

GlobalStat

Web Semantico - Primo elaborato

Boschi Francesco - 0000939879

26 marzo 2021

Indice

1	Introduzione	2
2	Obbiettivi	3
2.1	Presentazione del sito WEB	4
3	Metodologie e sorgenti dati	8
4	Architettura e tecnologie utilizzate	10
4.1	RDF, RDFS	10
4.2	OWL	13
5	Commenti finali	14

Capitolo 1

Introduzione

Nel 2010 l'Istituto Universitario Europeo [1] ha dato il via al progetto Global Governance Programme [2], il quale ha come obiettivo principale la nascita di una comunità di professori, studenti e ricercatori di alto livello che contribuiscano e siano di supporto alle future generazioni in ambito politico e non solo.

GlobalStat [3] nasce, grazie ad una prima collaborazione con l'istituto Fundação Francisco Manuel dos Santos [4] che è stato di fatto il primo finanziatore, come sottosezione del Global Governance Programme.

Dal momento della sua nascita ad oggi, GlobalStat ha stabilito e intensificato nuovi programmi di partnership, in quanto riconosciuto come risorsa fondamentale in ambito pubblico e accademico. Dal 2016 GlobalStat è infatti partner con l'Organizzazione per la Cooperazione e lo Sviluppo Economico (OCED) [5] con l'obiettivo di produrre nuove tecniche di gestione e visualizzazione dei dati necessari ad analizzare e favorire lo sviluppo economico. Anche lo stesso Servizio ricerca del Parlamento europeo [6] ha deciso di allacciare rapporti con GlobalStat per ottenere statistiche ed infografiche utili allo staff del Parlamento e a tutti i paesi dell'Unione Europea uniti ai loro partner globali.

Attualmente GlobalStat coopera con oltre 80 istituzioni internazionali, come Eurostat, l'Organizzazione per il Cibo e l'Agricoltura e l'Organizzazione Internazionale del Lavoro, con le quali vi è un continuo scambio di dati bidirezionale; le organizzazioni, oltre a sfruttare gli strumenti messi a disposizione da GlobalStats, forniscono dati e statistiche utili a mantenere l'ecosistema aggiornato.

Capitolo 2

Obbiettivi

Come afferma Gaby Umbach, direttore della fondazione GlobalStats [3], la statistica ricopre un ruolo fondamentale in numerosi ambiti della nostra vita, da quello politico a quello sociale.

Tale importanza si traduce in un aumento della domanda di dati statistici, che siano affidabili e pubblicamente disponibili, fattore spesso assente e fondamentale, specialmente nell'era della globalizzazione nella quale la velocità alla quale i dati vengono prodotti rende impossibile la loro gestione e analisi in maniera manuale.

Le fitte interconnessioni derivanti dalla globalizzazione portano ad avere forte impatto in ambito sociale, personale, culturale, politico, economico e ambientale, rendendo i dati che GlobalStat raccoglie vitali per poter gestire al meglio le risorse e le possibilità che ciascun paese e ciascun individuo possiedono.

GlobalStat si pone quindi come obbiettivo quello di soddisfare questa esigenza di informazioni trasparenti e disponibili pubblicamente, così che possano essere reperite e visualizzate in maniera semplice e intuitiva da chiunque ne abbia necessità.

Considerando l'obbiettivo alla base del progetto e il carattere multidimensionale della globalizzazione, GlobalStat presenta dati su una vasta gamma di argomenti, ad esempio:

- Dati demografici
- Dati relativi all'economia
- Dati relativi all'energia e alle risorse naturali
- Dati relativi all'ambiente e all'inquinamento
- Dati relativi alla produzione nei vari settori

- Dati relativi alla libertà, ai conflitti e ai pericoli
- Dati governativi
- Dati relativi alla salute e alle condizioni di vita
- Dati relativi a fattori etici e morali, come crimini, giustizie, diritti, educazione, uguaglianza e condizioni lavorative
- Dati relativi alle migrazioni
- Dati relativi allo sviluppo tecnologico
- Dati relativi a trend e mode

Emerge quindi come GlobalStat non miri solo a migliorare il lato strettamente economico di un paese, ma anche a informare sul modo in cui gli uomini vivono, le libertà di cui godono e le restrizioni che devono affrontare quotidianamente.

