

Why Machines Cannot Learn Mathematics, Yet

André Greiner-Petter¹, Terry Ruas², Moritz Schubotz¹,
Akiko Aizawa³, William Grosky², Bela Gipp¹

¹ University of Wuppertal, Wuppertal, Germany
{last}@uni-wuppertal.de

² University of Michigan-Dearborn, Dearborn, USA
{truas,wgrosky}@umich.edu

³ National Institute of Informatics, Tokyo, Japan
aizawa@nii.ac.jp

Abstract. Nowadays, Machine Learning (ML) is seen as the universal solution to improve the effectiveness of information retrieval (IR) methods. However, while mathematics is a precise and accurate science, it is usually expressed by less accurate and imprecise descriptions, contributing to the relative dearth of machine learning applications for IR in this domain. Generally, mathematical documents communicate their knowledge with an ambiguous, context-dependent, and non-formal language. Given recent advances in ML, it seems canonical to apply ML techniques to represent and retrieve mathematics semantically. In this work, we apply popular text embedding techniques to the arXiv collection of STEM documents and explore how these are unable to properly understand mathematics from that corpus. In addition, we also investigate the missing aspects that would allow mathematics to be learned by computers.

Keywords: Mathematical Information Retrieval, Machine Learning, Word Embeddings, Math Embeddings, Mathematical Objects of Interest

1 Introduction

Mathematics is capable of explaining complex concepts and relations in a compact, precise, and accurate way. Learning this idiom takes time and is often difficult, even to humans. The general applicability of mathematics allows a certain level of ambiguity in its expressions. This ambiguity is regularly mitigated by short explanations following or preceding these mathematical expressions, that serve as context to the reader. Along with context dependency, inherent issues of linguistics (e.g. ambiguity, non-formality) make it even more challenging for computers to understand mathematical expressions. Said that, a system capable of capturing the semantics of mathematical expressions automatically would be suitable for several applications, from improving search engines to recommender systems.

During our evaluations of MathMLBen [31], a benchmark for converting mathematical \LaTeX expressions into MathML, it is possible to notice several fundamental problems that generally affect prominent ML approaches to learn semantics of mathematical expressions. For instance, the first entry of the benchmark,

$$W(2, k) > 2^k / k^\varepsilon \tag{1}$$

is extracted from the English Wikipedia page about Van der Waerden’s theorem⁴. Without further explanation, the symbols W , k , and ε might have several possible meanings. Depending on which one is considered, even the structure of the formula may be different. If we consider W as a variable, instead of a function, it changes the interpretation of $W(2, k)$ to a multiplication operation.

Learning connections, such as between W and the entity ‘*Van der Waerden’s number*’, requires a large specifically labeled scientific database that contains these mathematical objects. Furthermore, a fundamental understanding of the mathematical expression would increase the performance during the learning process, e.g., that $W(2, k)$ and $W(n, k)$ contain the same function.

Word embedding techniques has received significant attention over the last years in the Natural Language Processing (NLP) community, especially after the publication of word2vec [18]. Recently, more and more projects try to adapt this knowledge for solving Mathematical Information Retrieval (MIR) tasks [4, 12, 36, 34]. While all of these projects follow similar approaches and obtain promising results, all of them fail to understand mathematical expressions because of the same fundamental issues. In this paper, we explore some of the main aspects that we believe are necessary to leverage the learning of mathematics by computer systems. We explain, with our evaluations of word embedding techniques on the arXMLiv 2018 [5] dataset, why current ML approaches are not applicable for MIR tasks, yet.

2 Background & Related Work

Understanding mathematical expressions essentially means comprehending the semantic value of its internal components, which can be accomplished by linking its elements with their corresponding mathematical definitions. Current MIR approaches [11, 32, 30] try to extract textual descriptors of the parts that compose mathematical equations. Intuitively, there are questions that arise from this scenario, such as (i) how to determine the parts which have their own descriptors, and (ii) how to identify correct descriptors over others.

Answers to (i) are more concerned in choosing the correct definitions for which parts of a mathematical expression should be considered as one mathematical object [10, 35, 31]. Current definitions, such as the content MathML

⁴ https://en.wikipedia.org/wiki/Van_der_Waerden's_theorem

3.0⁵ specification, are often imprecise⁶. For example, content MathML 3.0 uses `csymbol` elements for functions and specifies them as expressions that *refer to a specific, mathematically-defined concept with an external definition*⁷. However, it is not clear whether W or the sequence $W(2, k)$ (Equation 1) should be declared as a `csymbol`. Another example are content identifiers, which MathML specifies as *mathematical variables which have properties, but no fixed value*⁸. While content identifiers are allowed to have complex rendered structures (e.g. β_i^2), it is not permitted to enclose identifiers within other identifiers. Let us consider α_i , where α is a vector and α_i its i -th element. In this case, α_i should be considered as a composition of three content identifiers, each one carrying its own individualized semantic information, namely the vector α , the element α_i of the vector, and the index i . However, with the current specification, the definition of these identifiers would not be canonical. One possible workaround to represent such expressions with content MathML is to use a structure of four nodes, interpreting α_i as a function with the `csymbol` *vector-selector*. However, ML algorithms and MIR approaches would benefit from more precise definitions and a unified answer for (i). Most of the related work relies on these relatively vague definitions and in the analysis of content identifiers, focusing their efforts on (ii).

