

Using a Jupyter Notebook to perform a reproducible scientific analysis over semantic web sources

Alasdair J G Gray ([ORCID:0000-0002-5711-4872](http://orcid.org/0000-0002-5711-4872) (<http://orcid.org/0000-0002-5711-4872>))

Heriot-Watt University, Edinburgh, UK



(<http://creativecommons.org/licenses/by/4.0/>).

This work is licensed under a [Creative Commons Attribution 4.0 International License](http://creativecommons.org/licenses/by/4.0/)

(<http://creativecommons.org/licenses/by/4.0/>).

Abstract: In recent years there has been a reproducibility crisis in science. Computational notebooks, such as Jupyter, have been touted as one solution to this problem. However, when executing analyses over live SPARQL endpoints, we get different answers depending upon when the analysis in the notebook was executed. In this paper, we identify some of the issues discovered in trying to develop a reproducible analysis over a collection of biomedical data sources and suggest some best practice to overcome these issues.

Keywords: Reproducibility, Computational Notebooks

1. Introduction

In recent years there has been a reproducibility crisis in science ([Baker; 2016](#)). Open science advocates to overcome this crisis by openly publishing your work, included data, methods, and analysis ([Pearce; 2012](#)). Computational notebooks such as Jupyter ([Kluyver et al; 2016](#)) that combine analysis with narrative are one approach to publishing scientific analysis methods. The approach is known as Literate Programming ([Knuth; 1984](#)). A key idea of Literate Programming is that the documentation of the analysis code is kept up to date with the code, as they are edited in the same environment. When writing the analysis, the author is forced to think about how others can understand and interpret the analysis, thus writing analysis in a way to aid understanding of the code ([Piccolo, Frampton; 2016](#)). Ideally, the computation notebook becomes the publication, as this notebook attempts. Specifically, in this notebook we investigate the feasibility of publishing a reusable data analysis paper as a Jupyter Notebook in the Semantic Web community.

There is existing work on publishing computational workflows via notebooks ([Kluyver et al; 2016](#)), and works in several scientific disciplines that have also used Jupyter Notebooks, e.g. Astronomy ([Farr et al; 2017](#), [Medvedev et al; 2016](#)), Biology ([Grüning et al; 2017](#)), and Oceanography ([Turner and Gill; 2018](#)). In their paper Kluyver et al ([2016](#)) recognise some of the limitations of academic publishing within a Jupyter notebook, particularly in the lack of support for citations. In this work, we are looking at performing an analysis over pharmacology Linked Data. The contents of this notebook provide a good starting point for those who want to do an analysis of the content of the chemical substances known in some of the leading chemistry databases used to support early stage pharmacology research, and supports the understanding of the evolution of their content. The analysis will compare the compound count between ChEBI ([Hastings et al; 2016](#)), ChEMBL ([Gaulton et al; 2017](#)), DrugBank ([Wishart et al; 2018](#)), and Guide to Pharmacology ([Harding et al; 2018](#)).

In the rest of this notebook, we will first setup our computational environment in [Section 2](#). We will then perform a simple, but crucial, analysis in [Section 3](#) that does a crude comparison over four major pharmacology datasets. We will then discuss and conclude the paper in [Section 4](#).

2. Method

2.1 Computational Environment

This notebook was prepared using [Jupyter server \(http://jupyter.org/\)](http://jupyter.org/) version 5.0.0 running on Python 3.6.3 distributed by [Anaconda, Inc. \(https://anaconda.org/\)](https://anaconda.org/) (default, Oct 6 2017, 12:04:38) [GCC 4.2.1 Compatible Clang 4.0.1 (tags/RELEASE_401/final)].

2.2 Configuring the Environment

The following cells load various support libraries and define functions. The functions enable the running of queries over SPARQL endpoints and printing the results.