2.1 Presentazione del sito WEB

Il sito web si presenta con un'interfaccia tanto semplice quanto chiara e intuitiva, che vede sulla parte sinistra il menù suddiviso in categorie e sottocategorie con tutti gli indicatori disponibili, mentre in quella destra voci per ottenere informazioni sul progetto.

Una volta selezionato l'indicatore, di default verranno mostrate le statistiche in modalità mappa, dalla quale è possibile selezionare lo stato di interesse per visualizzare le informazioni specifiche, oltre che applicare i filtri messi a disposizione.

La visuale ad istogramma rappresenta gli stessi dati sotto forma differente.

La modalità di visualizzazione trend permette di avere una visione chiara dell'andamento dell'indicatore con il passare degli anni, a livello globale o nel paese selezionato tramite gli appositi filtri.

L'ultima modalità di visualizzazione è certamente quella più riutilizzabile per sviluppare progetti esterni a GlobalStat, in quanto permette non solo di consultare, ma anche di scaricare i dati relativi all'indicatore selezionato dettagliati per ciascun paese e anno.

Il download viene fatto in formato xlsx, quindi sotto forma di foglio di calcolo Excell. Proprio per questo motivo, in questa sezione non sono sfruttate particolari tecnologie relative al web semantico.

Ciascuna sezione analizzata fino ad ora presenta inoltre un la classica icona delle informazioni: alla sua pressione un popup permette all'utente di ottenere informazioni dettagliate sul significato dell'indicatore, note aggiuntive e soprattutto la fonte dalla quale è stato reperito, la quale spesso, come analizzeremo nel capitolo 4, fa ampio uso delle tecnologie del web semantico.

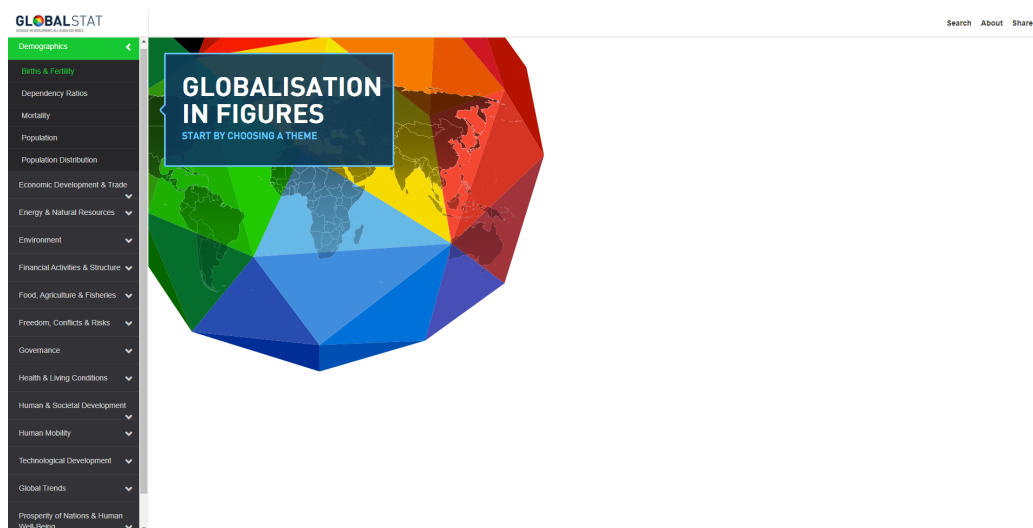


Figura 2.1: Homepage del sito GlobalStat.