In [32], an approach is presented for scoring pairs of identifiers and *definiens*⁹ by the number of words between them. Their approach is based on the assumption that correct definiens appear close to the identifier and to the complex mathematical expression that contains this same identifier. Kristianto et al. [11] introduce an ML approach, in which they train a Support Vector Machine (SVM) to consider sentence patterns and other characteristics as features (e.g. part-of-speech (POS) tags, parse trees). Later, [30] combine the aforementioned approaches and use pattern recognition based on the POS tags of common identifier-definiens pairs, the distance measurements, and SVM, reporting results for precision and recall of 48.60% and 28.06%, respectively. These results can be considered as a baseline for MIR tasks.

More recently, some projects try to use embedding techniques to learn patterns of the correlations between context and mathematics. In the work of [4], they embed single symbols and train a model that is able to discover similarities between mathematical symbols. Similarly to this approach, Krstovski and Blei [12] use a variation of word embeddings (briefly discussed in Section 3) to represent complex mathematical expressions as single unit tokens for IR. In 2019, M. Yasunaga et al. [34] explore an embedding technique based on recurrent neural networks to improve topic models by considering mathematical expres-

⁵ <https://www.w3.org/TR/MathML3/>

⁶ Note that OpenMath is another specification specifically designed for encoding semantics of mathematics. However, content MathML is an encoding of OpenMath and inherent problems of content MathML also apply for OpenMath (see <https://www.openmath.org/om-mml/>).

⁷ <https://www.w3.org/TR/MathML3/chapter4.html#contm.csymbol>

⁸ <https://www.w3.org/TR/MathML3/chapter4.html#contm.ci>

⁹ *definiens* is a phrase that defines an identifier or mathematical object. Considering equation (1), the correct definiens for W is the phrase '*Van der Waerden's number*'.

sions. They state their approach outperforms topic models that do not consider mathematics in text and report a topic coherence improvement of 0.012 over the LDA¹⁰ baseline. What all these embedding projects have in common is that they show promising examples and suppose a high potential, but do not evaluate their results for MIR.

Questions (i), (ii), and other pragmatic issues are already in discussion in a bigger context, as data production continues to rise and digital repositories seem to be the future for any archive structure. The National Research Council is making efforts to establish what they call the *Digital Mathematics Library* (DML)¹¹, a project under the International Mathematical Union. The goal of this future project is to take advantage of new technologies and help to solve the inability to search, relate, and aggregate information about mathematical expressions in documents over the web.

3 Machine Learning on Embeddings

The *word2vec* [18] technique computes real-valued vectors for words in a document using two main approaches: skip-gram and continuous bag-of-words (CBOW). Both produce a fixed length n -dimensional vector representation for each word in a corpus. In the skip-gram training model, one tries to predict the context of a given word, while CBOW predicts a target word given its context. In *word2vec*, context is defined as the adjacent neighboring words in a defined range, called a sliding window. The main idea is that the numerical vectors representing similar words should have close values if the words have similar context, often illustrated by the *king-queen* relationship¹².

Extending *word2vec*'s approaches, Le and Mikolov [13] propose *Paragraph Vectors* (PV), a framework that learns continuous distributed vector representations for any size of text segment (e.g. sentences, paragraphs, documents). This technique alleviates the inability of *word2vec* to embed documents as one single entity. This technique also comes in two distinct variations: Distributed Memory (DM) and Distributed Bag-of-Words (DBOW), which are analogous to the skip-gram and CBOW training models respectively. However, in both approaches, an extra feature vector representing the text segment, named paragraph-id, is included as another word. This paragraph-id is updated throughout the entire document, based on the current evaluated context window for each word, and is used to represent the whole text segment.

Recently, researchers have been trying to improve their semantic representations, producing multiple vectors (multi-sense embeddings) based on the word's sense, context, and distribution in the corpus [7, 28]. Another concern with traditional techniques is that they often neglect exploring lexical structures with valuable prior knowledge about the semantic relations, such as: WordNet [19], ConceptNet [15] and BabelNet [20]. These lexical structures offer a rich semantic

¹⁰ Latent Dirichlet Allocation

¹¹ <https://www.nap.edu/read/18619>

¹² $\mathbf{v}_{\text{king}} - \mathbf{v}_{\text{man}} \approx \mathbf{v}_{\text{queen}} - \mathbf{v}_{\text{woman}}$

environment that illustrate the word-senses, their use, and how they relate to each other. Some publications take advantage of the robustness provided by word embeddings approaches and lexical structures to combine them into multi-sense representations, improving their overall performance in many NLP downstream tasks [16, 29, 25].

