



IDENTIFIER NAMESPACES IN MATHEMATICAL NOTATION

MASTER THESIS

by

Alexey Grigorev

Submitted to the Faculty IV, Electrical Engineering and Computer
Science Database Systems and Information Management Group in
partial fulfillment of the requirements for the degree of

Master of Science in Computer Science

as part of the ERASMUS MUNDUS IT4BI programme

at the

TECHNISCHE UNIVERSITÄT BERLIN

July 31, 2015

Thesis Advisors:

Moritz SCHUBOTZ

Juan SOTO

Thesis Supervisor:

Prof. Dr. Volker MARKL

Eidesstattliche Erklärung

Ich erkläre an Eides statt, dass ich die vorliegende Arbeit selbstständig verfasst, andere als die angegebenen Quellen/Hilfsmittel nicht benutzt, und die den benutzten Quellen wörtlich und inhaltlich entnommenen Stellen als solche kenntlich gemacht habe.

Statutory Declaration

I declare that I have authored this thesis independently, that I have not used other than the declared sources/resources, and that I have explicitly marked all material which has been quoted either literally or by content from the used sources.

Berlin, July 31, 2015

Alexey GRIGOREV

Table of Contents

1	Introduction	5
1.1	Namespaces in Computer Science	5
1.2	Namespaces in Mathematics	7
1.3	Outcomes	9
1.4	Thesis Outline	10
2	Mathematical Definition Extraction	11
2.1	Formula Representation: MathML	12
2.2	Math-aware POS tagging	15
2.3	Extraction Methods	16
2.4	Performance Measures	20
3	Namespaces as Document Clusters	21
3.1	Namespaces in Mathematical Notation	21
3.2	Discovering Namespaces with Document Cluster Analysis	23
3.3	Vector Space Model	24
3.4	Identifier Space Model	27
3.5	Similarity Measures and Distances	28
3.6	Inverted Index	32
4	Document Clustering Techniques	34
4.1	Agglomerative clustering	34
4.2	K -Means	35
4.3	Extensions of K -Means	37
4.4	DBSCAN	38
4.5	Extensions of DBSCAN	40
5	Discovering Latent Semantics	41
5.1	Latent Semantic Analysis	41
5.2	Non-Negative Matrix Factorization	44
6	Implementation	47
6.1	Data set	47
6.2	Definition Extraction	47
6.3	Data Cleaning	50
6.4	Document Clustering	51
6.5	Building Hierarchy	54
6.6	Java Language Processing	56
7	Evaluation	59
7.1	Parameter Tuning	59

7.2	Result analysis	63
7.3	Building Hierarchy	64
7.4	Experiment Conclusions	64
8	Conclusions	65
8.1	Future Work	65
9	Bibliography	66

1 Introduction

namespaces are ... In this thesis we extend the notion of namespaces to mathematical formulae.

Review CS and extend in

1.1 Namespaces in Computer Science

In computer science, a *namespace* refers to a collection of terms that are managed together because they share functionality or purpose, typically for providing modularity and resolving name conflicts [1].

XML (eXtensible Markup Language) is a framework for defining markup languages. XML lets users define a set of tags to represent information in some specific domain [2]. For example, XHTML is an XML language for hypertext markup and MathML is a language for describing mathematical notation.

However, different XML languages may use the same names for elements and attributes. For example, consider two XML languages: XHTML for specifying the layout of web pages, and some XML language for describing furniture. Both these languages have the `<table>` elements there, in XHTML table is used to present some data in a tabular form (see listing 1), while the second one uses it to describe a particular piece of furniture in the database (see listing 2).

Listing 1: Describing tabular data in XHTML with the `<table>` tag

```
<table>
  <tr><th>China</th><th>Germany</th><th>Russia</th></tr>
  <tr><td>10</td><td>3</td><td>1</td></tr>
</table>
```

Listing 2: The `<table>` tag in an XML language for describing furniture

```
<table id="table31337">
  <width unit="cm">130</width>
  <length unit="cm">90</length>
  <color>Black</color>
</table>
```

The `<table>` elements have very different semantics in these languages and there should be a way to distinguish between these two elements. In XML this problem is solved with XML namespaces [3]: the namespaces are used to ensure the uniqueness of attributes and resolve ambiguity. It is done

by binding a short namespace alias with some uniquely defined URI (Unified Resource Identifier), and then appending the alias to all attribute names that come from this namespace. In the example above, we can bind an alias `h` with XHTML's URI <http://www.w3.org/TR/xhtml1> and then use `<h:table>` to refer to XHTML's table. Likewise, in the furniture database language the element names can be prepended with a prefix `d`, where `d` is bound to some URI, e.g. <http://www.furniture.de/2015/db> (see listing 3).

Listing 3: Namespaces in XML

```
<html xmlns:h="http://www.w3.org/TR/xhtml1">
  <h:body>
    <h:table>
      <h:tr><h:th>China</h:th><h:th>Germany</h:th><h:th>Russia</h:th></h:tr>
      <h:tr><h:td>10</h:td><h:td>3</h:td><h:td>1</h:td></h:tr>
    </h:table>
  </h:body>
</html>

<database xmlns:d="http://www.furniture.de/2015/db">
  <d:table id="table31337">
    <d:width unit="cm">130</d:width>
    <d:length unit="cm">90</d:length>
    <d:color>Black</d:color>
  </d:table>
</database>
```

The namespaces are also used in programming languages for organizing variables, procedures and other identifiers into groups and for resolving name collisions. In programming languages without namespaces the programmers have to take special care to avoid naming conflicts. For example, in the PHP programming language prior to version 5.3 [4] there is no notion of namespace, and the namespaces have to be emulated to ensure that the names are unique, and this often results in long names like `Zend_Search_Lucene_Analysis_Analyzer`¹.

Other programming languages have the notion of namespaces built in from the very first versions. For example, the Java programming language [5] uses packages to organize identifiers into namespaces. In Java, packages solve the problem of ambiguity. For example, in the standard Java API there

¹ http://framework.zend.com/apidoc/1.7/Zend_Search_Lucene/Analysis/Zend_Search_Lucene_Analysis_Analyzer.html

are two classes with the name `Date`: one in the package `java.util` and another in the package `java.sql`. To be able to distinguish between them, the classes are referred by their *fully qualified name*: an unambiguous name that uniquely specifies the class by combining the package name with the class name. Thus, to refer to a particular `Date` class in Java `java.util.Date` or `java.sql.Date` should be used.

It is not always convenient to use the fully qualified name in the code to refer to some class from another package. Therefore in Java it is possible to *import* the class by using the import statement which associates a short name alias with its fully qualified name. For example, to refer to `java.sql.Date` it is possible to import it by using `import java.sql.Date` and then refer to it by the alias `Date` in the class [5].

Although there is no strict requirement to organize the classes into well defined groups, it is a good software design practice to put related objects into the same namespace and by doing this achieve better modularity. There are design principles that tell software engineers how to best organize the source code: classes in a well designed system should be grouped in such a way that namespaces exhibit low *coupling* and high *cohesion* [6]. Coupling describes the degree of dependence between namespaces, and low coupling means that the interaction between classes of different namespaces should be as low as possible. Cohesion, on the other hand, refers to the dependence within the classes of the same namespace, and the high cohesion principle says that the related classes should all be put together in the same namespace.

1.2 Namespaces in Mathematics

Informally, a mathematical formula is a rule that shows the relationship between different variables. For example, $x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ is a formula for solving a quadratic equation $ax^2 + bx + c = 0$.

To give a more formal definition of formula, we first need to define a first-order language that contains primitive symbols such as (1) parentheses, brackets and other boundary symbols (2) *constants* (1, 2, \hbar , ...) and variables (x , y , ...) (3) *functions* (+, \cdot , ...) and (4) *predicates*, e.g. binary relation symbols ($=$, $<$, \geq , ...).

In this language, *constants* are symbols with pre-defined meaning from some alphabet and variables are symbols that can be assigned a value from this alphabet. Any symbol can be a variable, for example, x , y , \mathbf{w} , or it can be a symbol with subscripts, for example, x_1, x_2, \dots or even w_{slope} .

A *well-formed term* t (or just *term*) in this language is defined as

$$t \equiv c \mid x \mid f(t_1, t_2, \dots, t_n)$$

, which means that the term t can be a constant, a variable or an n -ary function $f(t_1, t_2, \dots, t_n)$. An *n -ary function* is an function that takes n terms t_1, t_2, \dots, t_n and produces some new term. An *n -ary predicate* (or an *n -ary relation symbol*) is typically a boolean-valued function that can be evaluated to **True** or **False** depending on the values it gets.

Then an *atomic well-formed formula* (or just *formula*) in this language is a n -ary predicate with n terms evaluated to **True** [7].

For example, $x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ is a formula because it represents an equation that always holds true for a quadratic equation $ax^2 + bx + c = 0$. The equality symbol “=” is a predicate that shows the relationship between variables x_1, x_2 and variables a, b, c .

In logic we can use any symbol for variables without changing the meaning of the formula. For example, the energy-mass equivalence relation $E = mc^2$ can be written as $x = yz^2$ and it still will still hold true and remain a valid formula. However, there are research communities in mathematics that have developed a special system of naming these variables, and these naming systems are called *mathematical notations*. For each symbol in a formula, the notation assigns a precise semantic meaning. Therefore, because of the notation, in Physics it is more common to write $E = mc^2$ rather than $x = yz^2$, because the notation assigns unambiguous meaning to the symbols “ E ”, “ m ” and “ c ”, and the meaning of these symbols is recognized among physicists.

However the notations may conflict. For example, while it is common to use symbol E to denote “Energy” in Physics, it also is used in Probability and Statistics to denote “Expected Value”, or in Linear Algebra to denote “Elimination Matrix”. We can compare the conflict of notations with the name collision problem in namespaces, and solve it by extending the notion of namespaces to mathematical notation.

Thus, let us define a *notation* \mathcal{N} as a set of pairs $\{(i, s)\}$, where i is a symbol and s is its semantic meaning, such that for any pair $(i, s) \in \mathcal{N}$ there does not exist a pair $(i', s') \in \mathcal{N}$ with $i = i'$. Let us call i *identifier* and s *definition*, and then (i, s) is an *identifier-definition pair*. We say that two notations \mathcal{N}_1 and \mathcal{N}_2 *conflict* if there exists a pair $(i_1, s_1) \in \mathcal{N}_1$ and a pair $(i_2, s_2) \in \mathcal{N}_2$ such that $i_1 = i_2$ and $s_1 \neq s_2$.

Then we can define *namespace* as a named notation, i.e. a pair (name, \mathcal{N}) where \mathcal{N} is a notation and “name” is a string that uniquely defines the nota-

tion. For convenience we can use the Java syntax to refer to specific entries of a namespace. If $(\text{name}, \mathcal{N})$ is a namespace and i is an identifier such that $(i, s) \in \mathcal{N}$ for some s , then “name. i ” is a *fully qualified name* of the identifier i that relates it to the definition s . For example, given a namespace $(\text{“Physics”}, \{(E, \text{“energy”}), (m, \text{“mass”}), (c, \text{“speed of light”})\})$, “Physics. E ” refers to “energy” – the definition of E in the namespace “Physics”.

Analogously to namespaces in Computer Science, formally a mathematical namespace can contain any set of identifier-definition pairs that satisfies the definition of the namespace, but typically namespaces of mathematical notation exhibit the same properties as well-designed Java packages: they have low coupling and high cohesion, meaning that all definitions come from the same area of mathematical knowledge and the definitions from different notations do not intersect heavily.

Additionally, we can introduce a document-centric view on the mathematical namespaces: suppose we have a collection of documents of n documents $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$ and a set of namespaces $\{(n_1, \mathcal{N}_1), (n_2, \mathcal{N}_2), \dots, (n_K, \mathcal{N}_K)\}$. A document d_j can use a namespace (n_k, \mathcal{N}_k) by *importing* identifiers from it. To import an identifier, the document uses an import statement where the identifier i is referred by its fully qualified name defined by this namespace.

TODO: Continue, re-introduce low coupling and high cohesion in terms of document collection.

However, in real-life scientific documents there are no import statements in the document preamble, and therefore ...

The goal is to automatically discover the set of mathematical namespaces given a collection of documents.

Manual labeling is time consuming

Want: approximation via ML algorithms

To do that need to be able to extract definition from text

1.3 Outcomes

TODO: taken from the proposal. The main idea here to describe what can happen if we solve the problem of namespace discovery in notation.

Once such namespaces are found, they can give good categorization of scientific documents based on formulae and notation used in them.

We believe that this may facilitate better user experience: for instance, it will allow users to navigate easily between documents of the same category

and see in which other documents a particular identifier is used, how it is used, how it is derived, etc. Additionally, it may give a way to avoid ambiguity. If we follow the XML approach [3] and prepend namespace to the identifier, e.g. “physics. E ”, then it will give additional context and make it clear that “physics. E ” means “energy” rather than “expected value”.

We also expect that using namespaces is beneficial for relating identifiers to definitions. Thus, as an application of namespaces, we would like to be able to use them for better definition extraction. It may help to overcome some of the current problems in this area, for example, the problem of *dangling identifiers* [8] - identifiers that are used in formulae but never defined in the document. Such identifiers may be defined in other documents that share the same namespace, and thus we can take the definition from the namespace and assign it to the dangling identifier.

1.4 Thesis Outline

This work is organized as follows: In chapter 2 we discuss how extract definitions for identifiers in texts with mathematical formulae; in chapter 3 the approaches to namespace extraction and we argue why the cluster analysis can be used for that; in chapter 4 we review cluster analysis methods and in chapter 5 we discuss how to extract latent information from data using different matrix factorization techniques. Finally, we describe how the techniques are implemented (chapter 6) and evaluated (chapter 7).

2 Mathematical Definition Extraction

Mathematical expressions are hard to understand without the natural language description, therefore we want to extract identifiers from mathematical expressions and then find their definitions from the surrounding text.

For example, given the sentence “The relation between energy and mass is described by the mass-energy equivalence formula $E = mc^2$, where E is energy, m is mass and c is the speed of light” the goal is to extract the following identifier-definition relations:

- (E , “energy”)
- (m , “mass”)
- (c , “the speed of light”)

Consider another example: “Let e be the base of natural logarithm”. We would like to extract (e , “the base of natural logarithm”).

Formally, a phrase that defines a mathematical expression consists of three parts [9]:

- *definiendum* is the term to be defined: it is a mathematical expression or identifier;
- *definiens* is the definition itself: it is the word or phrase that defines the definiendum in a definition.
- *definitior* is a relator verb that links definiendum and definiens.

In this work we are interested in the first two parts: *definiendum* and *definiens*. Thus we define a *relation* as a pair (definiendum, definiens). For example, (E , “energy”) is a relation where E is a definiendum, and “energy” is a definiens.

We have the following assumption about definition (i.e. definitions): Assumption: the definitions of mathematical expressions are always noun phrases

In general, a noun phrase can be

- a simple noun
- a compound noun (e.g. adjective + noun)
- a compound noun with a clause, prepositional phrase, etc

In this chapter we will discuss how the relations can be discovered automatically. The typical processing pipeline for definition extraction consists of the following steps [9] [8]:

- Read corpus of documents in some markup language, e.g. Wiki Markup or Latex
- Translate all formulas in the documents into MathML
- Process MathML formulas
- Replace formulas with some placeholder
- Annotate text using Math-Aware POS Tagging
- Find relations in the text

Thus, this chapter is organized as follows: first we introduce the Mathematical Markup Language (MathML) in section 2.1, then discuss the Math-Aware POS Tagging procedure in section 2.2 and finally review the extraction methods in section 2.3 and briefly discuss how the quality of extracted identifiers is evaluated in section 2.4.

