

Data visualization: Personal report

Joris LIMONIER

January 11, 2022

Contents

1	Sankey diagram by Joris LIMONIER	2
1.1	Introduction	2
1.1.1	Get the project	2
1.1.2	Background on Sankey diagrams	2
1.1.3	Code structure	2
1.2	Data pre-processing	2
1.2.1	Raw data	2
1.2.2	Clean the data	3
1.2.3	Transfer the data	4
1.2.4	Set up filtering	4
1.2.5	Group genres	4
2	Homework for 20/10/21	6
3	UX Protocol	7
3.1	Presentation & training	7
3.2	User test	8
3.3	Debriefing	8

1 Sankey diagram by Joris LIMONIER

1.1 Introduction

1.1.1 Get the project

My submission for the project of data visualization consists of a Sankey diagram. It is available [here](#)¹ and should be viewed in the browser (tested on Chrome), on a computer. The code for the project is posted [on GitHub](#).

1.1.2 Background on Sankey diagrams

Sankey diagrams are defined as follows: “Sankey diagrams are a type of flow diagram in which the width of the arrows is proportional to the flow rate.” [1] In our case, we have four columns that are linked by flows, representing the number of albums or the number of songs.

1.1.3 Code structure

The diagram was created with the following structure:

- Python [2] for data pre-processing (cleaning, formatting, writing).
- CSV for the data relative to each passage from one column to another (4 columns, therefore 3 CSV files).
- JSON to group the CSV data and put them in a format that can be fed to the D3.js library.
- JavaScript (Vanilla) to perform last-minute grouping and filter-relative updates
- D3.js [3] for the interactive data visualization.

Let us first detail the data preprocessing part.

1.2 Data pre-processing

1.2.1 Raw data

We grab the following raw data from the CSV Wasabi datasets that were downloaded.

- Album field
 - `_id` (album id)
 - `id.artist`

¹If the link is dead, go to <https://jorislimonier.github.io/>, navigate to “Projects”, then look for “Collaborative Data Visualization”

- genre
- Artist field
 - _id (artist id)
 - type (*e.g.* “Person”, “Orchestra”, “Group”, ...)
 - gender (male, female, unknown)
- Song field
 - id_album

1.2.2 Clean the data

The number of NaN values varies significantly from column to column. The Id columns don’t have any but some columns have many, making them sometimes tedious to work with. Overall however, the dataset is pretty clean. Barely a few Id’s are misformatted and some genres have unexpected characters.

Let us call “column transfer” (*CT*) the passage from one column to another. This is what we feed into D3.js and this is the part with the most impact on the diagram.

We clean the data in [sankey.py](#), which contains one class per *CT* :

- *TypeGender* : goes from the artist type to the gender of the artist.
- *GenderAlbums* : goes from the gender of the artist to the number of albums produced.
- *AlbumsSongs* : goes from the number of albums produced to the average number of songs per album.

The same file also contains the following class:

- *Sankey* : contains several functions to go from the individual CT classes to writing the data in the final format.

Each of the *CT* classes writes a CSV file. Usually, Sankey diagrams have three columns:

- source : the origin node for this flow
- target : the target node for this flow
- value : the flow quantity (*i.e.* the number of elements that go from a given source to a given target)

The *TypeGender* *CT* does indeed have these three columns. The *GenderAlbums* and *AlbumsSongs* *CT*’s on the other hand, are composed of a fourth column:

- genre : the genre of the album

This genre column allows to perform the filter operation.

1.2.3 Transfer the data

The *Sankey* class has a *write_final_data* method that takes the three CSV's and writes the format in a JSON file ([sankey-genre.json](#)), in the format demanded by D3.js. The format is as follows:

- nodes:
 - index
 - name
- links:
 - source
 - target
 - genre
 - value

where the set of the node indices is simply a bijection between the unique node names and $\mathbb{N}_{\geq 0}$.

1.2.4 Set up filtering

Filtering is performed in [sankey-filter.js](#). We allow for a filtering that only takes place for the two last *CT* 's (*i.e.* *GenderAlbums* and *AlbumsSongs*), because the *TypeGender CT* concerns artists, not albums (and obviously, artists don't have a genre).