The lack of solid references and applications that provide the same semantic structure of natural language for mathematical identifiers make their disambiguation process even more challenging. In natural texts, one can try to infer the most suitable word sense for a word based on the lemma¹³ itself, the adjacent words, dictionaries, thesaurus and so on. However, in the mathematical arena, the scarcity of resources and the flexibility of redefining their identifiers take this issue to a more delicate scenario. The context text preceding or following the mathematical equation is essential for its understanding.

More recently, [12] propose a variation of word embeddings for mathematical expressions. Their main idea relies on the construction of a distributed representation of equations, considering the word context vector of an observed word and its word-equation context window. They treat equations as single-unit words (EqEmb), which eventually appears in the context of different words. They also try to explore the effects of considering the elements of mathematical expressions separately (EqEmb-U). In this scenario, mathematical equations are represented using a Syntax Layout Tree (SLT) [37], which contains the spatial relationship between its symbols. While they present some interesting findings for retrieving entire equations, little is said about the vectors representing equation units and how they are described in their model. The word embedding techniques seem to have potential for semantic distance measures between complex mathematical expressions. However, they are not appropriate for extracting semantics of identifiers separately. This is an indication that the problems of representing mathematical identifiers are tied to more fundamental issues, which are explained in Section 5.

Since the overall performance of word embedding algorithms has shown superior results in many different NLP tasks, such as machine translation [18], relation similarity [9], word sense disambiguation [2], word similarity [21, 29], and topic categorization [26]. In the same direction, we also explore how well mathematical tokens can be embedded according to their semantic information. However, mathematical formulae are highly ambiguous and if not properly processed, their representation is jeopardized.

3.1 How to Embed Mathematics

There are two main standard formats in which to represent mathematics in science: \LaTeX and MathML. The former is used by humans for writing scientific documents. The latter, on the other hand, is popular in web representations of mathematics due to its machine readability and XML structure. There has been a major effort to automatically convert \LaTeX expressions to MathML [31]

¹³ canonical form, dictionary form, or citation form of a set of words

ones. However, neither \LaTeX nor MathML are practical formats for embeddings. Considering the equation embedding techniques in [12], we devise three main types of mathematical embeddings.

Mathematical Expressions as Single Tokens: EqEmb [12] uses entire mathematical expressions as one token. In this type, the inner structure of the mathematical expression is not taken into account. For example, Equation (1) is represented as one single token t_1 . Any other expression, such as $W(2, k)$ in the surrounding text of (1), is an entirely independent token t_2 . Therefore, this approach does not learn any connections between $W(2, k)$ and (1). While this approach seems to hold interesting results for comparing mathematical expressions, it fails in representing the semantic aspects of inner elements in mathematical equations.

Stream of Tokens: Instead of embedding mathematical expressions as a single token, we can represent them through a sequence of its inner tokens. For example, considering only the identifiers in Equation (1), we would have a stream of three tokens W , k , and ε . This approach has the advantage of learning all mathematical tokens. However, this method also has some drawbacks. Complex mathematical expressions may lead to long chains of elements, which can be especially problematic when the window size of the training model is too small. Naturally, there are approaches to reduce the length of chains. In Section 5 we show our own model which uses a stream of mathematical identifiers and cut out all other expressions. In [4], L. Gao et al. use a CBOW and embed all mathematical symbols, including identifiers and operands, such as $+$, $-$ or variations of equalities $=$. In [34], they do not cut out symbols and train their model on the entire sequence of tokens that the \LaTeX tokenizer generates. Considering Equation (1), it would result in a stream of 13 tokens. They use a long short-term memory (LSTM) architecture to handle longer chains of tokens and also to limit their length to 20 – 150 tokens. Usually, in word embeddings, such behaviour is not preferred since it increases the noise in the data¹⁴. We will see later in the paper (Section 3.1), that a typically trained model on mathematical embeddings is able to detect similarities between mathematical objects but do not perform well detecting connections to word descriptors. Therefore, we consider close relations of mathematical symbols to other mathematical symbols as noise. To mitigate this issue, we only work with mathematical identifiers and not any other symbols or structures.

Semantic Groups of Tokens: The third approach of embedding mathematics is only theoretical, and concerns the aforementioned problems related to the vague definitions of identifiers and functions in a standardized format (e.g. MathML). As previously discussed, current MIR and ML approaches would benefit from a basic structural knowledge of mathematical expressions, such that variations of function calls (e.g. $W(r, k)$ and $W(2, k)$) can be recognized as the same function. Instead of defining a unified standard, current techniques use their own ad-hoc interpretations of structural connections, e.g., α_i is one iden-

¹⁴ Noise means, the data consists of many uninteresting tokens that affect the trained model negatively.

tifier rather than three [31, 30]. We assume that an embedding technique would benefit from a system that is able to detect the parts of interest in mathematical expressions prior any training processes. However, such system still does not exist.