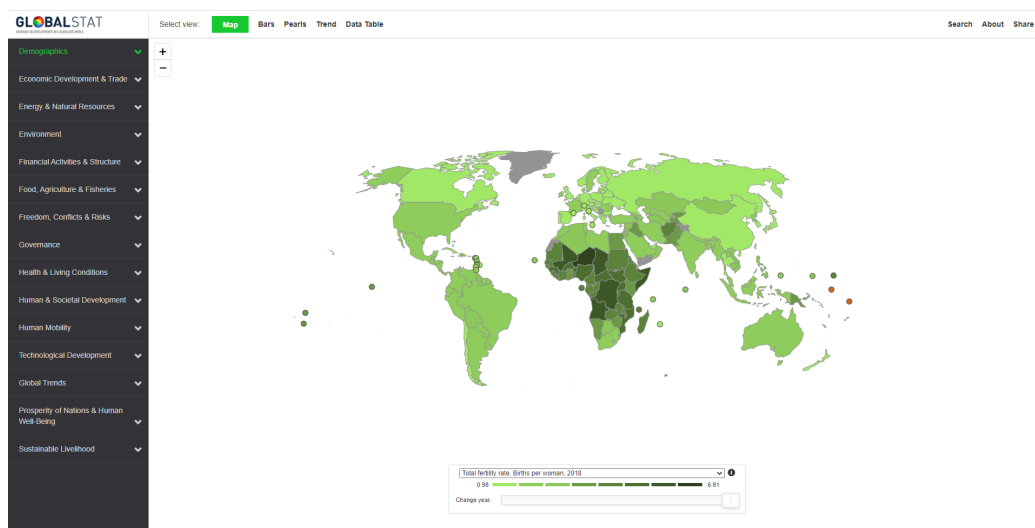


Figura 2.2: Mappa del sito GlobalStat.

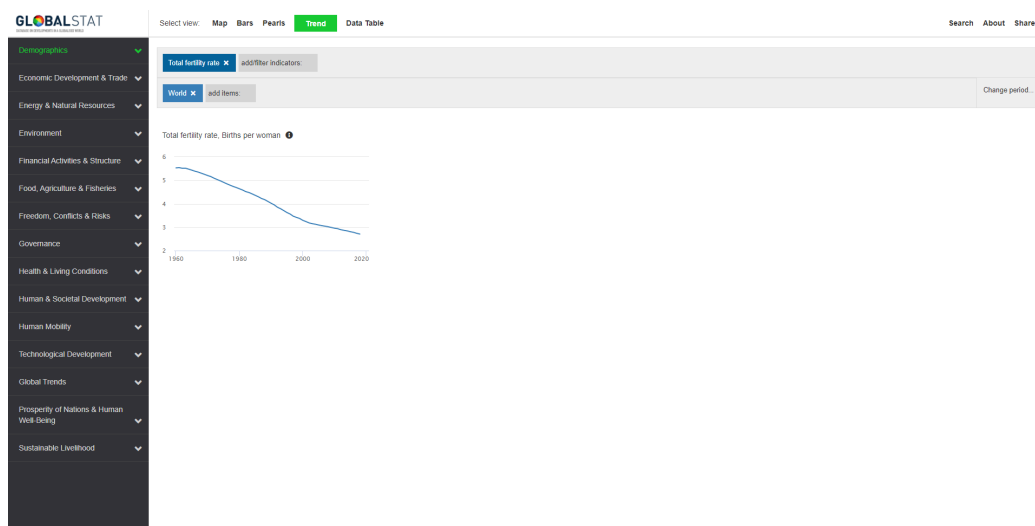


Figura 2.3: Trend del sito GlobalStat.

