

Understanding Scientific Documents with Synthetic Analysis on Mathematical Expressions and Natural Language

Takuto Asakura

Department of Informatics, SOKENDAI

1 Introduction

Converting Science, Technology, and Engineering (STEM) documents to formal expressions has a large impact on academic and industrial society. It enables us to construct databases of mathematical knowledge, search for formulae, and develop a system that generates executable codes automatically.

However, the conversion is an exceedingly ambitious goal. Mathematical expressions are commonly used in scientific communication in numerous fields such as mathematics and physics, and in many cases, they express key ideas in STEM documents. Despite the importance of mathematical expressions, formulae and texts are complementary to each other, and those in documents cannot be understood independently. Thus, deep synthetic analyses on natural language and mathematical expressions are necessary.

To date, a large number of efforts have been made for developing Natural Language Processing (NLP) techniques, including semantic parsing [4], but their targets are mostly ‘general’ texts. Naturally, conventional NLP techniques include only limited features to treat formulae and numerous linguistic phenomena specific to STEM documents [3].

Meanwhile, semantics on mathematical expressions also has been deeply investigated. Such results can be seen in logic theories, MathML specification [1], etc. However, there is a large space between formal expressions such as first-order logic and actual formulae in natural language texts.

2 Research Goals

There are a number of remaining works to achieve the conversion from STEM documents to a computational form (Figure 1). At first, we are going to focus on the two foundational parts for the synthetic analyses. The first one is token-level analyses on formulae. The main part of the analyses is associating formulae tokens to mathematical objects and text fragments (Section 2.1). This is a primal step for the conversion, but it is still almost untouched. The second one is the morphology of mathematical expression and semantics covering both formulae and texts (Section 2.2). Studying underlying theories is essential to deeply understand the structure of STEM document, and aim for the practical application by a bottom-up approach.

2.1 Associating Tokens in Formulae with Mathematical Objects and Their Descriptions in Texts

Tokens in formulae (e.g., x , ε , \times , \log) and their combination can refer to *mathematical objects*. We human beings are able to detect what each token or combination pointing to, by using common sense, domain knowledge, and referencing descriptions in the document or in the others. This detection is fundamental and should be one of the initial steps for understanding STEM documents, but unfortunately, it cannot be easily done by a machine. There are at least four factors which make the detection highly challenging: (1) ambiguity of tokens, (2) syntactic ambiguity of formulae, (3) necessity for common sense and domain knowledge, and (4) severe abbreviation. These difficulties often appear in formulae; giving an example for (1) as a representative, only in the first chapter of a book *Pattern Recognition and Machine Learning* (PRML) [2], a character **y** (letter ‘y’ in bold roman) is used in several meanings including a function, vectors, and a value (Table 1).

Copyright © by the paper’s authors. Copying permitted for private and academic purposes.

In: C. Kaliszyk, E. Brady, J. Davenport, W.M. Farmer, A. Kohlhase, M. Kohlhase, D. Müller, K. Pąk, and C. Sacerdoti Coen (eds.): Joint Proceedings of the FMM, LML, OpenMath Workshops, Doctoral Program and Work in Progress at the Conference on Intelligent Computer Mathematics 2019 co-located with the 12th Conference on Intelligent Computer Mathematics (CICM 2019), Prague, Czech Republic, July 8–12, 2019, published at <http://ceur-ws.org>

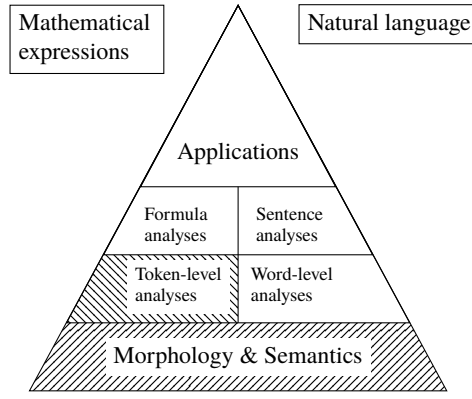


Figure 1: Overview of our task definitions. At first, we are tackling the token-level analyses on mathematical expressions (Section 2.1) and theories covering both formulae and texts (Section 2.2).

