

統計データの LOD 化とデータ間の関係の表現

Presentation of Statistical Data and their Relationship as LOD

武田 英明^{*1*2}
Hideaki Takeda

加藤 文彦^{*3}
Fumihito Kato

小出 誠二^{*3}
Seiji Koide

松村 冬子^{*4}
Fuyuko Matsumura

大向 一輝^{*1*2}
Ikki Ohmukai

小林 巖生^{*5}
Iwao Kobayashi

岩山 真^{*6}
Makoto Iwayama

浅野 優^{*6}
Yu Asano

濱崎 雅弘^{*7}
Masahiro Hamasaki

^{*1} 国立情報学研究所 ^{*2} 総合研究大学院大学 ^{*3} 情報・システム研究機構 ^{*4} 青山学院大学
National Institute of Informatics Graduate University for Advanced Studies Research Organization of Information and Systems Aoyama Gakuin University

^{*5} Open Community Data Initiative ^{*6} 日立製作所 中央研究所 ^{*7} 産業技術総合研究所
Open Community Data Initiative Central Research Laboratory, Hitachi Ltd. National Institute of Advanced Industrial Science and Technology

In this paper, we show Linked Open Data of the statistical data based on RDF Data Cube Vocabulary. Since the statistical data is important for the government, it should be open and shared as a part of Open Data by the government. As a part of Open Data Portal by Ministry of Economy, Trade and Industry, Japan, we translate and represent some parts of Census of Manufactures with RDF by using RDF Data Cube Vocabulary.

1. はじめに

オープンガバメントは現在、世界的潮流になっている。オープンガバメントのための施策は多岐にわたるが、一つの重要な点は政府の持つデータを公開し、再利用可能にする政府によるオープンデータ化である。現在、各国が争うようにオープンデータを公開するサイトを立ち上げて、データ公開を進めている。

日本政府も遅まきながら、2012年7月に電子行政オープンデータ戦略[IT 戦略本部 2012]を制定して、国としてオープンデータに取り組むことを決定している。その戦略に基づき、現在、総務省、経済産業省、内閣官房が関連する会議を開催して、戦略を具体化する作業をしている。総務省はオープンデータ流通推進コンソーシアム¹を立ち上げ、実証実験や各種会議を通じて官民連携のオープンデータの推進を進めている。経産省は IT 融合フォーラム/公共データワーキンググループを開催しつつ、自省のデータの公開するサイトであるオープンデータ実証用サイト「Open DATA METI」(β版)²を構築している。

一方、世界のオープンデータにおいて Linked Open Data (LOD)は先進的なデータ公開技術として認知されている。米国と英国においてはトップページから Semantic Web あるいは Linked Data というメニューが用意され、実際 RDF 化されたデータを入手することができる。日本のオープンデータ戦略においても RDF 化を視野にいれている。

本研究においては、前述の Open DATA METI において公開した工業統計調査の結果の一部の RDF 化について、どのようなことを行って来たかを述べる。

2. Open DATA METI

Open DATA METI はオープンデータ実証用サイトであり、前

連絡先: 武田英明, 国立情報学研究所, 千代田区一ツ橋2-1-2, 03-4212-2543, takeda@nii.ac.jp

¹ <http://www.opendata.gr.jp/>

² <http://datameti.go.jp/>

述の通り、経済産業省のデータを公開するサイトとして作成したものである。

このサイトには二つの主な機能がある。一つはデータカタログであり、経産省のもつ公開データをリストアップして掲載することで、簡単に検索して、原データにアクセスできるようにするというものである。もう一つは、RDF 化されたデータに関して、SPARQL Endpoint を提供することで、LOD 化されたデータの多様な利用を可能とするものである。

