

Investigating word interactions in texts. Application to text categorization in genomics

Martine Cadot, Michel Zitt, Gabriel Meurin, Alain Lelu

► **To cite this version:**

Martine Cadot, Michel Zitt, Gabriel Meurin, Alain Lelu. Investigating word interactions in texts. Application to text categorization in genomics. First SaarLorLux Workshop on Systems Biology 2009, Computational, Structural and Medical Approaches for Systems Biology, Dec 2009, Nancy, France. <inria-00442395>

HAL Id: inria-00442395

<https://hal.inria.fr/inria-00442395>

Submitted on 21 Dec 2009

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Investigating word interactions in texts. Application to text categorization in genomics.

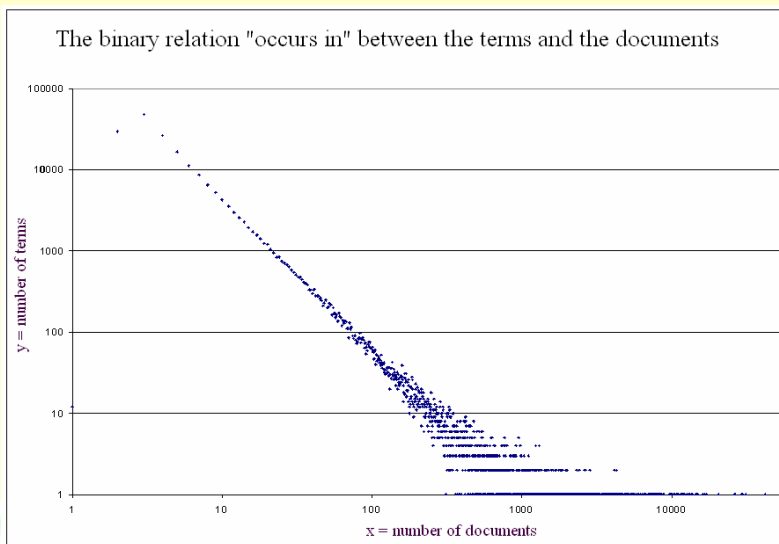
Martine Cadot (LORIA), Michel Zitt (INRA, OST), Gabriel Meurin (LORIA/INRA), Alain Lelu (LORIA, LASELDI)

Words interacting in a text may be compared, to a certain extent, to molecules interacting and building "complexes", i.e. phrases, named entities, or longer-range semantic or syntactic associations. We will call them "k-itemsets", k being their interaction level. We have shown (Cadot 06) that an adequately built subset of these k-itemsets is enough for describing the entirety of the relations at work in a corpus represented as a set of "bag-of-words" documents, whatever the level k of these relations.

Our objective is to reconstruct each category into which a corpus of scientific abstracts has been split, using a set of Boolean queries as a best compromise between conciseness and reproducibility of this categorization.

I - The corpus, its 50 categories.

- Scientific abstracts in genomics, pulled out from the Web of Science (Thomson Scientific ed.).
- This subset has been delineated using a hybrid method, based both on lexical queries and citation expansion/ shrinkage (Zitt et al. 2006) → 120,000 abstracts from 1999 to 2005.
- A vocabulary of 237,000 lemmatized words and phrases (>2 occurrences) has been pulled out and filtered (NeuroNav, www.diatopie.com). I.e.: *sequence, polymorphism, folded_structure, chromosome_4B, greenbug_resistance_gene,...*



This figure reads: e.g. 1047 terms occur each one exactly in 21 documents.

- 50 categories resulted from a clustering of the abstracts by the Axial K-means method (Bassecoulard et al. 2007).

M1/ Human_genome/ Human_genome_project	M17/ Map/ Linkage_maps/ Polymorphism	M33/ RNA- Virus
M2/ Translocation/ FISH/ leukemia	M18/ Population_genomics	M34/ PCR/ Methods/ applications
M3/ Plant_genomics/ Transgenic_plants	M19/ Repair/ DNA_damage	M35/ C-DNA/ Transcription/ C-DNA_library
M4/ DNA_sequence/ Satellite	M20/ Resistance/ Resistance_genes/ Plant & Fungi_resistance	M36/ Polymorphism
M5/ Strain/ Microbial_genomics	M21/ Hybrid/ Somatic_hybrids/ Fertility	M37/ Cell/ DNA_damage
M6/ Cell_identity & Gene_expression	M22/ Human/ C-DNA/ Gene_annotation	M38/ Genome/ Genome_sizes
M7/ Enzyme/ Escherichia_Coli	M23/ Exon/ Genomic_organization/ Gene_annotation	M39/ DNA/ Arrays/ Genomic_techniques
M8/ Alignment/ Bioinformatics	M24/ System/ Systems_biology/ Bioinformatics	M40/ QTL/ Trait/ Mapping/ Polymorphism
M9/ Genome	M25/ Patient/ Disease_genomics/ Biomarkers/ Pharmacogenomics	M41/ Signaling/ Kinase/ MAPK
M10/ Comparative_genomic_hybridization/ Tumor	M26/ Virus/ Nucleotide_sequence	M42/ Mutation/ Missense_mutation
M11/ SNPs/ Polymorphism	M27/ Evolution/ Evolutionary_genomics	M43/ Mouse/ Murine_genomics
M12/ Network/ Biological_networks/ Model	M28/ Cancer/ Genome & cancer	M44/ Expression/ Cell_identity & Gene_expression
M13/ Transcriptional/ Saccharomyces_cerevisiae/ Transcriptome	M29/ Promoter/ Transcription	M45/ LOD/ Linkage_analysis/ Polymorphism
M14/ Locus/ Microsatellite_locus/ Polymorphism	M30/ Mutant/ Mutagenesis	M46/ Human/ Primate/ Gene_annotation/ Comparative_genomics
M15/ Cell_line/ Tumor/ Genome & Cancer	M31/ LOH/ Tumor_suppressor/ Genome & Cancer	M47/ Species/ Phylogeny/ Evolutionary_genomics
M16/ Spectrometry/ Proteomics	M32/ Marker/ RAPD/ AFLP/ Polymorphism	M48/ C57BL/ Congenic_strains/ Murine_genomics
		M49/ Residue/ Amino_acid_sequence
		M50/ Virus/ Virus_replication/ Virus_recombination

