, a string, a logical (or arithmetic) operator, a class method, or indentation. The word itself is added to the document-term matrix if the token is a reserved word. Besides adding single tokens, this tokenizer also considers bigrams and trigrams. Although the document-term matrix does not consider the terms' order, this can be enforced by adding n-grams as a matrix feature. For example, the code snippet in Fig. 2(a) is first transformed to its augmented version (see Fig. 2(b)). Then, every single word, including the bigrams and the trigrams, are added as tokens to the document-term matrix. Table I presents some examples of the document-term matrix terms, and, in total, the document is tokenized into 75 terms.

TABLE I
EXAMPLE OF DOCUMENT-TERM MATRIX. A SET OF TERMS FROM THE COMPLETE DOCUMENT-TERM MATRIX AFTER AUGMENTING AND TOKENIZING THE DOCUMENT FROM FIG. 2

| Terms | Count |
|---|---|
| is_block | 3 |
| is_indent | 3 |
| is_number | 2 |
| is_op_logic | 3 |
| if | 1 |
| is_op_logic is_number | 2 |
| is_block is_indent | 3 |
| is_op_logic is_number is_op_logic | 1 |

After the document-term matrix creation, we applied some transformations to enhance document representation and decrease matrix sparsity. First, we removed tokens with document frequency below a fixed threshold to perform feature selection. This threshold was determined using a hyperparameter grid search ranging from 5% to 50% with a 5% step. Second, we decided how to count a token frequency: either the token is counted once per document (binary appearance) or every time it appears. Finally, some tokens may be more important than others. For example, term frequency by inverse document frequency (TF-IDF) [47], [48] recalculates the tokens' weights by balancing two factors: 1) a term that occurs in many
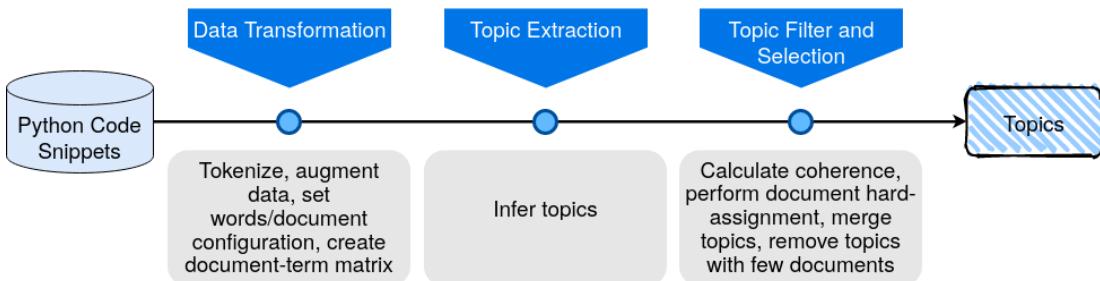
Fig. 1.  Code-clustering pipeline overview [34].

```
def light(switchA, switchB):
    if switchA == 1 and switchB == 1:
        return True
    else:
        return False
```

(a) Original code snippet

```
def light switchA switchB is_block
  is_indent if switchA is_op_logic is_number is_op_logic
            switchB is_op_logic is_number is_block
    is_indent return True
  is_dedent else is_block
    is_indent return False
  is_dedent
is_dedent
```

(b) Augmented code snippet

Fig. 2.  A code snippet and its augmented version. The augmented version will be processed in the augmented tokenizer, whereas the regular one will be processed in the standard tokenizer.

documents should not be as important as a more exclusive term, since it does not characterize documents well. Also, 2) a term that appears in a small number of documents may only be particular to those documents and not enough to distinguish a topic. Yan *et al.* [49] propose another way of recalculating the terms' weights. Their method (called NCut) comes from the normalized cut problem on term affinity graphs. This weighting scheme modifies terms' counts based on terms cooccurrence and not on document frequency. Their experiments show NMF's performance increase on short text clustering using the NCut weighting scheme.

*2) Topic Extraction:* As stated earlier, after document processing, a document-term matrix $D$ is generated. The matrix rows represent points in an $\mathbb{R}^n$ feature space, where $n$ is the total number of terms, and each term $w_j$ corresponds to a dimension. It becomes a classical clustering problem where we expect similar documents to be in surrounding regions in space. So, clustering algorithms like K-means, hierarchical clustering, and nearest neighbors are applicable here. However, for topic modeling tasks, algorithms like the NMF [2], [50] and the LDA [3] are effective since they interpret terms' counts as a set of visible variables generated from a set of hidden variables (topics) [51], [52]. Accordingly, the documents can be modeled as a distribution of topics and topics as a distribution of terms. We used two topic modeling techniques:

1) NMF: a matrix factorization technique with a particular property of only allowing non-negative values in its entries, which is well-suited for human-interpretability [2].
2) LDA: a generative probabilistic model that describes how to create documents in a collection. Once you have a dataset, a group of already written documents, we find the distributions that create these documents. The LDA algorithm tries to backtrack this probabilistic model to find a set of topics that are likely to have generated the dataset [4]. To generate a document, we sample from two distributions using the following iterative process:
   a) We sample a topic for the given document (a document is a distribution of topics).
   b) We sample a term from the topic sampled in step a) (a topic is a distribution of terms).

*3) Topic Filter and External Evaluation:* Given several document-term matrix creation options and two different topic modeling methods, we need to find the best set of hyperparameters. There are strategies in the literature to find a near-optimal set of models' hyperparameters, such as manual search, grid search, and random search [53]. Although random search demonstrates promising results in general machine learning tasks [53], Chuang *et al.* [54] and Wang and Blei [55] results were competitive using grid search in topic modeling tasks. We chose to use a grid search approach. There was a prior manual stage to define the regions in which grid search would act. Since the dataset is not large and the number of hyperparameters to try is not extensive, it is efficient to run an exhaustive search combining hyperparameters. In total, there are 1680 possible combinations: 10 minimum document frequencies (ranging from 5% to 50% with 5% step increment), 2 binary appearance options, 3 token weights (counts) transformation possibilities (none, TF-IDF, and NCut), 2 clustering methods (LDA and NMF), and 14 number of clusters (i.e., $10 \times 2 \times 3 \times 2 \times 14 = 1680$). The grid search was set to search between 2 and 15 clusters (the upper bound is based on the number of concepts from Table II in Section III-C).

To determine whether topics are well-defined, we can use topic coherence and pointwise mutual information (PMI) metrics, which correlated well with human-interpretability [8], [9], [56]. As explained in Section III-B2 (topic extraction), when using NMF or LDA, each topic is mapped to a list of top-N words that best define the topic. Topic coherence calculates the ratio between the cooccurrence of these top-N words and

their total occurrence. The assumption is that the words that best characterize a topic often appear together if a topic is well-defined. This paper applied two types of topic coherence metrics: UCI [7] and UMass [8]. The UCI metric based on PMI is calculated using an external validation source. The PMI can be substituted using normalized PMI (NPMI) to better correlate with humans' ratings [9]. The UMass metric uses the conditional probability of one word occurring given that one other high-ranked word occurred and can be measured using the modeled corpus, without depending on an external reference corpus. We used the UMass coherence to choose the best models since it is an internal validation metric (it only evaluates the clustered data). To assess the models, we used an external dataset with the UCI NPMI metric.