2.1 Formula Representation: MathML

MathML [10] stands for “Mathematical Markup Language” It is a standard for mathematical expressions defined by W3C that browsers should support to render math formulas. There are two types of MathML: Presentation MathML, which describes how mathematical expressions should be displayed, and Content MathML, which focuses on the meaning of mathematical expressions. In this section, we will discuss Presentation MathML.

A *token* in MathML is an individual symbol, name or number. Tokens are grouped together to form MathML expressions.

Tokens can be:

- identifier, variable or function names
- numbers
- operators (including brackets - so called “fences”)
- text and whitespaces

A “symbol” is not necessarily one character: it could be a string such as `<mi>sin</mi>` or `<mn>24</mn>`. In MathML they are treated as single tokens.

As in mathematics, MathML expressions are constructed recursively from smaller expressions or single tokens. Complex expressions are created with so-called “layout” constructor elements, while tokens are created with token elements.

Let us consider an example. A mathematical expression $(a + b)^2$ can be represented in MathML as follows:

Listing 4: TODO

```

<math xmlns="http://www.w3.org/1998/Math/MathML">
  <msup>
    <mrow>
      <mo>(</mo>
      <mrow>
        <mi>a</mi>
        <mo>+</mo>
        <mi>b</mi>
      </mrow>
      <mo>></mo>
    </mrow>
    <mn>2</mn>
  </msup>
</math>

```

It has the tree structure and recursive. If we take another mathematical expression $\frac{3}{(a+b)^2}$. It is a fraction and we see that its denominator is the same as the previous expression. This is also true for the MathML representation:

Listing 5: TODO

```

<math xmlns="http://www.w3.org/1998/Math/MathML">
  <mfrac>
    <mn>3</mn>
    <msup>
      <mrow>
        <mo>(</mo>
        <mrow>
          <mi>a</mi>
          <mo>+</mo>
          <mi>b</mi>
        </mrow>
        <mo>></mo>
      </mrow>
      <mn>2</mn>
    </msup>
  </mfrac>
</math>

```

Token Elements

Token elements are needed for representing tokens: the smallest units of mathematical notation that convey some meaning.

There are several token elements:

- `mi` identifier
- `mn` number
- `mo` operator, fence, or separator
- `mtext` text
- `mspace` space
- `ms` string literal

Often tokens are just single characters, like `<mi>E</mi>` or `<mn>5</mn>`, but there are cases when tokens are multi-character, e.g. `<mi>sin</mi>` or `<mi>span</mi>`.

In MathML `mi` elements represent some symbolic name or text that should be rendered as identifiers. Identifiers could be variables, function names, and symbolic constants.

Transitional mathematical notation often involve some special typographical properties of fonts, e.g. using bold symbols e.g. \mathbf{x} to denote vectors or capital script symbols e.g. \mathcal{G} to denote groups and sets. To address this, there is a special attribute “mathvariant” that can take values such as “bold”, “script” and others.

Numerical literals are represented with `mn` elements. Typically they are sequences of digits, sometimes with a decimal point, representing an unsigned integer or real number, e.g. `<mn>50</mn>` or `<mn>50.00</mn>`.

Finally, operators are represented with `mo` elements. Operators are ...

Layouts

Layout elements are needed to form complex mathematical expressions from simple ones. They group elements in some particular way. For example:

- `mrow` groups any number of sub-expressions horizontally
- `mfrac` form a fraction from two sub-expressions
- `msqrt` forms a square root (radical without an index)

Some layout elements are used to add subscripts and superscripts:

- `msub` attach a subscript to a base
- `msup` attach a superscript to a base
- `msubsup` attach a subscript-superscript pair to a base

And special kinds of scripts (TODO: describe in more details)

- `munder` attach an underscript to a base

- `mover` attach an overscript to a base
- `munderover` attach an underscript-overscript pair to a base

For example, \boldsymbol{v} will be rendered as

Listing 6: TODO

```
<math xmlns="http://www.w3.org/1998/Math/MathML">
  <mover>
    <mi>v</mi>
    <mo>&rarr;</mo>
  </mover>
</math>
```

This is how we would represent $\hat{\mathbf{x}}$ (a bold x with a hat) in MathML:

Listing 7: TODO

```
<math xmlns="http://www.w3.org/1998/Math/MathML">
  <mrow>
    <mover>
      <mrow>
        <mi mathvariant="bold">x</mi>
      </mrow>
      <mo>&#x005E;<!-- ^ --></mo>
    </mover>
  </mrow>
</math>
```

There are more complex elements such as `mtable`.

MathML presentation elements only suggest specific ways of rendering Math Entities

Certain characters are used to name identifiers or operators that in traditional notation render the same as other symbols or usually rendered invisibly.

entities `⁢` `→`

The complete list of MathML entities is described in [Entities].

2.2 Math-aware POS tagging

Part-of-Speech Tagging (POS Tagging) is a typical Natural Language Processing task which assigns a POS Tag to each word in a given text [11]. While the POS Tagging task is mainly a tool for text processing, it can also be applicable to scientific documents with mathematical expressions, and can be adjusted to dealing with formulae.

A *POS tag* is an abbreviation that corresponds to some part of speech. Penn Treebank POS Scheme [12] is a commonly used POS tagging scheme which defines a set of part-of-speech tags for annotating English words. For example, JJ is an adjective (“big”), RB as in adverb, DT is a determiner (“a”, “the”), NN is a noun (“corpus”) and SYM is used for symbols (“>”, “=”).

However the Penn Treebank scheme does not have special tags for mathematics, but it is flexible enough and can be extended to include additional tags. For example, we can include a math-related tag **MATH**. Usually it is done by first applying traditional POS taggers (like Stanford CoreNLP [13]), and then refining the results by re-tagging math-related tokens of text as **MATH** [14].

For example, consider the following sentence: “The relation between energy and mass is described by the mass-energy equivalence formula $E = mc^2$, where E is energy, m is mass and c is the speed of light”. In this case we will assign the **MATH** tag to “ $E = mc^2$ ”, “ E ”, “ m ” and “ c ”

However we can note that for finding identifier-definition relations the **MATH** tag alone is not sufficient: we need to distinguish between complex mathematical expressions and stand-alone identifiers - mathematical expressions that contain only one symbol: the identifier. For the example above we would like to be able to distinguish the expression “ $E = mc^2$ ” from identifier tokens “ E ”, “ m ” and “ c ”. Thus we extend the Penn Treebank scheme even more and introduce an additional tag **ID** to denote stand-alone identifiers.

Thus, in the example above “ $E = mc^2$ ” will be assigned the **MATH** tag and “ E ”, “ m ” and “ c ” will be annotated with **ID**.

In the next section we discuss how this can be used to find identifier-definition relations.

2.3 Extraction Methods

There are several ways of extracting the identifier-definition relations. Here we will review the following:

- Nearest Noun
- Pattern Matching
- Machine-Learning based methods
- Probabilistic methods

Nearest Noun Method

The Nearest Noun [15] [16] is the simplest definition extraction method. It assumes that the definition is a combination of ad It finds definitions by looking for combinations of adjectives and nouns (sometimes preceded by determiners) in the text before the identifier.

I.e. if we see a token annotated with ID, and then a sequence consisting only of adjectives (JJ), nouns (NN, NNS) and determiners (DET), then we say that this sequence is the definition for the identifier.

For example, given the sentence “In other words, the bijection σ normalizes G in ...” we will extract relation $(\sigma, \text{”bijection”})$.

Pattern Matching Methods

The Pattern Matching method [17] is an extension of the Nearest Noun method: in Nearest Noun we are looking for one specific pattern where identifier is followed by the definition, but we can define several such patterns and use them to extract definitions.

For example, we can define the following patterns:

- IDE DEF
- DEF IDE
- let|set IDE denote|denotes|be DEF
- DEF is|are denoted|defined|given as|by IDE
- IDE denotes|denote|stand|stands as|by DEF
- IDE is|are DEF
- DEF is|are IDE
- and many others

In this method IDE and DEF are placeholders that are assigned a value when the pattern is matched against some subsequence of tokens. IDE and DEF need to satisfy certain criteria in order to be successfully matched: like in the Nearest Noun method we assume that IDE is some token annotated with ID and DEF is a phrase containing adjective (JJ), nouns (NN) and determiners (DET). Note that the first pattern corresponds to the Nearest Noun pattern.

The patterns above are combined from two sources: one is extracted from a guide to writing mathematical papers in English ([18]) by **TODO**, and another is extracted from Graphs and Combinatorics papers from Springer by **TODO**.

The pattern matching method is often used as the baseline method for identifier-definition extraction methods [9] [19] [8].

Machine Learning Based Methods

The definition extraction problem can be formulated as a binary classification problem: given a pair (identifier, candidate-definition), does this pair correspond to real identifier-definition relation?

To do this we find all candidate pairs: identifiers are tokens annotated with ID, and candidate definitions are nouns and noun phrases from the same sentence as the definition.

Once the candidate pairs are found, we extract the following features [19] [16]:

- boolean features for each of the patterns from section 2.3 indicating if the pattern is matched
- indicator if there’s a colon or comma between candidate and identifier
- indicator if there’s another math expression between candidate and identifier
- indicator if candidate is inside parentheses and identifier is outside
- distance (in words) between the identifier and the candidate
- the position of candidate relative to identifier
- text and POS tag of one/two/three preceding and following tokens around the candidate
- text of the first verb between candidate and identifier
- many others

Once the features are extracted, a binary classifier can be trained to predict if an unseen candidate pair is a relation or not. For this task the popular choices of classifiers are Support Vector Machine classifier with linear kernel [19] [16] and Conditional Random Fields [19], but, in principle, any other binary classifier can be applied as well.

Probabilistic Approaches

In the Mathematical Language Processing approach [8] an identifier-definition relation is seen as a pair $\langle \text{identifier}, \text{probability distribution over definition candidates} \rangle$. Thus, for definition extraction the candidate definitions are ranked by their probability of defining the identifier, and then only most probable candidates are retained.

The main idea of this approach is that the definitions occur very closely to identifiers in sentences, and the closeness can be used to model the probability distribution over candidate definitions.

There are two assumptions:

- identifier and its definition can only occur in the same sentence, so the candidates are definitions are taken only from the sentence where the identifier occurs (as in the Machine Learning approach)
- definitions are more likely to occur closer to the formula where the identifier is used

With these assumptions in mind, the distribution over definition candidates can be modeled with:

- probability distribution $R_{\sigma_d}(\Delta(i, t))$ where $\Delta(i, t)$ is the distance between identifier i and candidate definition t , parameterized with σ_d ,
- probability distribution $R_{\sigma_s}(n(i, t))$ where $n(i, t)$ is the number of sentences between the formula where i is used and the sentence where t occurs, parameterized with σ_s , and
- term frequency $\text{tf}(t, s)$: how many times t occurs in the sentence s

All three elements can be combined together for ranking candidate definitions:

$$R(n, \Delta, t, d) = \frac{\alpha R_{\sigma_d}(\Delta) + \beta R_{\sigma_s}(n) + \gamma \text{tf}(t, s)}{\alpha + \beta + \gamma}$$

where α, β, γ are weighting parameters.

Instead of taking the raw distances, the Gaussian distribution is used to model the fact that the probability of being the definition does not decrease linearly as we move away from the identifier, but rather exponentially. Thus, the distances are modeled with

$$R_{\sigma}(\Delta) = \exp\left(-\frac{1}{2} \cdot \frac{\Delta^2 - 1}{\sigma^2}\right)$$

assuming that the probability to find the relation at $\Delta = 1$ is maximal.

There are several parameters in this model: σ_d and σ_s for controlling the widths of $R_{\sigma_d}(\Delta)$ and $R_{\sigma_s}(n)$ respectively; and α, β, γ are parameters for controlling the contribution of different ranking components.

The parameter σ_d is the standard deviation of Gaussian that models the distance to definition candidate. By examining several mathematical articles manually it has been established that $R_{\sigma_d}(1) \approx 2 \cdot R_{\sigma_d}(5)$, i.e. it is two times more likely to find the real definition at distance $\Delta = 1$ rather than at distance $\Delta = 5$, and thus $\sigma_d = \sqrt{12/\ln 2}$.

The parameter σ_s is the standard deviation of the Gaussian that models the distance (in the number of sentences) between the candidate definition and the formula where the identifier is used. Manual evaluation has shown that $\sigma_s = 2\sqrt{1/\ln 2}$.

Finally, for weights α, β, γ the following parameters were chosen in [8]: $\alpha = \beta = 1$ and $\gamma = 0.1$.

2.4 Performance Measures

The common way to evaluate the performance of an Information Retrieval system is to use Precision and Recall [20]. *Precision* is typically defined as is the fraction of retrieved documents that are relevant, while *recall* is defined as the fraction of relevant documents that are retrieved. These performance measures can be adapted to measure the quality of Formula Extraction systems [8] as:

- *Precision*: the fraction of definitions that are correct among extracted. It is calculated as the number of correctly extracted definitions divided by the total number of extracted definitions.
- *Recall*: the fraction of definitions correctly extracted among all correct definitions. It is calculated by dividing the number of correctly extracted definitions on the total number of identifiers with definition.

A common way of incorporating two measures into a single one is F_β -score defined as $F_\beta = \frac{(\beta^2+1)PR}{\beta^2P+R}$ where P and R are precision and recall respectively, and β is the trade-off parameter. $F_1 = \frac{2PR}{P+R}$ is the balanced F score when both precision and recall have equal weights.

3 Namespaces as Document Clusters

In this chapter we discuss how the process of namespace discovery can be automated.

First, in section 3.1 we describe identifier namespaces and we compare the namespace discovery with cluster analysis techniques applied to textual data and see how clustering algorithms can be useful. Next, in section 3.3 we review the Vector Space Model (VSM): the traditional way of representing a collection of documents as vectors, and then in section 3.4 we introduce the Identifier VSM - which is a way to represent identifier-definition relations in the vector space. Finally we go over common similarity and distance functions that are useful for document clustering in section 3.5 and discuss how similarity search can be made faster by using inverted index 3.6.

3.1 Namespaces in Mathematical Notation

TODO: rewrite a bit to keep the introduction in mind. Also refer to introduction

An *identifier namespace* is a coherent structure where each identifier is used only once and has a unique definition.

How to find a namespace? Namespaces can be constructed by manually labeling each identifier/definition pair with appropriate name. But this is very time consuming, and in this work we suggest a different approach: use Machine Learning techniques for discovering namespaces automatically from a collection of scientific documents containing mathematical formulae.

Many modern programming languages use namespaces for modularity. For example, in the Java programming language [5] namespaces are called “packages” and a class may refer to classes from other packages via the `import` statement.

Typically in a well-designed application, we can distinguish between two types of application packages [21]:

- *type 1*: domain-specific packages that deal with one particular concept, and
- *type 2*: packages that use many other packages of the first type

For example, we have an application `org.company.app` with several domain-specific packages: `org.company.app.domain.user` with classes related to users, `org.company.app.domain.account` with classes related to user accounts, and a system-related package `org.company.app.tools.auth`

that deals with authentication and authorization. Then we also have a package `org.company.app.web.manage`, which belongs to the type 2: it handles web requests while relying on classes from `user` and `account` to implement the business logic and on `auth` for making sure the requests are authorized.

