

Document Semantic Similarity

TIS Project

Alberto Pirovano Francesco Picciotti

Politecnico di Milano

22nd April 2017



Outline

- 1 State of art
- 2 Preprocessing
- 3 Word2Vec
- 4 Doc2Vec



L'attuale stato dell'arte per similitudine semantica tra documenti si suddivide in:

- NLP Tradizionale



L'attuale stato dell'arte per similitudine semantica tra documenti si suddivide in:

- NLP Tradizionale
- Geometrico



L'attuale stato dell'arte per similitudine semantica tra documenti si suddivide in:

- NLP Tradizionale
- Geometrico
- Deep Learning based



L'attuale stato dell'arte per similitudine semantica tra documenti si suddivide in:

- NLP Tradizionale
- Geometrico
- Deep Learning based



Questo approccio segue la letteratura del Natural Language Processing e consiste dei seguenti passi:

- Cleaning dei dati
- Pos-Tagging
- Stemming o Lemmatisation
- Parsing
- Ontologia

Tuttavia nel caso del nostro scopo presenta delle criticità, ovvero:

- Affidabilità del Pos-Tagger italiano di TreeTagger
- Reperire una Ontologia e un parsing toll nella lingua italiana



Questo approccio segue la letteratura del Natural Language Processing e consiste dei seguenti passi:

- Cleaning dei dati
- Pos-Tagging
- Stemming o Lemmatisation
- Parsing
- Ontologia

Tuttavia nel caso del nostro scopo presenta delle criticità, ovvero:

- Affidabilità del Pos-Tagger italiano di TreeTagger
- Reperire una Ontologia e un parsing toll nella lingua italiana

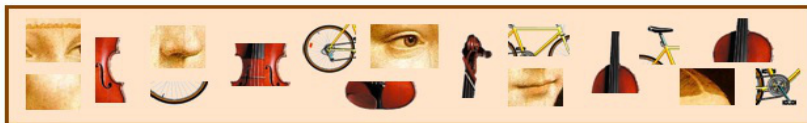
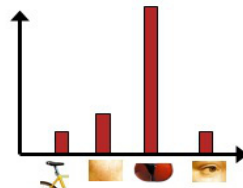
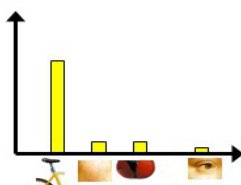
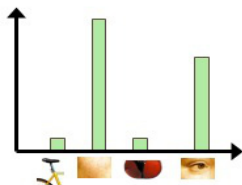


È basato sulla rappresentazione vettoriale del documento, dove le diverse dimensioni sono il Bag of Word di tutti i documenti. Questo approccio richiede dei seguenti steps:

- Cleaning dei dati
- Stemming o Lemmatisation
- Encoding del documento in vettore
- TF/IDF + LSA (Latent semantic analysis)



...limiti?



Per riassumere:

- Molto dipendente dal preprocessing del corpus
- Grande Bag Of World → bisogno di LSA, ma non è cos banale. Perché?

"LSA assumes that words that are close in meaning will occur in similar pieces of text"

- Un documento è un insieme **non ordinato** di parole



Negli ultimi anni il Deep Learning trova numerose applicazioni con ottimi risultati.

La creazione di Word2Vec da Google offre un modello che ha principalmente, i seguenti vantaggi e svantaggi.

Pros:

- Molto meno dipendente da un preprocessing
- Context-aware
- Combina il metodo *Geometrico* con quello *NLP Tradizionale*
- Non sfrutta una ontologia, ma la **crea**

Cons:

- Tecnica unsupervised
- Bisogno di un esperto per validare gradi di similarità
- Può risultare in GIGO system (Garbage In Garbage Out)



Outline

- 1 State of art
- 2 Preprocessing**
- 3 Word2Vec
- 4 Doc2Vec



"Preprocessing is 80% of NLP work"

Lev Konstantinovskiy

Qui spieghiamo il preprocessing e visualizzazione (volendo)



Outline

- 1 State of art
- 2 Preprocessing
- 3 Word2Vec**
- 4 Doc2Vec



Qui spaghiamo per bene come funziona word2vec



Outline

- 1 State of art
- 2 Preprocessing
- 3 Word2Vec
- 4 Doc2Vec



Qui spaghiamo per bene come funziona word2vec

