

Text mining steps

Step one	<p>Self-defined goals will be formatted as <i>comma-separated values</i>. This means, first self-defined goals will be parsed into their separate components using tokenizer algorithms as part of <i>nltk</i> library. For example, if a person writes “Walk outside 3-4 x a week in the morning for 15-20 min”, the sentence will be separated into tokens (words) as below:</p> <pre>['Walk', 'outside', '3-4', 'x', 'week', 'in', 'the', 'morning', 'for', '15-20', 'min']</pre> <p>Using <i>pandas</i> library, this information will be transferred into an in-memory data structures.</p>
Step two	<p>Having the goals separated into tokens will give the opportunity to: (i) identify the role of each word in the sentence; and (ii) evaluate if the word is part of the goal setting lexicon. Parts of speech (POS) tagging functions such as <i>pos_tag()</i> will be used to tag each token into their corresponding role (class). Each class has an associated code. To illustrate how the results of this step will look, the <i>pos_tag()</i> algorithm was applied to the goal example above:</p> <pre>[('Walk', 'VB'), ('outside', 'IN'), ('3-4', 'JJ'), ('x', 'JJ'), ('week', 'NN'), ('in', 'IN'), ('the', 'DT'), ('morning', 'NN'), ('for', 'IN'), ('15-20', 'JJ'), ('min', 'NN')]</pre> <p>However, the POS tagger does not always give the best results and may contain inaccurate results. To deal with this inaccuracy, different tagging algorithms available from <i>nltk</i> library¹ will be tried and the one with the highest accuracy will be selected.</p>
Step three	<p>At this stage, all the tokens in persons’ goals have been tagged. The next step is to evaluate tokens based on the previously gathered lexicon. Regular expression algorithms will be used to match tokens to the initial lexicon. If a match was found, that word is classified as a “matched” word. For example, in the case of our example above:</p> <pre>line = Walk outside 3-4 x a week in the morning for 15-20 min search = re.search(r'walk', line, re.M re.I) if search: print “Found!”</pre>

¹ <https://www.nltk.org/api/nltk.tag.html>