データカタログとしては、白書と統計データを中心に、現在、197 件のデータセット、13,282 件のリソースが登録されている。

LOD としては、今回その一つである工業統計調査の一部のデータを試験的に変換して、利用可能としている。

3. 工業統計調査データの LOD 化の論点

国勢調査のような調査活動の結果は、個票と呼ばれる調査結果から個人が特定されるような情報を除いて統計値が計算され、表形式にまとめられ、xls や csv のようなフォーマットで公開される。ここで具体的な統計数値の意味は、それに付随する各種属性のセットで表現されたと考える。たとえば、図1に示すような、「平成 22 年工業統計表『工業地区編データ 経済産業省大臣官房調査統計グループ』(平成 24 年 4 月 27 日公表)」の「第1表 都道府県別、産業中分類別統計表」において、Excel セル I10 の意味は、「食料品製造業」という産業分類における全国合計の「従業者数」であるが、このとき表側および表頭にある「食料品製造業」や「従業者数」はもちろんのこと、正確には「全国計」、「実数(人)」に加えて、他の表も同様な構造を持つために、「平成 22 年工業統計表『工業地区編データ 経済産業省大臣官房調査統計グループ』(平成 24 年 4 月 27 日公表)」や、「第1表 都道府県別、産業中分類別統計表」もこのセル I10 の数値の属性として考慮されなければならない。

一方、LOD における RDF (Resource Description Framework) のモデルは RDF グラフと呼ばれるグラフ構造である。ここでグラフというのは、二つのノードと両者を結合するリンクを最小単位とする構造のことである。RDF グラフはラベル付き有向グラフであ

平成22年工業統計表「工業地区別」データ 経済産業省大臣官房調査統計グループ1(平成24年4月27日公表)
[GO TO INDEX]

第1表 都道府県別 産業中分類別統計表

都道府県 産業分類	面種	事業所数	従業員数		金額 (百万円)	構成比 (%)	従業員1人 当たり金額 (千円)	産業別 特化係数	現金給与 総額 (百万円)	有形固定資産額 (従業員30人以上)		
			実数 (人)	人口比率 (%)						年末現在高 (百万円)	資本設備率 (千円)	
00 全国計	00 製造業計	372834	224403	7663847	6.034	289107683	100.0	36662	32719540	69568727	12377	
00 全国計	09 食品製造業	372834	30292	122307	0.094	24114367	8.3	30296	30296	5667585	4777	
00 全国計	10 飲料・たばこ・飼料製造業	372834	4391	102045	0.060	9613348	3.3	66617	422264	1501857	31127	
00 全国計	11 繊維工業	372834	15902	296827	0.234	3789828	1.3	12514	778520	1062883	6917	
00 全国計	12 木材・木製品製造業(家具を除く)	372834	6456	96045	0.076	2134101	0.7	21857	312668	429742	11367	
00 全国計	13 家具・器具製造業	372834	6610	99053	0.078	1575390	0.5	15604	351256	363296	8023	
00 全国計	14 パルプ・紙・紙加工品製造業	372834	6695	189907	0.149	7110758	2.5	36849	791344	3043217	23634	
00 全国計	15 印刷・同関連業	372834	13914	299038	0.235	6044642	2.1	19789	1182664	1665266	9207	
00 全国計	16 化学工業	372834	4742	344968	0.272	26212040	9.1	74823	1919273	7260091	23796	
00 全国計	17 石油製品・石炭製品製造業	372834	953	25387	0.020	14891705	5.2	487165	169144	2046332	119690	
00 全国計	18 プラスチック・ゴム・皮革製造業(靴を除く)	372834	14085	420179	0.331	10902553	3.8	25512	1584180	3008829	10652	
00 全国計	19 ゴム製品製造業	372834	2782	117176	0.092	3028976	1.0	25530	495365	780232	8679	
00 全国計	20 ぬいぐるみ・玩具・模型製造業	372834	1688	24761	0.019	361569	0.1	14338	69859	31139	2339	
00 全国計	21 窯業・土石製品製造業	372834	11055	249439	0.196	7101297	2.5	27993	1045271	2725592	19796	
00 全国計	22 鉄鋼業	372834	4486	219983	0.173	18148293	6.3	82146	1211008	8691488	37582	
00 全国計	23 非鉄金属製造業	372834	2909	143637	0.113	8911397	3.1	51551	712725	645124	20671	
00 全国計	24 金属製品製造業	372834	28974	578599	0.456	12292040	4.3	20828	2263374	3171679	10171	
00 全国計	25 はん用機械器具製造業	372834	7714	324636	0.256	10099831	3.5	30752	1643362	2568056	10038	
00 全国計	26 生産用機械器具製造業	372834	20118	543070	0.428	13645906	4.7	24926	2504766	3363739	9903	
00 全国計	27 業務用機械器具製造業	372834	4568	211834	0.167	6872908	2.4	32002	967737	1368773	8001	
00 全国計	28 電子部品・デバイス・電子回路製造業	372834	4907	452731	0.356	16633305	5.8	36489	2189256	5612339	13591	
00 全国計	29 電気機械器具製造業	372834	9673	485979	0.361	15119685	5.2	30674	2216163	2708903	6918	
00 全国計	30 情報通信機械器具製造業	372834	1984	212466	0.167	12584896	4.4	58785	1119001	991489	5059	
00 全国計	31 輸送用機械器具製造業	372834	11110	948824	0.747	54213562	18.8	57148	5178397	9996106	11757	
00 全国計	32 その他の製造業	372834	9415	156486	0.123	3907287	1.2	22711	568572	608993	7975	
01 北海道	00 製造業計	78459	5931	173973	3.151	5952064	100.0	32777	576683	1305163	11281	
01 北海道	09 食品製造業	78459	2065	82420	1.493	1084710	31.7	22942	3796	200846	386138	6407