Table 1: Usage of character \mathbf{y} in the first chapter of PRML (except exercises). Underlines by the author.

Text fragment from PRML Chap. 1	Meaning of \mathbf{y}
... can be expressed as a function $\underline{\mathbf{y}}(\mathbf{x})$ which takes ...	a function which takes an image as input
... an output vector $\underline{\mathbf{y}}$, encoded in ...	an output vector of function $\mathbf{y}(\mathbf{x})$
... two vectors of random variables \mathbf{x} and $\underline{\mathbf{y}}$...	a vector of random variables
Suppose we have a joint distribution $p(\mathbf{x}, \underline{\mathbf{y}})$ from ...	a part of pairs of values, corresponding to \mathbf{x}

The other part of the initial steps of understanding STEM documents is that connecting text fragments to the subjective mathematical objects. Our hypothesis is that for this step, general NLP approaches such as dependency parsing are more or less applicable. Of course, some tuning for STEM documents will be required, and also this process might need to be done interactively with the mathematical object detection for formulae.

2.2 Semantics and Morphology

Semantics on natural language and mathematical expressions have been studied separately. However, to understand STEM document, it is important to investigate a synthetic semantics covering both of texts and formulae.

Though morphology has been studied for natural languages, not so much for formulae. As a matter of fact, in terms of morphology, *words* also exist in formulae. For instance, a token M is a word in “Matrix M ”, but M is not a word in “An entry $M_{i,j}$ ” ($M_{i,j}$ is a word). Unlike morphemes in natural language, tokens in formulae do not have lexical categories, but some symbols (e.g., parentheses and equal sign) and positional information (e.g., superscript and subscript) have typical usages.

3 Completed and Remaining Research

For the beginning of our research, we simplified the detection task which we described in Section 2.1. Specifically, we are giving annotations on some research papers in the following manner:

1. Detecting minimal groups of tokens (we call them *chunks*) each of which refers to a mathematical object (chunking).
2. Categorizing chunks by the mathematical object they referring to.

This annotation (we call it *pilot annotation*) is the fundamental process to create the first gold dataset for associating tokens and mathematical objects. The annotated data will also be helpful for investigating the morphology on mathematical expressions.

In other words, we defined a classification task before annotating descriptions for tokens. Since there are many ways to describe a mathematical object, this classification can be done more coherently through the pilot annotation. Moreover, we are expecting that the classification is naturally rather easier to be automated than giving descriptions automatically for the first attempt.

Besides the pilot annotation, all the works which have to be done to achieve our goal are remaining. For the next step, we are planning to automate the annotation process by using features such as apposition nouns and syntactic information in formulae. At the same time, we have to decide the form of mathematical objects. For now, we can say that every

mathematical object should have a description and some attributes such as types (e.g., int and float). What attributes are necessary and sufficient is still not clear, and we will find it out after trying the annotation for several documents.

4 Publication Plans and Evaluation Plans

Currently, we are creating a new language resource as the pilot annotation, and we are planning to publish it for the community of language resources. For the further future, we will develop automation algorithms for mathematical object detection, which are works suitable for NLP and digital mathematical library community, including CICM. The analyses on underlying morphology and semantics are more like works in computational linguistics.

For the initiative dataset, it is better to make agreements among a few experts if possible. Following progress on developing algorithms and analyses on linguistic phenomena should be evaluated with our handmade gold datasets.

References

- [1] Ron Ausbrooks et al. *Mathematical Markup Language (MathML) 3.0 Specification*. World Wide Web Consortium (W3C), 2014.
- [2] Christopher M Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- [3] Michael Kohlhase and Mihnea Iancu. “Co-Representing Structure and Meaning of Mathematical Documents”. In: *Sprache und Datenverarbeitung, International Journal for Language Data Processing* 38.2 (2014).
- [4] Siva Reddy et al. “Transforming Dependency Structures to Logical Forms for Semantic Parsing”. In: *Transactions of the Association for Computational Linguistics* (2016).