We use JavaScript for the filtering part because we need to update the graph when genres are added/removed from the filter list. On each change of the selection dropbox (using an *eventListener*), we go through *sankey-genre.json* and check for each element in the "links" list if it is in the list of genres selected by the user. If so, we group it with all other elements going from the same source node to the same target node and sum their values.

1.2.5 Group genres

The grouping of genres is performed in [sankey.py](#). It contains a dictionary-like object that is of the following structure:

- keys: the new group name.
- values: a list of strings. If the old genre contains one of these strings, it will be assigned to the new group which is named after the *key*.

References

- [1] Sankey diagram. Page Version ID: 1062562883. URL: https://en.wikipedia.org/w/index.php?title=Sankey_diagram&oldid=1062562883.
- [2] Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009.
- [3] Michael Bostock, Vadim Ogievetsky, and Jeffrey Heer. D³: data-driven documents. 17(12).

2 Homework for 20/10/21

A paragraph describing the users.

Users could be people interested in the music industry, who want to compare men and women artists for sociological interpretation.

The list of visual tasks supported by users and the visualization goals.

User task	Details
Overview	Flow between columns
Zoom	TBD
Filter	Filter by genres

The list of (raw) attributes you will need from the WASABI dataset you are going to use.

The following attributes will be needed:

- Album field
 - `_id` (album id)
 - `id_artist`
- Artist field
 - `_id` (artist id)
 - `type` (*e.g.* “Person”, “Orchestra”, “Group”, “Choir”, “Other” or “”)
 - `gender` (male, female, unknown)
 - `members` (check this one, it may be the name of the members of a band)
- Song field
 - `id_album`
 - `genre`

The informal description of the processing of the raw data in order to make it to fit in the visualization technique. This might include calculated variables you must add in the process.

Clustering Artist - *type* variable to make *is_band* boolean variable

The name of visualization technique and the name of the member of the group who is going to implement it. Associate the visualization technique with the visual goal.

Sankey diagram with the following columns:

- Single artist or Band
- Male, Female, Unknown
- Number of Albums
- Number of Songs

A visual mapping of variables available in your data set (after data processing) and the visual variable available in the visualization technique you have chosen.

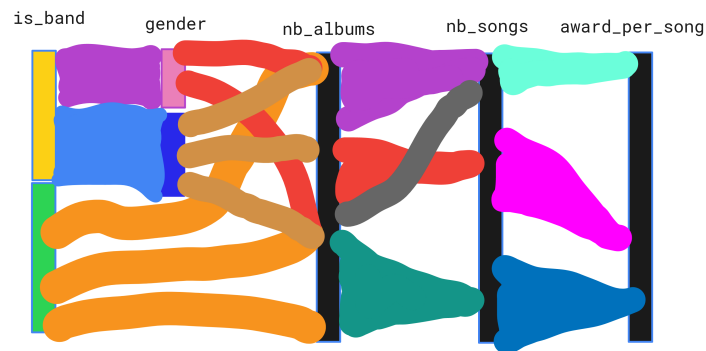


Figure 1: Representation of the Sankey diagram (JL)

3 UX Protocol

TODO:

1. Write text that will be read (to remove bias in the way I say things)
2. Test app on some people (minimum 3)

3.1 Presentation & training

1. Ask written consent for recording
2. Tell people why they are here and what we will do

3. Give the following information:
 - Age
 - Sex
 - Highest level of education reached
4. Present Sankey diagram and explain what a Sankey diagram is.
5. Is anything unclear?

3.2 User test

1. Task 1
 - (a) How many men made 5 albums.
 - (b) On a scale from 1 to 5 (1 means very easy, 5 means very difficult), how hard was that task?
2. Task 2
 - (a) How many rap solo artists are women.
 - (b) On a scale from 1 to 5 (1 means very easy, 5 means very difficult), how hard was that task?
3. Task 3
 - (a) Can you tell how many solo artists made 3 albums? Why/Why not?
 - (b) On a scale from 1 to 5 (1 means very easy, 5 means very difficult), how hard was that task?

3.3 Debriefing

1. What are your three favorite feature?
2. What is your three least favorite feature?
3. Would you recommend this application to a friend?
4. What would you do differently?