図1 工業統計調査の結果の表の例

る。すなわちリンクにも名前が付与されており、そのリンクには方向があるが双方向ではない。RDF のモデルはこのように単純なものであるため、表現方法として柔軟であり、根本的には表形式をこのようなグラフ構造に変換するのになんの問題もない。ただし、知識表現方法として原理的であればあるほど、どんな知識でも表現することはできてもその効率は悪くなる(表現に多くの記述を必要とする)というのが一般的であり、表形式から RDF グラフへの変換時の欠点として、表現効率の悪さが指摘できる。しかし、本来表形式に馴染まないデータで、表形式ではスペースになって0値ばかりが挿入されるようなデータでは、逆にグラフ構造のほうが表現効率がよくなる。

なぜ表形式になっているものをわざわざ RDF グラフにするのだろうか? その利点はなにか? 実はそれは RDF の利点を生かして RDF としての利用を考えない限り、表になっているものをわざわざ RDF に直す必要性はない。一般に RDF の利点として次のような点があげられる。

- a) URI を用いてグローバルに唯一な名前を定義し、その意味するところを定義できる。
- b) 背後に厳密なモデル論、意味論があるために、誤解の余地がない。またそのため、機械の助けをかりて様々な自動処理が可能となる。
- c) 様々なデータの関係を、組織を超えてリンクとして定義できるため、組織を超えたデータの利用が可能となる。たとえば、市町村の様々なデータがすでに LOD として提供されていた場合、企業など各組織の所有するデータを市町村データとリンクすることで、市町村データを経由して政府統計データと民間のデータなど互いに当初無関係であったデータ間の関係も抽出することが可能となる。

4. RDF Data Cube Vocabulary による統計データの LOD 化

統計データの LOD 化については、世界中で関心が高く、W3C の eGovernment Activity 中の Government Linked Data (GLD) Working Group で盛んに議論が行われている。その議論の中から、The RDF Data Cube Vocabulary[Cyganik,2013] が Working Draft として提案されている。今回は基本的にこの語彙に従ってスキーマを設計した。

4.1 RDF Data Cube Vocabulary の概要

この仕様では、ISO 標準である統計データの交換規約 SDMX の考え方を取り入れ、それを RDF フォーマットと LOD で扱うための方法を提案している。表を多次元空間でとらえるデータキューブという考え方、コードリスト、データフローなどという用語は、SDMX からきているものである。

「RDF データキューブ語彙」仕様書では、RDF とリンクデータの手法を用いて統計データを公開可能にすることについて次のような利点を挙げている。

- (1) 個々の観測値や観測値のグループが、(ウェブ)アドレス可能になる。それにより公開者と第三者がこのデータを注釈づけし(annotate)、リンク付けすることが可能となる。たとえば、ある報告書が詳細な出典のトレースバックを考慮した特定の図を参照することが可能となる。
- (2) データをデータセット横断的に、あるいは統計セットと非統計セットをフレキシブルに組み合わせることが可能になる(たとえば、宗教的寛容さに関連した国民的指標の高い値の国勢調査の領域で、すべての宗教的學校を発見するなど)。統計データはリンクデータのより広範なウェブの不可欠な一部となる。
- (3) リンクデータとして公開することで、現在静的なファイルのみを提供しているような公開者には、フレキシブルな、かつ非プロプライエタリな機械可読可能な公開の手段を提供することになり、プログラムからアクセス可能なすぐに使えるウェブ API をサポートすることになる。
- (4) 標準化されたツールやコンポーネントの再利用が可能となる。

