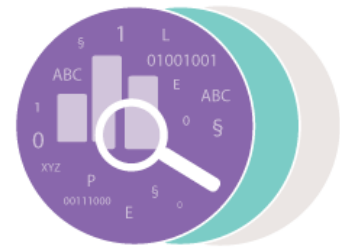




# FutureTDM

Explore . Analyse . Improve



REDUCING BARRIERS AND INCREASING UPTAKE OF TEXT AND DATA MINING FOR RESEARCH ENVIRONMENTS USING A COLLABORATIVE KNOWLEDGE AND OPEN INFORMATION APPROACH

## Deliverable D4.4

### Scientific Review Paper

## Project

|              |  |
|--------------|--|
| Acronym:     | <b>FutureTDM</b>   |
| Title:       | Reducing Barriers and Increasing Uptake of Text and Data Mining for Research Environments using a Collaborative Knowledge and Open Information Approach  |
| Coordinator: | SYNYO GmbH   |
| Reference:   | 665940   |
| Type:        | Collaborative project  |
| Programme:   | HORIZON 2020   |
| Theme:       | GARRI-3-2014 - Scientific Information in the Digital Age: Text and Data Mining (TDM)   |
| Start:       | 01. September, 2015  |
| Duration:    | 24 months  |
| Website:     | <a href="http://www.futuretdm.eu/">http://www.futuretdm.eu/</a>  |
| E-Mail:      | <a href="mailto:office@futuretdm.eu">office@futuretdm.eu</a>   |
| Consortium:  | <b>SYNYO GmbH</b> , Research & Development Department, Austria, (SYNYO)<br><b>Stichting LIBER</b> , The Netherlands, (LIBER)<br><b>Open Knowledge</b> , UK, (OK/CM)<br><b>Radboud University</b> , Centre for Language Studies The Netherlands, (RU)<br><b>The British Library Board</b> , UK, (BL)<br><b>Universiteit van Amsterdam</b> , Inst. for Information Law, The Netherlands, (UVA)<br><b>Athena Research and Innovation Centre in Information, Communication and Knowledge Technologies</b> , Inst. for Language and Speech Processing, Greece, (ARC)<br><b>Ubiquity Press Limited</b> , UK, (UP)<br><b>Fundacja Projekt: Polska</b> , Poland, (FPP) |

## Deliverable

|                      |   |
|----------------------|---|
| Number:              | <b>D4.4</b>   |
| Title:               | <b>Scientific Review Paper</b>  |
| Lead beneficiary:    | <b>RU</b>   |
| Work package:        | WP4: Fields of Application, Projects, Best Practices and Resources  |
| Dissemination level: | Public (PU)   |
| Nature:              | Report (RE)   |
| Due date:            | 30.04.2017  |
| Submission date:     | 02.05.2017  |
| Authors:             | <b>Maria Eskevich, RU</b><br><b>Antal van den Bosch, RU</b>   |
| Contributors:        | <b>Stelios Piperidis, ARC</b><br><b>Kanella Pouli, ARC</b><br><b>Maria Gavriilidou, ARC</b><br><b>Dimitris Galanis, ARC</b> |
| Review:              | <b>Burcu Akinci, SYNYO</b>  |

**Acknowledgement:** This project has received funding from the European Union's Horizon 2020 Research and Innovation Programme under Grant Agreement No 665940.

**Disclaimer:** The content of this publication is the sole responsibility of the authors, and does not in any way represent the view of the European Commission or its services.

This report by FutureTDM Consortium members can be reused under the [CC-BY 4.0](https://creativecommons.org/licenses/by/4.0/) license (<https://creativecommons.org/licenses/by/4.0/>).

# Table of Contents

- 1 Summary..... 5
- 2 Work on the scientific paper ..... 6
  - 2.1 Criteria for selection of potential publication venues..... 6
  - 2.2 Information on the submitted paper ..... 8
- 3 Work on the FutureTDM knowledge base enrichment..... 9
  - 3.1 Collection of Projects..... 9
  - 3.2 Collection of Organisations ..... 9
  - 3.3 Collection of Methods ..... 10
  - 3.4 Collection of Tools ..... 11
- 4 Conclusions ..... 15
- 5 Annex..... 16

# List of Tables

- Table 1: Total number of projects in the FutureTDM Knowledge Base Collection structured by economic sector ..... 9
- Table 2: Total number of projects in the FutureTDM Knowledge Base Collection structured by country ..... 10
- Table 3: FutureTDM text and data mining methods collection ..... 11
- Table 4: Total list of tools in the FutureTDM Knowledge Base Collection ..... 14

## 1 SUMMARY

This report provides information about the work on the scientific review paper that gives a critical review of the past, present and future of the state-of-the-art of the field of text and data mining (TDM). This publication aims to broaden perspective on the past, current, and future trends in TDM amongst audience, outlined via desktop research, consultation with the practitioners (Deliverables D.4.3. and D4.5.) and scientific communities. Report includes brief description of publication venue selection, text of the paper that can be further updated with the open accesses reviewers comments.

Additionally, the information about the update of FutureTDM knowledge base collection of relevant projects, organisations, methods and tools is provided including the current list of the collections in the Annex.

## 2 WORK ON THE SCIENTIFIC PAPER

### 2.1 Criteria for selection of potential publication venues

We have defined the two main criteria for the target publication venues:

- Diverse type of TDM applications already present in the journal
- The journal has to be Open Access, or at least to have the open access publication option

As the review of TDM work from the perspective of overall scientific development and insights on its visibility across different economic fields, that are formulated within the FutureTDM consortium, differ from the traditional type of publication content that is expected by journals, we took an initiative to contact a number of targeted journals in order to find a venue that would fit our criteria and to ensure that our content would be suitable. These were the journals to contact:

#### 1. Data Science Journal

- Link: <http://datascience.codata.org>
- Aims and scope of the journal: The CODATA Data Science Journal is a peer-reviewed, open access, electronic journal, publishing papers on the management, dissemination, use and reuse of research data and databases across all research domains, including science, technology, the humanities and the arts. The scope of the journal includes descriptions of data systems, their implementations and their publication, applications, infrastructures, software, legal, reproducibility and transparency issues, the availability and usability of complex datasets, and with a particular focus on the principles, policies and practices for open data. All data is in scope, whether born digital or converted from other sources.
  - *Reviews* can cover topics such as current controversies, the current “state of the art” or the historical development of studies as well as issues of regional or temporal focus. Papers should critically engage with the relevant body of extant literature. Review articles should be no longer than 3,000 words in length.
- Open access policy: multidisciplinary Open Access journal

#### 2. PLOSOne

- Link: <http://journals.plos.org/plosone/>
- Aims and scope of the journal: PLOS ONE features reports of original research from all disciplines within science and medicine. By not excluding research on the basis of subject area, PLOS ONE facilitates the discovery of connections between research whether within or between disciplines.

- We consider publishing systematic reviews only if the methods ensure the comprehensive and unbiased sampling of existing literature.
- Submissions describing methods, software, databases, or other tools. We consider submissions describing methods, software, databases, or other tools if they follow the appropriate reporting guidelines.
- Qualitative research. We consider publishing qualitative research only if it adheres to appropriate study design and reporting guidelines.
- Studies reporting negative results.
- Open access policy: multidisciplinary Open Access journal

### 3. DSH (Digital Scholarship for Humanities)

- Link: <https://academic.oup.com/dsh>
- Aims and scope of the journal: *DSH or Digital Scholarship in the Humanities* is an international, peer reviewed journal which publishes original contributions on all aspects of digital scholarship in the Humanities including, but not limited to, the field of what is currently called the Digital Humanities. Long and short papers report on theoretical, methodological, experimental, and applied research and include results of research projects, descriptions and evaluations of tools, techniques, and methodologies, and reports on work in progress. *DSH* also publishes reviews of books and resources.
- Open access policy: multidisciplinary Open Access journal

### 4. Science Communication (SAGE)

- Web-Link: <http://journals.sagepub.com/home/scx>
- Aims and scope of the journal: Science Communication is an international, interdisciplinary social science journal that examines such topics as the nature of scientific expertise as represented through communication and the processes or effects characterizing the communication of science in any context. Science is broadly defined to include environmental science, health science, and technology. Science Communication welcomes submissions of empirical research from authors in all relevant disciplines (including the social sciences, the humanities, and science itself). Both qualitative and quantitative research papers with a basis in theory are acceptable. Preference is given to articles that bridge the gap between theory and practice and that will be of interest across disciplines. In addition to peer-reviewed research, Science Communication publishes commentaries that analyse issues and trends in the field – whether scholarly, professional, or policy-related – and a periodic summary of new books in the field.
- Open access policy: OA option is available

### 5. IEEE Access

- Web-Link: <http://ieeaccess.ieee.org/learn-more-about-ieee-access/>

- Aims and scope of the journal: IEEE Access publishes articles that are of high interest to readers: original, technically correct, and clearly presented. The scope of this journal comprises all of IEEE's fields of interest, emphasizing applications-oriented and interdisciplinary articles. IEEE Access also accepts traditional technical articles, as well as reviews and surveys.
- Open access policy: Multidisciplinary Open Access Journal, fee \$1,750 per article

Eventually, we have chosen the “Data Science Journal” that is an open access multidisciplinary journal run by an SME publisher “Ubiquity Press”.

## 2.2 Information on the submitted paper

The details of the submission are the following:

- Authors: Maria Eskevich, Stelios Piperidis, Kanella Pouli, Maria Gavriilidou, Dimitris Galanis, Antal van den Bosch.
- Title: Text and Data Mining: Past, present, and future.
- Submission date: 01.05.2017.
- Text of the paper is in the Appendix



## 3 WORK ON THE FUTURETDM KNOWLEDGE BASE ENRICHMENT

### 3.1 Collection of Projects

FutureTDM collection of projects by the Deliverable submission date contains 140 items, 37 of which are already integrated into the platform. For 7 of the projects it is not possible to assign an economic sector, such as (+)Spaces , while the other 133 can be assigned to one or even a combination of sectors due to the multidisciplinary nature of the projects. Table 1 contains the overall statistics of the number of projects in the collection.

| Sector or combination of Sectors that that the project is assigned to | Number of projects |
|---|--------------------|
| Primary-Sector, Secondary-Sector, Tertiary-Sector, Quaternary-Sector  | 1                  |
| Primary-Sector, Quaternary-Sector                                     | 1                  |
| Secondary-Sector  | 8                  |
| Secondary-Sector, Tertiary-Sector                                     | 3                  |
| Secondary-Sector, Quaternary-Sector                                   | 38                 |
| Secondary-Sector, Tertiary-Sector, Quaternary-Sector                  | 6                  |
| Tertiary-Sector   | 13                 |
| Tertiary-Sector, Quaternary-Sector                                    | 24                 |
| Quaternary-Sector   | 38                 |
| Quinary-Sector  | 1                  |
| <b>TOTAL</b>  | <b>133</b>         |

**Table 1: Total number of projects in the FutureTDM Knowledge Base Collection structured by economic sector**

### 3.2 Collection of Organisations

FutureTDM collection of organisations by the Deliverable submission date contains 72 items, 59 of which are already integrated into the platform. 55 of these organisations have only one EU-country location, while three of them spread across more than one country within the EU, and 6 have locations outside EU in USA.

| Company locations details           | Country                                | Number    |
|-------------------------------------|--|-----------|
| Company is based in one country     | Austria                                | 1         |
|                                     | Belgium                                | 2         |
|                                     | Czech Republic                         | 2         |
|                                     | Denmark                                | 2         |
|                                     | Finland                                | 1         |
|                                     | France                                 | 4         |
|                                     | Germany                                | 10        |
|                                     | Greece                                 | 3         |
|                                     | Ireland                                | 1         |
|                                     | Italy                                  | 3         |
|                                     | Norway                                 | 1         |
|                                     | Slovenia                               | 1         |
|                                     | Spain                                  | 4         |
|                                     | Sweden                                 | 2         |
|                                     | Switzerland                            | 4         |
| The Netherlands                     | 5                                      |           |
| UK                                  | 9                                      |           |
| Company has offices internationally | Belgium, The Netherlands, Germany, USA | 1         |
|                                     | The Netherlands,                       | 1         |
|                                     | Germany/USA                            | 1         |
|                                     | Ireland/Greece                         | 1         |
|                                     | Spain/UK                               | 1         |
|                                     | UK/USA                                 | 2         |
|                                     | USA, many                              | 2         |
| Companies outside EU                | Japan                                  | 1         |
|                                     | USA                                    | 7         |
| <b>TOTAL</b>                        |  | <b>72</b> |

