

## Assignment# 2: Pre-processing

Due: October 30, 2021, in Blackboard

Given a text file (input.txt), write a program (in your language of choice) to pre-process the file. The program should

- Tokenize the lines of the file based on white-spaces.
- Case-fold each token to lowercase.
- Identify the stop-words. Use the stop words from the stop\_words.txt file.
- Remove all punctuation.
- Use Porter stemmer\* to stem the tokens.
- Write each token on a separate line in the output file.
- Write the type of token - for all non-empty token, next to the token, in the same line of the output file, write 'article' if the token is an article (i.e., a, an, the), 'stop word' if the token is a stop word, otherwise, write the *stemmed* version of the token.

Sample input file:

The less there is to justify a traditional custom, the harder it is to get rid of it.

Sample output file:

the	- article
less	- less
there	- stop word
is	- stop word
to	- stop word
justify	- justifi
a	- article
traditional	- tradit
custom	- custom
the	- article
harder	- harder
it	- stop word
is	- stop word
to	- stop word
get	- get
rid	- rid
of	- stop word
it	- stop word

**Deliverables:** 2 files

- The source code of the program. The name of the file should be lastname\_firstname.xyz (replace xyz with proper extension).
- The output txt file. The filename should be lastname\_firstname.txt

\*<https://tartarus.org/martin/PorterStemmer/>