4 Performance of Math Embeddings

The examples illustrated in [4, 12, 34] seem to be feasible as a new approach for distance calculations between complex mathematical expressions. While comparing mathematical expressions is essentially practical for search engines or automatic plagiarism detection systems, these approaches do not seem to capture the components of complex structure separately, which are necessary for other applications, such as automated reasoning. Another aspect to be considered is that in [12] they do not train mathematical identifiers, preventing their system from learning connections between identifiers and definiens (e.g., $W(2, k)$ and the definiens ‘*Van der Waerden number*’). Additionally, the connection between entire equations and definiens is, at some level, questionable. Entire equations are rarely explicitly named¹⁵. However, in the extension EqEmb-U [12], they use an SLT representation to tokenize mathematical equations and to obtain specific unit-vectors, which is similar to our *identifiers as tokens* approach.

In order to investigate the discussed approaches, we apply variations of a word2vec implementation to extract mathematical relations from the arXMLiv 2018 [5] dataset, an HTML collection of the arXiv.org preprint archive¹⁶, which is used as our training corpus. We also consider the subsets that do not report errors during the document conversion (i.e. *no_problem* and *warning*) which represent 70% of archive.org. There are other approaches that also produce word embeddings given a training corpus as an input, such as fastText [1], ELMo [23], and GloVe [22]. The choice for word2vec is justified because of its implementation, general applicability, and robustness in several NLP tasks [9, 8, 14, 16, 25, 29]. Additionally, in fastText they propose to learn word representations as a sum of the n -grams of its constituent characters (sub-words). This would incorporate a certain noise to our experiments. In ELMo, they compute their word vectors as the average of their characters representations, which are obtained through a two-layer bidirectional language model (biLM). This would bring even more granularity than fastText, as they consider each character in a word as having their own n -dimensional vector representation. Another factor that prevent us from using ELMo, for now, is its expensive training process¹⁷. Closer to the word2vec technique, GloVe [22] is also considered, but its co-occurrence matrix would escalate the memory usage, making its training for arXiv not possible at the moment. We also examine the recently published Universal Sentence En-

¹⁵ However, it is common for groundbreaking findings, such as *Pythagorean’s theorem* or the *energy-mass equivalence*.

¹⁶ <https://arxiv.org/>

¹⁷ <https://github.com/allenai/bilm-tf>

coder (USE) [3] from google, but their implementation does not allow one to use a new training corpus, only to access its pre-calculated vectors.

As a pre-processing step, mathematical expressions are represented using MathML¹⁸ notation. Firstly, we replace all mathematical expressions by the sequence of the identifiers it contains, i.e., $W(2, k)$ is replaced by ‘ $W k$ ’. Secondly, we remove all common English stopwords from the training corpus. Finally, we train a word2vec model (skip-gram) using the following hyperparameters configuration¹⁹: vector size of 300 dimensions, a window size of 15, minimum word count of 10, and a negative sampling of $1E - 5$.

The trained model is able to partially incorporate semantics of mathematical identifiers. For instance, the closest²⁰ 27 vectors to the mathematical identifier f are mathematical identifiers themselves and the fourth closest noun vector to f is $\mathbf{v}_{\text{function}}$. Inspired by the classic *king-queen* example, we explore which tokens perform best to model a known relation. Consider an approximation $\mathbf{v}_{\text{variable}} - \mathbf{v}_a \approx \mathbf{v} - \mathbf{v}_f$, where $\mathbf{v}_{\text{variable}}$ represents the word *variable*, \mathbf{v}_a the identifier a , and \mathbf{v}_f represents f . We are looking for \mathbf{v} that fits best for the approximation. We call this measure the *semantic distance* to f with respect to a given relation between two vectors. Table 1 shows the top 10 semantically closest results to f with respect to the relation between \mathbf{v}_a and $\mathbf{v}_{\text{variable}}$.

We also perform an extensive evaluation on the first 100 entries²¹ of the *MathMLBen* benchmark [31]. We evaluate the average of the *semantic distances* with respect to the relations between $\mathbf{v}_{\text{variable}}$ and \mathbf{v}_x , $\mathbf{v}_{\text{variable}}$ and \mathbf{v}_a , and $\mathbf{v}_{\text{function}}$ and \mathbf{v}_f . In addition, we consider only results with a cosine similarity of 0.7 or above to maintain a minimum quality in our results. The overall results were poor with a precision of $p = .0023$ and a recall of $r = .052$. For the identifier W (Equation (1)), the evaluation presents four semantically close results: *functions*, *variables*, *form*, and the mathematical identifier q . Even though expected, the scale of the presented results are astonishing.

Additionally, we also try the Distributed Bag-of-Words of Paragraph Vectors (DBOW-PV) [13] considering the approach of [30]. In [30], they analyze all occurrences of mathematical identifiers and consider the entire article at once. We assume this prevents the algorithm

Tokens	Cosine Distances
variables	0.7600
function	0.7154
appropriate	0.6925
independent	0.6789
instead	0.6784
defined	0.6729
namely	0.6719
continuous	0.6707
depends	0.6629
represents	0.6623

Table 1. Semantically closest 10 results to f with respect to the relation between \mathbf{v}_a and $\mathbf{v}_{\text{variable}}$.

