

OpenTED Browser: Insights into European Public Spendings

Yann-Aël Le Borgne¹, Adriana Homolova², and Gianluca Bontempi¹

¹ Machine Learning Group, Université Libre de Bruxelles
Boulevard du Triomphe, CP212
1050 Brussels, Belgium

² Independent Data Journalist, Amsterdam, Netherlands

Abstract. We present the *OpenTED browser*, a Web application allowing to interactively browse public spending data related to public procurements in the European Union. The application relies on Open Data recently published by the European Commission and the Publications Office of the European Union, from which we imported a curated dataset of 4.2 million contract award notices spanning the period 2006-2015. The application is designed to easily filter notices and visualise relationships between public contracting authorities and private contractors. The simple design allows for example to quickly find information about who the biggest suppliers of local governments are, and the nature of the contracted goods and services. We believe the tool, which we make Open Source, is a valuable source of information for journalists, NGOs, analysts and citizens for getting information on public procurement data, from large scale trends to local municipal developments.

Keywords: Public Procurements, European Union, Open Data, Government Transparency, Data Journalism

1 Introduction

Public procurement is the process whereby governments buy goods and services, such as office supplies, equipment, buildings, roads, and so forth. It represents about one third of total government expenditures in OECD countries [3]. This is usually a public-private deal in which the buyer is a public entity and the (winning) bidder is usually a privately owned company. Public procurements are mostly financed by public funds [7, 11].

The development of Tenders Electronic Daily (TED) [9] can be considered as the biggest EU-wide effort made so far to support procurement across borders. TED, which is managed by the Publications Office of the European Union [8], is the online version of the S Series of the Official Journal of the European Union (OJEU), which is a supplement to the Journal particularly focused on European public procurement. TED publishes over 1000 new over EU-threshold value procurement notices every day worth about EUR 400 billions a year [3–5]. Furthermore, TED also publishes other documents coming from funds to be spent on EU external aid and the procurement of EU institutions.

Along with TED, the Publications Office provides bulk downloads of its data. The raw dataset contains all contract notices for tender data since 1995 in XML format, and its size is around 100 GB (GigaBytes). The browsing and analysis of this dataset raises a number of complex challenges, due to varying quality of data between countries and years, missing values, variability in the naming of contracting authorities and entities, multilingual documents, and so forth.

The Directorate-General for Internal Market, Industry, Entrepreneurship and SMEs (DG-GROW) of the European Commission has recently initiated great efforts in the curation of the Publications Office data, and released in August 2015 a summary of all contract award notices (CANs) for European public procurements for the period spanning 2006-2015 [6]. The curated dataset provides an unprecedented source of information on public spending in the European Union and yields very valuable information on procurement data. In particular, it provides information related to the identity of contracting authorities and contractors, the nature of the supplies/works/services, the final contract values and the number of offers for a given contract notice. Such information not only makes it possible to provide insights in the network of public/private partnerships, but also to exhibit procurements patterns across all European countries, or to detect (and avoid) corruption [1, 4, 16, 18].

The relative large size of the dataset however still prevents its analysis from users without an analytics background: Data is provided as CSV files for each year of CAN (from 2006 to 2015, ten files in total), whose size ranges from around 100 MB (MegaBytes) to 300 MB, totalising about 2 GB (GigaBytes). While the size of the dataset can be stored without trouble on current laptop or desktop configurations, the opening of such files remains a challenge for standard spreadsheet applications such as Excel, and makes analyses spanning multiple years (i.e. filtering data from multiple files) very tedious. It must be stressed that most of the people interested in CAN data (journalists, entrepreneurs, citizen, ...) do not have the necessary analytics background to explore datasets of such size.

The OpenTED browser aims at bridging this gap, and our contributions are the following. First, data is stored on our OpenTED server, and does not need to be fully downloaded, making it easier for users with slow Internet connection to access the data. Users can furthermore filter the data they need (according to countries, years, type of goods, and so forth), and download only what they are interested in. Second, we provide a visualisation tool, based on Sankey diagrams, that represents as a graph the contract awards between public authorities and private contractors. The visualisation makes clear how much money the contracts are worth, and provides hyperlinks to the official TED award notices for further details. Finally, we make the code for both the data preprocessing (Python) and the Web application (R Shiny) open source. A docker container can be used to run the Web application on a local machine, making the interaction with the application even faster.