Table 2: Total number of projects in the FutureTDM Knowledge Base Collection structured by country

### 3.3 Collection of Methods

FutureTDM collection of projects splits scientific methods into 29 items, and all of them are integrated into the platform, as they are used to specify the organisations focus.

| TDM Methods             |                                   |                            |                          |
|-------------------------|-----------------------------------|----------------------------|--------------------------|
| Artificial Intelligence | Machine Learning                  | Statistics                 | Kernel Methods           |
| Association rules       | Machine Translation               | Summarisation              | Ensemble Learning        |
| Classification          | Multimedia Processing             | Term/Concept Extraction    | Dimensionality Reduction |
| Clustering              | Natural Language Processing       | Text Mining                | Decision Trees           |
| Data Mining             | Negation and Modality Detection   | Textual Entailment         | Deep Learning            |
| Graph Mining            | Predictive Analytics              | Artificial Neural Networks |                          |
| Information Extraction  | Regression                        | Bayesian Inference         |                          |
| Information Retrieval   | Sentiment Analysis/Opinion Mining | Instance-Based Learning    |                          |

**Table 3: FutureTDM text and data mining methods collection**

### 3.4 Collection of Tools

FutureTDM collection of projects by the Deliverable submission date contains 72 items, 59 of which are already integrated into the platform. Table contains the complete list of tools collected so far.

| Name                    | URL   |
|-------------------------|---|
| ALVEO                   | <a href="http://alveo.edu.au/">http://alveo.edu.au/</a>   |
| Alvis                   | <a href="https://migale.jouy.inra.fr/redmine/projects/alvisnlp">https://migale.jouy.inra.fr/redmine/projects/alvisnlp</a> |
| Angoss Knowledge STUDIO | <a href="http://www.angoss.com/">http://www.angoss.com/</a>   |
| AnnoMarket              | <a href="https://annomarket.eu/">https://annomarket.eu/</a>   |
| Apache cTAKES           | <a href="http://ctakes.apache.org">http://ctakes.apache.org</a>   |
| Apache OpenNLP          | <a href="http://opennlp.apache.org">http://opennlp.apache.org</a>   |
| Apache UIMA             | <a href="https://uima.apache.org">https://uima.apache.org</a>   |
| Argo                    | <a href="http://argo.nactem.ac.uk">http://argo.nactem.ac.uk</a>   |

|                               |   |
|-------------------------------|---|
| BioCatalogue                  | <a href="https://www.biocatalogue.org/">https://www.biocatalogue.org/</a>   |
| BiodiversityCatalogue         | <a href="https://www.biodiversitycatalogue.org">https://www.biodiversitycatalogue.org</a>   |
| Bluima                        | <a href="https://github.com/BlueBrain/bluima">https://github.com/BlueBrain/bluima</a>   |
| CLARIN-DK                     | <a href="http://clarin.dk/">http://clarin.dk/</a>   |
| ClearTK                       | <a href="https://cleartk.github.io/cleartk">https://cleartk.github.io/cleartk</a>   |
| Clementine                    | <a href="http://datamining.togaware.com/survivor/Clementine.html">http://datamining.togaware.com/survivor/Clementine.html</a>   |
| ContentMine                   | <a href="http://www.contentmine.org/">http://www.contentmine.org/</a>   |
| DKPro Core                    | <a href="https://dkpro.github.io/dkpro-core">https://dkpro.github.io/dkpro-core</a>   |
| ELKI                          | <a href="http://elki.dbs.ifi.lmu.de/">http://elki.dbs.ifi.lmu.de/</a>   |
| FICO Data Management          | <a href="http://www.fico.com/">http://www.fico.com/</a>   |
| Galaxy                        | <a href="https://galaxyproject.org/">https://galaxyproject.org/</a>   |
| GATE Embedded                 | <a href="https://gate.ac.uk">https://gate.ac.uk</a>   |
| Heart of Gold                 | <a href="http://heartofgold.dfki.de">http://heartofgold.dfki.de</a>   |
| IBM SPSS Predictive Analysis  | <a href="http://www-03.ibm.com/software/products/en/spss-predictive-analytics-enterprise">http://www-03.ibm.com/software/products/en/spss-predictive-analytics-enterprise</a>   |
| JCoRe                         | <a href="http://julielab.github.io">http://julielab.github.io</a>   |
| Jpylyzer                      | <a href="http://jpylyzer.openpreservation.org/">http://jpylyzer.openpreservation.org/</a>   |
| KAF                           | <a href="http://kyoto-project.eu/xmlgroup.iit.cnr.it/kyoto/index2091.html?option=com_content&amp;view=article&amp;id=504&amp;Itemid=173">http://kyoto-project.eu/xmlgroup.iit.cnr.it/kyoto/index2091.html?option=com_content&amp;view=article&amp;id=504&amp;Itemid=173</a> |
| Kepler                        | <a href="https://kepler-project.org/">https://kepler-project.org/</a>   |
| KNIME                         | <a href="http://www.knime.org/knime">http://www.knime.org/knime</a>   |
| Language Grid                 | <a href="http://langrid.org/en/index.html">http://langrid.org/en/index.html</a>   |
| LAPPS Grid                    | <a href="http://www.lappsgrid.org/">http://www.lappsgrid.org/</a>   |
| Massive Online Analysis (MOA) | <a href="http://moa.cms.waikato.ac.nz/">http://moa.cms.waikato.ac.nz/</a>   |
| Matchbox                      | <a href="http://matchbox.openpreservation.org/">http://matchbox.openpreservation.org/</a>   |
| Microsoft SQL Server          | <a href="https://technet.microsoft.com/en-">https://technet.microsoft.com/en-</a>   |

|                             |   |
|-----------------------------|---|
| Analysis Service (SSAS)     | <a href="http://us/library/ms175609%28v=sql.90%29.aspx">us/library/ms175609%28v=sql.90%29.aspx</a>  |
| Neural Designer             | <a href="https://www.neuraldesigner.com/">https://www.neuraldesigner.com/</a>   |
| NLTK                        | <a href="http://www.nltk.org">http://www.nltk.org</a>   |
| Open Calais                 | <a href="http://www.opencalais.com/">www.opencalais.com/</a>  |
| Oracle Data Miner GUI       | <a href="http://www.oracle.com/technetwork/database/options/odm/dataminerworkflow-168677.html">http://www.oracle.com/technetwork/database/options/odm/dataminerworkflow-168677.html</a> |
| Orange                      | <a href="http://orange.biolab.si/">http://orange.biolab.si/</a>   |
| Pagelyzer                   | <a href="http://pagelyzer.openpreservation.org/">http://pagelyzer.openpreservation.org/</a>   |
| Pegasus                     | <a href="https://pegasus.isi.edu">https://pegasus.isi.edu</a>   |
| Pentaho                     | <a href="http://www.pentaho.com/">http://www.pentaho.com/</a>   |
| Pipeline Pilot              | <a href="http://accelrys.com/products/collaborative-science/biovia-pipeline-pilot/">http://accelrys.com/products/collaborative-science/biovia-pipeline-pilot/</a>                       |
| QT21                        | <a href="http://qt21.metashare.ilsp.gr/">http://qt21.metashare.ilsp.gr/</a>   |
| Rapidminer Radoop           | <a href="http://www.rapidminer.com">www.rapidminer.com</a>  |
| Rapidminer Server           | <a href="http://www.rapidminer.com">www.rapidminer.com</a>  |
| Rapidminer Studio           | <a href="http://www.rapidminer.com">www.rapidminer.com</a>  |
| SAP Hana                    | <a href="https://hana.sap.com">https://hana.sap.com</a>   |
| SAS Data Mining             | <a href="http://www.sas.com">http://www.sas.com</a>   |
| SPSS                        | <a href="http://www.ibm.com/analytics/us/en/technology/spss/">http://www.ibm.com/analytics/us/en/technology/spss/</a>   |
| Taverna                     | <a href="http://www.taverna.org.uk/">http://www.taverna.org.uk/</a>   |
| Think Analytics             | <a href="http://thinkanalytics.com/">http://thinkanalytics.com/</a>   |
| Think Enterprise Data Miner | <a href="http://www.thinkanalytics.com/">www.thinkanalytics.com/</a>  |
| Triana                      | <a href="http://www.trianacode.org">http://www.trianacode.org</a>   |
| TTNWW                       | <a href="http://yago.meertens.knaw.nl/apache/TTNWW/">http://yago.meertens.knaw.nl/apache/TTNWW/</a>   |
| Weblicht                    | <a href="https://weblicht.sfs.uni-tuebingen.de/">https://weblicht.sfs.uni-tuebingen.de/</a>   |
| WEKA                        | <a href="http://www.cs.waikato.ac.nz/~ml/weka">http://www.cs.waikato.ac.nz/~ml/weka</a>   |
| xcorrSound                  | <a href="http://xcorrSound.openpreservation.org/">http://xcorrSound.openpreservation.org/</a>   |

|                     |   |
|---------------------|---|
| VisTrails           | <a href="http://www.vistrails.org">http://www.vistrails.org</a>                                   |
| Trendminer          | <a href="https://www.trendminer.com/">https://www.trendminer.com/</a>                             |
| Opendatasoft        | <a href="https://www.opendatasoft.com/">https://www.opendatasoft.com/</a>                         |
| Meta                | <a href="http://meta.com/">http://meta.com/</a>   |
| TensorFlow          | <a href="https://www.tensorflow.org/about/">https://www.tensorflow.org/about/</a>                 |
| Epinium             | <a href="http://epinium.com/">http://epinium.com/</a>   |
| Exclusivi           | <a href="http://exclusivi.com/hotels/">http://exclusivi.com/hotels/</a>                           |
| Eparkomat           | <a href="http://www.eparkomat.com/">http://www.eparkomat.com/</a>                                 |
| Data Scouts         | <a href="https://datascouts.eu/">https://datascouts.eu/</a>                                       |
| TopPlace            | <a href="http://www.avuxi.com/products/topplace">http://www.avuxi.com/products/topplace</a>       |
| Influency           | <a href="https://influency.com/en/">https://influency.com/en/</a>                                 |
| Plazi TreatmentBank | <a href="http://plazi.org/api-tools/api/#Plazi_API">http://plazi.org/api-tools/api/#Plazi_API</a> |
| SafeToNet           | <a href="https://safetonet.com/">https://safetonet.com/</a>                                       |
| Egas                | <a href="https://demo.bmd-software.com/egas/">https://demo.bmd-software.com/egas/</a>             |
| Wordfreak           | <a href="http://wordfreak.sourceforge.net/">http://wordfreak.sourceforge.net/</a>                 |
| Theano              | <a href="https://github.com/Theano/">https://github.com/Theano/</a>                               |

**Table 4: Total list of tools in the FutureTDM Knowledge Base Collection**

## 4 CONCLUSIONS

This Deliverable contains the full text of submitted TDM review paper, together with the explanation on the publication venue selection.

We provide quantitative updates on the progress with filling in the four main categories of the Knowledge Base Collection, while the Deliverables 6.1, 6.2 and 6.3<sup>1</sup> contains the qualitative description of the technology behind the FutureTDM Platform.

---

<sup>1</sup> [www.futuretdm.eu](http://www.futuretdm.eu)

## 5 ANNEX

In the following we present the scientific paper “**Text and Data Mining: Past, present, and future**” and the collections below described in the previous section. The list will be given in this order:

- Collection of Organisations
- Collection of Projects
- Collection of Tools



# Text and Data Mining: Past, present, and future

Maria Eskevich<sup>1</sup>, Stelios Piperidis<sup>2</sup>, Kanella Pouli<sup>2</sup>,  
Maria Gavriilidou<sup>2</sup>, Dimitris Galanis<sup>2</sup>, Antal van den Bosch<sup>1</sup>

<sup>1</sup>*Radboud University, Nijmegen, Netherlands,*

<sup>2</sup>*Institute for Language and Speech Processing (ILSP), "Athena" R.C., Greece*

---

## Abstract

In this paper we review text and data mining (TDM) as a set of data science techniques that detect information from massive amounts of data, present this information for further analysis, and support data-driven decision making. We summarise state-of-the-art TDM trends from a scientific perspective, and outline their application in industrial setting. We analyse the spread of TDM across diverse economic sectors via such parameters as publication and funding volumes in the context of research within the European Union as an example.