Defining $P(w_i)$ as the probability of the term $w_i$ occurring and $P(w_i, w_j)$ as the probability of terms $w_i$ and $w_j$ cooccurring, we calculated the coherence for a single topic $t_k$ using (1), (2), and (3). In this paper, the topic coherence for a single topic was calculated using top-5 and top-10 terms. After calculating each topic's coherence in a single hyperparameter combination, this combination's coherence was reported as the median of all topic coherence.

$$C_{UMass}(W) = \sum_{i=1}^{N} \sum_{j=1}^{N} log \frac{P(w_i, w_j) + \epsilon}{P(w_i)} \quad (1)$$

$$C_{UCI}(W) = \sum_{i=1}^{N} \sum_{j=1}^{N} NPMI(w_i, w_j) \quad (2)$$

$$NPMI(W) = \frac{log \frac{P(w_i, w_j) + \epsilon}{P(w_i) P(w_j)}}{-log(P(w_i, w_j) + \epsilon)} \quad (3)$$

where $W = (w_1, w_2, ..., w_N)$ are the top-N terms for calculating the coherence. An $\epsilon$ value of 0.01 was used to avoid taking a zero logarithm.

We performed hard-assignment to cluster documents by topic by assigning each document to the topic with the most relevance (weight) in the document-topic matrix. The hard-assignment was achieved with minimal loss of information when a topic strongly characterized a document. In addition to assigning documents to topic clusters, the set of features/terms that best characterize each cluster/topic were extracted for further analysis.

As a final step, we also merged semantically similar topics and removed those not containing relevant information.

### C. Topics Contextualization

To relate concepts and topics, we first defined the most commonly seen concepts in CS1 exercises. The following four references were used to create a list of concepts commonly used in CS1 courses:

1) Computer Science Curricula 2013 [57]: a document jointly built by the Association of Computer Machinery and the IEEE Computer Society. The document recommends curricular guidelines for computer science education, which we used as the main concept list. We used the papers in items 2, 3, and 4 to improve it.

2) Exploring programming assessment instruments: A classification scheme for examination questions [12]: creates a classification scheme characterizing exam questions by their concept areas, question style, and skills a student needs to solve them. We used the list of the proposed concepts as a second source to enhance the main list.

3) Reviewing CS1 exam question content [13]: this paper used nine experienced CS1 instructors to review the concepts required in examinations from different North American institutions. It created a list of concepts using the intersection of three other experiments [58]–[60].

4) Identifying challenging CS1 concepts in a large problem dataset [61]: this paper identifies the most challenging CS1 concepts for students based on 266,852 web-based code-writing student responses. Their concept list is based on the course structure from the CS1 class where they experimented.

Table II shows the final list of consolidated concepts.

TABLE II
LIST OF CS1 CONCEPTS

| 1. Syntax | 6. Logic | 11. Conditional |
|---|---|---|
| 2. Assignment | 7. Data type: string | 12. Loop |
| 3. Data type: number | 8. Data type: array | 13. Nested loop |
| 4. Data type: boolean | 9. Data type: tuple | 14. Function |
| 5. Math | 10. Data type: dict | 15. Recursion |

Then, to interpret the meaning of the topics, we asked 14 professors to perform three tasks. The professors (with 2 to 20 years of teaching experience) teach CS1 or other programming-related subjects.

- **Theme identification:** we present some code snippets belonging to the topic and found essential tokens for each topic. The professors were asked to label each topic with free-text descriptions. We tokenized the descriptions and counted the terms. We also created the topic titles based on the terms that appeared more frequently in the descriptions.

- **Concept identification:** each professor was asked to associate up to three concepts (from the 15 available in Table II) to 15 randomly assigned code snippets. Then, we aggregated the concepts related to each code to the associated topics. In this way, it is possible to understand the most relevant concepts in each topic and calculate how well the concepts inside a cluster agree to them. A similar approach was used by Lan *et al*. [62] to contextualize concept clusters.

- **Intruder identification:** Chang *et al*. [63] proposed quantitative methods to analyze the interpretability of the latent space created using topic modeling techniques. Their paper proposed a method to identify an intruder topic given a document. We adapted this method because we hard-assigned each document to a single topic. For our analysis, the professors were asked to identify the intruder document given a topic. This approach has been used by Mahmoud and Bradshaw [64] to assess the quality of their topic modeling approach on Java-based systems.

## IV. RESULTS AND DISCUSSION

We run each hyperparameter combination from the 1680 possibilities 10 times and calculated their average coherence and standard deviation. Next, the two best-ranked results are analyzed. They were calculated using Fagin's algorithm [11] for top-5 and top-10 terms UMass coherence. Table III shows the set of hyperparameters for each experiment.

TABLE III
SET OF HYPERPARAMETERS FOR EACH EXPERIMENT

|        | Min DF | Binary | Vectorizer | Method | # of clusters |
|--------|--------|--------|------------|--------|---------------|
| Exp. 1 | 0.35   | True   | NCut       | NMF    | 7             |
| Exp. 2 | 0.05   | True   | Count      | LDA    | 12            |

### A. Experiment 1

After the document-term matrix factorization, we hard-assigned each document to the topic with the highest relevance (highest weight in the document-topic distribution). Table IV shows the number of documents assigned per topic. After assigning each document to its related topic in this experiment, the documents are only assigned to four of the seven topics. Fig. 3 shows the documents projected to two dimensions using principal component analysis (PCA).

TABLE IV
NUMBER OF DOCUMENTS PER TOPIC IN EXPERIMENT 1

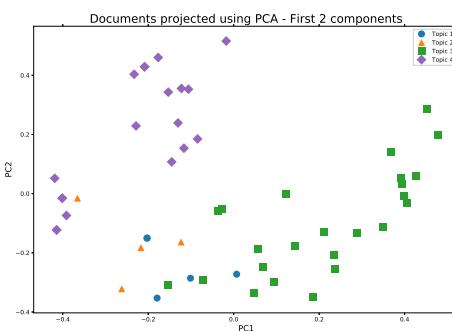| Topics          | 1 | 2 | 3  | 4  | 5 | 6 | 7 |
|-----------------|---|---|----|----|---|---|---|
| # of documents  | 5 | 4 | 26 | 19 | 0 | 0 | 0 |



Fig. 3. Documents projected into the first two dimensions using principal component analysis. Points with the same color and marker belong to the same cluster.

Using a minimum document frequency of 35% kept only 23 valid terms. Fig. 4 shows the essential terms per topic where the terms that are exclusively important for a single topic (a term is vital if it is above the 75th percentile of all weights) are denoted in green. In this plot, topics 3 and 4 share almost all terms. By adjusting the document-term matrix values using the NCut vectorizer, the factorization split topics 3 and 4 using the conditional *if* term. Topic 4 is exclusive for code snippets that are solved using conditional statements, whereas topic 3 comprises the opposite.
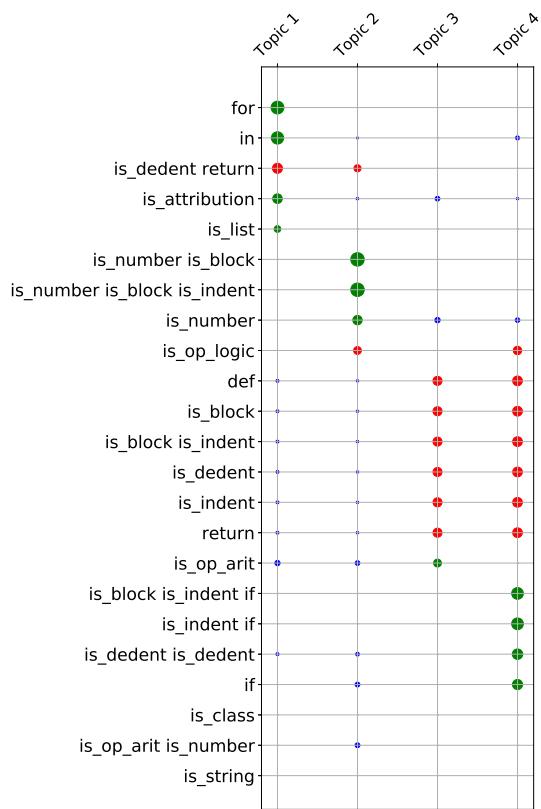


Fig. 4. Term importance for the four most populated topics. Each row represents a term, the size of the point corresponds to the term's weight for the topic, and red points are the points above the 75th percentile of all weights. The green points denote words above the 75th percentile limit on only one of the four topics.