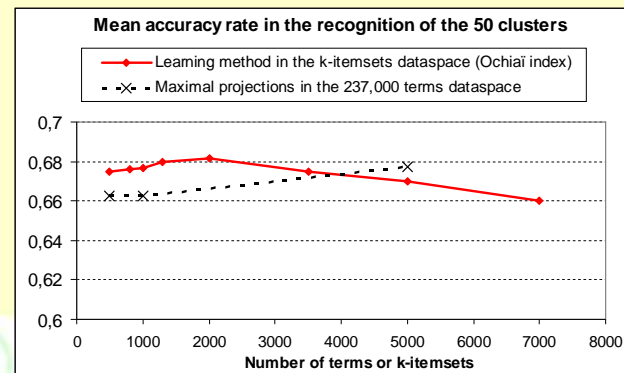
Ex. of the 14 first terms (words and phrases) most typical of the « Human Genome Project » cluster:

No	Phrase	Ochiai	8	human_genome_sequence	0,122
1	human_genome	0,467	9	primate	0,107
2	human	0,336	10	sequence	0,103
3	genome	0,257	11	chimpanzee	0,099
4	human_genome_project	0,238	12	completion	0,096
5	project	0,157	13	disease	0,094
6	draft	0,138	14	genomics	0,092
7	human_chromosome	0,125	15	human_gene	0,092

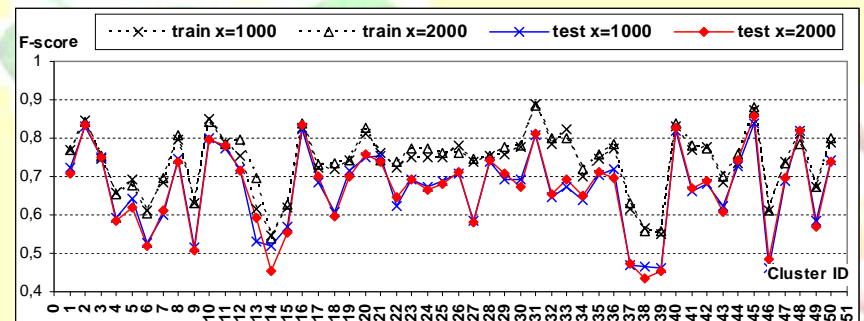
As most of data analysis methods do, this data partition takes into account the only « 2-itemsets » (a k-itemset of support s is an elementary association of k terms present in s documents).

II - Concise and reproducible representation of the 50 categories: using itemsets of order 1, 2, and higher order ones, which express complex interactions between terms in specific contexts.

- MIDOVA method (Cadot 2006) for mining ordered lists of informative itemsets specific to each category (train set: 1/10th of the corpus, test set: 9/10th) → ordered lists of simple Boolean queries for identifying the class of any document out of the corpus, and extending this categorization process to other databases (patents, ...).
- Results
 - At the same time as a 100% intrinsically exact reconstitution rate results from using the whole 237,000 terms, a maximum 68% rate results from using about x = 2000 k-itemsets (and then decreases), with a statistically-controlled generalization ability:



- 68% is a mean value, embedding many discrepancies:
 - >80%: clusters N°45 (LOD/ Linkage_analysis/ Polymorphism) and N°16 (Spectrometry/ Proteomics)
 - <50%: clusters N°14 (Locus/ Microsatellite_Locus/ Polymorphism) and N°38 (genome/ genome_size)



- example of class description: the 6 first k-itemsets of cluster 10:

comparative_genomic_hybridization, hybridization AND tumor AND genomic, hybridization AND tumor, CGH AND comparative_genomic_hybridization, losses AND genomic, tumor AND genomic.

As may be observed *comparative_genomic_hybridization* and *CGH* acronym appear together, integrally or partly, at the top of the list.

References

Bassecoulard E., Lelu A. and Zitt M. (2007). Mapping nanosciences by citation flows: a preliminary analysis, *Scientometrics*, vol 70, n°3, pp. 859-880.

Cadot M. (2006). *Extraire et valider les relations complexes en sciences humaines : statistiques, motifs et règles d'association*. Doctoral dissertation, University of Franche-Comté, France.

Laurens P., Zitt M. and Bassecoulard E. (to appear), "Delineation of the genomics field by hybrid citation-lexical methods: interaction with experts and validation process", *Scientometrics*.

Zitt M., Ramanana-Rahary S. and Bassecoulard E. (2006). Delineating complex scientific fields by a hybrid lexical-citation method: an application to nanosciences, *Information Processing and Management* Vol. 42-6, pp. 1513-1531.