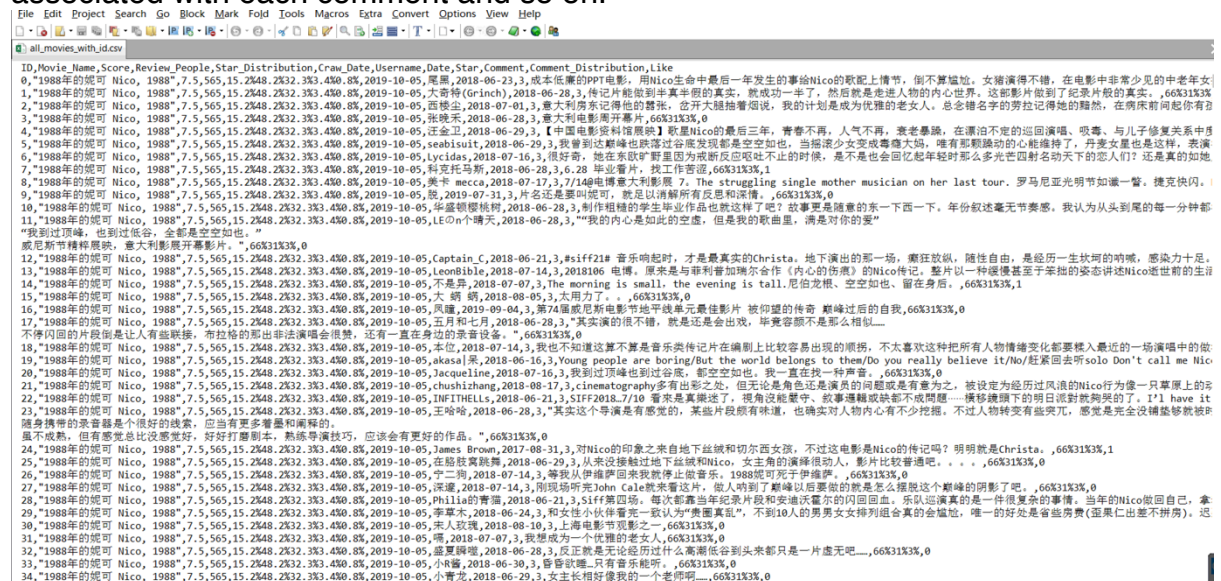# Status report

Group: DATA SMILE

- ## What we have done:

**Crawled** around 15000 movies and finally got 10269 unique movies after processing. All of the movie information and the corresponding comments have been written into one csv file. We upload this file to kaggle.com and it is now publicly available at: _https://www.kaggle.com/liujt14/dou-ban-movie-short-comments-10377movies_.   Each record in the dataset has 12 attributes, including movie name, overall score, stars associated with each comment and so on.



_All_movies_with_id.csv, 2.14GB_

**Processed dataset**. We have used common Chinese stopwords and punctuations list to remove less useful elements.  We utilize jieba segmentation tool to split each review as a list of words and find different parts of a Chinese sentence. With the calculation of TF-IDF, we found the most important words in one sentence and analyse its sentiment representation for the whole comment.

**Model Design.** The design of the overall framework is still in progress. Currently, we are discussing how to model the review word sequence. Typically, for a movie review containing several review sentences, there are two possible solutions: (1) regard the whole document as a long word sequence; (2) reserve the sentence boundaries and capture the hierarchical word-sentence & sentence-document relation. Also, in order to find more effective feature extractor, we are now comparing the performance of transformer with that of CNN/RNN on review representation learning.

- **What we will do:**

**Emotion analysis modelling**. Assume that there are few useless comments, which are issued by Internet marketers, inside Top 250 movies of Douban, we train one uniform-helpfulness model to analysis emotion of comments. The training dataset testing dataset are all Top 250 movies. The reference values will be the rank of each comment and the final rank. K-means clustering and dimension reduction algorithm(SVD) will be utilized to cluster movies according to reviewers' emotions per dimension. The main feature, or the bright spot of each movie, will also be presented in our output results.

**Helpfulness evaluation modelling** via RNN and CNN. After obtaining a sequence of review representations from our enhanced self-attention networks, we can borrow the recurrent neural networks (RNN) and convolutional neural networks (CNN) to calculate the helpfulness-aware overall representations. The training dataset and testing dataset include all movies information, i.e. ratings, reviews, descriptions and audiences, meanwhile, the related information with the long comments or from moviegoer reviewers will be considered as major reference for the evaluation of helpfulness. The useless comments will be given a lower weight and one new rank will be outputted. It is expected that this model can illustrate the "reality" of the rank, comparing the output of these two self-built models. The reality of Top 250 movies should approach 1 according to our assumptions.

**Baseline model testing.** Compare our emotion analysis results with other deep learning-based models and traditional models.

**Online real-time testing with new movies.** This is to demonstrate our ability to process big data with high velocity. Data streaming algorithm could be applied to analyse the new movie quality and propagation trend with the trained model. Ranks based on uniform-helpfulness model and helpfulness evaluation model will be given, compared to the final rank given by Douban.

**Data visualization.** The movies emotion map will be illustrated as a special heatmap to visualize the aggregated per-dimension emotion distribution. Clustering movies based on evoked emotion and emotion vector similarities might be also demonstrated by graphs. Besides, the estimated ratings, keywords and evaluation calculated by our models will be shown up to compare with original information.