Fig. 5 shows the topic distribution per document. As explained in Section III-B2 (topic extraction), distribution over topics describes a document. Darker cells imply that the topic characterizes a document better. As stated before, if a topic strongly characterizes a document, then we can hard-assign it to a single topic. However, Fig. 5 shows that most documents assigned to topics 1 and 2 (top part of the plot) spread throughout the topics. It suggests we have to combine the most important terms for each topic to interpret these code snippets.

Analyzing the code snippets from topic 1, they combine for-loops with conditional statements. Topic 2 is a mixture containing the code snippets that do not belong to any other topic.

Fig. 6, using the LDAVis tool [65], calculates the topics' distance and projects them to 2D using principal coordinate analysis. Topics 3 and 4 are located close to each other, and they correspond to 45% of the terms and 83% of the documents. Fig. 6 also validates that these topics are not that different when their crucial terms are analyzed. Still, the conditional statements that characterize topic 4 are enough to produce a linearly separable 2D data projection, except for a few outliers, as shown in Fig. 3.
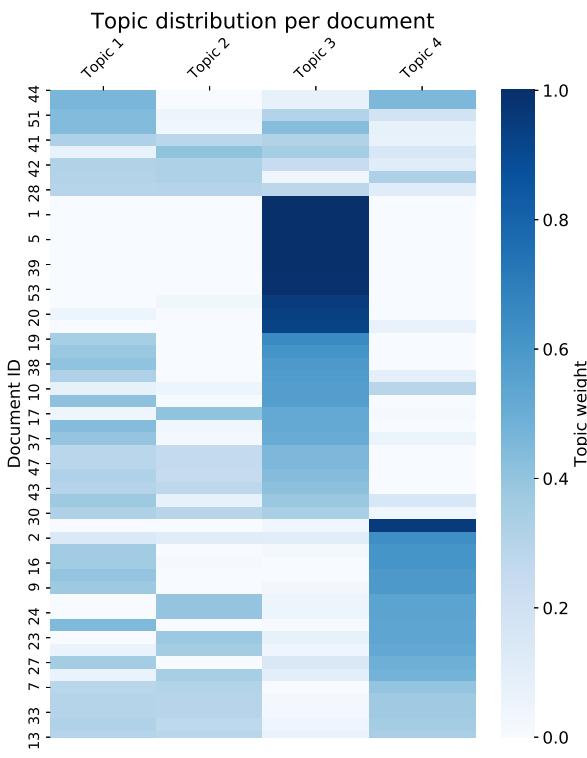
## Topic distribution per document



Fig. 5. Topic distribution per document. Darker cells indicate a better description of the document by the topic.
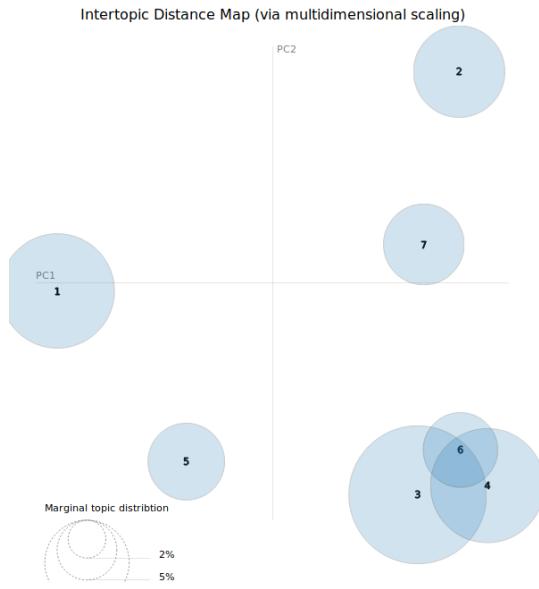


Fig. 6. Topics are represented as circles proportional to the number of terms whose weights are most associated with the topic and the distance between the circles is the intertopic distance.

### B. Experiment 2

Table V shows the number of assigned documents per topic with hyperparameters combination producing a more uniform grouping scheme than the previous one. Although

we initially set 12 clusters, two of them (topics 9 and 11) are empty after assigning each document to the topic with the highest relevance (weight). Topics 6, 8, 10, and 12 have the largest number of documents. Fig. 7 shows the topic per document distribution where the topics better characterize each document.

TABLE V
NUMBER OF DOCUMENTS PER TOPIC IN EXPERIMENT 2

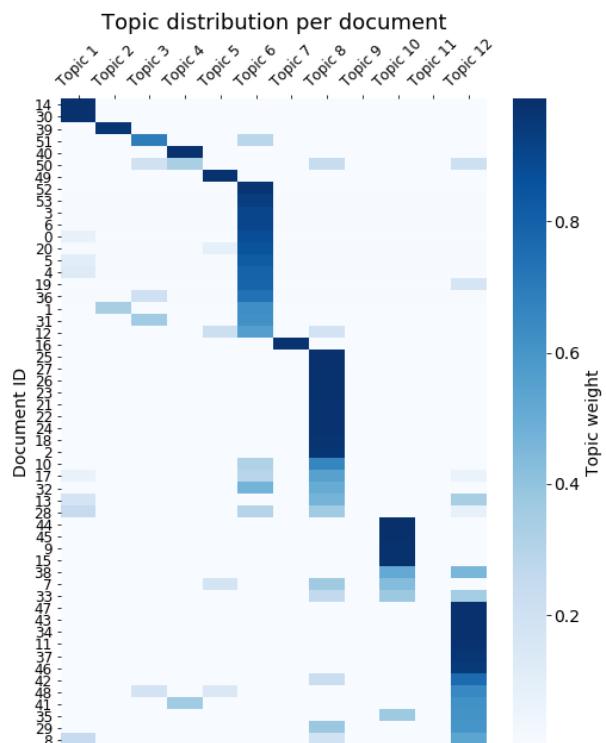| Topics | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # of documents | 2 | 1 | 1 | 2 | 1 | **13** | 1 | **14** | 0 | **7** | 0 | **12** |

## Topic distribution per document



Fig. 7. Topic distribution per document. Darker cells indicate a better description of the document by the topic.

The complete term per topic plot for this experiment is omitted because it is long and challenging to read. Using a minimum document frequency of 5% increased the number of terms to 236.

