

Text Normalization

Section 1: Descriptions of the System

1.1: The design

The tweets are one kind of special text corpus where we may have many expressions which are not standard expression. The design of the system is meant to handle the following characteristics of tweets:

- The special sequences: urls, the topic tags, the user tags, concatenation of several words, etc.
- The misspelling (deliberated or not) of words

For the two characteristics, what I do is to build a system of 2 main steps:

1. Preprocessing: remove the urls, detect and mark the topic tags and user tag, split the concatenated sequences.
2. Misspelling correction: For the words that are not correct (not in the dictionary), correct it using both context information and formal similarity information.

1.2: The building process

The main idea is to firstly build the overall pipeline, then fill the detailed implementations:

1. Build the “preprocess.py” file.
2. Build the pipeline where “system.py”.
 1. It does nothing to the preprocessed sentences: it just take the input as output, the purpose is to check whether the bash file works well.
 2. Add the function to detect the positions in the sentences where the words are not in the dictionary.
 3. Add the function to get the context vector of the target positions.
 4. Get the candidates replacements using both vector similarity and Levenshtein distance.
 5. Add the ‘STATE’ parameter to specify the way to select the candidates (manual or auto).

1.3: The components

The system is implemented in python.

1. For the preprocessing part, I used one existing library ekphrasis (<https://github.com/cbaziotis/ekphrasis>), the code in “preprocess.py” is a modified version of <https://github.com/cbaziotis/ekphrasis/blob/master/ekphrasis/examples/example.py>.
2. For the misspelling correction part, I used the library context2vec (<https://github.com/orenmel/context2vec>), the code in “system.py” is developed by myself, but I referred to the https://github.com/orenmel/context2vec/blob/master/context2vec/train/train_context2vec.py.

Section 2: Situations Handling

2.1: Situations well handled

For the simple misspelling, such as “survivor —> survivar”, the system can easily detect and correct it. The system considered both context information and formal similarity information, so the recovery of simple misspelled words in a relatively clear context is good.

Examples:

1.

- Original sentence:

RT @muhdnuri: Israel killing Muslims everydwy and no one bats an eye. Terrorist attack and Muslims got the blame? How shallow can you be? #...

- Manual mode:

```
-----Sentence 1:-----
```

```
b'rt <user> : israel killing muslims everydwy and no one bats an eye . terrorist attack and muslims got the blame ? how shallow can you be ? # \xe2\x80\xa6'
```

```
>> Extracting the words to be corrected ... ..
```

```
>> The following positions are not correct:
```

```
[6]
```

```
>> For the correction of word: everydwy
```

```
>>>> Calculating the candidates ... ..
candidate 1:      everyday , similarity: 0.3986249
Time: 2.09s
Please choose the replacement, enter the number of the candidate: (Enter 0 for no replacement)
>> 1
>>>> everyday      ---> everyday
```

- Modified sentence:

b'rt <user> : israel killing muslims everyday and no one bats an eye . terrorist attack and muslims got the blame ? how shallow can you be ? # \xe2\x80\xa6'<user> <user> that bitch stops a show because someone spilled water on stage but puts on a show when terrorists attacki'b'rt <user> : french president francois holland condemns the " terrorist attacks of unprecedented proportions . " <hashtag> pray for paris </hashtag>'

2.

- Original sentence:

RT @ABSCBNNews: **Franch** President **Frandois** Hollande condemns the "terrorist attacks of unprecedented proportions." #PrayForParis <https://t.c...>

- Manual mode:

```
-----Sentence 2:-----
b'rt <user> : franch president frandois hollande condemns the " terrorist attacks of unprecedented proportions . " <hashtag> pray for paris </hashtag>'
```

```
>> Extracting the words to be corrected ... ..
>> The following positions are not correct:
[3, 5, 6]
```

```
>> For the correction of word: franch
>>>> Calculating the candidates ... ..
candidate 1:      french , similarity:    0.5013202
candidate 2:      france , similarity:    0.45431128
candidate 3:      branch , similarity:    0.43451628
candidate 4:      ranch , similarity:     0.43080345
candidate 5:      franca , similarity:     0.42238182
Time: 1.41s
Please choose the replacement, enter the number of the candidate: (Enter 0 for no replacement)
>> 1
>>>> franch        ---> french
```

```
>> For the correction of word: frandois
>>>> Calculating the candidates ... ..
candidate 1:      francois , similarity: 0.468642
Time: 1.97s
Please choose the replacement, enter the number of the candidate: (Enter 0 for no replacement)
>> 1
>>>> frandois     ---> francois
```

```
>> For the correction of word: hollande
>>>> Calculating the candidates ... ..
candidate 1:      holland , similarity:   0.48591468
candidate 2:      hollander , similarity: 0.45005512
candidate 3:      hollands , similarity:  0.41884422
Time: 1.91s
Please choose the replacement, enter the number of the candidate: (Enter 0 for no replacement)
>> 0
```

- Modified sentence:

b'rt <user> : french president francois hollande condemns the " terrorist attacks of unprecedented proportions . " <hashtag> pray for paris </hashtag>'

2.2: Situations not properly handled

1. In the system, I didn't consider the impact of the keyboard: 'q' is more likely to be misspelled to 's' than 'p' because 'q' is closer to 's' in the 'qwerty' keyboard. One possible solution is to attach different penalty in the computation of Levenshtein distance: larger penalty for the replacement between two 'far-away' letters.
2. In the system, I didn't consider the frequent situation where people use the repeated vowels. For example, "liiiiiiiiike" means "like", but in our system, this transformation will be counted 9 units of Levenshtein distance. One possible solution is try to "merge" the characters (in this case, 'i') which repeat too many times.
3. I didn't have enough time to train the context2vec model using the given corpus, if the model is trained, we can have better result.