

Chapter 74

News Trends Processing Using Open Linked Data

Antonio Garrote

Universidad de Salamanca, Spain

María N. Moreno García

Universidad de Salamanca, Spain

ABSTRACT

In this chapter we describe a news trends detection system built with the aim of detecting daily trends in a big collection of news articles extracted from the web and expose the computed trends data as open linked data that can be consumed by other components of the IT infrastructure. Due to the sheer amount of data being processed, the system relies on big data technologies to process raw news data and compute the trends that will be later exposed as open linked data. Thanks to the open linked data interface, data can be easily consumed by other components of the application, like a JavaScript front-end, or re-used by different IT systems. The case is a good example of how open linked data can be used to provide a convenient interface to big data systems.

ORGANIZATION BACKGROUND

The organization involved in the project is an Internet company providing data analysis services for different kinds of web data. Integration of different data sources capable of generating more useful insights on clients' data is an essential task in the company strategy.

From the organizational point of view, the software development process is accomplished by small and highly autonomous engineering teams responsible for different projects and relying on services provided by other teams. Use of big data

technologies means that many times resources like computer clusters must be shared by different teams. This practice demands a high degree of cooperation between teams.

Within this context of small teams building a network of services that are used and combined by other teams, the use of open linked data makes it possible for the easy inter-operability between data resources as well as provides a shared vocabulary for the outcome data, processed by big data systems like Apache Hadoop.

Data management in big data systems is a hard problem. The organization undergoing the devel-

opment project described in this case generates and processes tera bytes of data on a daily basis. Most of these data have been so far stored as plain tab-separated files in the Hadoop Distributed File System (HDFS) (Shvachko et al., 2010).

Re-using and cataloging the available data sources have been traditionally an important issue, due to the distributed nature of the development teams in the organization. As a consequence, problems like finding if the right information is already available in some part of the cluster file system or if some particular data generation process is still in use have been hard to solve, usually involving a lot of communication overhead between members of different teams.

This situation was slightly improved when a more structured data storage technology like Apache Hive started to be used instead of direct access to plain HDFS files. Hive provides a data abstraction layer in the form of data tables with a certain data schema and a relation SQL-like data retrieval language that can be used on top of the map-reduce platform offered by Hadoop. The use of a schema and an easy interface to query the stored data made it easier for non technical users to retrieve information from the cluster as well as provided a better definition of the available data. However, the problem of finding available data in the cluster remained a problem.

When making available structured information about news trends started to be considered as development project, linked data appeared as a possible alternative to provide a more open interface to the available data stored in the cluster, as well as a mechanism to interlink isolated data sets using well known web technologies like URIs and hyper-links.

CASE DESCRIPTION

The main goal of the project was to make available daily news trends as a structured data source that could be used as an additional input in any data

analysis task being performed in the organization. Computation of the news trends was to be achieved in a series of steps involving:

- Crawling of news raw data from web sources.
- Classification of the news data by country, language, and topic.
- Extraction of trends using natural language processing techniques.
- Storage of the processed trends in the data cluster in a structured format compatible with Apache Hive.
- Building a data interface for the data available as a collection of web services that could be re-used by other applications without accessing directly the data stored in HDFS.
- Providing a web application exposing the news trend data through a user interface that could be used by non technical users.

The team assigned to the project consisted of two developers with a good knowledge of the statistical techniques for natural language processing as well as experience with the underlying Hadoop platform and web development skills.

The decision of using open linked data affected the later goals of the projects, like providing a web interface for the computed trends. A combination of linked data standards for data modeling like the Resource Description Language (RDF) embedded inside Java Script documents using the JSON-LD standard that could be made available through a series of simple RESTful (Fielding, 2000) web services was the intended solution for the web data interface.

Technology Concerns

From a technical point of view, the project showed important challenges. Crawling of web data is an error prone task, especially when no stable web

3 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/news-trends-processing-using-open-linked-data/125361

Related Content

An Electronic Contract Signing Protocol Using Fingerprint Biometrics

Harkeerat Bedi, Li Yang and Joseph M. Kizza (2013). *IT Policy and Ethics: Concepts, Methodologies, Tools, and Applications* (pp. 635-661).

www.irma-international.org/chapter/electronic-contract-signing-protocol-using/75050

Energy-Efficient MAC Protocols in Distributed Sensor Networks

Yupeng Huang and Rui Li (2013). *IT Policy and Ethics: Concepts, Methodologies, Tools, and Applications* (pp. 1776-1797).

www.irma-international.org/chapter/energy-efficient-mac-protocols-distributed/75099

Should Buyers Try to Shape IT Markets Through Non-Market (Collective) Action? Antecedents of a Transaction Cost Theory of Network Effects

Kai Reimers and Mingzhi Li (2005). *International Journal of IT Standards and Standardization Research* (pp. 44-67).

www.irma-international.org/article/should-buyers-try-shape-markets/2563

A Diffusion Model for Communication Standards in Supply Networks

Tim Stockheim, Michael Schwind and Kilian Weiss (2006). *International Journal of IT Standards and Standardization Research* (pp. 24-42).

www.irma-international.org/article/diffusion-model-communication-standards-supply/2576

Ensuring Users' Rights to Privacy, Confidence and Reputation in the Online Learning Environment: What Should Instructors Do to Protect Their Students' Privacy?

Louis B. Swartz, Michele T. Cole and David Lovejoy (2010). *Information Communication Technology Law, Protection and Access Rights: Global Approaches and Issues* (pp. 346-362).

www.irma-international.org/chapter/ensuring-users-rights-privacy-confidence/43504