



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Analysis and reproducibility of paper “A Word Embedding based Generalized Language Model for Information Retrieval”

D. Ganguly, D. Roy, M. Mitra, G.J.F. Jones

11 february 2020

Alessio Lazzaron, Matteo Marchiori, Andrea Oriolo, Fabio Piras



- 1 Introduzione
- 2 Aspetti riprodotti
- 3 Risultati e osservazioni

■ Language Model (LM)

- Sviluppo di modelli probabilistici in grado di prevedere un documento, in una collezione, data una query.

$$P(d|q) = \prod_{t \in q} \lambda \frac{tf(t, d)}{|d|} + (1 - \lambda) \frac{cf(t)}{cs}$$

Approcci di LM classici $\left\{ \begin{array}{l} \text{Latent Dirichlet allocation (LDA)} \\ \text{Latent semantic analysis (LSA)} \end{array} \right.$

Si propone un nuovo modello denominato
GENERALIZED LANGUAGE MODEL (GLM)

■ Word Embedding

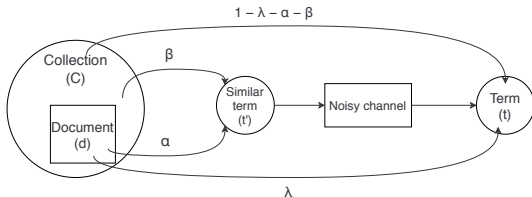
- Rappresentazione delle parole in un nuovo spazio, con lo scopo di memorizzarne le informazioni semantiche e sintattiche.

IDEA

Processo di trasformazione attraverso un «noisy channel» in cui un termine t' viene mutato in un altro termine t .

Tre approcci di trasformazione del termine:

- Direct term sampling
- Transformation via Document
- Transformation via Collection



$$P(t|d) = \lambda P(t|d) + \alpha \sum_{t' \in d} P(t, t'|d) P(t') \\ + \beta \sum_{t' \in N} P(t, t'|C) P(t') + (1 - \lambda - \alpha - \beta) P(t|C)$$

- Termine λ : Probabilità del termine t della query all'interno del documento d senza trasformazione.
- Termine α : Probabilità di trasformare il termine t della query in un termine t' appartenente a d .
- Termine β : Probabilità di trasformare il termine t della query in un termine t' appartenente alla collezione C .
- Termine $1 - \lambda - \alpha - \beta$: Probabilità del termine t della query di trovarsi all'interno della collezione C .

$$P(t, t'|d) = \frac{\text{sim}(t, t')}{\sum(d)} \frac{\text{tf}(t', d)}{|d|}$$

$$P(t, t'|C) = \frac{\text{sim}(t, t')}{\sum(Nt)} \frac{\text{cf}(t')}{cs}$$

- Per il calcolo della similarità tra t e t' entrambi i termini vengono rappresentati mediante Word Embedding.
- Una volta ottenuta la loro rappresentazione, per il calcolo dello score di similarità viene utilizzata la Cosine Similarity.

$$\text{sim}(A, B) = \cos(\vartheta) = \frac{A * B}{\|A\| * \|B\|} = \frac{\sum_{i=1}^n A_i * B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

$$P(t|d) = \lambda P(t|d) + \alpha \sum_{t' \in d} P(t, t'|d) P(t') \\ + \beta \sum_{t' \in N} P(t, t'|C) P(t') + (1 - \lambda - \alpha - \beta) P(t|C)$$

- Maggiore è il valore di similarità tra un termine t della query e i termini t' del documento, maggiore sarà il valore del termine in α .
- Maggiore è il valore di similarità tra un termine t della query e i termini t' della collezione, maggiore sarà il valore del termine in β .
- Il GLM generalizza il Language Model considerando non solo i singoli termini della query ma anche i termini semanticamente simili a essi, favorendo i documenti il cui contesto è simile a quello della query.



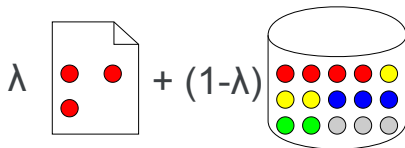
■ Lucene... Quale versione?

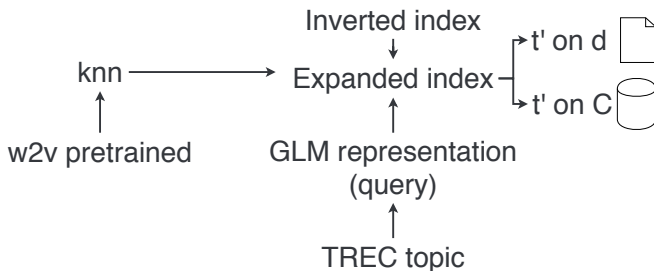


■ Indicizzazione... In che modo?

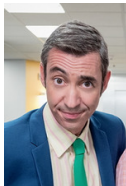
- Tokenizing
- Rimozione stopwords
- Indice trasposto

■ Base Language Model





- Primo approccio: chi fa da sè... Fa da sè!



- Secondo approccio: cerchiamo di prendere esempio e di migliorare

- Vincoli di tempo...

#	Score	Doc. Id	date	doc_number
0	0,6174	301084	940112	FT941-16044
1	0,6112	129708	930226	FT931-7040
2	0,6112	142295	930630	FT932-8
3	0,6112	171815	930715	FT933-14247
4	0,6112	187823	940329	FT941-706
5	0,6112	188903	940222	FT941-7958
6	0,6112	297797	940127	FT941-13091
7	0,5239	125326	930331	FT931-190

La più semplice query con Luke.
'from' è la parola con la frequenza più alta nella collezione.

Topic Set	MAP	GMAP	RECALL
<i>TREC-6</i>	0.2148 0.1612	0.0761 0.0843	0.4778 0.5382
<i>TREC-7</i>	0.1771 0.1734	0.0706 0.0843	0.4867 0.5444
<i>TREC-8</i>	0.2357 0.1899	0.1316 0.0906	0.5895 0.5484
<i>Robust</i>	0.2555 0.2021	0.1290 0.1016	0.7715 0.6052

Risultati di LM di base

- Indicizzazione con la stoplist standard di Lucene
- Parser standard di Lucene per ottenere query su cui trovare metriche
- Uso di uno stemmer? Improbabile...

- GLM migliore di LM
- Per ottenere i knn serve molto tempo
- Risultati non riprodotti
- Tempo sottostimato per implementare gli step necessari per riprodurre i risultati del paper
- È stato usato solo Lucene perché è quanto viene usato secondo il paper.



M. Agosti and G. Silvello, *Slides*, Padua, Italy: academic year 2019-2020.



D. Ganguly, M. Mitra, D. Roy and G. Jones, *A Word Embedding based Generalized Language Model for Information Retrieval*, Dublin, Ireland; Kolkata, India: SIGIR 2015.



D. Ganguly, *GLM GitHub repository*



V. Lavrenko, *Jelinek-Mercer Smoothing video*