MEMO

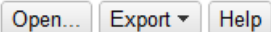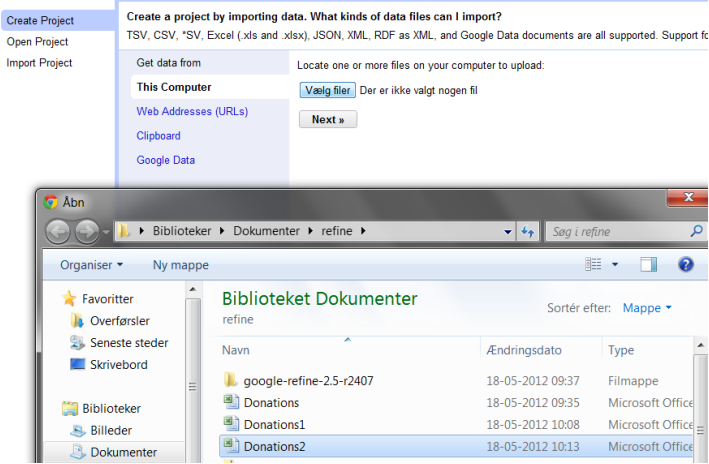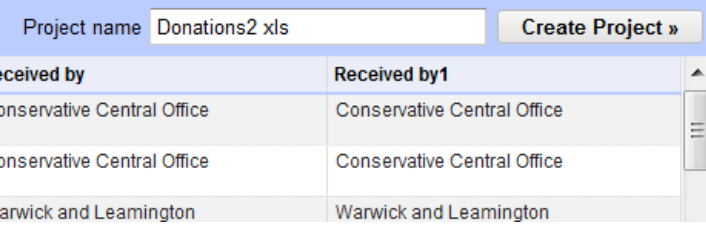| WHAT | WHY | HOW |
|---|---|---|
| | OpenRefine (ex-Google Refine) is a powerful tool for working with messy data, cleaning it, transforming it from one format into another, extending it with web services, and linking it to databases. | http://openrefine.org/ |
| **Who can use Refine?** | Everybody must use Refine. It's free. The latest totally stable version is Google Refine 2.5. It's free. We will use OpenRefine 2.6-rc2 Release Candidate 2, whish also is free and are working fine. | Download it to your PC, MAC or Linux. |
| **Open Refine** | OpenRefine opens in a browser in this url: http://127.0.0.1:3333/ It means it runs on your own machine. You are then ready to work on projects – either old or new after import of data. | Click on google-refine and run the program directly.  openrefine |

MEMO

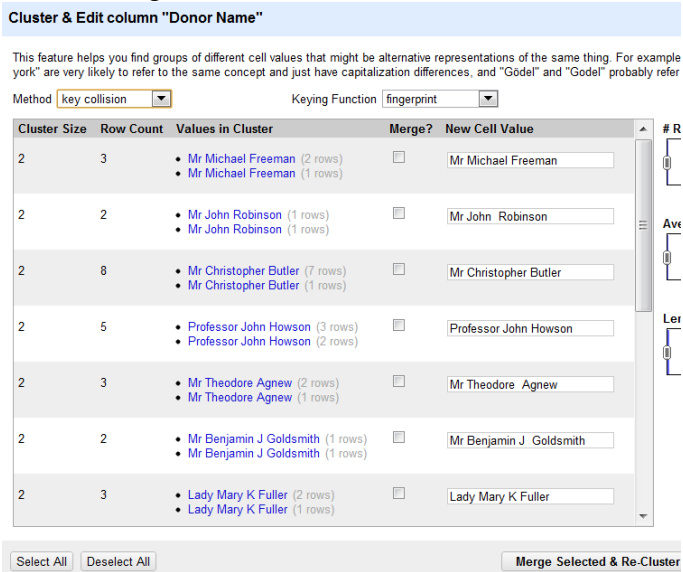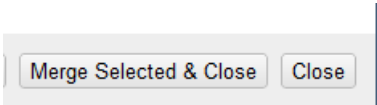| WHAT | WHY | HOW |
|---|---|---|
| **Prepare data** | Prepare data in a spreadsheet. Copy-paste the columns you will clean.<br>Then you in the end can compare the original and cleaned. | http://code.google.com/p/google-refine/ |
| **Open Refine** | Find the folder, where Google Refine is installed and start the program from here . | Click on google refine and run the program directly:<br><br>🔵 openrefine |
| **Start project** | If you should start a new project - click on the button "Open" in the upper right corner:<br><br>Open… \| Export ▾ \| Help | Click Open… |
| **Choose file** | <br>The project uses an excel file, Donations2, imported from the computer.<br>Note that you can also choose data from URL's and Google and clipboard.<br>URL's and Google provides unique opportunities to cleanse dirty data automatically and then send the cleaned data in to your own presentation.<br>We choose the basic model of cleaning a spreadsheet. | Click on:<br>Create Project<br>Choose Files<br>Mark the relevant file<br>Click Open<br>Click Next >> |
| **Create project** | The Import Wizard allows you to change the settings. When importing from Excel it is usually not necessary.<br> | Click on Create Project >> |

