

Topic Analysis in AceMap

Hao Lu

With group members Zhengtian Xu, Yiyi Shen, and Ziming Cheng

5130309109

June 24, 2016

Contents

Contents	1
1 Motivation	2
2 Principle	2
3 Model	2
3.1 Database	2
3.2 Distribution of papers over topics	2
3.3 Distribution of a topic over papers	4
3.4 Topic Strength and Dependencies	4
4 Demo	5
4.1 Topic Strength	5
4.2 Topic Relation	6
5 Future work	7

1 Motivation

- Exploit the idea of Data Mining: find what's under the explicit statistics, rather than just visualize plain data.
- Design an algorithm for calculating the trend of a topic and deciding the dependencies between topics.
- Deploy the algorithm in AceMap: most theoretical algorithms are not practical, thus we need a tradeoff between being pragmatic and precise.

2 Principle

The popular algorithms for topic analysis often use LDA models, which is used on a set of un-tagged papers. After iterative clustering, the papers will be clustered into different topics. As a result, we can get two kinds of distributions, that is, distribution of papers over topics, and distribution of topics over papers. As you might have expected, this is a rather time-consuming process.

3 Model

3.1 Database

The database from Microsoft is a powerful one, with topics already labeled among keywords.

PaperID	KeywordName	FieldOfStudy...
000000BE	vaporization	0BABB10F
0000024B	population dynarr	0206C527
00000569	planetary nebula	05AF48EC
000010BF	boundary layer	06A6EEEE
00001561	falsification	06ED6B1B
000017D0	three dimensiona	09633D62
00001B5B	educational techn	00FAE71A

However, we must acknowledge that the data is in a poor manner. The topics are in a tree-like hierarchy. Sometimes one subtopic may belong to multiple topics in different layers, with confidence or probability. This makes our work difficult. So we design a recursive procedure for classifying a topic's related, parent topics.

ChildFieldOfSt...	ChildFieldOfSt...	ParentFieldOf...	ParentFieldOf...	Confidence
0001BE57	L3	09892D7B	L1	0.5
0001DDB2	L3	0B0FEB68	L0	1
000226EC	L3	052C8328	L0	1
0002487D	L3	092C7EF1	L1	0.5467854356075009
0002C926	L3	065808E3	L2	0.5
00032F94	L3	064E5072	L1	1
0004676F	L3	07982D63	L0	0.6038416171751301
0005EC96	L3	007E3B49	L2	1

3.2 Distribution of papers over topics

For each paper, we get the distribution over topics in the following way:

1. Gather the papers that this paper has referred to.

2. For each reference, get the topics that the keywords map to.
3. For each topics, follow the hierarchy and get the topics in the layer we want.

In this way we'll get a probabilistic distribution.

```

1 def collectPhi(fieldID):
2     f = codecs.open(fieldID + ".phi",
3                     encoding="utf-8", mode="w")
4
5     # Need to fill in the user and passwd.
6     connection = MySQLdb.connect(...)
7
8
9     cursor = connection.cursor()
10
11     # Find all related topics
12     cursor.execute(
13         "SELECT ChildFieldOfStudyID, Confidence
14         FROM FieldOfStudyHierarchy
15         WHERE ParentFieldOfStudyID = {}".format(fieldID))
16     fields = cursor.fetchall()
17     print len(fields)
18     #for i in fields:
19     #    print i
20
21     # Find all papers in a specific topic
22     for field in fields + ((fieldID, 1),):
23         print field
24         cursor.execute(
25             "SELECT PaperID
26             From PaperKeywords
27             Where FieldOfStudyIDMappedToKeyword =
28             {}".format(field[0]))
29         result = cursor.fetchall()
30
31         # Find information of each paper
32         for paper in result:
33             cursor.execute
34                 ("SELECT PaperID,
35                 PaperPublishYear, PaperRank
36                 FROM Papers
37                 WHERE PaperID = {}".format(
38                 paper[0]))
39             arr = cursor.fetchone()
40             cursor.execute(
41                 "SELECT *
42                 FROM PaperKeywords
43                 WHERE PaperID = {}".format(paper[0]))
44             keywordNum = len(cursor.fetchall())
45             f.write(arr[0])
46             f.write(u"\t")

```

```

47         f.write(str(arr[1]))
48         f.write(u"\t")
49         f.write(str(arr[2] / 10000.0))
50         f.write(u"\t")
51         f.write(str(keywordNum))
52         f.write(u"\t")
53         f.write(str(field[1]))
54         f.write(u"\n")
55
56     connection.close()
57     f.close()

```

Code Listing 1: collectPhi()

3.3 Distribution of a topic over papers

For each topic in the layer we want, we generate its distribution in the following way:

1. Find all child topics.
2. For each child topic and the topic itself, find papers belonging to the topic.

The program for this is not completed by me, so I'll not post it here.