First we will pull in various libraries that we will use. At the time of publication, the following versions of these libraries were used:

- [SPARQLWrapper \(https://rdflib.github.io/sparqlwrapper/doc/1.8.1/\)](https://rdflib.github.io/sparqlwrapper/doc/1.8.1/): 1.8.1
- [JSON \(simplejson \(https://pypi.org/project/simplejson/\)\)](https://pypi.org/project/simplejson/): 3.15

The SPARQLWrapper dependency is captured in the [requirements.txt](https://github.com/AlasdairGray/SemSci2018/blob/master/requirements.txt) (<https://github.com/AlasdairGray/SemSci2018/blob/master/requirements.txt>) file on GitHub which enables this notebook to be run through the [Binder Service \(https://mybinder.org/\)](https://mybinder.org/).

launch binder

(<https://mybinder.org/v2/gh/AlasdairGray/SemSci2018/master?filepath=SemSci2018%20Publication.ipynb>)

In [1]:

```
from SPARQLWrapper import SPARQLWrapper, JSON
```

Ideally at this point we would print the version number of each module imported automatically. However the `print (SPARQLWrapper.__version__)` command does not work when using the `from SPARQLWrapper import` mechanism.

2.3 Configuring SPARQL endpoints

The endpoints used in this workbook:

- DrugBank: <http://bio2rdf.org/sparql> (<http://bio2rdf.org/sparql>)
- EBI: <https://www.ebi.ac.uk/rdf/services/sparql> (<https://www.ebi.ac.uk/rdf/services/sparql>)
- Guide to Pharmacology: <https://rdf.guidetopharmacology.org/sparql> (<https://rdf.guidetopharmacology.org/sparql>)

Note that the EBI SPARQL endpoint is used to query both ChEBI and ChEMBL ([Jupp et al; 2014](#)). DrugBank is not originally published as RDF, so we use the Bio2RDF conversion and endpoint ([Callahan et al; 2013](#)).

2.3.1 DrugBank Endpoint

In [2]:

```
#Define DrugBank SPARQL endpoint and function to run queries over it
dbSparql = SPARQLWrapper("http://bio2rdf.org/sparql")
dbSparql.setReturnFormat(JSON)
def queryDrugBank(query):
    dbSparql.setQuery(query)
    results = dbSparql.queryAndConvert()
    return results
```

While the Bio2RDF [documentation \(https://github.com/bio2rdf/bio2rdf-scripts/wiki/Bio2RDF-Dataset-Provenance#querying-the-provenance-graph\)](https://github.com/bio2rdf/bio2rdf-scripts/wiki/Bio2RDF-Dataset-Provenance#querying-the-provenance-graph) includes a query to extract the VOID metadata for the dataset, the query only provides results for the DisGeNet dataset, not DrugBank in which we are interested.

In [3]:

```
query = """
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX prov: <http://www.w3.org/ns/prov#>
PREFIX void: <http://rdfs.org/ns/void#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
SELECT * WHERE {
?dataset rdf:type void:Dataset .
?dataset rdfs:label ?datasetLabel .
?dataset dcterms:created ?creationDate .
?dataset dcterms:creator ?creationScript .
?dataset dcterms:rights ?license .
?dataset prov:wasDerivedFrom ?parentDataSet .
?dataset void:dataDump ?downloadLink .
?dataset void:sparqlEndpoint ?endpointLink .
}
"""
queryDrugBank(query)
```

Out[3]:

```
{'head': {'link': [],
'vars': ['dataset',
'datasetLabel',
'creationDate',
'creationScript',
'license',
'parentDataSet',
'downloadLink',
'endpointLink']},
'results': {'bindings': [], 'distinct': False, 'ordered': True}}
```

According to the Bio2RDF [website \(http://download.bio2rdf.org/files/release/3/release.html\)](http://download.bio2rdf.org/files/release/3/release.html), the version of DrugBank loaded dates from 2014-07-25. The original DrugBank version is not stated.

2.3.2 EBI SPARQL Endpoint

In [4]:

```
#Define EBI SPARQL endpoint and function to run queries over it
#EBI endpoint can be used to query a variety of datasets including ChEMBL and ChEBI
ebiSparql = SPARQLWrapper("https://www.ebi.ac.uk/rdf/services/sparql")
ebiSparql.setReturnFormat(JSON)
def queryEBI(query):
    ebiSparql.setQuery(query)
    results = ebiSparql.queryAndConvert()
    return results
```

The following cell executes a query to determine what version of ChEMBL is currently loaded into the EBI triplestore. At the time of writing, the version was 24.0, which was last modified on 5 January 2018.