The paper is organised as follows. Section 2 details the official TED CAN data, and how they were preprocessed for the OpenTED browser. Section 3

presents the OpenTED browser Web interface. Section 4 presents the lottery game ‘Public spending is fun! who wants to be a supplier?’, a ‘gamification’ of the search for contract award notices. Section 5 concludes the paper with open issues and perspectives.

2 Data and Methods

Raw data was retrieved from the Open Data portal of the European Union [6], slightly preprocessed for easier querying, and converted to Parquet format for efficiency reasons. We describe the data and preprocessing hereafter.

2.1 Data overview

The data comes from the European Economic Area, Switzerland, and the former Yugoslav Republic of Macedonia and covers the time period between 2006/01/01 and 2015/12/31. It is in comma separated value (CSV) format and is encoded as UTF-8. Generally, the data consists of notices above the procurement thresholds. However, publishing below threshold notes in TED is considered good practice, and thus a non-negligible number of below threshold notices is present as well [6].

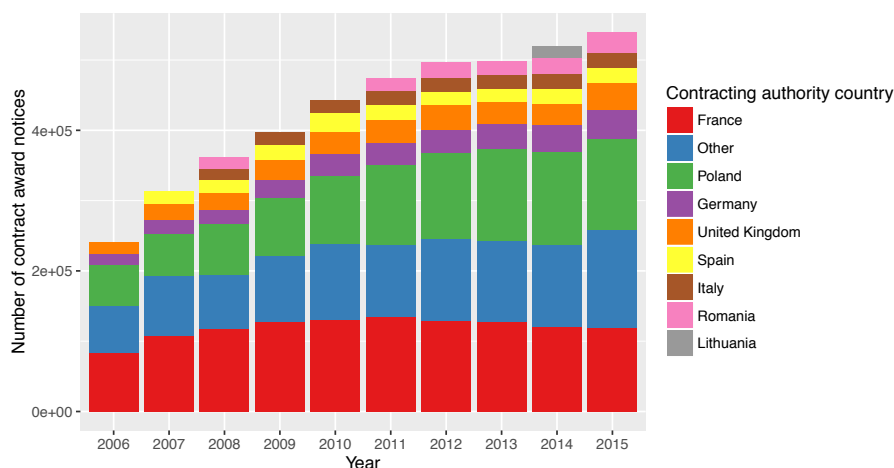


Fig. 1: Number of contract award notices (CANs) per country and per year, since 2006. Only countries with more than 15000 CANs per year are reported (number of CANs for other countries are summarized in ‘Other’).

The number of countries covered has increased throughout the years, generally in line with their accession to the single European market. This is illustrated in Fig. 1, which reports the number of CANs per country and per year, since

2006. Only countries with more than 15,000 CANs per year are reported for clarity reasons (number of CANs for other countries are summarized in ‘Other’).

The number of CANs increased from about 250,000 in 2006 to more than 500,000 in 2015, and a total of 4,283,986 CANs are available in the dataset. It is worth noting that the number of CANs available per country is quite imbalanced, and are the highest for France, Poland, Germany and the United Kingdom (UK).

The data includes 48 selected fields from CANs (Table 3 in the Annex), divided in six categories: *Notice metadata*, *Contracting authority or entity identification*, *Winning bidder identification*, *Various CAN level variables*, and *Various CA level variables* [6]. The size of the CSV file with all notices is slightly above 2GB.

2.2 Data preparation and conversion to Parquet

We remained as faithful as possible to the original data, and included all 4,283,986 records and 48 fields in the OpenTED browser. For making the filtering of the dataset more user-friendly, we however renamed nine fields, which are highlighted in the OpenTED browser (see Section 3). These fields consist in the date of the CAN, the identifier (ID), the name and country of the contracting authority, the name and country of the contractor, the value of the CAN, the number of offers and the CPV code. The renaming is detailed in Table 1.

Table 1: Renaming of highlighted fields in the OpenTED browser.

