



## TLN: text categorization con risorse semantiche: il metodo di Rocchio

Daniele Radicioni

Daniele Radicioni - TLN

## terminologia

- *documento* è un'unità di testo indicizzata nel sistema e disponibile per la ricerca.
- *collezione* è un insieme di documenti, utilizzati per soddisfare le richieste dell'utente.
- *termine* è l'elemento (dal singolo elemento lessicale a intere *phrases*) che occorre nella collezione.
- *interrogazione* è la richiesta dell'utente, espressa come insieme di termini.



Daniele Radicioni - TLN

2

## vector space model

- il valore di ogni feature è detto *term weight*, ed è una *funzione della frequenza del termine nel documento*.
  - e.g., in una ricetta di cucina sul *pollo fritto* trovata sul Web, i termini *pollo*, *fritto*, *olio*, *pepe* occorrono rispettivamente 8, 2, 7, 4 volte.
  - usando questa frequenza come *peso del termine*, il *vettore per questo documento* (documento *j*) sarà espresso come

$$\vec{d}_j = (8, 2, 7, 4)$$



Daniele Radicioni - TLN

3

## il vettore del documento

- in generale, il vettore di un documento  $d_j$  è rappresentato come

$$\vec{d}_j = (w_{1,j}, w_{2,j}, w_{3,j}, \dots, w_{n,j})$$

- dove l'elemento  $w_{ij}$  individua l'elemento  $i$  fra gli  $N$  complessivamente presenti nella *collezione*, e  $j$  è il  $j$ -esimo documento della collezione.



Daniele Radicioni - TLN

4

## il vettore della query

- analogamente, anche la query è rappresentata come un vettore; la query *pollo fritto* potrebbe essere

$$\vec{q} = (1, 1, 0, 0) \leftarrow (\text{pollo, fritto, olio, pepe})$$

- e in generale,

$$\vec{q} = (w_{1,q}, w_{2,q}, w_{3,q}, \dots, w_{n,q})$$



## un'altra ricetta

- consideriamo un'altra ricetta per il *pollo cotto*, contenuta nel documento  $d_k$ , che potrebbe essere rappresentata dal vettore

$$\vec{d}_k = (6, 1, 0, 0)$$

- vorremmo che la query *pollo fritto* restituisse il documento  $d_j$ , (la ricetta del pollo fritto), piuttosto che il documento  $d_k$  (pollo cotto)



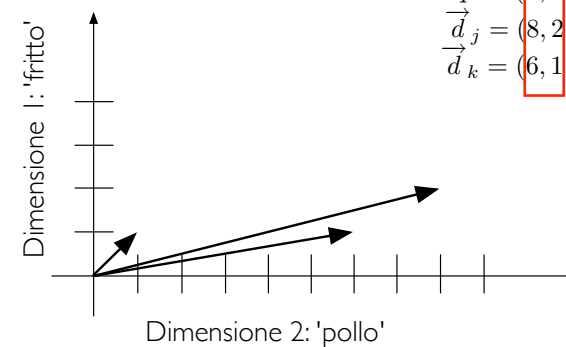
## FV Model

- è utile immaginare le *features* usate per rappresentare i documenti e le interrogazioni in questo modello come *dimensioni in uno spazio multidimensionale*, in cui i pesi delle features servono a individuare i documenti.
- una query utente tradotta in un vettore denota un punto in quello spazio.
- i documenti che sono vicini alla query sono in qualche modo più rilevanti (pertinenti?) ai fini della query, di quanto siano i documenti lontani.



## FV Model

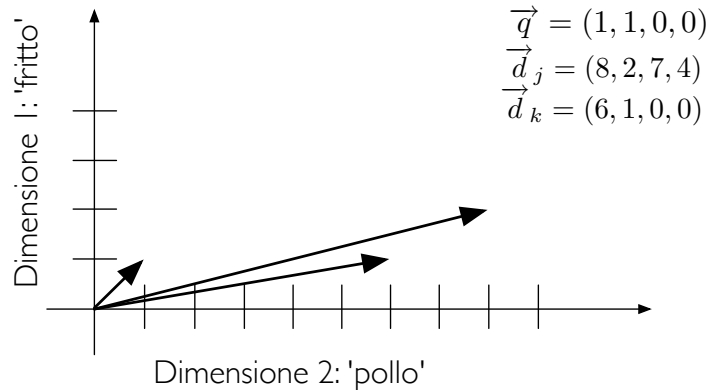
$$\begin{aligned}\vec{q} &= (1, 1, 0, 0) \\ \vec{d}_j &= (8, 2, 7, 4) \\ \vec{d}_k &= (6, 1, 0, 0)\end{aligned}$$



rappresentazione delle prime due dimensioni del vector model.