*Keywords:* Text and Data Mining, Data Science, Big Data, Review, Trends

---

## 1. Introduction

Text and Data Mining has been defined as “the discovery by computer of new, previously unknown information, by automatically extracting and relating information from different (...) resources, to reveal otherwise hidden meanings” [17], in other words, “an exploratory data analysis that leads to the discovery of heretofore unknown information, or to answers for questions for which the answer is not currently known” [17]. Nowadays, this umbrella term encompasses diverse techniques that allow for the interpretation of content of any type ranging from raw data, e.g. sensor data, text, images and multimedia, to processed content and structured data, e.g. diagrams, charts, references, maps, formulas, chemical structures, and metadata, on a large scale through the identification of patterns.

Components of TDM existed before the formal introduction of the term within the scientific community. Essentially TDM has grown as part of the general ‘data science’ field from a number of parent scientific disciplines, notably artificial intelligence and its subdisciplines such as machine learning, pattern recognition, and natural language processing, as well as mathematics, statistics, and information retrieval. It has benefited from multidisciplinary approaches and collaborations across these fields.

The growing quantity of digital content currently produced by society, as well as efforts on the digitization of archives, in combination with the potential ease of access to this data at any location via easy-to-use devices (with reliable Internet connection, and mobile devices) increase both the need for content mining technologies, and their proliferation in the fabric of modern society [29].

TDM is becoming increasingly ubiquitous on the global market, as each company that tracks its activities, products, and communications in some digital form has the potential to reuse this data for further analysis and improvement. For some sectors, countries and applications, it is already possible to estimate the potential profit or at least the size of the market [29, 9]. Other fields such as journalism, which is in the process of embracing the concept of Data Journalism, have to first agree on the conceptual changes required to their current work behaviour in order to merge the traditional work pattern with TDM [20].

This paper is structured as follows. We first survey the scientific trends and outline the most prominent scientific communities that develop the state-of-the-art algorithms and aim to advance beyond those (Section 2). Section 3 provides examples of the TDM scientific work in application across economic structure using the European Union as a case example. Section 4 then discusses future trends and their overall potential.

## **2. Scientific perspective: Research fields supporting TDM**

Text and data mining has gradually evolved into a large multidisciplinary field with complex approaches for a range of applications that can be found in all types of human activities. Historically, the development of mining technologies has its roots with the invention of mathematical models [14] and statistical analysis [8]. Due to the advent of computers and systems for storing larger quantities of data, the possibility to mine large datasets and filter relevant information has greatly expanded what used to be laborious

expert analysis work. Since the beginning of the 1990s, modern data mining tools have appeared on the market and in labs, allowing more advanced analysis, to discover hidden patterns and to extract implicit knowledge in any domain, as envisaged and pioneered for the life science domain by Swanson in [41].

With the invention and adoption of the World Wide Web [4], researchers expanded their efforts to develop text and data mining methods to search and explore the web. This led to the creation of search engines, starting from simple interface-based browsers such as *Mosaic* and *Netscape*, and developed into large-scale systems such as *Google* and *Bing*. This, together with the advances in natural language processing, formed the opening gambit towards more complex approaches of text mining to different fields (such as the biomedical domain, economics, social sciences, or humanities), and types of data and content (such as scientific publications and research data, social media, multimedia, or public sector information).

The core research fields that develop, test, and publish the algorithms for TDM implementations for each specific domain of use, are

- *machine learning (ML)*, which is the automated learning of tasks and the discovery of patterns and structure in data [1];
- *information retrieval (IR) and search*, which focuses on the optimisation of search methods that, based on an index of vast amounts of content, allow for the location and extraction of the most relevant facts answering any question and user type [37]; and
- *natural language processing (NLP)* techniques, which deal with human language in its textual and spoken form [21, 27].

Research progress within TDM-related communities can be traced across two main sources: a number of high-impact yearly conferences (such as Knowledge Discovery and Data Mining (KDD), International Conference on Data Engineering (ICDE), IEEE International Conference on Data Mining (ICDM) for ML; ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR) and International Conference on Information and Knowledge Management (CIKM) for IR; Annual meeting of the Association for Computational Linguistics (ACL) and Conference on Empirical Methods on Natural Language Processing (EMNLP) for NLP), and a set of prominent journals in the fields (such as IEEE Transactions on Knowledge and Data Engineering (TKDE), Information Processing Letters (IPL)

for ML; Information Retrieval Journal for IR; Transactions of the Association for Computational Linguistics (TACL) for NLP). With the speed of technology development and the race for implementation of new methods in the context of practical implementations, the publications at conferences, having the fastest turnaround, have a relatively high impact. There is also a noteworthy trend to pre-print publications via open access self-archiving services such as <https://arxiv.org/>, that originally was created for the field of physics, but gradually expanded to include mathematics, computer science, nonlinear sciences, quantitative biology, quantitative finance, and statistics. This pre-print culture has an impact on the publications acceptance at the traditional publication venues that needs to take into account these additional aspects when the reviewing process is organised [42].

### 2.1. Leading algorithms

A top 10 of algorithms was defined by its leading experts in 2006 [46], and they are still representative for the situation in 2017, with the notable exception of Deep Learning as a method that was dominant in the early 1990s [24, 45, 38, 23], was largely absent in the 2000s, and which returned in the 2010s [10]. We list the 2006 top 10 of methods below, and in Figure 1 we show the growth of citations for the main reference scientific articles that defined those methods using the information available via the Google Scholar service.<sup>1</sup> Across the board, numbers of citations have increased. Support vector machines [44] were already popular in 2006 and took the position as most frequently cited learning algorithm in 2017.

- Classification: C4.5 [36]; CART [6]; K Nearest Neighbours (kNN) [16]; Naive Bayes [49];
- Statistical Learning: SVM [44]; EM [30];
- Association Analysis: Apriori [3]; FP-tree [15];
- Link Mining: PageRank [7]; HITS [22];
- Clustering: k-Means [1967]; BIRCH [50];
- Bagging and Boosting: AdaBoost [12];

---

<sup>1</sup><https://scholar.google.com/>

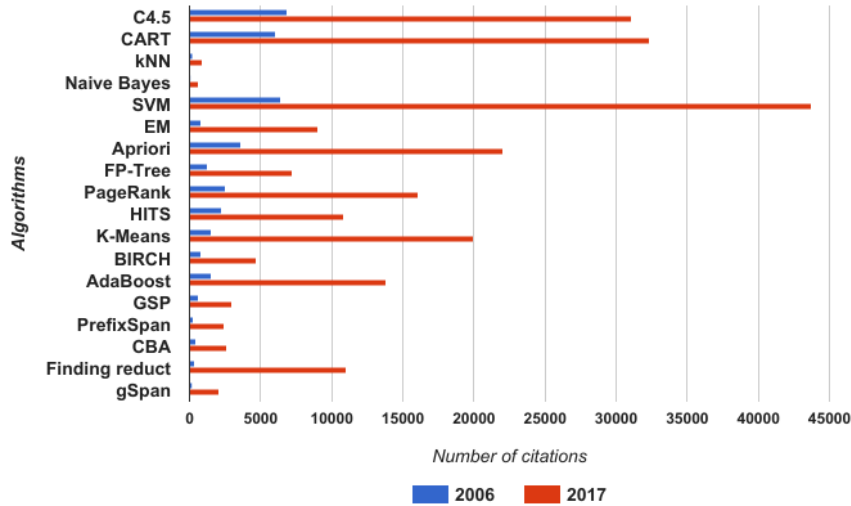


Figure 1: Comparison of the top 10 TDM papers citations: 2006 vs 2017 according to GoogleScholar.

- Sequential Patterns: GSP [39]; PrefixSpan [34];
- Integrated Mining: CBA [25];
- Rough Sets: Finding reduct [33];
- Graph Mining: gSpan [48].

## 2.2. NLP and IR perspective: main used algorithms and tasks

As IR has gradually developed from the techniques targeting simple keyword search [26] to more sophisticated use of language modeling [18] that requires speech and text pre-processing, and overall language modeling, the NLP and IR fields converge into sharing the algorithms used. Below we enlist the leading techniques used across those fields:

- Latent Dirichlet allocation (LDA) [5] for topic modeling and text classification;

- Hidden Markov Models (HMMs) [28], recurrent neural networks (RNNs) and Long short-term memory (LSTMs) for statistical language modeling as a component in speech recognition [19], machine translation [47], and information retrieval [18];
- Word2Vec [31] and GloVe [35] for text classification, information retrieval [13], machine translation [32];
- CRFs (conditional random fields) [40] and other sequence learners for part-of-speech tagging, Named Entity Recognition;
- CF (collaborative filtering) for recommendation [2];
- TF-IDF, BM25 for information retrieval [37, 27].

### **3. Tracing TDM application in different scientific areas and fields of activity**

The growth of data-driven business, services and research, has led to an introduction of an extended model for the economic structure that reflects a novel type of knowledge-based economy, where knowledge and data become a separate valuable commodity (OCDE, 1996). Thus recently, the research that supports knowledge sharing and growth, as well as education of the professionals that can carry out these activities, have been termed a quaternary economic sector, to be added to the traditional primary, secondary, and tertiary ones. Moreover, the high-level decision makers in governments, large industry companies, and education, having the potential to shape the future of the entire sectors with their vision and decisions, can arguably be placed in a separate, quinary sector.

We focus our investigation on the current and potential presence of TDM across all economic sectors, and thus we envisage these relations in Figure 2. The primary, secondary, and tertiary sectors follow roughly the same type of interaction with TDM technologies. These sectors represent sources of all possible types of data that are available for mining, and at the same time, have tasks and challenges that require smart solutions. The quaternary sector is particularly central to TDM development and implementation, as it produces the TDM experts and advances the technology itself. Experts at the decision-maker level use TDM tools to test and prove their vision of

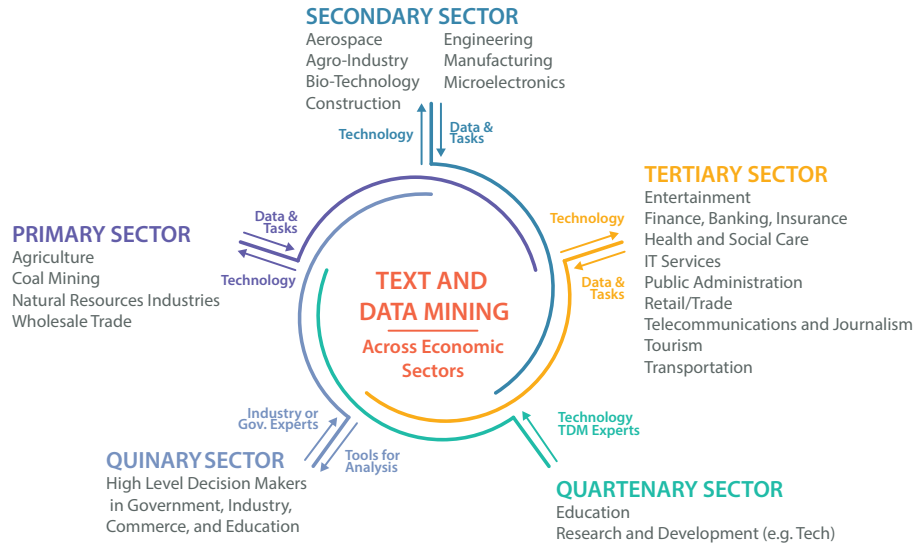


Figure 2: General economic structure and connection between TDM and all economic sectors

economic development, as well as to hypothesize where a specific decision will lead to, and to measure and assess its progress.

With the economic sectors and their relationship with TDM defined for any type of knowledge-based society, we take an example of the European Union infrastructure, and use TDM techniques to outline the scientific span of the TDM work and applications. We employ two axes, namely (i) EU funded research projects and infrastructures, and (ii) scientific articles published in journals and conferences. These two axes are considered as reflecting the strategic decisions of the EU and the policies adopted as regards the selection of specific research domains for funding, as well as the scientific trends attested in the publications.