Original name	Description	New name
DT_DISPATCH	The date when the buyer dispatched (sent) the notice	Dispatch_Date
ID_NOTICE_CAN	Unique identifier of the contract award notice for publication to TED	Award_Notice_Id_Link
CAE_NAME	Official name	Contracting_Authority_Name
ISO_COUNTRY_CODE	Country	Contracting_Authority_Country
WIN_NAME	Official name	Contractor_Name
WIN_COUNTRY_CODE	Country	Contractor_Country
VALUE_EURO	CAN value, in EUR, without VAT. If the value was not present, the lowest bid is included	Contract_Value_Euros
NUMBER_OFFERS	Number of offers received	Number_Offers_Received
CPV	The main Common Procurement Vocabulary code of the main object of the contract	CPV_Code

Besides the renaming, we also reformatted fields involving dates, countries, and values. In the original data, date formats are *Day-Months (3 first letters)-Year (2 last digits)*, e.g., ‘31-DEC-13’. We reformatted it as *Year (4 digits)-Month (2 digits)-Day*, e.g., ‘2013-12-31’. Such conversion makes it more suitable to define filtering intervals on dates, using operator such as greater or less than (e.g. date higher than ‘2013-06-01’ and less than ‘2013-12-31’ to get CANs from the last six months of 2013). Country related fields were converted from the ISO code (2 letters, e.g., ‘FR’) to their full names (e.g., ‘France’). All award

value fields were converted from floats to integers. Finally, award notice IDs were converted to hyperlinks linking to the CAN page on the official TED Web site [9].

Finally, we converted the CSV data to the Apache Parquet format [2, 15]. Apache Parquet is a columnar data storage format designed to support very efficient compression and encoding schemes. Besides compression, Parquet also allows data to be queried from the files using SQL like syntax, a feature which we use in the browser. Data types were associated to CAN fields in order to perform filtering and SQL queries, which we detail in Table 3. All fields related to numbers were associated an *Integer* data type. Fields involving strings were associated a *String* data type, or *factor* when the number of values was less than 300 hundreds. The use of the factor type allows to present users with a list of choices in the TED browser filtering tool. After Parquet conversion, the dataset size was reduced to 315MB.

We make available the code for data preparation and Parquet conversion as an IPython notebook [13].

3 OpenTED browser

The OpenTED browser is a Web application that provides a user-friendly access to the CANs. Its main features are a filtering tool for extracting subsets of CANs, and a visualisation tool that displays the relationships between contracting authorities and contractors by means of a Sankey diagram. The OpenTED browser is made available online at [14].

3.1 Filtering tool

The filtering tool allows to filter CANs by setting conditions on the content of any of the CAN fields listed in Tables 1 and 3. A snapshot of the tool is given in Fig. 2. All 48 fields can be filtered. The filtering operators provided for a field depends on the data type of the field. A summary of available filtering operators for a given data type is provided in Table 2.

Table 2: Available operators for the different data types.

Data type	Available operators
String	equal,not_equal,less, less_or_equal, greater,greater_or_equal, between, in, not_in,begins_with, ends_with, is_null, is_not_null
Factor	equal,not_equal,is_null, is_not_null
Integer	equal,not_equal,less, less_or_equal, greater,greater_or_equal, between,in, not_in,is_null, is_not_null

Conditions can be combined by logical conjunction and disjunction operators, and may also be nested thanks to the grouping option. An example of filtering is given in Fig. 2. The filter retrieves CANs for which: (i) the contracting authority country is ‘Belgium’, (ii) the CPV code either begins with ‘301’ (Office

Get to know Tender better, play the lottery!

TED Award Notices 2006-2015 Sankey diagram CPV codes What is this interface?

AND OR + Add rule + Add group

Contracting_Authority_Country equal Belgium Delete

AND OR + Add rule + Add group Delete

CPV_Code begins with 301 Delete

CPV_Code begins with 302 Delete

Contract_Value_Euros greater or equal 1000000 Delete

Apply filters

Download selection (CSV)

Select additional variables to display

Show 10 entries

Showing 1 to 10 of 128 entries

Contracting Authority Country	Contracting Authority Name	Dispatch Date	CPV Code Meaning	Contractor Country	Contractor Name	Contract Value Euros	Number Offers Received	CPV Code	Award Notice Id Link
Belgium	Centre de recherche en aéronautique ASBL (Cenaero)	2015-11-30	Super computer	France	Serviware SAS	1496201	2	30211100	427971-2015
Belgium	Resa SA	2015-09-23	Agency fuel cards	Belgium	Total Belgium SA	7062845	3	30163100	340064-2015
Belgium	Universiteit Gent	2015-04-07	Computer equipment and supplies	Belgium	Dimension Data	1250000	1	30200000	123188-2015

Fig. 2: OpenTED browser: Filtering tool. The user can filter CANs by setting conditions on CAN field values. A table displays the corresponding subset of CANs.

machinery, equipment and supplies except computers, printers and furniture) or ‘302’ (Computer equipment and supplies) and (iii) the contract value in euros is more than one million.

