

2:W35: Open Refine

Upload your answers to these questions:

2.1) Danish monarchs, create spreadsheet

1. Create a spreadsheet listing the names of Danish monarchs with their birth- and death-date and start and end year of reign. Make it *tidy*! They should be sortable by year of birth. Suitable source websites are *here* and *here*, but you can also use another source, provided you reference it. (Group collaboration is expected and welcome. Remember to attach this spreadsheet to Brightspace submission)

SEE ATTACHED FILE [monarchs.csv](#) IN THE FOLDER [hw_w5_2](#)

2.2) Does OpenRefine alter the raw data during sorting and filtering?

Not during sorting and filtering. However, when the spreadsheet is uploaded to OpenRefine, the data types changes. They can quickly be transformed into a different format by clicking on the column arrow > Edit cells > Common transforms.

2.3) Fixing the interviews dataset and look for water-deprived months

Fix the *interviews* dataset in OpenRefine enough to answer this question: "Which two months are reported as the most water-deprived/driest by the interviewed farmer households?"

I'll use the variable *months_no_water* to answer the question: *Which two months are reported as the most water-deprived/driest by the interviewed farmer households?"*.

I go to the column menu, choose the menu item « Split into several columns... »

7. Real-Data-Challenge: What are the 10 most frequent occupations (erhverv) among unmarried men and women in 1801 Aarhus? (hint: some expert judgement interpretation is necessary, look at the *HISCO classification* "Historical International Standard of Classification of Occupations" on *Dataverse* if ambitious)

First, I want to split the cells that contains 2 months into separate cells. I see that the cell contains hard brackets, the symbol ' and semicolon, and the semicolon is what separates the months. I therefore split the multi-valued cells **by separator ;**.

I here click the following:

- the small arrow next to the column **months_no_water**
- **Edit cells**
- **Split multi-valued cells..**
- then I fill **;** into the **Separator** field and click **OK**.

I now want to get rid of all the special characters. I do that by following these steps:

- Edit Cells
- Transform
- To remove all left square brackets and replace with nothing: ```value.replace("[", "")`
- To remove all right square brackets and replace with nothing: `value.replace("]", "")`
- To remove all ' signs I write: ``value.replace("'", "")`
- I then go to the left pane and **Sort by count**.

I see that the two months reported as the most water-deprived/driest are November (count = 41) and October (count = 38).

The screenshot shows a data visualization interface with a facet titled "months_no_water". The facet is sorted by count, showing the following data:

Month	Count
NULL	45
Nov	41
Oct	38
Sept	37
Aug	31
Sept	27
Oct	25
Dec	11
Oct	9
Nov	7
Sept	6

2.4) 10 most frequent occupations

Real-Data-Challenge: What are the 10 most frequent occupations (erhverv) among unmarried men and women in 1801 Aarhus? (hint: some expert judgement interpretation is necessary, look at the *HISCO classification* "Historical International Standard of Classification of Occupations" on *Dataverse* if ambitious)

Data:

<https://raw.githubusercontent.com/aarhusstadsarkiv/datasets/master/censuses/1801/census-1801-normalized.csv>

HISCO classification: <https://github.com/cedarfoundation/hisco>

I first load the data set into OpenRefine by copy pasting the link.

I have to do some cleaning since some cells with **erhverv** contains more than one value. I therefore split by `og, '`

I want to filter by **civilstand = ugift** and then group by men and women. I'll use **Facet** which allows me to group and also to filter the data by these values.

First I'll do a **Text Facet** which will group by identical text values in a specific column and then list unique values with the number of records it appears in.

So first, I scroll over the `civilstand` column. Then I click the down arrow and choose `Facet > Text facet`. In the left panel, I now see a box with every unique value in the `civilstand` column along with a number representing how many times that value occurs in the column. I click `ugift` to exclude all other categories.

×

—

civilstand

change invert reset

5 choices

Sort by: name count

Cluster

enke 1527

gift 8296

separeret 4

skilt 3

ugift 12282

(blank) 136

Facet by choice counts

exclude

I then do the same with `koen`. So first, I scroll over the `koen` column. Then I click the down arrow and choose Facet > Text facet. In the left panel, I now see a box with every unique value in the `koen` column along with a number representing how many times that value occurs in the column. I click `kvinde` to first have a look at the most frequent occupations for women.

I then do the same with the `erhverv` column. Then I click Sort by: count to see each count for each occupation for the women (because I selected `kvinde`). I then click on `mand` in the `koen` pane and compare the counts.

Gender	Female	Male
--------	--------	------

Gender	Female	Male
Occupations sorted by count	<div> <div> <div>koen</div> <div>change invert reset</div> </div> <div> <div>2 choices</div> <div>Sort by: name count</div> <div>Cluster</div> </div> <div> <div>kvinde 12282</div> <div>mand 13156</div> <div>(blank) 2</div> <div>Facet by choice counts</div> </div> </div> <div> <div> <div>erhverv</div> <div>change</div> </div> <div> <div>338 choices</div> <div>Sort by: name count</div> <div>Cluster</div> </div> <div> <div>Tienestepige 61</div> <div>hospitalslem 21</div> <div>Væverske 18</div> <div>lever af at spinde 17</div> <div>Inderste 14</div> <div>tjener faderen 13</div> <div>lever af almisse 12</div> <div>spinder 12</div> <div>Vanfør 10</div> <div>Almisselem 9</div> <div>indsidder 8</div> </div> </div>	<div> <div> <div>koen</div> <div>change invert reset</div> </div> <div> <div>2 choices</div> <div>Sort by: name count</div> <div>Cluster</div> </div> <div> <div>kvinde 12282</div> <div>mand 13156</div> <div>(blank) 2</div> <div>Facet by choice counts</div> </div> </div> <div> <div> <div>erhverv</div> <div>change</div> </div> <div> <div>714 choices</div> <div>Sort by: name count</div> <div>Cluster</div> </div> <div> <div>National Soldat 96</div> <div>soldat ved 1. Jyske Inf. Reg. 94</div> <div>nationalsoldat 61</div> <div>Tienestekarl 47</div> <div>læredreng 42</div> <div>Nationalsoldat 36</div> <div>Bonde og Gaardbeboer 32</div> <div>Tienestedræng 32</div> <div>Væver 32</div> <div>gårdskarl 30</div> <div>Soldat 27</div> </div> </div>

OBS: Scrolling down the facet panes, it becomes evident that the same occupations occur multiple times (as separate occupations) due to differences in spelling. The duplicates that I could choose to get rid of by combining them into one category are among others:

- different types of *soldier*
- *stuepige*/*Stuepige*
- *vanvittig*/*Vanvittig*
- *gaardbeboer*/*gårdbeboer*
- *_bonde*/*Bonde*

The different categories can be searched for and replaced using RegEx (for instance soldier category (`([\w-]+)(oldat)$`), but I just press in the left "erhverv" pane and edit manually. I see that the most frequent occupations for unmarried men is **soldat** (some sort of soldier) and the most frequent occupations (erhverv) among unmarried men and women in 1801 is **Tienestepige**.