### 3.1. Experimental framework

In order to estimate the distribution of scientific publications relating to Text and Data Mining per economic sector, we used 1,195,499 abstracts extracted from the latest dumps that have been made available at the CORE

(COncecting REpositories) site<sup>2</sup>. From the 1,195,499 abstracts in our dataset, 11504 are related to Text and Data Mining, i.e. 0,96% of the total. This finding is commensurate to the findings of [11] and [43]. These abstracts were analysed using LDA, a probabilistic Bayesian model of text generation. We used the MALLET<sup>3</sup> implementation of LDA and learnt a model of 100 topics from the 1,195,499 abstracts. Each topic  $t$  (in our analysis) is represented by a small set of words; the ones that are more likely to belong to  $t$  according to the estimated LDA probabilities (word-topic distribution). Based on the word set, each topic was manually assigned an application area label from the respective classification that is depicted in Figure 2.

From all abstracts, we kept (for our analysis) only the ones in which a TDM-related term occurs (e.g. Text Mining, Text Analysis, Text Analytics, Data Mining, Natural Language Processing)<sup>4</sup>. We used the topic-document probabilities that are returned from LDA and indicate which topics are discussed in each abstract. The sector label (e.g. Aerospace) of the LDA topic with the highest topic-document probability is assigned to each abstract.

### 3.2. Primary sector

The primary sector includes all raw materials and natural resources as well as the methods used for their mining (in the prototypical sense of the word) and/or production.

The main objectives of TDM research in the primary sector include, among others:

- standardisation, understanding and forecasting of natural phenomena (e.g. earthquakes) and climate changes interacting with the lives of humans and other species;
- prediction of better or worse/problematic crops or certain climate phenomena which affect them;
- correlation of pesticide use and animal/human diseases;
- better use of water resources/materials for production optimization.

---

<sup>2</sup><https://core.ac.uk/>

<sup>3</sup><http://mallet.cs.umass.edu/>

<sup>4</sup>The full list of terms is in the Appendix 8



### 3.2.1. Funding

A total amount of approximately €98,2M, with a European Commission maximum contribution of €75,3M, has been invested in 24 projects concerning TDM in the primary sector. Most projects focus on the collection and processing of data from natural resources, oceanic data, climate data or agricultural data. The following have been selected among the top twenty projects, based on their total investment, as representative cases for the primary sector:

- MARS: Managing Aquatic ecosystems and water Resources under multiple Stress;
- ADAPTAWHEAT: Genetics and physiology of wheat development to flowering: tools to breed for improved adaptation and yield potential;
- SUSTAINMED: Sustainable agri-food systems and rural development in the Mediterranean Partner Countries;

### 3.2.2. Publications

The majority of the publications of TDM interest (retrieved from the CORE repository) which fall in the primary sector deal with the areas of Agriculture and Environment; other areas (e.g. Natural Resources) are comparatively less researched into while some others (e.g. Wholesale Trade) are not represented in the publications' set, see Figure 3.

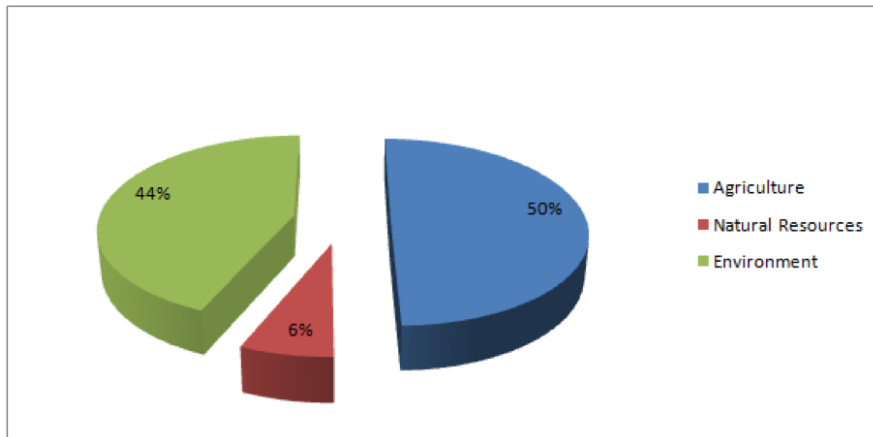


Figure 3: Publications in the primary sector.

### 3.3. Secondary sector

The secondary sector includes all the methods/techniques/tools/machinery used in order to transform the primary sector materials or resources into goods and products. All human activities related to manufacturing, processing and construction pertain to the secondary sector.

#### 3.3.1. Funding

A total amount of €806,4M, with a European Commission maximum contribution of €633,8M, has been invested on 253 projects concerning TDM in the secondary sector. The majority of projects focus on the collection and processing of biological and medical data for biotechnology applications and the subsequent provision of better health care services. Many projects concern the processing of energy data and the development of green or smart energy systems. The following projects are representative cases for the secondary sector:

- GEN2PHEN: Genotype-To-Phenotype Databases: A Holistic Solution;
- LinkedDesign: Linked Knowledge in Manufacturing, Engineering and Design for Next-Generation Production;
- Fortissimo 2: Factories of the Future Resources, Technology, Infrastructure and Services for Simulation and Modelling 2.

#### 3.3.2. Publications

Half the publications of TDM interest which fall in the secondary sector concern Engineering. The other areas of interest cover Aerospace and Energy followed by Construction and Microelectronics as depicted in Figure 4.

### 3.4. Tertiary sector

The tertiary sector includes all the services provided to individuals or organized groups such as businesses. All human activities related to retail and wholesale sales, transportation, entertainment, tourism, insurance, and many more pertain to the tertiary sector.

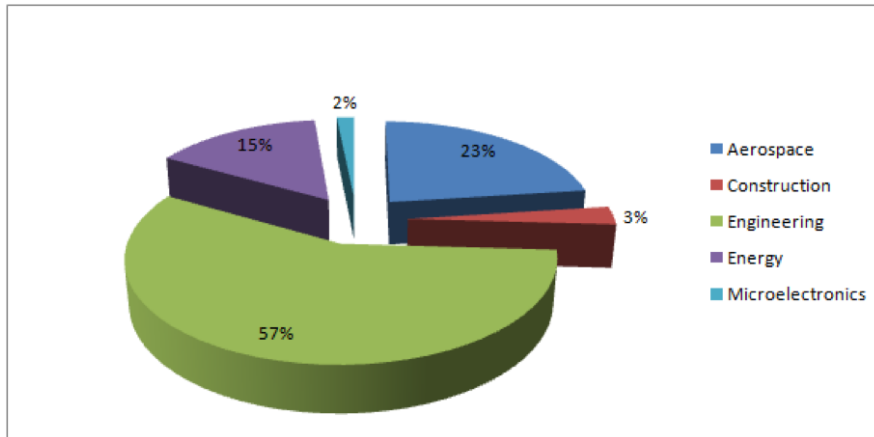


Figure 4: Publications in the secondary sector.

#### 3.4.1. Funding

A total amount of €1,1B, with a European Commission maximum contribution of €852,4M, has been invested on 313 projects concerning TDM in the tertiary sector. Most projects focus on the collection and processing of data from medical records and biology databases. The next category in frequency is that of IT services followed by the business application area. The following projects are representative cases for the tertiary sector:

- SENSEI: Integrating the Physical with the Digital World of the Network of the Future;
- SIIP: Speaker Identification Integrated Project;
- APO-SYS: “Apoptosis systems biology applied to cancer and AIDS. An integrated approach of experimental biology, data mining, mathematical modelling, biostatistics, systems engineering and molecular medicine”;
- mPlane: mPlane an Intelligent Measurement Plane for Future Network and Application Management;

#### 3.4.2. Publications

The retrieved publications of TDM interest for the tertiary sector cover a variety of fields as depicted in Figure 5. The majority of publications (60%) pertain to the fields of Medicine, Health Care and Social Care, while 30%

are related to generic IT services. The remaining publications are related to Finance, Public Administration, Entertainment and Transportation.

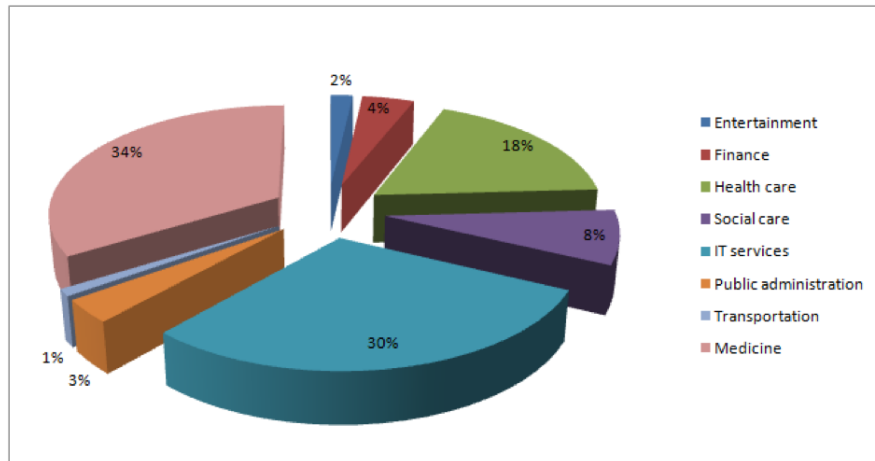


Figure 5: Publications in the tertiary sector.

### 3.5. Quaternary sector

The quaternary sector includes all services dealing with knowledge and information such as research, development and education among others. Research and development are terms applying to all scientific fields; therefore the following sections include projects, infrastructures, language resources and publications which could also be identified as belonging to any of the first three economic sectors but are presented here because the focus lies on the innovation, introduction, and improvement of products and processes of the various disciplines and application areas.

#### 3.5.1. Funding

A total amount of €402M, with a European Commission maximum contribution of €319,9M, has been invested on 273 projects concerning TDM in the quaternary sector. This sector includes all these projects that are classified by the EC as scientific research under RTD horizontal topics without further distinguishing, in most cases, the particular scientific area, the development of which the particular action will affect. As a result strong infrastructural projects like EUDAT2020 and projects requiring substantial financial investment in the area of Medicine and Health Sciences, e.g. the

METACARDIS project, are grouped together. The following are the representative examples for the tertiary sector:

- METACARDIS: Metagenomics in Cardiometabolic Diseases;
- SCY: Science Created by You;
- NEXT-TELL: Next Generation Teaching, Education and Learning for Life.

### 3.5.2. Publications

The publications concerning TDM for the quaternary sector, as depicted in Figure 6, fall into two big categories: the first being research, which is a generic term including many scientific fields, and the second one being education. The research publications presented here are those for which no specific scientific field was denoted from the set of words of topic  $t$ . Otherwise, the classification of the topic was done in accordance to the application areas of the primary, secondary or tertiary sector.

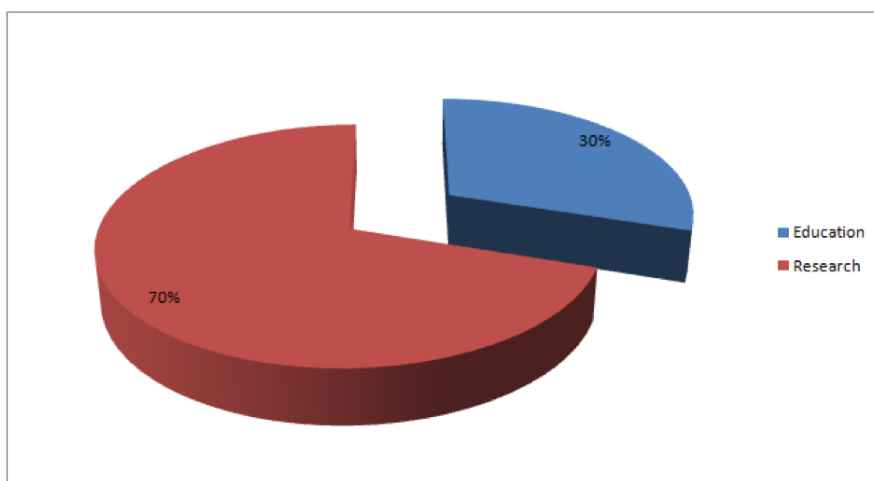


Figure 6: Publications in the quaternary sector.