Applying the filter returns the set of CANs matching the conditions. In the example above, 128 CANs are returned, and displayed as a table below the filtering widget. The table is interactive, and the user can reorder columns by increasing or decreasing order of values. Award notice IDs are hyperlinks that connect to the page of the CAN on the official TED Web site [9]. Finally, the set of filtered CANs can be downloaded as a CSV using the ‘Download selection’ button.

3.2 Sankey diagram

Sankey diagrams visualise the magnitude of flow between nodes in a network. They provide a useful visualisation tool for contract award notices, where con-

tracting authorities and contractors can be seen as nodes of a network, and the contract values as ‘flows’.

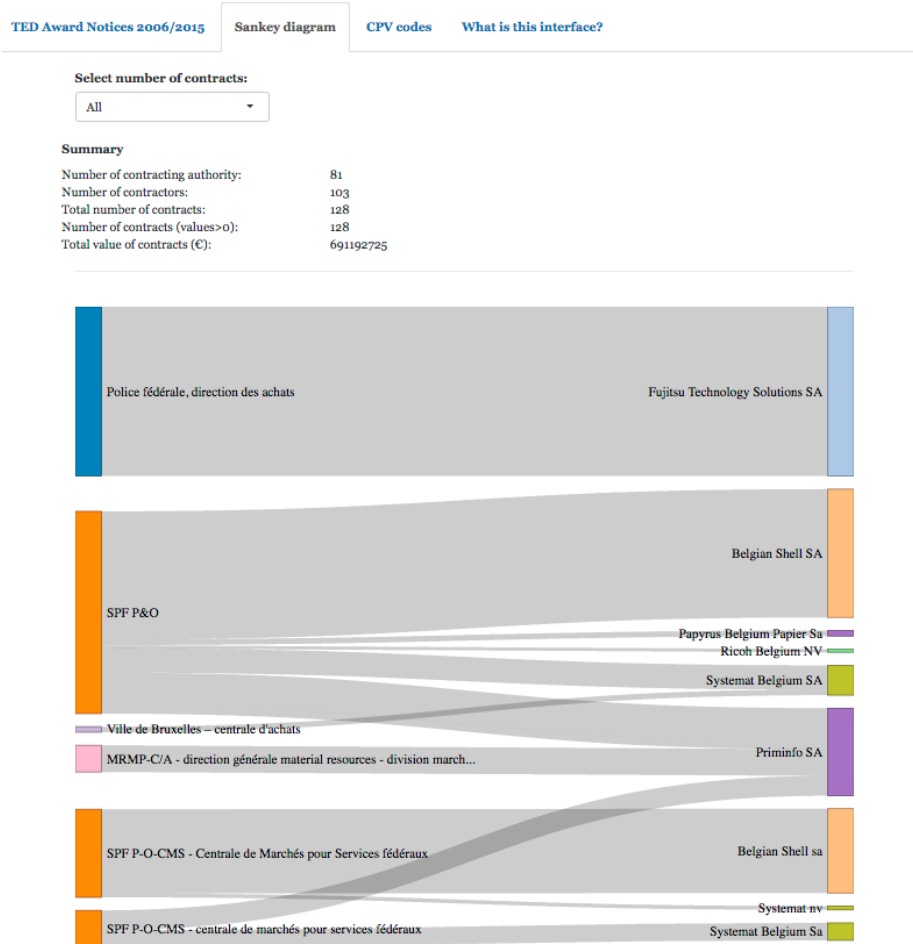


Fig. 3: OpenTED browser: Sankey diagram for the visualisation of CANs. Contracting authorities and contractors are represented on the left and right sides of the network, respectively. The thickness of flows is proportional to the sum of award values between two parties.

The set of contracting authorities are represented on the left side of the network, and the set of contractors on the right side. The thickness of the flow is proportional to the sum of contract values between contracting authorities and contractors.