In [5]:

```
#Query to determine the version of ChEMBL loaded
query = """
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX pav: <http://purl.org/pav/>

SELECT ?version ?dateOfLastUpdate
where {
    <http://rdf.ebi.ac.uk/dataset/chembl> pav:hasCurrentVersion ?currentVersion.
    ?currentVersion pav:version ?version ;
                    dcterms:hasDistribution ?dist.
    ?dist pav:lastUpdateOn ?dateOfLastUpdate
}
"""
queryEBI(query)
```

Out[5]:

```
{'head': {'vars': ['version', 'dateOfLastUpdate']},
 'results': {'bindings': [{'dateOfLastUpdate': {'datatype': 'http://www.w3.org/2001/XMLSchema#dateTime',
 'type': 'typed-literal',
 'value': '2018-01-05T00:00:00Z'},
 'version': {'type': 'literal', 'value': '24.1'}}]}}
```

We now determine the version of ChEBI that is loaded in the EBI triplestore. Note that slightly different metadata is available for ChEBI so we retrieve the date issued as the last update date is not available. However, the version URL may indicate that the dataset has been updated more recently than the last issue date. When executing in July 2018 the version is <http://rdf.ebi.ac.uk/dataset/chebi/03-07-2018> while the issued date is 2018-01-01.

In [6]:

```
#Query to determine the version of ChEBI loaded
query = """
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX pav: <http://purl.org/pav/>

SELECT  ?version ?dateIssued
where {
  <http://rdf.ebi.ac.uk/dataset/chebi> pav:hasCurrentVersion ?currentVersion.
  ?currentVersion pav:version ?version ;
                  dcterms:issued ?dateIssued
}
"""
queryEBI(query)
```

Out[6]:

```
{'head': {'vars': ['version', 'dateIssued']},
 'results': {'bindings': [{'dateIssued': {'datatype': 'http://www.w3.org/2001/XMLSchema#date',
      'type': 'typed-literal',
      'value': '2018-01-01'},
    'version': {'type': 'literal',
      'value': 'http://rdf.ebi.ac.uk/dataset/chebi/03-07-2018'}}]}}
```

2.3.3 Guide to Pharmacology SPARQL Endpoint

In [7]:

```
#Define Guide to Pharmacology SPARQL endpoint and function to run queries over it
gtpSparql = SPARQLWrapper("https://rdf.guidetopharmacology.org/sparql")
gtpSparql.setReturnFormat(JSON)
def queryGtPdb(query):
    gtpSparql.setQuery(query)
    results = gtpSparql.queryAndConvert()
    return results
```

Currently, the Guide to Pharmacology does not support querying of their metadata. A VoID file is available through their [website](http://www.guidetopharmacology.org/download.jsp#rdf) (<http://www.guidetopharmacology.org/download.jsp#rdf>), but this has not yet been loaded into their triplestore. The latest version of the Guide to Pharmacology data is listed as 2018.2.

2.3.4 Discussion

It would be ideal to query each of the endpoints to programmatically determine the version of the dataset that was loaded at the time of execution, as we did with the EBI endpoint. This can be achieved using the [VoID](https://www.w3.org/TR/void/) (<https://www.w3.org/TR/void/>) vocabulary, although acquiring the knowledge to the structure of the metadata and the graph in which it is loaded is often reliant on the documentation of the data.

3. Analysis: Compare the number of compound structures known in GtoPdb, ChEMBL, DrugBank, and ChEBI

This question arose from discussions with [Chris Southan](https://scholar.google.co.uk/citations?user=y1DsHJ8AAAAAJ&hl=en) (<https://scholar.google.co.uk/citations?user=y1DsHJ8AAAAAJ&hl=en>), University of Edinburgh, whilst creating the Guide to Pharmacology RDF data. Chris currently performs his investigation on [PubChem](https://pubchem.ncbi.nlm.nih.gov/) (<https://pubchem.ncbi.nlm.nih.gov/>). I recreated Chris's search queries (15/3/2018) and generated the output shown in the figure below:

Query	Items found
Search "ChEBI"[SourceName]	91407
Search "ChEMBL"[SourceName]	1729327
Search "DrugBank"[SourceName]	9789
Search "IUPHAR/BPS Guide to Pharmacology"[SourceName]	6969
Search (("IUPHAR/BPS Guide to Pharmacology"[SourceName] AND "DrugBank"[SourceName]) AND "ChEMBL"[SourceName]) AND "ChEBI"[SourceName]	1523