### 3.6. Quinary sector

The quinary sector is the highest economic sector encompassing decision making for policy guidelines in Industry, Government, Science and Technology, having a profound impact on economy in general. It is the special

characteristics of this sector and the interrelation with all the other economic sectors, which make it challenging to isolate resources and material pertaining exclusively to the quinary sector itself. Decision makers need all the available data from various types of infrastructures as well as all the TDM software tools and technologies in order to analyse, deploy them, and finally make decisions affecting different application areas within all previous sectors. Taking these peculiarities into consideration, in this paper we present only the EU projects which have been funded under certain RTD Horizontal Topics, specifically those which aim at the formulation or identification of future developments in RTD and long-term strategic options.

### *3.6.1. Funding*

A total amount of €192,1M, with a European Commission maximum contribution of €149M, has been invested on 116 projects concerning TDM in the quinary sector. Most projects focus on the evaluation or application of laws and regulations, policies, policy strategies or plans of action for research and development in science and technology. The following examples show decision making projects that are the most representative for the quinary sector:

- KHRESMOI: Knowledge Helper for Medical and Other Information users
- LOD2: LOD2 - Creating Knowledge out of Interlinked Data
- FIRST: Large scale information extraction and integration infrastructure for supporting financial decision making

## **4. Conclusion**

In this paper we reviewed the most important algorithms that define the state-of-the-art and potential for advancement for the scientific development of TDM. Using LDA analysis on the CORE database, we highlighted the state of TDM research and funding coverage in context of the European Union test case, showcasing the spread of publications per sector and topics that are intertwined with TDM.

## 5. Acknowledgements

This research is funded by H2020 Project FutureTDM (call GARRI-3-2014, grant agreement 665940).

## 6. Competing interests

The authors have no competing interests.

## 7. References

- [1] Glossary of terms. *Mach. Learn.*, 30(2-3):271–274, February 1998.
- [2] Gediminas Adomavicius and Alexander Tuzhilin. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING*, 17(6):734–749, 2005.
- [3] Rakesh Agrawal and Ramakrishnan Srikant. Fast algorithms for mining association rules. In *Proc. of 20th Intl. Conf. on VLDB*, pages 487–499, 1994.
- [4] Tim Berners-Lee. The world wide web - past, present and future. *J. Digit. Inf.*, 1(1), 1997.
- [5] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003.
- [6] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, Monterey, CA, 1984.
- [7] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30(1):107 – 117, 1998.
- [8] Kenneth P. Burnham and David R. Anderson. *Model Selection and Multimodel Inference*. Springer-Verlag, 2002.
- [9] Frank Buytendijk. *Hype Cycle for Big Data*. Stamford: Gartner, 2014.
- [10] Li Deng and Dong Yu. Deep learning: Methods and applications. *Found. Trends Signal Process.*, 7(3&#8211;4):197–387, June 2014.

- [11] S. Filippov. Mapping text and data mining in academic and research communities in europe. Technical report, Nuffield College, Oxford, UK, Discussion paper, 2014.
- [12] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1):119–139, August 1997.
- [13] Debasis Ganguly, Dwaipayan Roy, Mandar Mitra, and Gareth J.F. Jones. Word embedding based generalized language model for information retrieval. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '15, pages 795–798, New York, NY, USA, 2015. ACM.
- [14] N. Gershenfeld. *The Nature of Mathematical Modeling*. Cambridge University Press, 1998.
- [15] Jiawei Han, Jian Pei, and Yiwen Yin. Mining frequent patterns without candidate generation. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, SIGMOD '00, pages 1–12, New York, NY, USA, 2000. ACM.
- [16] Trevor Hastie and Robert Tibshirani. Discriminant adaptive nearest neighbor classification. *IEEE Trans. Pattern Anal. Mach. Intell.*, 18(6):607–616, June 1996.
- [17] Marti A. Hearst. Untangling text data mining. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics on Computational Linguistics*, ACL '99, pages 3–10, Stroudsburg, PA, USA, 1999. Association for Computational Linguistics.
- [18] Djoerd Hiemstra. *Using language models for information retrieval*. Univ. Twente, 2001.
- [19] Xuedong Huang, Yasuo Ariki, and Mervyn Jack. *Hidden Markov Models for Speech Recognition*. Columbia University Press, New York, NY, USA, 1990.
- [20] A. B. Jones and J. M. Smith. Big data and journalism: Epistemology, expertise, economics, and ethics. *Digital Journalism*, 3:447466, 2014.



- [21] Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 1st edition, 2000.
- [22] Jon M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, September 1999.
- [23] Yann LeCun and Yoshua Bengio. The handbook of brain theory and neural networks. chapter Convolutional Networks for Images, Speech, and Time Series, pages 255–258. MIT Press, Cambridge, MA, USA, 1998.
- [24] Yann Lecun, Lon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324, 1998.
- [25] Bing Liu, Wynne Hsu, and Yiming Ma. Integrating classification and association rule mining. In *KDD*, pages 80–86, 1998.
- [26] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA, 2008.
- [27] Christopher D. Manning and Hinrich Schütze. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA, USA, 1999.
- [28] A. Markov. Extension of the limit theorems of probability theory to a sum of variables connected in a chain. In R. Howard, editor, *Dynamic Probabilistic Systems (Volume I: Markov Models)*, chapter Appendix B, pages 552–577. John Wiley & Sons, Inc., New York City, 1971.
- [29] V. Mayer-Schönberger and K. Cukier. *Big Data: A Revolution That Will Transform How We Live, Work and Think*. John Murray Publishers, 2013.
- [30] G. J. McLachlan and D. Peel. *Finite mixture models*. Wiley Series in Probability and Statistics, New York, 2000.
- [31] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.

- [32] Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. Exploiting similarities among languages for machine translation. *CoRR*, abs/1309.4168, 2013.
- [33] Zdzislaw Pawlak. *Rough Sets: Theoretical Aspects of Reasoning About Data*. Kluwer Academic Publishers, Norwell, MA, USA, 1992.
- [34] Jian Pei, Jiawei Han, Behzad Mortazavi-Asl, Helen Pinto, Qiming Chen, Umeshwar Dayal, and Meichun Hsu. Prefixspan: Mining sequential patterns by prefix-projected growth. In *Proceedings of the 17th International Conference on Data Engineering*, pages 215–224, Washington, DC, USA, 2001. IEEE Computer Society.
- [35] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, volume 14, pages 1532–1543, 2014.
- [36] J. Ross Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993.
- [37] C. J. Van Rijsbergen. *Information Retrieval*. Butterworth-Heinemann, Newton, MA, USA, 2nd edition, 1979.
- [38] P. Smolensky. Tensor product variable binding and the representation of symbolic structures in connectionist systems. *Artif. Intell.*, 46(1-2):159–216, November 1990.
- [39] Ramakrishnan Srikant and Rakesh Agrawal. Mining sequential patterns: Generalizations and performance improvements. In *Proceedings of the 5th International Conference on Extending Database Technology: Advances in Database Technology*, EDBT '96, pages 3–17, London, UK, UK, 1996. Springer-Verlag.
- [40] Charles Sutton and Andrew McCallum. An introduction to conditional random fields. *Found. Trends Mach. Learn.*, 4(4):267–373, April 2012.
- [41] D.R. Swanson. Fish oil, raynauds syndrome, and undiscovered public knowledge. *Perspect. Bio. Med.*, 30:7–18, 1986.
- [42] Andrew Tomkins, Min Zhang, and William D. Heavlin. Single versus double blind reviewing at WSDM 2017. *CoRR*, abs/1702.00502, 2017.

- [43] Hsu-Hao Tsai. Global data mining: An empirical study of current trends, future forecasts and technology diffusions. *Expert Syst. Appl.*, 39(9):8172–8181, July 2012.
- [44] Vladimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag New York, Inc., New York, NY, USA, 1995.
- [45] Alexander Waibel, Toshiyuki Hanazawa, Geoffrey Hinton, Kiyohiro Shikano, and Kevin J. Lang. Readings in speech recognition. chapter Phoneme Recognition Using Time-delay Neural Networks, pages 393–404. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1990.
- [46] Xindong Wu, Vipin Kumar, J. Ross Quinlan, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J. McLachlan, Angus Ng, Bing Liu, Philip S. Yu, Zhi-Hua Zhou, Michael Steinbach, David J. Hand, and Dan Steinberg. Top 10 algorithms in data mining. *Knowl. Inf. Syst.*, 14(1):1–37, December 2007.
- [47] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, 2016.
- [48] Xifeng Yan and Jiawei Han. gspan: Graph-based substructure pattern mining. In *Proceedings of the 2002 IEEE International Conference on Data Mining*, ICDM ’02, pages 721–, Washington, DC, USA, 2002. IEEE Computer Society.
- [49] Keming Yu and David J. Hand. Idiot’s bayesnot so stupid after all? *International Statistical Review*, 69(3):385–398, 2001.
- [50] Tian Zhang, Raghu Ramakrishnan, and Miron Livny. BIRCH: An efficient data clustering method for very large databases. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data*, SIGMOD ’96, pages 103–114, New York, NY, USA, 1996. ACM.

8. Appendix: List of TDM related terms and their frequencies in FP7 and HORIZON2020 projects

| Keywords in OBJECTIVE     | FP7 | keywords in OBJECTIVE     | Horizon 2020 |
|---------------------------|-----|---------------------------|--------------|
| data analysis             | 240 | big data                  | 66           |
| machine learning          | 156 | data analysis             | 61           |
| data mining               | 88  | machine learning          | 51           |
| large data                | 65  | data analytics            | 21           |
| language processing       | 39  | data mining               | 14           |
| big data                  | 38  | large data                | 14           |
| information retrieval     | 37  | linked data               | 6            |
| information extraction    | 30  | business intelligence     | 5            |
| linked data               | 30  | data science              | 5            |
| text mining               | 17  | language processing       | 5            |
| data science              | 19  | information retrieval     | 4            |
| data analytics            | 18  | predictive analytics      | 4            |
| information access        | 17  | information access        | 3            |
| summarization             | 12  | sentiment analysis        | 2            |
| business intelligence     | 11  | text and data mining      | 2            |
| knowledge discovery       | 11  | information extraction    | 1            |
| multimedia processing     | 8   | knowledge discovery       | 1            |
| conversation analysis     | 7   | summarization             | 1            |
| opinion mining            | 7   | TDM                       | 1            |
| sentiment analysis        | 5   | abbreviation detection    | 0            |
| semantic mining           | 4   | competitive intelligence  | 0            |
| text analytics            | 3   | concept extraction        | 0            |
| content mining            | 3   | content classification    | 0            |
| text understanding        | 3   | content mining            | 0            |
| trend detection           | 2   | conversation analysis     | 0            |
| graph mining              | 2   | datafication              | 0            |
| predictive analytics      | 2   | graph mining              | 0            |
| sense disambiguation      | 2   | linguistic identification | 0            |
| content classification    | 2   | modality detection        | 0            |
| textual entailment        | 1   | multimedia processing     | 0            |
| datafication              | 0   | negation detection        | 0            |
| concept extraction        | 0   | opinion mining            | 0            |
| modality detection        | 0   | semantic mining           | 0            |
| negation detection        | 0   | sense disambiguation      | 0            |
| competitive intelligence  | 0   | text analytics            | 0            |
| linguistic identification | 0   | text mining               | 0            |
| abbreviation detection    | 0   | text understanding        | 0            |
| text and data mining      | 0   | textual entailment        | 0            |
| TDM                       | 0   | trend detection           | 0            |