¹⁸ The source T_EX file has to use mathematical environments for its expressions.

¹⁹ Non mentioned hyperparameters are used with their default values as described in the Gensim API [27]

²⁰ Considering cosine similarity.

²¹ Same entries used in [30]

from finding the right descriptor in the text, since later or prior occurrences of an identifier might appear in a different context, and therefore potentially introduce different meanings. Instead of using the entire document, we consider the algorithm of [30] only in the input paragraph and similar paragraphs given by our DBOW-PV model. Unfortunately, the obtained variance within the paragraphs brings a high number of false positives to the list of candidates, which affects negatively our performance.

We also experiment other hyperparameters when training our word embeddings model to see if it is possible to improve the overall results. However, while the performance decreases, no drastic structural changes appear in the model. Figure 1 illustrates a t-SNE plot²² of the model trained with 400 dimensions, a window size of 25, and minimum count of 10 words, without any filters applied to the text. The plot is similar to the visualized model presented in [4], even though they use a different embedding technique. Compared to [4], we provide a bigger picture of the model that reveals some dense clusters for numbers at $(11, -1)$ with the math token for *invisible times* nearby, equation abbreviations such as, *eq1*, at $(-8, -9)$, and logical operators at $(-9, -7)$. We highlight mathematical tokens in the model in red and word tokens in blue. The plot in Figure 1 illustrates that mathematical tokens are close to each other.

Based on the presented results, one can still argue that more settings should be explored (e.g. different embedding techniques, parameters) for the embeddings phase and different pre-processing steps (e.g. stemming and lemmatization) be adopted. This would probably solve some minor problems, such as removing the inaccurate first hit in Table 1. Nevertheless, the overall results would not be improved to a point of being comparable to [30] findings, which report a precision of $p = 0.48$. The main reason for this is that, mathematics as a language is highly customizable. Many of the defined relations between mathematical concepts and their descriptors are only valid in a local scope. Consider, for example, an author that notes his algorithm by π . This does not change the general meaning of π , even though it effects the meaning in the scope of the article. Current ML approaches only learn patterns of most frequently used combinations, e.g., between f and *function*, as seen in Table 1. Furthermore, we assume this is a general problem that different embedding techniques and tweaks of settings, such as those illustrated in Section 2, would not solve. Therefore, in the following section, we present some concepts that we believe can help ML algorithms to better understand mathematics.

5 Make Math Machine Learnable

A case study [33] has shown that 70% of mathematical symbols are explicitly declared in the context. Only four reasons are causing an explicit declaration in

²² Note that t-SNE plots may misleadingly create clusters that do not exist in the model. To overcome this issue we create several plots with different settings. The results remain similar to the plot we show which is an indication that the visual clusters exist also in the model.

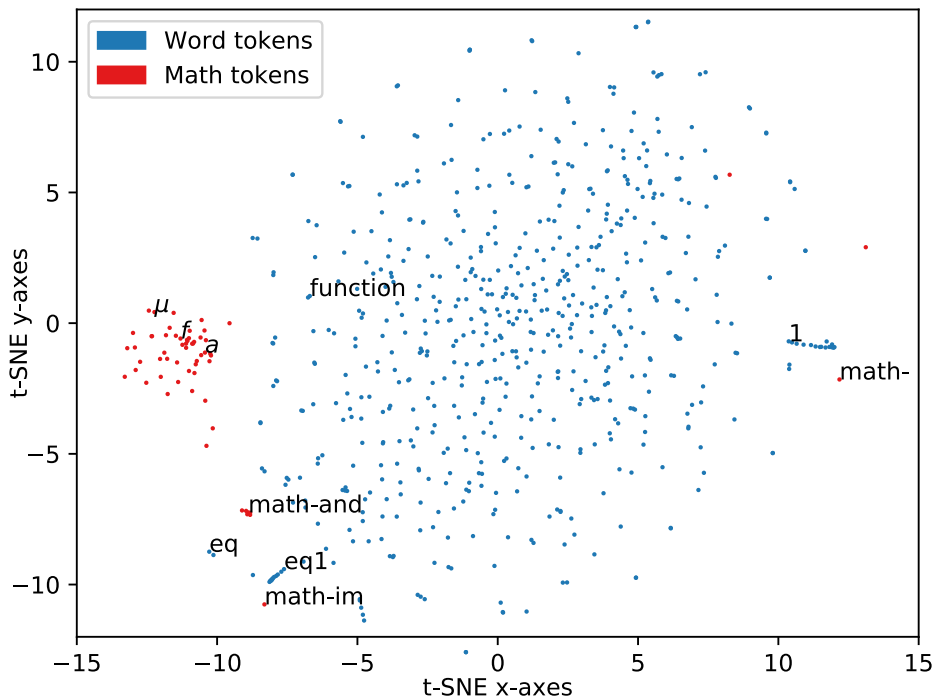


Fig. 1. t-SNE plot of top 1000 closest vectors of the identifier f . For this plot, we used a perplexity of 80 and the default values of the TSNE python package for other settings.

