Date: May 2014

# Open Data Ireland:

# Open Data Publication Handbook

*Authors: Deirdre Lee, Richard Cyganiak, Stefan Decker*

*Insight Centre for Data Analytics, NUI Galway*

# Contents

# Glossary

API        Application Programming Interface

CKAN       Comprehensive Knowledge Archive Network

CSO        Central Statistics Office of Ireland

CSV        Comma-Separated Values

DCAT       Data Catalog Vocabulary

HIQA       Health Information and Quality Authority

IETF       Internet Engineering Task Force

ISDE       Irish Spatial Data Exchange

ISDI       Irish Spatial Data Infrastructure

ISO        International Organization for Standardization

OASIS      Organization for the Advancement of Structured Information Standards

OGC        Open Geospatial Consortium

OSi        Ordnance Survey Ireland

SPARQL     SPARQL Protocol and RDF Query Language

URI        Unique Reference Identifier

W3C        World Wide Web Consortium

# 1 Introduction

The Open Data Publishing Handbook is designed as a step-by-step guide for public bodies wishing to publish Open Data on the National Open Data Portal (data.gov.ie). It is complemented by the Open Data Ireland Best Practice Handbook (Lee & Cyganiak 2014), a detailed guide to the key elements of an Open Data Ecosystem with a focus on the Irish context, and the Open Data Ireland Roadmap (Lee et al. 2014). The Open Data Publishing Handbook is aimed at public body employees who will be responsible for the publication of Open Data. Some technical knowledge is assumed.

# 2 Step-by-Step Guide to Open Data Publishing

## Step 1 Carry out a Data Audit

Before deciding what data to publish as Open Data, public bodies need to have an overview of what data they currently manage, and could therefore potentially release as Open Data. This can be a challenging task, as data in large organisations is typically dispersed over multiple websites, databases, shared storage, and personal computers. A data audit is carried out to establish an inventory of what datasets currently exist.

A data audit should produce a complete list of datasets that fall under the scope and terms of the audit. "Datasets" in the context of Open Data include databases, registers, ongoing data collections and surveys, geospatial datasets, and so on. Pure text-based information, such as emails, memos, reports and press releases, are not considered as "data" for this purpose.

The results should be collected according to a standard template structured with the target open data portal's metadata schema in mind.

### Existing data catalogues

A public body may already manage a data catalogue, or contribute metadata to a cross-sector catalogue such as the Irish Spatial Data Exchange. These sources should be reviewed first to establish a starting point. A data audit may already have been performed to provide content for these catatalogues.

### Online data audit

Websites managed by the public body can be reviewed to identify datasets that are already released on the Web. Relevant keywords to look for are terms such as 'data', 'database', 'download', 'API', 'export', 'statistics', 'reports'. Relevant file formats to search for include:

- Microsoft Excel (XLS, XLSX)
- Comma-Separated Values (CSV, TSV)
- Extensible Markup Language (XML)
- Keyhole Markup Language and Geography Markup Language (KML, GML)
- Resource Description Framework (RDF)
- ESRI Shapefiles (SHP)
- Archived forms of the formats above (ZIP, GZ, BZ2, etc.)

General web searches can be performed using the advanced search of a search engine, such as Google. Search results can be restricted to a particular domain and/or to a particular file type. The formats and keywords are similar to the web crawl. The syntax to use is:

| site: | filetype: |
| --- | --- |

For example:

| site:http://www.per.gov.ie/ | filetype:xls |
| --- | --- |

A more advanced technique that can be used to deal with largeer collections of datasets is to use tools such as the *Scrapy* screen-scraping toolkit[1] and the *Xidel* HTML/XML data extraction tool[2], to extract links to data files and relevant metadata.

---

[1] http://scrapy.org/

### Internal Data Audit

An internal data audit should be carried out to identify datasets that are not currently released on websites. This can be carried out by the Open Data officer assigned by each public body. An internal data audit may involve staff interviews and reviews of technical assets, such as relational databases, shared file storage, or personal computer drives.