GLOBALSTAT

Select view: Map Bars Pearls Trend **Data Table**

Search About Share

Total fertility rate

add filter items

Change year

Total fertility rate, Births per woman

download selection

Country	1960 Q	1961 Q	1962 Q	1963 Q	1964 Q	1965 Q	1966 Q	1967 Q	1968 Q	1969 Q	1970 Q	1971 Q	1972 Q	1973 Q	1974 Q	1975 Q	1976 Q	1977 Q	1978 Q	1979 Q	1980 Q	1981 Q
Rwanda	8.19	8.19	8.2	8.2	8.2	8.2	8.2	8.2	8.21	8.22	8.23	8.25	8.28	8.31	8.34	8.37	8.4	8.43	8.45	8.46	8.45	8.4
Kenya	7.95	8	8.04	8.08	8.1	8.12	8.13	8.13	8.12	8.1	8.08	8.05	8.01	7.96	7.91	7.84	7.77	7.7	7.62	7.54	7.45	7.3
Côte d'Ivoire	7.69	7.72	7.75	7.78	7.81	7.84	7.87	7.89	7.91	7.93	7.94	7.94	7.94	7.93	7.91	7.88	7.83	7.76	7.68	7.59	7.4	
Jordan	7.69	7.8	7.9	7.98	8.03	8.06	8.05	8.03	8	7.97	7.93	7.88	7.82	7.75	7.68	7.6	7.52	7.46	7.39	7.33	7.26	7.1
Samoa	7.65	7.65	7.63	7.6	7.57	7.52	7.46	7.4	7.33	7.27	7.19	7.12	7.04	6.95	6.86	6.76	6.66	6.55	6.43	6.32	6.2	6.0
Dominican Republic	7.56	7.49	7.4	7.3	7.19	7.05	6.9	6.74	6.56	6.38	6.18	5.98	5.78	5.58	5.39	5.2	5.02	4.84	4.68	4.53	4.38	4.2
Algeria	7.52	7.57	7.61	7.65	7.67	7.68	7.68	7.67	7.67	7.66	7.64	7.62	7.6	7.56	7.51	7.43	7.34	7.23	7.11	6.96	6.79	6.6
Syria	7.47	7.5	7.52	7.54	7.56	7.56	7.57	7.57	7.57	7.57	7.57	7.57	7.56	7.54	7.51	7.47	7.42	7.36	7.29	7.2	7.09	6.9
Honduras	7.46	7.45	7.44	7.44	7.44	7.44	7.44	7.42	7.38	7.33	7.27	7.2	7.11	7.03	6.94	6.84	6.75	6.65	6.54	6.43	6.31	6.1
Niger	7.45	7.47	7.49	7.51	7.52	7.53	7.54	7.54	7.55	7.56	7.57	7.58	7.6	7.62	7.64	7.67	7.7	7.74	7.78	7.81	7.84	7.8
Afghanistan	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.45	7.4
Nicaragua	7.37	7.31	7.25	7.19	7.13	7.07	7.03	6.98	6.94	6.9	6.86	6.82	6.77	6.71	6.65	6.57	6.5	6.41	6.33	6.23	6.14	6.0
Tonga	7.36	7.35	7.3	7.23	7.12	6.97	6.8	6.59	6.37	6.15	5.94	5.76	5.62	5.51	5.45	5.43	5.44	5.47	5.51	5.54	5.55	5.5
Madagascar	7.3	7.3	7.3	7.31	7.31	7.31	7.31	7.3	7.3	7.29	7.27	7.25	7.22	7.19	7.15	7.1	7.04	6.98	6.9	6.81	6.73	6.6
Somalia	7.25	7.25	7.26	7.26	7.26	7.26	7.26	7.25	7.23	7.21	7.18	7.15	7.12	7.09	7.06	7.03	7.02	7	7	7	7.01	7.0
Oman	7.25	7.25	7.25	7.26	7.26	7.27	7.28	7.28	7.29	7.29	7.31	7.35	7.41	7.5	7.62	7.75	7.89	8.02	8.14	8.23	8.3	8.3
Kuwait	7.24	7.27	7.31	7.35	7.38	7.41	7.42	7.4	7.36	7.28	7.17	7.02	6.84	6.63	6.4	6.18	5.96	5.77	5.6	5.45	5.32	5.1
Saint Vincent and the Grenadines	7.22	7.16	7.07	6.98	6.88	6.76	6.63	6.49	6.35	6.19	6.01	5.83	5.63	5.42	5.2	4.97	4.75	4.54	4.34	4.15	3.99	3.8
Saudi Arabia	7.22	7.23	7.24	7.25	7.26	7.26	7.26	7.26	7.27	7.27	7.28	7.29	7.3	7.31	7.31	7.31	7.31	7.3	7.28	7.25	7.21	7.1
Libya	7.2	7.24	7.31	7.4	7.52	7.65	7.78	7.91	8.01	8.09	8.13	8.15	8.14	8.1	8.08	7.97	7.87	7.74	7.58	7.41	7.22	7.0
Vanuatu	7.2	7.12	7.03	6.94	6.84	6.73	6.63	6.53	6.43	6.35	6.27	6.2	6.13	6.07	6	5.93	5.86	5.79	5.72	5.65	5.58	5.5

