

NLP - LEZIONE 19

DEL 12/12/2019

RAPPR. DISTRIBUITA I

SIMILARITÀ TRA FRASI

Siamo interessati a CALCOLARE la SIMILARITÀ tra due date frasi

$$P_1 = w_1 \dots w_m$$

$$P_2 = k_1 \dots k_m$$

Per fare questo possiamo iniziare cercando una eventuale relazione di PARAFRASI tra P_1 e P_2 . Questa relazione significa che P_1 e P_2 stanno indicando la STESSA INFORMAZIONE.

Notiamo però che il nostro di lavorare con le frasi intere non è più utile lavorare con le SOTTOFRASI delle due frasi. Il concetto di SINONIMIA è infatti più CHIARO rispetto a quello della PARAFRASI.

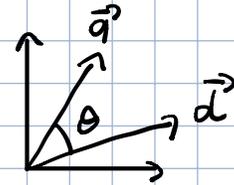
PARAFRASI := Relazione tra due FRASI che hanno lo stesso significato.

SINONIMIA := Relazione tra due PAROLE che hanno lo stesso significato.

Motivazioni che studiamo IR (INFORMATION RETRIEVAL)
per calcolare la similarità tra due frasi otteniamo

i) Rappresentare le FRASI come VETTORI in uno spazio vettoriale.

ii) Utilizziamo il PRODOTTO SCALARE per calcolare l'angolo θ tra i due vettori.



Con il passo i) otteniamo trasformate nella RAPPRESENTAZIONE DISCRETA ($P = w_1 \dots w_m$) in una RAPPRESENTAZIONE CONTINUA ($\vec{d} = (v_1, \dots, v_m)$). Le rappresentazioni CONTINUE sono anche dette DISTRIBUTED REPRESENTATIONS.

Q: Che RELAZIONI intercorrono tra le rapp. DISCRETE e quelle DISTRIBUITE?

Nella trasformazione della rapp. DISCRETA \rightarrow DISTRIBUITA utilizzate in IR si PERDONO le informazioni relative all'ORDINE delle parole nelle frasi.

In generale si passa da rapp. DISCRETA a quelle DISTR. per poter lavorare su degli SPAZI METRICI. Molti algoritmi di APPRENDIMENTO funzionano su spazi metrici.

Andiamo quindi a vedere come si passa da word .
DISCRETE a word . DISTRIBUITE nel contesto del NLP.

BAG OF WORDS ~ (14:00 min)

Sia $P = w_1 \dots w_m$. Vogliamo capire come ottenere

$$F(P) = F(w_1 \dots w_m) = \vec{P}$$

Nel modello BAG OF WORDS non possiamo utilizzare la
COMPOSIZIONALITÀ. In particolare definiamo invece
una serie di VETTORI ONE-HOT per ogni PAROLA DISTINTA

$$F(w_1) = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

\vdots

$$F(w_m) = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

ONE-HOT VECTORS,

formano una BASE ORTONORMALE
della matrice vettoriale.

(LOCAL DISTRIBUTED REPRESENTATION,
PLATE 1995)

Dato una frase $P = w_1 \dots w_m$ possiamo come segue:

$$F(w_1 \dots w_m) = \sum_{i=1}^m F(w_i)$$

È poi possibile assegnare un PESO al contributo di
ogni parola (vedi tf-idf ranking in IR).

Come detto prima, il modello BAG OF WORDS perde informazioni sull'ordine delle parole. Se vogliamo CODIFICARE in una rapp. DISCRETA la frase $P = w_1 \dots w_m$ potremmo procedere come segue

$$F(w_1 \dots w_m) = \underbrace{F(w_1)}_{\text{CODIFICA } w_1} \oplus F(w_2) \oplus \dots \oplus F(w_m)$$

OPERAZIONE DI
CONCATENAZIONE

Con questa rapp. mantengo l'ordine delle parole in quanto mantengo SEPARATE le rapp. delle varie PARTI. Motivato inoltre che se devo rapp. n simboli, con questa rapp. necessito $\log n$ bits per simbolo. Con la rapp. distribuite invece necessito n bits per simbolo. Questo tipo di rapp. è proprio quello che viene utilizzato per MEMORIZZARE i dati nei CALCOLATORI moderni.

La rapp. simbolica però non è METRICA, e quindi non posso calcolare la DISTANZA tra due ELEMENTI.