Each of the queries can be rerun by the following links:

- ChEBI: "ChEBI"[SourceName]_([https://www.ncbi.nlm.nih.gov/pccompound?
term=%22ChEBI%22%5BSourceName%5D&cmd=DetailsSearch](https://www.ncbi.nlm.nih.gov/pccompound?term=%22ChEBI%22%5BSourceName%5D&cmd=DetailsSearch)).
- ChEMBL: "ChEMBL"[SourceName]_([https://www.ncbi.nlm.nih.gov/pccompound?
term=%22ChEMBL%22%5BSourceName%5D&cmd=DetailsSearch](https://www.ncbi.nlm.nih.gov/pccompound?term=%22ChEMBL%22%5BSourceName%5D&cmd=DetailsSearch)).
- DrugBank: "DrugBank"[SourceName]_([https://www.ncbi.nlm.nih.gov/pccompound?
term=%22DrugBank%22%5BSourceName%5D&cmd=DetailsSearch](https://www.ncbi.nlm.nih.gov/pccompound?term=%22DrugBank%22%5BSourceName%5D&cmd=DetailsSearch))
- Guide to Pharmacology: "IUPHAR/BPS Guide to Pharmacology"[SourceName]
([https://www.ncbi.nlm.nih.gov/pccompound?term=%22IUPHAR/BPS Guide to
Pharmacology%22%5BSourceName%5D&cmd=DetailsSearch](https://www.ncbi.nlm.nih.gov/pccompound?term=%22IUPHAR/BPS Guide to Pharmacology%22%5BSourceName%5D&cmd=DetailsSearch)).
- Intersection of all sources: ([\("IUPHAR/BPS Guide to Pharmacology"\[SourceName\] AND "DrugBank"
\[SourceName\]\) AND "ChEMBL"\[SourceName\] AND "ChEBI"\[SourceName\]](https://www.ncbi.nlm.nih.gov/pccompound?term=%28%28%22IUPHAR/BPS%20Guide%20to%20Pharmacology%22%5BSourceName%5D%20AND%20%22DrugBank%22%5BSourceName%5D%20AND%20%22ChEMBL%22%5BSourceName%5D%20AND%20%22ChEBI%22%5BSourceName%5D%29%29%5BSourceName%5D&cmd=DetailsSearch)
[\(\[https://www.ncbi.nlm.nih.gov/pccompound?
term=%28%28%22IUPHAR/BPS%20Guide%20to%20Pharmacology%22%5BSourceName%5D%20AND%20%22DrugBank%22%5BSourceName%5D%20AND%20%22ChEMBL%22%5BSourceName%5D%20AND%20%22ChEBI%22%5BSourceName%5D%29%29%5BSourceName%5D&cmd=DetailsSearch\]\(https://www.ncbi.nlm.nih.gov/pccompound?term=%28%28%22IUPHAR/BPS%20Guide%20to%20Pharmacology%22%5BSourceName%5D%20AND%20%22DrugBank%22%5BSourceName%5D%20AND%20%22ChEMBL%22%5BSourceName%5D%20AND%20%22ChEBI%22%5BSourceName%5D%29%29%5BSourceName%5D&cmd=DetailsSearch\)](https://www.ncbi.nlm.nih.gov/pccompound?term=%28%28%22IUPHAR/BPS%20Guide%20to%20Pharmacology%22%5BSourceName%5D%20AND%20%22DrugBank%22%5BSourceName%5D%20AND%20%22ChEMBL%22%5BSourceName%5D%20AND%20%22ChEBI%22%5BSourceName%5D%29%29%5BSourceName%5D&cmd=DetailsSearch)).

We will now recreate the same data by querying the SPARQL endpoints.