Figura 2.4: Tabella del sito GlobalStat.

Age-specific fertility rates | 25-29 years

Definition Notes Data source

Age-specific fertility rates (UNDESA)

Number of births to women in a particular age group, divided by the number of women in that age group. The age groups used are: 15-19, 20-24,.....45-49. The data refer to five-year periods running from 1 July to 30 June of the initial and final years. (metadata - UNDESA)

Additional information

Countries or areas listed individually are only those with 90,000 inhabitants or more in 2017.

Data refer to five-year periods running from 1 July of the initial year to 30 June of the final year. The dataset reports the final year.

For a detailed explanation of methodology used by UN for population estimates and projections, see the pdf file: United Nations, Department of Economic and Social Affairs, Population Division (2014). [World Population Prospects: The 2012 Revision, Methodology](#) of the United Nations Population Estimates and Projections.

<https://population.un.org/wpp/Download/Standard/Population/>

(metadata – UNDESA)

Figura 2.5: Info del sito GlobalStat.

Capitolo 3

Metodologie e sorgenti dati

I dati raccolti sono sensibili al carattere di vasta portata del progetto; per questo motivo è stato necessario provvedere apposite metodologie atte a produrre statistiche effettivamente utili e soprattutto il più possibile affidabili. GlobalStats raccoglie dati sui 193 stati nazionali sovrani riconosciuti membri delle Nazioni Unite (ONU) [7] e, sulla base di questi, fornisce una panoramica sulle prestazioni dei singoli stati, dei continenti, di undici comunità di cooperazione e integrazione regionale e di organizzazioni internazionali. Questa forte cooperazione con le comunità è il pilastro principale dell'aggregazione dei dati effettuata da GlobalStat, in quanto tramite esse gli stati si impegnano in una collaborazione regionale per sostenersi a vicenda e migliorare lo sviluppo e le condizioni di vita.

Sebbene tendenzialmente i confini tra stati siano statici nel tempo, per tenere traccia di eventuali separazioni o formazione di nuovi territori, GlobalStat calcola in maniera dinamica i dati aggregati adattandosi al formato con cui vengono resi disponibili. Esempio lampante in cui questa metodologia ha permesso di avere dati coerenti nonostante le vicende politiche è la dissoluzione della Cecoslovacchia avvenuta nel 1993.

Anche in fase di visualizzazione, quindi, i dati rispecchieranno di anno in anno la composizione degli stati facenti parte dell'ONU, rendendo più complesso il calcolo di medie e la possibilità di effettuare comparazioni nel tempo.