Fig. 3 gives a snapshot of the Sankey diagram obtained for the subset of CANs matching the conditions given in Section 3.1, and illustrates that the diagram gives a clear overview of the relationships existing between the different parties. A few statistics at the top of the page summarise the number of contracting authorities, contractors, contracts and the sum of contract values in the Sankey diagram.

Tooltips and hyperlinks are tied to the diagram edges, giving the amount of the contract value, and linking to the page of the contract award notice on the official TED Web site, respectively.

3.3 Common Procurement Vocabulary

The subjects of procurement contracts are encoded using Common Procurement Vocabulary (CPV) [10]. CPV is based on a tree structure comprising codes of up to 9 digits (an 8 digit code plus a check digit) associated with a wording that describes the type of supplies, works or services forming the subject of the contract. The first two digits identify the general division, and subsequent digits refine the division into more specific categories. For example, all codes starting with *30* are from the general division of *Office and computing machinery, equipment and supplies except furniture and software packages*, while the code *3012453* refer more specifically to *Scanner transparency adapters*.

TED Award Notices 2006/2015 Sankey diagram **CPV codes** What is this interface?

Meanings of Common Procurement Vocabulary (CPV) codes

Original code	CPV plain code	NUM-digits	Real Code	Category content
<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text" value="3"/>	<input type="text" value="30"/>	<input type="text" value="computer"/>
30100000-0	30100000	3	301	Office machinery, equipment and supplies except computers, printers and furniture
30200000-1	30200000	3	302	Computer equipment and supplies

Show 10 entries Previous 1 Next

Showing 1 to 2 of 2 entries (filtered from 9,454 total entries)

Fig. 4: OpenTED browser: Meanings of Common Procurement Vocabulary (CPV) codes.

The meaning of CPV codes can be found from the EU Publications office [10]. In order to facilitate the use of CPV codes in the browser, we added a table with all current CPV codes (9454 in total), that can be searched and filtered. The table is available in the *CPV codes* tab, a snapshot of which is given in Fig. 4. The user may fix the number of digits in the results in order to restrict the search to more general categories, and may sort the columns by increasing or decreasing order of values. It should be noted that we only include the current nomenclature

(valid since 2008), and that the integration of the previous nomenclature (before 2008) is part of future work.

3.4 Implementation

The Web application was developed using R Shiny [17], and is made Open Source at [13]. It is worth mentioning that thanks to the expressiveness of R Shiny, the application is rather compact, i.e. less than 500 lines of code in total, making it easily reusable and adaptable. We furthermore provide a Docker container for facilitating its deployment on a different server, or run it on a local machine for faster interaction with the browser [14].

4 Lottery game: Who wants to be a supplier?

The lottery game is a third-party Web site aimed at promoting the OpenTED browser [12]. The Web site gives the player a quest on tender data for finding the biggest suppliers of some goods or services in a given country. Fig. 5 gives a snapshot of the page, where the quest is to find suppliers for ‘Insurance and pension in Sweden in 2013’.

HOME FAQ RULES FEEDBACK & SHARE

Welcome! Would you like to find out who are the biggest government suppliers?

Dive with us into **tender data!**

Choose a country:

Any

Find the biggest suppliers for:

Insurance and pension services in Sweden in 2013

Let's play! That's not fun enough Show me the solution

Tweet your results!

Fig. 5: Lottery game: Public spending is fun! Who wants to be a supplier?

If the player accepts (“Let’s play” button), she is redirected to the OpenTED browser where the goal is to find the subset of CAN corresponding to the quest.

She may also get the solution by selecting ‘Show me the solution’, which will also redirect to the OpenTED browser, but the filtering tool will be prefilled with the set of conditions answering the quest. If the user does not like the quest, she can ask for another one by selecting the ‘Not fun enough’ button.

5 Conclusion and perspectives

The OpenTED browser provides an intuitive online gateway for filtering contract award notices (CANs) of the European procurement system, and for getting insights into the business relationships existing between contracting authorities and entities. We believe that the tool, thanks to its simplicity and ease of use, can be of significant interest for a number of users.