In [8]:

```
# Define function to extract count result from JSON SPARQL result set
def extract_count(results):
    """
    Extract the count result from the JSON format
    """
    for result in results["results"]["bindings"]:
        return result["count"]["value"]

#Initialise counts dictionary to store count for each dataset in
counts = {}

# Find the number of InChI Keys stored in ChEBI
query = """
SELECT (count(?inchikey) as ?count)
WHERE {
    ?ligand <http://purl.obolibrary.org/obo/chebi/inchikey> ?inchikey.
}
"""
results = queryEBI(query)
counts['ChEBI'] = extract_count(results)

# Find the number of InChI Keys stored in ChEMBL
query = """
SELECT (COUNT(?substance) as ?count)
FROM <http://rdf.ebi.ac.uk/dataset/chembl>
WHERE {
    ?substance a <http://semanticscience.org/resource/CHEMINF_000059>
}
"""
results = queryEBI(query)
counts['ChEMBL'] = extract_count(results)

#Find the number of InChI Keys stored in DrugBank
query = """
SELECT (COUNT(?s) as ?count)
WHERE {?s a <http://bio2rdf.org/drugbank_vocabulary:InChIKey>}
"""
results = queryDrugBank(query)
counts['DrugBank'] = extract_count(results)

# Find the number of InChI Keys stored in GtPdb
query = """
PREFIX gtpo: <http://rdf.guidetopharmacology.org/ns/gtpo#>
SELECT (count(?inchikey) as ?count)
WHERE {
    ?ligand gtpo:inChIKey ?inchikey.
}
"""
results = queryGtPdb(query)
counts['GtPdb'] = extract_count(results)

# Print results
print (counts)

{'ChEBI': '90510', 'ChEMBL': '1820035', 'DrugBank': '6810', 'GtPdb':
'7146'}
```

The table below captures the counts of the number of InChI keys in PubChem in March 2018, and then the

SPARQL queries which have been executed in June and July of 2018.

Dataset	PubChem (2018-03-15)	SPARQL (2018-06-08)	SPARQL (2018-07-24)
ChEBI	91,407	184,393	90,510
ChEMBL	1,729,327	1,820,035	1,820,035
DrugBank	9,789	6,810	6,810
Guide to Pharmacology	6,969	7,065	7,146
Intersection	1,523	--	--

The ChEMBL and Guide to Pharmacology datasets contain more data in the SPARQL endpoint than the PubChem version. This is to be expected as the datasets available through the SPARQL endpoints are newer and the datasets are adding more content. The Guide to Pharmacology change between June and July corresponds with the release of version 2018.3 (<https://blog.guidetopharmacology.org/2018/06/28/database-release-2018-3/>) on 28 June 2018. The exception to this is DrugBank, which has seen a decrease. This is due to the Bio2RDF endpoint offering an older version of the dataset.

On PubChem we were able to calculate the intersection of the datasets, i.e. we could count the number of substances that are contained in all the datasets. This has not been possible with the SPARQL endpoints due to the distributed nature of the endpoints and the unreliability of federated queries. There are also complications due to the way that the substance structures (InChI Keys) are represented in the various datasets. The following queries demonstrate the output of the first 5 results of each dataset.

In [9]:

```
# Show 5 InChI Keys stored in ChEBI
query = """
SELECT ?inchikey
WHERE {
    ?ligand <http://purl.obolibrary.org/obo/chebi/inchikey> ?inchikey.
} LIMIT 5
"""
queryEBI(query)
```

Out[9]:

```
{'head': {'vars': ['inchikey']},
 'results': {'bindings': [{'inchikey': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'MHWLWQUZZRMNGJ-UHFFFAOYSA-N'}},
  {'inchikey': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'DURULFYMVIFBIR-UHFFFAOYSA-N'}},
  {'inchikey': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'MYKUKUCHPMASKF-VIFPVBQESA-N'}},
  {'inchikey': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'IEDVJHCEMCRBQM-UHFFFAOYSA-N'}},
  {'inchikey': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'ZZUFCTLCJUWOSV-UHFFFAOYSA-N'}}]}}
```

In [10]:

```
# Show 5 InChI Keys stored in ChEMBL
query = """
SELECT ?inchiKey
FROM <http://rdf.ebi.ac.uk/dataset/chembl>
WHERE {
  ?substance a <http://semanticscience.org/resource/CHEMINF_000059>;
    <http://semanticscience.org/resource/SIO_000300> ?inchiKey .
} LIMIT 5
"""
queryEBI(query)
```

