

Try to fetch information from statistics.gov.scot using SPARQL

statistics.gov.scot (<http://statistics.gov.scot>) makes many useful datasets available as linked open data.

Let's execute a [SPARQL](https://www.w3.org/TR/sparql11-query/) (<https://www.w3.org/TR/sparql11-query/>) query to retrieve information about all of those datasets...

In [3]:

```
1 ; Add code libraries
2
3 %classpath add mvn org.clojure data.csv 1.0.0
4 (require '[clojure.data.csv :as csv])
5
6 %classpath add mvn clj-http clj-http 3.10.1
7 (require '[clj-http.client :as http])
8
9 (require '[clojure.string :as str])
10
11 (import 'java.net.URLEncoder)
12 (import 'java.time.LocalDate)
13 (import 'com.twosigma.beakerx.chart.xychart.TimePlot
14         'com.twosigma.beakerx.chart.xychart.plotitem.Line)
```

Out[3]:

```
class com.twosigma.beakerx.chart.xychart.plotitem.Line
```

In [4]:

```
1 ; Define convenience functions
2
3 ; Convert the CSV structure to a list-of-maps structure.
4 (defn to-maps [csv-data]
5   (map zipmap (->> (first csv-data)
6                   (map keyword)
7                   repeat)
8         (rest csv-data)))
9
10 ; Ask statistic.gov.scot to execute the given SPARQL query
11 ; and return its result as a list-of-maps.
12 (defn exec-query [sparql]
13   (->> (http/post "http://statistics.gov.scot/sparql"
14                   {:body (str "query=" (URLEncoder/encode sparql))
15                    :headers {"Accept" "text/csv"
16                              "Content-Type" "application/x-www-form-urlencoded"
17                              :debug false})
18         :body
19         csv/read-csv
20         to-maps))
```

Out[4]:

```
#'beaker_clojure_shell_d969f8d8-804b-4a28-a997-e83a4343705f/exec-q
uery
```

In [5]:

```

1 ; Query for information about all the datasets
2
3 (def sparql "
4
5 PREFIX dcterms: <http://purl.org/dc/terms/>
6 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
7 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
8 SELECT ?dataset ?name ?publisherLabel
9
10 WHERE {
11     ?dataset rdf:type <http://publishmydata.com/def/dataset#Dataset>.
12     ?dataset rdfs:label ?name.
13     ?dataset dcterms:publisher ?publisher.
14
15     ?publisher rdfs:label ?publisherLabel.
16 }
17 ")
18
19 (def all-datasets (exec-query sparql))
20
21 (str (count all-datasets) " datasets")

```

Out[5]:

290 datasets

My [previous report](https://stir.sharepoint.com/:w:/r/sites/DataCommonsScotland/_layouts/15/Doc.aspx?sourcedoc=%7B5DD2D579-68D4-4C21-91CD-CCF6BE8BC276%7D&file=2020-04-16-progress-report-ash.docx&action=default&mobileredirect=true) (https://stir.sharepoint.com/:w:/r/sites/DataCommonsScotland/_layouts/15/Doc.aspx?sourcedoc=%7B5DD2D579-68D4-4C21-91CD-CCF6BE8BC276%7D&file=2020-04-16-progress-report-ash.docx&action=default&mobileredirect=true) used the Household Waste dataset from SEPA and the Population dataset from NRS.

Let's find those two datasets in the statistics.gov.scot information...

In [6]:

```

1 ; Filter to find SEPA's Household Waste and NRS' Population datasets
2
3 (->> all-datasets
4     (filter #(or
5         (and (= "SEPA" (:publisherLabel %))
6             (str/includes? (:name %) "Household Waste"))
7         (and (= "National Records of Scotland" (:publisherLabel %))
8             (str/includes? (:name %) "Population")
9             (str/includes? (:name %) "Current"))))
10     (sort-by :name))

```

index	:dataset	
0	http://statistics.gov.scot/data/household-waste	C
1	http://statistics.gov.scot/data/population-estimates-current-geographic-boundaries	P

