

Linked Open Data as the fuel for Smarter Cities

Mikel Emaldi, Jon Lázaro, Oscar Peña, Diego López-de-Ipiña, Sacha Vanhecke

No Institute Given

Abstract. In the last decade big efforts have been carried out in order to move towards the Smart City concept, both from the academic and industrial points of view, encouraging researchers and data stakeholders to find new solutions on how to cope with the huge amount of generated data.

Open Data has arisen as a way to share data in order to be consumed freely without restrictions from copyright, patents or other mechanisms of control. Nowadays Open Data is an achievable concept thanks to the World Wide Web, and has been re-defined for its application in different domains.

Regarding public administrations, the concept of Open Government has found an ally in Open Data concepts, defending citizens' right to access data, documentation and proceedings of the governments.

We propose the use of Linked Open Data, a set of best practices to publish data on the Web proposed by the W3C, in a new data life-cycle management model, allowing governments and individual to handle better their data, easing the consumption by anybody, including both companies and third parties interested in the exploitation of the data, and citizens as end users receiving relevant curated information and reports about their city.

In summary, Linked Open Data uses the previous Openness concepts to evolve from an infrastructure thought for humans, to an architecture for the automatic consumption of big amounts of data, providing relevant and high quality data to end users with low maintenance costs. Smart data can now be achievable in smart cities.

1 Introduction

In the last decade, cities have been sensorized, protocols are constantly refined to deal with the possibilities that new hardware offers, communication networks are offered in all flavours, etc., generating lots of data everyday that needs to be dealt with. Public administrations usually are not able to process all the data they have in an efficient way, resulting in large amounts of data going un-analyzed and affecting the end-users that could benefit from them.

Citizens are also being encouraged to adopt the role of linked open data providers. User-friendly Linked Data apps should allow citizens to easily contribute with new trustable data that can be linked to already existing published (more static generally) Linked Open Data provided by city councils. In addition, people-centric mobile sensing, empowered by the technology inside actual smartphones, should progress into continuous people-centric enriched Linked Data. Linked Open Data also encourages the linkage to other resources described formally through structured vocabularies, allowing the discovery of related information and the possibility to make inferences, resulting in higher quality data.

Data management is becoming one of the greatest challenges of the 21st century. Regarding urban growth, experts predict that global urban population will double by the year 2050, meaning that nearly 70% of the whole planet's inhabitants will be living in a major town or city.

This prediction arises the need to deal with the huge amounts of data generated by cities, enabling the possibility to manage their resources in an efficient way. The *Smart Data* term has been coined to address the data that makes itself understandable, by extracting relevant information and insights from big data and presenting the conclusions as human-friendly visualizations.

The problems of managing data are moving to a new level. It's not only a matter of caring about *more* data, but how we can use it efficiently in our processes. It's about how we can deal with increasing volumes of data (from standalone databases to real *Big Data*) and integrate them to our advantage, making it useful and digestible in order to make better decisions.

In the last few years, the *Smart city* concept has been adopted to refer to those cities aware of their citizens' life quality, worried about the efficiency and trustworthiness of the services provided by governing entities and businesses.

Smart data can help cities reach a *Smart City* status, analysing the generated data streams and providing useful information to their users: citizens, council managers, third parties, etc.

Although, efficient data lifecycle management processes need to be adopted as best practices, avoiding to convert input data in non-sense noise that can not be used to improve council's services.

Thus, our approach relies is based on an actual review of the state of the art regarding data lifecycle management, proposing our own model as a more refined approach to the existing ones. We also encourage the adoption of Linked Open Data principles to publish both the whole generated data and the processed data, in order to allow further research on

the area by third parties and the development of new business models relying on public access data. Similar proposals regarding Linked Data are defended by [?].

2 Background and definitions

3 Data Life Cycle

Throughout the literature, a variety of different definitions of data life cycle models can be found. Although they have been developed for different actuation domains, we describe here some of them which we think that can be applied for generic data, independently of its original domain.