データキューブでは、統計データセットは何らかの論理空間中の点とされる観測値の集まりから構成されると考える。一つのキューブは次元、属性、測度の集まりとして定義される。これらの各要素はデータキューブのコンポーネントと呼ばれる。

- (1) 次元コンポーネントは観測値を同定するものである。次元コンポーネントの値の集合は一個の観測を同定する。たとえば一つの観測値には観測された時間や観測がカバーする地理学上の領域が含まれる。
- (2) 測度コンポーネントは計測された値であり観察された現象を表現する。

- (3) 属性コンポーネントは観測された値を限定し、解釈することを可能にする。それは測度の単位やスケーリングファクタを指定することを可能にし、どんなスケーリングファクタや観測値の状態(推測値あるいは暫定値)のようなメタデータも指定することもできる。この語彙においては、

4.2 LOD 化の例

特定の領域における観測値の次元コンポーネント、測度コンポーネント、属性コンポーネントは、各々 qb:DimensionProperty、qb:MeasureProperty、qb:AttributeProperty のインスタンスであるプロパティを定義して利用する。例えば、あるノードが観測値を表現するものであり、次元コンポーネントとして産業分類を持ち、その次元の値が産業中分類「食料品製造業」であるとき、このノードからこの次元へのリンクをつけるために ktsh:refSangyoChuBunrui という名前のプロパティを次のように定義する。

```
ktsh: refSangyoChuBunrui a qb: DimensionProperty ;
rdfs: label "日本標準産業分類(中分類)"@j a ;
rdfs: range jsic: JsicConcept .
```

ここから ktsh:refSangyoChuBunrui は次元コンポーネントのためのプロパティであり、それは観測値の次元として jsic:JsicConcept のインスタンスを持つということがわかる。産業中分類「食料品製造業」を <http://datameti.go.jp/scheme/jsic/2007/C09> (jsic:C09) という URI を持つノードとしたとき、これは jsic:JsicConcept に型付けされる。

同様に、次元コンポーネントとして都道府県というものがあり、その次元の値として「全国計」を持つとき、ブランクノードからこの次元へのリンクをつけるために、ktsh-sac:refPrefecture という名前の次元コンポーネント・プロパティを次のように定義する。

```
sac: refPrefecture a qb: DimensionProperty ;
rdfs: label "reference area (prefecture)"@en ;
rdfs: label "都道府県"@j a ;
rdfs: subPropertyOf sdmx-dimension: refArea ;
rdfs: range sac: Prefecture ;
qb: concept sdmx-concept: refArea .
```

ここで sac:refPrefecture は sdmx-dimension:refArea のサブプロパティであることが述べてある。SDMX 標準には「内容指向のガイドライン (COG)」が含まれる。それは統計上の共通の概念とコードリストを定義しているが、RDF データキューブにおいても、この SDMX の概念を再利用可能とすることを目的に、以下のものが定義される。

- sdmx-concept: COG 定義の各概念に対する SKOS 概念
- sdmx-code: COG 定義の各コードリストに対する SKOS 概念とそのスキーマ
- sdmx-dimension: 次元として用いられる各 COG 概念に相当するコンポーネント・プロパティ
- sdmx-attribute: 属性として用いられる各 COG 概念に相当するコンポーネント・プロパティ
- sdmx-measure: 測度として用いられる各 COG 概念に相当するコンポーネント・プロパティ

そこで観測値から数値(リテラル)へのリンクのプロパティとして ktsh:numberOfEmployee を qb:MeasureProperty の実現として次のように定義する。

```
ktsh: numberOfEmployees a qb: MeasureProperty ;
rdfs: label "従業員数(人)"@j a ;
rdfs: subPropertyOf sdmx-measure: obsValue ;
sdmx-attribute: unitMeasure
ktsh: UnitOfPerson ;
```

```
rdfs: range xsd: integer .
```

ただし単位はプロパティ ktsh:numberOfEmployees に定義してあることに注意されたい。もし「従業員数(人口比率)」のプロパティがほしければ、単位を変えたプロパティを別途用意する必要がある。

