

The Structure of Mathematical Expressions

An ARXIV Case Study

Deyan Ginev and Bruce R. Miller

National Institute of Standards and Technology

March 30, 2012



Contents

Contents	2
1 Introduction	3
1.1 Motivation	3
1.2 Related Resources	3
2 Methods	5
2.1 Training Corpus	5
2.2 Structural Annotation	5
2.3 Annotation Vocabulary	5
3 A Study of Mathematical Syntax	9
3.1 Basics	9
3.2 Discrete math	9
3.3 Continuous math	11
3.4 Other fields	12
4 Discussion	13
5 Conclusion	15

Introduction

In this study, we survey the notational diversity of present-day mathematical expressions, in order to uncover their linguistic phenomena. A practical motivation for this study is to provide a foundation for determining the boundary between syntactic and semantic phenomena in said expressions, from the perspective of language modeling. The ultimate goal of this project is to construct a grammar of mathematical expressions, which captures all relevant syntactic properties established in this study, and allows for the semantic analysis necessary to model and observe the semantic relationships.

1.1 Motivation

We want to enable machine-reading of formulas, in order to provide a variety of user-assistance services, such as semantic search, text-to-speech synthesis, semantic interactions (definition lookup), as well as computer algebra support (“evaluate subexpressions on demand”) and ultimately computer verification (“does that proof step really hold?”).¹

EdN:1

1.2 Related Resources

Notation census, beginnings of study are in Deyan’s thesis, Naproche and FMathL have examples, but no real systematic study.²

EdN:2

¹EdNOTE: expand

²EdNOTE: expand

Methods

2.1 Training Corpus

The primary corpus on which we base this investigation is the Cornell pre-print archive “ARXIV”³, consisting of over 700,000 articles in 37 scientific subfields.

EdN:3

arXiv Sandbox

4

EdN:4

As a secondary resource, we we will also consult entry-level literature on highschool mathematics, in order to exhibit basic phenomena, as well as to demonstrate phenomena apriori known to the authors.⁵

EdN:5

2.2 Structural Annotation

As one of the goals of our study is to establish a first guess of an underspecified operator tree⁶, any annotation must at its core mark up the applicative logical structure of the mathematical expression. This process will build up a formula tree, the collection of which can later be used as a gold standard for developing a grammatical model of the language of symbolic mathematics.

EdN:6

7 8

EdN:7

EdN:8

2.3 Annotation Vocabulary

Another core goal is to discover and describe interesting linguistic phenomena that occur naturally in our corpus. Examples of what we consider “interesting” are phenomena that

³EdNOTE: cite here

⁴EdNOTE: Say that, on the ARXIV front, we first start with the train sandbox from Deyan’s thesis

⁵EdNOTE: Wikipedia? PEMDAS?

⁶EdNOTE: make sure the concepts are introduced and/or rephrase

⁷EdNOTE: I’m currently thinking of rendering the annotations as trees (tikz,pstricks...custom tree drawing package?), so that the annotator can proofread the annotations in an intuitive manner.

⁸EdNOTE: In the XHTML, I’m thinking of ContentMML+SVG rendering, all of this figured out by the binding, maybe a custom stylesheet?

Train1	Differential Geometry http://arxmliv.kwarc.info/files/9609/dg-ga.9609012
Train2	Quantum Physics http://arxmliv.kwarc.info/files/0910/0910.5733/
Train3	High Energy Physics - Theory http://arxmliv.kwarc.info/files/9407/hep-th.9407125/
Train4	Commutative Algebra http://arxmliv.kwarc.info/files/0809/0809.4873/
Train5	Statistics Theory http://arxmliv.kwarc.info/files/0905/0905.1486/
Train6	General Relativity and Quantum Cosmology http://arxmliv.kwarc.info/files/0807/0807.2507/
Train7	Cosmology and Extragalactic Astrophysics http://arxmliv.kwarc.info/files/0908/0908.2548
Train8	Exactly Solvable and Integrable Systems http://arxmliv.kwarc.info/files/0905/0905.2033
Train9	Geometric Topology http://arxmliv.kwarc.info/files/0809/0809.4477
Train10	Algebraic Geometry http://arxmliv.kwarc.info/files/0704/0704.0537

Table 2.1: Sandbox of Ten Random ARXIV Papers from Diverse Scientific Subfields

induce ambiguity, or legitimize what would typically be ungrammatical fragments. Cases of ambiguity are well-known to follow from semantic overloading of symbols, implicit argument scopes of operations or eliding syntax, leaving the reader with the task of guessing the “invisible” dynamics. Use of custom shorthands, however, as well as custom notations in general, expands the grammar of symbolic mathematics, often in completely non-standard ways that can only be grasped through a deep understanding of the document at hand.

