

Text Classification: the Geometrical approach. Vector models, and similarity

R. Basili

Corso di *Web Mining e Retrieval*
a.a. 2022-23

March 8, 2023

Outline

Outline

- 1 *Overview*
- 2 *Vector Spaces*
 - Inner Product, Norms and Distances
- 3 *Distance, similarity and classification*
 - The Rocchio TC model
 - Memory Based Learning
 - Distances and similarities
 - Distances and similarities: Discussion
 - Other Distance Metrics
 - Discussion
- 4 *A digression: IT*
- 5 *Probabilistic Norms*
 - Mutual Information
 - Probabilistic Norms
- 6 *References*



Real-valued Vector Space

Vector Space definition:

A vector space need to satisfy the following axioms:

Real-valued Vector Space

Vector Space definition:

A vector space need to satisfy the following axioms:

Sum

To every pair, \underline{x} and \underline{y} , of vectors in V there corresponds a vector $\underline{x} + \underline{y}$, called the sum of \underline{x} and \underline{y} , in such a way that:

- 1 sum is commutative, $\underline{x} + \underline{y} = \underline{y} + \underline{x}$
- 2 sum is associative,
 $\underline{x} + (\underline{y} + \underline{z}) = (\underline{x} + \underline{y}) + \underline{z}$
- 3 there exist in V a unique vector Φ (called the origin) such that
 $\underline{x} + \Phi = \underline{x} \forall \underline{x} \in V$
- 4 $\forall \underline{x} \in V$ there corresponds a unique vector $-\underline{x}$ such that $\underline{x} + (-\underline{x}) = \Phi$

Real-valued Vector Space

Vector Space definition:

A vector space need to satisfy the following axioms:

Sum

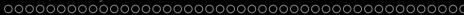
To every pair, \underline{x} and \underline{y} , of vectors in V there corresponds a vector $\underline{x} + \underline{y}$, called the sum of \underline{x} and \underline{y} , in such a way that:

- 1 sum is commutative, $\underline{x} + \underline{y} = \underline{y} + \underline{x}$
- 2 sum is associative,
 $\underline{x} + (\underline{y} + \underline{z}) = (\underline{x} + \underline{y}) + \underline{z}$
- 3 there exist in V a unique vector Φ (called the origin) such that
 $\underline{x} + \Phi = \underline{x} \forall \underline{x} \in V$
- 4 $\forall \underline{x} \in V$ there corresponds a unique vector $-\underline{x}$ such that $\underline{x} + (-\underline{x}) = \Phi$

Scalar Multiplication

To every pair α and \underline{x} , where α is a scalar and $\underline{x} \in V$, there corresponds a vector $\alpha \underline{x}$, called the product of α and \underline{x} , in such a way that:

- 1 associativity $\alpha(\beta \underline{x}) = (\alpha\beta)\underline{x}$
- 2 $1\underline{x} = \underline{x} \quad \forall \underline{x} \in V$
- 3 mult. by *scalar* is distributive wrt. vector addition $\alpha(\underline{x} + \underline{y}) = \alpha\underline{x} + \alpha\underline{y}$
- 4 mult. by *vector* is distributive wrt. scalar addition $(\alpha + \beta)\underline{x} = \alpha\underline{x} + \beta\underline{x}$



Vector Operations

Sum of two vector \underline{x} and \underline{y}

$$\underline{x} + \underline{y} = |\underline{x}\rangle + |\underline{y}\rangle = \begin{pmatrix} x_1 + y_1 \\ \cdot \\ \cdot \\ \cdot \\ x_n + y_n \end{pmatrix}$$



Vector Operations

Sum of two vector \underline{x} and \underline{y}

$$\underline{x} + \underline{y} = |\underline{x}\rangle + |\underline{y}\rangle = \begin{pmatrix} x_1 + y_1 \\ \cdot \\ \cdot \\ \cdot \\ x_n + y_n \end{pmatrix}$$

Multiplication by scalar α

$$\alpha \underline{x} = \alpha |\underline{x}\rangle = \begin{pmatrix} \alpha x_1 \\ \cdot \\ \cdot \\ \cdot \\ \alpha x_n \end{pmatrix}$$

Linear combination

$$\underline{y} = c_1 \underline{x}_1 + \cdots + c_n \underline{x}_n$$

or

$$|\underline{y}\rangle = c_1 |\underline{x}_1\rangle + \cdots + c_n |\underline{x}_n\rangle$$