以上のコンポーネントを用いて、一つの表が定義される。図2に一つの表の定義の例を示す。

4.3 コード体系の LOD 化

統計調査においては、調査対象の分類の体系は重要である。統計調査ではそれらの分類にコード(英数字からなる文字列)を割り振ることが多く、コード体系(あるいはコードリスト)と呼ばれる。コード体系は特定の調査のみに出現する系もあれば、特定の統計調査とは別に定義されることもある。工業統計においては、前者の例は工業地区のコード体系であり、後者としては日本標準産業分類や標準地域コードが例である。

Data Cube Vocabulary においてもコード体系は別に定義する。コード体系は基本的に SKOS(Simple Knowledge Organization System)を用いて表現する。すなわち、skos:broader/narrower で階層的な概念を関係づける。コード体系の要素(概念)は Data Cube Vocabulary の次元コンポーネントのプロパティの値となっている。コード体系の要素のあるクラスのインスタンス(4.2 節の例であれば jsic:JsicConcept)と宣言することで、それらの要素が次元コンポーネントのプロパティの値の候補になる。

5. 工業統計調査の LOD 化

上記の方針の元に工業統計調査の LOD 化を行った。

工業統計調査とは、統計法に基づき行政機関が実施する統計調査のうち、重要なものとして総務大臣が指定した基幹統計調査の一つで、国が我が国の製造業の実態を把握するためにを行っているものである。今回はこのうち、平成 22 年(2010 年)の調査結果を利用した。

5.1 対象データ

工業統計調査の結果は多数の表として公開されている。それぞれは調査票で集計したデータを様々な切り口で集計・集約したものである。大分類で品目編、産業編、用地・用水編、市区町村編、工業地区編、産業細分類別、企業統計編に別れ、その中に多数の表形式のデータがある。今回のその中の以下の4つの表を例題として取り上げた。これらの表の名称と次元と測度を示す¹。

- (1) 産業細分類別統計表 都道府県別産業細分類別統計表 (kougyo_h22-k8-data-j-1003.ttl)

次元: 都道府県、産業細分類

測度: 事業所数、従業者数、現金給与総額、原材料使用額等、製造品出荷額等、生産額付加価値額、有形固定資産投資総額

- (2) 市区町村編 市区町村別、産業中分類別統計表 (kougyo_h22-k6-data-j-2000.ttl)

次元: 市区町村、産業中分類

測度: 事業所数、従業者数、現金給与総額、原材料使用額等、製造品出荷額等、粗付加価値額、有形固定資産年末現在高

¹ 測度については簡略的に表記している。実際にはより細かく定義されている。

```
kougyo:k6-data-j-2000t a qb:DataStructureDefinition ;
  rdfs:label "工業統計表「市区町村編」データ(経済産業省大臣官房調査統計グループ) 2. 市区町村別、産業中分類別統計表(スキーマ)"@ja ;
  # dimension
  qb:component [qb:dimension ktsh:refMunicipality; qb:order 1] ;
  qb:component [qb:dimension ktsh:refSangyoChuBunrui; qb:order 2] ;
  qb:component [qb:dimension ktsh:refYear; qb:order 3] ;
  # measure
  qb:component [qb:measure ktsh:numberOfEstablishments] ;
  qb:component [qb:measure ktsh:numberOfEstablishments_withBetween30To299Employees] ;
  qb:component [qb:measure ktsh:numberOfEstablishments_with300OrMoreEmployees] ;
(中略)
# attributes
qb:component [qb:attribute sdmx:attribute:unitMeasure; qb:componentAttachment qb:DataSet] ;
```

図 2 : RDFS による表の定義の例

- (3) 産業編 都道府県別、東京特別区・政令指定都市別統計表 (2) 従業者 30 人以上の事業所に関する統計表 ②産業中分類別の在庫額、有形固定資産額及びリース契約による契約額及び支払額 (kougyo_h22-k3-data-j-3220.ttl)

次元: 都道府県、産業中分類

測度: 在庫額、有形固定資産額、リース契約による契約額及び支払額

- (4) 用地・用水編 第1部 事業所数、従業者数、製造品出荷額等、事業所敷地面積、建築面積及び延べ建築面積表 4. 工業地区別、産業中分類別統計表(kougyo_h22-k4-data-j-1400.ttl)