the context: (a) a new mathematical symbol is defined, (b) a known notation is changed, (c) used symbols are present in other contexts and require specifications to be properly interpreted, or (d) authors declarations were redundant (e.g. for improving readability). We assume (d) is a rare scenario compared to (a-c), unless in educational literature. Current math-embedding techniques can learn semantic connections only in those 70%, where the definiens are available. Besides (d), the algorithm would learn either rare notations (in case of (a)) or ambiguous notations (in cases (b-c)). The flexibility that mathematical documents allow to (re)define used mathematical notations further corroborates to the complexity of learning mathematics.

How would be possible for machines to learn mathematics? One of the major problems is the ambiguity of mathematical expressions (a-c). Natural languages also consist of ambiguous and context-sensitive words. A typical workaround for this problem is to consider the use of lexical databases (e.g. WordNet [19]) to identify the most suitable word senses for a given word. However, mathematics lacks of such system, which makes its adoption not feasible at the moment. In [35] they propose the use of tags, similarly to the POS tags in linguistics, but for tagging mathematical \TeX tokens, bringing more information to the tokens considered. As a result, a lexicon containing several meanings for a large set of

mathematical symbols is developed. Such lexicons might enable the disambiguation approaches in linguistics to be used in mathematical embeddings in the near future.

Furthermore, learning algorithms would benefit from a literature focused in (a) and (d), instead of (b-c). Similar to students who start to learn mathematics, ML algorithms have to consider the structure of the content they learn. It is hard to learn mathematics only considering arXiv documents without prior or complementary knowledge. Usually, these documents represent state-of-the-art findings containing new and unusual notations and lack of extensive explanations (e.g. due to page limitations). In contrast, educational books carefully and extensively explain new concepts. We assume better results can be obtained if ML algorithms are to be trained in multiple stages. A basic model trained on educational literature should capture standard relations between mathematical concepts and descriptors. This model should also be able to capture patterns independently how new or unusual the notations are present in the literature. In 2014, Matsuzaki et al. [17] present some promising results to automatically answer mathematical questions from Japanese university entrance exams. While the approach involves many manual adjustments and analysis, the promising results illustrate the different levels of knowledge that is still required for understanding arXiv documents and university entrance level exams. A well-structured digital mathematical library that distinguishes the different levels of progress in articles (e.g. introductions vs. state-of-the-art publications) would also benefit mathematical machine learning tasks.

Another problem in recent publications, is the lack of standards for properly evaluating MIR algorithms, leading to several publications that present promising results without an extensive evaluation [4, 12, 34]. While ML algorithms in NLP benefit from available extensive training and testing datasets, ongoing discussions about interpretations of mathematical expressions [31], and imprecise standards thwarts research progress in MIR. A common standard for interpreting semantic structures of mathematics would help to overcome the issues of different evaluation techniques. Numerous applications (e.g., search engines) would benefit directly from unified interpretations and representations of mathematical semantics. Therefore, we introduce *Mathematical Objects of Interest* (MOI). Currently, there are three approaches to tokenize and interpret semantics of mathematical expressions: (1) tokenize the mathematical \TeX string and tag the tokens with semantic information [35], (2) analyze the elements of presentational MathML [37], and (3) analyze elements of content MathML [30, 31]. The Part-of-Math (POM) tagger [35] proposes a multi-scan approach for mathematical \TeX strings that incorporates more semantic information into a parse tree in each iteration. The available first scan creates a parse tree and tags each node with further information about the symbol of each node (similar to POS-tags in linguistics). In [37], they generate custom tree representations of presentational MathML (SLT) that allow wildcards for certain positions in the trees and improves search engine results. Applications of (3) are discussed in Section 2.

The goal of MOIs is to combine the advantages of concepts (1-3) and propose a unified solution for interpreting mathematical expressions. We suggest MOIs as a recursive tree structure in which each node is an MOI itself. The current workaround of the problematic example of α_i as an element of the vector α in content MathML is vague and inappropriate for content specific tasks. As an MOI, this expression would contain three nodes, with α_i as the parent node of two leaves α and i . While it first seems non-intuitive that α , as the vector, is a child node of its own element, this structure is able to incorporate all three components of semantic information of the expression. Hence, an MOI structure should not be misinterpreted as a logical network explaining semantic connections between its elements, but as a highly flexible and lightweight structure for incorporating semantic information of mathematical expressions.

We believe that a new standard, such as MOIs, has to be extensively studied and discussed by specialists on the field before its acceptance. Therefore, the introduction of MOIs in this paper is fairly broad and still need to be further investigated. However, we suggest the use of MOIs would assist ML algorithms to simplify and unify their evaluations in MathIR projects.

