

Introducci' on the T'ecnicas Miner'ia Text



Maria-Amparo Vila
vila@decsai.ugr.es

Investigaci' on group Databases and Intelligent
Systems Informaci' on <https://idbis.ugr.es/>

Department of Sciences Computaci'
on artificial Intelligence University of
Granada

Outline other slideshow

one. The problem Miner'ia Text

two. I preprocessing

2.1 sint'actico preprocessing

2.2 sem'antico preprocessing

3. Reduccion of t'erminos (variables)

3.1 direct T'ecnicas

3.2 T'ecnicas based on measurements

3.3 T'ecnicas based on major components. latent semantics

Four. Regrouping and Text Mining

5. Classi f'icaci'on and Text Mining

6. Asociaci'on and Text Mining



The problem Miner'ia Text

The concept of "Miner'ia Text" (TM) is something that is in discusi'on today'ia. It podr'ia de fi ned as:

"Extracci'on process knowledge or patterns, previously unknown, non-trivial and interesting (potentially 'utiles) and understandable by users from of unstructured text documents. "

Text Mining is a data mining extensi'on where the discovery was made from unstructured databases.



The problem Miner'ia Text

should not be confused with Text Mining Recuperaci'on of Informaci'on from textual databases.

The recuperaci'on of informaci'on looking "documents" according to some requirements.

TM seek:

- unknown knowledge
- Understandable by users
- nontrivial
- Interesting



The problem Miner'ia Text

It is considered that the TM is' texts:

last step in processing



The problem Miner'ia Text

Podr'ia think is that you can directly apply the DM cl'asicas t'ecnicas the textual informaci'on. Nothing is further from reality . The DM works with databases with known schema. Each text document is an ordered colecci'on of words and signs of separaci'on meaning associated with situaci'on in the text which is determined by restrictions sint'actico and sem'antico type. There are semi-structured texts such as documents written in XML. In TM, the data include:

- inherently unstructured
 - impl'icita structure
 - Much greater wealth than in structured cases
- ambiguous
- multi-lingual



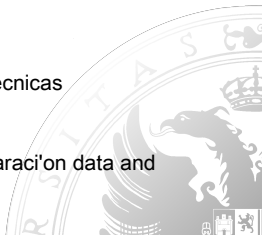
The problem Miner'ia Text

?This lack of structure is the biggest problem of the TM and involves the need to preprocess the texts, passing them to a *intermediately*

- Bags (bags) of t'erminos
- matrix structures (datasets)
- sem'anticas conceptual graphs or networks.
- Type structures "ontolog'ia"

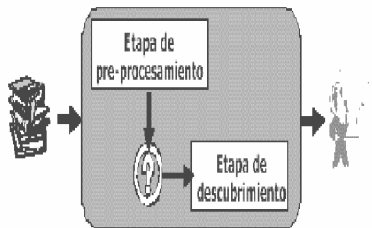
Preparaci'on phase data inherent in any process of extracci'on knowledge, it is crucial.

- The obtenci'on of the "parison" in turn involves use of extracci' t'ecnicas on knowledge,
- In many cases, it is not clear then, where the process ends preparaci'on data and begins to miner'ia.



The problem Miner's Text

A simple scheme of this search process:



The problem Miner'ia Text

?Some authors simply define the TM as an interdisciplinary field that includes elements of:

- Recover of information
- Extracci'on of informaci'on by LING computational uistica
- Grouping (clustering)
- Categorizaci'on,
- Other t'ecnicas DM



The problem Miner'ia Text

? All authors agree that est'an process TM includes the following phases:

one. I preprocessing

two. Miner'ia (proper)

3. Visualizaci'on In

funci' on the complexity of the preprocessing we find the possibility that this step occurs or not a complex intermediate form.



Text preprocessing Miner'ia

preprocessing Sint'actico

b'asica Idea

It is free text processing so that the output can be treated in an automated manner.

We must move from unstructured data structure data. At first **T'erminos bag**, possibly annotated

Stages in the processing sint'actico

one. Tokenizaci'on

two. Recognition / Eliminaci'on signs of puntuaci'on

3. Recognition / Eliminaci'on of "vac'ias words" (stopwords)

Four. M'recognition of words multiple. (N-grams)

5. Sint'acticos recognition types (POS)

6. Lematizaci'on (Steaming)



Text preprocessing Miner'ia

preprocessing Sint'actico

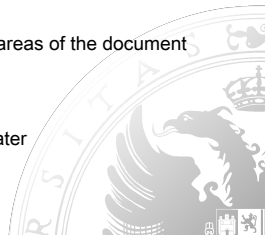
Tokenizaci'on

- This is the text from "tokens" strings that represent words.
- Several "tokens" represent the same t'ermino or word.
- In the end we will have bags of words with its occurrence in each t'ermino

- Common problems:

Prepare the document to extract text (PDF, HTML, XML). Different areas of the document may have to be treated differently.