Within [this](https://github.com/ash-mcc/dcs/blob/master/linked-data-experiment/plot-info-graphics.ipynb) (<https://github.com/ash-mcc/dcs/blob/master/linked-data-experiment/plot-info-graphics.ipynb>) previous lab book, we queried its triplestore for the waste tonnage generated per council citizen per year by executing the [Datalog](https://en.wikipedia.org/wiki/Datalog) (<https://en.wikipedia.org/wiki/Datalog>) query:

```

    [:find ?c ?y ?t2 :where
l, year & tonnage
    [?e1 :waste-type/name "Subtotal"]
    [?e2 :waste-process/name "Generated"]
    [?e3 :waste-tonnes-cytp/waste-type ?e1] ; fixed at Subt
otal
    [?e3 :waste-tonnes-cytp/waste-process ?e2] ; fixed at Gene
rated
    [?e3 :waste-tonnes-cytp/tonnes ?t1] ; tonnage for a
council for a year
    [?e3 :waste-tonnes-cytp/council ?e4]
    [?e3 :waste-tonnes-cytp/year ?e5]
    [?e4 :council/name ?c] ; used to join
waste- & population- facts
    [?e5 :year/value ?y] ; used to join
waste- & population- facts
    [?e6 :population-cy/council ?e4]
    [?e6 :population-cy/year ?e5]
    [?e6 :population-cy/population ?p] ; population fo
r a council for a year
    [(/ ?t1 ?p) ?t2]] ; calculate the
tonnage per citizen

```

In theory, we ought to be able to achieve the same result by executing an equivalent SPARQL query against the datasets held in statistics.gov.scot.

Let's try building up such a query step-by-step...

In [7]:

```

1  ; Query for the waste tonnage generated per council per year
2
3  (def sparql "
4
5  PREFIX qb: <http://purl.org/linked-data/cube#>
6  PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
7  PREFIX pdmx: <http://purl.org/linked-data/sdmx/2009/dimension#>
8  PREFIX sdmx: <http://statistics.gov.scot/def/dimension/>
9  PREFIX snum: <http://statistics.gov.scot/def/measure-properties/>
10
11  SELECT ?c ?y ?t
12  WHERE {
13
14      ?tonnage qb:dataSet <http://statistics.gov.scot/data/household-waste> .
15      ?tonnage pdmx:refArea ?area .
16      ?tonnage pdmx:refPeriod ?period .
17      ?tonnage sdmx:wasteCategory ?wastecategory .
18      ?tonnage sdmx:wasteManagement ?wastemanagement .
19      ?tonnage snum:count ?t .
20      ?area rdfs:label ?c .
21      ?period rdfs:label ?y .
22
23      ?wastecategory rdfs:label \"Total Waste\" .
24      ?wastemanagement rdfs:label \"Waste Generated\" .
25  }
26  ")
27
28  (def tonnage-generated-per-council-per-year (exec-query sparql))
29
30  (sort-by (juxt :c :y) tonnage-generated-per-council-per-year)

```

index	:c	:y	:t
0	Aberdeen City	2011	97184
1	Aberdeen City	2012	97230
2	Aberdeen City	2013	94117
3	Aberdeen City	2014	96130
4	Aberdeen City	2015	95241
5	Aberdeen City	2016	96123
6	Aberdeen City	2017	87786
7	Aberdeen City	2018	85540
8	Aberdeenshire	2011	148052
9	Aberdeenshire	2012	134708
10	Aberdeenshire	2013	131809
11	Aberdeenshire	2014	131389
12	Aberdeenshire	2015	130249
13	Aberdeenshire	2016	131863
14	Aberdeenshire	2017	127632
15	Aberdeenshire	2018	120519
16	Angus	2011	62200
17	Angus	2012	60803

18	Angus	2013	54258
19	Angus	2014	58892
20	Angus	2015	57609
21	Angus	2016	58847
22	Angus	2017	56278
23	Angus	2018	54619
24	Argyll and Bute	2011	46301