## Step 2    Select what Data to Publish

There are three main avenues to selecting what data to publish:

   a)  Upgrade data already in the public domain to 'Open Data'
   b)  Follow international best practice
   c)  Demand-drive

### Upgrade Data Already in the Public Domain to 'Open Data'

As demonstrated in the Best Practice Handbook, there are already a number of Irish data catalogues, such as the CSO's statistical data catalogue StatCentral[3] and Irish Spatial Data Exchange[4]. The first step for any public body is to ensure that any data catalogue they manage is published as Open Data, i.e. data that is freely available in a machine-readable format under an Open License including commercial reuse (see steps below). Even if the data is not in a data catalogue, but is uploaded as part of a standard website, the data should be associated with the national Open Data License, which will be available on data.gov.ie

### Follow International Best Practice

There are a number of international indicators of high-value datasets: the G8 Open Data Charter, the Open Data Barometer, and the Open Data Census. In the Best Practice Handbook, we have identified common high-value datasets across all available G8 Charter National Action Plans, the Danish Basic-Data Registers, the Open Data Barometer and the Open Data Census. These datasets are shown in Table 1 and should be prioritised for publication as Open Data in Ireland.

---

[2] http://videlibri.sourceforge.net/xidel.html
[3] http://statcentral.ie/
[4] http://catalogue.isde.ie/

**Table 1: Common High-Value Datasets**

| G8 Open Data Charter Category | Common High-Value Datasets |
|---|---|
| Companies | • Company register<br>• Insolvency and bankruptcy records |
| Crime and Justice | • Crime statistics<br>• Justice statistics<br>• Justice spending |
| Earth Observation | • Meteorological<br>• Fishing/Hunting levels<br>• Agriculture |
| Education | • School attendee<br>• Post-education<br>• School locations |
| Energy and Environment | • Pollution<br>• Water quality<br>• Air quality<br>• Natural resources<br>• Waste<br>• Energy consumption |
| Finance and contracts | • Government budgets<br>• Government spending<br>• Tenders/procurement |
| Geospatial | • National maps<br>• Thematic geo-information<br>• Environmental geo-information<br>• Local/administrative boundaries<br>• Topographical geo-information<br>• Postcodes and addresses |
| Global Development | • Development aid<br>• International assistance |
| Government Accountability and Democracy | • Government structures and contacts<br>• Government salaries and pay-scales<br>• Legislation<br>• Hospitality/gift<br>• Election results |
| Health | • Health performance Drug/prescription<br>• Restaurant hygiene |
| Science and Research | • Research |
| Social Mobility and Welfare | • Housing<br>• Employment/unemployment<br>• Social security/welfare |
| Statistics | • National statistics<br>• Census |
| Transport and Infrastructure | • Public transport schedules<br>• Public transport stops<br>• Road network<br>• Road traffic accidents |

### Demand-driven Dataset Selection

The selection of what data to publish should be demand-driven, as it is ultimately the use of Open Data that drives its value. While there are many datasets that are internationally recognised as being high-value, other datasets will be of interest based on the particular jurisdiction or stakeholder group. The most efficient way to decide what data to release is to ask potential data-users for suggestions.

This can be facilitated on the Open Data portal with a simple suggestion form or through wider consultations as described below. To prioritise what datasets are most in-demand, users could rate/like suggested datasets. The most important element with such an approach is the responsiveness of the data-publisher. While it may not be possible to publish all data requested, it is best practice to acknowledge each request and comment either, if and when it will be released, or why it is not possible to release it.

## Step 3    Ensure Data Protection Laws are Adhered to

Open Data should not include personal or sensitive data that could be linked back to an individual. In all cases of Open Data the citizen's fundamental right to privacy must be protected and the data publisher must comply with data protection principles set out by the Office of the Data Protection Commissioner[5]. Statistical data that has its origins in information concerning individuals, such as A&E admittance numbers, crime rates, or household deprivation statistics, may be published as Open Data. However, this data must be anonymised/aggregated using recognised statistical methods. If required, guidance on statistical methods should be sought from the in-house statistician, or the CSO. If there are any concerns in relation to data privacy, guidance should be sought from the Data Protection Commissioner.

## Step 4    Associate Data with an Open License

What makes Open Data 'open' is that it is free to be used, including for commercial use. Associating an Open License with Open Data is necessary to ensure the legal grounding for its potential reuse. For a data user (individual/organisation/company/etc.) wishing to use and build on top of public data, they require assurance of what they legally can and can't do with the data. If no license is specified, each data-user must contact the data publisher on a case-by-case basis.

As outlined in the Open Data Ireland Roadmap, the Steering and Implementation Group will identify a standard Open License to associated with all Irish Open Data, based on the recommendations of the Best Practice Handbook. This Open Data License will be available on data.gov.ie. Until a standard Open License is defined, Open Data should be associated with the Irish PSI License[6] or Creative Commons 4.0 License[7].