Words "composite" requires preprocessing sem'antico we will see later



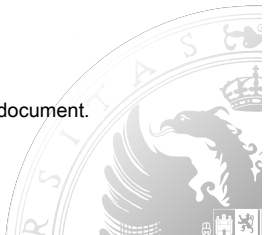
Text preprocessing Miner'ia

preprocessing Sint'actico

Recognition / eliminaci'on signs of puntuaci'on

Problems:

- It is detected when a sign is truly puntuaci'on
 - () ? ! , " Are delimiters and can be tokens.
 - ' ; , Y . They may form part of or be t'ermino delimiters. Habr'a likely to make an in-depth an'alisis if there are many signs.
- It is very dependent on language
- At the end of the process they are removed from the set of t'erminos document.



Text preprocessing Miner'ia

preprocessing Sint'actico

Recognition / Eliminaci'on of "vac'ias words" (stopwords)

- They are words that do not provide informaci'on from a point of view not LING
u'isitico
- They are essentially functional role
- Seg' are removed from the text a A given list
- It is very dependent on language

<http://www.navigla.es/posicionamiento-seo/palabrasstopwords-seo-espanol/>

You can find the list of words vac'ias espa~
Google's engine B' earch

not using



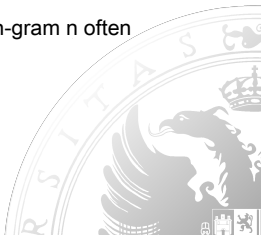
Text preprocessing Miner'ia

preprocessing Sint'actico

M'recognition of words

ultiple.

- There are groups of words that have meaning for themselves: Data base, etc. Operating System these groups should be treated as t'erminos. They can be regarded as a particular case of n-grams (n-grams)
- His autom'atica detecci'on is based on the idea that a ser'am'as n-gram n often together any words.



Text preprocessing Miner's

preprocessing Sint'actico

Recognition of words

multiple. To detect ngram one

possible algorithm:

one. Maximum set a value of n (2,3, ...) of words to be taken together. Setting a threshold measurement aparición n-gram: Minmes

two. For $k = 2..N$

2.1 k-sequences of consecutive words are discussed in the text T_k .

2.2 those are selected T_k such that:

$$AM(T_k) = k (\log_{10} \frac{\text{freq}(T_k)}{\text{freq}(\text{word}_i)}) \text{ minmes}$$

2.3 If there alg' a T_{kone} , $k > two$ let's be included in T_k remove $T(kone)$
the set of n-grams selected

Text preprocessing Miner'ia

preprocessing Sint'actico

Recognition / Eliminaci'on of sint'acticos types (Part of Speech, POS)

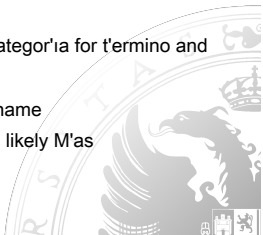
- For certain applications it is necessary to recognize the grammatical functions of t'erminos name, own name, verbs in different times etc.
- There are algorithms that allow you to tag the t'erminos in this regard.

LING algorithms based on rules u'isticas. M'as are old.

Autom'atico learning algorithms. They are based on recognizing a categor'ia for t'ermino and assign the next word his M'as categor'ia likely seg'

one language. PE in English ANN (adjective name
name) is a likely sequence in espa~ nol NNA is likely M'as

- Obviously they depend on the language



Text preprocessing Miner'ia

preprocessing Sint'actico

Recognition / Eliminaci'on of sint'acticos types (Part of Speech, POS)

Some categories of POS in ingl'es

Tag	Description
CC	Coordinating conjunction
CD	Cardinal number
DT	Determiner
EX	Existential there
FW	Foreign word
IN	Preposition or subordinating conjunction
JJ	Adjective
JJR	Adjective, comparative
JJS	Adjective, superlative
LS	List item marker
MD	Modal
NN	Noun, singular or mass
NNS	Noun, plural
POS	Possessive ending
UH	Interjection
VB	Verb, base form
VBD	Verb, past tense
VBG	Verb, gerund or present participle
VBN	Verb, past participle
VBP	Verb, non-3rd person singular present
WDT	Wh-determiner

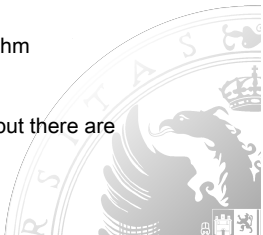


Text preprocessing Miner'ia

preprocessing Sint'actico

Stemming (Steaming)

- The different ways the same word usually problem'aticas for an'alisis word, as they have different ortograf'ia and similar meaning (for example, learn, learn, learning)
- we mean by *Lematizaci'on (Steaming or Lematization)* the process transform a word in his "standarizada root"
- For Ingl'es is not a big problem, there are algorithms available p'ublicamente that give good results. M'as is known Porter algorithm
- for espa~ nol no algorithm generally recognized as the Porter but there are several available algorithms.



