

Technical Blog Exploratory Analyst:

As we chose the dataset about books, I have faced several technical issues to resolve the problem definition. Before starting talking about the technical issues we will mention at first the resolution of the easiest one that could be implemented easy.

As I thought at first the data was already clean while trying to execute simple codes. I had to remove Null values in dataset books in order to compute results and getting an efficient answer. Most of the Null values come from column Year. While uploading to GitHub it doesn't show that which I couldn't resolve the issue. In the same column Year as it's type was float we needed to change it into integer in order to conduct a more realistic number through years rather than having 2009.0.

In the entire dataset we didn't don't have type of genre books. What I had to do is to create a function that combine dataset tags and add new column genre by merging on tag_id with book_tags and tags. In result it created a new dataset that is called main_tags. However, the function created wasn't accurate as 100% it only took the first word of tag_name. I have created a new dataset manually of most know genre that I called genres. What the function did is to match between genres and tag_name and choose the first word to place it in the new columns genre. If it couldn't find a match it returned 0. It also created a problem and duplicate in the genre between original_title as every goodreads_book_id has the can have the same title but not the same goodreads_book_id. In order to place it as worldcloud we had to remove the 0 and to stick with the duplicate.

To ease the process as time is limited for this project I have chosen a simple code that can be technical and business response to this dataset by mentioning how many books in the dataset and providing the min rating, avg rating and max rating for this book dataset.