We can observe that the type 1 packages are mostly self-contained and not highly coupled between each other, but type2 packages mostly use other packages of type 1: they depend on them.

We can extend this idea to scientific documents and identifier namespaces. A document can be seen as a class that uses concepts defined in other documents. Then the documents can be grouped such that some groups are of *type 1*: they define the namespaces. In some sense the documents of type 1 are “pure” - they contain information about closely related concepts and not highly coupled with other document groups. But some documents are of *type 2* and they are not pure: they draw from different concepts.

This intuition allows us to have the following assumptions:

1. documents are “mixtures” of namespaces: they take identifiers from several namespaces
2. there are some documents are more “pure” than others: they either take identifiers exclusively from one namespace or from few very related namespaces

With these assumptions we can refer to “pure” groups as *namespace defining* groups. These groups can be seen as “type 1” packages: they define namespaces that are used by other “type 2” document groups.

Additionally, we can assume that there is a strong correlation between identifiers in a document and the namespace of the document, and this correlation can be exploited to categorize documents into groups.

Thus by combining these assumptions we can conclude that it should be approximate the process of namespace discovery by discovering groups of namespace defining documents, and this can be done by applying cluster analysis techniques to documents, represented by identifiers they contain.

In the next section we will argue why we can use traditional document clustering techniques and what are the characteristics that texts and identifiers have in common.

3.2 Discovering Namespaces with Document Cluster Analysis

We believe that cluster analysis techniques developed for text documents should work for identifiers.

Let us consider the characteristics of text data:

There are many distinct words in natural language. For example, if \mathcal{V} is a set of all possible words, then usually $|\mathcal{V}| \approx 10^5$, but each individual document may contain only 500 distinct words, or sometimes even less if we consider sentences or small documents (e.g. tweets)

number of words across different document may vary a lot

word distributions follow Power Laws (e.g. Zipf's law)

The identifiers have the same properties!

Natural languages suffer from lexical problems of variability and ambiguity, and the two main problems are synonymy and polysemy [22] [23]:

- two words are *synonymous* if they have the same meaning (for example, “word” and “term” are synonyms),
- a word is *polysemous* if it can have multiple meanings (for example, “trunk” can refer to a part of elephant or a part of a car).

We can note that identifiers have the same problems. For example, in Information Theory, the Shannon Entropy is usually denoted by “ H ”, but sometimes it is also denoted by “ I ” or by “ S ”, thus these identifiers may be seen as synonyms. Also, “ E ” can stand both for “Energy” and “Expected value”, so “ E ” is polysemous.

These problems have been studied in Information Retrieval and Natural Language Processing literature. One possible solution for the polysemy problem is Word Sense Disambiguation [11]: either replace a word with its sense [24] or append the sense to the word, for example if the polysemous word is “bank” with meaning “financial institution”, then we replace it with “bank_finance”. The same idea can be used for identifiers, for example “ E ” can be replaced with “ E_{energy} ”.

Document clustering techniques usually use Vector Space Models [25] [26] to represent documents. We can define “Identifier Spaces” analogous to Vector Space Models. We assume that Identifier Spaces exhibit the same characteristics as traditional Vector Space Models, and thus

In the next section we review the Vector Space Model, and then introduce the Identifier VSM in the chapter 3.4

Then we can apply cluster analysis techniques to document-identifier matrices.

3.3 Vector Space Model

Vector Space Model is a statistical model for representing documents in some (very high dimensional) vector space. It is an Information Retrieval model [20], but it is also used for various Text Mining tasks such as Document Classification [27] and Document Clustering [25] [26].

The process of transforming a text to its vector representation is called “vectorization”. But before documents can be vectorized they are preprocessed. The preprocessing usually consists of the following steps:

- tokenization: extracting individual words from the text;
- stop words removal: removes functional words that have no discriminative power;
- word normalization (includes stemming or lemmatization): reduces words to some common form;

There are two assumptions made about the data:

- *Bag of Words assumption*: the order of words is not important, only word counts;
- *Independence assumption*: we treat all words as independent.

Both assumptions are quite strong, but nonetheless this method often gives good results.

Document-Term Matrix - representation of a document for text analysis each row of the matrix - is a “document vector” each component of the document vectors is a concept, a key word, or a term, but usually it’s terms documents don’t contain many distinct words, so the matrix is sparse

Notation: let $\mathcal{D} = \{d_1, \dots, d_n\}$ be a set of m documents and let $\mathcal{V} = \{t_1, \dots, t_m\}$ be a set of n terms (the vocabulary). Each document is set of weighed terms $d_i = \{w_1, \dots, w_m\}$ where w_j is the weight of term t_j .

There are following term weighting schemes [20]:

- binary: 1 if a term is present, 0 otherwise
- term frequency (TF): number of occurrences of the term in a document
- document frequency (DF): number of documents containing the term
- TF-IDF: combination of TF and inverse DF

Term Frequency (TF) weights terms by local frequency in the document. That is, the term is weighed by how many times it occurs in the document. We can define TF formally as

$$\text{tf}(t, d) = |\{t' \in d : t' = t\}|$$

Sublinear TF: sometimes the term is used too often in a document and we want to reduce its influence compared to other less frequent tokens. This can be done by applying some sublinear transformation to TF, for instance, a squared root $\sqrt{\text{tf}(w, d)}$ or a logarithm $\log \text{tf}(w, d)$.

Document Frequency (DF) weights terms by their global frequency in the collection, which is the number of documents that contain the token. Formally it can be defined as

$$\text{df}(t, \mathcal{D}) = |\{d \in \mathcal{D} : t \in d\}|$$

Inverse Document Frequency (IDF): more often we are interested in domain specific words than in neutral words, and these domain specific words tend to occur less frequently and they usually have more discriminative power: that is, they are better in telling one document apart from another. So IDF should give more weights to rare words rather than to frequent words. Typically IDF is defined as follows:

$$\text{idf}(t, \mathcal{D}) = \log \frac{|\mathcal{D}|}{\text{df}(t, \mathcal{D})}$$

A good weighting system gives the best performance when it assigns more weights to terms with high TF, but low DF [28]. This can be achieved by combining both TF and IDF schemes. TF is good for getting high frequency words, but using just TF is not enough if high frequency words are contained in many documents, thus need a collection dependent factor that favors terms that are contained in fewer documents: IDF.

Usually a sublinear TF is used to avoid the dominating effect of words that occur too frequently. As the result, terms appearing too rarely or too frequently are ranked low.

So, we can combine TF and IDF then multiply:

$$\text{tf-idf}(t, d | \mathcal{D}) = (1 + \log \text{tf}(t, d)) \cdot \log \frac{|\mathcal{D}|}{\text{df}(t, \mathcal{D})}$$

Once the weighting scheme is established, documents can be represented by vectors $d_i = (w_1, \dots, w_m)$ where w_j is the weight of term t_j .

Then these vectors can be put together to form a matrix. Let D be a $m \times n$ matrix, where rows of D are indexed by terms t_i , columns of D are indexed by documents d_j , and element $(D)_{ij}$ is a weight w_i of term t_i in document d_j . Then such matrix D is called a *term-document matrix*.

Alternatively, D can be $n \times m$ matrix with rows being indexed by documents d_j and columns - by terms t_i . Then such D is called *document-term matrix*. Note that if D is a term-document matrix, then D^T is a document-term matrix. In Information Retrieval literature term-document matrices are used more often, than document-term matrices, but in some applications, like Clustering, it is more convenient to use document-term matrix.

Let us consider the column space of the term-document matrix D . The column of D are documents from the corpus \mathcal{D} , so the column space of D contains document vectors where dimensions are terms t_1, t_2, \dots, t_m . This vector space is called *Document VSM*

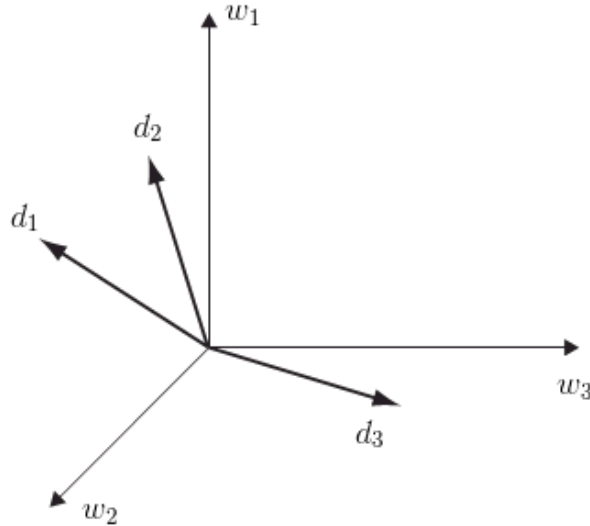


Fig. 1: **TODO redraw in vector**

Alternatively, we can consider the row space of D . This is called the *Term VSM*: the vectors in this space are terms and they are indexed by documents d_1, d_2, \dots, d_n .

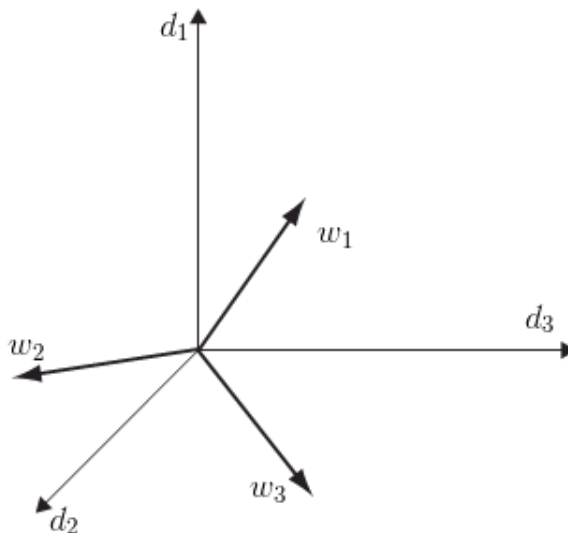


Fig. 2: **TODO** redraw in vector

3.4 Identifier Space Model

The Vector Space Model gives a good foundation. In this work documents are not represented by words they contain, but rather by identifiers they have. Because words and identifiers exhibit similar properties (see section 3.2) we can replace the Term VSM by *Identifier VSM*: a vector space where documents are represented as vectors indexed by identifiers they contain. Thus, the “vocabulary” of this space is $\text{id}_1, \dots, \text{id}_k$, and documents are represented as $d_j = (w_1, \dots, w_k)$ where w_i is the weight of identifier id_i .

Then we can define an identifier-document matrix D as $k \times n$ matrix where columns are documents and rows are identifiers.

The Identifier VSM also suffers from the problems of polysemy and synonymy (see section 3.2). They can be solved by extracting definitions for all the identifiers and incorporating these definitions into the Identifier VSM.

There are three ways of adding the definition information into Identifier VSM:

- use only identifier information, and do not include the definitions;
- use “weak” identifier-definition association: include identifiers and definitions as separate dimensions;
- use “strong” association: append definition to identifier.

To illustrate how it is done, consider two relations (λ , “regularization”) and (w , “weight vector”)

- no definitions: dimension of the Identifier VSM are (λ , w)
- “weak” association: dimensions are (λ , w , regularization, weight vector)
- “strong” association: dimensions are (λ _regularization, w _weight vector)

3.5 Similarity Measures and Distances

Once the documents are represented in some vector space, we need to define how to compare these documents to each other. There are two ways of doing this: using a similarity function that tells how similar two objects are (the higher values, the more similar the objects), or using a distance function, sometimes called “dissimilarity function”, which is the opposite of similarity (the higher the values, the less similar the objects).

We will consider the following similarity and distance functions:

- Euclidean distance
- Inner product (or dot product)
- Cosine similarity
- Jaccard coefficient

Euclidean Distance

The Euclidean distance function is the most commonly used distance function in vector spaces. Euclidean distance corresponds to the geometric distance between two data points in the vector space. For example, if we have two points \mathbf{x} and \mathbf{z} , then the Euclidean distance is the length of the line that connects these two points. It is defined as

$$\|\mathbf{x} - \mathbf{z}\| = \sqrt{(\mathbf{x} - \mathbf{z})^T (\mathbf{x} - \mathbf{z})} = \sqrt{\sum_i (x_i - z_i)^2}$$

This distance is also often called “ L_2 distance”.

Let us rewrite the expression for calculating the Euclidean distance:

$$\|\mathbf{x} - \mathbf{z}\|^2 = (\mathbf{x} - \mathbf{z})^T (\mathbf{x} - \mathbf{z}) = \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{z} + \mathbf{z}^T \mathbf{z} = \|\mathbf{x}\|^2 + \|\mathbf{z}\|^2 - 2\mathbf{x}^T \mathbf{z}$$

A useful way to think about Euclidean distance for document spaces is to consider all three elements of this distance: individual lengths of both vectors and the inner product between them.

Consider two document vectors \mathbf{d}_1 and \mathbf{d}_2 with lengths l_1 and l_2 respectively. If these two documents have no terms in common, then the inner product term $\mathbf{d}_1^T \mathbf{d}_2$ between them is 0, and therefore the (squared) distance is $\|\mathbf{d}_1\|^2 + \|\mathbf{d}_2\|^2 = l_1^2 + l_2^2$. However if we consider two documents of the same lengths l_1 and l_2 that share several terms in common, their inner product, then their inner product is no longer 0, and therefore these documents become closer. Thus, the more terms the documents have in common, the closer they are.

While Euclidean distance is useful for low-dimensional data, it does not always work very well in high dimensions, especially with sparse vectors such as document vectors [29]. The problem is usually caused by the individual length components of the distance function: if two documents have no words in common, but they are not long, they can have the same distance as two long documents that have plenty of words in common.

Consider the following example (from [29]): There are 4 data points $\mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3, \mathbf{d}_4$ in 10-dimensional space indexed by terms A_1, \dots, A_{10}

	A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8	A_9	A_{10}
\mathbf{d}_1	3	0	0	0	0	0	0	0	0	0
\mathbf{d}_2	0	0	0	0	0	0	0	0	0	4
\mathbf{d}_3	3	2	4	0	1	2	3	1	2	0
\mathbf{d}_4	0	2	4	0	1	2	3	1	2	4

The distance between \mathbf{d}_1 and \mathbf{d}_2 is $\|\mathbf{d}_1 - \mathbf{d}_2\| = 5$, and the distance between \mathbf{d}_3 and \mathbf{d}_4 is also $\|\mathbf{d}_3 - \mathbf{d}_4\| = 5$. But when we think of these data points as documents, intuitively \mathbf{d}_3 and \mathbf{d}_4 should be closer because they have 7 terms in common, whereas \mathbf{d}_1 and \mathbf{d}_2 have none.

Thus if the data is sparse it is better to use different measures of distance/similarity, preferably ones that ignore dimensions where both vector have 0 values.

Inner product

The Euclidean distance consists of three components: lengths and the inner product. Since the lengths may cause some problems, it is possible to take only the inner product part and treat it as a similarity function: the more similar two vectors are, the larger is their inner product.

Geometrically the inner product between two vectors \mathbf{x} and \mathbf{z} is defined as

$$\mathbf{x} \cdot \mathbf{z} = \|\mathbf{x}\| \|\mathbf{z}\| \cos \theta$$

where θ is the angle between vectors \mathbf{x} and \mathbf{z} [30]. If the vectors are perpendicular, then $\cos \theta = 0$, and the inner product is also 0. The perpendicular vectors are called *orthogonal*.

