# Text Mining For Definitions in the arXiv

Luis Berlioz

lab232@pitt.edu

Formal Methods in Mathematics / Lean Together 2020

January 10, 2020

## Objectives and Outline

### Objective

*Create a machine learning system that can find the definitions and the terms being defined in large collections of mathematical texts.*

The problem is broken down into two parts:

The Classifier: Tells if a given paragraph is a definition or not

A Named Entity Recognition system: given a definition, returns the term that is being defined (definiendum).

For each part I will describe how to:

- ▶ Get and process the relevant data.
- ▶ Train the machine learning system.
- ▶ Take a look at the results.

# arXiv Website Bulk Download

All the LaTeX source files can be downloaded from an Amazon S3 bucket

- ▶ About 1 Terabyte of .tar files.
- ▶ Each .tar file is about 500 Megabytes.
- ▶ Download without affecting the website's traffic.
- ▶ LaTeX source is converted to a more structured format.

# LaTeXML
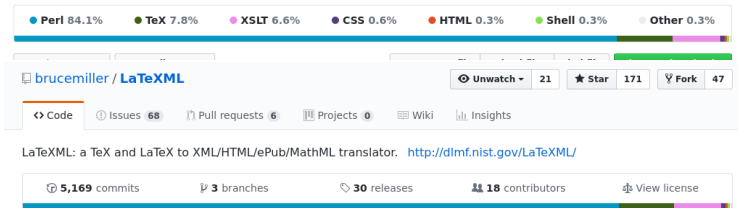
Process each article to get a more structured format

# Obtaining and Classifying Definitions

- Sometimes the author of an article uses a LaTeX macro to label a definition. These are our positive labels:

```
</para>
<theorem class="ltx_theorem_definition" inlist="thm theorem:definition" xml:id="Thmdefinition1">
  <tags>
    <tag>Definition 1</tag>
    <tag role="refnum">1</tag>
    <tag role="typerefnum">Definition 1</tag>
```

- To get examples of non-definitions, we pick paragraphs at random and assume they are not definitions.
- This has the drawback that some of the non-definitions are wrong.
- There are 1,707 articles in 2015 math.AG, we go from 5,229 labeled definitions to 71,067 "probable" definitions.

# Some Classification Results

▶ Results using SVC (Support Vector Classifier) in **scikit-learn**

|     | precision | recall | f1-score | support |
|-----|-----------|--------|----------|---------|
| 0.0 | 0.79      | 0.91   | 0.85     | 2358    |
| 1.0 | 0.95      | 0.88   | 0.91     | 4520    |

▶ Sanity check:

```
Def = ['a banach space is defined as a complete vector space.',
       'This is not a definition honestly. even if it includes technical words like scheme and cohomology',
       'There is no real reason as to why this classifier is so good.',
       'a triangle is equilateral if and only if all its sides are the same length.']
vdef = count_vect.transform(Def)
clf.predict(vdef)
```

```
array([1., 0., 0., 1.])
```



Distribution of appearance of definitions

# Extracting the Definienda

Obtaining the data for Named Entity Recognition system



- ▶ Go through every of wikipedia article looking for a Definition section that contains the title.
- ▶ We obtain a pair: (**Definienda**, Definition).
- ▶ Just 5,321 matches out of almost 6 million articles.
- ▶ Other websites/datasets:
  - ▶ All types of wikis, e.g. ProofWiki, GroupProps (500)
  - ▶ The Stacks project (3,000)
  - ▶ PlanetMath (1,500)

# Training and Evaluating the NER System

Results of the IOB parser using the ChunkParserI method in the **nltk** library

| Input | | Output |
|-------|-----|--------|
| Token | POS | NER |
| We | PRP | O |
| define | VBP | O |
| a | DT | O |
| Banach | NNP | B–DFNDUM |
| space | NN | I–DFNDUM |
| as | IN | O |
| a | DT | O |
| complete | JJ | O |
| vector | NN | O |
| space | NN | O |

```
ChunkParse score:
    IOB Accuracy:   91.2%%
    Precision:      32.0%%
    Recall:         68.0%%
    F-Measure:      43.5%%
```
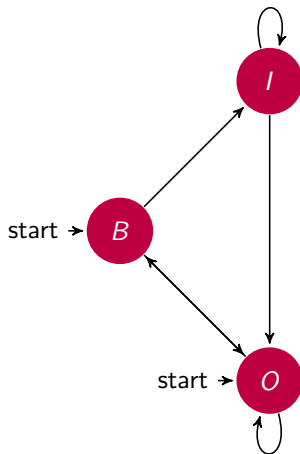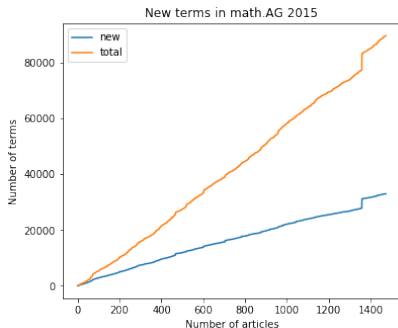
Some definitions found in the 2015 math.DG articles:

## Ex. Things we like

An induced generalized Kähler structure on _inline_math_ is a Lie algebraic generalized Kähler structure with _inline_math_. It is a canonical generalized Kähler structure if _inline_math_.

## Ex. Things we don't like

Suppose _inline_math_ is a vector space. The only connection on the graded manifold _inline_math_ is the trivial connection.



New terms in math.AG 2015

# Conclusions and Future Work

- We think that we have collected enough evidence to believe that a robust collector of definitions is possible.
- A lot of interesting work ahead:
    - Organize the definitions in a *dependency tree structure.*
    - Produce word embedding with math tokens
      (e.g. where *Banach space* is just one token).
    - Apply disambiguation and polysemy detection techniques.

# Text Mining For Definitions in the arXiv

Luis Berlioz
lab232@pitt.edu

Formal Methods in Mathematics / Lean Together 2020

January 10, 2020