Fig. 8 shows the intertopic map. The main topics 6, 8, 10, and 12 correspond to 85% of the documents and 77.4% of the terms and we do not observe any main topic overlap in this plot. The next subsections analyze these main topics in detail. Topics 2 and 4 will also be analyzed since they occupy a different space on the map. This step belongs to the topic filter and selection phase from the code-clustering pipeline depicted in Fig. 1. After hard-assigning the documents to the clusters (removing topics 9 and 11), merging topics 2 and 4, and removing topics with a few documents (less than three

documents per topic: topics 1, 3, 5, and 7), it resulted in five conceptual clusters (six from the original topics in total) to be analyzed in detail.
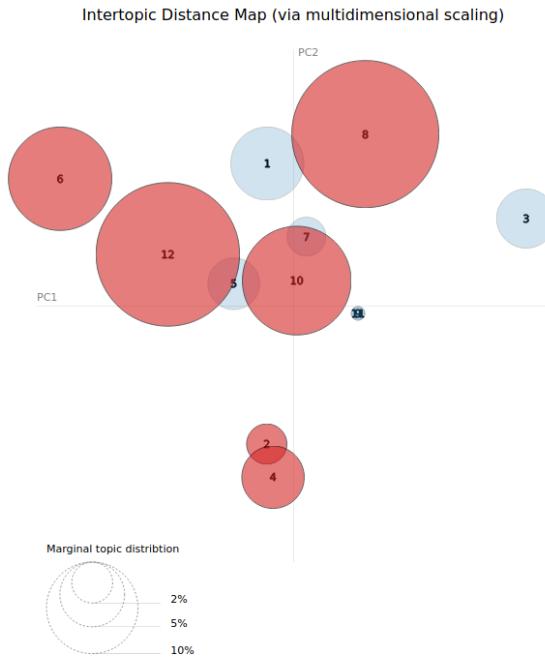


Fig. 8. Intertopic distance map for experiment 2. Topics are represented as circles proportional to the number of terms whose weights are most associated with the topic and the distance between the circles is the intertopic distance. Topics in red are further discussed in detail.

Using the Sievert and Sheley relevance metric implemented in the LDAVis tool [65], the top-30 most relevant terms per topic were extracted. Table VI shows the five most relevant terms. We wrote a description for each topic after analyzing the essential terms and code snippets from each class. Fig. 9 shows some examples of code snippets for each topic, highlighting the terms used to define them.

- Topic 8 is strongly characterized by conditional statements, logical operators, and Boolean values.
- Topic 6 does not seem to have a clear definition by just inspecting its terms. Indentation terms appear among the most relevant ones. By analyzing the code snippets, topic 6 comprises code with one indentation structure (simple coding structures), sometimes without assigning variables to solve the exercise.
- Topic 12 is characterized by for-loops, especially range loops combined with arithmetic operations.
- Topics 2 and 4 present for-loops, conditionals, and lists, and these tools are used to perform string operations. A strong mark is the presence of the auxiliary functions split and join.
- Topic 10 is about lists and their usual operations: for-loops, conditionals, and appending elements.

To better understand the topics' differences, we performed a topic dissimilarity analysis shown in Fig. 10, representing the real distances of the topics, not reduced to two dimensions as in Fig. 8. In this analysis, topic 8 (conditional) is the

TABLE VI
FIVE MOST RELEVANT TERMS FOR EACH OF THE ANALYZED TOPICS. THE TERMS STARTING WITH 'IS' ARE THE SPECIAL ANNOTATED TERMS EXPLAINED IN SECTION III-B1 (DATA TRANSFORMATION)

| Topics 2 & 4 | Topic 6 | Topic 8 | Topic 10 | Topic 12 |
|---|---|---|---|---|
| split | is_op_arit | is_op_logic | append | range |
| is_string | is_indent | if | is_list | is_op_arit |
| join | def | else | for | for |
| for | is_number | True | is_attribution | is_number |
| if | len | False | if | is_attribution |

furthest from all others. Topic 10 is central, being slightly closer to topic 4 than to topic 12. This result is reasonable as it contains elements between these two topics (numeric and range loops vs. loops with strings). Table VI shows their differences as being the data types, which can be observed in the most important terms ("range" and "is_number" for topic 12, "append" and "is_list" for topic 10 and "split", and "is_string" for topics 2 and 4).

We also analyzed how the test dataset fits into the six (including topics 2 and 4) valid topics to understand if the topics are representative of the possible concepts present in an unseen code. We assigned each code to the topic with the highest weight as we did for the training set. Table VII shows the number of assigned documents per topic. Except for two documents, all the others belong to one of the six valid topics. It confirms that the different topics (the ones considered invalid) detect specific code traits and not their general concepts. It is important to notice that topic modeling is a soft clustering technique: a document has a probability of belonging to each topic and can be associated with more than one. So, a document can be related to the main topic with its specificity related to minor ones.

TABLE VII
NUMBER OF TEST SET DOCUMENTS PER TOPIC FOR EXPERIMENT 2

| Topics | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # of documents | 2 | 4 | 0 | 1 | 0 | 23 | 0 | 17 | 0 | 3 | 0 | 15 |

### C. Coherence Evaluation

Both experiments were analyzed using the UCI coherence metric with NPMI [9], as described in Section III-B3 (topic filter and external evaluation), to validate how well the proposed methodology performs in an external dataset. Although AST trees have been used to cluster code, they do not provide an intuitive way to analyze the important features besides reading it. We compare our results with a K-means clustering method using the proposed augmented tokenizer and logistic regression to extract the important features per cluster as a baseline. We also compared our best results using the standard tokenizer instead of our proposed tokenizer. We used $k = 5$ for K-means since there were five main conceptual clusters found in the LDA. We ran each method 100 times and averaged their UCI coherence metric. Statistical difference was measured using the Mann–Whitney U test [66], and all the results were statistically significant with $p < 0.001$. Table VIII reports

```python
def sort_csv(csv):
    items = csv.split(', ')
    items.sort()
    return ", ".join(items)
```

(a) Topic 2

```python
def sort_dedupe(words):
    items = words.split(' ')
    items_dedupe = []
    for word in items:
        if word not in items_dedupe:
            items_dedupe.append(word)
    items_dedupe.sort()
    return " ".join(items_dedupe)
```

(b) Topic 4

```python
import math
def formula(D):
    C = 50
    H = 30
    Q = round(math.sqrt(2*C*D/float(H)))
    return Q
```

(c) Topic 6

```python
def is_prime(number):
    '''Returns True for prime numbers, False otherwise'''
    #Edge Cases
    if number == 1:
        prime = False
    elif number == 2:
        prime = True
    #All other primes
    else:
        prime = True
        for check_number in range(2, int(number/2)+1):
            if number % check_number == 0:
                prime = False
                break
    return prime
```

(d) Topic 8

```python
def dedupe(dup_list):
    nodup_list = []
    for i in dup_list:
        if i not in nodup_list:
            nodup_list.append(i)
    return nodup_list
```

(e) Topic 10

```python
def fatorial(number):
    total = 1

    for i in range(number, 1, -1):
        total = total * i

    return total
```

(f) Topic 12

Fig. 9. Example of code snippets belonging to each topic. In red are the relevant terms for each topic.
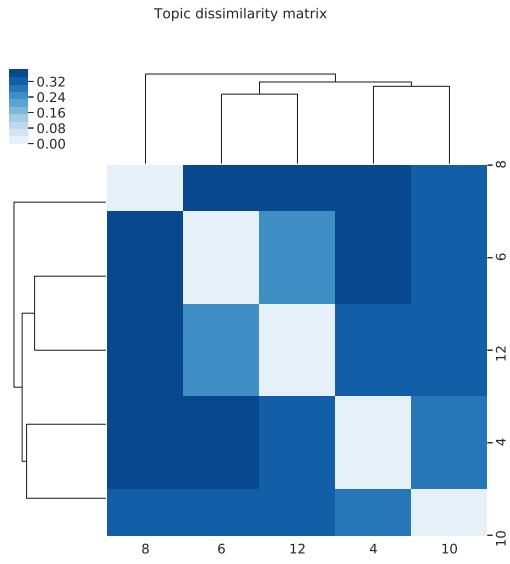


Fig. 10. Topic dissimilarity matrix. The darker the cell more distant the topics are from each other.

the mean and standard deviation for each experiment. In the UCI coherence with NPMI metric, the values are bounded between 1 and -1, where 1 means that the top words only occur together, zero means that they are distributed as expected under independence, and -1 means that they only occur separately. The UCI coherence for the standard tokenizer using the top-10 terms could not be measured because there were no important top-10 term pairwise combinations in this setting that appeared in at least one document. The NMF experiment with the augmented tokenizer considering the top-10 terms demonstrated the best UCI occurrence metric, followed by the LDA experiment, which had the best performance considering the top-5 terms.