The geometric definition makes sense from the Law of Cosine point of view. Consider two vectors \mathbf{x} and \mathbf{z} , and let $\mathbf{y} = \mathbf{x} - \mathbf{z}$. Then the Law of Cosines states that

$$\|\mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{z}\|^2 - 2 \|\mathbf{x}\| \|\mathbf{z}\| \cos \theta$$

(see fig. 3). By replacing $\|\mathbf{x}\| \|\mathbf{z}\| \cos \theta$ with the inner product $\mathbf{x} \cdot \mathbf{z}$ and \mathbf{y} with $\mathbf{x} - \mathbf{z}$ we get the Euclidean distance

$$\|\mathbf{x} - \mathbf{z}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{z}\|^2 - 2 \mathbf{x} \cdot \mathbf{z} \cos \theta$$

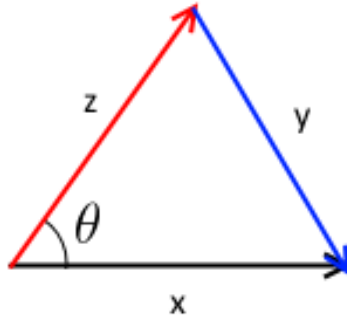


Fig. 3: **TODO redraw in vector** The Law of Cosine: $\|\mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{z}\|^2 - 2 \|\mathbf{x}\| \|\mathbf{z}\| \cos \theta$

In Linear Algebra, however, the inner product (also called “Dot product”) is defined differently as a sum of element-wise products of two vectors: given two vectors \mathbf{x} and \mathbf{z} , the inner product is $\mathbf{x}^T \mathbf{z} = \sum_{i=1}^n x_i z_i$ where x_i and z_i are i th elements of \mathbf{x} and \mathbf{z} , respectively. The geometric and algebraic definitions are equivalent [30].

Thus, for two documents \mathbf{d}_1 and \mathbf{d}_2 , the more terms these documents have in common, the bigger their dot product is, and the more similar they are.

However the magnitude of each individual vector still matters. If one document vector is particularly long compared to other vectors and it is not orthogonal to them (i.e. $\cos \theta \neq 0$), then it is likely to be one of the most similar vectors to others – only because of its length.

Cosine Similarity

The important drawback of the Inner product is that it is very sensitive to lengths of the vectors. Therefore it may make sense to consider only the angle between them: the angle does not depend on the magnitude, but still acts as a very good indicator of vectors being similar or not.

The angle between two vectors can be calculated from the geometric definition of inner product. Recall that given two vectors \mathbf{x} and \mathbf{z} , the inner product is defined as $\mathbf{x} \cdot \mathbf{z} = \|\mathbf{x}\| \|\mathbf{z}\| \cos \theta$. By rearranging the terms we get

$$\cos \theta = \frac{\mathbf{x} \cdot \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|}$$

and thus

$$\theta = \arccos \frac{\mathbf{x} \cdot \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|}$$

Note, however, that it is not necessary to invert the cosine: the angle θ between two vectors ranges from 0 to π , and cosine is monotonically decreasing on this interval. Thus, $\cos \theta \in [-1, 1]$, and the smaller the angle between two vectors, the bigger the cosine of this angle. Document vectors typically do not have negative components (all weights are positive), and thus, the angle between two vectors can be at most π , and therefore $\cos \theta \in [0, 1]$.

Using this, we define *cosine similarity* between two documents \mathbf{d}_1 and \mathbf{d}_2 as

$$\text{cosine}(\mathbf{d}_1, \mathbf{d}_2) = \frac{\mathbf{d}_1 \cdot \mathbf{d}_2}{\|\mathbf{d}_1\| \|\mathbf{d}_2\|}$$

If the documents have unit lengths, then cosine similarity is the same as dot product:

$$\text{cosine}(\mathbf{d}_1, \mathbf{d}_2) = \frac{\mathbf{d}_1 \cdot \mathbf{d}_2}{\|\mathbf{d}_1\| \|\mathbf{d}_2\|} = \mathbf{d}_1 \cdot \mathbf{d}_2$$

. Thus we can unit-normalize all document vectors (i.e. compute $\mathbf{d}' = \mathbf{d}/\|\mathbf{d}\|$ for all $\mathbf{d} \in \mathcal{D}$) and then simply use the dot product to compute similarity between two documents. This normalization is often called “cosine normalization” in Information Retrieval.

Cosine similarity gives the similarity score between two document vectors. What if instead we need a distance function? The maximal possible cosine is 1 for two identical documents. Therefore we can define *cosine distance* between two vectors \mathbf{d}_1 and \mathbf{d}_2 as

$$d_c(\mathbf{d}_1, \mathbf{d}_2) = 1 - \text{cosine}(\mathbf{d}_1, \mathbf{d}_2)$$

The cosine distance is not a proper metric [31], but it is nonetheless useful.

Finally, there is a connection between the cosine distance and the Euclidean distance [31]. Consider two unit-normalized vectors $\mathbf{d}'_1 = \mathbf{d}_1 / \|\mathbf{d}_1\|$ and $\mathbf{d}'_2 = \mathbf{d}_2 / \|\mathbf{d}_2\|$. The Euclidean distance between them is

$$\|\mathbf{d}'_1 - \mathbf{d}'_2\|^2 = \|\mathbf{d}'_1\|^2 - 2\mathbf{d}'_1{}^T \mathbf{d}'_2 + \|\mathbf{d}'_2\|^2 = 2 - 2\mathbf{d}'_1{}^T \mathbf{d}'_2$$

Since the vectors are unit-normalized, we know that $\text{cosine}(\mathbf{d}'_1, \mathbf{d}'_2) = \mathbf{d}'_1{}^T \mathbf{d}'_2$, so we have

$$\|\mathbf{d}'_1 - \mathbf{d}'_2\|^2 = 2(1 - \text{cosine}(d'_1, d'_2)) = 2d_c(d'_1, d'_2)$$

Thus we can use Euclidean distance on unit-normalized vectors and interpret it as Cosine distance.

Jaccard Coefficient

Jaccard similarity, jaccard distance

Other Similarity Measures

Other commonly used similarity measures include:

- Dice Coefficient: $\frac{d_1^T d_2}{\|\mathbf{d}_1\|^2 + \|\mathbf{d}_2\|^2}$
- Pearson correlation coefficient
- etc

3.6 Inverted Index

In databases, indexing techniques are used to make queries faster and avoid full table scan. Inverted Index is an Information Retrieval technique for achieving better retrieval speed by reducing the number of documents that the system needs to examine [20]. This technique is also useful for other Text Mining tasks such as computing document-document similarity matrix, where each document is compared to every other document in the document collection.

The main idea is to use the fact that documents usually contain only a small portion of terms. For example, a number of distinct terms in the corpus

may be 10^5 , but typically individual documents do not contain more than 500 distinct terms. Thus, when these documents have very sparse document vectors, and many similarity functions used for documents ignore zero entries (e.g. cosine similarity or Jaccard coefficient).

For the cosine to be non-zero, two documents need to share at least one term. Therefore to find documents similar to document d , we use the index to restrict the list of documents to compare to those that have at least one term in common.

TODO: Describe the algorithm of building the index and using it for the similarity search 1.

Algorithm 1 Inverted Index for similarity computations

```

function BUILD-INDEX(documents  $\mathcal{D}$ )
   $V \leftarrow$  all distinct terms from  $\mathcal{D}$ 
  for each  $w_i \in V$  do
     $\text{index}[w_i] \leftarrow \emptyset$ 
  for each  $d \in \mathcal{D}$  do
    for each  $w_i \in d$  do
       $\text{index}[w_i] \leftarrow \text{index}[w_i] \cup \{d\}$ 
  return index

function RETURN-CANDIDATES(document  $d$ , index)
  for each  $w_i \in d$  do
     $D_i \leftarrow \text{index}[w_i] - d$   $\triangleright$  all documents containing  $w_i$  except  $d$ 
  return  $\bigcup_i D_i$ 

```

In the next section we will review some algorithms for document clustering, and Inverted Index can be used to speed up similarity computations needed for these clustering algorithms.

4 Document Clustering Techniques

Cluster analysis is a set of techniques for organizing collection of items into coherent groups. In Text Mining clustering is often used for finding topics in a collection of document. In Information Retrieval clustering is used to assist the users and group retrieved results into clusters.

There are several types of clustering algorithms:

- Hierarchical (Agglomerative and Divisive)
- Partitioning
- Density-based
- and others

In this chapter we will review some of the clustering techniques: in section 4.1 we will discuss Agglomerative clustering. We discuss K-Means, a partitioning algorithms, in section 4.2 and its extensions in section 4.3. Finally, a density-based algorithm DBSCAN is explained in section 4.4 along with its extensions in section 4.5.

4.1 Agglomerative clustering

The general idea of agglomerative clustering algorithms is to start with each document being its own cluster and iteratively merge clusters based on best pair-wise cluster similarity.

Thus, a typical agglomerative clustering algorithms consists of the following steps:

- let each document be a cluster on its own
- compute similarity between all pairs of clusters and store the results in a similarity matrix
- merge two most similar clusters
- update the similarity matrix
- repeat until everything belongs to the same cluster

These algorithms differ only in the way they calculate similarity between clusters.

It can be:

- **Single Linkage (SLINK)**: The clusters are merged based on the closest pair. It can be very efficient, but it encourages chaining behavior.

- **Complete Linkage (CLINK)**: The clusters are merged based on the worst-case similarity - the similarity between the most distant objects on the clusters. It's very expensive computationally, but it avoids chaining altogether.
- **Group-Average Linkage**: Similarity between clusters is calculated as average pair-wise similarity between all objects in the clusters, and the most similar clusters are merged.
- **Ward's Method**: The clusters to merge are chosen such that the within-cluster error (e.g. sum of squares) between each object and its centroid is minimized.

Among these algorithms only SLINK is computationally feasible for large data sets, but it doesn't give good results compared to other agglomerative clustering algorithms. Additionally, these algorithms are not always good for document clustering because they tend to make mistakes at early iterations that are impossible to correct afterwards [32].

4.2 K-Means

Unlike agglomerative clustering algorithms, K-Means is an iterative algorithm, which means that it can correct the mistakes made at earlier iterations.

Lloyd's algorithm is the most popular way of implementing K-Means [33]: given a desired number of clusters k , it iteratively improves the Euclidean distance between each data point and the centroid, closest to it.

Let $\mathcal{D} = \{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_n\}$ be the document collection, where documents \mathbf{d}_i are represented in some document vector space \mathbb{R}^m and k is the desired number of clusters. Then we define k cluster centroids $\boldsymbol{\mu}_j$ that are also in the same document vector space \mathbb{R}^m . Additionally for each document \mathbf{d}_i we maintain the assignment variable $c_i \in \{1, 2, \dots, k\}$, which specifies to what cluster centroid $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k$ the document \mathbf{d}_i belongs.

The algorithm consists of three steps: (1) seed selection step, where each $\boldsymbol{\mu}_j$ is randomly assigned some value, (2) cluster assignment step, where we iterate over all document vectors \mathbf{d}_i and find its closest centroid, and (3) move centroids step, where the centroids are re-calculated. Steps (2) and (3) are repeated until the algorithm converges. The pseudocode for K-Means is presented in the listing 2.

Usually, K-Means shows very good results for document clustering, and in several studies it (or its variations) shows the best performance [34] [32].

Algorithm 2 Lloyd’s algorithm for K -Means

```

function K-MEANS(no. clusters  $k$ , documents  $\mathcal{D}$ )
  for  $j \leftarrow 1 \dots k$  do                                     ▷ random seed selection
     $\mu_j \leftarrow \text{random } \mathbf{d} \in \mathcal{D}$ 
  while not converged do
    for each  $\mathbf{d}_i \in \mathcal{D}$  do                                     ▷ cluster assignment step
       $c_i \leftarrow \arg \min_j \|\mathbf{d}_i - \mu_j\|^2$ 
    for  $j \leftarrow 1 \dots k$  do                                     ▷ move centroids step
       $\mathcal{C}_j \leftarrow \{\mathbf{d}_i \text{ s.t. } c_i = j\}$ 
       $\mu_j \leftarrow \frac{1}{|\mathcal{C}_j|} \sum_{\mathbf{d}_i \in \mathcal{C}_j} \mathbf{d}_i$ 
  return ( $c_1, c_2, \dots, c_n$ )

```

However for large document collections Lloyd’s classical K -Means takes a lot of time to converge. The problem is caused by the fact that it goes through the entire collection many times. Mini-Batch K -Means [35] uses Mini-Batch Gradient Descent method, which is a different optimization technique that converges faster. The pseudocode for Mini-Batch K -Means is presented in listing 3.

Algorithm 3 MiniBatch K -Means

```

function MINIBATCH-K-MEANS(no. clusters  $k$ , no. iterations  $t$ , batch size  $b$ , documents  $\mathcal{D}$ )
  for  $j \leftarrow 1 \dots k$  do                                     ▷ random initialization
     $\mu_j \leftarrow \text{random } \mathbf{d} \in \mathcal{D}$ 
  repeat  $t$  times
     $\mathcal{M} \leftarrow b$  random examples from  $\mathcal{D}$ 
    for each  $\mathbf{d}_i \in \mathcal{M}$  do
       $\text{centroids}[\mathbf{d}_i] \leftarrow \arg \min_j \|\mathbf{d}_i - \mu_j\|^2$            ▷ cache the centroid nearest to  $\mathbf{d}_i$ 
    for each  $\mathbf{d}_i \in \mathcal{M}$  do
       $c_i \leftarrow \text{centroids}[\mathbf{d}_i]$                                ▷ the centroid index of document  $\mathbf{d}_i$ 
       $v[c_i] \leftarrow v[c_i] + 1$                                ▷ counts per centroid  $c_i$ 
       $\eta \leftarrow 1/v[c_i]$                                      ▷ per-centroid learning rate
       $\mu_{c_i} \leftarrow (1 - \eta) \cdot \mu_{c_i} + \eta \cdot \mathbf{d}_i$        ▷ gradient step
  until converged
  return ( $c_1, c_2, \dots, c_n$ )

```

Note that K means uses Euclidean distance, and Euclidean distance does not always behave well in high-dimensional sparse vector spaces like Document VSMs (see section 3.5). However, as discussed in section 3.5, if document vectors are normalized, the Euclidean distance and Cosine dis-

tance are related, and therefore Euclidean K -means is the same as “Cosine Distance” K -Means.

K -Means is the most popular clustering algorithms and there are many extensions to this algorithm. In the next section we will discuss some of the extensions related to document clustering.

4.3 Extensions of K -Means

There are several extensions of K -Means.

For example, Bisecting K -Means [32] is a combination of partitioning and hierarchical (divisive) algorithms. It’s a variant of K -Means that gradually splits the document space in halves until the desired number of clusters is obtained. Bisecting K -Means can achieve good performance while giving the user additional information about ... ?

Algorithm:

- start with a single cluster
- choose a cluster to split (for example, the largest one)
- apply K -means to this cluster with $K = 2$ to split it
- repeat until have desired number of clusters

TODO: pseudocode

Scatter/Gather is another popular variation of K -means, but initially used for clustering retrieved documents for Information Retrieval systems [36]. This variation includes: special smart seed selection procedure (applying hierarchical cluster on a subset of document vectors to initialize the centroids at the initialization step) and several cluster refinement operations. Additionally, in Scatter/Gather cluster centroids are concatenations of all terms in the cluster documents, not a mean value; and the cosine similarity is used instead of Euclidean distance.