Collection of Organisations

| Name   | URL   | Country                                | Zip        | City            | Street                                 | Type     | Data Sources   | Sector   | Application Field Sector I  | Application Field Sector II | Application Field Sector III | Application Field Sector IV   | Application Field Sector V  |
|--|---|--|------------|-----------------|--|----------|--|--|-----------------------------|-----------------------------|------------------------------|---|---|
| The Hub  | <a href="http://thehub.com">http://thehub.com</a>   | The Netherlands                        | 1088XK     | Amsterdam       | Science Park 402                       | Company  | Business, Internet   | Tertiary Sector, Quaternary Sector                   | N/A                         | N/A                         | N/A                          | Finance, Banking, Insurance   | Research and Development (e.g. Tech)  |
| Accenture  | <a href="https://www.accenture.com">https://www.accenture.com</a>                                     | USA, many interns                      | N/A        | N/A             | N/A                                    | Company  | Business, Internet   | N/A  | N/A                         | N/A                         | N/A                          | Health and Social Care  | Research and Development (e.g. Tech)  |
| ADAPT  | <a href="http://adaptcentre.ie">http://adaptcentre.ie</a>   | Ireland                                | Dublin 2   | Dublin          | Trinity College, College Green         | Research | All  | Tertiary Sector, Quaternary Sector                   | N/A                         | N/A                         | N/A                          | Entertainment, Finance  | Education, Research and Development (e.g. Tech)                               |
| Agri   | <a href="http://agri.tycho.org/tycho-search/index.do">http://agri.tycho.org/tycho-search/index.do</a> | Greece                                 | 10776      | Athens          | Acharon Street 7                       | Company  | Scientific   | Primary Sector                                       | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| AgriKnow   | <a href="http://www.agriknow.com/agriknow/">http://www.agriknow.com/agriknow/</a>                     | Greece                                 | 152 35     | Athens          | Grimoussi 17                           | Company  | Scientific   | Primary Sector                                       | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Algorithmia  | <a href="http://www.algorithmia.com/">http://www.algorithmia.com/</a>                                 | The Netherlands                        | 1075       | Amsterdam       | Orange Nassaubaan 62-2                 | Company  | All  | Secondary Sector                                     | N/A                         | N/A                         | N/A                          | Engineering   | N/A   |
| Alteos Limited   | <a href="http://www.alteos.com/">http://www.alteos.com/</a>   | Spain                                  | 37900      | Salamanca       | Carretera de Madrid 13                 | Company  | Business, Personal   | Secondary Sector, Tertiary Sector, Quaternary Sector | N/A                         | N/A                         | N/A                          | Manufacturing   | Entertainment, Finance, Banking, Insurance, Retail/Trade, Telecommunications  |
| Artificial Intelligence Techniques Ltd                                   | <a href="http://www.aitech.com/">http://www.aitech.com/</a>   | UK                                     | CB2 0WT    | Cambridge       | Cambridge Business Park, Cowley Park   | Company  | All  | Secondary Sector, Tertiary Sector                    | N/A                         | N/A                         | N/A                          | Manufacturing   | Health and Social Care  |
| Autonomy   | <a href="http://www.autonomy.com/">http://www.autonomy.com/</a>                                       | Spain/UK                               | N/A        | N/A             | N/A                                    | Company  | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | Finance, Banking, Insurance   | Health and Social Care, Public Administration                                 |
| AVUKI  | <a href="http://www.avuki.com/">http://www.avuki.com/</a>   | USA                                    | N/A        | Columbus, Ohio  | N/A                                    | Company  | Business, Scientific   | Secondary Sector, Quaternary Sector                  | N/A                         | N/A                         | N/A                          | Bio-technology, Manufacturing   | Research and Development (e.g. Tech)  |
| Bailela  | <a href="http://www.bailela.com/index.html">http://www.bailela.com/index.html</a>                     | Austria                                | 1040       | Vienna          | Schleierhofgasse 7                     | Company  | Business, Public Sector  | Secondary Sector, Tertiary Sector, Quaternary Sector | N/A                         | N/A                         | N/A                          | Finance, Banking, Insurance   | Education   |
| BGS Stat   | <a href="http://www.bgsstat.com/">http://www.bgsstat.com/</a>   | Denmark                                | 2300       | Copenhagen      | Njalsgade 136, building 27             | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Research and Development (e.g. Tech)  |
| Centre for Language Technology (CST)                                     | <a href="http://csl.uu.se/csl/">http://csl.uu.se/csl/</a>   | Sweden                                 | 7000       | Uppsala         | Uppsala University                     | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Research and Development (e.g. Tech)  |
| ClearPath Clearance Center, Inc.   | <a href="http://www.clearpath.com/">http://www.clearpath.com/</a>                                     | Switzerland                            | CH-4603    | Basel           | Murgengrabenstrasse 47                 | Company  | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | Bio-technology  | Research and Development (e.g. Tech)  |
| CrossRef   | <a href="http://www.crossref.org">http://www.crossref.org</a>   | UK/USA                                 | OX1 1BF    | Oxford          | New Road                               | Research | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | Agro-industry, Bio-technology   | Finance, Banking, Insurance   |
| CrossRef   | <a href="http://www.crossref.org">http://www.crossref.org</a>   | Czech Republic                         | 170 00     | Prague          | U Prichonu 22 / 466                    | Company  | All  | Tertiary Sector                                      | N/A                         | N/A                         | N/A                          | Health and Social Care  | Research and Development (e.g. Tech)  |
| DeepMind   | <a href="https://www.deepmind.com/">https://www.deepmind.com/</a>                                     | UK/USA                                 | N/A        | London          | N/A                                    | Company  | Business   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Research and Development (e.g. Tech)  |
| Deutsche Forschungszentrum für Künstliche Intelligenz gmbh (DFKI)        | <a href="http://www.dfki.de">http://www.dfki.de</a>   | Germany                                | D-66113    | Saarbrücken     | Stuhlnissenweg 3                       | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Deezer   | <a href="http://www.deezer.com/">http://www.deezer.com/</a>   | N/A                                    | N/A        | N/A             | N/A                                    | Company  | All  | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Dutch-Flemish Human Language Technology Agency ePronomat                 | <a href="http://www.epronomat.com/">http://www.epronomat.com/</a>                                     | The Netherlands, Belgium               | N/A        | N/A             | N/A                                    | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| ePronomat  | <a href="http://www.epronomat.com/PrdctyIndex.html">http://www.epronomat.com/PrdctyIndex.html</a>     | Czech Republic                         | N/A        | N/A             | N/A                                    | Company  | Business   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| EpiGene  | <a href="http://www.epigenome.com/">http://www.epigenome.com/</a>                                     | Spain                                  | N/A        | N/A             | N/A                                    | Company  | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Exclusiv Ltd.  | <a href="http://www.exclusiv.com/">http://www.exclusiv.com/</a>                                       | Ireland/Greece                         | N/A        | N/A             | N/A                                    | Company  | Business   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Expert Systems S.P.A.  | <a href="http://www.expert-systems.com">http://www.expert-systems.com</a>                             | Italy                                  | 38122      | Boziano         | Via Santa Croce 77                     | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | Bio-technology  | Entertainment, Finance, Banking, Insurance, Public Administration             |
| Fondazione Bruno Kessler (FBK)   | <a href="http://www.fbk.it/">http://www.fbk.it/</a>   | Italy                                  | 38122      | Boziano         | Via Santa Croce 77                     | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| German National Library of Economics (ZBW)                               | <a href="http://www.zbw.eu/">http://www.zbw.eu/</a>   | Germany                                | 68159      | Manheim         | Quadrat 82, 1                          | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Research and Development (e.g. Tech)  |
| GESIS Leibniz-Institut für die Sozialwissenschaften                      | <a href="http://www.gesis.org/">http://www.gesis.org/</a>   | Germany                                | 68159      | Manheim         | Quadrat 82, 1                          | Research | Personal, Scientific   | Tertiary Sector                                      | N/A                         | N/A                         | N/A                          | N/A   | Health and Social Care  |
| Heinrich Heine University of Technology                                  | <a href="http://www.aahp.fi/en/about/history/">http://www.aahp.fi/en/about/history/</a>               | Finland                                | 2150       | Helsinki        | Espoo                                  | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Institut national de recherche en informatique et en automatique (INRIA) | <a href="http://www.inria.fr/">http://www.inria.fr/</a>   | France                                 | 78153      | Roquecournot    | BP 105                                 | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Institut Universitari de Lingüística Aplicada (IULA)                     | <a href="http://www.iula.com/">http://www.iula.com/</a>   | Spain                                  | 8018       | Barcelona       | Universitat Pompeu Fabra Campus del    | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Research and Development (e.g. Tech)  |
| Institute for Language and Speech Processing (ILSP / "Athens" R.C.)      | <a href="http://www.ilsp.gr/">http://www.ilsp.gr/</a>   | Greece                                 | 151 25     | Athens          | Marousi, Artemidos 6                   | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Intituto di Linguistica Computazionale "Antonio Zampolli" (ILC - CNR)    | <a href="http://www.ils.cnr.it/">http://www.ils.cnr.it/</a>   | Italy                                  | 86124      | Pisa            | Italian National Research Council/Pisa | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| JISC   | <a href="http://www.jisc.ac.uk/">http://www.jisc.ac.uk/</a>   | UK                                     | 852 6BA    | Bristol         | One Castlepark, Tower Hill             | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Karlsruhe Institute of Technology  | <a href="http://www.kit.edu/">http://www.kit.edu/</a>   | Germany                                | 76131      | Karlsruhe       | Kaiserstraße 12                        | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Komatsu  | <a href="http://www.komatsu.com/en/aboutus/">http://www.komatsu.com/en/aboutus/</a>                   | Japan                                  | N/A        | Toyko           | N/A                                    | Company  | N/A  | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Kudoss   | <a href="http://www.kudoss.nl">http://www.kudoss.nl</a>   | The Netherlands                        | N/A        | N/A             | N/A                                    | Company  | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| LinguaMatics   | <a href="http://www.lingumatics.com/">http://www.lingumatics.com/</a>                                 | UK                                     | CB8 0WG    | Cambridge       | Milton Road                            | Company  | Personal, Scientific   | Secondary Sector, Tertiary Sector, Quaternary Sector | N/A                         | N/A                         | N/A                          | Bio-technology  | Health and Social Care  |