次元: 工業地区、産業中分類

測度: 事業所数、従業者数、製造品出荷額等、事業所敷地面積、事業所建築面積、事業所延べ建築面積

このうち、次元においては「都道府県」(1,3)、「工業地区」(4)、「市区町村」(2)は階層的な関係であり、「産業中分類」(2,3,4)と「産業細分類」(1)も同様である。これらの次元はそれぞれコード体系として表の表現とは別に用意される。

測度についてもいくつかの項目は共通である。例えば(1)と(2)は測度は同じである。(1)(2)と(4)でも共通の測度(製造品出荷額等)があるが、単位が異なる(前者が万円単位、後者が百万円単位)なので測度コンポーネントは別に定義する必要がある。

5.2 RDF 化とクエリ

上記の 4 つの表を RDFS を用いて定義して、データを RDF として表現した。スキーマの定義を含めて、全体で 2,827,017 トリプルとなった。このデータは Open Data METI サイトでダウンロードあるいは SPARQL Endpoint¹を通じてアクセス可能である。

このデータに対して行う SPARQL Query の例を図3に示す²。この例では異なる表のデータを一つのクエリの中で参照して、統合した解を返すようにしている。(1)の表から北海道における産業細分類別の従業員数のデータをとり、産業中分類に集約すると共に、それに対応する有形固定資産土地の値を(3)の表から集めている。

6. 考察

今回、工業統計調査を例に取り、Data Cube Vocabulary の方法に則って統計データの RDF 化を行った。Data Cube Vocabulary の方法では個別のデータそのものをノードと表現し、そのプロパティとして次元や測度を表現する。このため表にまたがって次元

を指定して検索したりといった柔軟性の高い利用方法が可能となる。例えばコード体系が共通なら、一つの統計に限らず複数の統計からのデータをつなげて使うことも可能になる。ただし、このためには次元や測度、コード体系を共通に利用できるようにきちんと定義する必要がある。

現実的な課題としては定義のためのコストが高いことおよびデータ量を大きくすることがあげられる。このため、柔軟性の高い利用(横断的利用など)が期待されるデータを中心に RDF 化を行うことが実際のな方策だと思われる。

参考文献

- [IT 戦略本部 2012] 高度情報通信ネットワーク社会推進戦略本部, 電子行政オープンデータ戦略, 2012. <http://www.kantei.go.jp/jp/singi/it2/denshigyousei.html>.
 [Cyganiak,2013] R. Cyganiak and D. Reynolds (Eds.), The RDF Data Cube Vocabulary, W3C Working Draft 12 March 2013, W3C, 2013. <http://www.w3.org/TR/2013/WD-vocab-data-cube-20130312/>

```
#北海道の産業中分類別、有形固定資産土地(百万円)と従業員数
PREFIX ktsh:<http://datameti.go.jp/scheme/kougyou-toukei-schema/>
PREFIX kougyo: <http://datameti.go.jp/lof/kougyou-toukei/>
PREFIX qb: <http://purl.org/linked-data/cube/#>
select distinct ?sanchu_label ?total_jugyoin ?landprice
where
{
  {
    select distinct ?sanchu (SUM(?jugyoin) AS ?total_jugyoin)
    where
    {
      ?cell1 qb:dataSet kougyo:h22-k8-data-j-1003 .
      ?cell1 ktsh:refSangyoSaiBunrui ?sansai .
      ?sansho skos:narrower ?sansai .
      ?sanchu skos:narrower ?sansho .
      ?cell1 ktsh:refPrefecture
      <http://datameti.go.jp/scheme/standard-area-code/C01> .
      ?cell1 ktsh:numberOfEmployees ?jugyoin .
    } Group by ?sanchu
  }
  ?cell2 qb:dataSet kougyo:h22-k3-data-j-3220 .
  ?cell2 ktsh:refSangyoChuBunrui ?sanchu .
  ?cell2 ktsh:refPrefecture
  <http://datameti.go.jp/scheme/standard-area-code/C01> .
  ?cell2
  ktsh:valueOfTangibleFixedAssets_purchase_lands_byMillionYear
  ?landprice .
  ?sanchu rdfs:label ?sanchu_label .
}
```

図 3 : SPARQL Query の例

¹ <http://datameti.go.jp/sparql>

² <http://bit.ly/ZMerUz>