There are two cluster refinement operations: split and join.

The split operation is used to continue splitting clusters, and it’s applied only to the clusters that are not coherent enough. Essentially, the split operation splits the non-coherent clusters in the same way as Bisecting K -Means. The coherence is measured via *self-similarity* of a cluster, which is the mean similarity of all documents in the cluster to its centroid, or the mean pair-wise similarity between all documents of the cluster.

The join operation merges the clusters that are very similar to each other. The similarity is measured by computing “typical” terms for each cluster (usually the most frequent terms of the centroid) and examining which clusters have significant overlaps between their typical terms.

However, when there are many documents, the centroids tend to contain a lot of words, which leads to a significant slowdown. A solution to this problem is a center adjustment method, called vector average dumping **TODO** [37]. Alternatively, some terms of the centroid can be truncated. There are several possible ways of truncating the terms: for example, we can keep only the top c terms, or remove the least frequent words such that at least 90% (or 95%) of the original vector norm is retained [38].

4.4 DBSCAN

DBSCAN is a clustering algorithm that can discover clusters of complex shapes based on the density of data points [39].

The *density* associated with a data point is obtained by counting the number of points in a region of specified radius ε around the point. If a point has a density of at least some user defined threshold MinPts , then it is considered a *core point*. The clusters are formed around these core points, and if two core points are within the radius ε , then they belong to the same cluster. If a point is not a core point itself, but it belongs to the neighborhood of some core point, then it is a *border point*. But if a point is not a core point and it is not in the neighborhood of any other core point, then it does not belong to any cluster and it is considered *noise*.

DBSCAN works as follows: it selects an arbitrary data point p , and then finds all other points in ε -neighborhood of p . If there are more than MinPts points around p , then it is a core point, and it is considered a cluster. Then the process is repeated for all points in the neighborhood, and they all are assigned to the same cluster, as p . If p is not a core point, but it has a core point in its neighborhood, then it’s a border point and it is assigned to the same cluster and the core point. But if it is a noise point, then it is marked as noise or discarded.

The DBSCAN algorithm uses the Euclidean distance, but can be adapted to use any other distance or similarity function. For example, to modify the algorithm to use the Cosine similarity (or any other similarity function) the **REGION-QUERY** has to be modified to return $\{x : \text{similarity}(x, p) \geq \varepsilon\}$.

The details of implementation of **REGION-QUERY** are not specified, and it can be implemented differently. For example, it can use Inverse Index (see

Algorithm 4 DBSCAN

```

function DBSCAN(database  $\mathcal{D}$ , radius  $\varepsilon$ , MinPts)
  result  $\leftarrow \emptyset$ 
  for all  $p \in \mathcal{D}$  do
    if  $p$  is visited then
      continue
    mark  $p$  as visited
     $\mathcal{N} \leftarrow \text{REGION-QUERY}(p, \varepsilon)$   $\triangleright \mathcal{N}$  is the neighborhood of  $p$ 
    if  $|\mathcal{N}| < \text{MinPts}$  then
      mark  $p$  as NOISE
    else
       $\mathcal{C} \leftarrow \text{EXPAND-CLUSTER}(p, \mathcal{N}, \varepsilon, \text{MinPts})$ 
      result  $\leftarrow \text{result} \cup \{\mathcal{C}\}$ 
  return result

function EXPAND-CLUSTER(point  $p$ , neighborhood  $\mathcal{N}$ , radius  $\varepsilon$ , MinPts)
   $\mathcal{C} \leftarrow \{p\}$ 
  for all  $x \in \mathcal{N}$  do
    if  $x$  is visited then
      continue
    mark  $x$  as visited
     $\mathcal{N}_x \leftarrow \text{REGION-QUERY}(x, \varepsilon)$   $\triangleright \mathcal{N}_x$  is the neighborhood of  $x$ 
    if  $|\mathcal{N}_x| \geq \text{MinPts}$  then
       $\mathcal{N} \leftarrow \mathcal{N} \cup \mathcal{N}_x$ 
   $\mathcal{C} \leftarrow \mathcal{C} \cup \{x\}$ 
  return  $\mathcal{C}$ 

function REGION-QUERY(point  $p$ , radius  $\varepsilon$ )
  return  $\{x : \|x - p\| \leq \varepsilon\}$   $\triangleright$  all points within distance  $\varepsilon$  from  $p$ 

```

section 3.6, and listing 1 for the pseudocode) to make the similarity search faster.

4.5 Extensions of DBSCAN

As discussed, DBSCAN can be extended to use any distance or similarity function. Shared Nearest Neighbors Similarity (SNN Similarity) [29] is a special similarity function that is particularly useful for high-dimensional spaces. And this similarity function is applicable to document clustering and topic discovery [40].

SNN Similarity is specified in terms of the K nearest neighbors. Let $\text{NN}_{K,\text{sim}}(p)$ be a function that returns top K closest points of p according to some similarity function sim . Then the SNN similarity function is defined as

$$\text{snn}(p, q) = |\text{NN}_{K,\text{sim}}(p) \cup \text{NN}_{K,\text{sim}}(q)|$$

The extension of DBSCAN that uses the SNN Similarity is called SSN Clustering algorithm. The user needs to specify the SSN similarity function by setting parameter K and choosing the base similarity function $\text{sim}(\cdot, \cdot)$ (typically Cosine, Jaccard or Euclidean). The algorithm itself has the same parameters as DBSCAN: radius ε (such that $\varepsilon < K$) and the core points density threshold MinPts . The **REGION-QUERY** function is modified to return $\{q : \text{snn}(p, q) \geq \varepsilon\}$. For pseudocode, see the listing 5.

Algorithm 5 SNN Clustering Algorithm

```

function SNN-CLUSTER(database  $\mathcal{D}$ ,  $K$ , similarity function  $\text{sim}$ , radius  $\varepsilon$ ,  $\text{MinPts}$ )
  for all  $p \in \mathcal{D}$  do                                      $\triangleright$  Pre-compute the  $K$ NN lists
     $\text{NN}[p] \leftarrow \text{NN}_{K,\text{sim}}(p)$ 
  for all  $(p, q) \in (\mathcal{D} \times \mathcal{D})$  do                        $\triangleright$  Pre-compute the SNN similarity matrix
     $A[p, q] \leftarrow |\text{NN}[p] \cup \text{NN}[q]|$ 
  return DBSCAN( $A, \varepsilon, \text{MinPts}$ )

```

The algorithm's running time complexity is $O(n^2)$ time, where $n = |\mathcal{D}|$, but it can be sped up by using the Inverted Index.

5 Discovering Latent Semantics

In chapter 3 we have discussed the lexical variability and ambiguity problems in natural language: synonymy and polysemy. We can treat these problems as “statistical noise” and apply dimensionality reduction techniques to find the optimal dimensionality for the data and thus reduce the amount of noise there. In this chapter we discuss two approaches for doing this: Latent Semantic Analysis in section 5.1 and Non-Negative Matrix Factorization in section 5.2.

5.1 Latent Semantic Analysis

The Vector Space Model discussed in section 3.3. The solution to this problem is to assume that there exists some optimal document vector space where the document vectors do not suffer from the ... This vector space can be found by finding the best k -rank approximation to the Term-Document matrix using Singular Value Decomposition (SVD). This technique is called Latent Semantic Analysis [41] or Latent Semantic Indexing in the context of Information Retrieval [22]. It is also a popular Text Mining technique for reducing the dimensionality of text data and it is often used for document clustering [26] [42].

There are three major steps in Latent Semantic Analysis [43]:

- Preprocess documents
- Construct a Term-Document matrix D using the Vector Space Model
- Reduce dimensionality of D by using SVD

The first two steps are the same as for traditional Vector Space Models: consider a set of document $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$ with the vocabulary $\mathcal{V} = \{t_1, t_2, \dots, t_m\}$, then D is a $m \times n$ Term-Document Matrix. If the matrix D has rank r , then the SVD of D is $D = U\Sigma V^T$, where:

- U is an $m \times r$ orthogonal matrix, i.e. $UU^T = I$;
- Σ is a diagonal $r \times r$ matrix with singular values ordered by their magnitude;
- V is an $n \times r$ orthogonal matrix, $VV^T = I$.

The dimensionality reduction is done by finding the best k -rank approximation of D , which is obtained by keeping only the first k singular values of Σ and setting the rest to 0. Typically, not only Σ is truncated, but also U

and V , and therefore, the k -rank approximation of D using SVD is written as $D \approx D_k = U_k \Sigma_k V_k^T$ where U_k is an $m \times k$ matrix with first k columns of U , Σ_k is an $k \times k$ diagonal matrix with singular values, and V_k is an $n \times k$ matrix with first k columns of V . This decomposition is called *rank-reduced* SVD and when applied to text data it reveals the “true” latent semantic space. The parameter k corresponds to the number of “latent concepts” in the data.

Let us illustrate LSA with an example (from [22] and [41]). Consider the following nine documents:

- c_1 : “Human machine interface for ABC computer applications”
- c_2 : “A survey of user opinion of computer system response time”
- c_3 : “The EPS user interface management system”
- c_4 : “System and human system engineering testing of EPS”
- c_5 : “Relation of user perceived response time to error measurement”
- m_1 : “The generation of random, binary, ordered trees”
- m_2 : “The intersection graph of paths in trees”
- m_3 : “Graph minors IV: Widths of trees and well-quasi-ordering”
- m_4 : “Graph minors: A survey”

Let the vocabulary be $\mathcal{V} = \{ \text{human, interface, computer, user, system, response, time, EPS, survey, trees, graph, minors} \}$. Then, the term-document matrix D with terms weighed by TF is

$$D = \begin{bmatrix} & c_1 & c_2 & c_3 & c_4 & c_5 & m_1 & m_2 & m_3 & m_4 \\ \text{human} & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ \text{interface} & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{computer} & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{user} & 0 & 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ \text{system} & 0 & 1 & 1 & 2 & 0 & 0 & 0 & 0 & 0 \\ \text{response} & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ \text{time} & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ \text{EPS} & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ \text{survey} & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ \text{trees} & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ \text{graph} & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ \text{minors} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

We can observe that the row vectors for “human” and “user” are orthogonal: their inner produce is zero, but these terms are similar and the inner product should be positive.

Let us apply SVD to find the best rank-2 approximation for D : $D \approx D_2 = U_2 \Sigma_2 V_2^T$. The rank-2 approximation D_2 is

$$D_2 = \begin{bmatrix} & c_1 & c_2 & c_3 & c_4 & c_5 & m_1 & m_2 & m_3 & m_4 \\ \text{human} & 0.16 & 0.4 & 0.38 & 0.47 & 0.18 & -0.05 & -0.12 & -0.16 & -0.09 \\ \text{interface} & 0.14 & 0.37 & 0.33 & 0.4 & 0.16 & -0.03 & -0.07 & -0.1 & -0.04 \\ \text{computer} & 0.15 & 0.51 & 0.36 & 0.41 & 0.24 & 0.02 & 0.06 & 0.09 & 0.12 \\ \text{user} & 0.26 & 0.84 & 0.61 & 0.7 & 0.39 & 0.03 & 0.08 & 0.12 & 0.19 \\ \text{system} & 0.45 & 1.23 & 1.05 & 1.27 & 0.56 & -0.07 & -0.15 & -0.21 & -0.05 \\ \text{response} & 0.16 & 0.58 & 0.38 & 0.42 & 0.28 & 0.06 & 0.13 & 0.19 & 0.22 \\ \text{time} & 0.16 & 0.58 & 0.38 & 0.42 & 0.28 & 0.06 & 0.13 & 0.19 & 0.22 \\ \text{EPS} & 0.22 & 0.55 & 0.51 & 0.63 & 0.24 & -0.07 & -0.14 & -0.2 & -0.11 \\ \text{survey} & 0.1 & 0.53 & 0.23 & 0.21 & 0.27 & 0.14 & 0.31 & 0.44 & 0.42 \\ \text{trees} & -0.06 & 0.23 & -0.14 & -0.27 & 0.14 & 0.24 & 0.55 & 0.77 & 0.66 \\ \text{graph} & -0.06 & 0.34 & -0.15 & -0.3 & 0.2 & 0.31 & 0.69 & 0.98 & 0.85 \\ \text{minors} & -0.04 & 0.25 & -0.1 & -0.21 & 0.15 & 0.22 & 0.5 & 0.71 & 0.62 \end{bmatrix}$$

What is the effect of dimensionality reduction here? The frequencies of words have changed, and some of the entries even become negative. Consider two cells: (“survey”, m_4) and (“trees”, m_4). In the original document we have $D[\text{survey}, m_4] = 1$ and $D[\text{trees}, m_4] = 0$, but in the reduced space we have $D_2[\text{survey}, m_4] = 0.42$ and $D_2[\text{trees}, m_4] = 0.66$. Note that the count for “survey” went down while the count for “trees” went up. The reason for this is that the document m_4 contains “graph” and “minor”, which are graph related, so another graph related word “trees” got additional score. On the other hand the term “survey” was not expected in this context, so the count went down.

Additionally, we can compute the similarity in the reconstructed space as well. For example, the terms “user” and “human” are no longer orthogonal and have positive dot product. Thus, it tells us that the terms are similar even though they never co-occur together in the original input space.

LSA can be used for clustering as well, and this is usually done by first transforming the document space to the LSA space and then doing applying transitional cluster analysis techniques there [38].

However these is not need to reconstruct the rank-reduced matrix to apply clustering, and in many cases it is not possible: the original input space is very sparse, but the rank-reduced reconstructed matrix becomes very dense. Therefore we do not reconstruct the entire matrix, but instead

keep only the low dimensional representation $V_k \Sigma_k$, which is enough for many clustering algorithms.

For example, consider the inner product. Document-document similarity in the original space is calculated as DD^T (the columns of D are the document $\mathbf{d}_1, \dots, \mathbf{d}_n$), and by applying SVD we have $D^T D = V \Sigma^T U^T U \Sigma V^T = V \Sigma^T \Sigma V^T = (V \Sigma)(V \Sigma)^T$. Thus, to compute the similarity between document \mathbf{d}_i and \mathbf{d}_j we compute the inner product between i th and j th rows of $V \Sigma$. Likewise, to compute the similarity between documents i and j in the reduced representation, we compute the inner product between the respective rows of $V_k \Sigma_k$.

Additionally, the Euclidean distance in the reduced space can also be computed directly on the rows of $V_k \Sigma_k$. Recall that the Euclidean distance can be expressed as an inner product $\|\mathbf{d}_i - \mathbf{d}_j\|^2 = \mathbf{d}_i^T \mathbf{d}_i - 2 \mathbf{d}_i^T \mathbf{d}_j + \mathbf{d}_j^T \mathbf{d}_j$, and since we know how to compute the inner product in the semantic space, we can compute the distance in this space as well by using the rows of $V_k \Sigma_k$.

This means that we can apply any clustering algorithm, including K -means, on the rows of $V_k \Sigma_k$ and without having to reconstruct the entire term-document matrix.