| Ludwig Maximilians Universität   | <a href="http://www.lmu-muenchen.de/index">http://www.lmu-muenchen.de/index</a>                       | Germany                                | 80339      | Munich          | Geschwister Scholl Platz 1             | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Max Planck Computing and Data Facility                                   | <a href="http://www.rig.mpg.de/">http://www.rig.mpg.de/</a>   | Germany                                | 85748      | Garching        | Gießenbachstraße 2                     | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Research and Development (e.g. Tech)  |
| Mata   | <a href="http://www.mata.com/">http://www.mata.com/</a>   | USA                                    | N/A        | Washington, DC  | N/A                                    | Company  | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| National Centre for Text Mining (NuCTeM)                                 | <a href="http://www.nuitem.com/">http://www.nuitem.com/</a>   | UK                                     | M1 7DN     | Manchester      | Princess Street 111                    | Company  | Business, Internet, Personal, Secondary Sector, Tertiary Sector, Quaternary Sector | N/A  | N/A                         | N/A                         | N/A                          | Health and Social Care  | Education, Research and Development (e.g. Tech)                               |
| OpenDataSoft   | <a href="http://www.opendatasoft.com/">http://www.opendatasoft.com/</a>                               | USA/International                      | N/A        | N/A             | N/A                                    | Company  | Business, Internet   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Paperkit   | <a href="http://www.paperkit.com/">http://www.paperkit.com/</a>                                       | Switzerland                            | N/A        | N/A             | N/A                                    | Company  | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Piati  | <a href="http://www.piati.com/">http://www.piati.com/</a>   | Norway                                 | 855        | Oslo            | Sognsvien 70A                          | Company  | Personal, Scientific   | Tertiary Sector, Quaternary Sector                   | N/A                         | N/A                         | N/A                          | N/A   | Health and Social Care  |
| PuGene   | <a href="http://www.pugene.com/">http://www.pugene.com/</a>   | Belgium                                | 91150      | Brussels        | Avenue Roger Vandendriessche 9         | Company  | Business, Personal   | Tertiary Sector                                      | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Python predictions   | <a href="http://www.pythonpredictions.com/">http://www.pythonpredictions.com/</a>                     | Germany                                | 44227      | Dortmund        | Stoekumer Str. 475                     | Company  | Business, Internet   | F, Secondary Sector, Tertiary Sector                 | N/A                         | N/A                         | N/A                          | Entertainment, Finance, Banking, Insurance, IT Services, Retail/Trade | N/A   |
| RapidMiner   | <a href="http://www.rapidminer.com/">http://www.rapidminer.com/</a>                                   | Sweden                                 | 114 28     | Stockholm       | Brinellvägen 8                         | Company  | Business, Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| Raytheon Institute of Technology (KIT)                                   | <a href="http://www.kit.com/">http://www.kit.com/</a>   | USA                                    | 69190      | Waldorf         | Hasso Plattner Ring 7                  | Company  | All  | Primary Sector                                       | Natural Resources Mining, W | Engineering, Manufacturing  | Entertainment, Finance       | Research and Development (e.g. Tech)                                  | N/A   |
| Salesforce Ltd   | <a href="http://www.salesforce.com/">http://www.salesforce.com/</a>                                   | USA                                    | N/A        | Washington, DC  | N/A                                    | Research | N/A  | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| SAP  | <a href="http://www.sap.com/">http://www.sap.com/</a>   | The Netherlands                        | 10980H     | Amsterdam       | Science Park 402                       | Company  | Business, Internet   | F, Tertiary Sector, Quaternary Sector                | N/A                         | N/A                         | N/A                          | Health and Social Care, R   | Research and Development (e.g. Tech)  |
| Scholarly Publishing and Academic Resources Coalition (SPARC)            | <a href="http://www.sparc.org/">http://www.sparc.org/</a>   | USA                                    | 20036      | Washington, DC  | N/A                                    | Research | N/A  | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Slovenian Social Science Data Archives                                   | <a href="http://www.ssdas.com/">http://www.ssdas.com/</a>   | Slovenia                               | 1000       | Ljubljana       | Kongresni trg 12                       | Research | Personal, Scientific   | Tertiary Sector, Quaternary Sector                   | N/A                         | N/A                         | N/A                          | Health and Social Care  | Research and Development (e.g. Tech)  |
| Spanish  | <a href="http://www.spanish.com/">http://www.spanish.com/</a>   | USA                                    | N/A        | N/A             | N/A                                    | Company  | Business, Scientific   | Secondary Sector, Quaternary Sector                  | N/A                         | N/A                         | N/A                          | Bio-technology  | Education, Research and Development (e.g. Tech)                               |
| Språkbanken, The Swedish language bank                                   | <a href="http://www.sprakbanken.se">http://www.sprakbanken.se</a>                                     | Sweden                                 | N/A        | N/A             | N/A                                    | Research | Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Squares on Blue  | <a href="http://www.squaresonblue.com/">http://www.squaresonblue.com/</a>                             | France                                 | 78000      | Versailles      | Rue Saint Honoré 13                    | Company  | All  | Tertiary Sector                                      | N/A                         | N/A                         | N/A                          | Finance, Banking, Insurance   | Health and Social Care, Public Administration                                 |
| Statistik  | <a href="http://www.statistik.com/">http://www.statistik.com/</a>                                     | Switzerland                            | 3018       | Bern            | Morgenstrasse 129                      | Company  | Business, Internet, Personal, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | Finance, Banking, Insurance   | Health and Social Care  |
| Strategic Feed   | <a href="http://www.strategicfeed.com/">http://www.strategicfeed.com/</a>                             | France                                 | 92 088     | Paris           | Place de la Pyramide 5                 | Company  | Business, Personal   | Tertiary Sector, Quaternary Sector                   | N/A                         | N/A                         | N/A                          | N/A   | Health and Social Care  |
| Suhoi Centre of Expertise in the Social Sciences                         | <a href="http://www.suhoi.com/">http://www.suhoi.com/</a>   | Switzerland                            | CH-1015    | Lausanne        | Bâtimet Gloggi                         | Research | Personal, Scientific   | Tertiary Sector, Quaternary Sector                   | N/A                         | N/A                         | N/A                          | N/A   | Health and Social Care  |
| Synthetic Partners   | <a href="http://www.syntheticpartners.com/en/">http://www.syntheticpartners.com/en/</a>               | Spain                                  | 28006      | Madrid          | Calle Velázquez 92                     | Company  | Business, Internet   | F, Tertiary Sector                                   | N/A                         | N/A                         | N/A                          | Finance, Banking, Insurance   | Health and Social Care, IT Services, Public Admin                             |
| TeamGig  | <a href="http://www.teamgig.com/">http://www.teamgig.com/</a>   | The Netherlands                        | NL-1072 AB | Amsterdam       | Nieuwendammerdijk 28A-17               | Company  | Business, Internet   | F, Secondary Sector, Tertiary Sector                 | N/A                         | N/A                         | N/A                          | Manufacturing   | Finance, Banking, Insurance, IT Services, Public Administration, Retail/Trade |
| They say   | <a href="http://www.theysay.com/">http://www.theysay.com/</a>   | UK                                     | EC2V 3ND   | London          | Canal 24                               | Company  | Business, Internet   | Secondary Sector                                     | N/A                         | N/A                         | N/A                          | N/A   | Finance, Banking, Insurance, Retail/Trade                                     |
| Translight   | <a href="http://www.translight.com/">http://www.translight.com/</a>                                   | Germany                                | D-01307    | Dresden         | Tatzberg 47.51                         | Company  | Business, Internet   | F, Tertiary Sector, Quaternary Sector                | N/A                         | N/A                         | N/A                          | Health and Social Care  | Research and Development (e.g. Tech)  |
| Transliminer   | <a href="http://www.transliminer.com/">http://www.transliminer.com/</a>                               | Belgium, The Netherlands, Germany, USA | N/A        | N/A             | N/A                                    | Company  | Business, Scientific   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| UK Data Archive  | <a href="http://www.ukda.ac.uk/">http://www.ukda.ac.uk/</a>   | UK                                     | CO4 3UJ    | Essex           | Wivenhoe Park                          | Research | All  | Tertiary Sector, Quaternary Sector                   | N/A                         | N/A                         | N/A                          | N/A   | Finance, Banking, Insurance   |
| University of Bristol  | <a href="http://www.bristol.ac.uk/">http://www.bristol.ac.uk/</a>                                     | UK                                     | BS8 1TH    | Bristol         | Tyndall Avenue                         | Research | Scientific   | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Education, Research and Development (e.g. Tech)                               |
| USRD   | <a href="http://www.usrd.com/">http://www.usrd.com/</a>   | Denmark                                | 1190       | Brussels        | Boulevard de l'Atomité                 | Company  | Business   | N/A  | N/A                         | N/A                         | N/A                          | N/A   | N/A   |
| Verisk   | <a href="http://www.verisk.com/">http://www.verisk.com/</a>   | USA                                    | NI 07310   | 188 Jersey City | 545 Washington Boulevard               | Company  | All  | N/A  | N/A                         | N/A                         | N/A                          | N/A   | Finance, Banking, Insurance   |
| Verisk   | <a href="http://www.verisk.com/">http://www.verisk.com/</a>   | USA                                    | UT 84113   | 95555 Lake City | N/A                                    | Company  | All  | Quaternary Sector                                    | N/A                         | N/A                         | N/A                          | N/A   | Research and Development (e.g. Tech)  |
| VeriTrails   | <a href="http://www.veritrails.com/">http://www.veritrails.com/</a>                                   | USA                                    | 38240      | Meriden         | 6 chemin de Maupertuis                 | Research | Business, Internet   | F, Tertiary Sector                                   | N/A                         | N/A                         | N/A                          | N/A   | Finance, Banking, Insurance   |
| Xerox Lab  | <a href="http://www.xerox.com/">http://www.xerox.com/</a>   | USA                                    | 38240      | Meriden         | 6 chemin de Maupertuis                 | Research | Business, Internet   | F, Tertiary Sector                                   | N/A                         | N/A                         | N/A                          | N/A   | Finance, Banking, Insurance   |



