

## Group Project Diary

| Assigned To     | Tasks            | Description  |
|-----------------|------------------|--|
| Ahn<br>Joo-Hyun | Topic discussion | - Suggested a topic which are airline delay data set, crash accidents data sets by airport and airlines.   |
|                 | Data wrangling   | - Replaced or deleted some dirty data what we do not use in the future in Julia.<br>- Bound every the-numbers dataset (7 files with moviedl after using API for dataset from numbers.com) to do join with IMDB data in R.<br>- Joined data between IMDB dataset (basic, rating, crew) in R.<br>- Created tidy version CSV files to the next step in Julia.<br>- Create new variables from given columns for analysing movie data such as counting words in the title, calculating genre complexity, Rescaling a financial data for easy to read. |
|                 | Plotting         | - Made graphs for visualisation in Julia.<br>- Analysing the relationship between some factors for example budget, words count, rating and revenue   |
|                 | Presentation     | - Presented visualization part of presentation and conclusion.   |
|                 | Organisation     | - Created a google drive and diary for the group to share ideas and relevant materials.<br>- Draw an overall outline how to proceed to the group project.<br>- Coordinated with a team and conducted regular meetings.   |
|                 | Prototyping      | - Active discussion. Offered ideas about how to combine two different datasets and visualization.  |
|                 | Report           | - Wrote report part related to plots, visualization, conclusion.<br>- Reviewed codes regarding the wrangling and plotting in Julia part.<br>- Summarised R and Julias codes for wrangling and visualisation and made comments for each steps.<br>- Updated project diary.<br>- Consolidated final codes in R and Julia by collecting each parts from group members.  |
| Jia<br>Zhiying  | Topic discussion | - Offered two different potential topics, one is to combine weather information with traffic data; another is to scrape and analysis the replies of a popular post.  |
|                 | Data wrangling   | - Optimized the data model from the perspective of tidy dataset and suggest to delete columns with the similar attributes.   |
|                 | Web scrapping    | - Completed web scraping of box office data from boxofficemojo.com, but the amount of data is less, only 790.<br>- The project team decided to scrape more box office data from the-numbers.com.<br>- Worked with Moon to figure out how to get the changed page links.  |
|                 | Prototyping      | - Active discussion, offered ideas about visualization.  |
|                 | Presentation     | - Prepared the presentation and slides about the goals and data sources.   |
|                 | Report           | - Wrote report part related to goals and data sources.<br>- Wrote the project report, consolidation in one file of all parts written by members.<br>- Updated project diary.   |

## Group Project Diary

| Assigned To       | Tasks            | Description   |
|-------------------|------------------|---|
| Changmin Moon     | Topic discussion | <ul style="list-style-type: none"> <li>- Presented a topic about movie datasets which can lead to prediction of specific upcoming movies.</li> <li>- Regarding the topic, found two main data set sources: imdb.com and the-numbers.com.</li> </ul>   |
|                   | Web scrapping    | <ul style="list-style-type: none"> <li>- Scraped the financial information from the-numbers.com and managed to collect the data about 5786 movies as a single data frame.</li> </ul>  |
|                   | Data wrangling   | <ul style="list-style-type: none"> <li>- Wrangled the movie data collected from the-numbers.com such as omitting unnecessary units.</li> <li>- Wrangled date column into 'yyyy-mm-dd' format - Once stuck in wrangling date column but found the way by using 'ludridate' package.</li> </ul>   |
|                   | Prototyping      | <ul style="list-style-type: none"> <li>- Active discussion. Offered ideas about what unique key should be and what to analyse and present.</li> </ul>   |
|                   | Presentation     | <ul style="list-style-type: none"> <li>- Edited and collected whole PPT material for the presentation. Made the beginning part of presentation.</li> <li>- Presented introduction of the presentation.</li> </ul>   |
|                   | Report           | <ul style="list-style-type: none"> <li>- Summarised R codes for web-scrapping, and made comments for each steps.</li> <li>- Wrote report part related to introduction, web-scrapping, and wrangling about the-numbers.com data parts.</li> <li>- Summarised the whole group project into a brief txt file to publish the group project to Github.</li> <li>- Uploaded relevant files of the group project to Github, and shared with other members how to work.</li> <li>- Edited the group project report and added supplementary explanation. And updated project diary.</li> </ul> |
|                   | Organisation     | <ul style="list-style-type: none"> <li>- Proposed guidelines for presentation and defined progress.</li> </ul>  |
| Komlev Viacheslav | Topic discussion | <ul style="list-style-type: none"> <li>- Offered some topics: Quality of life by cities, Quality of life by countries, energy production and contribution in NZ or worldwide in relation with CO2 emission.</li> <li>- Found data with movies revenues, links saved in common folder on google drive, file DataSetLinks.doc.</li> </ul>   |
|                   | Data wrangling   | <ul style="list-style-type: none"> <li>- Uploaded some tsv files (title_basic, title rating) from IMDB and found the way to read data without errors.</li> <li>- Made some wrangling (filter by year, excluding values which are not used in the project).</li> <li>- Wrote a function which mutate a new column with imdbID.</li> <li>- Created new title_genres table (save data into files title_genres.csv, title_rating.csv).</li> </ul>   |
|                   | Using API        | <ul style="list-style-type: none"> <li>- Wrote an API function which gets movieID from omdbapi.com for financial data from the-numbers.com.</li> <li>- Created several CSV files for further processing by using API function.</li> </ul>   |
|                   | Prototyping      | <ul style="list-style-type: none"> <li>- Active discussion, offered ideas about visualization, API function for joining data.</li> </ul>  |
|                   | Presentation     | <ul style="list-style-type: none"> <li>- Preparation for presentation, data model, wrangling, and API function parts.</li> </ul>  |
|                   | Report           | <ul style="list-style-type: none"> <li>- Added comments regarding the codes for API function and relevant wrangling in R.</li> <li>- Project diary aggregation.</li> <li>- Updated project diary.</li> <li>- Wrote about parts regarding data model, project steps, data wrangling techniques, API function.</li> </ul>   |
|                   | Organisation     | <ul style="list-style-type: none"> <li>- Arising questions about necessity of discussion related to project. (solving problems, topic searching).</li> </ul>  |