A generic LSA-based clustering algorithm therefore consists of the following steps:

- Build a term-document matrix D from the document collection
- Select number of latent concepts k and apply rank-reduced SVD on D to get $V_k \Sigma_k$
- Apply the cluster algorithm on the rows of $V_k \Sigma_k$

LSA has some drawbacks. Because SVD looks for orthogonal basis for the new document space, there are negative values that are harder to interpret. Additionally, with negative values in the reconstructed space can cause the cosine to take negative values as well. However, it does not affect significantly the properties of the cosine distance: it still will always give the results larger than 0. This problem can be solved by using a different matrix decomposition technique instead of SVD, and we discuss one of them in the next section.

5.2 Non-Negative Matrix Factorization

Apart from SVD there are many other different matrix decomposition techniques that can be applied for document clustering and for discovering the latent structure of the term-document matrix [44], and one of them is Non-Negative Matrix Factorization (NMF) [45].

NMF is a matrix decomposition technique. When it is applied to non-negative data, NMF produces non-negative rank-reduced approximations. Since term-document matrices do not have negative values, it makes NMF a good candidate to replace SVD in LSA. The main conceptual difference between SVD and NMF is that SVD looks for orthogonal directions to represent document space, while NMF does not require orthogonality. As the result, SVD often produces semantic spaces with negative values, but NMF does not [46] (see fig. 4).

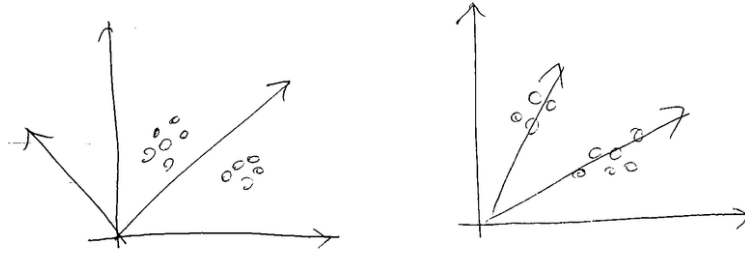


Fig. 4: **TODO redraw** Directions found by SVD (on the left) vs directions by NMF (on the right)

The NMF of an $m \times n$ term-document matrix D is $D \approx D_k = UV^T$ where U is an $m \times k$ matrix, V is an $n \times k$ matrix and k is the number of semantic concepts in D . Non-negativity of elements in D_k is very good for interpretability: it ensures that documents can be seen as a non-negative combination of the key concepts.

Additionally, NMF is useful for clustering: the results of NMF can be directly interpreted as cluster assignment and there is no need to use separate clustering algorithms [46].

What is more, NMF is a co-clustering algorithms: it clusters both rows of D and columns of D at the same time. For a term-document matrix $D \approx UV^T$, where U defines the reduced vector space for terms and V defines the reduced vector space for documents. Since all elements are non-negative, it can have the following interpretation: elements $(U)_{ij}$ of U represent the degree to which terms i belongs to cluster j , and elements $(V)_{ij}$ represent the degree to which document i belongs to cluster j .

The document clustering using NMF consists of the following steps [46]:

- Construct the term-document matrix D ;
- Perform NMF on D to get U and V ;
- Normalize rows \mathbf{v}_i of V by using $\mathbf{v}'_i = \mathbf{v}_i \cdot \|\mathbf{u}_i\|$ and rows \mathbf{u}_i of U with $\mathbf{u}'_i = \mathbf{u}_i / \|\mathbf{u}_i\|$;
- To determine the cluster assignment for document \mathbf{d}_i , examine \mathbf{v}'_i (the i th row of V) and find the largest component of this vector. That is, i th document belongs to cluster x if $x = \arg \max_j v_{ij}$ where v_{ij} are components of \mathbf{v}_i ;
- If the desired number of clusters K is larger than the rank k of the reduced matrix D_k , the clustering can be performed directly on the rows of V , for example, by using K -Means.

6 Implementation

In section 6.1 we describe the data set that we use, then we describe how we extract identifiers from this dataset (section 6.2) and how this dataset is cleaned (section 6.3). Next, the implementation of clustering algorithms is described in the section 6.3. After the clusters are found, we combine them into a hierarchy in the section 6.5.

Finally, in the section 6.6 we explore how the same set of techniques can be applied to source code in Java.

6.1 Data set

In this work we apply the discussed techniques to the English version of Wikipedia [47]. It's a big web encyclopedia where articles are written and edited by the community. For our work wikipedia is interesting because there are many math pages.

At present (July 12, 2015) English wikipedia contains about 4.9 mln articles² and it is 1.5 Gb in the compressed form. However, only a small portion of these articles are math related: only about 30.000 pages contain at least one `<math>` tag.

The math information is enriched with semantic information by MediaWiki and we use this augmented data representation. 30.000 math pages with augmented math tags occupy around 1.5 Gb in uncompressed form.

Apart from the text data and formulas wikipedia articles have information about categories, which can also be exploited. It is hard to extract category information from the raw wikipedia mark up, but this information is available in a structured form in DBpedia [48].

6.2 Definition Extraction

Before we can proceed to discovering identifier namespaces, we need to extract identifier-definition relations. For this we use the probabilistic approach, discussed in the section 2.3. The extraction process is implemented using Apache Flink [49].

But before the original dataset can be preprocessed, it is enriched with augmented MathML (see section 2.1), and then the dataset is filtered such that only articles with the math tag are retained.

In the wikipedia dataset each document is represented using wiki XML. It makes it easy to extract title and content, and then, the all the formulas

² <https://en.wikipedia.org/wiki/Wikipedia:Statistics>

are extracted from the content. The formulas are extracted by looking for `<math>` tags. However some formulas are typed without the tag, but only with unicode characters. Such formulas are not easy to detect and therefore in this work we choose not to process them. Once all `<math>` tags are found, they (along with the content) are replaced with a special placeholder `FORMULA_%HASH%`, where `%HASH%` is MD5 hash [50] of the tag’s content represented as a hexadecimal string. After that the content of the tags is kept separately from the document content.

Once formulas are retrieved, we extract the definitions from them. We are not interested in the semantics of a formula, only in the identifiers it contains. Hence we need only to look for all `<ci>` tags. There are two types of identifiers: simple identifiers such as “ t ”, “ C ”, “ μ ”; and complex identifiers with subscripts such as “ x_1 ”, “ ξ_i ” or even “ β_{slope} ”. We do not process superscripts because they are usually powers (for example, x^2), and therefore they are not interesting for this work. There are exceptions to this, for example, “ σ^2 ” is an identifier, but these cases are rare and can be ignored.

Since MathML is XML, the identifiers are extracted with XPath queries:

- `//m:mi[not(ancestor::m:msub)]/text()` for all `<ci>` tags that are not subscript identifiers
- `//m:msub` for subscript identifiers

Once the identifiers are extracted, the rest of the formula is discarded. As the result, we have a “Bag of Formulas”.

The content of a wiki document is structured and authored with a special markup language for specifying document layout elements such as headers, lists, text formatting and tables: Wiki markup. Thus the next step is to process the Wiki markup and extract the textual content of an article, and this is done using a Java library “Mylyn Wikitext” [51]. Almost all annotations are discarded at this stage, and only inner-wiki links are kept: they can be useful as candidate definitions.

Once the markup annotations are removed and the text content of an article is extracted, we then apply Natural Language Processing (NLP) techniques. Thus, the next step is the NLP step, and for NLP we use StanfordNLP [13]. The first part at this stage is to tokenize the text and also split it by sentences. Once it is done, we then apply Math-aware POS tagging (see section 2.2). For identifiers and math formulas we introduce two new POS classes: “ID” and “MATH”, respectively. These classes are not a part of the standard Penn Treebank POS Scheme [12] used by StanfordNLP, therefore we need to label all the instances of these tags ourselves during the

additional post-processing step. If a token starts with “**FORMULA_**”, then we recognize that it is a placeholder for a math formula, and therefore we annotate it with the “**MATH**” tag. Additionally, if this formula contains only one identifier, this placeholder token is replaced by the identifier and it is tagged with “**ID**”. Additionally, we keep track of all identifiers found in the document and then for each token we check if this token is in the list. If it is, then it is re-annotated with “**ID**” as well.

At the Wikipedia markup processing step we discard almost all markup annotations, but keep only inter-wiki links, because these links are good definition candidates. To use them, we introduce another POS Tag: “**LINK**”. To detect all inner-wiki links, we first find all token subsequences that start with `[[` and end with `]]`. Then these subsequences are concatenated and tagged as “**LINK**”.

Also we are interested in all sequences of successive nouns (both singular and plural) possibly modified by an adjective. We concatenate all such sequences into one token tagged with “**NOUN_PHRASE**”.

Next we select the most probably identifier-definition pairs. At this stage we are interested only in tokens annotated with “**LINK**” and “**NOUN_PHRASE**”: these tokens are definition candidates, and we rank each token by a score that depends how far it is from the identifier of interest and how far is the closest formula that contains this identifier (see section 2.3). The output of this step is a list of identifier-definition pairs along with the score. Only pairs with scores above the user specified threshold are retained.

The following is the list of the most common identifier-definition pairs:

- t : “time” (1086)
- m : “mass” (424)
- θ : “angle” (421)
- T : “temperature” (400)
- r : “radius” (395)
- v : “velocity” (292)
- ρ : “density” (290)
- G : “group” (287)
- V : “volume” (284)
- λ : “wavelength” (263)
- R : “radius” (257)
- n : “degree” (233)
- r : “distance” (220)
- c : “speed of light” (219)
- L : “length” (216)

- n : “length” (189)
- n : “order” (188)
- n : “dimension” (185)
- n : “size” (178)
- M : “mass” (171)

6.3 Data Cleaning

The Natural Language data is famous for being noisy and hard to clean [52]. The same is true for mathematical identifiers and scientific texts with formulas. In this section we describe how the data was preprocessed and cleaned at different stages of Definition Extraction (section 6.2).

Often identifiers contain additional semantic information visually conveyed by special diacritical marks or font features. For example, the diacritics can be hats to denote “estimates” (e.g. “ \hat{w} ”), bars to denote the expected value (e.g. “ \bar{X} ”), arrows to denote vectors (e.g. “ \vec{x} ”) and others. As for the font features, boldness is often used to denote vectors (e.g. “ \mathbf{w} ”) or matrices (e.g. “ \mathbf{X} ”), calligraphic fonts are used for sets (e.g. “ \mathcal{H} ”), double-struck fonts often denote spaces (e.g. “ \mathbb{R} ”), etc. Unfortunately there is no common notation established across all fields of mathematics and there is a lot of variance. For example, a vector can be denoted by “ \vec{x} ”, “ \mathbf{x} ” or “ \mathbf{x} ”, and a real line by “ \mathbb{R} ”, “ \mathbf{R} ” or “ \mathfrak{R} ”. Therefore we discard all this additional information, such that “ \bar{X} ” becomes “ X ”, “ \mathbf{w} ” becomes “ w ” and “ \mathfrak{R} ” becomes “ R ”.

The diacritic marks can easily be discarded because they are represented by special MathML instructions that easily can be ignored (see the section 2.1 for details). But, on the other hand, the visual features are encoded directly on the character level: the identifiers use special unicode symbols to convey font features such as boldness or Fraktur, so it needs to be normalized by converting characters from special “Mathematical Alphanumeric Symbols” unicode block [53] back to the standard ASCII positions (“Basic Latin” block).

Additionally, there is a lot of noise on the annotation level in MathML formulas: many non-identifiers are captured as identifiers inside `<ci>` tags. Among them there are many mathematic-related symbols like “ \wedge ”, “ $\#$ ”, “ \vee ”, “ \int ”; miscellaneous symbols like “ \diamond ” or “ \circ ”, arrows like “ \rightarrow ” and “ \Rightarrow ”, and special characters like “ \lceil ”.

To filter out these one-symbol false identifiers we fully exclude all characters from the following unicode blocks: “Spacing Modifier Letters”, “Mis-

cellaneous Symbols”, “Geometric Shapes”, “Arrows”, “Miscellaneous Technical”, “Box Drawing”, “Mathematical Operators” (except “ ∇ ” which is sometimes used as an identifier) and “Supplemental Mathematical Operators” [53]. Some symbols (like “=”, “+”, “~”, “%”, “?”, “!”) belong to commonly used unicode blocks which we cannot exclude altogether. For these symbols we manually prepare a stop list for filtering them.

It also captures multiple-symbol false positives: operators and functions like “sin”, “cos”, “exp”, “max”, “trace”; words commonly used in formulas like “const”, “true”, “false”, “vs”, “iff”; English stop words like “where”, “else”, “on”, “of”, “as”, “is”; units like “mol”, “dB”, “mm”. These false identifiers are excluded by a stop list as well: if a candidate identifier is in the list, it is filtered out.

Then, at the next stage, the definitions are extracted. However many shortlisted definitions are either not valid definitions or too general. For example, some identifiers become associated with “if and only if”, “alpha”, “beta”, “gamma”, which are not valid definitions. Other definitions like “element”, “number” or “variable” are valid, but they are too general and not descriptive. We maintain a stop list of such false definitions and filter them out from the result.

The next stage is using identifier/definition pairs for document clustering. We can note that if some definition is used only once throughout the entire data set, it is not useful because it does not have any discriminative power. Therefore all such definitions are excluded.

6.4 Document Clustering

At the Document Clustering stage we want to find cluster of documents that are good namespace candidates.

Before we can do this, we need to vectorize our dataset: i.e. build the Identifier Space (see section 3.4) and represent each document in this space.

There are three choices for dimensions of the Identifier space:

- identifiers alone
- “weak” identifier-definition association
- “strong” association: use identifier-definition pairs

In the first case we are only interested in identifier information and discard the definitions altogether.

In the second and third cases we keep the definitions and use them to index the dimensions of the Identifier Space. But there is some variability in

the definitions: for example, the same identifier “ σ ” in one document can be assigned to “Cauchy stress tensor” and in other it can be assigned to “stress tensor”, which are almost the same thing. To reduce this variability we perform some preprocessing: we tokenize the definitions and use individual tokens to index dimensions of the space. For example, suppose we have two pairs (σ , “Cauchy stress tensor”) and (σ , “stress tensor”). In the “weak” association case we have will dimensions (σ , Cauchy, stress, tensor), while for the “strong” association case we will have (σ _Cauchy, σ _stress, σ _tensor).

Additionally, the effect of variability can be decreased further by applying a stemming technique for each definition token. In this work we use Snowball stemmer for English [54] implemented in NLTK [55]: a python library for Natural Language Processing.

Each document is vectorized (converted to a vector form) by using `TfidfVectorizer` from scikit-learn [56]. We use the following settings:

- `use_idf=True, min_df=2`
- `use_idf=False, min_df=2`
- `use_idf=False, sublinear_tf=True, min_df=2`

In the first case we use inverse document frequency (IDF) to assign additional collection weight for “terms” (see section 3.3), while in second and in third we use only term frequency (TF). In the second case we apply a sublinear transformation to the TF component to reduce the influence of frequently occurring words. In all three cases we keep only “terms” that are used in at least two documents.

The output is a document-identifier matrix (analogous to “document-term”): documents are rows and identifiers/definitions are columns. The output of `TfidfVectorizer` is row-normalized, i.e. all rows has unit length.

Once we the documents are vectorized, we can apply clustering techniques to them. We use K -Means (class `KMeans` in scikit-learn) and Mini-Batch K -Means (class `MiniBatchKMeans`) [56]. Note that if rows are unit-normalized, then running K -Means with Euclidean distance is equivalent to cosine distance (see section 4.2).