3.4 Topic Strength and Dependencies

With the two kinds of distributions, we can perform many valuable calculations. For example, we can collect all papers related to a topic, thus forming its strength or evolution pattern. Also, we can get the dependency between two topics, thus generating a dependency graph.

```

1 def calculateStrength(paperYear,
2   paperRank, paperKeywordNum, paperConfidence):
3     strength = [0] * 120
4     num = [0] * 120
5     connection = MySQLdb.connect(...)
6     cursor = connection.cursor()
7     cnt = 0
8     try:
9         for paper in paperYear.keys():
10            year = paperYear[paper]
11            rank = paperRank[paper]
12            keywordNum = paperKeywordNum[paper]
13            Confidence = paperConfidence[paper]
14            paperStrength = rank / keywordNum * Confidence
15            strength[year - baseYear] += paperStrength
16            num[year - baseYear] += 1
17            cursor.execute(
18                "SELECT PaperID
19                FROM PaperReferences
20                WHERE PaperReferenceID={}".format(paper))
21            citations = cursor.fetchall()
22
23            for citation in citations:

```

```

24         if citation[0] in paperYear:
25             citationStrength =
26                 paperStrength * alpha / len(citations)
27                 strength[paperYear[citation[0]] - baseYear]
28                 += citationStrength
29             cnt += 1
30             #if cnt % 100 == 0:
31                 print cnt
32     except Exception:
33         pass
34     finally:
35         connection.close()
36     return strength, num

```

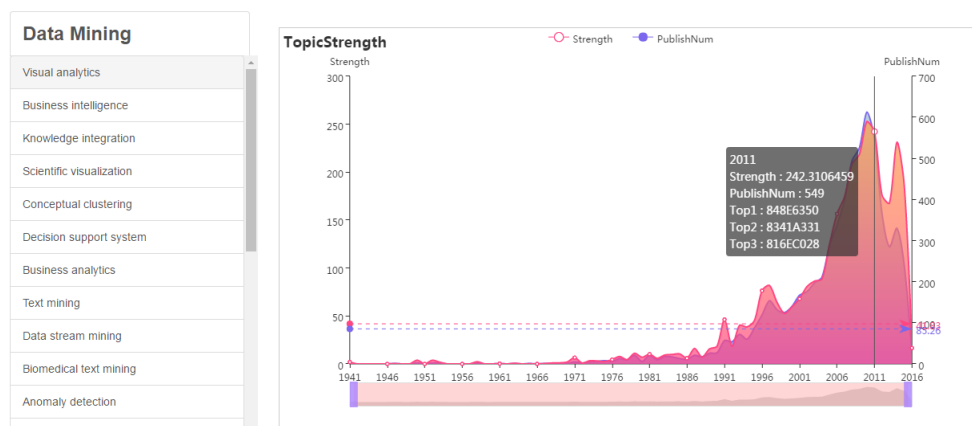
Code Listing 2: calculateStrength()

4 Demo

The visualization is mainly finished by the two girls in our group. I and Zhengtian Xue also contributes in some aspects.

4.1 Topic Strength

Here I select two topics under Data Mining and show the results of ours compared to the results now on AceMap.