TABLE VIII
THE MEAN AND STANDARD DEVIATION OF UCI COHERENCE USING NPMI

| | $C_{UCI}$ (NPMI) | |
|---|---|---|
| | Top-5 terms | Top-10 terms |
| NMF (Experiment 1 with augmented tokenizer) | 0.73 (0.03) | **0.83 (0.02)** |
| LDA (Experiment 2 with augmented tokenizer) | **0.75 (0.06)** | 0.76 (0.04) |
| K-Means (with augmented tokenizer) | 0.55 (0.09) | 0.53 (0.03) |
| LDA (best result for standard tokenizer) | 0.53 (0.03) | — |

### D. Discussion about Experiments 1 and 2

Experiments 1 and 2 are the two best-ranked results on the top-5 and top-10 UMass coherence metric. Both experiments have different hyperparameters. Experiment 1 uses a minimum document frequency of 35%, NCut to weight the terms, and then perform NMF to extract the topics. Experiment 2 uses a minimum document frequency of 5%, regular count of words, and LDA to extract the topics.

We found both experiments to have their main concepts in a few clusters (two main clusters in Experiment 1 and six main clusters in Experiment 2). The remaining clusters are associated with code specificity. In the case of Experiment 1, using NCut and a high document frequency threshold, the topic modeling from Experiment 1 focused on finding structures with high volume and cooccurrence rates, resulting in separation of the *if/else* structure from the rest. The conditional structure was first separated from a hierarchical perspective, and the remaining structures were all grouped in a cluster. In Experiment 2, the conditional topic (topic 8) is also the furthest from the other topics. As shown in the hierarchical clustering of Fig. 9, this topic is the last one to be aggregated (or the first one to be separated). The common code snippets between the conditional clusters in each experiment also validate this result. From the 14 code snippets associated with the conditional topic (topic 8) in Experiment 2, 11 of them (79%) belong to the conditional topic in Experiment 1 (topic 4). Therefore, Experiment 2 demonstrates more granularity than Experiment 1.

### E. Topics Contextualization

Experiment 2 shows a better distribution of exercises and less overlap among topics than Experiment 1; we followed up the analysis with this experiment, which was the second-best result considering the $C_{UCI}$ with NPMI metric.

- **Theme identification:** after tokenizing the labels given by the professors, we separated the most frequent tokens and manually created the topic titles, as follows. All the tokens in the presented titles appeared in at least 50% of the labels.
  1) Topics 2 & 4: string manipulation,
  2) Topic 6: math functions,
  3) Topic 8: conditional structure,
  4) Topic 10: list loops,
  5) Topic 12: math loops.
- **Concept identification:** each professor was asked to associate up to three concepts (from the 15 available in Table II) to each presented code. Four professors analyzed each code. In 37 of the 54 code snippets, there was at least one concept in common between all four professors. In 53 of the 54 code snippets, at least one concept was common between three out of the four professors (75%). Therefore, we decided to use the 75% threshold of the agreement to relate the exercises' concepts. The concepts in each topic were aggregated to provide an overview of the main concepts needed to solve the cluster's problems, as summarized in Table IX. These results validate the topic themes defined in the previous task: the professors also elected each topic's main concept as a word in free-text labels. Notice that we performed the two tasks independently.
- **Intruder identification:** four code snippets were presented in each of the five groups, three belonging to the same topic (randomly chosen from the topic pool) and an intruder (also randomly chosen from another topic). Fig. 11 shows a confusion matrix, presenting how well the

TABLE IX
THE RELATION BETWEEN THE FOUND TOPICS AND THE MAIN CS1
CONCEPT RELATED TO THE TOPIC

| Topic | Main Concept | Agreement inside cluster |
|---|---|---|
| 2 & 4 - String manipulation | 7. Data type: string | 66.7% |
| 6 - Math functions | 14. Function | 50% |
| 8 - Conditional structure | 11. Conditional | 86.7% |
| 10 - List loops | 12. Loop | 71.4% |
| 12 - Math loops | 12. Loop | 75% |

professors could distinguish the intruder in each topic. In this figure, each row sums to 1 and represents how often the professors correctly guess the intruder (the diagonal values) and how often the professors confuse the intruder (the intruder cluster is depicted in the columns). We can draw some insights from this analysis:

- The "conditional structure" topic performed well, with the intruder code being identified 79% of the time, meaning that identifying a code snippet from a different cluster can be done 4 out of 5 times.
- The intruder code inside the "math loops" topic was identified 2 out of 3 times (64%), being confused with "list loops" the last third of the time. These topics work on the same main concept, as seen in the concept identification task.
- The same behavior is not seen for the "list loops" topic. This topic and the "string manipulation" topic present very similar behavior. They can be distinguished frequently (2x better than the 25% random baseline), but we do not see a confusion pattern.
- On the other hand, the "math functions" topic seemed to confuse the professors, being switched half of the time with the "math loops" topic. By backtracking the previous two tasks and analyzing them with this point of view, this topic's main concept was function, present in all exercises. As stated in Section IV-B, this topic does not have well-defined terms; indentation terms appeared among the most relevant. Although the indentations when programming in Python can indicate the difficulty of an exercise, they are not a natural human way of splitting them. Therefore, it is hard for humans to interpret this kind of clustering scheme.
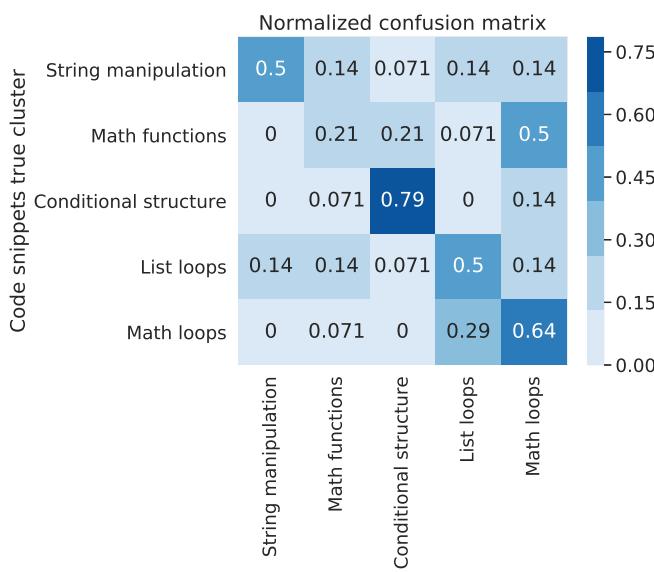
Fig. 11. Normalized confusion matrix for the intruder-identification task.

## V. Conclusions and Future Work

Based on the evaluation metric, our proposed method found semantically related code-clustering schemes suitable for human-interpretability with minimal supervision, giving support for RQ1. The method is expected to provide semantics for large amounts of unannotated code.