Bisecting K -Means (see section 4.3) was implemented on top of scikit-learn: at each step we take a subset of the dataset and apply K -Means with $K = 2$ to this subset. If the subset is big (with number of documents $n > 2000$), then we use Mini-Batch K -means with $K = 2$ because it converges much faster.

Scatter/Gather, an extension to K -means (see section 4.3), was implemented manually using `scipy` [57] and `numpy` [58] because `scikit-learn`’s implementation of K -Means does not allow using user-defined distances.

DBScan (section 4.4) and SNN Clustering (section 4.5) algorithms were also implemented manually: available DBScan implementations usually take distance measure rather than a similarity measure. The similarity matrix created by similarity measures are typically very sparse, because usually only a small fraction of the documents are similar to some given document. Similarity measures can be converted to distance measures, but in this case the matrix will no longer be sparse, and we would like to avoid that. Additionally, available implementations are usually general purpose implementations and do not take advantage of the structure of the data: in text-like data clustering algorithms can be sped up significantly by using an inverted index (section 3.6)

Dimensionality reduction techniques are also important: they not only reduce the dimensionality, but also help reveal the latent structure of data. In this work we use Latent Semantic Analysis (LSA) (section 5.1) which is implemented using randomized Singular Value Decomposition (SVD) [59]. The implementation of randomized SVD is taken from `scikit-learn` [56] - method `randomized_svd`. Non-negative Matrix Factorization is an alternative technique for dimensionality reduction (section 5.2). Its implementation is also taken from `scikit-learn` [56], class `NMF`.

To assess the quality of produced clusters we use wikipedia categories. It is quite difficult to extract category information from raw wikipedia text, therefore we use DBPedia [48] for that: it provides machine-readable information about categories for each wikipedia article. Additionally, categories in wikipedia form a hierarchy, and this hierarchy is available as a SKOS ontology.

A cluster is said to be “pure” if all documents have the same category. Using categories information we can find the most frequent category of the cluster, and then we can define purity as

$$\text{purity}(C) = \frac{\max_i \text{count}(c_i)}{|C|}$$

(**TODO:** Add backlink to purity definition).

Then we can calculate the overall purity of a cluster assignment and use this to compare results of different clustering algorithms. However it is not enough just to find the most pure cluster assignment: because as the number

of clusters increases the overall purity also grows. Thus we can also optimize for the number of clusters with purity p of size at least n .

When the number of clusters increase, the purity always grows (see fig. 5), but at some point the number of pure clusters will start decreasing (see fig. 6).

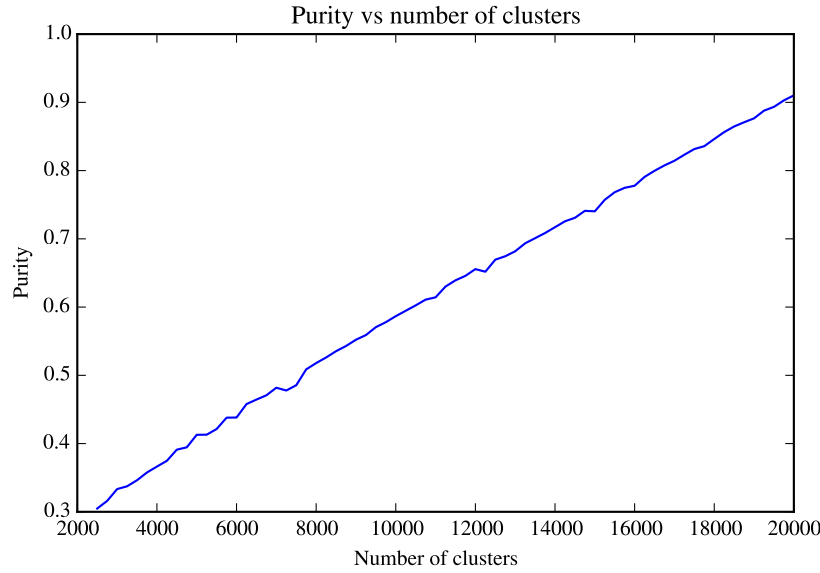


Fig. 5: K in K -Means vs overall purity of clustering: the purity increases linearly with K

6.5 Building Hierarchy

Once

AMS Mathematics Subject Classification (2010) [60] Excluded all sub-categories those code end with '99': they are usually 'Miscellaneous topics' or 'None of the above, but in this section'. top level categories 'General', 'History and biography', and 'Mathematics education' were also excluded. Additionally we exclude the following:

- Quantum theory → Axiomatics, foundations, philosophy
- Quantum theory → Applications to specific physical systems
- Quantum theory → Groups and algebras in quantum theory

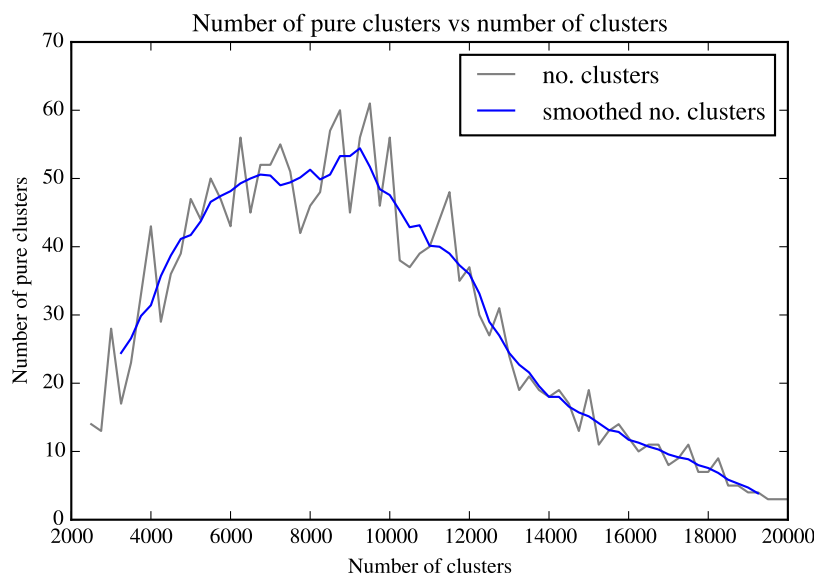


Fig. 6: K in K -Means vs the number of pure clusters: it grows initially, but after $K \approx 1000$ starts to decrease

- Partial differential equations → Equations of mathematical physics and other areas of application
- Statistics → Sufficiency and information
- Functional analysis → Other (nonclassical) types of functional analysis
- Functional analysis → Miscellaneous applications of functional analysis

So these categories do not interfere with PACS.

APS Physics and Astronomy Classification Scheme (2010) [61]

We remove the “GENERAL” top-level category. In PACS there are 3 levels of categories, but we merge all 3-rd level categories into 2nd level.

ACM Classification Scheme [62] available as a SKOS [63] ontology at their website [64]. The SKOS ontology graph was processed with RDFLib [65]

We keep the following top level categories: “Hardware”, “Computer systems organization”, “Networks”, “Software and its engineering”, “Theory of computation”, “Information systems”, “Security and privacy”, “Human-centered computing”, “Computing methodologies”.

After obtaining the data and parsing, all categories, the hierarchies are merged into one and then we try to match the found namespaces with second-level category in the hierarchy.

This is done by keywords matching: we extract all words from the category (this includes top level category name, subcategory name and all sub-sub categories concatenated). From the cluster we also extract the category information. Then we try to do keyword matching using cosine similarity between the cluster and each category. The cluster is assigned to the category with the best cosine.

If the cosine score is low (below 0.2) or there is only one keyword matched, then the cluster is assigned to the “OTHERS” category.

6.6 Java Language Processing

TODO: also refer back to the introduction

In chapter 3 we have compared the identifier namespaces with namespaces in programming languages and with packages in the Java programming language in particular. We motivated the assumption that there exist “namespace defining” groups of documents by arguing that these groups also exist in programming languages. Thus, the same set of techniques for namespace discovery should be applicable to source code as well.

If a programming language is statically typed, like Java or Pascal, usually it is possible to know the type of a variable from the declaration of this variable. Therefore we can see variable names as “identifiers” and variable types as “definitions”. Clearly, there is a difference between variable types and identifier definitions, but we believe that this comparison is valid because the type carries additional semantic information about the variable and in what context it can be used – like the definition of an identifier.

The information about variables and their types can be extracted from a source code repository, and each source file can be processed to obtain its Abstract Syntax Tree (AST). By processing the ASTs, we can extract the variable declaration information. Thus, each source file can be seen as a document, which is represented by all its variable declarations.

In this work we process Java source code, and for parsing it and building ASTs we use a library `JavaParser` [66]. The Java programming language was chosen because it requires the programmer to always specify the type information when declaring a variable. It is different for other languages when the type information is usually inferred by the compilers at compilation time.

In Java a variable can be declared in three places: as an inner class variable (or a “field”), as a method (constructor) parameter or as a local variable inside a method or a constructor. We need to process all three types of variable declarations and then apply additional preprocessing, such as converting the name of the type from short to fully qualified using the information from the import statements. For example, `String` is converted to `java.lang.String` and `List<Integer>` to `java.util.List<Integer>`, but primitive types like `byte` or `int` are left unchanged. Secondly,

Consider an example in the listing 8. There is a class variable `threshold`, a method parameter `in` and two local variables `word` and `posTag`. The following relations will be extracted from this class: (“threshold”, `double`), (“in”, `domain.Word`), (“word”, `java.lang.String`), (“posTag”, `java.lang.String`). Since all primitives and classes from packages that star with `java` are discarded, at the end the class `WordProcessor` is represented with only one relation (“in”, `domain.Word`).

Listing 8: A java class

```
package process;

import domain.Word;

public class WordProcessor {

    private double threshold;

    public boolean isGood(Word in) {
        String word = in.getWord();
        String posTag = in.getPosTag();
        return isWordGood(word) && isPosTagGood(posTag);
    }

    // ...

}
```

In the experiments we applied this source code analysis to the source code of Apache Mahout 0.10 [67], which is an open-source library for scalable Machine Learning and Data Mining.

As on 2015-07-15, this dataset consists of 1560 java classes with 45878 variable declarations. After discarding declarations from the standard Java API, primitives and types with generic parameters, only 15869 declarations were retained.

The following is top-15 variable/type declarations:

- ("conf", org.apache.hadoop.conf.Configuration), 491 times
- ("v", org.apache.mahout.math.Vector), 224 times
- ("dataModel", org.apache.mahout.cf.taste.model.DataModel), 207 times
- ("fs", org.apache.hadoop.fs.FileSystem), 207 times
- ("log", org.slf4j.Logger), 171 times
- ("output", org.apache.hadoop.fs.Path), 152 times
- ("vector", org.apache.mahout.math.Vector), 145 times
- ("x", org.apache.mahout.math.Vector), 120 times
- ("path", org.apache.hadoop.fs.Path), 113 times
- ("measure", org.apache.mahout.common.distance.DistanceMeasure), 102 times
- ("input", org.apache.hadoop.fs.Path), 101 times
- ("y", org.apache.mahout.math.Vector), 87 times
- ("comp", org.apache.mahout.math.function.IntComparator), 74 times
- ("job", org.apache.hadoop.mapreduce.Job), 71 times
- ("m", org.apache.mahout.math.Matrix), 70 times

7 Evaluation

In section 7.1 we describe how we select the best clustering algorithm.

7.1 Parameter Tuning

We have the following parameters

Way to incorporate definition information (3 ways): no identifier, soft association, hard association

Weighting schemes: TF-IDF, TF, log TF

Clustering algorithm

DBSCAN, SNN Clustering: params: min_pts, epsilon K-Means, Bisecting K-Means: params: k

Distance and similarity measures used Euclidean distance, cosine similarity, jaccard similarity, SNN Similarity

Ways to reduce dimensionality SVD $D \approx D_k = U_k \Sigma_k V_k^T$, param: k what's the rank of the approximation matrix NMF $D \approx D_k = U_k V_k^T$ param: k number of columns in U and V - also the rank of D_k

The approach for finding the best parameter set is a grid search: different combination are tries and the best result is kept.

Agglomerative: Ward linkage: takes forever never finished

Only identifiers

Usual K-Means

some clusters are useful, but most of them aren't

For example

Article	Identifiers
APL (programming language)	n, O, R
Binary search tree	O, n
Boolean satisfiability problem	O, n
Complexity	O, n
Earley parser	O, n
Heapsort	O, n, Ω
Lisp (programming language)	O, n
Priority queue	O, n
Sieve of Eratosthenes	O, n
Smoothsort	O, n
Comb sort	Ω, n, p, O
Divide and conquer algorithm	O, n, Ω
Stack (abstract data type)	O, n, t
Skip list	n, p, O
Graph minor	O, n, h
Flex lexical analyser	O, n
Gift wrapping algorithm	O, n, h
Perlin noise	n, O
Pseudo-polynomial time	n, O, m
Hirzebruch surface	O, n, m, p
Beap	n, O
Pairing heap	O, n, Ω
Cost efficiency	O, n, p

(54 documents in total) These articles appear to relate to

DBSCAN SNN

$k = 10$ dist = jaccard

$\varepsilon=7$ points MinPts=5 points

Article	Identifiers
Epsilon Eridani in fiction	M_{\odot}
Solar mass	M_{\odot}
Orders of magnitude (mass)	M_{\odot}
Carbon-burning process	M_{\odot}
Baryonic dark matter	M_{\odot}
PSR J16142230	M_{\odot}
KennicuttSchmidt law	M_{\odot}
Portal:Star/Selected article/19	M_{\odot}
NGC 6166	M_{\odot}
Celestial Snow Angel	M_{\odot}
Huge-LQG	M_{\odot}
High-velocity cloud	M_{\odot}
NGC 4845	M_{\odot}
Pulsating white dwarf	M_{\odot}
Robust associations of massive baryonic objects	M_{\odot}
Black Widow Pulsar	M_{\odot}
Betelgeuse	M_{\odot}
Andromeda Galaxy	M_{\odot}

In general doesn't give clusters

TODO add some numbers and graphs

WEAK ASSOCIATION

K-Means weak association

DBSCAN k=15, eps=8, min_pts=5

Article	Identifiers
Papyrus 66	<i>P</i> papyrus
Alexandrian text-type	<i>P</i> , papyrus
Western text-type	<i>P</i> , papyrus
Codex Ephraemi Rescriptus	<i>P</i> , papyrus
Bodmer Papyri	<i>P</i> , papyrus
Categories of New Testament manuscripts	<i>P</i> , papyrus
Papyrus 4	<i>P</i> , papyrus
Papyrus 75	<i>P</i> , papyrus
Uncial 0308	<i>P</i> , <i>M</i> , papyrus, 47
Codex Athous Lavrensis	<i>P</i> , papyrus
Papyrus 92	<i>P</i> , papyrus
Papyrus 90	<i>P</i> , papyrus
Papyrus 9	<i>P</i> , papyrus
Papyrus 15	<i>P</i> , papyrus
Papyrus 16	<i>P</i> , papyrus
Papyrus 20	<i>P</i> , papyrus
Papyrus 39	<i>P</i> , papyrus
Papyrus 49	<i>P</i> , papyrus
Papyrus 65	<i>P</i> , papyrus
Papyrus 111	<i>P</i> , papyrus
Uncial 0243	<i>P</i> , papyrus
Minuscule 1739	<i>P</i> , papyrus
Minuscule 88	<i>P</i> , papyrus
Authorship of the Epistle to the Hebrews	<i>P</i> , papyrus
Egerton Gospel	<i>P</i> , papyrus
Rylands Library Papyrus P52	<i>P</i> , papyrus
Codex Vaticanus	<i>P</i> , papyrus

K-Means

TODO: think of different ways to represent it - maybe truncate some definitions?

