DATA ON THE WEB

Capacitación
**Mejores Prácticas de Datos en la Web**

# Mini curso: Open Refine

Newton Calegari
@newtoncalegari

ceweb.br

| title | author | year |
|---|---|---|
| Things Fall Apart | Chinua Achebe | |
| Fairy tales | Hans Christian Andersen | 1835-37 |
| The Divine Comedy | Dante Alighieri | 1308-1321 |
| The Epic Of Gilgamesh | Unknown | 18th - 17th ce |
| The Book Of Job | Unknown | 7th - 4th cent |
| One Thousand and One Nights | Unknown | 700-1500 |
| Njál's Saga | Unknown | 13th century |
| Pride and Prejudice | Jane Austen | |
| Le Père Goriot | Honoré de Balzac | |
| Molloy |  Malone Dies |  The Unnama |
| The Decameron | Giovanni Boccaccio | 1349-53 |
| Ficciones | Jorge Luis Borges | 1944-86 |
| Wuthering Heights | Emily Brontë | |
| The Stranger | Albert Camus | |
| Poems | Paul Celan | |
| Journey to the End of the Night | Louis-Ferdinand Céline | |
| Don Quijote De La Mancha | Miguel de Cervantes | 1605 (part 1) |
| The Canterbury Tales | Geoffrey Chaucer | 14th century |
| Stories | Anton Chekhov | |
| Nostromo | Joseph Conrad | |
| Great Expectations | Charles Dickens | |
| Jacques the Fatalist | Denis Diderot | |

```turtle
@prefix schema: <http://schema.org/> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix dct: <http://purl.org/dc/terms/> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
@prefix owl: <http://www.w3.org/2002/07/owl#> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .

<https://en.wikipedia.org/wiki/Things_Fall_Apart> a schema:Book ;
    schema:name "Things Fall Apart"@en ;
    schema:author "Chinua Achebe" ;
    schema:datePublished "1958"^^xsd:date ;
    schema:locationCreated "Nigeria" ;
    schema:inLanguage "English" ;
    schema:numberOfPages "209" .

<https://en.wikipedia.org/wiki/Fairy_Tales_Told_for_Children._First
    schema:name "Fairy tales"@en ;
    schema:author "Hans Christian Andersen" ;
    schema:datePublished "1836"^^xsd:date ;
    schema:locationCreated "Denmark" ;
    schema:inLanguage "Danish" ;
    schema:numberOfPages "784" .

<https://en.wikipedia.org/wiki/Divine_Comedy> a schema:Book ;
```

127.0.0.1:3333

Pesquisar

Google refine *A power tool for working with messy data.*

**Create Project**

Open Project

Import Project

**Create a project by importing data. What kinds of data files can I import?**

TSV, CSV, *SV, Excel (.xls and .xlsx), JSON, XML, RDF as XML, and Google Data documents are all supported. Support for other formats can be added with Google Refine extensions.

Get data from

**This Computer**

Web Addresses (URLs)

Clipboard

Google Data

Locate one or more files on your computer to upload:

Selecionar arquivo…    Nenhum arquivo selecionado.

**Next »**

Version 2.5 [r2407]

Help
About

**Top 100 books - CSV**

newtoncalegari.com.br/dwbp-costa-rica/

| | |
|---|---|
| title | schema:name |
| author | schema:author |
| year | schema:datePublished |
| country | schema:locationCreated |
| language | schema:inLanguage |
| pages | schema:numberOfPages |

RDF Schema Alignment

The RDF schema alignment skeleton below specifies how the RDF data that will get generated from your grid-shaped data. The cells in each record of your data will get placed into nodes within the skeleton. Configure the skeleton by specifying which column to substitute into which node.

**Base URI:** http://localhost:3333/ edit

**RDF Skeleton** | RDF Preview

Available Prefixes: schema rdfs foaf dct xs

wikipedia_link URI

definir el recurso como un URI

✗schema:Book
   add rdf:type

✗ ➤schema:author→   □   **author** cell

✗ ➤schema:datePublished→   □   **year2** cell

✗ ➤schema:locationCreated→   □   **country** cell

✗ ➤schema:inLanguage→   □   **language** cell

✗ ➤schema:numberOfPages→   □   **pages** cell

add property

Add another root node     Save

javascript:() OK   Cancel

definir el recurso como un URI
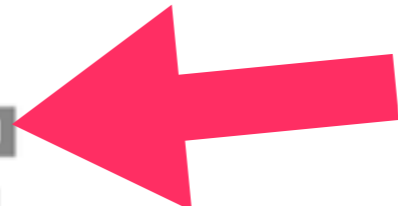
# RDF Schema Alignment

The RDF schema alignment skeleton below specifies how the RDF data that will get generated from your grid-shaped data. The cells in each record of your data will get placed into nodes within the skeleton. Configure the skeleton by specifying which column to substitute into which node.