3.1 Data Documentation Initiative

The first model to be analysed is the model proposed by Data Documentation Initiative (DDI). The DDI introduced a Combined Life Cycle Model for data managing [1]. As Figure 1 shows, this model has eight elements or steps which can be summarized as follows, according to [2]:

- **Study concept.** At this stage, apart from choosing the research question and the methodology for collecting the data, the processing and analysis step of the needed data to answer the question is planned.
- **Data collection.** This model proposes different methods to collect data, like surveys, health records, statistics or Web-based collections.
- **Data processing.** At this stage, the collected data is processed to answer the proposed research question. The data may be recorded in both machine-readable and human-readable form.
- **Data archiving.** Both data and metadata should be archived to ensure long-term access to them, guaranteeing confidentiality.
- **Data distribution.** This stage involves the different ways in which data is distributed, as well as questions related to the terms of use of the used data or citation of the original sources.
- **Data discovery.** Data may be published in different manners, through publications, web-indexes, etc.
- **Data analysis.** Data can be used by others to achieve different goals.
- **Repurposing.** Data can be used outside of their original framework, restructuring or combining it to satisfy diverse purposes.

3.2 Australian National Data Service

In late 2007, the *Australian National Data Service* (ANDS) was founded with the objective of create a national data management environment. ANDS established a set of verbs, denominated *Data Sharing Verbs*, that describe the entire life cycle of the data [3]:

- **Create.** *Create* (or *collect* for disciplines with an observational focus) is about the kinds of metadata that could be collected and the tools to fulfill this collection task.
- **Store.** This *Data Sharing Verb* remarks the need for stable and web-accessible storage, taking care about the appropriate storing of data.
- **Describe.** The more information inside the storage, the more difficult its discovery, access and exploit is. Annotating the data with the proper metadata solves this issue.



Fig. 1. Combined Life Cycle Model (ownership: DDI Alliance).

- **Identify.** The application of this verb implies the proper identification of each data resource, assigning a persistent identifier to each of them.
- **Register.** This Verb pertains to record the descriptions of the different data collections with one or more public catalogues.
- **Discover.** To improve data-reusing, ANDS suggests to enable different discovery services.
- **Access.** To guarantee the appropriate access to data, ANDS advises to provide a suitable search engine to retrieve these data. If data is not electronically available, ANDS recommends to provide contact details to get the data in conventional formats.
- **Exploit.** *Exploit*, the final *Data Sharing Verb*, comprises the tools, methodologies and support actions to enable reutilisation of data.

3.3 Ecoinformatics data life cycle

Michener and Jones define in [4] the concept of “ecoinformatics”: *a framework that enables scientists to generate new knowledge through innovative tools and approaches for discovering, managing, integrating, analysing, visualizing and preserving relevant biological, environmental, and socio-economic data and information.* To manage these data, the following data life cycle has been defined, as can be seen at Figure 2:

- **Plan.** This step involves the confection of a data management planning.
- **Collect.** This step considers both manual (hand-written data sheets) and automatic (sensor networks) data-gathering methods.
- **Assure.** Quality assurance and quality control (QA/QC), an issue addressed in previously mentioned models is not taken into account. Michener and Jones proposal is based on developing methods to guarantee the integrity of data. Quality assurance can also include the definition of standards for formats, codes, measurement units, metadata, etc.

- **Describe.** As other data life cycle models, this model remarks the value of the metadata to answer questions about *who*, *when*, *where*, *how* and *why*.
- **Preserve.** Data preservation implies the storage of the data and metadata, ensuring that these data can be verified, replicated and actively curated over time.
- **Discover.** The authors describe the data discovering process as *one of the greatest challenges*, as many data are not immediately available because they are stored in individual laptops. The main challenges to publish the data in a proper way are related to the creation of catalogues and indexes, and about the implementation of the proper search engines.
- **Integrate.** Integrating data from different and heterogeneous sources can become a difficult task, as it requires *understanding methodological differences, transforming data into a common representation, and manually converting and recording data to compatible semantics before analysis can begin*.
- **Analyze.** As well as the importance of a clear analysis step, this models remarks the importance of documenting this analysis with sufficient detail to enable its reproduction in different research frameworks.