Text preprocessing Miner'ia

preprocessing Sint'actico

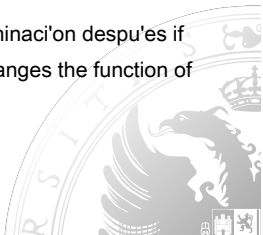
Stemming (Steaming)

- Is not always Useful carry out the lematizaci'on.

Overall for cases in which they work frequently t'erminos relizarlo interesting because it summarizes several t'erminos a s'olo and increases frequency of this.

But it can cause problems if a subsequent sem'antico preprocessing is done

- In any case the lematizaci'on habr'a to carry out labeling and eliminaci'on despu'es if the sint'acticas forms (POS) unwanted because lematizaci'on changes the function of the terms.



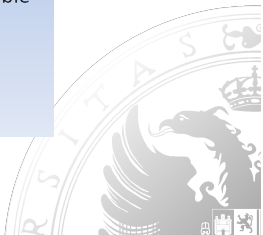
Text preprocessing Miner'ia

preprocessing Sint'actico

Stemming (Steaming)

Lematizaci'on rules in ingl'es

■ ATIONAL -> ATE	relational -> relate
■ TIONAL -> TION	conditional -> condition
■ ENCI -> ENCE	valenci -> valence
■ ANCI -> ANCE	hesitanci -> hesitance
■ IZER -> IZE	digitizer -> digitize
■ ABLI -> ABLE	conformabli -> conformable
■ ALLI -> AL	radicalli -> radical
■ ENTLI -> ENT	differentli -> different
■ ELI -> E	vileli -> vile
■ OUSLI -> OUS	analogousli -> analogous



Text preprocessing Miner'ia

preprocessing Sem'antico

b'asica Idea

Once clean and labeled the t'erminos relations sem'antico type are used to reduce them again.

Possible relationships between t'erminos

- **Synonymy** : differently, the same meaning (class lecci'on)
- **homonymy** : same meaning different form (bank financial institution, sit site)
- **Polysemy** : Similarly, different signi fi ed related (Bank Blood Bank)
- **hyponymy** A word is a subclass of another (dog, animal)



Text preprocessing Miner'ia

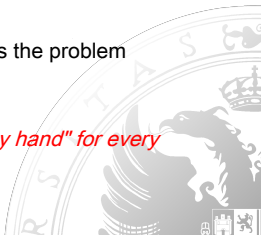
preprocessing Sem'antico

Desambiguaci'on

Process by which is assigned to multiple t'erminos s'olo one that has the mOsmol meaning that all of them

- There are tools that help you work with sin'onimos and hip'onimos, reaching assign each set of terms to a class semantics (desambiguaci'on).
- The famous M'as is Wordnet; but there is no algorithm that solves the problem completely and less in several languages.

The desambiguaci'on is an open problem and has to be made "by hand" for every problem, if needed



Reducci' on variables

b'asicas ideas

problems

- Colecci'on given a document each est'aj characterized by t'erminos.
- T'erminos assembly is called a colecci'on **dictionary** and may have several thousand items.
- The REPRESENTATION of each document is based on the dictionary
- T'ecnicas of reducci'on the aim of t'erminos *reduce n'umber of t'erminos Dictionary*



Reducci' on variables

b'asicas ideas

direct T'ecnicas

They are used directly in the preprocessing. Reduce Dictionary

one. The eliminaci'on of vac'ias words

two. M'recognizing words multiple

3. The eliminaci'on types of sint'acticos no signi fi cant

Four. the lematizaci'on

5. the desambiguaci'on



Reducci' on variables

Tecnicas based on measurements

b'asico process

- For each t'ermino k and each document i a "measure of importance" is set w_{ik} the t'ermino k in the document i .
- To measure the importance of a t'ermino in a document colecci'on set $w_k = \text{Agg}_n$

$i = \text{one } w_{ik}$ where Agg is a

AGGREGATION measure may be the m'aximo, m'inimo, average, sum etc.

- the seg' t'erminos are ordered w_k and selected the m first or a percentage of the total (ie 75% better) etc.



Reducci' on variables

Tecnicas based on measurements

Measures based on frequencies

- Be $n_{ik} = n_i'$ umber of times the t'ermino appears k in the document i , Y n_i
 n_i' umber of associated t'erminos i
- The frequency of each t'ermino k in i It is given by $F_{ik} = n_{ik}$
- If we work with frequencies take measures $w_{ik} = F_{ik}$. w_k
can be: $w_k = t_k = P_{ni=one} F_{ik}$ Total or frequency
 $w_k = P_{ni=one} F_{ik} / n$ average frequency.