**Base URI:** http://localhost:3333/ edit

**RDF Skeleton** | RDF Preview

Available Prefixes:     schema rdfs foaf dct xsd owl rdf ➕add prefix ⚙manage prefixes

**wikipedia_link** URI      ☐  ✕ ➤—schema:name→          ☐  **title** cell
✕schema:Book                ✕ ➤—schema:author→        ☐  **author** cell
add rdf:type               ✕ ➤—schema:datePublished→  ☐  **year2** cell
                           ✕ ➤—schema:locationCreated→ ☐  **country** cell
                           ✕ ➤—schema:inLanguage→      ☐  **language** cell
                           ✕ ➤—schema:numberOfPages→   ☐  **pages** cell

definir el formato del contenido (value)

Add another root node                                            Save

javascript:{} OK   Cancel

**RDF Schema Alignment**

The RDF schema alignment skeleton below specifies how the RDF data that will get generated from your grid-shaped data. The cells in each record of your data will get placed into nodes within the skeleton. Configure the skeleton by specifying which column to substitute into which node.

**Base URI:** http://localhost:3333/

RDF Skeleton | RDF Preview

Available Prefixes:

wikipedia_link
×sch
add rdf:type

**RDF Node**
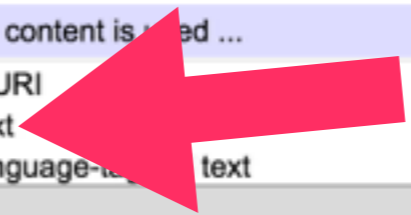
Use content from cell...
- ( row index)
- title
- ● author
- year2
- country
- language
- pages
- wikipedia_link
- Constant Value

The cell's content is used ...
- as a URI
- ● as text
- as language-... text
- as integer number
- as non-integer number
- as date (YYYY-MM-DD)
- as dateTime (YYYY-MM-DD HH:MM:SS)
- as boolean
- as custom datatype (specify type URI)
- as a blank node

Use custom expression...

value

preview/edit

Add another root node

Save

OK | Cancel

**definir el formato del contenido (value)**

definir las propiedades (predicados)

# RDF Schema Alignment

The RDF schema alignment skeleton below specifies how the RDF data that will get generated from your grid-shaped data placed into nodes within the skeleton. Configure the skeleton by specifying which column to substitute into which node.

**Base URI:** http://localhost:3333/ edit

RDF Skeleton    **RDF Preview**

This is a sample `Turtle` representation of (up-to) the *first 10* rows

```
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix dct: <http://purl.org/dc/terms/> .
@prefix owl: <http://www.w3.org/2002/07/owl#> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .


<https://en.wikipedia.org/wiki/Things_Fall_Apart> a schema:Book ;
        schema:name "Things Fall Apart"@en ;
        schema:author "Chinua Achebe" ;
        schema:datePublished "1958"^^xsd:date ;
        schema:locationCreated "Nigeria" ;
        schema:inLanguage "English" ;
        schema:numberOfPages "209" .

<https://en.wikipedia.org/wiki/Fairy_Tales_Told_for_Children._First_Collection.> a schema:Book
        schema:name "Fairy tales"@en ;
            author "
```

https://en.wikipedia.org/wiki/Things_Fall_Apart

contenido en Turtle/RDF

# RDF Schema Alignment

The RDF schema alignment skeleton below specifies how the RDF data that will get generated from your grid-shaped dat placed into nodes within the skeleton. Configure the skeleton by specifying which column to substitute into which node.