Out[10]:

```
{'head': {'vars': ['inchiKey']},
 'results': {'bindings': [{'inchiKey': {'type': 'literal',
    'value': 'YAPYGARUFVRIHX-ACCUITESSA-N'}},
  {'inchiKey': {'type': 'literal', 'value': 'IDDOHJSCCLKCAC-UHFFFAOYSA-N'}},
  {'inchiKey': {'type': 'literal', 'value': 'JJVMPIBZWGJKIV-GLGDLAMFSA-N'}},
  {'inchiKey': {'type': 'literal', 'value': 'MZITZBZTMCXZQS-FVNDCLGZSA-N'}},
  {'inchiKey': {'type': 'literal', 'value': 'KPNTUIZGUAJMMK-DWKDZXJCSA-N'}}]}}
```

In [11]:

```
# Show 5 InChI Keys stored in DrugBank
query = """
    SELECT ?inchi
    WHERE {?s a <http://bio2rdf.org/drugbank_vocabulary:InChIKey>;
           <http://bio2rdf.org/drugbank_vocabulary:value> ?inchi.
    } LIMIT 5
    """
queryDrugBank(query)
```

Out[11]:

```
{'head': {'link': [], 'vars': ['inchi']},
 'results': {'bindings': [{'inchi': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'InChIKey=PMATZTZNYRCHOR-IMVLJIQENA-N'}},
 {'inchi': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'InChIKey=SFKQVVDKFKYTNA-YVGXZPIDNA-N'}},
 {'inchi': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'InChIKey=DEQANNDTNATYII-RRCPSPWKPSA-N'}},
 {'inchi': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'InChIKey=NGVDGCNIFYWLIFO-UHFFFAOYSA-N'}},
 {'inchi': {'datatype': 'http://www.w3.org/2001/XMLSchema#string',
    'type': 'typed-literal',
    'value': 'InChIKey=SEKGMJVHSSBBHRD-WZHZPDAFSA-M'}}]},
'distinct': False,
'ordered': True}}
```

As can be seen from the results, additional string processing would be required to compare the InChI Key string from DrugBank with the other datasets. This would be a costly, but not impossible operation. However, one that would be best suited to centralising all the data first and then performing the SPARQL queries.

4. Conclusions

In this notebook, we have conducted a very simple analysis, counting the number of compounds in different pharmacology datasets. The purpose of publishing this as a notebook (c.f. open science) is to enable the repetition of this computation, since the datasets are continuously evolving, and to share our methods in a way that allows others to reproduce our results whilst also reading our narrative. However, this has raised the following questions:

1. How should you present the computational environment used?
2. How should you present the results of the original computation against those of the live results when the content of datasets is evolving?

For the first question, you would ideally use the computational environment to report about itself. However, this is not always straightforward. Libraries don't necessarily report their versions, and SPARQL endpoints all have different approaches (or perhaps none) for providing metadata about their datasets.

The second question focuses on recognising that sources change over time, often gaining more data, although as shown with the ChEBI data sometimes containing errors too. The publication of this notebook captures a point in time, or indeed three. The first is the time when the query was run over the PubChem database (15

March 2018), the second when the SPARQL queries were executed for the peer review version of this notebook (8 June 2018), and the third when the camera ready version was prepared (24 July 2018). This allowed us to compare the results at the different timepoints. In this notebook we chose to copy the results of the computation for publication into the text in order to allow others reusing the notebook to be able to compare their answers. We hope that this proves to be a useful approach. We note that the reuse of this notebook relies on the continued availability of the SPARQL endpoints. Newer computational notebook environments are being developed that overcome some of the limitations of Jupyter, e.g. Wrattler ([Petricek et al; 2018](#)), although these still need to mature and gain wider uptake by the community.

We have deliberately not used any Jupyter extensions in this notebook, believing that this will aid reproducibility into the future. However, this complicated the generating of the paper in this notebook, as there is no inbuilt support for academic references. If the use of computational notebooks is to grow, then referencing needs to become a core feature. In the future, it would be interesting to investigate the additional effort required in making the notebook itself a semantic resource. There have already been several discussions on embedding metadata in Jupyter Notebooks. However, there is no UI support for entering this metadata, no standard approach for representing authorship metadata, and no support for creating this as semantic annotations.