$\frac{F_{ik}}{n_i}$



Reducing on variables

Técnicas based on measurements

Discrimination measures

issue

*Are frequent terms those who best represent a document collection ?
Terms no better represent certain documents appearing in some and not
in others, **discriminate** some of the documents?*

As $tf \times idf$ Yes d_k It is the n'

Number of documents in which it appears

k , the larger less discriminates k . It is defined then:

$$idf_k = \log_{two}(N / A_k) + one \quad Y \quad w_{ik} = tf_{ik} \times idf_{ik} = F_{ik} idf_k$$

Reducci' on variables

Tecnicas based on measurements

Discriminaci'on measures

Noise measurement For each t'ermino it is defined:

$$n_k = \sum_i (F_{ik} / t_k) \log_{\text{two}} (F_{ik} / t_k) \quad s_k = \log_{\text{two}} t_k n_k$$

$$\text{so } w_{ik} = F_{ik} s_k$$

Measures based on similarity of documents For colecci'on

document it defines the average similarity between them

, for each t'ermino k is defined s_k as the *average similarity suppressing t'ermino k documents*.

$s_k = k$ discriminator measures the power k .

$$w_{ik} = F_{ik} s_k$$

Reduccion variables

Técnicas based on principal component selection

Initial model: The vector model documents / terms

Every document is represented as a vector: $d = (w_1, \dots, w_M)$ where each w_i represents the weight that the term i in document d .

issue

Although it has made a *cleaning process* "empty words", elimination of synonyms etc ..., and *process frequency reduction* we have hundreds of variables, ie M can be very large. Much weights can be zero in a particular document, that is, the vectors are very large and "empty"

To do a process of reduction dimensions

The process reduction dimensional transforms a vector

(w_1, \dots, w_M) where M is very large in another (v_1, \dots, v_m) where m may be fixed.



Reducci' on variables

Tecnicas based on principal component selecci'on

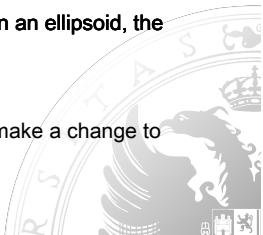
The problem

**Given a variable n dimensional $X = (x_{one}, \dots, x_n)$ Y m values
the same x_{ij} , i two $\{one, \dots, m\}$, j two $\{one, \dots, n\}$, find standard linear transformations
(SLC) that "summarize" the best possible data, capturing most variance thereof.**

intuitive idea

If the items are considered as a point cloud R^n , all can be enclosed in an ellipsoid, the average center, whose matrix is the covariance matrix.

The axes of the ellipsoid are a rectangular coordinate system, if we make a change to this coordinate system, the point spread along axis



Reducci' on variables

Tecnicas based on principal component selection

The model matem'atico

Be $\bar{\mu}$ the average $\bar{X}^T Y$ ^ its covariance matrix, is
 find a linear transformaci'on $\bar{y} = o(\bar{x} - \bar{\mu})$ such that new
 coordinate axes are the axes of the ellipsoid. It is proved to be a matrix such that:

$$O^T \Sigma O = \Lambda$$

where,

0BBBBBBB @

1CCCCCCCCA

one

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

n



Reducci' on variables

Tecnicas based on principal component selection

Vector values $\bar{x}_1, \dots, \bar{x}_n$ and are the
 "Eigenvalue" of the covariance matrix and associated with each λ_j there is a
 "eigenvector" \bar{y}_j that is the j th column of the matrix and verify c'andose:

$$\text{for } j \in \{1, \dots, n\} \text{ and } \lambda_j = \bar{y}_j^T (X - \bar{X}) (X - \bar{X})^T \bar{y}_j / \bar{y}_j^T \bar{y}_j$$

\bar{y}_j It is called j th principal component.

$$\text{for } k, j \in \{1, \dots, n\} \text{ Cov}(\bar{y}_j, \bar{y}_k) = 0 \text{ Var}(\bar{y}_j) = \lambda_j$$

$$\text{Var}(\bar{y}_1) \dots \text{Var}(\bar{y}_n)$$



Reducci' on variables

Tecnicas based on principal component selecci'on

Proporci'on explained variance

?The proporci'on of variance explained by k factors is

$(1 + \dots + k) / (1 + \dots + n)$ and it allows us to reduce the dimensionality of space. That is express fen'omeno less variables.

How many components make ?:

- At least 90% of variance explained
- All eigenvalues that are greater than the average of the same.
- If the matrix is used instead Correlation of the covariance matrix eigenvalues greater than 1.

The proporci'on of variaci'on explained variable j by k component is given by $\rightarrow (x_j, Y_k) = r_{jk}$, with these values can identify components and identify one semantics for them.