In [8]:

```

1 ; Query for the population per council per year (attempt #1)
2
3 (def sparql "
4
5 PREFIX qb: <http://purl.org/linked-data/cube#>
6 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
7 PREFIX pdmx: <http://purl.org/linked-data/sdmx/2009/dimension#>
8 PREFIX sdmx: <http://statistics.gov.scot/def/dimension/>
9 PREFIX snum: <http://statistics.gov.scot/def/measure-properties/>
10 PREFIX uent: <http://statistics.data.gov.uk/def/statistical-entity#>
11
12 SELECT ?c ?y ?p
13 WHERE {
14
15     ?population qb:dataSet <http://statistics.gov.scot/data/population-esti
16     ?population pdmx:refArea ?area .
17     ?population pdmx:refPeriod ?period .
18     ?population sdmx:age ?age .
19     ?population sdmx:sex ?sex .
20     ?population snum:count ?p .
21
22     ?area rdfs:label ?c .
23     ?area uent:code ?y .
24     ?period rdfs:label ?periodLabel .
25
26     ?areaCode rdfs:label \"Council Areas\" .
27     ?age rdfs:label \"All\" .
28     ?sex rdfs:label \"All\" .
29
30     VALUES ?periodLabel {\"2011\" \"2012\" \"2013\" \"2014\" \"2015\" \"2016\"
31 }
32 ")

```

Out[8]:

```
#'beaker_clojure_shell_d969f8d8-804b-4a28-a997-e83a4343705f/sparql
```

But statistics.gov.scot fails (with a timeout error message) when the above SPARQL query is executed against it.

Bill Roberts (from Swirrl IT whose PublishMyData platform underlies the statistics.gov.scot site) helped me figure out the root cause of the failure: basically, that query results in too many *joins* to compute within the 30 seconds limit of their public SPARQL endpoint. Bill suggested replacing `label` s with 'fixed' facts and performing the `VALUES` filter outside of SPARQL.

So lets try again...

In [9]:

```

1 ; Query for the population per council per year (attempt #2)
2
3 (def sparql "
4
5 PREFIX qb: <http://purl.org/linked-data/cube#>
6 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
7 PREFIX pdmx: <http://purl.org/linked-data/sdmx/2009/dimension#>
8 PREFIX sdmx: <http://statistics.gov.scot/def/dimension/>
9 PREFIX snum: <http://statistics.gov.scot/def/measure-properties/>
10 PREFIX uent: <http://statistics.data.gov.uk/def/statistical-entity#>
11
12 SELECT ?c ?y ?p
13 WHERE {
14
15     ?population qb:dataSet <http://statistics.gov.scot/data/population-esti
16     ?population pdmx:refArea ?area .
17     ?population pdmx:refPeriod ?period .
18     ?population sdmx:age <http://statistics.gov.scot/def/concept/age/all> .
19     ?population sdmx:sex <http://statistics.gov.scot/def/concept/sex/all> .
20     ?population snum:count ?p .
21
22     ?area uent:code <http://statistics.gov.scot/id/statistical-entity/S12>
23     ?area rdfs:label ?c .
24     ?period rdfs:label ?y .
25 }
26 ")
27
28 (def population-per-council-per-year
29   (->> (exec-query sparql)
30         (filter #(>= (Integer/parseInt (:y %)) 2011))))
31
32 (sort-by (juxt :c :y) population-per-council-per-year)

```

index	:c	:y	:p
0	Aberdeen City	2011	222460
1	Aberdeen City	2012	224910
2	Aberdeen City	2013	227070
3	Aberdeen City	2014	228920
4	Aberdeen City	2015	230350
5	Aberdeen City	2016	229840
6	Aberdeen City	2017	228800
7	Aberdeen City	2018	227560
8	Aberdeenshire	2011	253650
9	Aberdeenshire	2012	255560
10	Aberdeenshire	2013	257770
11	Aberdeenshire	2014	260530
12	Aberdeenshire	2015	261960
13	Aberdeenshire	2016	262190
14	Aberdeenshire	2017	261800
15	Aberdeenshire	2018	261470
16	Angus	2011	116200

17	Angus	2012	116220
18	Angus	2013	116290
19	Angus	2014	116740
20	Angus	2015	116900
21	Angus	2016	116520
22	Angus	2017	116280
23	Angus	2018	116040
24	Argyll and Bute	2011	88930

Now let's use the above two SPARQL queries in a single query SPARQL query that will find the waste tonnage generated per council citizen per year (i.e. be equivalent to the Datalog query that was discussed earlier)...

