

Intermediate Report

Project 1 for CS421- Natural Language Processing

Shobhit Lamba - slamba4@uic.edu

Mohit Adwani - madwan2@uic.edu

For Essay Autograder, we had to dive into various techniques learnt in class and implement most, if not all, of them to grade different aspects of the an essay. Here are the different techniques used for each of the four tasks:

1. Counting number of sentences

We took a multi-level approach to count the length of an essay:

- First, we get the sentences by applying nltk sentence tokenizer on the essay. Then we split all sentences which have a newline character (\n).
- We saw a pattern in the essay in which there wasn't a space after period (which denotes the end of a sentence) and hence the sentence tokenizer didn't split them into 2 sentences. Hence, we split the sentences into two if the character after a period is alpha and before the period isn't period(.) since some essays had two or more consecutive periods to denote continuation. Also, the sentence after the period should at least have 3 characters.
- While looking at multiple finite verbs in the sentence(hint given in project_part1.pdf), if the sentence didn't have coordinate or subordinate clause then we saw a pattern in the parse tree, denoting the finite verb phrase as '(SBAR (S' – where SBAR denotes 'clause introduced by a (possibly empty) subordinating conjunction'. If the sentence has a subordinating conjunction then it is denoted by '(SBAR (IN that) (S'.

2. Spelling mistakes

Again, we took a multilevel approach for counting number of spelling mistakes:

- First, we simply tried to find the spelling mistakes using wordnet. But wordnet falsely reported simple words like "how", "you" etc as wrong spellings. So we used it as our first pass of spellcheck and its output was used as an input for the next step.
- After that, I employed Peter Norvig's spell correction code to further filter out the words. Now to a great accuracy, we can say that the words left in the list are probably the only words that are incorrect in the essay.

3. Grammar

Grammar is the hardest task we had to tackle in part 1 of the project. Here, we had to handle the following:

- Subject-Verb Agreement
- Missing Verbs
- Incorrect Verbs
- Verb Tense

For this, we had to employ POS tagging. After every word of the essay was tagged, the task was subdivided as c.i and c.ii

For c.i (Subject-Verb Agreement), we defined certain rules that a subject and verb combination shouldn't follow, which were basically the singular and plural agreements. If the POS-tag sequences were found in the list of rules, error count is incremented.

For c.ii (Rest of the conditions), we again defined rules, but this time the list included rules that are correct. So, if the sequence does not match any of the defined rules, error count is incremented.

After the list of scores were generated for a, b, c.i and c.ii, we normalized it using the following equation-

$$x = (x - \min(x)) / (\max(x) - \min(x))$$

This gave us values normalized to a range of 0-1.

To get the spelling errors in range of 0-4, we multiplied every value with 4.

All the other scores were multiplied by 4 and subtracted from 5 to get the final scores in the range of 1-5. This ensured that lower the number of errors, higher the score obtained by the essay.

All values finally were rounded off to nearest integer value.

Final score was calculated for each essay using the given equation which is:

$$\text{Final score} = 2 * a - b + c.i + c.ii$$

Since we did not calculate c.iii, d.i and d.ii yet, they were given a value of zero and excluded from calculations.

The results are finally written into a text file.

★ Pattern of Errors

Every part of the project had some issues that we faced.

While calculating length of an essay, we found that our function wrongly separated sentences when a period was encountered with an abbreviation (like in “e.x.”). But since the number of such cases were minimal, we decided to ignore it. We also found that using capitalization as a characteristic for new line did not help. So we excluded that.

For part b, we found that wordnet had a very weird way of finding out existence of word in the corpus. It uses a function called “synset” which finds synonyms of words in the corpus. Since many words do not have synonyms, they were falsely marked as wrong spellings. That is why we implemented a second layer of scrutiny which employed probabilistic methods to find out if a word is misspelled. This really boosted our performance in part b.

For part c, the main problem was generalizing the code. Right now it is not that general, since english grammar has a lot of rules to define grammatical errors pertaining to verbs.. So, we had to separately jot down rules for both c.i and c.ii. Additionally, c.i does not have rules pertaining to adverbs so it might still have a certain degree of inaccuracy on the higher end. Since the nltk POS tagger sometimes doesn't tag the verbs correctly, some verb formations, even though correct, are tagged as incorrect.

★ Variations in Final Equation

The current final score equation is as follows (for part 1):

$$\text{Final score} = 2 * a - b + c.i + c.ii$$

Intuitively, we observed the following things:

1. Number of sentences should not get such a high weight in final score. An essay with 50 sentences might not be better than an essay with 25 sentences. It may so happen that someone just writes one sentence, “Pollution is bad.” over and over 50 times while someone else wrote an essay that was coherent and flowing but only had 25 sentences. Our autograder might wrongly grade the first essay higher with higher weight given to the number of sentences.

2. Spelling mistakes should have more weight proportionately. More number of spellings mean the writer is not well versed in writing in english and consequently writes poorly.
3. While grammar holds huge importance, our model for subject verb agreement only detected a maximum of 2 errors while most of the scores were zero. Scaling them to a scale of 1-5, zero errors got a 5, 1 error got a 3 while 2 got 1. But 2 errors is not that bad compared to 1 or zero, yet score difference is an anomaly. What if someone wrote an essay with an error in each line. He scores a 1 too, but in contrast, his essay was way more poorly written than the one with 2 errors.
4. Our scores for c.ii had a wide array of results and we believe c.ii should have more weight in the final score as well.

Running a linear regression analysis on the coefficients, we found that even though the coefficients had a varying result in each run, they stayed in the following range for most part:

Coefficient of a \rightarrow 0.40-0.45

Coefficient of b \rightarrow -(0.25-0.32)

Coefficient of c.i \rightarrow (0.07-0.10)

Coefficient of c.ii \rightarrow (0.34-0.38)

We believe the values will change with introduction of the other 3 variables in part 2 of the project.
