自然語言處理 HW2

組員 :

110590450 歐佳昀(70%) => 作業架構、流程、分析、改善

110590452 莊于潔(30%) => 作業模型挑選、參數調整

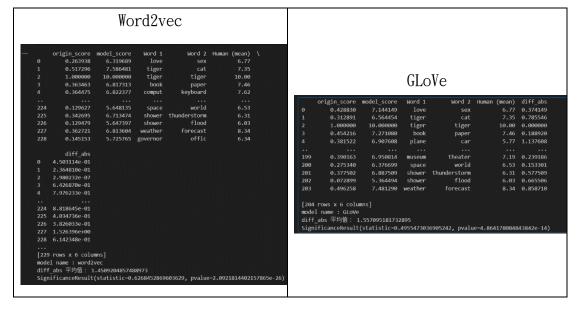| 環境 : 使用 Jupyter Notebook | 語言 : python |
|---|---|

- Results
- Word embeddings
- Word similarity
  - Correlation on WordSim353
- Analogy prediction
  - Accuracy for BATS categories

## Word embeddings



## Word similarity

WCorrelation on WordSim353

## Analogy prediction
### Accuracy for BATS categories

| GLoVe | Word2vec |
|---|---|
| ```
data/BATS\1_Inflectional_morphology :  category  [ noun - plural_reg ]
total accuracy    total :  35 ps :  35 lm :  35 data :  50
total :0.7
Stemming  :0.7
lemmatization  : 0.7
``` | ```
data/BATS\1_Inflectional_morphology :  category  [ noun - plural_reg ]
total accuracy    total :  32 ps :  35 lm :  32 data :  50
total :0.64
Stemming  :0.7
lemmatization  : 0.64
``` |
| ```
verb_ving ved verb_ving ved
data/BATS\1_Inflectional_morphology :  category  [ verb_ving - ved ]
total accuracy    total :  38 ps :  41 lm :  38 data :  50
total :0.76
Stemming  :0.82
lemmatization  : 0.76
``` | ```
49    data/BATS\1_Inflectional_morphology :  category  [ verb_ving - ved ]
50    total accuracy    total :  38 ps :  43 lm :  38 data :  50
51    total :0.76
52    Stemming  :0.86
53    lemmatization  : 0.76
``` |
| ```
139
140    country capital country capital
141    data/BATS\3_Encyclopedic_semantics :  category  [ country - capital ]
142    total accuracy    total :  46 ps :  43 lm :  46 data :  50
143    total :0.92
144    Stemming  :0.86
145    lemmatization  : 0.92
``` | ```
data/BATS\1_Inflectional_morphology :  category  [ verb_3psg - ved ]
total accuracy    total :  27 ps :  30 lm :  27 data :  50
total :0.54
Stemming  :0.6
lemmatization  : 0.54
``` |

### *SVD - word2vec

```
data/BATS\1_Inflectional_morphology :  category  [ noun - plural_reg ]
total accuracy    total :  11 ps :  11 lm :  11 data :  50
total :0.22
Stemming  :0.22
lemmatization  : 0.22
```

```
data/BATS\1_Inflectional_morphology :  category  [ noun - plural_irreg ]
total accuracy    total :  7 ps :  7 lm :  0 data :  50
total :0.14
Stemming  :0.14
lemmatization  : 0.0
```

```
data/BATS\1_Inflectional_morphology :  category  [ verb_inf - 3psg ]
total accuracy    total :  3 ps :  7 lm :  3 data :  50
total :0.06
Stemming  :0.14
lemmatization  : 0.06
```