6 Conclusion and future Work

In this paper, we explore how text embedding techniques (e.g. word2vec, PV) are unable to represent mathematical expressions adequately. After experimenting with popular mathematical representations in MIR, we expose fundamental problems that prevent ML algorithms from learning mathematics. We also discover the same problems in several related research projects. Many of these projects show promising examples without extensive evaluations, motivating more researchers to follow the same idea with equivalent fundamental issues.

We present concepts for enabling ML algorithms to learn mathematical expressions. Some of these concepts are generally time consuming, such as a lexical database for mathematics. For a concrete contribution we propose the MOIs, a unified solution for interpreting mathematical expressions semantically. We think this is a necessary first step to enable unified evaluations and libraries.

In future work, we plan to explore the element distributions of mathematical expressions and compare them with known distributions in linguistics [24] for leveraging MOIs'. Preliminary results have shown that distributions of mathematical objects also following Zipf's law, similar to word distributions in natural language. Thus, we are currently analyzing term frequencies and inverse document frequencies, commonly known as tf-idf, of mathematics. This research should help to discover existing meaningful mathematical structures that are already in use in scientific publications. Such structures should be interpreted as MOIs. Therefore, this research will build a base to constructively discuss MOIs.

Acknowledgments This work was supported by the German Research Foundation (DFG grant GI-1259-1). We thank Howard Cohl who provided insights and expertise.

References

- [1] P. Bojanowski et al. “Enriching Word Vectors with Subword Information”. In: *Transactions of the Association for Computational Linguistics* 5 (2017).
- [2] J. Camacho-Collados et al. “A Unified Multilingual Semantic Representation of Concepts”. In: *Proc. 53rd Annual Meeting of the Association for Computational Linguistics (ACL), Beijing, China*. ACL, 2015.
- [3] D. Cer et al. “Universal Sentence Encoder for English”. In: *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Ed. by E. Blanco and W. Lu. Association for Computational Linguistics, 2018.
- [4] L. Gao et al. “Preliminary Exploration of Formula Embedding for Mathematical Information Retrieval: can mathematical formulae be embedded like a natural language?” In: *CoRR* abs/1707.05154 (2017). arXiv: 1707.05154.
- [5] D. Ginev. *arXMLiv:08.2018 dataset, an HTML5 conversion of arXiv.org*. SIGMathLing – Special Interest Group on Math Linguistics. 2018. URL: <https://sigmathling.kwarc.info/resources/arxmliv/>.
- [7] E. H. Huang et al. “Improving Word Representations via Global Context and Multiple Word Prototypes”. In: *Proc. 50th Annual Meeting of the Association for Computational Linguistics (ACL) Vol. 1*. Jeju Island, Korea: ACL, 2012.
- [8] I. Iacobacci et al. “Embeddings for Word Sense Disambiguation: An Evaluation Study”. In: *Proc. 54th Annual Meeting of the Association for Computational Linguistics (ACL) Vol. 1, Berlin, Germany*. ACL, 2016.
- [9] I. Iacobacci et al. “SensEmbed: Learning Sense Embeddings for Word and Relational Similarity”. In: *Proc. 53rd Annual Meeting of the Association for Computational Linguistics (ACL) Vol. 1, Beijing, China*. ACL, 2015.
- [10] M. Kohlhase. “Math Object Identifiers - Towards Research Data in Mathematics”. In: *Lernen, Wissen, Daten, Analysen (LWDA) Conference Proceedings, Rostock, Germany, September 11-13, 2017*. Ed. by M. Leyer. Vol. 1917. CEUR-WS.org, 2017.
- [11] G. Y. Kristianto et al. “Extracting Textual Descriptions of Mathematical Expressions in Scientific Papers”. In: *D-Lib Magazine* 20.11/12 (2014).
- [12] K. Krstovski and D. M. Blei. “Equation Embeddings”. In: *CoRR* abs/1803.09123 (2018). arXiv: 1803.09123.
- [13] Q. Le and T. Mikolov. “Distributed Representations of Sentences and Documents”. In: *Proceedings of the 31th International Conference on Machine Learning, ICML, Beijing, China*. Beijing, China: JMLR.org, 2014.
- [14] J. Li and D. Jurafsky. “Do Multi-Sense Embeddings Improve Natural Language Understanding?” In: *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP), Lisbon, Portugal*. Lisbon, Portugal: Association for Computational Linguistics, 2015.
- [15] H. Liu and P. Singh. “ConceptNet &Mdash; A Practical Commonsense Reasoning Tool-Kit”. In: *BT Technology Journal* 22.4 (Oct. 2004).

