Spell Corrector

Egemen Göl 2016400009

Cmpe 493 - Introduction to Information Retrieval Assignment 1

Implementation

Learn

The spelling corrector starts by learning from corpus.txt which contains assumed-correct English text by tokenizing it with a regular expression and counting every token.

Then using the spell-errors.txt file, the program constructs a dictionary where on every misspelled word there, it notes the corrected word and increments this conversion's count. This is just like the corpus Counter object but it counts for every mis-typed word by itself.

Correct words from spell-errors are added to the corpus, in a weighted manner, because it is assumed that this file is important from the test cases. Current weight is 3, meaning every correct spell-error word is counted three times in the corpus.

Correct

Correction works by asking corpus and spell-errors objects for their correction suggestions, and taking the weighted best among them. Current weight favors spell-error by 30.

If the word exists in the corpus, no correction takes place.

Corpus calculates the probability in two ways.

- If simple mode is selected, probability is defined by (the count of the word in corpus) / (corpus size).
- If smooth mode is selected, probability is defined by (the count of the word in corpus
 + 1) / (corpus size + 1 * (nof unique words+1))

If there are multiple words with the same probability, for the sake of unbiasedness, the preferred word is selected randomly, in a pre-seeded manner.

1-edit distance is calculated by constructing deletes, transposes, replaces and inserts of every combination and then these are feeded to corpus probability function to calculate the respective frequencies.

Spell-error calculates the probability by choosing the most common corrections, randomly selects one from them, just like corpus.

If they have no suggestions, they return empty string with probability zero. This allows easy choosing between these two suggestions and returning empty string.

Measure

Taking measurements was the heaviest part of the project, since we decided to have measurement mechanism to act as if the correction system is a black box, and build statistics from there.

Two types of statistics are calculated:

- 1. Confusion matrices
- 2. Accuracy metrics.

Confusion Matrices

Statistics engine assumes the "correct correction" takes place by suggesting a "1-edit-distance" correction. If it is correct in this suggestion, with respect to the reference word, the respective confusion matrix is updated.

Matrix indexing works like following:

```
For REPLACE x -> y ==> (x, y)

For TRANSPOSE xy -> yx ==> (x, y)

For INSERT x -> xy ==> (x, y)

For DELETE xy -> x ==> (x, y)

For INSERT and DELETE, gives '_' for x if y is at the beginning of the word.
```

Then the confusion alphabet is '_abcdefghijklmnopqrstuvwxyz'

Matrice files are in .csv format, with index and header included.

Confusion matrix calculation, since it uses 1-distance editing assumption, only uses the 'spell-errors.txt' file as a corpus consisting of only the corrected words, rather than a dictionary of suggestions.

Accuracy Metrics

Accuracy metrics take 'spell-errors.txt' file into consideration, since it directly works on the output of the correction system and comparing it to the reference. It tries to find out the method used in correction, and it classifies this method into types:

- EDIT
- TABLE
- NO_OP

EDIT means 1-edit-distance correction, TABLE means the suggestions from the 'spell-errors.txt' are used, NO_OP means the word was already correct.

For measuring the 'spell-errors.txt' file, reference word is actually a set of words from the file, since it frequently contradicts itself. The set consists of the most-frequent suggestions from the file. If the system corrects the word into any of these reference words, it is counted as a correct correction.

Command-Line Usage

There are two python executables:

- 1. Correction executable 'main.py'
- 2. Measurement executable 'measure.py'

main.py

```
$ python3 main.py -corpus ./data/corpus.txt -spell-errors ./data/spell-errors.txt simple <
./data/test-words-misspelled.txt > output.txt
```

Uses stdin and stdout for correction. 'smooth' or 'simple' must be selected by the user.

`-corpus` and `-spell-errors` can be skipped since they default to `./data/corpus.txt` and `./data/spell-errors.txt` files, respectively.

```
$ python3 main.py smooth < ./data/test-words-misspelled.txt > output.txt
```

measure.py

```
$ python3 measure.py all
```

Produces confusion matrix and accuracy metric files into the ./measurements/ folder.

Takes a bit of time, ~15 seconds on the author's computer.