Ongoing improvements concern the correction of inconsistencies in the data, related to CPV codes and to the naming of contracting authorities and entities. The nomenclature for CPV codes changed in 2008, which requires to adapt queries to two nomenclatures when searching for CANs before and after 2008. We plan to address this issue shortly. The second type of inconsistency relates to variations in the naming of bidders and buyers. These may be caused by typos, but are also due to different ways of naming an entity (e.g., Siemens, Siemens A.G., Siemmens A.G., and so forth). While efforts to provide unique identifiers instead of names are being promoted, name inconsistencies currently remain an important challenge to properly group CANs according to their contracting authorities and entities.

In a larger perspective, a wide range of avenues exist for improving the browsing of TED notices, and for providing better insights into tender data. To name a few, possible extensions concern the integration of contract notices, also recently made available by the Publications Office as curated Open Data, and make use of advanced analysis techniques such as clustering, graph analysis, or outlier detection, to investigate questions such as what makes a bidder successful, what are the tendering patterns among EU countries, or to identify indicators of fraudulent behaviours.

Acknowledgments. The authors acknowledge the support of “BruFence: Scalable machine learning for automating defense system” (RBC/14 PFS-ICT 5), a project funded by the Institute for the Encouragement of Scientific Research and Innovation of Brussels (INNOVIRIS, Brussels Region, Belgium), Journalismfund.eu for organising the DataHarvest/European Investigative Journalism Conference, the OpenTED working group (<http://ted.openspending.org>), the European Union Publications Office, <http://ted.europa.eu>, 1998–2016, and the European Commission, Directorate-General for Internal Market, Industry, Entrepreneurship and SMEs (DG-GROW) for providing the data, and Jáchym Hercher, Policy Officer at DG-GROW, for providing feedback and improvements on this article.

References

1. Alvarez, J.M., Labra, J.E., Cifuentes, F., Alor-Hernández, G., Sánchez, C., Luna, J.A.G.: Towards a pan-european e-procurement platform to aggregate, publish and search public procurement notices powered by linked open data: the moldeas approach. *International Journal of Software Engineering and Knowledge Engineering* 22(03), 365–383 (2012)
2. Apache Software Foundation: Parquet, <https://parquet.apache.org/>, (Viewed June 2016)
3. Cernat, L., Kutlina-Dimitrova, Z.: International public procurement: From scant facts to hard data. Robert Schuman Centre for Advanced Studies Research Paper No. RSCAS (2015)
4. De La Iglesia, J.L.M.: Alternative estimation of public procurement advertised in the official journal as of GDP official indicator using open government data. *Computers in Industry* 65(5), 905–912 (2014)
5. DG GROW G4 - Innovative and e-Procurement: 2014 public procurement indicators. (February 2016), <http://ec.europa.eu/DocsRoom/documents/15421/>
6. DG Internal Market, Industry, Entrepreneurship, and SMEs, European Commission, Brussels: Ted csv dataset (2006-2015), tenders electronic daily, supplement to the official journal of the european union, <https://open-data.europa.eu/cs/data/dataset/ted-csv>, version 2.0. Accessed on 2016-06-25.
7. DIGIWHIST - Deliverable 1.1.: Towards a comprehensive mapping of information on public procurement tendering and its actors across europe. (August 2015), http://digiwhist.eu/wp-content/uploads/2016/01/DIGIWHIST_D1_1-AccessToTenderInfo.pdf
8. European Commission: The publications office of the european union (2016), <http://publications.europa.eu>, (Viewed June 2016)
9. European Commission: TED. tenders electronic daily (2016), <http://ted.europa.eu>, (Viewed June 2016)
10. European Union Publications Office: Common Procurement Vocabulary (2008), uRL: <https://simap.ted.europa.eu/cpv>. Viewed June 2016.
11. Hoekman, B.: International cooperation on public procurement regulation. Robert Schuman Centre for Advanced Studies Research Paper No. RSCAS 88 (2015)
12. Homolova, A.: Lottery game: Who wants to be a supplier. (2015), uRL: <http://supplier.tenders.exposed/>. Viewed June 2016.
13. Le Borgne, Y.A.: IPython notebook for data preparation and Parquet conversion (2016), <https://github.com/Yannael/OpenTED>, viewed June 2016.
14. Le Borgne, Y.A.: Opented browser web site (2016), <http://yleborgne.net/opented>, viewed June 2016.
15. Melnik, S., Gubarev, A., Long, J.J., Romer, G., Shivakumar, S., Tolton, M., Vasilakis, T.: Dremel: Interactive analysis of web-scale datasets. In: Proc. of the 36th Int'l Conf on Very Large Data Bases. pp. 330–339 (2010)
16. Miroslav, M., Miloš, M., Velimir, Š., Božo, D., Jorje, L.: Semantic technologies on the mission: Preventing corruption in public procurement. *Computers in Industry* 65(5), 878–890 (2014)
17. RStudio, Inc: Easy web applications in R. (2013), uRL: <http://www.rstudio.com/shiny/>. Viewed June 2016.
18. Uyarra, E.: Review of measures in support of public procurement of innovation. Report within the MIOIR-NESTA Compendium of Evidence on Innovation Policy. London/Manchester (2013)