- [16] M. Mancini et al. “Embedding Words and Senses Together via Joint Knowledge-Enhanced Training”. In: *Proc. 21st Conference on Computational Natural Language Learning (CoNLL), Vancouver, Canada*. Association for Computational Linguistics, 2017.
- [17] T. Matsuzaki et al. “The Most Uncreative Examinee: A First Step toward Wide Coverage Natural Language Math Problem Solving”. In: *Proc. Twenty-Eighth AAAI Conference on Artificial Intelligence, Québec City, Québec, Canada*. Ed. by C. E. Brodley and P. Stone. AAAI Press, 2014.
- [18] T. Mikolov et al. “Distributed Representations of Words and Phrases and Their Compositionality”. In: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*. Lake Tahoe, Nevada: Curran Associates Inc., 2013.
- [19] G. A. Miller. “WordNet: A Lexical Database for English”. In: *Commun. ACM* 38.11 (1995).
- [20] R. Navigli and S. P. Ponzetto. “BabelNet: The Automatic Construction, Evaluation and Application of a Wide-Coverage Multilingual Semantic Network”. In: *Artificial Intelligence* 193 (2012).
- [21] A. Neelakantan et al. “Efficient Non-parametric Estimation of Multiple Embeddings per Word in Vector Space”. In: *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar*. Association for Computational Linguistics, 2014.
- [22] J. Pennington et al. “Glove: Global Vectors for Word Representation.” In: *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar*. Vol. 14. Association for Computational Linguistics, 2014.
- [23] M. Peters et al. “Deep Contextualized Word Representations”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics, 2018.
- [24] S. T. Piantadosi. “Zipf’s word frequency law in natural language: A critical review and future directions”. In: *Psychonomic Bulletin & Review* 21.5 (Mar. 2014).
- [25] M. T. Pilehvar and N. Collier. “De-Conflated Semantic Representations”. In: *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)*. The Association for Computational Linguistics, 2016.
- [26] M. T. Pilehvar et al. “Towards a Seamless Integration of Word Senses into Downstream NLP Applications”. In: *CoRR* abs/1710.06632 (2017).
- [27] R. Řehůřek and P. Sojka. “Software Framework for Topic Modelling with Large Corpora”. English. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. <http://is.muni.cz/publication/884893/en>. Valletta, Malta: ELRA, May 2010.
- [28] J. Reisinger and R. J. Mooney. “Multi-prototype Vector-space Models of Word Meaning”. In: *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Com-*

- putational Linguistics*. Los Angeles, California: Association for Computational Linguistics, 2010.
- [29] T. Ruas et al. “Multi-sense Embeddings through a Word Sense Disambiguation Process”. Pre-print.
 - [30] M. Schubotz et al. “Evaluating and Improving the Extraction of Mathematical Identifier Definitions”. In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction - 8th International Conference of the CLEF Association, CLEF 2017, Dublin, Ireland, September 11-14, 2017, Proceedings*. Ed. by G. J. F. Jones et al. Vol. 10456. Springer, 2017.
 - [31] M. Schubotz et al. “Improving the Representation and Conversion of Mathematical Formulae by Considering their Textual Context”. In: *Proceedings of the 18th ACM/IEEE on Joint Conference on Digital Libraries, JCDL 2018, Fort Worth, TX, USA, June 03-07, 2018*. Ed. by J. Chen et al. ACM, 2018.
 - [32] M. Schubotz et al. “Semantification of Identifiers in Mathematics for Better Math Information Retrieval”. In: *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval*. Pisa, Italy: ACM, 2016.
 - [33] M. Wolska and M. Grigore. “Symbol Declarations in Mathematical Writing”. In: *Towards a Digital Mathematics Library*. Ed. by P. Sojka. Paris, France: Masaryk University Press, 2010.
 - [34] M. Yasunaga and J. Lafferty. “TopicEq: A Joint Topic and Mathematical Equation Model for Scientific Texts”. In: *CoRR* abs/1902.06034 (2019). arXiv: 1902.06034.
 - [35] A. Youssef. “Part-of-Math Tagging and Applications”. In: *Proc. CICM*. Ed. by H. Geuvers et al. Cham: Springer International Publishing, 2017.
 - [36] A. Youssef and B. R. Miller. “Deep Learning for Math Knowledge Processing”. In: *Proc. CICM*. Ed. by F. Rabe et al. Vol. 11006. Springer, 2018.
 - [37] R. Zanibbi et al. “Multi-Stage Math Formula Search: Using Appearance-Based Similarity Metrics at Scale”. In: *Proc. 39th International ACM SIGIR Conference on Research and Development in Information Retrieval, Pisa, Italy*. Ed. by R. Perego et al. ACM, 2016.

```
1 @inproceedings{GreinerPetter2019,  
2   author    = {Greiner-Petter, Andr\' { e } and Ruas, Terry  
3             and Schubotz, Moritz and Aizawa, Akiko and Grosky,  
4             William and Gipp, Bela},  
5   location  = {Paris, France},  
6   booktitle = {Submitted to 4th Joint Workshop on  
7             Bibliometric-enhanced Information Retrieval and Natural  
             Language Processing for Digital Libraries colocated at  
             the 42nd International ACM SIGIR Conference},  
8   date      = {2019-07},  
9   title     = {Why Machines Cannot Learn Mathematics, Yet},  
10  }
```

Listing 1.1. Use the following BibTeX code to cite this article