MEMO

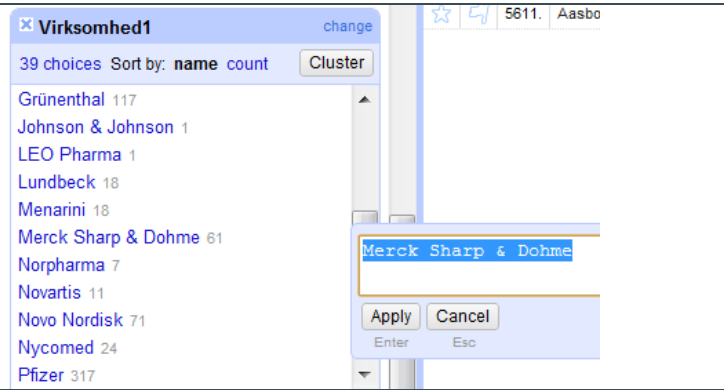| WHAT | WHY | HOW |
|---|---|---|
| **Trim cells** | Refine contains a wealth of cleaning methods . We begin to remove spaces at the beginning and end of cells:  Note also that we are working with Donor Name1 , while Donor Name is unchanged throughout the process. Google Refine also contains a good log system, so the process can be rolled back, yet it is wise to preserve the original content in a column. | Click the dropdown ▼ next to the heading of the column you want to clean. |
| **Facet** | Facet is the name of the main cleaning method. Select Text Facet, which offers a variety of mathematical methods to compare text. When you click , you get the following image on the left side of the screen:  Out of 2,126 rows are donor name in 1,503 different versions. | Select Facet - Text facet Click Cluster |
| **Cluster** | Under Cluster you are offered various methods to compare and link the content of the cells. You can choose everything at once and deselect one by | Select method and function Select dovetailing Click Merge Selected & Re - |

MEMO

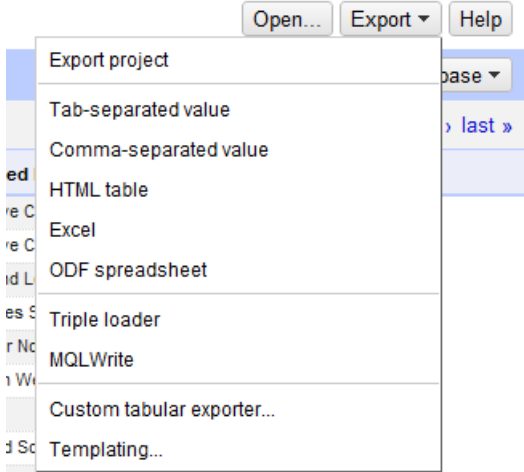| | | Cluster |
|---|---|---|
| | one. Or you can select it one by one.<br>There are two basic methods :<br>Key collision<br>Nearest neighbor<br><br>**Cluster & Edit column "Donor Name"**<br><br>This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer<br><br>Method [ key collision ▼ ]   Keying Function [ fingerprint ▼ ]<br><br>| Cluster Size | Row Count | Values in Cluster | Merge? | New Cell Value |<br>| 2 | 3 | • Mr Michael Freeman (2 rows)<br>• Mr Michael Freeman (1 rows) | ☐ | Mr Michael Freeman |<br>| 2 | 2 | • Mr John Robinson (1 rows)<br>• Mr John Robinson (1 rows) | ☐ | Mr John  Robinson |<br>| 2 | 8 | • Mr Christopher Butler (7 rows)<br>• Mr Christopher Butler (1 rows) | ☐ | Mr Christopher Butler |<br>| 2 | 5 | • Professor John Howson (3 rows)<br>• Professor John Howson (2 rows) | ☐ | Professor John Howson |<br>| 2 | 3 | • Mr Theodore Agnew (2 rows)<br>• Mr Theodore Agnew (1 rows) | ☐ | Mr Theodore  Agnew |<br>| 2 | 2 | • Mr Benjamin J Goldsmith (1 rows)<br>• Mr Benjamin J Goldsmith (1 rows) | ☐ | Mr Benjamin J  Goldsmith |<br>| 2 | 3 | • Lady Mary K Fuller (2 rows)<br>• Lady Mary K Fuller (1 rows) | ☐ | Lady Mary K Fuller |<br><br>[ Select All ] [ Deselect All ]                                [ Merge Selected & Re-Cluster ] | |
| **Key collision** | Key collision looking for convergence of important words. It is based on the idea of bringing together the core content of a text string. The method is subdivided into these three :<br><br>**Fingerprint:** Looking for identical characters and match hereinafter "John Smith" to "Smith, John". This method removes all special characters, set everything to lowercase, splits in words separated by spaces and removes duplicates. It also converts special versions of letters to their ascii-based value, ie ö becomes o. Fingerprint find the various ways in which the significant words are put together.<br><br>**N- gram Fingerprint :** Using the same principles as fingerprint , but allows variation of the word. N stands for the number of variations , permitting . That is, the higher n is , the greater is the difference between the words in the match . Even high N gives relatively few false. Try to vary.<br>1 - gram allows, for example, match "Krzysztof" ,"Kryzysztof" and "Krzystof" as identical.<br>N-grams find many typos.<br><br>**Phonetic Fingerprint :** Looking for the same sounds and finds, for example, "Horowitz" and "Horowicz". **Metaphone3** used for English, while **Cologne-phonetic** is German sounds. | **See more here:**<br>https://github.com/OpenRefine/OpenRefine/wiki/Clustering-In-Depth |