Seguendo la prassi di statistica internazionale, GlobalStat calcola e rende disponibile dati aggregati se sono soddisfatte due condizioni:

1. L'indicatore in questione è disponibile per più della metà dei membri del gruppo
2. Se la prima condizione è soddisfatta, il valore aggregato viene calcolato solo se la popolazione totale dei paesi per i quali sono disponibili i dati rappresenta almeno $2/3$ della popolazione totale del gruppo

Indipendentemente da tali condizioni, seguendo un approccio prudentiale, i dati per Asia e Europa non vengono calcolati fino al 1992, anno dello scioglimento dell'URSS.

Questa scelta è dovuta al fatto che i confini geografici dell'URSS erano condivisi tra due continenti, per i quali non è possibile dedurre il contributo dell'Unione in maniera isolata.

Se le condizioni 1 e 2 sono soddisfatte, GlobalStat applica metodologie di calcolo degli aggregati differenti dipendentemente dalla tipologia di dati:

- Gli aggregati dei dati dei singoli paesi, se espressi come valori assoluti, sono calcolati semplicemente tramite la somma ignorando i valori mancanti
- Gli aggregati dei dati dei singoli paesi, se espressi come valori rapporti (tassi/proporzioni/percentuali), sono calcolati come medie ponderate sulla base del conteggio della popolazione sul totale, ignorando i valori mancanti

GlobalStat si affida a ben 113 fonti differenti di dati dipendentemente dagli indicatori necessari e dai paesi di riferimento.

Pertanto, unendo questo fattore alle metodologie sopra elencate, quando un utente consulta GlobalStat, dovrebbe essere consapevole di questa peculiarità multi-fonte ed essere quindi cauto nel confrontare dati e indicatori differenti, poiché potrebbero potenzialmente variare e presentare incoerenze dovute all'intervallo temporale, alle pratiche di raccolta e dalla loro specificità.

Capitolo 4

Architettura e tecnologie utilizzate

GlobalStat, come analizzato in precedenza, funge principalmente da aggregatore di numerosissimi dati provenienti da sorgenti differenti, le quali possono essere a loro volta aggregatori piuttosto che semplici banche dati.

Nel primo caso GlobalStat produce le proprie statistiche a partire dai file riassuntivi tabulari, resi disponibili per esempio in formato .xlsx, facendo uso ridotto delle tecnologie del web semantico in maniera diretta essendo precedentemente sfruttate dall'aggregatore di riferimento.

In tutti gli altri casi, invece, nei quali le stime devono essere effettuate fondendo dati provenienti da fonti differenti, è fondamentale l'utilizzo di tecnologie per disambiguare il significato delle risorse e interconnetterle all'interno del web.

4.1 RDF, RDFS

RDF è la tecnologia alla base della sezione open data del progetto, fondamentale per descrivere e mettere in relazione tra loro le risorse all'interno del web tramite l'utilizzo di triple.

Tendenzialmente, in aggiunta ad RDF, viene utilizzato RDFS, che fornisce un ulteriore livello di conoscenza connettendo i concetti forniti tramite RDF alle primitive di corrispondenza (Class, Subclass, Property ecc.), permettendo così di fare inferenza e interrogare il dominio di appartenenza. Le seguenti analisi prenderanno come riferimento una risorsa messa a disposizione da World Bank [8], un'istituzione finanziaria internazionale che fornisce dati a paesi di tutto il mondo ed è quindi una delle principali sorgenti dati sfruttate