OSS: La rapp. DISTRIBUITA è chiamata così perché l'informazione dei SIMBOLI DISCRETI non si trova più in un POSTO UNICO ma è SPARSA in tanti posti.

FUNCTIONAL COMPOSITIONALITY

~ (35:00 min)

La FUNCTIONAL COMPOSITIONALITY è ottenuta applicando una funzione alle parti e poi mettendo tutto assieme. La funzione applicata sulle PARTI e la funzione applicata sul TUTTO potrebbero essere DIVERSE, come nel modello BAG OF WORDS.

La CONCATENATIVE COMPOSITIONALITY è ottenuta concatenando la CODIFICA di ogni parola. Con questo tipo di rapp. NON PERDIAMO informazioni sui singoli simboli utilizzati.

Notiamo che la composizionalità utilizzata nel BAG OF WORDS model non è concatenativa, in quanto perdiamo l'ordine delle parole.

Ci poniamo quindi i seguenti QUESITI:

Q1) Quanto CONCATENATIVA può essere una composizionalità FUNZIONALE?

Q2) Utilizzando una rapp. DISTRIBUITA, come facciamo a RIDURRE la DIMENSIONE dei vettori che utilizziamo?

OSS: L'ANALISI delle differenze tra i vari tipi di COMPOSIZIONALITÀ fu introdotta inizialmente dalle persone che hanno sviluppato la teoria delle RETI NEURALI.

Si era interessati a capire se le rapp. DISTRIBUITE contengono più o meno INFORMAZIONI rispetto a quelle DISCRETE.

Alcuni (FODOR e PYLSHYN) sostenevano che non c'era alcuna differenza tra le due rapp. Un altro (?) diceva che l'unica diff. è il fatto che nel mondo DISTRIBUITO le informazioni non sono ONTICAMENTE e questo ci permette di fare una SINGOLA operazione MATRICIALE che prende TUTTE le informazioni ASSIEME.

SEMANTICA DISTRIBUZIONALE (PRESENTAZIONE FEDERICO & PAOLO) ~ (49:00)

Lo SPAZIO METRICO DISTRIBUZIONALE è stato introdotto per cercare di DIMINUIRE la DIMENSIONE dei vettori con cui si lavora. È basato sulla seguente IPOTESI LINGUISTICA, dette anche IPOTESI DISTRIBUZIONALE:

due PAROLE non SIMILI tendono a comparire negli STESSI DOCUMENTI.

OSS: Lo SPAZIO METRICO DISTRIBUZIONALE è DIVERSO da quello DISTRIBUITO, e trova le sue origini nelle IDEE di due LINGUISTI: FIRTH e HARRIS.

FIRTH → "Il SIGNIFICATO delle PAROLE è definito dal CONTESTO in cui vengono utilizzate."

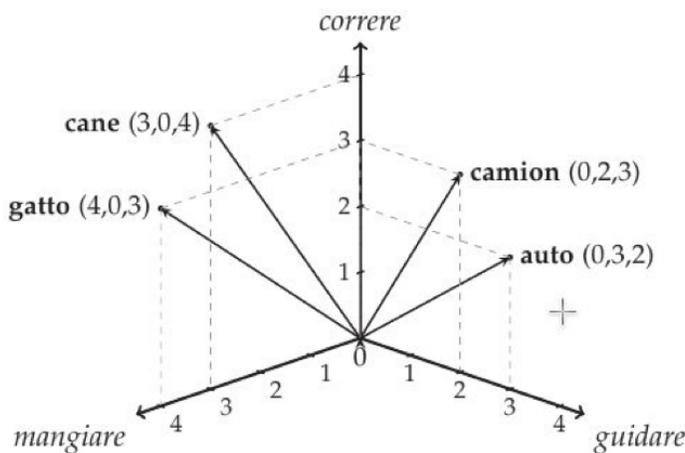
HARRIS → "Due PAROLE NON SIMILI se si PRESENTANO in CONTESTI SIMILI."

L'idea dietro all'IPOTESI DISTRIBUZIONALE è che la PRESENZA di una parola INFLUENZA le parole vicine alla parola: "Il CANE abbaia" vs "Il GATTO miagola".

Lo spazio metrico DISTRIBUZIONALE funziona quindi definendo un insieme di CONTESTI, che formano la base dello spazio vettoriale. Si vuole quindi associare ad ogni PAROLA un VETTORE.