## similarity



- se per misurare la somiglianza fra i vettori **consideriamo gli angoli fra i vettori**, la query è più simile al documento  $j$ .



Daniele Radicioni - TLN

9

## metriche di similarità

- invece dell'angolo si usa **il coseno dell'angolo fra i vettori**, calcolato come:

$$\text{sim}(\vec{q}, \vec{d}_j) = \frac{\sum_{i=1}^N w_{i,q} \times w_{i,j}}{\sqrt{\sum_{i=1}^N w_{i,q}^2} \times \sqrt{\sum_{i=1}^N w_{i,j}^2}}$$

- il coseno fra due documenti identici sarà 1, mentre il coseno fra due documenti *ortogonali* (privi di termini in comune) sarà 0.



Daniele Radicioni - TLN

10

## normalizzazione del dot product

- il coseno può essere pensato come *prodotto scalare*

$$\text{sim}_{\text{dot-product}}(\vec{v}, \vec{w}) = \vec{v} \cdot \vec{w} = \sum_{i=1}^N v_i \times w_i$$

- normalizzato per la lunghezza dei vettori, definita come

$$|\vec{v}| = \sqrt{\sum_{i=1}^N v_i^2}$$

- quindi

$$\text{sim}_{\text{cosine}}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\sum_{i=1}^N v_i \times w_i}{\sqrt{\sum_{i=1}^N v_i^2} \times \sqrt{\sum_{i=1}^N w_i^2}}$$



Daniele Radicioni - TLN

11

## pesatura dei termini

- finora abbiamo assunto che i **pesi dei termini** fossero determinati semplicemente con il conto della **frequenza dei vari termini nei documenti**.
- in realtà due fattori sono critici nell'attribuzione del giusto 'peso' ai termini:
  1. la *frequenza*;
  2. la frequenza in un documento ma la non-frequenza a livello dell'intera collezione



Daniele Radicioni - TLN

12

## IDF: inverse document frequency

- l'idea è che i termini che sono presenti in pochi documenti sono utili per **distinguere quei documenti dal resto della collezione** (i termini che ricorrono in tutti i documenti non sono utili: e.g., *stop words*).
- **pertanto l'IDF è definito per mezzo del rapporto  $N/n_i$** , dove  $N$  è il numero totale di documenti nella collezione, e  $n_i$  è il numero di documenti in cui il termine  $i$  occorre.



## inverse document frequency

- **quanti meno sono i documenti in cui un termine occorre, tanto più alto è il peso**. il peso minimo è 1, corrispondente a termini che occorrono in tutti i documenti.
- dato l'elevato numero di documenti in varie collezioni, il rapporto è 'schiacciato' con una funzione logaritmica.

$$\text{idf}_i = \log \left( \frac{N}{n_i} \right)$$



## inverse document frequency

- combinando  $tf$  e  $idf$  otteniamo il seguente schema:

$$w_{i,j} = \text{tf}_{i,j} \times \text{idf}_{i,j}$$

- utilizzando lo schema  $tf \cdot idf$  il peso del termine  $i$  nel vettore del documento  $j$  è calcolato come il prodotto della sua frequenza totale in  $j$  per il logaritmo della sua  $idf$  nella collezione

$$\text{sim}(\vec{q}, \vec{d}) = \frac{\sum_{w \in q, d} \text{tf}_{w,q} \text{tf}_{w,d} (\text{idf}_w)^2}{\sqrt{\sum_{q_i \in q} (\text{tf}_{q_i,q} \text{idf}_{q_i})^2} \times \sqrt{\sum_{d_i \in d} (\text{tf}_{d_i,d} \text{idf}_{d_i})^2}}$$



## The Rocchio Method



## profilo di una classe di documenti

- Some linear classifiers consist of an *explicit profile* (or *prototypical document*) of the category.
- Learning a linear classifier is often preceded by local *term space reduction*; in this case, a profile of  $c_i$  is a weighted list of the terms whose presence or absence is most useful for discriminating  $c_i$ .



