

Visualization Project: Tag Refinery

Christoph Kralj*

University of Vienna

1 MOTIVATION

In today's world, tags are everywhere. Movies, music, papers and many more things have tags attached. Those should help us find what we want or help us finding new interesting things. I distinguish between lectured and social tags. Lectured are added by experts in the given field and usually of high quality but very reduced and maybe biased. Social tags on the other hand are very diverse and rich on information. This comes with the drawbacks of: spelling error, context specific language, non-informative words and many more problems. Tag Refinery focuses on social tags and tackles there problems by providing a visual interface for the underlying data wrangling algorithm. The final goal is a cleaned and reduced tag space in respect to the user needs.

1.1 Tasks

The main target is a visual interface which let the user easily manipulate the cleaning pipeline and gives instant feedback about the changes. I try to achieve this by solve the following tasks:

Get an overview of the tags Exploring the tag space (Filtering, zooming) Set the parameters of the underlying algorithm (Interaction) Reflecting over the changes with help of graphs (Brushing & linking)

1.2 Users

The user is every person which likes to make use of social tags. An example would be creating a music recommendation app based on emotions. In this case the user can clean the tag space in a way that emotional words are improved.

1.3 Algorithm

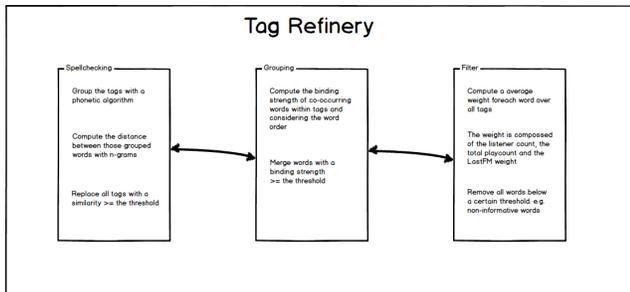


Figure 1: Algorithm overview

The underlying algorithm is part of my master thesis. The overall concept is shown in Figure 1. It consists of three modules: Spellchecking, grouping and filtering. The spellchecking module groups the tags with help of a phonetic algorithm and replaces within each group all words above a given similarity with the highest listeners count word in that group. Grouping is achieved by

*e-mail: christoph.kralj@gmail.com

finding co-occurring words and grouping them in respect to their relative frequency. The last module computes an importance score for each tag by using a weighted-average mean of the Playcount, Listeners and the LastFMWeight. This importance score is used to filter non-informative tags by removing all tags below a certain threshold.

1.4 Data

I use music tags from a subset of the million songs database¹ as example dataset. The tags are mined from LastFM². The format of the csv file looks like this:

SongID,SongName,Listeners,Playcount,TagID,TagName,TagWeight

2 APPROACH

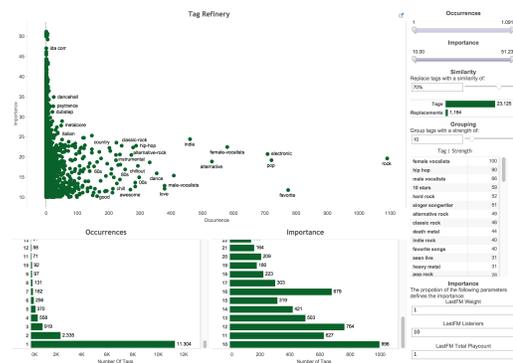


Figure 2: Final interface

Figure 2 shows the final interface. On the top left corner is the scatter plot (x-axis: occurrences; y-axis: importance score) which let the user explore the tag space by using the two filter sliders on the top right corner. The two histograms on the bottom side show the distribution of the the tag occurrences and the importance score. On the right side below the two filter sliders is the control element for the spellchecking module. The user chooses at which percent similarity words will be replaced. Directly below the element is a graph with two bars which show the total replacements in respect to the total tag count. This should give the user a hint how strong the effect of the chosen parameter is. Below the bars is the next input element, this time for grouping. The user can select a binding-strength threshold which is used by the algorithm to decide which words should be grouped. A lower values means more words will be grouped. This can be seen in the list below that element which shows all groups and the corresponding binding-strength. In the lower right corner are the input elements for the importance calculation. Each input is used as weight for the weighted-average computation.

¹<http://labrosa.ee.columbia.edu/millionsong/>

²<http://www.last.fm/>

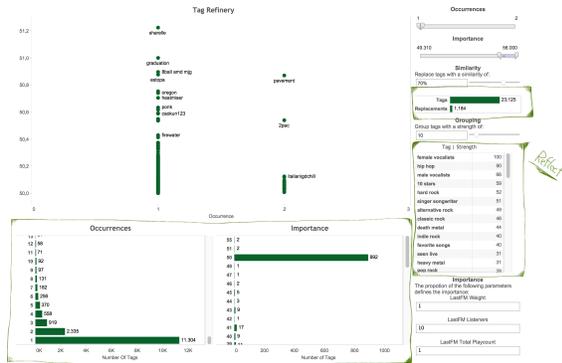


Figure 5: The graphs for reflection are highlighted in green.

5 DISCUSSION

The overall interface works in a nice and smooth way. I have removed everything which I think is not necessary for finding the best parameter combination for a given problem. Still there are a few points which can be improved. Especially the scatterplot which is good but not perfect may be a good starting point. Here are a list of strengths and weaknesses of the interface:

Strengths:

- Tag Refinery does not distract the user from his tasks by providing a simple and easy to use interface.
- It gives a good overview of the tag space and lets the user explore it interactively.
- This in combination helps the user find his optimal parameter combination.

Weaknesses:

- The scatterplot gets cluttered in the lower left corner.
- Importance calculation is hard to understand.
- All graphs show only overview.

5.1 Lessons learned

I learned a lot about how important it is how you name things to reduce ambiguity. Also I did not know the idea of mockups and they helped me a lot designing the first interface and getting my ideas onto paper. From this point I realised that fancy visualization techniques (heatmap, word cloud) should be used with care and at the end may be a simple visualization the better choice. From the technical side I am not aware of any better way of creating a visual prototype than Tableau. But on the other hand I think its not suitable for creating a end-user app. I think a good way of creating a visual tool is to create a prototype interface in tableau and implementing it afterwards in something like D3.

ACKNOWLEDGEMENTS

I wish to thank Torsten Miller, Moheesen Kamalza, the Vis class and my non-computer scientist users.