Reducci' on variables

Tecnicas based on principal component selection

New REPRESENTATION of documents y'terminos

The expresi'on:

$$\bar{y} = o(\bar{x} - \mu)$$

We can get:

$$\bar{x} = \kappa - o$$

So that, if κ matrix represents the new "factors", then

REPRESENTATION represents each document in the new space factor and the REPRESENTATION each t'ermينو in funci' on of the factors. The same occurs if the first k elements are taken κ , the first k columns and the first k columns.

Reducci' on variables

T'ecnicas based on principal component selecci'on

An'alisis latent semantics

In our case the model identi fi ca c'omo:

- The items are documents
- T'erminos weights in each document are the variables
- New combinations of factors are interrelated terms
- Representations of documents funci'on new factors can be used in many different ways



Reducci' on variables

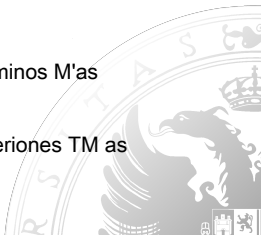
T'ecnicas based on principal component selecci'on

An'alisis latent semantics

- the fi c gr'a representations according factors can give an idea of a first cluster of documents
- T'erminos representations in funci'on factors can give an idea of the semantics of them. To the extent that some authors speak of new "concepts" identi fi ed by the factors.

Previous experiences show that:

- They can be reduced to 200 factors, several thousand spaces t'erminos M'as retaining 90% of the variance
- The reducci'on dimensions substantially not affect processes posteriores TM as clustering



Regrouping and TM: some approaches

cl'asico model

- Input data: *vector or vector model reduced*
- Distance used: *The measure cosine*

Since each document is a vector is calculated cosine'angulo forming. Yes $t_1 = (w_{eleven...} w_{one d})$ Y $t_2 = (w_{twenty-one...} w_{two d})$ are two vectors, then:

$$\cos(t_{one}, t_2) = (t_{one} \cdot t_2) / |t_1| |t_2|$$

where \cdot represents the scalar product and $|\cdot|$ the modulus, ie:

$$\cos(t_{one}, t_2) = \frac{\sum_{j=1}^n w_{one j} w_{two j}}{\sqrt{\sum_{j=1}^n w_{one j}^2} \sqrt{\sum_{j=1}^n w_{two j}^2}}$$

- If a factor model used can be used without problems eucl'idea distance

Regrouping and TM: some approaches

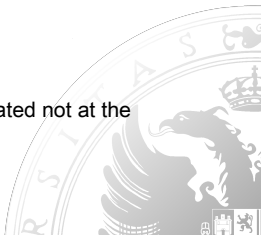
?It is very common to use the 's method k-means with cosine distance. some variants that allow improved m'etodo are used:

-Selecci' on centroids by

- An initial cluster jer'arquico a set peque~ documents. do not give
- An'alisis in previous and subsequent components REPRESENTATION
- Other t'ecnicas fractionation set of documents

- Utilizaci'on improvement of the k-means:

- 'S method k-means continuous continuously centroids are calculated not at the end of each step



Regrouping and TM: some approaches

- Utilization improvement of the k-means:

- 'S method k-means in two phases, where a refinement of the process is performed at each stage of the algorithm to avoid falling into "local optimum" of the function of coherence.
- 'S method k-means by bipartition.

? There are approaches that recommend the use of more advanced grouping techniques such as Medoids (CLARANS) or methods based on density, advanced hierarchical methods (BIRCH).



Regrouping and TM: discovery events

issue

We are surrounded by texts "have events", ie "news". The discovery of events in the context of news is the identification of "stories" that correspond to new events or previously identified. The detection of events can be done in two ways:

- Retrospectively
- on line

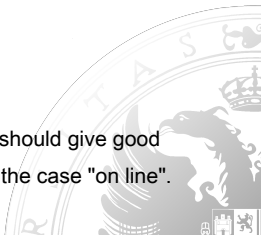


Regrouping and TM: discovery events

New events occur continuously and news relating to them are usually in the form of "blowout" or "explosions" with the following characteristics:

- one. You are grouped in time,
- two. The "explosion" associated with new events are separated for a certain amount of time,
3. Every new event is usually accompanied by a change in the type of terminos and vocabulary used to describe the associated news,
- Four. The events usually take a time window of between one and duration four weeks.

These four observations seem to suggest that any technical grouping should give good results in detection new events, both in the retrospective case as in the case "on line".



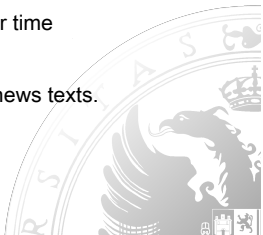
Regrouping and TM: discovery events

Model

- Each news is an "object" characterized by a set of attributes (t'erminos)
- Each event is a "prototype" of a set of similar news funci' on their t'erminos.