## il centroide

- Let  $D$  be the *set of documents* and  $T$  be the *dictionary* i.e. the set of all different terms occurring in  $D$ .

$$T = \{t_1, \dots, t_m\}$$

- then the *absolute frequency* of term  $t \in T$  in document  $d \in D$  is given by  $tf(d, t)$ .
- We denote the term vectors

$$\vec{t}_d = (tf(d, t_1), \dots, tf(d, t_m))$$



## il centroide

- il *centroide* di un insieme di  $X$  vettori di termini è definito come

$$\vec{t}_X := \frac{1}{|X|} \sum_{\vec{t}_d \in X} \vec{t}_d$$



## metodo di Rocchio: intuizione

- A classifier built by means of the Rocchio method rewards
  - the closeness of a test document to the centroid of the positive training examples, and
  - its distance from the centroid of the negative training examples.



## metodo di Rocchio

$$\vec{c}_i = \langle f_{1i}, \dots, f_{|T|i} \rangle$$

- for each class  $c_i$  (with  $i$  varying over classes) we compute the weight of the  $k$ -th feature  $f$  as

$$f_{ki} = \beta \cdot \sum_{d_j \in POS_i} \frac{w_{kj}}{|POS_i|} - \gamma \cdot \sum_{d_j \in NEG_i} \frac{w_{kj}}{|NEG_i|}$$

$$POS_i = \{d_j \in Tr | \Phi(d_j, c_i) = T\}$$

$$NEG_i = \{d_j \in Tr | \Phi(d_j, c_i) = F\}$$

- $\beta$  and  $\gamma$  are **control parameters** that allow setting the relative importance of positive and negative examples
- if  $\beta$  is set to 1 and  $\gamma$  to 0 the profile of  $c$  is the centroid of its positive training examples.



Daniele Radicioni - TLN

21

## raffinamento: i near positives

$$f_{ki} = \beta \cdot \sum_{d_j \in POS_i} \frac{w_{kj}}{|POS_i|} - \gamma \cdot \sum_{d_j \in NPOS_i} \frac{w_{kj}}{|NPOS_i|}$$

- **Open issue:** whether the set  $NEG_i$  should be considered in its entirety, or whether a well- chosen sample of it, such as the set  $NPOS_i$  of near-positives (defined as “the most positive among the negative training examples”).
- Near-positives are the most difficult documents to tell apart from the positives.



Daniele Radicioni - TLN

22

## consegna

1. sono dati due set di documenti (italiano vs. inglese; 20 docs x 10 classi, vs. 20 docs x 20 classi) con caratteristiche diverse. sceglierne uno.
2. scrivere un programma che utilizzi una rappresentazione basata sul *feature vector model*;
3. il programma deve costruire i profili (Rocchio) relativi alle classi cui appartengono i documenti, e
4. deve classificare i nuovi documenti (utilizzando come metrica di distanza la *cosine similarity* o una delle metriche derivate) in una delle classi date.
  - per ciascuna classe prendiamo il 90% dei documenti per costruire i profili, e il 10% per testare
  - **opzionale:** partizionare in 10 split il dataset, ripetendo la procedura 10 volte ('training' su 90% e test su 10%), calcolando la media dell'accuratezza così ottenuta.



Daniele Radicioni - TLN

23

## consegna - scelte progettuali

- a proposito del punto 1:
  - quali feature scegliere? sperimentare con i lemmi dei termini in input e con i relativi BabelNet IDs (nel secondo caso è necessario predisporre una qualche forma di WSD).
  - implementando entrambe le alternative, quali differenze nei risultati ci possiamo attendere?
- nell'implementazione del metodo di Rocchio partire assegnando a  $\beta$  e  $\gamma$  i valori di 1 e 4, rispettivamente.
- algoritmo per calcolare i documenti  $NPOS$ ?



Daniele Radicioni - TLN

24