DATA ON
THE WEB

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

# Mini curso: Open Refine

Newton Calegari
@newtoncalegari

ceweb.br

openrefine.org

Pesquisar

Follow us on: **Github** **Twitter**

Google™ Custom Search **Search** ×

# OPEN
# Refine

A free, open source, powerful tool
for working with messy data

**Home**

**Download**

**Documentation**

**Community**

**Post archive**

OpenRefine News:
Spring 2016

OpenRefine News:
December 2015

OpenRefine News:
November 2015

OpenRefine News:
October 2015

OpenRefine News:
September 2015

OpenRefine News:
August 2015

#Welcome!

OpenRefine (formerly Google Refine) is a powerful tool for working with messy data: cleaning it; transforming it from one format into another; and extending it with web services and external data.

Please note that since October 2nd, 2012, Google is not actively supporting this project, which has now been rebranded to OpenRefine. Project development, documentation and promotion is now fully supported by volunteers. Find out more about the history of OpenRefine and how you can help the community.

## Using OpenRefine - The Book

**Using OpenRefine**, by Ruben Verborgh and Max De Wilde, offers a great introduction to OpenRefine. Organized by recipes with hands on examples, the book covers the following topics:

1. Import data in various formats
2. Explore datasets in a matter of seconds
3. Apply basic and advanced cell transformations
4. Deal with cells that contain multiple values
5. Create instantaneous links between datasets
6. Filter and partition your data easily with regular expressions
7. Use named-entity extraction on full-text fields to automatically identify topics
8. Perform advanced data operations with the General Refine Expression Language

## Introduction to OpenRefine

## 1. Explore Data

**Using OpenRefine**

openrefine.org

**refine.deri.ie**

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

## Mapping vocabulary to data

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

Uploading to RDF Store

200.160.6.177:8080/fuseki/

**SPARQL - basics**

```
SELECT ?s  ?p  ?o
WHERE {
   ?s  ?p  ?o
}
```

## SPARQL - basics

```sparql
PREFIX schema: <http://schema.org/>
SELECT ?title ?author ?country
WHERE {
  ?book schema:author        ?author .
  ?book schema:name          ?title .
  ?book schema:locationCreated ?country .
  FILTER(?country = 'Brazil')
}
```