In the case of retrospective detecci'on "Explosions" are looking at a hist'orico file by:

- ↑ Sudden changes in distribuci'on of t'erminos over time
- ↑ L'exicas similarities and proximity in time of the news texts.



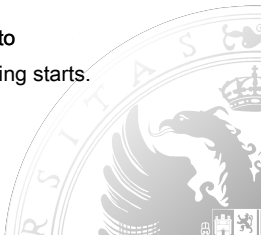
Regrouping and TM: discovery events

In the detection "on line" the idea is to process the news as they arrive. The algorithms have been developed for this problem are "boolean" answer:

one. Each story is compared, when it appears, with "Clustering" already detected.

two. If the news can join this grouping
It is assumed to correspond to the same event.

3. If the news can not be incorporated is supposed to
It corresponds to a new event and a new grouping starts.



Classi fi caci'on text

b'asicas ideas

general model classi fi caci'on

- Let's consider *The vector model documents / t'ermynos*

Every document is represented as a vector: $d = (w_{one}, .. w_N)$

where each w_i represents the weight that the t'ermyno i in document d .

- Suppose we have M documents $\{d_i, i = one.. M\}$
- A set of H classes provided for documents.
- The problem of classifying documents (categorize) *It does not seem very different from the problem of classifying General caci'on*



Classi fi caci'on text

b'asicas ideas

t'ipicos problems classified text caci'on

Document Categorizaci'on

We have a document colecci'on classi fi ed by t'opicos and we want to train a categorizer way that allows us to classify new documents. They tend to be "long" documents.

problems

- Car initially classified documents. Usually you must become an expert.
- The t'opicos documents often change over time.



Classi fi caci'on text

b'asicas ideas

t'ipicos problems classified text caci'on

M'as problems with short texts that include a word catgorizador .

- Classi fi er news (Google News for example). Est'a detecci'on related to the event. In some cases it is pre-clustering and working with jerarqu'ias
- Electr'onicos spam filters emails
- An'alisis feelings on social networks.
- Studies short texts categorized database (descriptive an'alisis):

m'edicas descriptions. Diagn'osticos etc. Opini'on surveys with free text.

Miner'ia Web usage and content (an'alisis of t'opicos on Twiter yh'abitos of use etc ..)



Classi fi caci'on text

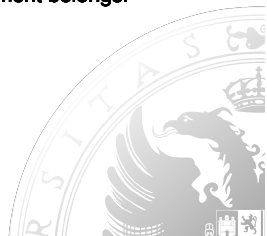
b'asicas ideas

general model classi fi caci'on

- We dataset:

items variables	t_{one}	t_{two}	\dots	t_N	C
d_{one}	W_{eleven}	W_{12}	\dots	$W_{one\ N}$	C_{one}
\vdots	\vdots	\vdots	\dots	\vdots	\vdots
d_M	$W_{M\ one}$	$W_{M\ two}$	\dots	$W_{M\ N}$	C_M

- It is predicted, with the values of $\{w_{ij}\}$ class c_i to which the document belongs.
- In principle the process of classi fi caci'on not vary:
 - one.** We have a training set and a set of test
 - two.** We apply a model classi fi caci'on the training set
 - 3.** We apply the model to predict the test set.
 - Four.** We measure how well the model



Classificazione text

specific problems classification

text on

issue

Vector REPRESENTATION documents:

- one.** In principle it contains many variables. It is impossible to think, principle in discrete techniques as spanning trees of decision
- two.** The matrix documents / terms is scattered.

Necessary reduce drastically terms space

Solution

- one.** Reduction apply principal component.
- two.** Establishing processes classify binary classification and reduce dictionary the representative Matrix as terms class to consider



Classification text

specific problems classification

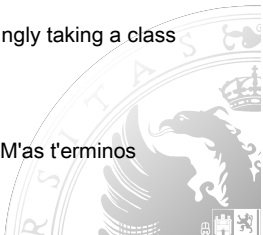
text on

Application problems of major components

- interpretation the semantics of the classification is lost. All methods wouldn't be predictive in this case.
- Classification of each document involves a process prior transformation. *REPRESENTATION is necessary to use the document in the space factor*

Processes binary classification

- If there is class H H processes classified develop classification increasingly taking a class as positive examples and the rest as negative.
- Every time the dictionary of the process considering the relevant terms class we are studying is changed.



Classi fi caci'on text

speci fi c problems classi fi caci'

text on

Processes binary classi fi caci'on

- The process can lead to classi fi ed each document belongs to a class M'as:

belonging to several classes remains (B'

Index search t'opicos)

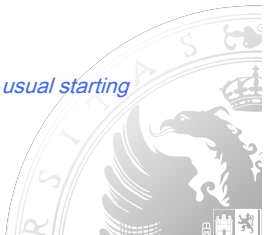
It is assigned somehow the "winning" class (probability measures, distance measurements etc.)