Collection of TDM Tools

| Name                       | URL   | Type                                    | Organisation/Producer                             | Country       | Data Sources     | Sector                                      | Application/Field                           | Application/Field  | Application/Field   | Application/Field | Application/Field | Methods | Projects | Experts |
|----------------------------|---|---|---|---------------|------------------|---|---|--|---|-------------------|-------------------|---------|----------|---------|
| ALVID                      | <a href="http://alvid.edu.au/">http://alvid.edu.au/</a>   | Infrastructure, Platform                | National Collaborative Research Infrastructure    | Australia     | Business, Interf | Tertiary-Sector, Quaternary-Sector          | Entertainment, Research and Development (e. | Data Mining, Text Mining                                     | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Alix                       | <a href="https://axois.org/axois/fr/redmine/projects/">https://axois.org/axois/fr/redmine/projects/</a>   | Open Source, Component Collection, Work | N/A   | N/A           | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Angoss Knowledge STUDIO    | <a href="http://www.angoss.com/">http://www.angoss.com/</a>   | Commercial                              | Angoss Software Corporation                       | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Multimedia Processing                           | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Anomarkit                  | <a href="https://www.anomarkit.eu/">https://www.anomarkit.eu/</a>   | Commercial, Web Service Registry        | University of Sheffield                           | UK            | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Text Mining                                     | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Apache CTAKES              | <a href="http://ctakes.apache.org/">http://ctakes.apache.org/</a>   | Open Source, Component Collection, Work | Children's Hospital Boston, Mayo Clinic           | USA           | Business, Scien  | Secondary-Sector, Tertiary-Sec              | Bio-technology Health and Social Care       | Information Extraction, Natural Language Processing, Text    | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Apache OpenNLP             | <a href="http://opennlp.apache.org/">http://opennlp.apache.org/</a>   | Open Source, Component Collection, Work | IBM, Apache Software Foundation                   | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | Natural Language Processing                                  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Apache UIMA                | <a href="http://uima.apache.org/">http://uima.apache.org/</a>   | Open Source, Component Collection, Work | IBM, Apache Software Foundation                   | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | Information Extraction, Text Mining                          | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| App                        | <a href="https://app.socran.us/uk/">https://app.socran.us/uk/</a>   | Open Source, Component Collection, Work | The National Centre for Text Mining, The Univer   | USA           | N/A              | N/A   | N/A   | Text Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| BioCatalogue               | <a href="http://www.biocatalogue.org/">http://www.biocatalogue.org/</a>   | Open Source, Web Service Registry       | The University of Manchester                      | UK            | Scientific       | Secondary-Sector, Quaternary-Bio-technology | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| BiodiversityCatalogue      | <a href="http://www.biodiversitycatalogue.org/">http://www.biodiversitycatalogue.org/</a>   | Open Source, Web Service Registry       | The University of Manchester                      | UK            | Scientific       | Secondary-Sector, Quaternary-Bio-technology | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Blumen                     | <a href="http://www.blumen.ch/">http://www.blumen.ch/</a>   | Open Source, Component Collection, Work | EPFL - Research Institute                         | Switzerland   | Business, Scien  | Secondary-Sector, Quaternary-Bio-technology | Research and Development (e.                | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| CLARIN-DK                  | <a href="http://clarin.dk/">http://clarin.dk/</a>   | Webapplication, Workflow                | N/A   | EU            | All              | Quaternary-Sector                           | Research and Development (e.                | Natural Language Processing, Term/Concept Extraction         | CLARIN  | N/A               | N/A               | N/A     | N/A      | N/A     |
| ClearTK                    | <a href="http://clear.tkhub.io/clear/">http://clear.tkhub.io/clear/</a>   | Open Source, Component Collection, Work | Center for Computational Language and Education   | USA           | N/A              | N/A   | N/A   | Machine Learning   | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Clementine                 | <a href="https://www.ibm.com/software/ai/clementine/">https://www.ibm.com/software/ai/clementine/</a>   | Commercial                              | IBM   | USA           | All              | Quaternary-Sector                           | Research and Development (e.                | Artificial Intelligence, Data Mining, Information Extraction | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| ContentMine                | <a href="http://www.contentmine.com/">http://www.contentmine.com/</a>   | Open Source, Webapplication, Pipeline   | ContentMine (non-profit)                          | UK/EU         | Scientific       | Tertiary-Sector, Quaternary-Sector          | Education, Research and Devel               | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Data Mouts                 | <a href="http://datamouts.eu/">http://datamouts.eu/</a>   | Open Source                             | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| DeFlow Core                | <a href="https://github.com/DeFlow/DeFlow-Core">https://github.com/DeFlow/DeFlow-Core</a>   | Open Source, Component Collection, Work | Linguistics Knowledge Processing Lab (LTKP)       | Te Germany    | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Egis                       | <a href="https://www.egis.com/egis/">https://www.egis.com/egis/</a>   | Open Source, Standalone                 | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| EUKI                       | <a href="http://euki.de/limu.de/">http://euki.de/limu.de/</a>   | Open Source, Workflow                   | Ludwig Maximilian University of Munich            | Germany       | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Equipomat                  | <a href="http://www.equipomat.com/">http://www.equipomat.com/</a>   | Commercial                              | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Epinium                    | <a href="http://epinium.com/">http://epinium.com/</a>   | Commercial, Workflow, Standalone        | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Exclusiv                   | <a href="http://www.exclusiv.com/">http://www.exclusiv.com/</a>   | Open Source, Workflow Environment       | N/A   | USA           | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| FIRO Data Management       | <a href="http://www.firo.com/">http://www.firo.com/</a>   | Commercial                              | Fair Isaac Corporation                            | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Multimedia Processing                           | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Galaxy                     | <a href="https://galaxyproject.org/">https://galaxyproject.org/</a>   | Open Source, Workflow Environment       | Center for Comparative Genomics and Bioinform     | USA           | Business, Scien  | Secondary-Sector, Quaternary-Bio-technology | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| GAZE Embedded              | <a href="http://www.gaze.com/">http://www.gaze.com/</a>   | Open Source, Component Collection, Work | The University of Sheffield                       | UK            | Business, Inter  | Secondary-Sector, Quaternary-Bio-technology | Education, Research and Devel               | Information Extraction, Natural Language Processing, Text    | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Heart of Gold              | <a href="http://www.heartofgold.de/">http://www.heartofgold.de/</a>   | Open Source, Component Collection, Work | DPI GmbH, Language Technology Lab                 | Germany       | Business         | Quaternary-Sector                           | Education, Research and Devel               | Natural Language Processing                                  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| IBM SPSS Predictive Analyt | <a href="http://www.ibm.com/software/products/ibm-spss-predictive-analytics/">http://www.ibm.com/software/products/ibm-spss-predictive-analytics/</a>                                     | Commercial                              | IBM   | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Influenza                  | <a href="http://influenza.com/en/">http://influenza.com/en/</a>   | Webapplication, Workflow                | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| JColte                     | <a href="http://www.jcolte.de/">http://www.jcolte.de/</a>   | Open Source, Component Collection, Work | Friedrich Schiller University, Department of cog  | Germany       | Business, Scien  | Secondary-Sector, Quaternary-Engineering    | Education, Research and Devel               | Information Extraction, Information Retrieval, Natural Lan   | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Jpylyzer                   | <a href="http://www.jpylyzer.com/">http://www.jpylyzer.com/</a>   | Open Source, Pipeline, Standalone       | N/A   | EU            | All              | Quaternary-Sector                           | Research and Development (e.                | Multimedia Processing  | Scope   | N/A               | N/A               | N/A     | N/A      | N/A     |
| KAF                        | <a href="http://kaf-project.eu/kafgroup/it-err-11/kaf/">http://kaf-project.eu/kafgroup/it-err-11/kaf/</a>   | Open Source, Standalone                 | N/A   | EU            | All              | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Information Extraction, Term/Concept Extra      | KYOTO   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Kapler                     | <a href="http://www.kapler.com/">http://www.kapler.com/</a>   | Open Source, Workflow Environment       | LX Davis Computer Science                         | USA           | Business, Scien  | Secondary-Sector, Quaternary-Engineering    | Research and Development (e.                | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| KNIME                      | <a href="http://www.knime.com/knime/">http://www.knime.com/knime/</a>   | Open Source, Workflow                   | KNIME.com AG                                      | Switzerland   | Business, Scien  | Secondary-Sector, Tertiary-Sec              | Finance, Banking, Insurance                 | Data Mining, Text Mining                                     | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Language Grid              | <a href="http://www.languagegrid.org/">http://www.languagegrid.org/</a>   | Infrastructure, Platform                | The Language Grid Association                     | Japan         | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | Association rules, Multimedia Processing, Natural Lan        | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| LAPPS Grid                 | <a href="http://www.lappsgrid.org/">http://www.lappsgrid.org/</a>   | Infrastructure, Platform                | Vassar College, et al                             | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | Natural Language Processing                                  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Massive Online Analysis (M | <a href="http://www.massive-online-analysis.org/">http://www.massive-online-analysis.org/</a>   | Open Source                             | University of Waikato                             | New Zealand   | All              | Quaternary-Sector                           | Research and Development (e.                | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Matchbox                   | <a href="http://matchbox.openpreserve.org/">http://matchbox.openpreserve.org/</a>   | Open Source, Pipeline, Standalone       | N/A   | EU            | All              | Quaternary-Sector                           | Research and Development (e.                | Multimedia Processing  | Scope   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Meta                       | <a href="http://meta.com/">http://meta.com/</a>   | Commercial                              | N/A   | USA           | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Microsoft SQL Server Anal  | <a href="https://technet.microsoft.com/en-us/library/924648c1-3100-407c-b068-746104c17133.aspx">https://technet.microsoft.com/en-us/library/924648c1-3100-407c-b068-746104c17133.aspx</a> | Commercial                              | Microsoft Technet                                 | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Multimedia Processing                           | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Neural Designer            | <a href="http://www.neuraldesigner.com/">http://www.neuraldesigner.com/</a>   | Commercial                              | Artelnic  | Spain         | Business         | Tertiary-Sector, Quaternary-Sector          | Health and Soc                              | Research and Development (e.                                 | Data Mining, Machine Learning, Predictive Analytics       | N/A               | N/A               | N/A     | N/A      | N/A     |
| NITE                       | <a href="http://www.nite.org/">http://www.nite.org/</a>   | Open Source, Component Collection, Work | Team NITE   | Sweden/USA/Au | Business, Inter  | Secondary-Sector, Quaternary-Sector         | Education, Research and Devel               | Artificial Intelligence, Information Retrieval, Machine Lea  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Open Calais                | <a href="http://www.open-calais.com/">http://www.open-calais.com/</a>   | Commercial, Open Source                 | Thomson Reuters                                   | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| OpenDataSoft               | <a href="http://www.opendatasoft.com/">http://www.opendatasoft.com/</a>   | Commercial, Open Source                 | IBM   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Oracle Data Miner GUI      | <a href="http://www.oracle.com/technetwork/data/oracle-dataminer-gui-134846.pdf">http://www.oracle.com/technetwork/data/oracle-dataminer-gui-134846.pdf</a>                               | Commercial                              | Oracle  | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Orange                     | <a href="http://orange.biolab.si/">http://orange.biolab.si/</a>   | Open Source                             | N/A   | N/A           | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Pageflyzer                 | <a href="http://pageflyzer.openpreserve.org/">http://pageflyzer.openpreserve.org/</a>   | Open Source, Standalone                 | N/A   | EU            | Internet         | Quaternary-Sector                           | Research and Development (e.                | Multimedia Processing  | Scope   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Papyrus                    | <a href="http://papyrus.nl.edu/">http://papyrus.nl.edu/</a>   | Open Source, Workflow                   | Information Sciences Institute, The University of | USA           | Scientific       | Secondary-Sector, Quaternary-Aerospace      | Research and Development (e.                | Data Mining, Text Mining                                     | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Pentaho                    | <a href="http://www.pentaho.com/">http://www.pentaho.com/</a>   | Commercial, Open Source                 | Pentaho Corporation                               | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Pipeline Pilot             | <a href="http://www.pipeline-pilot.com/collaborative-software/">http://www.pipeline-pilot.com/collaborative-software/</a>   | Workflow                                | Accelrys  | USA           | Scientific       | Secondary-Sector, Quaternary-Bio-technology | Research and Development (e.                | Data Mining, Graph Mining, Multimedia Processing, Text       | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Pilot TreatmentBank        | <a href="http://pilot.openpreserve.org/">http://pilot.openpreserve.org/</a>   | Webapplication, Workflow                | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| QT11                       | <a href="http://www.qt11.metababa.be/">http://www.qt11.metababa.be/</a>   | Infrastructure, Platform                | N/A   | EU            | Scientific       | Quaternary-Sector                           | Research and Development (e.                | Machine Translation  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Rapidminer Radoop          | <a href="http://www.rapidminer.com/">http://www.rapidminer.com/</a>   | Commercial, Open Source                 | RapidMiner  | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Predictive Analytics, Statistics                | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Rapidminer Server          | <a href="http://www.rapidminer.com/">http://www.rapidminer.com/</a>   | Commercial, Open Source                 | RapidMiner  | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Predictive Analytics, Statistics                | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Rapidminer Studio          | <a href="http://www.rapidminer.com/">http://www.rapidminer.com/</a>   | Commercial, Open Source                 | RapidMiner  | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Predictive Analytics, Statistics                | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| SafeTNet                   | <a href="http://www.safetnet.com/">http://www.safetnet.com/</a>   | Open Source, Component Collection, Work | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| SAP Hana                   | <a href="http://www.sap.com/">http://www.sap.com/</a>   | Commercial                              | SAP America, Inc.                                 | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| SAS Data Mining            | <a href="http://www.sas.com/">http://www.sas.com/</a>   | Commercial                              | SAS Institute                                     | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| SPSS                       | <a href="http://www.ibm.com/ibm/analytic/us/en/ibm/spss/">http://www.ibm.com/ibm/analytic/us/en/ibm/spss/</a>   | Commercial, Workflow, Standalone        | IBM   | USA           | All              | Tertiary-Sector, Quaternary-Sector          | Finance, Banki                              | Education, Research and Devel                                | Data Mining, Information Extraction, Predictive Analytics | N/A               | N/A               | N/A     | N/A      | N/A     |
| Taverna                    | <a href="http://www.taverna.org.uk/">http://www.taverna.org.uk/</a>   | Open Source, Workflow Environment       | School of Computer Science, University of Man     | UK            | Scientific       | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Text Mining                                     | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| TensorFlow                 | <a href="http://www.tensorflow.org/about/">http://www.tensorflow.org/about/</a>   | Commercial                              | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Theano                     | <a href="http://www.theano.org/">http://www.theano.org/</a>   | Open Source, Workflow                   | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Think Analytics            | <a href="http://www.thinkanalytics.com/">http://www.thinkanalytics.com/</a>   | Commercial, Open Source                 | ThinkAnalytics Ltd                                | USA           | Business         | Quaternary-Sector                           | Research and Development (e.                | Data Mining  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Think Enterprise Data Min  | <a href="http://www.thinkanalytics.com/">http://www.thinkanalytics.com/</a>   | Open Source                             | Think Analytics                                   | USA           | Business, Scien  | Quaternary-Sector                           | Research and Development (e.                | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| TopPlace                   | <a href="http://www.topplace.com/products/topplace/">http://www.topplace.com/products/topplace/</a>   | Open Source, Workflow                   | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Trendminer                 | <a href="http://www.trendminer.com/">http://www.trendminer.com/</a>   | Commercial, Open Source                 | N/A   | USA           | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Triana                     | <a href="http://www.triana-odi.org/">http://www.triana-odi.org/</a>   | Open Source, Workflow                   | Cardiff   | UK            | Scientific       | Quaternary-Sector                           | Research and Development (e.                | Data Mining, Multimedia Processing, Text Mining              | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| TTNWWW                     | <a href="http://www.ttnwww.com/">http://www.ttnwww.com/</a>   | Webapplication, Workflow                | N/A   | EU            | All              | Quaternary-Sector                           | Research and Development (e.                | Natural Language Processing, Term/Concept Extraction         | CLARIN  | N/A               | N/A               | N/A     | N/A      | N/A     |
| ViTrials                   | <a href="http://www.vitrials.com/">http://www.vitrials.com/</a>   | Commercial, Open Source                 | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| Weblicht                   | <a href="http://www.weblicht.de/">http://www.weblicht.de/</a>   | Webapplication, Workflow                | N/A   | EU            | All              | Quaternary-Sector                           | Research and Development (e.                | Natural Language Processing, Term/Concept Extraction         | CLARIN  | N/A               | N/A               | N/A     | N/A      | N/A     |
| WEKA                       | <a href="http://www.cs.waikato.ac.nz/ml/weka/">http://www.cs.waikato.ac.nz/ml/weka/</a>   | Open Source, Component Collection, Work | University of Waikato                             | New Zealand   | All              | Quaternary-Sector                           | Research and Development (e.                | Classification, Clustering, Data Mining, Predictive Analyt   | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| WordRank                   | <a href="http://www.wordrank.com/">http://www.wordrank.com/</a>   | Open Source, Workflow                   | N/A   | EU            | N/A              | N/A   | N/A   | N/A  | N/A   | N/A               | N/A               | N/A     | N/A      | N/A     |
| xcorSound                  | <a href="http://www.xcor.com/">http://www.xcor.com/</a>   | Open Source, Standalone                 | N/A   | EU            | All              | Quaternary-Sector                           | Research and Development (e.                | Multimedia Processing  | Scope   | N/A               | N/A               | N/A     | N/A      | N/A     |