**Base URI:** http://localhost:3333/ edit

RDF Skeleton    **RDF Preview**

This is a sample `Turtle` representation of (up-to) the *first 10* rows

```
@prefix rdrs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix dct: <http://purl.org/dc/terms/> .
@prefix owl: <http://www.w3.org/2002/07/owl#> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .


<https://en.wikipedia.org/wiki/Things_Fall_Apart> a schema:Book ;
        schema:name "Things Fall Apart"@en ;
        schema:author "Chinua Achebe" ;
        schema:datePublished "1958"^^xsd:date ;
        schema:locationCreated "Nigeria" ;
        schema:inLanguage "English" ;
        schema:numberOfPages "209" .

<https://en.wikipedia.org/wiki/Fairy_Tales_Told_for_Children._First_Collection.> a schema:Book
        schema:name "Fairy tales"@en ;
        author "
```

https://en.wikipedia.org/wiki/Things_Fall_Apart

contenido en Turtle/RDF

presione OK cuando finalize

exporte los datos en RDF

```turtle
@prefix schema: <http://schema.org/> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix dct: <http://purl.org/dc/terms/> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
@prefix owl: <http://www.w3.org/2002/07/owl#> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .

<https://en.wikipedia.org/wiki/Things_Fall_Apart> a schema:Book ;
    schema:name "Things Fall Apart"@en ;
    schema:author "Chinua Achebe" ;
    schema:datePublished "1958"^^xsd:date ;
    schema:locationCreated "Nigeria" ;
    schema:inLanguage "English" ;
    schema:numberOfPages "209" .

<https://en.wikipedia.org/wiki/Fairy_Tales_Told_for_Children._First
    schema:name "Fairy tales"@en ;
    schema:author "Hans Christian Andersen" ;
    schema:datePublished "1836"^^xsd:date ;
    schema:locationCreated "Denmark" ;
    schema:inLanguage "Danish" ;
    schema:numberOfPages "784" .

<https://en.wikipedia.org/wiki/Divine_Comedy> a schema:Book ;
```

datos en RDF (Turtle)

Este arquivo XML não parece ter qualquer informação de estilo associado a ele. A estrutura do documento é mostrada abaixo.

```xml
- <rdf:RDF>
  - <rdf:Description rdf:about="https://en.wikipedia.org/wiki/Things_Fall_Apart">
      <rdf:type rdf:resource="http://schema.org/Book"/>
      <schema:name xml:lang="en">Things Fall Apart</schema:name>
      <schema:author>Chinua Achebe</schema:author>
      <schema:datePublished rdf:datatype="http://www.w3.org/2001/XMLSchema#date">1958</schema:datePublished>
      <schema:locationCreated>Nigeria</schema:locationCreated>
      <schema:inLanguage>English</schema:inLanguage>
      <schema:numberOfPages>209</schema:numberOfPages>
    </rdf:Description>
  - <rdf:Description rdf:about="https://en.wikipedia.org/wiki/Fairy_Tales_Told_for_Children._First_Collection.">
      <rdf:type rdf:resource="http://schema.org/Book"/>
      <schema:name xml:lang="en">Fairy tales</schema:name>
      <schema:author>Hans Christian Andersen</schema:author>
      <schema:datePublished rdf:datatype="http://www.w3.org/2001/XMLSchema#date">1836</schema:datePublished>
      <schema:locationCreated>Denmark</schema:locationCreated>
      <schema:inLanguage>Danish</schema:inLanguage>
      <schema:numberOfPages>784</schema:numberOfPages>
    </rdf:Description>
  - <rdf:Description rdf:about="https://en.wikipedia.org/wiki/Divine_Comedy">
      <rdf:type rdf:resource="http://schema.org/Book"/>
      <schema:name xml:lang="en">The Divine Comedy</schema:name>
      <schema:author>Dante Alighieri</schema:author>
      <schema:datePublished rdf:datatype="http://www.w3.org/2001/XMLSchema#date">1315</schema:datePublished>
      <schema:locationCreated>Italy</schema:locationCreated>
      <schema:inLanguage>Italian</schema:inLanguage>
      <schema:numberOfPages>928</schema:numberOfPages>
    </rdf:Description>
  - <rdf:Description rdf:about="https://en.wikipedia.org/wiki/Epic_of_Gilgamesh">
      <rdf:type rdf:resource="http://schema.org/Book"/>
      <schema:name xml:lang="en">The Epic Of Gilgamesh</schema:name>
      <schema:author>Unknown</schema:author>
      <schema:datePublished rdf:datatype="http://www.w3.org/2001/XMLSchema#date">-1700</sc
      <schema:locationCreated>Sumer and Akkadian Empire</schema:locationCreated>
```

datos en RDF / XML