Linear dependence

Conditions for linear dependence

A set of vectors $\{\underline{x}_1, \dots, \underline{x}_n\}$ are *linearly dependent* if there a set constant scalars c_1, \dots, c_n exists, not all 0, such that:

$$c_1 \underline{x}_1 + \dots + c_n \underline{x}_n = \underline{0}$$

Linear dependence

Conditions for linear dependence

A set of vectors $\{\underline{x}_1, \dots, \underline{x}_n\}$ are *linearly dependent* if there a set constant scalars c_1, \dots, c_n exists, not all 0, such that:

$$c_1 \underline{x}_1 + \dots + c_n \underline{x}_n = \underline{0}$$

Conditions for linear independence

A set of vectors $\{\underline{x}_1, \dots, \underline{x}_n\}$ are *linearly independent* if and only if the *linear condition* $c_1 \underline{x}_1 + \dots + c_n \underline{x}_n = \underline{0}$ is satisfied only when $c_1 = c_2 = \dots = c_n = 0$



Basis

Definition:

A *basis* for a space is a set of n linearly independent vectors in a n -dimensional vector space V_n .

Inner Product

Definition:

Is a real-valued function on the cross product $V_n \times V_n$ associating with each pair of vectors $(\underline{x}, \underline{y})$ a unique real number.

The function (\cdot, \cdot) has the following properties:

- 1 $(\underline{x}, \underline{y}) = (\underline{y}, \underline{x})$
- 2 $(\underline{x}, \lambda \underline{y}) = \lambda (\underline{x}, \underline{y})$
- 3 $(\underline{x}_1 + \underline{x}_2, \underline{y}) = (\underline{x}_1, \underline{y}) + (\underline{x}_2, \underline{y})$
- 4 $(\underline{x}, \underline{x}) \geq 0$ and $(\underline{x}, \underline{x}) = 0$ **iff** $\underline{x} = \underline{0}$

Standard Inner Product

$$(\underline{x}, \underline{y}) = \sum_{i=1}^n x_i y_i$$



Norm

Geometric interpretation

Geometrically the *norm* represent the length of the vector

Norm

Geometric interpretation

Geometrically the *norm* represent the length of the vector

Definition

The *norm* id a function $\|\cdot\|$ from V_n to \mathbb{R}

Euclidean Norm:

$$\|\underline{x}\| = \sqrt{(\underline{x}, \underline{x})} = \sqrt{\sum_{i=1}^n x_i^2} = (x_1^2 + \dots + x_n^2)^{1/2}$$

Properties

- 1 $\|\underline{x}\| \geq 0$ and $\|\underline{x}\| = 0$ if and only if $\underline{x} = 0$
- 2 $\|\alpha \underline{x}\| = |\alpha| \|\underline{x}\|$ for all α and \underline{x}
- 3 $\forall \underline{x}, \underline{y}, |(\underline{x}, \underline{y})| \leq \|\underline{x}\| \|\underline{y}\|$ (Cauchy-Schwartz)

A vector $\underline{x} \in V_n$ is a *unit vector*, or *normalized*, when $\|\underline{x}\| = 1$



From Norm to distance

In V_n we can define the distance between two vectors \underline{x} and \underline{y} as:

$$d(\underline{x}, \underline{y}) = \|\underline{x} - \underline{y}\| = \sqrt{(\underline{x} - \underline{y}, \underline{x} - \underline{y})} = ((x_1 - y_1)^2 + \dots + (x_n - y_n)^2)^{1/2}$$

These measure, noted sometimes as $\|\underline{x} - \underline{y}\|_2^2$, is also named *Euclidean distance*.

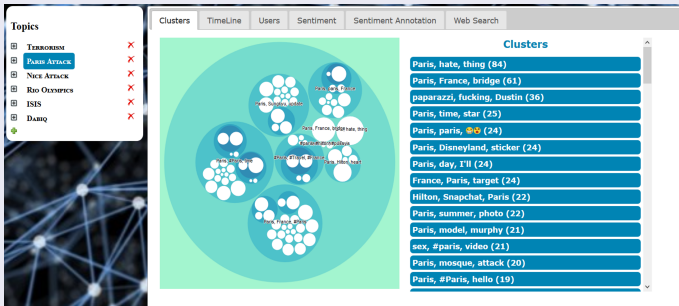
Similarity

Looking to texts as points a n -dimensional space

A structure for organizing large bodies of texts for efficient searching and browsing can be the notion of metric space.

Internet search engines may suitably exploit cluster analysis to documents in order to organize them visually.