---

[5] https://www.dataprotection.ie/
[6] http://psi.gov.ie/files/2010/03/PSI-Licence.pdf
[7] http://creativecommons.org/version4

## Step 5    Publish Data as 3- to 5-star Open Data

In order to ensure Open Data is as easy as possible for potential data users to reuse, it should be available in an open (non-proprietary), machine-readable format. This not only facilitates processing and analysis of each dataset, but it also supports the integration of multiple datasets. Tim Berners-Lee, the inventor of the Web and Linked Data initiator, suggested a 5-star deployment scheme for Open Data[8]. The greater the number of stars, the more reusable the data is, and the easier it is to reuse and interconnect data.

★  Publish data on the Web under an Open License

★ ★  Publish data in a machine-readable, structured format

★ ★ ★  Publish data in a non-proprietary format

★ ★ ★ ★  Use URIs to identify things, so that people can point at your stuff

★ ★ ★ ★ ★ Link your data to other data to provide context

Open Data published to data.gov.ie should be at least 3-star, meaning that the data is associated with an Open License, is machine-readable and non-proprietary, as shown in Figure 1. CSV is one of the most widely accepted formats for publishing Open Data. Data can be made available in multiple formats.
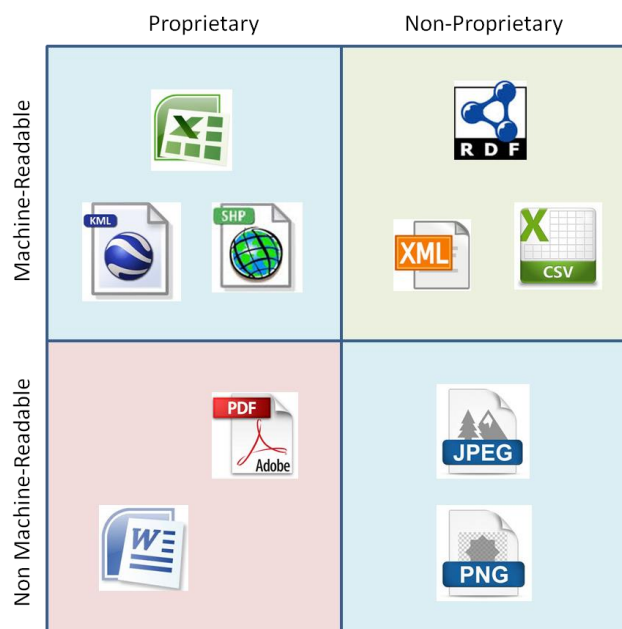


**Figure 1: Machine-Readable and Non-Proprietary Data Formats**

Some business applications can directly export data in CSV format. Software tools such as FME Server[9] or Talend[10] can be used for exporting SQL database query results as CSV. Excel spreadsheets can be exported as CSV, but care should be taken to re-structure presentation-oriented spreadsheets into re-use oriented 'tidy data' first where possible (Wickham 2014).

---

[8] http://www.w3.org/DesignIssues/LinkedData.html
[9] http://www.safe.com/fme/fme-server/
[10] https://www.talend.com/

## Step 6    Associate Data with Standardised Metadata

Metadata provides structured information about a dataset, so that a potential user can understand more about its context, for example, date, authors, purpose, scope. Metadata is an integral element of data management, as it facilitates the discovery and efficient use of the data. Metadata should be created along with the dataset, published alongside the data and updated whenever the dataset is updated. It is essential to embed standardised metadata production and management into all legacy and new data processes. For Open Data published via data.gov.ie, the metadata should be created when the data is being uploaded. CKAN provides a standard metadata upload mechanism. There is a complete and functional mapping of the CKAN dataset schema to Linked Data formats, including the W3C recommended DCAT vocabulary. To access Linked Data formats, the CKAN API can be called in the usual way but with the format you would like to be returned specified. More information can be found in the CKAN Documentation[11].

## Step 7    Use Data Standards

Data standards help give a common meaning to data. This is especially important when data is being used by a third-party, data is being integrated from different sources, or when data is being shared across public bodies. Data standards not only define the meaning of certain concepts, but also how concepts relate to each other. Data standards should have both a human-readable and machine-readable representation. When publishing Open Data, first try and reuse international standards defined by reputable standards organisations, such as ISO, the European Commission, W3C, IETF, OGC and OASIS. If international standards are unavailable or unsuitable, use national standards. For specific topics such as geospatial, statistics, or health, promote national standards defined by the responsible organisation (OSI, CSO, HIQA, etc.). A list of recommended data standards for use by Irish public bodies should be available on data.gov.ie. For guidance on data standards, contact the CSO.