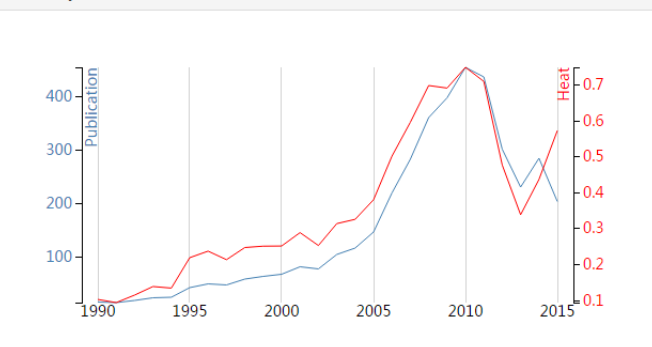


Visual Analytics

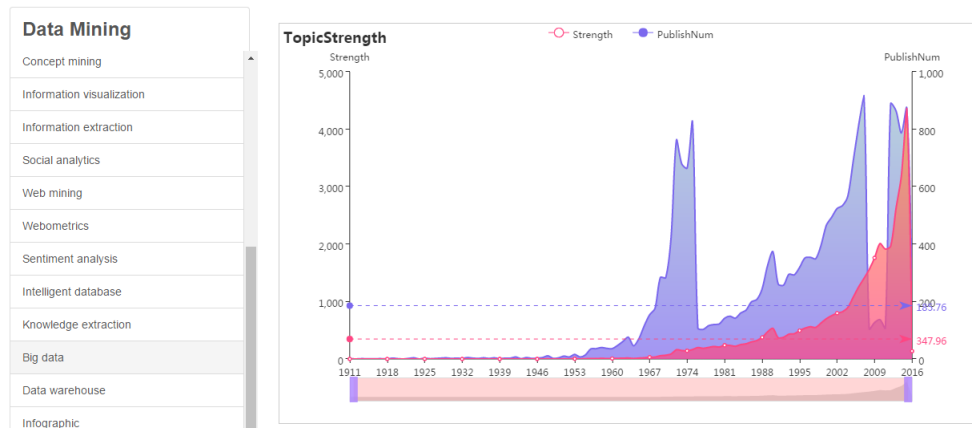
♥ Related Fields



Development and Trends



You'll see we've got a smoother and more reasonable strength. For example, you can see that in recent years, although the number of publication falls (which may be because of an incomplete database), the strength doesn't fall that much. This corresponds to our expectation.

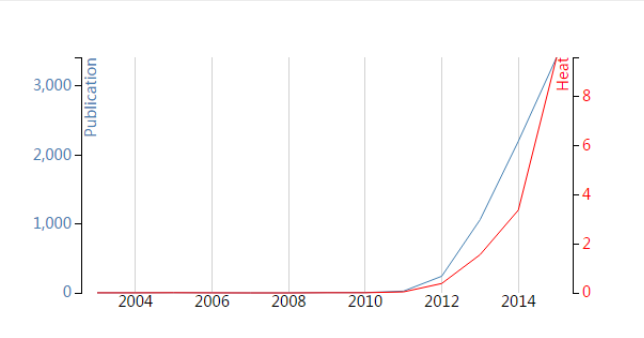


Big Data

♥ Related Fields



Development and Trends

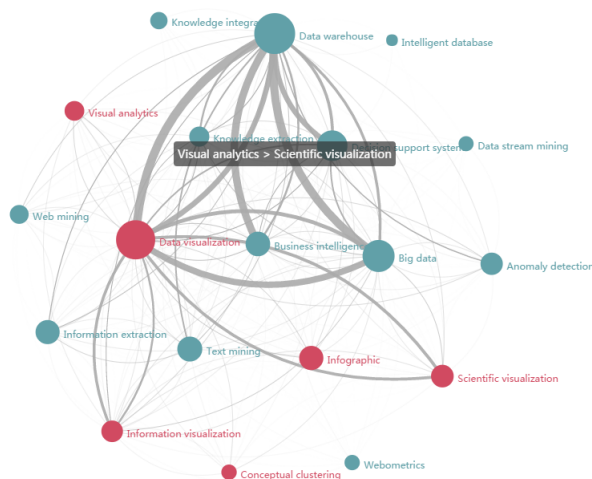


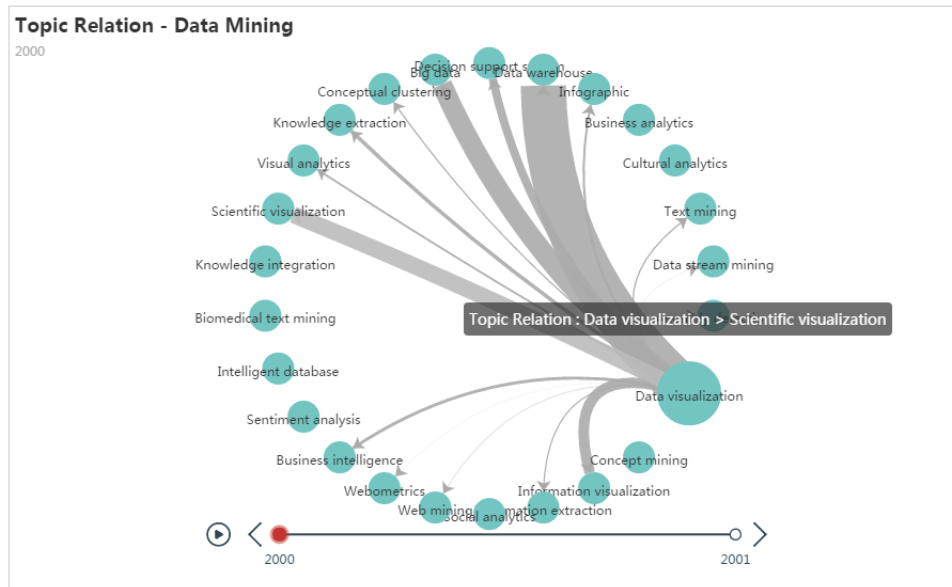
Here we'll get another view. You'll see that maybe due to the false classification of the database, the number of papers changes very abnormally. However, strength here is a good indicator, since it's very smooth and grows to our expectation.

4.2 Topic Relation

We can also show the relation among topics by their dependencies.

Data Mining 2000





5 Future work

- The program is still not fast enough. The very bottleneck is the cost of accessing the database.
- Many internal procedures are not integrated.
- Further analysis of the data we collect is available such as timeline and evolution of topics.
- We'll try to cover all CS topics in the future.