Fig. 2. Data life cycle in ecoinformatics. Taken from [4].

3.4 UK Data Archive

Another data life cycle model is the one proposed by *UK Data Archive*¹. This model is oriented to help researchers publish their data in a manner that allows other researchers to continue their work independently. In Figure 3, the following stages can be observed:

- **Creating data.** Creating the data involves the design of the research question, planning how data are going to be managed and their sharing strategy. If we want to reuse existing data, we have to locate existing data and collect them. Whether data is new or existing, at this stage the metadata has to be created.
- **Processing data.** Like in other models, at this stage the data is translated, checked, validated and cleaned. In the case of confidential data, data needs to be “anonymized”. The UK Data Archive recommends the creation of metadata at this stage too.
- **Analysing data.** At this stage data are interpreted and derived into visualizations or reports. In addition, the data are prepared for preservation, as mentioned in the following stage.
- **Preserving data.** To preserve data properly, they are migrated to the best format and stored in a suitable medium. In addition to the previously created metadata, the creating, processing, analysis and preserving processes are documented.
- **Giving access to data.** Once the data is stored, we have to distribute our data. Data distribution may involve controlling the access to them and establish a sharing license.
- **Re-using data.** At last, the data can be re-used enabling new research topics.

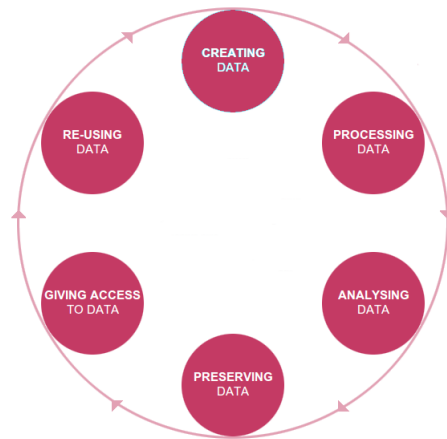


Fig. 3. Data life cycle proposed by UK Data Archive.

¹ <http://www.data-archive.ac.uk/create-manage/life-cycle>

3.5 The LOD2 Stack Data Life Cycle

The last analyzed data life cycle has been developed under the LOD2² project. This project proposes a technological and methodological stack which supports the entire life cycle of Linked Data [?]. As Figure 4 shows, the proposed life cycle phases are the following:



Fig. 4. Linked Data life cycle from LOD2 stack.

- **Storage.** As RDF data presents more challenges than relational data, they propose the collaboration between known and new technologies, like columnstore technology, dynamic query optimization, adaptive caching of joins, optimized graph processing and cluster/cloud scalability.
- **Authoring.** LOD2 provides provenance about data collected through distributed social, semantic collaboration and networking techniques.
- **Interlinking.** At this phase, LOD2 offers approaches to manage the links between different data sources.
- **Classification.** This stage deals with the transformation of raw data into Linked Data. This transformation implies the linkage and integration of data with upper level ontologies.
- **Quality.** Like other models, LOD2 develops techniques for assessing quality based on different metrics.

² <http://lod2.eu/>

- **Evolution/Repair.** At this stage, LOD2 deals with the dynamism of the data from the Web, managing changes and modifications over the data.
- **Search/Browsing/Exploration.** This stage is focused on offering Linked Data to final users through different search, browsing, exploration and visualization techniques.

3.6 A common data life cycle for smart cities

Based on these data life cycle models, we proposed a common data life cycle for managing data in smart cities. As can be seen at Figure 5, the different stages of mentioned models have been aggregated, forming our proposed model.



Fig. 5. Proposed model.

The different stages of this model, which are going to be explained widely in following sections, are:

- **Capture.** The first step of our model consists of capturing the data. In a smart city environment, there are a lot of alternatives to capture data, like sensors, data published by public administration, social networks or in more traditional way like surveys.
- **Process.** Once the required data is captured, they are prepared to be stored and in need of proper methods to explore them. This processing involves the analyzing, refining, cleaning, formatting and transformation of the data. At this stage we also suggest to search links and relationships with other data.
- **Store.** The storage of data is, probably, the most delicate action in the life cycle. Above the storage all the analysis tools are build, and is the “final endpoint” when someone requests our data. A suitable

storage should have indexing, replication, distribution and backup capabilities, among other services.