In [10]:

```

1  ; Query for the waste tonnage generated per council citizen per year
2
3  (def sparql "
4
5  PREFIX qb: <http://purl.org/linked-data/cube#>
6  PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
7  PREFIX pdmx: <http://purl.org/linked-data/sdmx/2009/dimension#>
8  PREFIX sdmx: <http://statistics.gov.scot/def/dimension/>
9  PREFIX snum: <http://statistics.gov.scot/def/measure-properties/>
10 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
11
12 SELECT ?c ?y ?t2
13 WHERE {
14
15     ?tonnage qb:dataSet <http://statistics.gov.scot/data/household-waste> .
16     ?tonnage pdmx:refArea ?area .
17     ?tonnage pdmx:refPeriod ?period .
18     ?tonnage sdmx:wasteCategory ?wastecategory .
19     ?tonnage sdmx:wasteManagement ?wastemanagement .
20     ?tonnage snum:count ?t1 .
21
22     ?wastecategory rdfs:label \"Total Waste\" .
23     ?wastemanagement rdfs:label \"Waste Generated\" .
24
25     ?population qb:dataSet <http://statistics.gov.scot/data/population-esti
26     ?population pdmx:refArea ?area .
27     ?population pdmx:refPeriod ?period .
28     ?population sdmx:age <http://statistics.gov.scot/def/concept/age/all> .
29     ?population sdmx:sex <http://statistics.gov.scot/def/concept/sex/all> .
30     ?population snum:count ?p .
31
32     ?area rdfs:label ?c .
33     ?period rdfs:label ?y .
34     BIND((xsd:integer(?t1)/xsd:integer(?p)) AS ?t2) .
35 }
36 " )
37
38 (def tonnage-generated-per-council-citizen-per-year (exec-query sparql))
39
40 (sort-by (juxt :c :y) tonnage-generated-per-council-citizen-per-year)

```

index	:c	:y	:t2
0	Aberdeen City	2011	0.436860559201654229973928
1	Aberdeen City	2012	0.432306255835667600373483
2	Aberdeen City	2013	0.414484520192011274056458
3	Aberdeen City	2014	0.419928359252140485759217
4	Aberdeen City	2015	0.413462122856522682873888
5	Aberdeen City	2016	0.418217020536025060911939
6	Aberdeen City	2017	0.38368006993006993006993
7	Aberdeen City	2018	0.375900861311302513622781
8	Aberdeenshire	2011	0.583686181746501084171102
9	Aberdeenshire	2012	0.527109093754891219283143
10	Aberdeenshire	2013	0.511343445707413585754743

11	Aberdeenshire	2014	0.5043142824242889494492
12	Aberdeenshire	2015	0.497209497633226446785769
13	Aberdeenshire	2016	0.502929173500133490979824
14	Aberdeenshire	2017	0.487517188693659281894576
15	Aberdeenshire	2018	0.46092859601483917849084
16	Angus	2011	0.535283993115318416523236
17	Angus	2012	0.523171571158148339356393
18	Angus	2013	0.466574941955456187118411
19	Angus	2014	0.504471475072811375706699
20	Angus	2015	0.492805816937553464499572
21	Angus	2016	0.505037761757638173704085
22	Angus	2017	0.483986928104575163398693
23	Angus	2018	0.470691140985866942433644
24	Argyll and Bute	2011	0.520645451478691105363769

Now let's put that data into an infographic...