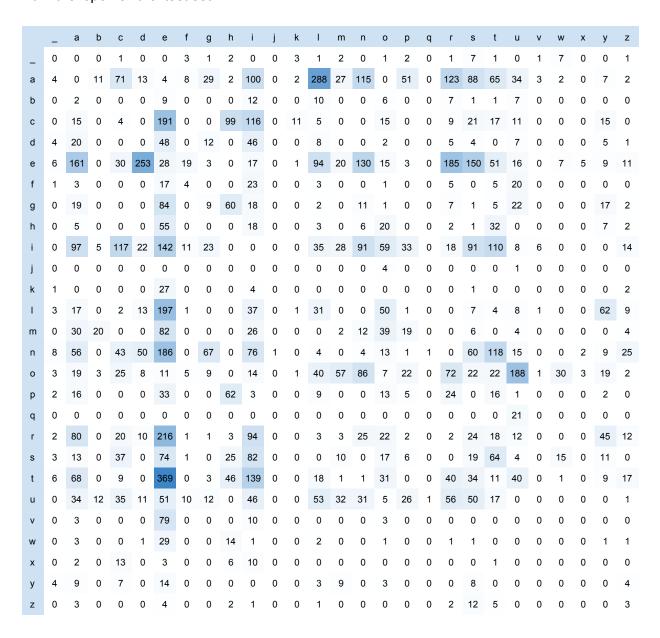
Measurement Outputs

Confusion Matrices

They can be obtained by running:

\$ python3 measure.py all

Since there are 2*2*4=16 tables, we will include only the insert matrix from simple insertion from the 'spell-errors' test set.



Accuracy Metrics

Test, Simple	Test, Smooth	Spell, Simple	Spell, Smooth
<pre>{ "true": { "edit": 353, "table": 12, "no_op": 0 }, "false": { "edit": 13, "table": 1, "no_op": 5 }, "null": 0 }</pre>	<pre>{ "true": { "edit": 353, "table": 12, "no_op": 0 }, "false": { "edit": 13, "table": 1, "no_op": 5 }, "null": 0 }</pre>	<pre>{ "true": { "edit": 22735, "table": 11176, "no_op": 17 }, "false": { "edit": 879, "table": 0, "no_op": 2316 }, "null": 0 }</pre>	<pre>{ "true": { "edit": 22735, "table": 11177, "no_op": 17 }, "false": { "edit": 878, "table": 0, "no_op": 2316 }, "null": 0 }</pre>
~95%	~95%	~91%	~91%

We can observe much from this table.

First, we need to take these percentages with a grain of salt since the system is trained over these datasets, a high percentage is expected.

Second, in our implementation, smoothing makes no difference at all. Even worsens the score in the spell-errors dataset.

"null" means "no-correction since I have no idea". Our program is trained over these datasets, it always has an idea.

Incorporating the 'spell-error' file made little difference, in the test_set the contribution is 12/384 == +3%, not that much. Since other is the training set itself, it is meaningless.

We did not optimize the weight we put into merging 'spell-errors' into corpus (3x), since it is not measured. The effect is indirect, hard to measure.

The 'no_op' errors in 'spell-errors' test is surprising, we are not sure why there are so much errors there. System just leaves the word as is, this may seem like stemming from the weighted merge operation, that move uses only the 'correct' or 'target' words from that file. Maybe this is because of the inconsistencies within the 'spell-error' file.

Maybe the biggest reason for this is the assumed-correct situation in which if the word exists in the corpus (and the corpus includes the target words from the 'spell-errors' file'), the system does not correct it. This is a choice of the authors, since user must be respected. Excluding this may improve our percentages, but this is good practice in our judgement.

Criticism and Possible Improvements

As mentioned above, weighted merge operation's effect is unknown.

Also the high error rate from no_op must be investigated further, it is thought that the inconsistency within the 'spell-error' file may be the culprit.

The system is slow, this may be improved by a re-write in a faster language.

Since 1-edit distance correction is the main requirement, other techniques are not discussed.

The keyboard layout for mobile and desktop, incurs some statistical bias to the errors, this may be exploited.

Using a bigger (and more consistent) training set would definitely help.

The choice of not correcting when the word is in the corpus, feels right by the user. There is no need to be a *know-all* in the eyes of the user.

In hindsight, I am satisfied with the approach I took, the architecture I chose and the implementation order I decided. If I would do this assignment once more, there would be few things that I would change. It was a fun and challenging project, especially in terms of testing. Thank you.