

LaSer : search engine for mathematical formulae

Prithwish Mukherje, Biswajit Paria, Agnivo Saha, Haritabh Singh, Anurag
Anand, Sandesh C, and Mayank Singh*

Department of Computer Science, IIT, Kharagpur

E-mail: mayank4490@gmail.com

Abstract

We present a search engine for mathematical formulae. The MathWebSearch system harvests the web for content representations (currently MathML) of formulae and indexes them using unigrams, bigrams and trigrams. The query for now is limited to Latex format only and can and will be extended to free form queries.

Introduction

As the world of information technology grows, being able to quickly search data of interest becomes one of the most important tasks in any kind of environment, be it academic or not. This paper addresses the problem of searching mathematical formulae from a semantic point of view, i.e. to search for mathematical formulae not via their presentation but their structure and meaning.

*To whom correspondence should be addressed

Related Work

We have seen that early work has been done in this field, allowing users to search for mathematical equations. There are multiple sites available such as <http://latexsearch.com/>, <http://uniquation.com/en/>, <http://www.searchonmath.com/> that provide equation searching features. Most of the work done involves converting the equations available into xmls, substitution trees and indexing via various open source software. On the other hand we are extracting various features other than simple xmls of the equations that help in increasing the efficiency of the system, as described in the following sections.

Outline of workflow

1 Convert Latex Formula to xml

The conversion from Latex formulae to xml has been performed using MathML

e.g : $\lambda = 1$ is converted to the following xml

```
<?xml version="1.0" encoding="UTF-8"?>
<m:math xmlns:m="http://www.w3.org/1998/Math/MathML" display="block">
<m:mrow>
<m:mi>  $\lambda$  </m:mi>
<m:mo>=</m:mo>
<m:mn>1</m:mn>
</m:mrow>
</m:math>
```

2 Simplifying using equations and generating additional xmls.

The simplification has been done using python library sympy. The MathML's are converted into expressions and sympy's simplify routine is called and the corresponding MathML is generated. We work on both the original MathML and the simplified MathML. e.g : $((a + b) + c)$ is simplified to $a + b + c$ and then MathML is generated as above. If the operands are constants, then the equation is simplified by evaluating it. e.g : $2 + 3.5$ is replaced by 5.5 .

3 Number Normalization and Unicode Normalization.

3.1 Number Normalization

Firstly, we are taking the MathML and taking the decimal numbers obtained inside $< mn >$ tags. We are adding all the numbers obtained by reducing precision. For eg. - Original MathML :-

```
< mathxmlns = "http://www.w3.org/1998/Math/MathML" display = "block" >
< mrow >
< mi > x < /mi >
< mo > + < /mo >
< mn > 2.90646 < /mn >
< /mrow >
< /math >
```

Converted MathML :-

```
< mathxmlns = "http://www.w3.org/1998/Math/MathML" display = "block" >
< mrow >
< mi > x < /mi >
< mo > + < /mo >
```

< mn > 2.90646 < /mn > < mn > 3.0 < /mn > < mn > 2.9 < /mn > < mn > 2.91 < /mn > < mn > 2.906 < /mn > < mn > 2.9065 < /mn >
< /mrow >
< /math >

3.2 Unicode Normalization

Here we are converting all unicode variants of a letter back to the same letter to increase search results.

Add description and examples.

4 Operator Grouping

Add description

5 Unigram + Bigram + Trigram Feature Extraction

Here the xml is treated simply as string and the unigram, bigram and trigram features are extracted. The corresponding postings list (inverted indexing of the equations) are also computed.

Examples (The format is {feature : [(eqn_id, term_frequency)]} :-

Unigram Features :- {< mi > ZN < /mi > : [(2586, 1)], < mn > x2147; < /mn > : [(1338, 1), (4122, 1)]}

Bigram Features :- {'< m : mo > 2887 < /m : mo >', '< m : mi > u03b2 < /m : mi >'} : [(1807, 1), (1808, 3)]}

Trigram Features :- {'< m : mn > 56 < /m : mn >', '< m : mo > < /m : mo >', '< m : mi > f < /m : mi >'} : [(2143, 1)]}

6 TF-IDF weights calculation

Each equation is converted into a vector representation using the extracted unigram, bigram and trigram features

The $tf_{weight}(eqn, feature) = 1 + \log_{10}tf$, $idf_{weight}(feature) = 1 + \log_{10}(N/df)$, here $tf_{weight}(eqn, feature)$ stands for the term frequency or the number of times the “feature” occurs in “eqn” and $idf_{weight}(feature)$ stands for how rarely/frequently the term occurs in the equations.

7 Query Normalization

Latex query given as input is also normalised as mentioned above to reduce unicode variants. The resulting equation is also simplified. Finally the query is converted into a high dimensional vector space using the extracted features mentioned above.

8 Cosine Similarity

Currently cosine similarity is being used as the metric to determine the similarity between query vector and the indexed vectors.

9 Backend and frontend

The Backend is running on a python server and webpage is used as the frontend

The LaSer Application

Add snapshots of UI here

Conclusions and Future Work

- Support free form queries as input.

- Retrieve on basis of structural similarity, for example: $x + y = 1$ and $a + b$ and $(2 * X - 1) + (2 * Y)$ are all structurally similar.
- Retrieve on simplification for example $x * x = 1$ and $x^2 = 1$ are same after simplification.
- Support search by names of popular equations.

References

1. A Search Engine for Mathematical Formulae, Michael Kohlhase and Ioan A. Sucan
2. $5e^{x+y}$: A Math Aware Search Engine(for CDS), Master Thesis Project, Arthur Oviedo,
Supervisors: CERN: Nikolaos Kasioumis, EPFL: Karl Aberer