Acknowledgements

This work was conducted through the ELIXIR-EXCELERATE project that is funded by the European Commission within the Research Infrastructures programme of Horizon 2020, grant number 676559.

The author would also like to thank the reviewers – [Idafen Santana-Perez](http://orcid.org/0000-0001-8296-8629) (<http://orcid.org/0000-0001-8296-8629>), [Michel Dumontier](https://orcid.org/0000-0003-4727-9435) (<https://orcid.org/0000-0003-4727-9435>), and [Oscar Corcho](https://orcid.org/0000-0002-9260-0753) (<https://orcid.org/0000-0002-9260-0753>) – for their helpful [comments](https://semsci.github.io/SemSci2018/openReview/Gray.txt) (<https://semsci.github.io/SemSci2018/openReview/Gray.txt>) that have been used to improve the published version of this notebook.

References

Baker, M. 1,500 scientists lift the lid on reproducibility. *Nature* 533, 452–454 (2016).
DOI: [10.1038/533452a](https://doi.org/10.1038/533452a) (<https://doi.org/10.1038/533452a>)

Callahan A., Cruz-Toledo J., Ansell P., Dumontier M. Bio2RDF Release 2: Improved Coverage, Interoperability and Provenance of Life Science Linked Data. *ESWC 2013*: 200-212. (2013)
DOI: [10.1007/978-3-642-38288-8_14](https://doi.org/10.1007/978-3-642-38288-8_14) (https://doi.org/10.1007/978-3-642-38288-8_14)

Farr, W. M., Stevenson, S., Coleman Miller, M., Mandel, I., Farr, B., and Vecchio, A. Distinguishing spin-aligned and isotropic black hole populations with gravitational waves. *Nature* 548, 426 (2017).
DOI: [10.0.4.14/nature23453](https://doi.org/10.0.4.14/nature23453) (<https://doi.org/10.0.4.14/nature23453>)

Gaulton A, Hersey A, Nowotka M, Bento AP, Chambers J, Mendez D, Mutowo P, Atkinson F, Bellis LJ, Cibrián-Uhalte E, Davies M, Dedman N, Karlsson A, Magariños MP, Overington JP, Papadatos G, Smit I, Leach AR. The ChEMBL database in 2017. *Nucleic Acids Research*, 45(D1) D945-D954. (2017)
DOI: [10.1093/nar/gkw1074](https://doi.org/10.1093/nar/gkw1074) (<https://doi.org/10.1093/nar/gkw1074>)

Grüning BA, Rasche E, Rebolledo-Jaramillo B, Eberhard C, Houwaart T, Chilton J, Coraor N, Backofen R, Taylor J, and Nekrutenko A. Jupyter and Galaxy: Easing entry barriers into complex data analyses for biomedical researchers. *PLoS Comput Biol* 13(5): e1005425 (2017).
DOI: [10.1371/journal.pcbi.1005425](https://doi.org/10.1371/journal.pcbi.1005425) (<https://doi.org/10.1371/journal.pcbi.1005425>)

Harding SD, Sharman JL, Faccenda E, Southan C, Pawson AJ, Ireland S, Gray AJG, Bruce L, Alexander SPH, Anderton S, Bryant C, Davenport AP, Doerig C, Fabbro D, Levi-Schaffer F, Spedding M, Davies JA; NC-IUPHAR. The IUPHAR/BPS Guide to PHARMACOLOGY in 2018: updates and expansion to encompass the new guide to IMMUNOPHARMACOLOGY. *Nucleic Acids Research*, 46 (Issue D1): D1091-D1106 (2018). DOI: [10.1093/nar/gkx1121](https://doi.org/10.1093/nar/gkx1121) (<https://doi.org/10.1093/nar/gkx1121>).

Hastings J, Owen G, Dekker A, Ennis M, Kale N, Muthukrishnan V, Turner S, Swainston N, Mendes P, Steinbeck C. ChEBI in 2016: Improved services and an expanding collection of metabolites. *Nucleic Acids Research* (2016). DOI: [10.1093/nar/gkv1031](https://doi.org/10.1093/nar/gkv1031) (<https://doi.org/10.1093/nar/gkv1031>).