Clustering of texts for browsing



Text Classification in the Vector Space Model

Text Classification: Definition

Given:

- a set of target categories, $C = \{C_1, \dots, C_n\}$:
- the set T of documents,

define a function: $f : T \leftarrow 2^C$

Vector Space Model (Salton89)

Features are dimensions of a Vector Space.

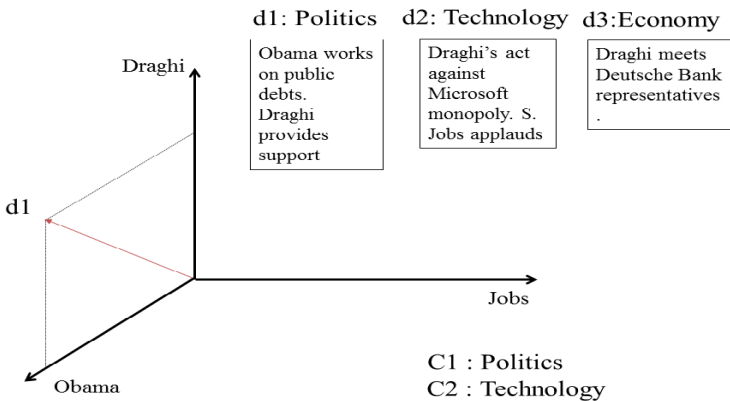
Documents d and Categories C_i are mapped to vectors of feature weights (\underline{d} and \underline{C}_i , respectively).

Geometric Model of $f()$:

A document d is assigned to a class C_i if $(\underline{d}, \underline{C}_i) > \tau_i$

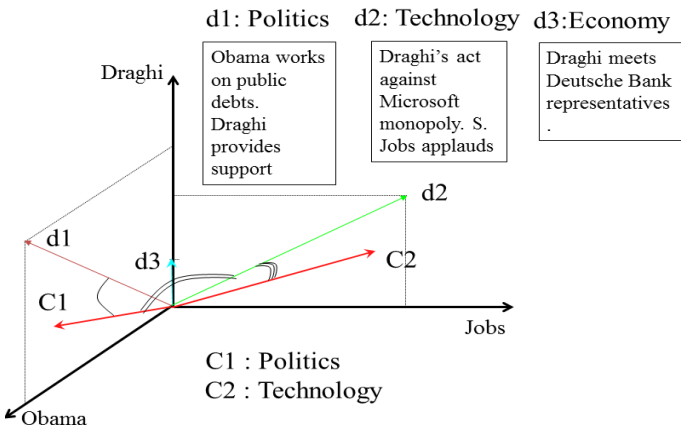
Text Classification: Vector Space Modeling

In Vector Space Model documents words corresponds to the space (orthonormal) basis, and individual texts are mapped into vectors ...



Text Classification: Classification Inference

Categories are also vectors and cosine similarity measures can support the final inference about category membership, e.g. $d1 \in C1$ and $d2 \in C2$:





A simple model for Text Classification

Motivation

Rocchio's is one of the first and simple models for *supervised text classification* where:

- *document vectors* are weighted according to a standard function, called $tf \cdot idf$,
- *category vectors*, $\underline{C}_1, \dots, \underline{C}_n$, are obtained by *averaging* the behaviour of the training examples.

We thus need to define a weighting function: $\omega(w, d)$ for individual words w in documents d and a method to design a category vector, i.e. a profile, as a linear combination of document vectors.

Similarity

Once vectors for documents and Category profiles (\underline{C}_i) are made available than the standard cosine similarity is adopted for inferencing, i.e. again a document d is assigned to a class C_i if $(\underline{d}, \underline{C}_i) > \tau_i$

k-NN: the algorithm

For each each training example $\langle x, c(x) \rangle \in D$
 Compute the corresponding TF-IDF vector, \underline{x} , for document x .

Test instance y :
 Compute TF-IDF vector \underline{y} for document y .
 For each $\langle x, c(x) \rangle \in D$

$$s_x = \text{cosSim}(\underline{y}, \underline{x}) = \frac{(\underline{y}, \underline{x})}{\|\underline{x}\| \cdot \|\underline{y}\|}$$

Sort examples $x \in D$ by decreasing values of s_x .
 Let kNN be the set of the closest (i.e. first) k examples in D .

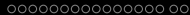
RETURN the majority class of examples in kNN .



Similarity

The role of similarity among vectors

In most of the examples above, document data are expressed as high-dimensional vectors, characterized by very sparse term-by-document matrices with positive ordinal attribute values and a significant amount of outliers.



