

# README

## Environment

1. python 3.6
2. python packages:  
NLTK, re, pandas, json.
3. java package:  
stanford-corenlp-full-2017-06-09/stanford-corenlp-3.8.0.jar,  
stanford-corenlp-full-2017-06-09/stanford-corenlp-3.8.0-models.jar

## Code Explanation

As stated in report, the whole process includes data cleaning, target extraction

## Data Cleaning

the codes are sotred in directory *data\_process*

1. *clean\_data.py* : this code is used to remove meaningless marks  
input: *raw\_data.csv*, output: *step1\_data.csv*
2. *remove\_non\_english\_sents.py* : this code is used to remove reviews not written in English  
input: *step1\_data.csv*, output: *step2\_data.csv*
3. *groupData\_base\_on\_id.py*: group data based on course id.  
input: *step2\_data.csv*. output: 36 course review files

All data used in above are stored in directory *data*, and 36 course review files are stored in the subdirectory *course* under *data*. These 36 files are named in form of *course\_"course\_id".csv*

## Target Extraction

The codes are stored in directory *extract\_target*.

1. *extract\_target\_list.py*: this code is to grow target and opinion list by propagation  
input: *course\_"course\_id".csv*, output: *course\_"course\_id"\_target\_list.txt*

2. *extract\_target\_phrase.py*: this code is to extract target word or phrase for each review.

input: *course\_"course\_id".csv*, *course\_"course\_id"\_target\_list.txt*

output: *course\_"course\_id"\_transaction.csv*

3. *filter\_target.py*: this code is used to prune the targets.

input: *course\_"course\_id"\_transaction.csv*

output: *course\_"course\_id"\_target\_filtered.txt*

All output files produced by this part are stored under the directory *result*.

## **Result**

Final result stored in the files names as *course\_"course\_id"\_target\_filtered.txt*.