da GlobalStat.

```
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:content="http://purl.org/rss/1.0/modules/content/"
  xmlns:dc="http://purl.org/dc/terms/"
  xmlns:dcat="http://www.w3.org/ns/dcat#"
  xmlns:sioc="http://rdfs.org/sioc/ns#">

  <rdf:Description rdf:about="https://datacatalog.worldbank.org/
    ↪ dataset/global-economic-prospects">
    <rdf:type rdf:resource="http://rdfs.org/sioc/ns#Item"/>
    <rdf:type rdf:resource="http://xmlns.com/foaf/0.1/Document"/>
    <content:encoded>&lt;p&gt;The Global Economic Prospects</
      ↪ content:encoded>
    <dc:creator></dc:creator>
    <dc:Frequency>3</dc:Frequency>
    <dcat:granularity></dcat:granularity>
    <dc:license>notspecified</dc:license>
    <dc:accessRights></dc:accessRights>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:Distribution rdf:resource="https://datacatalog.wb.org/
      ↪ dataset/global-economic/resource/ID"/>
    <dcat:theme rdf:resource="https://datacatalog.wb.org/
      ↪ taxonomy_term/2031"/>
    <dc:temporal>2018 - 2022</dc:temporal>
    <dc:title>Global Economic Prospects</dc:title>
    <dc:date rdf:datatype="http://www.w3.org/2001/XMLSchema#
      ↪ dateTime">2017-01-31T00:28:40-05:00</dc:date>
    <dc:created rdf:datatype="http://www.w3.org/2001/XMLSchema#
```

```
    ↪ dateTime">2017-01-31T00:28:40-05:00</dc:created>
<dc:modified rdf:datatype="http://www.w3.org/2001/XMLSchema#
    ↪ dateTime">2021-02-09T00:14:01-05:00</dc:modified>
<sioc:num_replies rdf:datatype="http://www.w3.org/2001/
    ↪ XMLSchema#integer">0</sioc:num_replies>
</rdf:Description>
</rdf:RDF>
```

Il file XML in questione fa riferimento ad una specifica risorsa disponibile nel catalogo di WordBank, raggiungibile tramite l'apposito URI presente in **rdf:about**.

A partire dal tag **rdf:type** sono elencate tutte le proprietà della risorsa. Queste ultime rappresentano delle specifiche caratteristiche che, per essere riconosciute in maniera standard e inconfutabile, sono connesse ad uno specifico vocabolario che ne determina l'ontologia. Nel caso specifico vengono sfruttati tre dizionari differenti, ciascuno relativo ad un ambito ed in grado di disambiguare al meglio il significato della proprietà.

Il sistema di metadati Dublin Core [9] (identificabile dalla proprietà **dc:**) definisce i quindici elementi core per rappresentare le risorse, vale a dire "Creator", "Date", "Description", "Format" ecc.

Il sistema Data Catalog Vocabulary [10] (identificato da **dcat:**) è invece utilizzato per descrivere i dataset e i servizi all'interno di un catalogo utilizzando un modello e un vocabolario standard, così da facilitare l'aggregazione (fondamentale all'interno di GlobalStat) e la consultazione.

L'ultima specifica di ontologie utilizzata è la Semantically-Interlinked Online Communities [11] (**sioc:**), che fornisce i concetti e le proprietà principali per descrivere informazioni relative a comunità online.

Questo esempio era solo un caso specifico selezionato tra le centinaia di sorgenti alle quali GlobalStat fa riferimento e alla miriade di risorse in esse contenute.

XML non è quindi l'unico formato disponibile; ciascuna sorgente sfrutterà i formati che ritiene più consoni, WordBank stesso permette di reperire i dati anche sotto forma di JSON, venendo così il più possibile incontro alle esigenze e alle preferenze dell'utilizzatore.

In maniera analoga anche la scelta dei vocabolari per le ontologie varia e viene fatta dipendentemente dagli elementi che devono essere rappresentati all'interno della risorsa e al grado di specificità richiesto dall'ente.

4.2 OWL

Sia DCAT che SIOC sono definiti sfruttando **OWL**, attualmente disponibile alla versione 2.

OWL 2 Web Ontology Language [12], abbreviato OWL 2, è un linguaggio per le ontologie utilizzato nel web semantico. Le ontologie forniscono classi, proprietà e valori che sono memorizzati sotto forma di documenti del web semantico.

Le ontologie fornite da OWL possono essere combinate tra di loro oltre che essere utilizzate per arricchire informazioni contenuto in documenti RDF.