Distance/similarity functions that have not a geometrical origin.

The role of probability

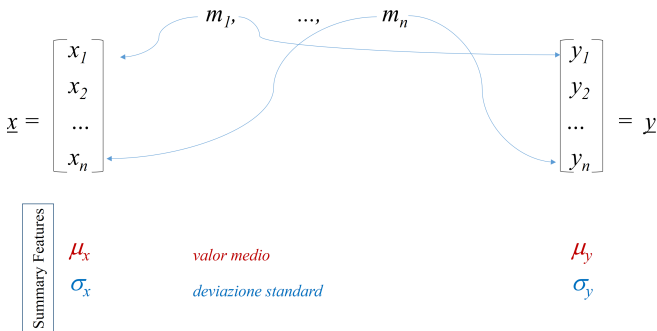
Very often objects in machine learning are described statistically, i.e. through the notion of distribution of probability that characterizes them: it serves to establish expectations about the values assumed by the object properties (e.g. how likely is 20 as the *age* of the instance of a “*young person*”).

Distances are this required to account for the likelihood that a value (e.g. 20) has with respect to others, and amplify (or decrease) the estimates according to such trends: this implies that non linear operators may arise and euclidean distances are not enough. Probability Theory and Information theory thus play a role in establishing some metrics that are useful in some Machine Learning tasks.

Pearson Correlation:

objects ($\underline{x}, \underline{y}$), mentors (m_i) and features (x_i, y_i)

Vectors \underline{x} and \underline{y} are derived from mentor judgments as follows.



As a consequence, summary features are other useful descriptors of collective attitudes of mentors towards objects x and y .



Pearson Correlation

Normalized Pearson Correlation

The [0,1]-normalized Pearson correlation can also be seen as a probabilistic measure as in:

$$\begin{aligned} nS^{(P)}(\underline{x}, \underline{y}) &\triangleq r_{xy} \triangleq \frac{\sum x_i y_i - n\mu_x \mu_y}{\sqrt{(\sum x_i^2 - n\mu_x^2)} \sqrt{(\sum y_i^2 - n\mu_y^2)}} \\ &= \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{(n-1)\sigma_x \sigma_y}, \end{aligned}$$

where μ_y denotes the average feature value of \underline{x} over all dimensions, and σ_x and σ_y are the standard deviations of \underline{x} and \underline{y} , respectively.

Jaccard Similarity

Binary Jaccard Similarity

The *binary Jaccard coefficient* measures the degree of overlap between two sets and is computed as the ratio of the number of shared features of \underline{x} AND \underline{y} to the number possessed by \underline{x} OR \underline{y} .

Example

For example, given two sets' binary indicator vectors $\underline{x} = (0, 1, 1, 0)^T$ and $\underline{y} = (1, 1, 0, 0)^T$, the cardinality of their intersect is 1 and the cardinality of their union is 3, rendering their Jaccard coefficient $1/3$.

The binary Jaccard coefficient it is often used in retail market-basket applications.



Extended Jaccard Similarity

Extended Jaccard Similarity

The *extended Jaccard coefficient* is the generalized notion of the binary case and it is computed as:

$$s^{(J)}(\underline{x}, \underline{y}) = \frac{\underline{x}^T \underline{y}}{\|\underline{x}\|_2^2 + \|\underline{y}\|_2^2 - \underline{x}^T \underline{y}}$$

Dice coefficient

Dice coefficient

Another similarity measure highly related to the extended Jaccard is the *Dice coefficient*:

$$s^{(D)}(\underline{x}, \underline{y}) = \frac{2\underline{x}^T \underline{y}}{\|\underline{x}\|_2^2 + \|\underline{y}\|_2^2}$$

The Dice coefficient can be obtained from the extended Jaccard coefficient by adding $\underline{x}^T \underline{y}$ to both the numerator and denominator.



