

Hands-on Data Analysis on Messy Datasets

June 25, 2018

Motivation of The Lab Session




A Real-World Scenario: The Rodents dataset

- ▶ data on rodents during a survey.
- ▶ useful for studying population dynamics and species interactions.
- ▶ size: 35549 rows and 35 columns.
- ▶ each row denotes the information collected on an individual rodent.
- ▶ data is provided by Ernest et al. (2018)
- ▶ meta-data is also available at ¹.

¹http://esapubs.org/archive/ecol/E090/118/Portal_rodent_metadata.htm

Data Wrangling Challenges in the Rodents Dataset

- ▶ NOTE: I WILL GIVE SOME EXAMPLES HERE.
- ▶ **missing data.**
- ▶ **format variabilities:** typos, abbreviations, leading and trailing whitespace.
- ▶ some issues are addressed in ².

²<http://www.datacarpentry.org/python-ecology-lesson/> 

A Data Wrangling Tool: OpenRefine

- ▶ a tool for working with messy datasets.
- ▶ see Verborgh and De Wilde (2013) for details.
- ▶ useful links: the software ³ and the documentation ⁴.

³<http://openrefine.org>

⁴[https://github.com/OpenRefine/OpenRefine/wiki/Documentation-For-](https://github.com/OpenRefine/OpenRefine/wiki/Documentation-For-Users)

Installation

- ▶ detailed installation instructions ⁵.
- ▶ download the file depending on the OS at ⁶.
- ▶ install OpenRefine as follows:
 - ▶ Linux: extract.
 - ▶ Mac: open, drag icon into the Applications folder.
 - ▶ Windows: unzip.

⁵<http://openrefine.org/download.html>

⁶<https://github.com/OpenRefine/OpenRefine/releases/tag/2.8>

Running and Loading Data

- ▶ run OpenRefine depending on the operating system:
 - ▶ Linux: `./refine` in your installation folder
 - ▶ Mac: OpenRefine in your Applications folder
 - ▶ Windows: `.exe` file in your installation folder
- ▶ get the dataset:
 - ▶ clone the git repository at ⁷.
 - ▶ use the file in
 `datasets/Portal_rodents_19772002_scinameUUIDs.csv`.
- ▶ import the data:
 - ▶ click “Create Project”.
 - ▶ click “Choose Files”.
 - ▶ select `Portal_rodents_19772002_scinameUUIDs.csv`.
 - ▶ click “Next”.

⁷<https://github.com/tahaceritli/acm-summer-school>

Data Preview

- ▶ configuration page for importing.
- ▶ a subset of the data is shown.
- ▶ use the defaults.
- ▶ click “Create Project” on the top-right corner.

Checking for Unique Values

- ▶ click drop-down arrow in the “survey_id” column.
- ▶ select “Facet>Customized Facet>Duplicates Facet”.
- ▶ results in a binary facet of “true” or “false”.
- ▶ “true” facet denotes rows with unique values.

View Range of Values

- ▶ click the drop-down arrow in the “scientificName” column.
- ▶ select “Facet>Text Facet”.
- ▶ lists the values and their counts.
- ▶ any problems with the data?

Updating Cell Values

- ▶ notice the spelling errors, e.g. “Amphespiza bilineata” for “Amphispiza bilineata”.
- ▶ hover over the former and select “edit” to update its value.

Filtering Rows

- ▶ click the drop-down arrow in the “scientificName” column.
- ▶ select “Text Filter”.
- ▶ type “bai”, which lists 48 matching rows.

Clustering

- ▶ data often contains inconsistencies due to data collection procedures.
- ▶ “Clustering” helps to find cells in a column, that refers to the same entity with different values.
- ▶ various methods to determine clusters.

Clustering

- ▶ click the drop-down arrow in the “scientificName” column.
- ▶ select “Edit cells>Cluster and edit...”.
- ▶ change the method to nearest neighbor.
- ▶ you can now check boxes and merge the clusters.

Trimming Whitespace

- ▶ click the drop-down arrow in the “scientificName” column.
- ▶ select “Edit cells>Common transforms>Trim leading and trailing whitespace”.

Deliverables

- ▶ **NOTE: I WILL UPDATE THIS PART. COULDN'T FIND ENOUGH TIME.**
- ▶ **missing data:**
 - ▶ report missing data encodings used in each column.
 - ▶ check whether they are actually missing. you may want to check meta-data!
 - ▶ hint: commonly used encodings: “NA”, “N/A”, “Null”, “-1”, “-99”, etc.
- ▶ **format variabilities:**
 - ▶ report how many unique values exist in the “country” column of the original dataset.
 - ▶ is there any problems with the data? if so, explain how did you solve them?
 - ▶ report how many unique values you have after fixing the potential inconsistencies.
 - ▶ hint: “clustering” feature of OpenRefine could be useful for this task.

Submissions

- ▶ NOTE:I WILL EXPLAIN THIS PART A BIT MORE.
- ▶ 1 page of either .txt or .pdf (not a .doc file!).

References I

Ernest, M., Brown, J., Valone, T., and White, E. P. (2018). Portal Project Teaching Database. https://figshare.com/articles/Portal_Project_Teaching_Database/1314459.

Verborgh, R. and De Wilde, M. (2013). *Using OpenRefine*. Packt Publishing Ltd.

- ▶ NOTE: I WILL ALSO GIVE REFERENCE to FREIRE'S PRESENTATIONS.