Homework 2

1. Does OpenRefine alter the raw data during sorting and filtering?

It changes only the raw data you retrieves in the program, not in the file save on your desktop. You can save a new edition, with the refines.

2. 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?"

Step 1; Find a column with the right information 'Months_no_water' (I belive).

Step 2; Was to use regex, to replace the squarebrackets, citationssigns, and spaces, with nothing

value.replace("[","") value.replace("[","") value.replace("[",",","")

Step 3; After this, we entered custom text facet, and made the list split on semicolon

Value.Split(";")



October and September is the driest months.

3. What were the 10 most frequent occupations (erhverv) among unmarried men and women in 1801 Aarhus?

Step 1; So I simply just toke 'cevilstand', made a text facet, and chose "ugift"

Step 2; toke 'amt', made a text facet, and chose "Århus"

Step 3; Then I clustered as many as I could, with my knowledge, and the ones that seemed to be alike.

Step 4; Then I chose count, to find the top 10

1. Nationalsoldat
2. Soldat ved 1. Jyske inf reg.
3. Tienestepige
4. Tienestekarl
5. landsoldat
6. læredreng
7. hospitalslem
8. tienestedræng
9. gårdskarl
10. tjener faderen

4; maybe if you clustered more occupations, the list would have different. However this would be a matter of judgement.