As multiple interesting observations can be made for a single large mathematical formula, it is natural to annotate multiple relevant subexpressions. More concretely, for each phenomenon of interest, we annotate the greatest common subtree (GCT) of all participating subtrees. In case we find a long-range relationship in a large formula, the annotation would hence be placed on the formula root.

The annotations can be utilized for different purposes - browsing by specific phenomena, syntactic feature or lemma, training a classifier, etc. Thus, we take a compositional, standardized approach to providing labels from a fixed vocabulary for the relevant ontological classes of structural properties.

EdN:9

9

⁹EDNOTE: Additional tokens: super, sub, fenced

Property	Keywords
Fixity	over, under, prefix, infix, postfix, superfix, subfix, circumfix, transfix, nofix ¹
Role (Symbols)	separator, modifier, relation, operator, metarelation, binder
Role (Objects)	factor, term, statement, variable, constant, modified
Role (Structure)	tuple, sequence, expression, shorthand, template, language
Composition	invisible, atom, complex, chained
Shallow Semantics	type, function, constructor, other
Linguistic	ellipsis, metonymy, ambiguity, vagueness, anaphora
Math Practices	framing

Table 2.2: Keyword Vocabulary for Syntactic Properties

Chapter 3

A Study of Mathematical Syntax

3.1 Basics

Foundations

10 11 12

EdN:10

EdN:11

EdN:12

High School

13 14

EdN:13

EdN:14

3.2 Discrete math

Set Theoretic Notations

15 16

EdN:15

EdN:16

Logical Operators

17

EdN:17

Combinatorics

18 19

EdN:18

EdN:19

¹⁰EdNOTE: arithmetic, grouping fences and equality

¹¹EdNOTE: basic relations and orderings

¹²EdNOTE: arithmetic and algebraic sequences?

¹³EdNOTE: geometry here, otherwise a separate geometry subsection

¹⁴EdNOTE: trigonometry, complex and rational numbers

¹⁵EdNOTE: elementhood, inclusions, set constructors, overloaded arith ops

¹⁶EdNOTE: also maps : domains -> codomains, xRy notations

¹⁷EdNOTE: classic logic, HOL, type theories

¹⁸EdNOTE: Infinite sums

¹⁹EdNOTE: binomials, combinations, permutations,

Expression	Denotation	Annotation
1. $W \in \mathcal{P} \cap \mathcal{Z}$	set membership	
2. $\nu : \times^n \mathbb{V} \rightarrow \mathbb{R}$	a map	
3. $\mathcal{Z}^* = \{X \in \mathcal{V} \mid \omega(X, W) \in \mathbb{Z}, \text{ for all } W \in \mathcal{Z}\}$	definition to set	
4. $\text{span}_{\mathbb{R}}\{W_1, \dots, W_g\}$	span of a set	
Discussion: set operators can take fenced yet not simply <i>grouped</i> arguments, [Train1]		

Table 3.1: Set Theory Notations, Part 1

Number Theory

20 21 22 23

EdN:20

EdN:21

Graph Theory

EdN:22

EdN:23

24 25 26

EdN:24

EdN:25

EdN:26

Algebra

27 28 29 30

EdN:27

EdN:28

Functions Theory

EdN:29

EdN:30

31

EdN:31

3.3 Continuous math

Calculus

32

EdN:32

Probability

33 34

EdN:33

EdN:34

Interval Notation and Arithmetic

35

EdN:35

Topology

36

EdN:36

²⁰EdNOTE: modulo modifiers

²¹EdNOTE: tuples

²²EdNOTE: divisibility notations $a \mid b$ and b/a

²³EdNOTE: DLMF sneaky notations

²⁴EdNOTE: edge and vertex notations

²⁵EdNOTE: incidence and adjacency notations

²⁶EdNOTE: Wiki is very nice: http://en.wikipedia.org/wiki/Glossary_of_graph_theory

²⁷EdNOTE: vectors

²⁸EdNOTE: maps and complements

²⁹EdNOTE: groups

³⁰EdNOTE: lattices

³¹EdNOTE: talk about associativity of application and composition, “;” and “o” as notation variants, discuss complex examples

³²EdNOTE: differentials, integrals, limits, remember brownian motion integral notations!

³³EdNOTE: Bayes formula with multiple denotations of P

³⁴EdNOTE: Various conditional and joint probability notations

³⁵EdNOTE: introduce interval notations, then move to interval arithmetic

³⁶EdNOTE: manifold constructors and notations

3.4 Other fields

Quantum Physics

EdN:37 ³⁷ ³⁸ ∴
EdN:38

³⁷EDNOTE: Bra-ket notation

³⁸EDNOTE: computer science, biology, chemistry...

Chapter 4

Discussion

Conclusion