Although code clustering in the CS1 context has been widely applied using the AST trees, the advantages of working with topic modeling are the terms per topic results that may help experts better assess each cluster's contents. The methodology has also been shown to overcome the small-sized code snippets challenge by extending the tokenizer to augment the corpus with the code structure. The standard tokenizer could not create semantically related topics, but adding structural information, as features: indents and data types, and enforcing the order using n-grams enriched the code representation and found topics suitable for human-interpretability. For example, augmenting the corpus with structural information as indents/blocks (in Python, indents indicate how deep a block of code is; other languages like C++ and Java could count the number of "{" and "}") helps to separate single loops from nested loops. Combining trigrams (to enforce order) with structural information can distinguish subtle differences in precondition and postcondition loops. Notice that postcondition loops do not exist in Python, so we could not verify this specific assumption. In our dataset, we expect trigram tokens to be enough to capture these varieties because a typical CS1 solution does not have more than three or four nested structures. Still, it may limit our model in identifying large nested structures on more complex code. Also, even though there is a recursion concept in the concepts list, there was no exercise using this technique in our dataset to verify how it would be clustered.

In the second experiment, we set the best number of topics to 12, but only six (considering the merge of topics 2 and 4) were valid. Standard topic modeling techniques with the augmented tokenizer yielded the best results. The LDA-based clustering demonstrates better interpretable results, as shown by our detailed topic analysis. When tested with an unseen dataset, the six topics comprised the main concepts present in the test set. Except for a few exceptions, most of them relate to the six main topics, and a few of them account for code specificity.

To understand how humans could read and interpret the extracted relationships (RQ2), we asked 14 professors to contextualize CS1 concepts. Our results showed that professors could relate topics and concepts with 54 observations in the training set. Each cluster's main topic was explicit in the theme label task, and the clusters contained different concepts per topic. The professors can use this information to understand how often concepts cooccur while solving exercises. A future research direction is necessary to investigate if more specific clusters could be produced by increasing the number of samples or using more complex models.

Although there are potentials to use these experiments in a real-world environment, they require further research to effectively and practically integrate them into professors' working tools. These tools should support other programming languages as well.

## References

[1] M. C. Desmarais, "Mapping question items to skills with non-negative matrix factorization," *SIGKDD Explorations Newslett.*, vol. 13, no. 2, pp. 30–36, May 2012, doi: 10.1145/2207243.2207248.

[2] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, pp. 788–791, Oct. 1999, doi: 10.1038/44565.

[3] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, Jan. 2003.

[4] M. Steyvers and T. Griffiths, "Probabilistic topic models," in *Handbook of Latent Semantic Analysis*, 1st ed., T. Landauer, S. D. McNamara, and W. Kintsch, Eds. New York, NY, USA: Psychology Press, 2007, ch. 21, pp. 427–448, doi: 10.4324/9780203936399.

[5] T. Hofmann, "Probabilistic latent semantic analysis," in *Proc. 15th Conf. Uncertainty Artificial Intelligence*, K. B. Laskey and H. M. Prade, Eds., Stockholm, Sweden, 30 Jul.–1 Aug. 1999, pp. 289–296.

[6] S. Lloyd, "Least squares quantization in PCM," *IEEE Trans. Inf. Theory*, vol. 28, no. 2, pp. 129–137, Mar. 1982, doi: 10.1109/TIT.1982.1056489.

[7] D. Newman, J. H. Lau, K. Grieser, and T. Baldwin, "Automatic evaluation of topic coherence," in *Human Language Technologies: 2010 Annu. Conf. North American Chapter Association Computational Linguistics*, Los Angeles, CA, USA, 1–6 Jun. 2010, pp. 100–108.

[8] D. Mimno, H. M. Wallach, E. Talley, M. Leenders, and A. McCallum, "Optimizing semantic coherence in topic models," in *Proc. Conf. Empirical Methods Natural Language Processing*, R. Barzilay and M. Johnson, Eds., Edinburgh, Scotland, United Kingdom, 27–29 Jul. 2011, pp. 262–272.

[9] N. Aletras and M. Stevenson, "Evaluating topic coherence using distributional semantics," in *Proc. 10th Int. Conf. Computational Semantics*, A. Koller and K. Erk, Eds., Potsdam, Germany, 19–22 Mar. 2013, pp. 13–22.

[10] M. Röder, A. Both, and A. Hinneburg, "Exploring the space of topic coherence measures," in *Proc. 8th ACM Int. Conf. Web Search and Data Mining*, Shanghai, China, 31 Jan.–6 Feb. 2015, pp. 399–408, doi: 10.1145/2684822.2685324.

[11] R. Fagin, "Combining fuzzy information from multiple systems," in *Proc. 15th ACM SIGACT-SIGMOD-SIGART Symp. Principles Database Systems*, R. Hull, Ed., Montreal, Canada, 3–5 Jun. 1996, pp. 216–226, doi: 10.1145/237661.237715.

[12] J. Sheard *et al.*, "Exploring programming assessment instruments: A classification scheme for examination questions," in *Proc. 7th Int. Computing Education Research Workshop*, Providence, RI, USA, 8–9 Aug. 2011, pp. 33–38, doi: 10.1145/2016911.2016920.

[13] A. Petersen, M. Craig, and D. Zingaro, "Reviewing CS1 exam question content," in *Proc. 42nd ACM Technical Symp. Computer Science Education*, T. J. Cortina, E. L. Walker, L. A. S. King, and D. R. Musicant, Eds., Dallas, TX, USA, 9–12 Mar. 2011, pp. 631–636, doi: 10.1145/1953163.1953340.

[14] T. Fei, W. J. Heng, K. C. Toh, and T. Qi, "Question classification for e-learning by artificial neural network," in *Proc. 2003 Joint Conf. 4th Int. Conf. Information, Communications and Signal Processing and 4th Pacific Rim Conf. Multimedia*, Singapore, 15–18 Dec. 2003, pp. 1757–1761, doi: 10.1109/ICICS.2003.1292768.

[15] S. Supraja, K. Hartman, and A. W. H. Khong, "Toward the automatic labeling of course questions for ensuring their alignment with learning outcomes," in *Proc. 10th Int. Conf. Educational Data Mining*, X. Hu, T. Barnes, A. Hershkovitz, and L. Paquette, Eds., Wuhan, China, 25–28 Jun. 2017, pp. 56–63.

[16] A. Godea, D. Tulley-Patton, S. Barbee, and R. Nielsen, "Classifying educational questions based on the expected characteristics of answers," in *Int. Conf. Artificial Intelligence Education*, C. Penstein Rosé, R. Martínez-Maldonado, H. U. Hoppe, R. Luckin, M. Mavrikis, K. Porayska-Pomsta, B. McLaren, and B. du Boulay, Eds., London, The United Kingdom, 20 Jun. 2018, pp. 104–108, doi: 10.1007/978-3-319-93846-2_20.

[17] P. González, E. Gibaja, A. Zapata, V. H. Menéndez, and C. Romero, "Towards automatic classification of learning objects: Reducing the number of used features," in *Proc. 10th Int. Conf. Educational Data Mining*, X. Hu, T. Barnes, A. Hershkovitz, and L. Paquette, Eds., Wuhan, China, 25–28 Jun. 2017, pp. 394–395.

[18] S. Nunn, J. T. Avella, T. Kanai, and M. Kebritchi, "Learning analytics methods, benefits, and challenges in higher education: A systematic literature review," *Online Learning*, vol. 20, no. 2, Jun. 2016, doi: 10.24059/olj.v20i2.790.

[19] A. Dutt, M. A. Ismail, and T. Herawan, "A systematic review on educational data mining," *IEEE Access*, vol. 5, pp. 15 991–16 005, Jan. 2017, doi: 10.1109/ACCESS.2017.2654247.

[20] M. W. Rodrigues, S. Isotani, and L. E. Zárate, "Educational data mining: A review of evaluation process in the e-learning," *Telematics and Inform.*, vol. 35, no. 6, pp. 1701–1717, Sep. 2018, doi: 10.1016/j.tele.2018.04.015.