Article	Identifiers
Direct shear test	<i>angle, φ, friction</i>
Truncated dodecadodecahedron	<i>golden, ratio, ϕ</i>
Golden triangle (mathematics)	<i>golden, section, θ, ϕ</i>
Petrophysics	<i>percentageφ, symbol, S_w</i>
Greedy algorithm for Egyptian fractions	<i>golden, terms, d, ϕ, denominator, possible, expa</i>
Pi Josephson junction	<i>π, junction, φ</i>
Snub dodecahedron	<i>golden, ratio, τ, 2ξ, $-$, V, ξ</i>
Lucas number	<i>golden, terms, ϕ, m, L, number, values, ratio, L</i>
54 (number)	<i>golden, φ, ratio</i>
Special right triangles	<i>π, c, b, ratio, c., ϕ, m, golden, radians, n, line, in</i>
Universal code (data compression)	<i>golden, code, ratio, power, ϕ, l, n, q, p</i>
Existential instantiation	<i>csymbolφvariable</i>
Perles configuration	<i>goldenratioϕ</i>
Wythoff array	<i>goldenratioϕcolumnmnumber$n$$\varphi$fib</i>
Golomb sequence	<i>a_n, n, golden, ratio, ϕ</i>
Almost integer	<i>constantπ, gel fonds, φ, exampl</i>
16:10	<i>golden, ratio, ϕ</i>
Leonardo number	<i>golden, ratio, ϕ, computations, ψ, L, n</i>
RogersRamanujan continued fraction	<i>golden.functionsratio$G\phi Hq$modu</i>
Great rhombic triacontahedron	<i>goldenratioϕ</i>
108 (number)	<i>goldenφratio</i>
Random Fibonacci sequence	<i>golden.BratioM_nprobabilityϕf_nsequencenrandominc</i>
Rhombic triacontahedron	<i>goldenratioϕr_iSedger$_mV$</i>
Exact trigonometric constants	<i>functionπratioϕimagegoldenvaluesxV</i>
Bilunabirotanda	<i>goldenratioϕ</i>
Feigenbaum constants	<i>mapzratioplaces$f\phi$goldenc$_n\alpha\delta$varia</i>

Found some interesting clusters but in general doesn't show good results.

Need to use semantic menthols

Batch, $k = 2500 \dots 10000$ with step 50 SVD with $n=600$

HEre results

7.2 Result analysis

Hierarcical methods are too slow, and SLINK is not good. Bisecting K -Means is good for explaining steps but not very practical

MiniBatch K means is preferred to usual K Means: fast but same results

NMF takes a lot of time to decompose a matrix with large k

$k = 100$ 30 min, but with results inferior to SVD $k = 250$ 2 hours, with results comparable to SVD

The complexity of NMF is $O(kn)$

The best definition embedding technique is soft association The best clustering algorithm is K -Means with $K = 10000$ on the semantic space produced by rank-reduced SVD with $k = 200$

7.3 Building Hierarchy

How to evaluate???

7.4 Experiment Conclusions

8 Conclusions

The results are super.

8.1 Future Work

the work is done assuming that document imports only from one namespace but it can import from several. can solve that by dividing the document in parts (e.g. by paragraphs) and then applying the same analysis independently to each paragraph - instead of each document.

Can use additional information from wiki articles. For example, can extract some keywords from the article and use it in clustering

Or interwiki pages.

Pages that describe certain namespaces may be quite interconnected. There are link-based clustering methods e.g. Botafogo and Schneiderman 1991

Can extract wiki graph and use this for clustering . There are hybrid approaches that use both usual textual representation + links [25]

It can be interesting to apply these techniques to a larger dataset, for example, arXiv.

Other dim red techniques for LSA, e.g. Local NMF [68] There should also be randomized NMF that works faster.

Try other clustering techniques: spectral clustering [69] other ways to embed identifiers like word2vec [70] or GloVe [71]

How to extend this method to situations when no additional information about document category is known. I.e. need to replace the notion of purity with some other objective for discovering namespaces and namespace-defining clusters

9 Bibliography

References

1. Erik Duval, Wayne Hodgins, Stuart Sutton, and Stuart L Weibel. Metadata principles and practicalities. *D-lib Magazine*, 8(4):16, 2002.
2. Anders Møller and Michael I Schwartzbach. *An introduction to XML and Web Technologies*. Pearson Education, 2006.
3. Henry Thompson, Tim Bray, Dave Hollander, Andrew Layman, and Richard Tobin. Namespaces in XML 1.0 (third edition). W3C recommendation, W3C, December 2009. <http://www.w3.org/TR/2009/REC-xml-names-20091208/>.
4. Kevin McArthur. Whats new in php 6. *Pro PHP: Patterns, Frameworks, Testing and More*, pages 41–52, 2008.
5. James Gosling, Bill Joy, Guy Steele, Gilad Bracha, and Alex Buckley. *The Java 8 Language Specification, Java SE 8 Edition*. Addison-Wesley Professional, 2014.
6. Craig Larman. *Applying UML and patterns: an introduction to object-oriented analysis and design and iterative development*. Pearson Education India, 2005.
7. Jon Barwise, John Etchemendy, Gerard Allwein, Dave Barker-Plummer, and Albert Liu. *Language, proof and logic*. CSLI publications, 2000.
8. Robert Pagael and Moritz Schubotz. Mathematical language processing project. *arXiv preprint arXiv:1407.0167*, 2014.
9. Giovanni Yoko Kristianto, MQ Ngien, Yuichiroh Matsubayashi, and Akiko Aizawa. Extracting definitions of mathematical expressions in scientific papers. In *Proc. of the 26th Annual Conference of JSAI*, 2012.
10. David Carlisle, Robert R Miner, and Patrick D F Ion. Mathematical markup language (MathML) version 3.0 2nd edition. W3C recommendation, W3C, April 2014. <http://www.w3.org/TR/2014/REC-MathML3-20140410/>.
11. Dan Jurafsky and James H Martin. *Speech & language processing*. Pearson Education India, 2000.
12. Beatrice Santorini. Part-of-speech tagging guidelines for the penn treebank project (3rd revision). 1990.
13. Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J Bethard, and David McClosky. The stanford CoreNLP natural language processing toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60, 2014.
14. Ulf Schöneberg and Wolfram Sperber. POS tagging and its applications for mathematics. In *Intelligent Computer Mathematics*, pages 213–223. Springer, 2014.
15. Mihai Grigore, Magdalena Wolska, and Michael Kohlhase. Towards context-based disambiguation of mathematical expressions. In *The Joint Conference of ASCM*, pages 262–271, 2009.
16. Keisuke Yokoi, Minh-Quoc Nghiem, Yuichiroh Matsubayashi, and Akiko Aizawa. Contextual analysis of mathematical expressions for advanced mathematical search. In *Prof. of 12th International Conference on Intelligent Text Processing and Computational Linguistics (CICLing 2011), Tokyo, Japan, February*, pages 20–26, 2011.
17. Minh Nghiem Quoc, Keisuke Yokoi, Yuichiroh Matsubayashi, and Akiko Aizawa. Mining coreference relations between formulas and text using Wikipedia. In *23rd International Conference on Computational Linguistics*, page 69, 2010.
18. Jerzy Trzeciak. *Writing mathematical papers in English: a practical guide*. European Mathematical Society, 1995.
19. Giovanni Yoko Kristianto, Akiko Aizawa, et al. Extracting textual descriptions of mathematical expressions in scientific papers. *D-Lib Magazine*, 20(11):9, 2014.

20. Christopher D Manning, Prabhakar Raghavan, Hinrich Schütze, et al. *Introduction to information retrieval*, volume 1. Cambridge university press Cambridge, 2008.
21. Eric Evans. *Domain-driven design: tackling complexity in the heart of software*. Addison-Wesley Professional, 2004.
22. Scott C. Deerwester, Susan T Dumais, Thomas K. Landauer, George W. Furnas, and Richard A. Harshman. Indexing by latent semantic analysis. *JAsIs*, 41(6):391–407, 1990.
23. Alfio Gliozzo and Carlo Strapparava. *Semantic domains in computational linguistics*. Springer Science & Business Media, 2009.
24. Christopher Stokoe, Michael P Oakes, and John Tait. Word sense disambiguation in information retrieval revisited. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pages 159–166. ACM, 2003.
25. Nora Oikonomakou and Michalis Vazirgiannis. A review of web document clustering approaches. In *Data mining and knowledge discovery handbook*, pages 921–943. Springer, 2005.
26. Charu C Aggarwal and ChengXiang Zhai. A survey of text clustering algorithms. In *Mining Text Data*, pages 77–128. Springer, 2012.
27. Fabrizio Sebastiani. Machine learning in automated text categorization. *ACM computing surveys (CSUR)*, 34(1):1–47, 2002.
28. Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval. *Information processing & management*, 24(5):513–523, 1988.
29. Levent Ertöz, Michael Steinbach, and Vipin Kumar. Finding clusters of different sizes, shapes, and densities in noisy, high dimensional data. In *SDM*, pages 47–58. SIAM, 2003.
30. Deborah Hughes-Hallett, William G. McCallum, Andrew M. Gleason, et al. *Calculus: Single and Multivariable, 6th Edition*. Wiley, 2013.
31. Tuomo Korenius, Jorma Laurikkala, and Martti Juhola. On principal component analysis, cosine and euclidean measures in information retrieval. *Information Sciences*, 177(22):4893–4905, 2007.
32. Michael Steinbach, George Karypis, Vipin Kumar, et al. A comparison of document clustering techniques. In *KDD workshop on text mining*, volume 400, pages 525–526. Boston, 2000.
33. Rui Xu, Donald Wunsch, et al. Survey of clustering algorithms. *Neural Networks, IEEE Transactions on*, 16(3):645–678, 2005.
34. Mark Hall, Paul Clough, and Mark Stevenson. Evaluating the use of clustering for automatically organising digital library collections. In *Theory and Practice of Digital Libraries*, pages 323–334. Springer, 2012.
35. David Sculley. Web-scale k-means clustering. In *Proceedings of the 19th international conference on World wide web*, pages 1177–1178. ACM, 2010.
36. Douglass R Cutting, David R Karger, Jan O Pedersen, and John W Tukey. Scatter/Gather: A cluster-based approach to browsing large document collections. In *Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 318–329. ACM, 1992.
37. Bjornar Larsen and Chinatsu Aone. Fast and effective text mining using linear-time document clustering. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 16–22. ACM, 1999.
38. Hinrich Schütze and Craig Silverstein. Projections for efficient document clustering. In *ACM SIGIR Forum*, volume 31, pages 74–81. ACM, 1997.
39. Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, volume 96, pages 226–231, 1996.
40. Levent Ertöz, Michael Steinbach, and Vipin Kumar. Finding topics in collections of documents: A shared nearest neighbor approach. pages 83–103, 2004.
41. Thomas K Landauer, Peter W Foltz, and Darrell Laham. An introduction to latent semantic analysis. *Discourse processes*, 25(2-3):259–284, 1998.

42. Stanisław Osiński, Jerzy Stefanowski, and Dawid Weiss. Lingo: Search results clustering algorithm based on singular value decomposition. In *Intelligent information processing and web mining*, pages 359–368. Springer, 2004.
43. Nicholas Evangelopoulos, Xiaoni Zhang, and Victor R Prybutok. Latent semantic analysis: five methodological recommendations. *European Journal of Information Systems*, 21(1):70–86, 2012.
44. Stanisław Osiński. Improving quality of search results clustering with approximate matrix factorisations. In *Advances in Information Retrieval*, pages 167–178. Springer, 2006.
45. Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
46. Wei Xu, Xin Liu, and Yihong Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pages 267–273. ACM, 2003.
47. Wikimedia Foundation. English wikipedia XML data dump, 2015. <http://dumps.wikimedia.org/enwiki/latest/>, accessed: TODO.
48. Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören Auer, and et al. DBpedia – a crystallization point for the web of data. *Web Semant.*, 7(3):154–165, September 2009.
49. Apache Software Foundation. Apache Flink 0.8.1. <http://flink.apache.org/>, accessed: 2015-01-01.
50. Ronald Rivest. The MD5 message-digest algorithm. 1992.
51. Eclipse Foundation. Mylyn WikiText 1.3.0, 2015. <http://projects.eclipse.org/projects/mylyn.docs>, accessed: 2015-01-01.
52. Daniel Sonntag. Assessing the quality of natural language text data. In *GI Jahrestagung (1)*, pages 259–263, 2004.
53. Julie D Allen et al. *The Unicode Standard*. Addison-Wesley, 2007.
54. Martin F Porter. Snowball: A language for stemming algorithms, 2001.
55. Steven Bird. NLTK: the natural language toolkit. In *Proceedings of the COLING/ACL on Interactive presentation sessions*, pages 69–72. Association for Computational Linguistics, 2006.
56. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, et al. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
57. Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001–. <http://www.scipy.org/>, accessed: 2015-02-01.
58. S. van der Walt, S.C. Colbert, and G. Varoquaux. The numpy array: A structure for efficient numerical computation. *Computing in Science Engineering*, 13(2):22–30, March 2011.
59. A Tropp, N Halko, and PG Martinsson. Finding structure with randomness: Stochastic algorithms for constructing approximate matrix decompositions. Technical report, 2009.
60. American Mathematical Society. AMS mathematics subject classification 2010, 2009. <http://msc2010.org/>, accessed: 2015-06-01.
61. American Physical Society. PACS 2010 regular edition, 2009. <http://www.aip.org/publishing/pacs/pacs-2010-regular-edition/>, accessed: 2015-06-01.
62. Bernard Rous. Major update to ACM’s computing classification system. *Commun. ACM*, 55(11):12–12, November 2012.
63. Alistair Miles, Brian Matthews, Michael Wilson, and Dan Brickley. SKOS Core: Simple knowledge organisation for the web. In *Proceedings of the 2005 International Conference on Dublin Core and Metadata Applications: Vocabularies in Practice*, DCMI ’05, pages 1:1–1:9. Dublin Core Metadata Initiative, 2005.
64. Association for Computing Machinery. ACM computing classification system, 2012. <https://www.acm.org/about/class/2012>, accessed: 2015-06-21.
65. Daniel Krech. RDFLib 4.2.0. <https://rdflib.readthedocs.org/en/latest/>, accessed: 2015-06-01.

66. Sreenivasa Viswanadha, Danny van Bruggen, and Nicholas Smith. JavaParser 2.1.0, 2015. <http://javaparser.github.io/javaparser/>, accessed: 2015-06-15.
67. Apache Software Foundation. Apache Mahout 0.10.1. <http://mahout.apache.org/>, accessed: 2015-06-15.
68. Stan Z Li, Xin Wen Hou, HongJiang Zhang, and QianSheng Cheng. Learning spatially localized, parts-based representation. In *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, volume 1, pages I–207. IEEE, 2001.
69. Andrew Y Ng, Michael I Jordan, Yair Weiss, et al. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 2:849–856, 2002.
70. Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *Proceedings of Workshop at ICLR*, 2013.
71. Jeffrey Pennington, Richard Socher, and Christopher D Manning. GloVe: Global vectors for word representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12:1532–1543, 2014.