Appendix

Table 3: Fields present in contract award notices CSV files

Field	Description	Type
Notice metadata		
ID_NOTICE_CAN	Unique identifier of the contract award notice	String
YEAR	Year of publication of the notice	Integer
ID_TYPE	Standard form number	Factor
DT_DISPATCH	The date when the buyer dispatched (sent) the notice for publication to TED	String
XSD_VERSION	Version of the XML schema definition [ADDED]	Factor
CANCELLED	1 = this notice was later cancelled [ADDED]	Factor
Contracting authority or entity identification		
CAE_NAME	Official name	String
CAE_NATIONALID	“National ID” e.g. VAT number for utilities	String
CAE_ADDRESS	Postal address	String
CAE_TOWN	Town	String
CAE_POSTAL_CODE	Postal code	String
ISO_COUNTRY_CODE	Country	Factor
Winning bidder identification		
WIN_NAME	Official name	String
WIN_ADDRESS	Postal address	String
WIN_TOWN	Town	String
WIN_POSTAL_CODE	Postal code	String
WIN_COUNTRY_CODE	Country	Factor
Various CAN level variables		
CAE_TYPE	Type of contracting authority	Factor
MAIN_ACTIVITY	The classification corresponds to COFOG divisions	String
B_ON_BEHALF	This indicates either a central purchasing body or several buyers buying together	Factor
TYPE_OF_CONTRACT	Type of contract	Factor
TALLOCATION_NUTS	The Nomenclature of Territorial Units for Statistics (NUTS) code placement	String
B_FRA_AGREEMENT	The notice involves the establishment of a framework agreement	Factor
B_DYN_PURCH_SYST	The notice involves contract(s) based on a dynamic purchasing system	Factor
CPV	The main Common Procurement Vocabulary code of the main object of the contract	String
ADDITIONAL_CPV1-4	The first four CPV listed in the notice.	String
B_GPA	The contract is covered by the Government Procurement Agreement	Factor
VALUE_EURO	CAN value, in EUR, without VAT. If the value was not present, the lowest bid is included	Integer
VALUE_EURO_FIN_1	CAN value, in EUR, without VAT. If the value was not present, second estimate	Integer
VALUE_EURO_FIN_2	CAN value, in EUR, without VAT. If the value was not present, third estimate	Integer
TOP_TYPE	Type of procedure	Factor
CRIT_CODE	Award criteria	Factor
CRIT_CRITERIA	Information on award criteria.	String
CRIT_WEIGHTS	Information on award criteria weighing	String
B_ELECTRONIC_AUCTION	An electronic auction has been used	Factor
NUMBER_AWARDS	The number of CAs for a given CAN. [ADDED]	Integer
Various CA level variables		
ID_AWARD	Unique contract award identifier	String
CONTRACT_NUMBER	Contract No	String
LOT_NUMBER	An identifier of a lot	String
TITLE	Title	String
NUMBER_OFFERS	Number of offers received	Integer
NUMBER_OFFERS_ELECTR	Number of offers received by electronic means	Integer
AWARD_EST_VALUE_EURO	Estimated CA value, in EUR, without VAT	Integer
AWARD_VALUE_EURO	Total final CA value, in EUR, without VAT. If the value was not present, the lowest bid is included	Integer
AWARD_VALUE_EURO_FIN_1	CA value, in EUR, without VAT. If a value field is missing, second estimate	Integer
B_SUBCONTRACTED	The contract is likely to be subcontracted	Factor
B_EU_FUNDS	The contract is related to a project and / or programme financed by European Union funds	Factor
DT_AWARD	Date of contract award	String