OWL è un linguaggio computazionale basato sulla logica: la conoscenza espressa tramite esso può essere dedotta da software i quali potranno inoltre verificare la consistenza e fare inferenza di nuove informazioni.

Per quanto concerne DC, invece, essendo uno dei primi progetti nati in ambito ontologie, non è definito mediante OWL.

La differenza nel linguaggio rende più difficoltosa l'integrazione di DC con i dizionari più recenti; inoltre con il passare degli anni sono emerse numerose limitazioni come l'utilizzo di "LastName" e "FirstName" per identificare il creatore di una risorsa.

Possono infatti essere presenti numerosi casi di omonimia all'interno di un singolo database, oltre ad avere spesso difficoltà a disambiguare il nome dal cognome.

Per questi motivi, negli ultimi anni, stanno nascendo numerosi progetti atti a costruire versioni di DC basate su OWL, facilitandone così l'integrazione con ontologie definite ad hoc e disambiguando ulteriormente i termini [13].

Nonostante le numerose limitazioni analizzate, vista l'ampia diffusione dovuta alla longevità di Dublin Core, è tutt'ora uno dei vocabolari più diffusi per identificare i creatori di risorse.

Capitolo 5

Commenti finali

GlobalStat è sicuramente un progetto ambizioso, viste le difficoltà e i problemi che costantemente deve affrontare per portare statistiche affidabili, oltre che ammirevole in quanto si pone come obbiettivo principale quello di migliorare nel complesso la vita delle persone.

Quest'ultima, tra tutte le motivazioni, è stata quella che mi ha maggiormente spinto ad approfondire questo lavoro.

Attualmente GlobalStat funge principalmente da aggregatore: grazie alla sua folta rete di collaborazioni ottiene dati su scala globale, spesso a loro volta già precedentemente processati, e produce tramite essi statistiche che mette gratuitamente a disposizione degli utenti, principalmente in modalità visuale.

Non vi è quindi la possibilità di reperire i dati direttamente in maniera programmata, ma il sito stesso effettua un reindirizzamento alla sorgente originale, la quale adotterà poi delle specifiche politiche e tecnologie per l'accesso ai dati.

Sebbene per gli enti governativi e le organizzazioni internazionali i dati siano già accessibili in maniera automatica tramite un agente software, la speranza è che in futuro tale possibilità venga estesa a chiunque, permettendo così una diffusione di queste importanti informazioni su scala ancora più ampia.

Inoltre, come tendenzialmente accade in statistica, l'importanza e il valore cresce insieme alla quantità di dati disponibili: proprio per questo, negli anni a venire, i vantaggi che GlobalStat sarà in grado offrire emergeranno maggiormente ed è quindi fondamentale continuino ad essere accessibili in maniera semplice e automatizzata.

Bibliografia

- [1] European university institute. Online; consultato il 16/03/2021.
- [2] Global governance programme. Online; consultato il 16/03/2021.
- [3] Globalstat. Online; consultato il 16/03/2021.
- [4] Fundação francisco manuel dos santo. Online; consultato il 16/03/2021.
- [5] Organizzazione per la cooperazione e lo sviluppo economico. Online; consultato il 16/03/2021.
- [6] Servizio ricerca del parlamento europeo. Online; consultato il 19/03/2021.
- [7] Stati membri delle nazioni unite. Online; consultato il 17/03/2021.
- [8] Datacatalog worldbank. Online; consultato il 25/03/2021.
- [9] Dublin core (dc). Online; consultato il 26/03/2021.
- [10] Data catalog vocabulary (dcat). Online; consultato il 26/03/2021.
- [11] Semantically-interlinked online communities (sioc). Online; consultato il 26/03/2021.
- [12] Web ontology language (owl). Online; consultato il 26/03/2021.
- [13] Dublin core and owl. Online; consultato il 26/03/2021.