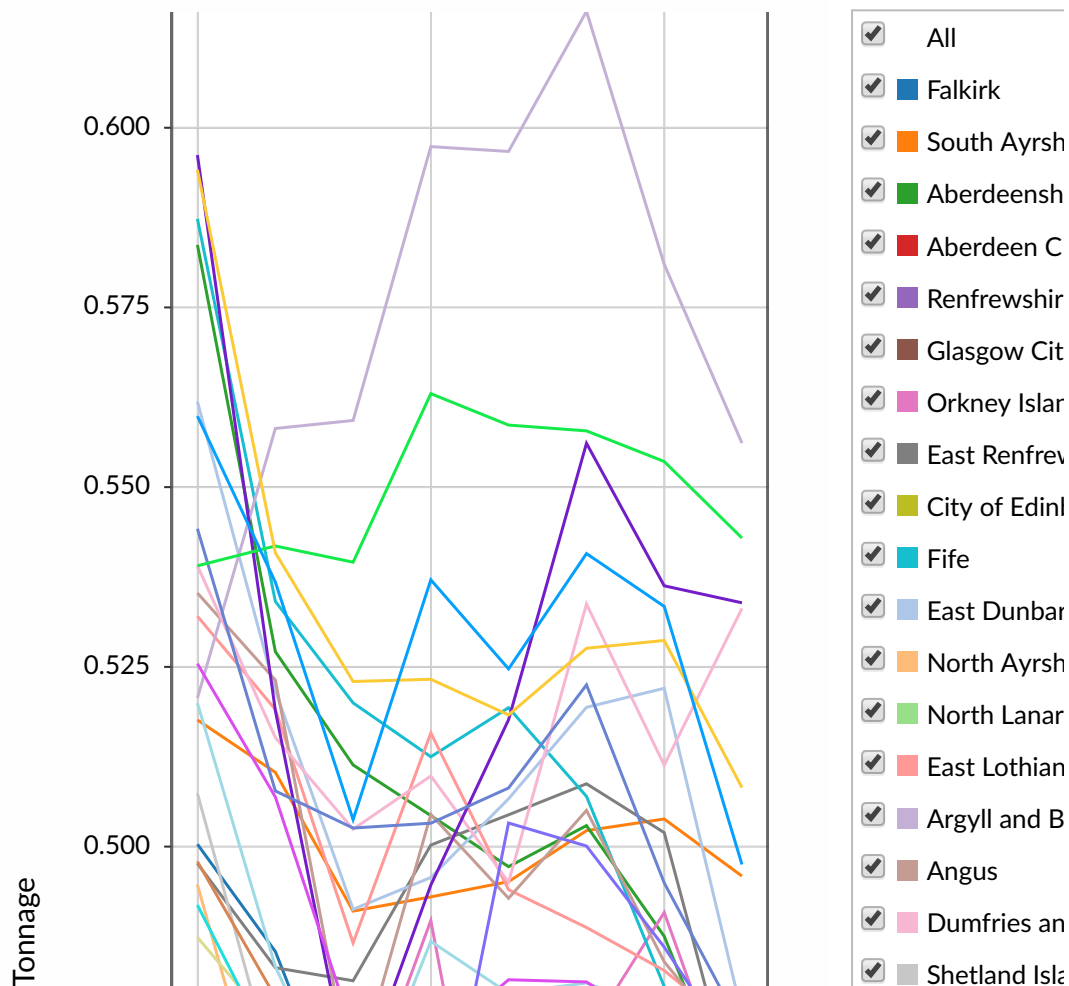
In [11]:

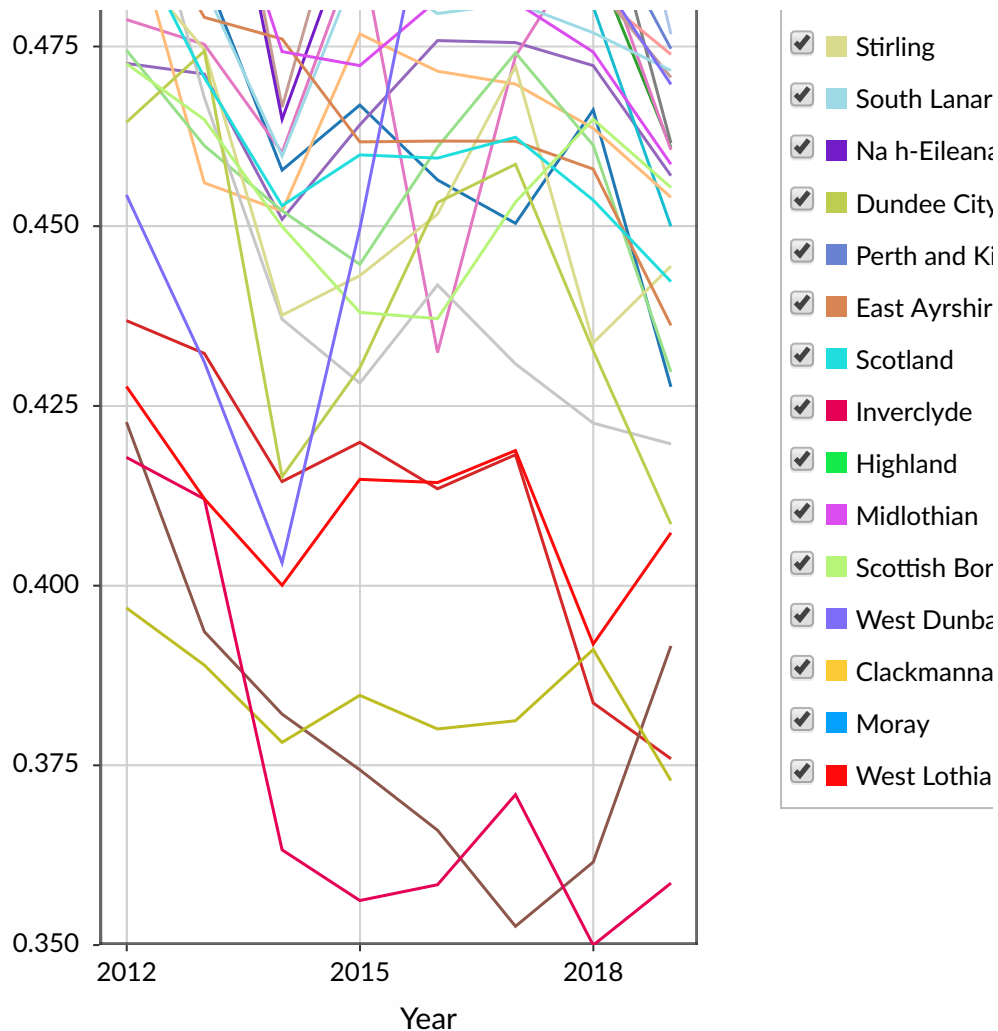
```

1 ; Plot an info-graphic the waste tonnage geneareted per council citizen per
2
3 (def lines
4   (->> tonnage-generated-per-council-citizen-per-year
5     (group-by :c)
6     (map (fn [[c coll1]]
7       (let [coll2 (sort-by :y coll1)]
8         { :label c
9           :x-series (->> coll2
10             (map :y)
11             (map (fn [yyyy] (LocalDate/parse (str yyyy "-12-31"
12           :y-series (->> coll2
13             (map :t2))))))))))
14
15 (def time-plot
16   (doto (TimePlot.)
17     (.setTitle "Waste generated per council-citizen per year")
18     (.setXLabel "Year")
19     (.setYLabel "Tonnage")))
20 (doseq [line lines]
21   (.add time-plot (doto (Line.)
22     (.setDisplayName (:label line))
23     (.setX (:x-series line))
24     (.setY (:y-series line)))))
25 time-plot

```

Waste generated per council-citizen per year





Conclusions

- I'm impressed at statistics.gov.scot's platform - it manages to curate, link and make accessible a sizeable number of datasets from several sources.
Its datasets include some (e.g. [Household Waste](https://statistics.gov.scot/data/household-waste) (<https://statistics.gov.scot/data/household-waste>) and [Population](https://statistics.gov.scot/data/population-estimates-current-geographic-boundaries) (<https://statistics.gov.scot/data/population-estimates-current-geographic-boundaries>)) that are core to providing data about waste in Scotland.
- Its SPARQL interface is powerful and has shown itself to be useful. However, some care must be taken to ensure that queries against it, are performant.
For example, we might have to split one query (requiring many *joins*) into several smaller queries (requiring less *joins*) that are needing together outside of SPARQL.
- For the general public, statistics.gov.scot is not "easy to use" - I have seen several references to this in Slack workspaces (OpenDataScotland and SODU).
Whereas, statisticians/coders may accept that it is reasonably simple given that it is needing together diverse datasets into a coherent knowledge.
(Certainly it has taken me over one work-day to understand enough about the structures within statistics.gov.scot, just to be able to construct the queries above.)
One of the objectives of the Data Commons Scotland project is to address this "easy of use" issue.