Jupp S, Malone J, Bolleman J, Brandizi M., Davies M., Garcia L., Gaulton A., Gehant S., Laibe C., Redaschi N., Wimalaratne S.M., Martin M., Le Novère N., Parkinson H., Birney E. and Jenkinson A.M. The EBI RDF Platform: Linked Open Data for the Life Sciences Bioinformatics 30 1338-1339. 2014. DOI: [10.1093/bioinformatics/btt765](https://doi.org/10.1093/bioinformatics/btt765) (<https://doi.org/10.1093/bioinformatics/btt765>).

Kluyver T., Ragan-Kelley B., Pérez F., Granger B.E., Bussonnier M., Frederic J., Kelley K., Hamrick J.B., Grout J., Corlay S. and Ivanov P. Jupyter Notebooks-a publishing format for reproducible computational workflows. In *ELPUB*, 87-90 (2016). DOI: [10.3233/978-1-61499-649-1-87](https://doi.org/10.3233/978-1-61499-649-1-87) (<https://doi.org/10.3233/978-1-61499-649-1-87>).

Knuth, D. E. Literate Programming. *The Computer Journal* 27, 97–111 (1984). DOI: [10.1093/comjnl/27.2.97](https://doi.org/10.1093/comjnl/27.2.97) (<https://doi.org/10.1093/comjnl/27.2.97>).

Medvedev D, Lemson G, and Rippin M. SciServer Compute: Bringing Analysis Close to the Data. In *Proceedings of the 28th International Conference on Scientific and Statistical Database Management (SSDBM '16)*, Article 27, 4 pages (2016). DOI: [10.1145/2949689.2949700](https://doi.org/10.1145/2949689.2949700) (<https://doi.org/10.1145/2949689.2949700>).

Pearce, J. M. Open Source Research in Sustainability. *Sustainability: The Journal of Record*, 5, 238–243 (2012). DOI: [10.1089/SUS.2012.9944](https://doi.org/10.1089/SUS.2012.9944) (<https://doi.org/10.1089/SUS.2012.9944>).

Petricek, T., Geddes, J. and Sutton, C. Wrattler: Reproducible, live and polyglot notebooks. in 10th USENIX Workshop on The Theory and Practice of Provenance (TaPP 2018) (2018). <https://www.usenix.org/conference/tapp2018/presentation/petricek> (<https://www.usenix.org/conference/tapp2018/presentation/petricek>).

Piccolo S.R. and Frampton M. B. Tools and techniques for computational reproducibility. *Gigascience* 5, 30 (2016). DOI: [10.1186/s13742-016-0135-4](https://doi.org/10.1186/s13742-016-0135-4) (<https://doi.org/10.1186/s13742-016-0135-4>).

Turner C., and Gill, I. Developing a Data Management Platform for the Ocean Science Community. *Marine Technology Society Journal* 52, no. 3: 28-32 (2018). DOI: [10.4031/MTSJ.52.3.8](https://doi.org/10.4031/MTSJ.52.3.8) (<https://doi.org/10.4031/MTSJ.52.3.8>).

Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR, Sajed T, Johnson D, Li C, Sayeeda Z, Assempour N, Iynkkaran I, Liu Y, Maciejewski A, Gale N, Wilson A, Chin L, Cummings R, Le D, Pon A, Knox C, Wilson M. DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic Acids Research*, 2017. DOI: [10.1093/nar/gkx1037](https://doi.org/10.1093/nar/gkx1037) (<https://doi.org/10.1093/nar/gkx1037>).

Version History

- 2018-07-31 Camera ready version ([e7d1dce](https://github.com/AlasdairGray/SemSci2018/commit/e7d1dcee4fba3b64038761dadcf4c173cac8f9e4)) (<https://github.com/AlasdairGray/SemSci2018/commit/e7d1dcee4fba3b64038761dadcf4c173cac8f9e4>)

- 2018-06-08 Version submitted to Semantic Science 2018 workshop ([8130dda](https://github.com/AlasdairGray/SemSci2018/commit/8130ddaa41238d4fc65a2a198fde2953689db3e0)
(<https://github.com/AlasdairGray/SemSci2018/commit/8130ddaa41238d4fc65a2a198fde2953689db3e0>))