[21] A. Hernández-Blanco, B. Herrera-Flores, D. Tomás, and B. Navarro-Colorado, "A systematic review of deep learning approaches to educational data mining," *Complexity*, vol. 2019, May 2019, doi: 10.1155/2019/1306039.

[22] H. Aldowah, H. Al-Samarraie, and W. M. Fauzy, "Educational data mining and learning analytics for 21st century higher education: A review and synthesis," *Telematics and Inform.*, vol. 37, pp. 13–49, Apr. 2019, doi: 10.1016/j.tele.2019.01.007.

[23] S. Trivedi, Z. Pardos, G. Sárközy, and N. Heffernan, "Spectral clustering in educational data mining," in *Proc. 4th Int. Conf. Educational Data Mining*, M. Pechenizkiy, T. Calders, C. Conati, S. Ventura, C. Romero, and J. Stamper, Eds., Eindhoven, the Netherlands, 6–8 Jul. 2011, pp. 129–138.

[24] A. P. d. B. B. R. Figueira, "A repository with semantic organization for educational content," in *2008 8th IEEE Int. Conf. Advanced Learning Technologies*, Santander, Spain, 1–5 Jul. 2008, pp. 114–116, doi: 10.1109/ICALT.2008.60.

[25] J. Huang, C. Piech, A. Nguyen, and L. Guibas, "Syntactic and functional variability of a million code submissions in a machine learning MOOC," in *Proc. Workshops 16th Int. Conf. Artificial Intelligence Education 2013*, E. Walker and C.-K. Looi, Eds., vol. 1, Memphis, TN, USA, 9–13 Jul. 2013, pp. 25–32, doi: 10.1007/978-3-642-39112-5.

[26] A. Nguyen, C. Piech, J. Huang, and L. Guibas, "Codewebs: Scalable homework search for massive open online programming courses," in *Proc. 23rd Int. Conf. World Wide Web*, Seoul, Korea, 7–11 Apr. 2014, pp. 491–502, doi: 10.1145/2566486.2568023.

[27] B. Paassen, B. Mokbel, and B. Hammer, "Adaptive structure metrics for automated feedback provision in intelligent tutoring systems," *Neurocomputing*, vol. 192, pp. 3–13, Jun. 2016, doi: 10.1016/j.neucom.2015.12.108.

[28] T. Price, R. Zhi, and T. Barnes, "Evaluation of a data-driven feedback algorithm for open-ended programming," in *Proc. 10th Int. Conf. Educational Data Mining*, X. Hu, T. Barnes, A. Hershkovitz, and L. Paquette, Eds., Wuhan, China, 25–28 Jun. 2017, pp. 192–197.

[29] B. Mokbel, S. Gross, B. Paassen, N. Pinkwart, and B. Hammer, "Domain-independent proximity measures in intelligent tutoring systems," in *Proc. 6th Int. Conf. Educational Data Mining*, S. K. D'Mello, R. A. Calvo, and A. Olney, Eds., Memphis, TN, USA, 6–9 Jul. 2013, pp. 334–335.

[30] Y. Lu and I.-H. Hsiao, "Modeling semantics between programming codes and annotations," in *Proc. 29th Hypertext and Social Media*, D. Lee, N. Sastry, and I. Weber, Eds., Baltimore, MD, USA, 9–12 Jul. 2018, pp. 101–105, doi: 10.1145/3209542.3209578.

[31] T. H. Chen, S. W. Thomas, and A. E. Hassan, "A survey on the use of topic models when mining software repositories," *Empirical Softw. Eng.*, vol. 21, no. 5, pp. 1843–1919, Oct. 2016, doi: 10.1007/s10664-015-9402-8.

[32] D. Azcona, P. Arora, I.-H. Hsiao, and A. Smeaton, "User2code2vec: Embeddings for profiling students based on distributional representations of source code," in *Proc. 9th Int. Learning Analytics and Knowledge Conf.*, Tempe, AZ, USA, 4–8 Mar. 2019, pp. 86–95, doi: 10.1145/3303772.3303813.

[33] Q. Chen, L. Yao, and J. Yang, "Short text classification based on LDA topic model," in *2016 Int. Conf. Audio, Language and Image Processing*, Shanghai, China, 11-12 Jul. 2016, pp. 749–753, doi: 10.1109/ICALIP.2016.7846525.

[34] N. Abolhassani and L. Ramaswamy, "Extracting topics from semi-structured data for enhancing enterprise knowledge graphs," in *Int. Conf. Collaborative Computing: Networking, Applications and Worksharing*, X. Wang, H. Gao, M. Iqbal, and G. Min, Eds., vol. 292, London, United Kingdom, 19–22 Aug. 2019, pp. 101–117, doi: 10.1007/978-3-030-30146-0_8.

[35] I.-H. Hsiao and P. Awasthi, "Topic facet modeling: Semantic visual analytics for online discussion forums," in *Proc. 5th Int. Conf. Learning Analytics Knowledge*, J. Baron, G. Lynch, N. Maziarz, P. Blikstein, A. Merceron, and G. Siemens, Eds., Poughkeepsie, NY, USA, 16–20 Mar. 2015, pp. 231–235, doi: 10.1145/2723576.2723613.

[36] I.-H. Hsiao and Y.-L. Lin, "Enriching programming content semantics: An evaluation of visual analytics approach," *Comput. Human Behav.*, vol. 72, pp. 771–782, Jul. 2017, doi: 10.1016/j.chb.2016.10.012.

[37] Zhao *et al.*, "Comparing Twitter and traditional media using topic models," in *Proc. 33rd European Conf. Advances Information Retrieval*, P. Clough, C. Foley, C. Gurrin, G. J. Jones, W. Kraaij, H. Lee, and V. Murdoch, Eds., Dublin, Ireland, 18–21 Apr. 2011, vol. 6611, pp. 338–349, doi: 10.1007/978-3-642-20161-5_34.

[38] M. Steyvers, P. Smyth, M. Rosen-Zvi, and T. Griffiths, "Probabilistic author-topic models for information discovery," in *Proc. 10th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, R. Kohavi, J. Gehrkeand, W. DuMouchel, and J. Ghosh, Eds., Seattle, WA, USA, 22-25 Aug. 2004, pp. 306–315, doi: 10.1145/1014052.1014087.

[39] M. Rosen-Zvi, C. Chemudugunta, T. Griffiths, P. Smyth, and M. Steyvers, "Learning author-topic models from text corpora," *ACM Trans. Inf. Syst.*, vol. 28, no. 1, pp. 1–38, Jan. 2010, Art. no. 4, doi: 10.1145/1658377.1658381.

[40] S. Li, J. Li, and R. Pan, "Tag-weighted topic model for mining semi-structured documents," in *Proc. 23rd Int. Joint Conf. Artificial Intelligence*, F. Rossi, Ed., Beijing, China, 3–9 Aug. 2013, pp. 2855–2861.

[41] J. Weng, E.-P. Lim, J. Jiang, and Q. He, "TwitterRank: Finding topic-sensitive influential twitterers," in *Proc. 3rd ACM Int. Conf. Web Search and Data Mining*, B. D. Davison, T. Suel, N. Craswell, and B. Liu, Eds., New York, NY, USA, 3–6 Feb. 2010, pp. 261–270, doi: 10.1145/1718487.1718520.

[42] S. Syed and M. Spruit, "Full-text or abstract? Examining topic coherence scores using latent Dirichlet allocation," in *2017 IEEE Int. Conf. Data Science and Advanced Analytics*, Tokyo, Japan, 19–21 Oct. 2017, pp. 165–174, doi: 10.1109/DSAA.2017.61.

