# An application of Spectral Clustering on Named Entity Recognition

Ali Josue Limon, Michele Cerú

New York University

*ajl649@nyu.edu; mc3784@nyu.edu*

December 7, 2016

# Project description

- Goal:
- Purpose:

# Dataset

- bleble
- bleble

# Model parameters

- Size of the word embedding space: $d_E$
- technique used for the word embedding training: cbow or skip-gram
- Number of words embedding considered around the target word: $n_w$
- Kernell used for the spectral clustering: cosine similarity
    - Cosine similarity
    - Knearest neighbours
    - $k(x, y) = exp(-\gamma ||x - y||_1)$

# Evaluation Criteria

Adjusted Rand index: accuracy of the predicted labels given the ground truth.

$S_{11} = \{$pairs in the same cluster in C and C'$\}$

$S_{00} = \{$pairs in different cluster both in C and C'$\}$

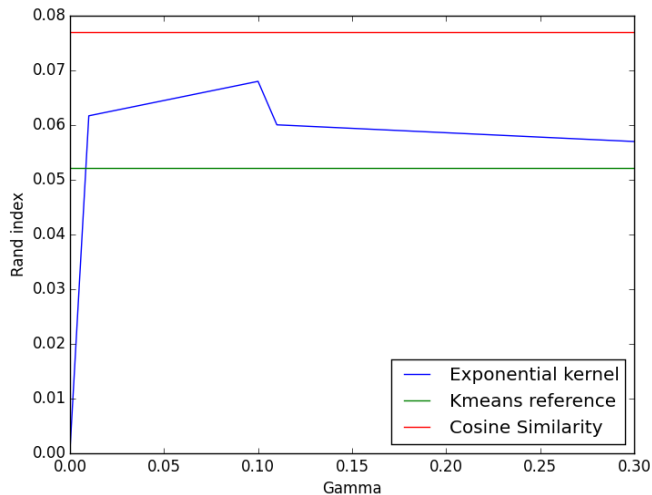$S_{10} = \{$pairs in the same cluste C but not in C'$\}$

$S_{01} = \{$pairs in the same cluste C' but not in C$\}$

Rand index:

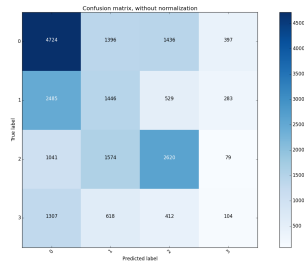$$R = \frac{n_{11} + n_{00}}{n_{11} + n_{10} + n_{01} + n_{00}} = \frac{n_{11} + n_{00}}{\binom{n}{2}}$$

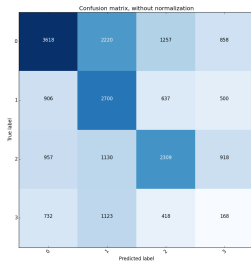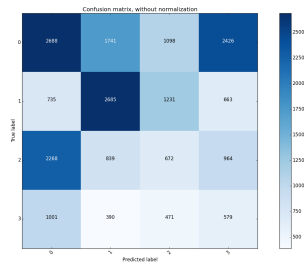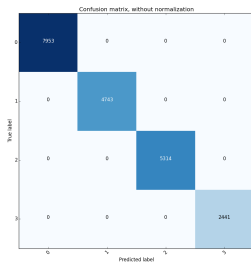Adjusted Rand index: Center the random clustering to 0

# result

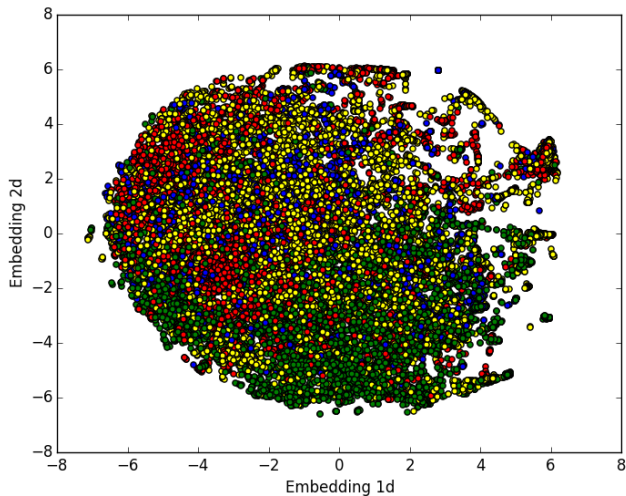# Data samples distribution ground truth

| Cluster Label | Name Entity | Number of sample |
|:---:|:---:|:---:|
| 0 | Organization | 7953 |
| 1 | Location | 4743 |
| 2 | Person | 5314 |
| 3 | Miscellaneous | 2441 |

# Confusion matrix

# Data visualisation in 2d

# Conclusion

- bleble
- bleble