Usually working with classi fi cations binary. The adaptaci'on processes classified caci'on take into account this fact

Processes classified caci'on word used M'as t'ecnicas taking the usual starting the vector model documents / t'erminos. Seg'

one type is usually t'ecnica

*use the presence / absence or tf * idf model*



Classificaci'on text

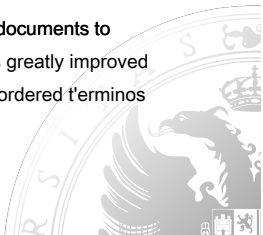
Aplicaci' t'ecnicas on the classifi caci' on a word

Using the KNN t'ecnica

- Overall similarity measures are used instead of distances:

If it is a binary pattern similarity measures if it is a model $tf * idf$ measure cosine.

- **The main problem is computational aplicaci'on** : For each document must calculate the similarity with everyone.
- **Soluci' on** : Aplicaci'on of t'ecnicas of recuperaci'on of informaci'on All documents to classify a "consultation" k-documents that best satisfy the process is greatly improved with the use of inverted indices are obtained (t'ermino / documents) ordered t'erminos

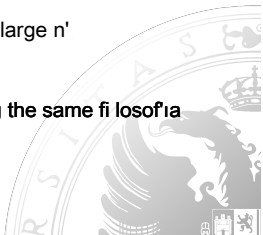


Classificaci'on text

Aplicaci' t'ecnicas on the classifi caci' on a word

Using rule-based t'ecnicas

- As explanatory m'etodos est'an M'as get oriented rules describing the classes to classify.
- Commonly used binary data
- You should try to work with simple rules that prevent overlearning. This involves using t'ecnicas pruning.
- Experimental data indicate that it is not necessary to work with a large n'umber of starting t'erminos
- It is also possible to work with **spanning trees of decisi'on** following the same fi losof'ia above.



Classification text

Applicazioni tecniche on the classification on a word

Using Probabilistic Techniques

Naive Bayes Method

We want to predict a binary class variable N random binary

(t_1, \dots, t_N), or N -dimensional binary variable

the conditional probability, lead to:

t . The expressions

$$Pr(C | t = x) = \frac{Pr(C) \prod_{j=1}^N Pr(t_j = x_j | C)^{w_j x_j + b}}{Pr(t = x) \exp(X_j w_j x_j + b)}$$

where w_j and b They can be calculated from the estimated probabilities. $Pr(t_j = 1 | C)$, $Pr(C)$ etc.

There are two models of computation: Bernoulli and multinomial

Classification text

Applications techniques on the classification on a word

Using Techniques Probabilistic

Regression on Logistic

Since we are classified in binary form you can use tf * idf and work with models
logistic. Logistic regression coefficients of regression

Logistic regression can give us an idea of the weight of the classified terms classification.
Especially if it has done in steps.



Classification text

Aplicaci' t'ecnicas on the classi fi caci' on a word

GLM .

For each document D to classify its puntuaci'on is calculated to belong to a given class of linear form:

$$\text{punt}(D) = \sum_j w_j x_j + b$$

x_j is 0 or 1 in the document and w_j They are the "weights of each t'ermino in class". The process of classifying the documents caci'on is very r'apido with this m'etodo:

one. the t'erminos are arranged with their weights (inverted list)

two. D for each list is scrolled by adding the weights of t'erminos

It is in D, $\text{sum}(D)$

3. $\text{sum}(D) > \epsilon > 0$? the document is assigned to the class.



Classificaci'on text

Aplicaci' t'ecnicas on the classifi caci' on a word

GLM .

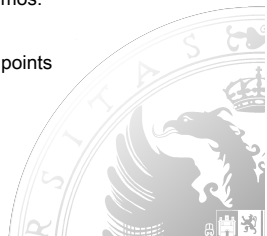
- The process of classifying M'as r'apido caci'on is possible.
- They may extend lists weights using, sin'onimos, n-grams etc.

- **Problem: learn the weights**

Direct optimizaci'on leads to a non-convex porblema with local m'inimos.

It works to find the hyperplane separating the positive and negative points

It is normally used an SVM



Association rules and TM

An interesting example

- T : *términos in a document*
- D : *word set*
- $\{pain\} \rightarrow \{death\}$

Basic Idea

A set of documents associated with términos can be viewed as a transactional database:

items variables	t_{one}	t_{two}	t_N
d_{one}	1	0	one
\vdots	\vdots	\vdots	\vdots
d_M	0	0	one

Association rules and TM

Associations and co-occurrences between a text characteristics

Formulation

$T = t_1, t_2, \dots, t_n$ set of terms

$D = d_1, d_2, \dots, d_m$ set of documents indexed by terms

Each document d_i generates a subset $d_i(t) \subseteq T$.