[43] M. Pratusevich, "Practice python," www.practicepython.org, 2017, [Online; accessed 07-Jan-2021].

[44] S. Sentance and A. McNicol, "Python school," https://pythonschool.net/, 2016, [Online; accessed 07-Jan-2021].

[45] J. Hu, "Python programming exercises," https://github.com/zhiwehu/Python-programming-exercises, 2018, [Online; accessed 07-Jan-2021].

[46] W3Resource, "W3Resource," https://www.w3resource.com/python/python-tutorial.php, 2018, [Online; accessed 07-Jan-2021].

[47] G. Salton and M. J. McGill, *Introduction to Modern Information Retrieval*. New York, NY, USA: McGraw-Hill, 1986.

[48] W. Zhang, T. Yoshida, and X. Tang, "A comparative study of TF*IDF, LSI and multi-words for text classification," *Expert Syst. Appl.*, vol. 38, no. 3, pp. 2758–2765, Mar. 2011, doi: 10.1016/j.eswa.2010.08.006.

[49] X. Yan, J. Guo, S. Liu, X.-q. Cheng, and Y. Wang, "Clustering short text using NCut-weighted non-negative matrix factorization," in *Proc. 21st ACM Int. Conf. Information and Knowledge Management*, Maui, HI, USA, 29 Oct.–2 Nov. 2012, pp. 2259–2262, doi: 10.1145/2396761.2398615.

[50] A. Cichocki and A. H. Phan, "Fast local algorithms for large scale nonnegative matrix and tensor factorizations," *IEICE Trans. Fundam. Electron., Commun. and Comput. Sci.*, vol. E92.A, no. 3, pp. 708–721, Mar. 2009, doi: 10.1587/transfun.E92.A.708.

[51] D. O'Callaghan, D. Greene, J. Carthy, and P. Cunningham, "An analysis of the coherence of descriptors in topic modeling," *Expert Syst. Appl.*, vol. 42, no. 13, pp. 5645–5657, Aug. 2015, doi: 10.1016/j.eswa.2015.02.055.

[52] Y. Chen, H. Zhang, R. Liu, Z. Ye, and J. Lin, "Experimental explorations on short text topic mining between LDA and NMF based schemes," *Knowledge-Based Syst.*, vol. 163, pp. 1–13, Jan. 2019, doi: 10.1016/j.knosys.2018.08.011.

[53] J. Bergstra and Y. Bengio, "Random search for hyperparameter optimization," *J. Mach. Learn. Res.*, vol. 13, pp. 281–305, Feb. 2012.

[54] J. Chuang, S. Gupta, C. Manning, and J. Heer, "Topic model diagnostics: Assessing domain relevance via topical alignment," in *Proc. 30th Int. Conf. Machine Learning*, S. Dasgupta and D. McAllester, Eds., 16–21 Jun. 2013, pp. 612–620.

[55] C. Wang and D. M. Blei, "Collaborative topic modeling for recommending scientific articles," in *Proc. 17th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, San Diego, CA, USA, 21–24 Aug. 2011, pp. 448–456, doi: 10.1145/2020408.2020480.

[56] J. H. Lau, D. Newman, and T. Baldwin, "Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality," in *Proc. 14th Conf. European Chapter Association Computational Linguistics*, S. Wintner, S. Goldwater, and S. Riezler, Eds., Gothenburg, Sweden, 26–30 Apr. 2014, pp. 530–539, doi: 10.3115/v1/E14-1056.

[57] Joint Task Force on Computing Curricula, Association for Computing Machinery (ACM) and IEEE Computer Society, *Computer Science Curricula 2013: Curriculum Guidelines for Undergraduate Degree Programs in Computer Science*. New York, NY, USA: Association for Computing Machinery, Dec. 2013, doi: 10.1145/2534860.

[58] C. Schulte and J. Bennedsen, "What do teachers teach in introductory programming?" in *Proc. 2nd Int. Workshop Computing Education Research*, R. Anderson, S. A. Fincher, and M. Guzdial, Eds., 9–10 Sep. 2006, pp. 17–28, doi: 10.1145/1151588.1151593.

[59] A. Robins, P. Haden, and S. Garner, "Problem distributions in a CS1 course," in *Proc. 8th Australasian Conf. Computing Education*, D. Tolhurst and S. Mann, Eds., vol. 52, Hobart, Australia, Jan. 2006, pp. 165–173.

[60] Goldman *et al.*, "Setting the scope of concept inventories for introductory computing subjects," *ACM Trans. Comput. Educ.*, vol. 10, no. 2, pp. 1–29, Jun. 2010, Art. no. 5, doi: 10.1145/1789934.1789935.

[61] Y. Cherenkova, D. Zingaro, and A. Petersen, "Identifying challenging CS1 concepts in a large problem dataset," in *Proc. 45th ACM Technical Symp. Computer Science Education*, J. Dougherty, K. Nagel, A. Decker, and K. Eiselt, Eds., Atlanta, GA, USA, 5–8 Mar. 2014, pp. 695–700, doi: 10.1145/2538862.2538966.

[62] A. S. Lan, A. E. Waters, C. Studer, and R. G. Baraniuk, "Sparse factor analysis for learning and content analytics," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1959–2008, 2014.

[63] J. Chang, S. Gerrish, C. Wang, J. L. Boyd-graber, and D. M. Blei, "Reading tea leaves: How humans interpret topic models," in *Advances in Neural Information Processing Systems 22, Proc. 23rd Annu. Conf. Neural Information Processing Systems*, Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, Eds., Vancouver, British Columbia, Canada, 7–10 Dec. 2009, pp. 288–296.

[64] A. Mahmoud and G. Bradshaw, "Semantic topic models for source code analysis," *Empirical Softw. Eng.*, vol. 22, no. 4, pp. 1965–2000, Aug. 2017, doi: 10.1007/s10664-016-9473-1.

[65] C. Sievert and K. Shirley, "LDAvis: A method for visualizing and interpreting topics," in *Proc. Workshop Interactive Language Learning, Visualization, and Interfaces*, Baltimore, MD, USA, 27 Jun. 2014, pp. 63–70, doi: 10.3115/v1/W14-3110.

[66] H. B. Mann and D. R. Whitney, "On a test of whether one of two random variables is stochastically larger than the other," *Ann. Math. Statist.*, vol. 18, no. 1, pp. 50–60, Mar. 1947, doi: 10.1214/aoms/1177730491.

**Laura O. Moraes** Laura O. Moraes received her master degree in Systems and Computing Engineering from UFRJ, Rio de Janeiro, Brazil, 2016. She is currently pursuing a Ph.D. degree in the same department. In 2018, she worked as a fellow at the Data Science for Social Good Europe, organized by the University of Chicago with Nova SBE. Her research interests include text classification, item to skill mapping and the relationship between skills and student performance.

**Carlos Eduardo Pedreira** Prof. Carlos Eduardo Pedreira holds Bachelor (1975) and MSc degrees (1981) in electrical engineering from the Catholic University of Rio de Janeiro. In 1987, he received a Ph.D. degree from Imperial College of Science, Technology and Medicine, University of London. Presently, he is a Professor at the Federal University of Rio de Janeiro where he is the head of the AI sector at the Systems and Computing Department at COPPE. He is a visiting researcher at the University of Salamanca, Spain since 2002. His main research interests are Pattern Classification, Machine Learning and its applications in health. He was the founding president of the Brazilian Society of Computational Intelligence. Received the Santander Bank Award of Science and Innovation in 2006.