- **Publish.** Most of the previously mentioned models prioritize the analysing stage over the publication stage. In our model, we defend the opposite approach for a very simple reason: when you consume your data before the publication of them, and using different processes as the rest of the people who is going to consume them, you are not making enough emphasis on publishing these data correctly. Everybody has ever met a research paper or an application in which accessing the data was difficult, or, when once the data was collected it became totally incomprehensible. To avoid this issue, we propose to publish the data before consuming them, and consume them through the same way as the rest of the people does.
- **Consume.** Once the data is published, we use the provided methods to consume the data. This data consumption involves the data mining, analytics or reasoning.
- **Visualize.** To understand the data properly, designing suitable visualizations is essential to show correlations between data and the conclusions of the data analysis in a human-understandable way.
- **Discovery.** This step is about discovering data from external sources, or about our data being discovered by others. This step enables a new execution of the entire life cycle, avoiding the capturing stage, since we are using third parties data. Otherwise, the data life of different datasets can elapse concurrently, being merged at processing stage.

4 Identified challenges

Taking into consideration the large amounts of data present at smart cities, data management's difficulty can be described in terms of:

- Volume
- Variety
- Velocity

These three variables can also be found in *Big Data*-related articles (also known as the *Big Data's Vs*) [?] [?], so it's not surprising at all that smart cities are going to deal with Big Data problems in the near future (if they are not dealing with them right now).

Data scientists need to take into account the following three variables, which could overlap in certain environments. Should this happen, each scenario will determine the most relevant factors of the process, generating un-desired drawbacks on the other ones.

4.1 Volume

The high amount of data used and generated by cities nowadays needs to be properly analysed, processed, stored and eventually accessible. This means conventional IT structures need to evolve, enabling scalable storage technologies, distributed querying approaches and massively parallel processing algorithms and architectures.

However, big amounts of data should not be seen as a drawback attached to smart cities. The larger the datasets, the better analysis algorithms can perform, so deeper insights and conclusions should be expected as an outcome. These could ease the decision making stage.

As management consultant Peter Drucker once said: "*If you can't measure it, you can't manage it*", thus leaving no way to improve it either. This adage manifests that if you want to take care of some process, but you are not able to measure it or you can't access the data, you will not be able to manage that process. That being said, the higher amounts of data available, the greater the opportunities of obtaining useful knowledge will become.

4.2 Variety

Data is rarely found in a perfectly ordered and ready for processing format. Data scientists are used to work with diverse sources, which seldom fall into neat relational structures: embedded sensor data, documents, media content, social data, etc. Variety in data sources, in storage systems, in data-types to get together in a unified analytic...

There is also an increasing concern on data trustworthiness. As pointed out by [?]. *data provenance is fundamental to understanding data quality*. They also highlight that established information storage systems may not be adequate to keep semantic sense of data.

In a previous research [?], we introduced a provenance data model to be used in user-generated Linked Data datasets, which follow W3C's PROV-O ontology³.

³ <http://www.w3.org/TR/prov-o/>

Several efforts are trying to convert existing data in high quality data, providing an extra confidence layer in which data analysts can rely.

4.3 Velocity

Finally, we must assume that data generation is experiencing an exponential growth. That forces our IT structure to not only tackle with volume issues, but with processing rates. A widely spread concept among data businesses is that sometimes you can not rely on five-minute-old data for your business logic.

That's why *streaming data* has moved from academic fields to industry to solve velocity problems. There are two main reasons to consider streaming processing:

- Sometimes, input data is too fast to store in their entirety without rocketing costs.
- If applications mandate immediate response to the data, batch processes are not suitable. Due to the rise of smartphone applications, this trend is increasingly becoming a common scenario.

