# CRC cards

| WebScraping | |
|---|---|
| Instance variables: referrer (http://google.com) userAgent (Mozilla) | Collaborators: NLP Analysis (feed input) |
| Methods: sanitizeURL | Make sure URL is in proper form for feeding into JSOUP parser |
| getConnection | Establish connection to website, throw exceptions as necessary if unable; return html if connection |
| cleanContent | Get only relevant html tags, clean content so that content cannot be malicious to application |
| readHTML | Read in relevant tags to string array for output |
| Runner | Solicit user input, run methods on input url, return for passing into NLP |

| Model Training | |
|---|---|
| Instance variables: *still learning how this works; will train for keyword and topic tagging on a dataset 1. Acquire data 2. Clean data 3. Train the model 4. Evaluate the model | Collaborators: Dataset NLP Analysis |

| NLPAnalysis | |
|---|---|
| Instance variables:<br>userWords (string of words from user's URL)<br><br>InputStream tokenModelIn<br><br>InputStream posModelIn<br><br>InputStream dictLemmatizer (?)<br><br><br>tokensArray (string array of tokens)<br><br>tagsArray (string array of tags)<br><br>HashMap<String, String><br>tokenToPOSTagMap<br><br>HashMap<String, Integer><br>tokenToCountMap | Collaborators:<br>WebScraping (output)<br><br>File: token model (pre-trained)<br><br>File: pos model (pre-trained)<br><br>File: a lemma dictionary (from OpenNLP)<br>-likely needed for KeywordAnalysis<br><br><br><br><br><br><br><br><br>Sentiment Analysis (input)<br>KeywordAnalysis |
| Methods:<br>Constructor - takes in string from web scraping output (userWords)<br><br>Lematize (helper)<br>Return tokensArray<br><br>createTokenToPOSTagMap<br>Return tokenToPOSTagMap<br><br><br>createTokenToCountMap<br>Return TokenToCountMap<br><br><br><br><br>nlpLemmatize<br>Return lemmaArray | Responsibilities:<br>Creates the NLPAnalysis object<br><br><br>Lematizes the string input from the user<br><br><br>Tokenizes, POS tags, stores those key-value pairs in hashmap<br><br><br>Put all words (minus function words) and their frequency of occurrence (from user's URL) into a hashmap for the sentiment analysis<br><br>Create string array of lemma of each word |

| | |
|---|---|
| getTokenToCountMap | Getter for Sentiment Analysis and Keyword Analysis to access |

| Keyword Analysis | |
|---|---|
| Instance variables:<br>keywordArray<br><br>InputStream keywordModelIn | Collaborators:<br>NLPAnalysis (input)<br>User Interface (output)<br><br>File: text keyword model (from our Model Training) |
| Methods:<br>Constructor takes a String of the words from the user's URL (userWords)<br><br>nlpKeywordTag<br>Return keywordArray | Responsibilities:<br>WebScraping (input)<br><br><br>Use our trained model to find keywords and thus identify the topic(s) of the text |

| Sentiment Analysis | |
|---|---|
| Instance variables:<br>negativeWordCount<br>positiveWordCount<br>scoreOutput | Collaborators:<br>NLP Analysis (input) |
| Methods:<br><br>createDictionary<br><br><br>wordCounter | Responsibilities:<br><br>Create a dictionary of positive and negative words<br><br>Count the positive and negative words in the text |

| | |
|---|---|
| scoreCounter | Label the positive words with grade 1 and the negative words with grade -1 |
| scoreDisplay → go to user interface?? | Output a positivity score of the text |

| User Interface | |
|---|---|
| Instance variables: | Collaborators:<br>NLP Analysis & Sentiment Analysis |
| Methods:<br>Simple GUI? Print a file? Console? | Responsibilities:<br>Returns to user the topic (several keywords) of the text from their URL input<br><br>Optional (if time): recommend other URLs based on that |