MEMO

| | | |
|---|---|---|
| | Normally , you start to cluster with key collision and takes these methods from the top down. | |
| **Nearest neighbor** | Nearest neighbor is the second basic method. While key collisions is very fast and well defined, it can sometimes be either too strict or loose and end up giving you much manual work to clean. Here it is an idea to try to finetune the selection of nearest neighbor that makes it possible to provide the parameters (like you can in n- gram). **Levenshtein:** The method looks at how many edits that is needed to get two strings match. "Paris" and "paris" has an edit distance of 1 because P must be changed to p "New York" and "newyork" has an edit distance of 3 In Refine you can define both the radius (distance) and Block Chars. Distance is absolutely central to work with. The higher the more difference there may be in text blocks while they are still perceived as equal. **PPM ( Prediction by Partial Matching):** The method examines how much difference there is in two strings. It looks like Levenshtein , but uses a different mathematical model. With Radius and Block Chars you can ask about how much difference you will accept and still consider it as a match . Matches typically "Johnson" and "Johnsons" . | |
| **End Cluster** | Once you are happy with your match, terminate the cluster. | Click Close |
| **Edit manually** | Although there are many algorithms to find duplicate names, so they will not find everything. Merck Sharp & Dohme may for example be as abbreviated MSD . It must be edited manually by simply typing in the cell. | |

MEMO



| | |
|---|---|
| Calculation | Refine has its own editing language. This can be used for calculations . |
| | Often it is easier to export to Excel and carry out the calculations here. |

MEMO

| WHAT | WHY | HOW |
|---|---|---|
| **Export** | After cleaning of data, you would like to continue to analyze them in Excel.<br><br><br><br>Choose Excel. Google Refine sends automatical an Excel-version to the browser download. Then you can work with it in Excel. | Click on Export<br>Choose Excel |

Exercise

| WHAT | WHY | HOW |
|---|---|---|
| **Clean titles** | Download file. It might open directly in Excel, then save it.<br>Import file to Open Refine. Clean by text facet, and check the result. Export the data to Excel. | data.kaasogmulvad.dk/unv/2016/refine/prof.csv |
| **Clean company spellings** | Prepare data in Excel. Add an extra column for the data you want to clean – not to spoil your original data. Import Excel to Open Refine. Work on the copied column and do text facet, and check the result. Export the data to Excel. | data.kaasogmulvad.dk/unv/2016/refine/defendants.xlsx |
| **Clean data on doctors** | In Denmark Doctors shall provide information on sidelines. They spell the name of the medical companies and their own names in very different ways. You need a clean version. | |
| **Make data ready in Excel** | Remember to copy and add a column for the data to be cleaned, Company, so the data are included in the spreadsheet twice. Second column can be called Company1. | data.kaasogmulvad.dk/unv/2016/refine/doctors20160919.xlsx<br><br>Original data:<br>http://ext.laegemiddelstyrelsen.dk/tilladelselaegertandlaeger/tilladelse_laeger_tandlaeger_full_soeg.asp?vis=hele |
| **Clean data** | Clean the column, Company1. Consider how it shall be cleaned out of a journalistic purpose. | |
| **Eksport data to Excel** | | |
| **Group data** | In Excel – make a list of the companies that have most doctors to work for them in Denmark. | |