## Step 8    Use Unique Identifiers

Unique identifiers enable specific resources to be distinguished across multiple datasets. Unique identifiers are especially important for reference datasets that are used by multiple agencies, e.g. school roll numbers, road numbers, company registration numbers, statute book numbers. Universal Resource Identifiers (URIs) are unique identifiers that are used on the World Wide Web. It is best practice to publish Open Data with URIs, so that the data is accessible online. Creating a common URI structure is important, so that they are persistent, i.e. resilient to change.

The Steering and Implementation Group will develop and adopt a national URI strategy for Ireland. All URIs created by public bodies should be created according to the national strategy and be documented on data.gov.ie.

---

[11] http://docs.ckan.org/en/1117-start-new-test-suite/linked-data-and-rdf.html

## Step 9    Provide Access to the Data

It is important that Open Data is published in a timely and accessible manner to preserve the value of the data and to ensure data is available to the widest range of users and for the widest range of purposes. There are a couple of different methods to facilitate access to Open Data: as bulk data (data-dump), via an Application Programming Interface (API), as a feed, via a SPARQL endpoint, etc. The most commonly used and useful method is as bulk data, meaning that the complete dataset should be available in downloadable form. Bulk data should not be seen as stale data and should be kept up-to-date by the data providers. Real-time data can be published as feeds. If there is a demand from users for an API, public bodies should consider introducing one. However public bodies should (i) use existing API standards whenever possible, e.g. the OGC web services or SPARQL, (ii) before creating a new API, collaborate with potential users on its structure, and (iii) provide complete documentation for each API.

## Step 10  Publish Data on the National Open Data Portal, data.gov.ie

Open Data should be easy to find and easy to access. A data catalogue provides a registry or listing of all existing datasets and a pointer (URL) to where the data can be accessed. A data catalogue of all available Open Data datasets is usually the key component of an Open Data Portal, along with social, news and community elements. Ireland's National Open Data Portal is data.gov.ie[12], which is built using CKAN[13], Open Source data portal software.
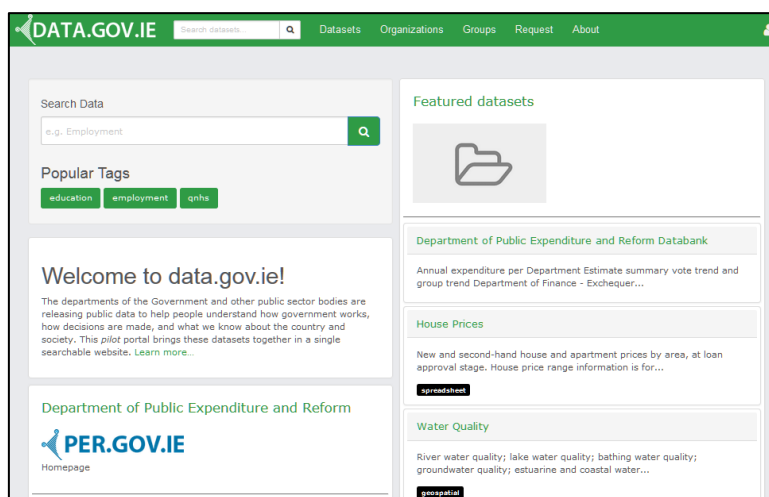


**Figure 2: National Open Data Portal data.gov.ie**

For a detailed guide on how to organize, publish and find data to data.gov.ie, see CKAN's online user guide[14]. CKAN also has an API for developers who want to write code that interacts with CKAN sites and their data. The API is documented in CKAN's API guide[15]. Adding and updating information in data.gov.ie requires a user account and appropriate permissions[16].

---

[12] http://data.gov.ie/
[13] http://ckan.org/
[14] http://docs.ckan.org/en/latest/user-guide.html
[15] http://docs.ckan.org/en/latest/api/index.html
[16] http://data.gov.ie/contact

# 3 Bibliography

Lee, D. & Cyganiak, R., 2014. *Open Data Ireland: Best Practice Handbook*,

Lee, D., Decker, S. & Cyganiak, R., 2014. *Open Data Ireland Roadmap*,

Wickham, H., 2014. Tidy Data. *Journal of Statistical Software*, VV(Ii). Available at: http://vita.had.co.nz/papers/tidy-data.pdf.