Conditional-entropy

Conditional Entropy

the *conditional entropy* $H[\xi|\eta]$ of ξ and η is defined as:

$$\begin{aligned} H[\xi|\eta] &= - \sum_{j=1}^L p(y_j) \sum_{i=1}^M p(x_i|y_j) \ln p(x_i|y_j) = \\ &= - \sum_{j=1}^L \sum_{i=1}^M p(x_i, y_j) \ln p(x_i|y_j) \end{aligned}$$

Conditional and joint entropy

Conditional and Joint Entropy

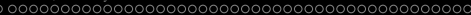
The conditional and joint entropies are related just like the conditional and joint probabilities:

$$H[\xi, \eta] = H[\eta] + H[\xi|\eta]$$

Conveyed Information

The *information conveyed* by η , denoted $I[\xi|\eta]$, is the reduction in entropy of ξ by finding out the outcome of η . This is defined by:

$$I[\xi|\eta] = H[\xi] - H[\xi|\eta]$$



Mutual Information

Given two random variable ξ and η :

Mutual Information

The *mutual information* between ξ and η is defined as:

$$\begin{aligned} MI[\xi, \eta] &= E\left[\ln \frac{P(\xi, \eta)}{P(\xi) \cdot P(\eta)}\right] = \\ &= \sum_{(x,y) \in \Omega_{(\xi, \eta)}} f_{(\xi, \eta)}(x, y) \ln \frac{f_{(\xi, \eta)}(x, y)}{f_{\xi}(x) \cdot f_{\eta}(y)} \end{aligned}$$

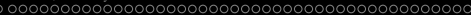
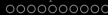


Mutual Information

Mutual Information measures the amount of information about a random variable ξ an observer receives when the outcome of a random variable η is known, in fact:

Mutual Information

$$\begin{aligned} MI[\xi, \eta] &= H[\xi] - H[\xi|\eta] = \\ &= \sum_{(x,y) \in \Omega_{(\xi,\eta)}} f_{(\xi,\eta)}(x,y) \ln \frac{f_{(\xi,\eta)}(x,y)}{f_{\xi}(x) \cdot f_{\eta}(y)} \end{aligned}$$



Pointwise Mutual Information

Another way to look to mutual information is about the individual values (i.e. outcomes) $\xi = x_i$ and $\eta = y_j$.

Pointwise Mutual Information

Given the two random variable ξ and η : the *pointwise mutual information* between $\xi = x_i$ and $\eta = y_j$ is defined as:

$$MI[x_i, y_j] = f_{(\xi, \eta)}(x_i, y_j) \ln \frac{f_{(\xi, \eta)}(x_i, y_j)}{f_{\xi}(x_i) \cdot f_{\eta}(y_j)} = P(x_i, y_j) \ln \frac{P(x_i, y_j)}{P(x_i) \cdot P(y_j)}$$

Pointwise Mutual Information

Pointwise Mutual Information (pmi)

$$MI[x_i, y_j] = P(x_i, y_j) \ln \frac{P(x_i, y_j)}{P(x_i) \cdot P(y_j)}$$

Use of the pmi

If $MI[x_i, y_j] \gg 0$, there is a strong correlation between x_i and y_j

If $MI[x_i, y_j] \ll 0$, there is a strong negative correlation.

When $MI[x_i, y_j] \approx 0$ the two outcomes are almost independent.

Cross-entropy as a Norm

Cross-entropy

$$H_p[q] = - \sum_{x \in \Omega_\xi} p(x) \ln q(x)$$

Relative Entropy (or Kullback-Leibler distance)

$$D[p||q] = \sum_{x \in \Omega_\xi} p(x) \ln \frac{p(x)}{q(x)} = H_p[q] - H[p]$$

Cross-entropy and Norms

Relative Entropy (or Kullback-Leibler distance)

$$D[p||q] = \sum_{x \in \Omega_{\xi}} p(x) \ln \frac{p(x)}{q(x)} = H_p[q] - H[p]$$

KL distance as a norm?

Unfortunately, as

$$D[p||q] \neq D[q||p]$$

the KL distance is *not* a valid metric in the classical terms. It is a *measure of the dissimilarity* between p and q .



Norms, Similarity and Learning

Why ranking probability distributions is necessary?

- During a learning process we need to figure out the circumstances (i.e. the state of affairs of the world) under which a certain concept/class/property manifest.
- This make a direct reference to the probability of some (stochastic) event. Stochastic events are used to describe circumstances and properties.
- Moreover, learning proceeds from experience, i.e. known facts or previous classified examples, to rules, i.e. probability joint distributions over *decisions* and *circumstances*
- Learning in general means **to induce the proper probability distributions from the known examples**. There are several many ways to do it!!!



Probability and Information References

Elementary Information Theory

- in (Krenn & Samuelsson, 1997), Brigitte Krenn, Christer Samuelsson, *The Linguist's Guide to Statistics Don't Panic*, Univ. of Saarlandes, 1997.

URL: <http://nlp.stanford.edu/fsnlp/dontpanic.pdf>