5 Open Linked data as a viable approach

In the previous section, we identified some of the challenges smart cities will need to face in the following years. The data lifecycle model proposed at Figure 5 relies on Linked Open Data principles to try to solve these issues, reducing costs and enabling third parties to develop new business models on top of Linked Open Data.

Next we describe how Linked Open Data principles could help in the model's stages:

5.1 Capture

5.2 Process

5.3 Store

5.4 Publish

5.5 Discovery

5.6 Consume

At this stage, we focus on consuming data for our logic processes, should they involve data mining algorithms, analytics, reasoning, etc.

Whereas complex processing algorithms can be used independently of the dataset format, Linked Open Data can greatly help at reasoning purposes. Linked Open Data allows to describe entities using constraints and restriction rules (belonging, domain, range, etc.), favoring the inference of new information from the existing one. Thanks to Linked Data, we are not feeding our algorithms with raw data (numbers, strings, values...), but with semantically meaningful data (height in cm, world countries, company names...).

Using Linked Data we can make use of semantics to enrich input data and the processing algorithms, resulting in higher quality outputs.

While dealing with Linked Data streams, stream-querying languages such as CQELS's⁴ language (an extension of the declarative SPARQL 1.1 language using the EBNF notation) can greatly help in the task. CQELS[?] (Continuous Query Evaluation over Linked Stream) is a native and adaptive query processor for unified query processing over Linked Stream Data and Linked Data developed at DERI Galway⁵.

5.7 Visualize

In order to make meaning from data, humans have developed a great ability to understand visual representations. The main objective of data visualization is to communicate information in a clean and effective way through graphical means. It's also suggested that visualization should also encourage users engagement and attention.

⁴ <https://code.google.com/p/cqels/>

⁵ <http://www.deri.ie/>

The "*A picture is worth a thousand words*" saying reflects the power images and graphics have when expressing information.

As Linked Data is based on subject-predicate-object triples, graphs are a natural way to represent triple stores, where subject and object nodes are inter-connected through predicate links. When further analysis is applied on triples, a diverse variety of representations can be chosen to show processed information: charts, infographics, flows, etc.

6 Evaluation

7 Lessons learned

8 Further research

References

1. Bizer, C., Boncz, P., Brodie, M.L., Erling, O.: The meaningful use of big data: four perspectives – four challenges. *SIGMOD Rec.* **40**(4) (January 2012) 56–60
2. Initiative, D.D.: Overview of the DDI version 3.0 conceptual model (April 2008)
3. Ball, A.: Review of data management lifecycle models. (February 2012)
4. Burton, A., Treloar, A.: Designing for discovery and re-use: the ‘ANDS data sharing verbs’ approach to service decomposition. *International Journal of Digital Curation* **4**(3) (2009) 44–56
5. Michener, W.H., Jones, M.B.: *Ecoinformatics: supporting ecology as a data-intensive science.* (2012)
6. Auer, S., Bühmann, L., Dirschl, C., Erling, O., Hausenblas, M., Isele, R., Lehmann, J., Martin, M., Mendes, P.N., Van Nuffelen, B., et al.: Managing the life-cycle of linked data with the LOD2 stack. In: *The Semantic Web–ISWC 2012*, Springer (2012) 1–16
7. Zikopoulos, P., Eaton, C., et al.: *Understanding big data: Analytics for enterprise class hadoop and streaming data.* McGraw-Hill Osborne Media (2011)
8. Russom, P.: *Big data analytics.* TDWI Best Practices Report, Fourth Quarter (2011)
9. Buneman, P., Davidson, S.B.: *Data provenance—the foundation of data quality* (2013)
10. Emaldi, M., Pena, O., Lázaro, J., Vanhecke, S., Mannens, E.: To trust, or not to trust: Highlighting the need for data provenance in mobile apps for smart cities
11. Le-Phuoc, D., Dao-Tran, M., Parreira, J.X., Hauswirth, M.: A native and adaptive approach for unified processing of linked streams and linked data. In: *The Semantic Web–ISWC 2011.* Springer (2011) 370–388

9 Acronyms and terms