Be $OR \subseteq T$. The set of all documents D such that d

$OR \subseteq d(T)$ It is called *covering assembly* U and $[U]$ shows. In its simplest form a rule of association is an implication of the form:

$U) \text{ or with } OR \subseteq YOU$

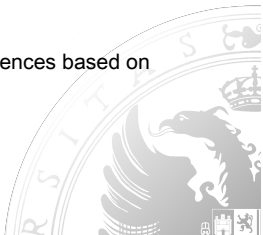


Association rules and TM

There are several approaches to the problem depending on how T is defined:

- T is the set of all documents. $d(T)$ is the document REPRESENTATION t as a "pocket" word. This approach is not very efficient and computationally expensive
- T is a set of keywords in documents. $d(T)$ is the document INDEXING T through a set of keywords.

In this approach the problem is to obtain T . Techniques of co-occurrences based on natural language processing can be employed



Association rules and TM

- The elements of T are grouped according to classification one taxonomic. (Thesaurus)
This approach allows you to discover partnerships with different level of granularity
- T elements are phrases or full sentences. In this case it is discovered co-occurrences of sentences in documents collection.
- T elements are "episodes" (complete sets of vector characteristics).

In this case the so-called "rules of episodes" are obtained.



Extracci' on attributes of semantics in short texts

issue

There attributes database (relational or not) corresponding to short text with a restricted semantics.

Examples:

- Databases M'edicas (diagn'osticos, descripci'on interventions etc.)
- Fields "observations" in various situations (surveys, expert etc.)
- Document Abstract

It is desired to obtain a REPRESENTATION of the attribute (intemedia shape) such that it:

- Collect semantics
- allow the query as a model ROO (where applicable)
- allows performing TM on said structure



Extracci' on attributes of semantics in short texts

b'asicas ideas

- Once "clean" text, sets / sequences repeated words in n'
UMBER sufficient to constitute tuples
"Sentences" or "concepts" system.
- It is possible that some "sentences" are repeated incomplete, corresponding to less specialized concepts, so that we have a subset of t'erminos reticulum.
- This leads us to the concept of "itemset (seq) frequent" in transactional database of texts. Remember that every frequent itemset has the property "a priori" Then

If you get frequent itemset base texts:

- *The maximal itemset give us complete sentences*
- *The subret'iculo gives us the full semantics structure of the database.*



Extracci' on attributes of semantics in short texts

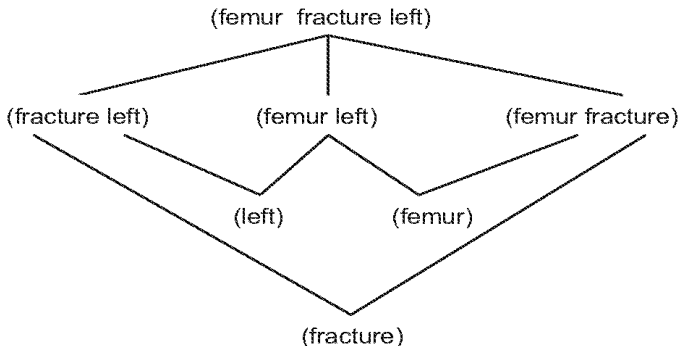
Example: database ED

Field 1	Field 2	Field 3	...	Support
right	fracture			4.26
right	femur			1.826
back	sebaceous	cyst		1.304
back	cyst			1.304
right	femur	fracture		1.173
femur	left			1.043
femur	fracture	left		1.043
.	.			.
.	.			.
.	.			.

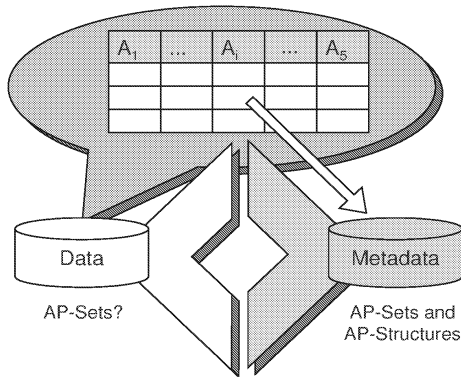


Extracci' on Semantics in short text attributes

Example: database ED



Extracci' on Semantics in short text attributes



Extracci' on Semantics in short text attributes

↑ The overall semantics of a set of short texts is collected in a structure derived from whole frequent itemset

↑ The value of each tuple seg' a text containing a substructure global

↑ Are de fi ned properties and operations ap-structures that allow:

- Get the value of each tuple in the attribute as substructure
- Approximate perform queries using this new data structure

↑ This structure is implemented as a TDA in a BDROO

↑ We have developed TM with these structures, studying models:

- Data Warehousing
- And obtaining ordered itemset GENERATION autom'atica tag-clouds