In questo esempio i contesti sono

$B = \{ \text{MANGIARE, CORRERE, GUIDARE} \}$



Formalmente lo SPAZIO SEMANTICO DISTRIBUZIONALE è definito come quadrupla $S = \langle T, B, M, F \rangle$, con:

- $T := \{ \text{PAROLE TARGET } w \}$
- $B :=$ la BASE che contiene i CONTESTI LINGUISTICI nei quali viene volute la SIMILARITÀ.
- $M :=$ Matrice di CO-OCCORRENZA
- $F :=$ METRICA utilizzata per misurare la DISTANZA tra i vettori.

L'aspetto FONDAMENTALE nella costruzione di questi spazi è la SCELTA dei CONTESTI, ciascuno con la rispettiva GRANULARITÀ, che utilizziamo per COSTRUIRE i vari VETTORI.

OSS: Motiemi che utilizzando CONTESTI altrettanto "grandi" potrei trovare SIMILI parole OPPOSITE come "SALIRE LE SCALE", "SCENDERE LE SCALE".

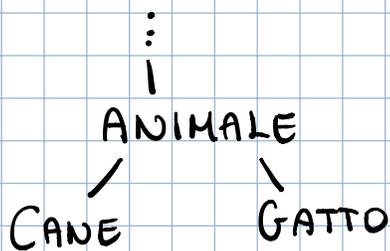
Nei SPAZI DISTRIBUZIONALI devo quindi stare ATTENTO a:

- i) \mathcal{S} VETTORI che rapp. le varie PAROLE.
- ii) La FINESTRA utilizzata per CONTESTUALIZZARE le varie parole. Più GRANDE è la finestra, e più le parole saranno SIMILI per ragioni diverse dalle loro PROPRIETÀ.

WORD EMBEDDING

Il WORD-EMBEDDING è la TRASPOSIZIONE nelle RETI NEURALI del fatto che la similarità tra due parole DIPENDE dalla LUNGHEZZA DEL CONTESTO e del CORPUS che utilizzo per fare TRAINING.

COTOPY



CANE e GATTO possono essere considerate simili in quanto non la SPECIALIZZAZIONE della STESSA CLASSE SUPERIORE.

Rimetto a quanto fatto prima con lo SPAZIO METRICO DISTRIBUITO, con lo SPAZIO METRICO DISTRIBUZIONALE abbiamo cambiato i VETTORI e abbiamo DIMINUITO la DIMENSIONE dello SPAZIO.

Iniziamo con

$$F(w) = \vec{w} \in \mathbb{R}^m \quad (\text{ONE-HOT VECTOR})$$

Poniamo adesso utilizzare la MATRICE DI CO-OCCORRENZA $W_{d \times m}$ definite nello spazio metrico DISTRIBUZIONALE come segue

$$W_{d \times m} \cdot \vec{w} = \vec{d} \in \mathbb{R}^d \quad \left(\begin{array}{l} \text{MODELLO} \\ \text{DISTRIBUZIONALE} \end{array} \right)$$

dove,

- $m := \#$ di PAROLE
- $d := \#$ di CONTESTI
- $(w_{i,s}) := \#$ di volte in cui la PAROLA s appare nel CONTESTO i .

Notiamo che $\vec{d} = W \cdot \vec{w}$ è un esempio di WORD EMBEDDING in quanto stiamo PASSANDO da un vettore all'altro tramite una TRASFORMAZIONE LINEARE.

Notiamo altresì che NON SEMPRE ad OGNI DIMENSIONE è associato un SIGNIFICATO PRECISO. Più accade infatti che ad una DIMENSIONE associamo una COMBINAZIONE LINEARE di più parole.

Per cercare di DIMINUIRE d l'idea è studiare quelle di trovare le righe della MATRICE W che non LINEARMENTE DIPENDENTI. Questa idea è implementata nella PCA (PRINCIPAL COMPONENT ANALYSIS)

La PCA forma un insieme di tecniche che utilizzano gli AUTOVETTORI della MATRICE per ridurre il $\#$ di DIMENSIONI.

La SVD (SINGULAR VALUE DECOMPOSITION) è la più famosa tecnica di PCA.

OSS: Andando a diminuire il # di dim. non necessariamente PERDERE l'IDENTITÀ UNIVUCA dei SIMBOLI.

Ad esempio i SINONIMI potrebbero essere mappati nello STESSO ELEMENTO.

OSS: Alcuni esempi di RETI NEURALI che utilizzano WORD-EMBEDDING sono:

- ELMO
- BART
- CONCEPT-NET