


Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 00000 Conclusion

Strategies in Identifying Issues Addressed in Legal Reports

Gilbert Ritschard¹ Matthias Studer¹ Vincent Pisetta²

¹Dept of Econometrics, University of Geneva
<http://mephisto.unige.ch>
²ERIC Laboratory, University of Lyon 2

Compstat 2008, Porto, Portugal, August 24 - 29



20/8/2008gr 1/45

Introduction ●000 The Text Mining Process 000000 Text representation 00000000 Learning 00000 Conclusion

Aim of presentation

- Automatic identification of issues reported in texts.
- Experience with **reports on application of ILO Conventions**.
- Describing the text mining process
 - Quantitative representation of the texts.
 - Learning rules for predicting issues reported by any given text.

20/8/2008gr 4/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 00000 Conclusion

Project on Social Dialogue Regimes

- Financially supported by the Geneva International Academic Network (GIAN).
- Joint project between
 - Depts of Econometrics and of Sociology (U. of Geneva),
 - ERIC (U. of Lyon 2)
 - International Institute of Labour Studies (ILO, Geneva).
- Analysis of the determinants and socioeconomic correlates of Social Dialogue Regimes (SDR)
 - Sociopolitical regimes in which workers have freedom to establish organizations of their own choosing, negotiate collectively over working conditions, and participate through their associations in the design and implementation of policies that affect their lives

20/8/2008gr 6/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 00000 Conclusion

The Objectives of Text Mining

- Focus on CEACR comments from 1991 to 2002.
 - CEACR = Committee of Experts on the Application of Conventions and Recommendations
- Creation of synthetic indicators of legal rights.
- Use of **indicators for aggregate data analysis**, aimed at exploring how particular violations are (or are not) linked to others, as well as relationships with socioeconomic indicators.

20/8/2008gr 7/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 00000 Conclusion

Why resorting to text mining ?

- Extract useful information from huge number of expert comments (~1200 reports).
- Two main goals
 - Assist legal experts in the search of relevant information, speed up the process.
 - Produce **indicators for synthetic aggregate analysis**.

20/8/2008gr 8/45

Introduction 0000 **The Text Mining Process** ●00000 Text representation 00000000 Learning 00000 Conclusion

What is text mining ?

- Process of analyzing text to extract information useful for particular purposes.
- More than indexing or search engine.
- Aims at discovering knowledge about :
 - Content
 - Structure
 - Semantic
 - Ontology : typical terminology grouped into concepts and organized into conceptual hierarchies
 - ...

20/8/2008gr 11/45

Introduction 0000 The Text Mining Process 000000 Text representation 0000000 Learning 00000 Conclusion

Different text mining usages

- Statistical analysis of used words or sentence structure.
- Relationships between texts (Which texts are similar to a given one?)
- Summarizing automatically articles or documents.
- Unsupervised text categorization (clustering).
- **Supervised categorization** (detecting spam).
- Technological watch.
- Building ontologies (finding typical terminology and organizing it into conceptual hierarchies).
- Retrieving concepts from texts.
- **Information Retrieval.**
- ...

20/8/2008gr 12/45

Introduction 0000 The Text Mining Process 000000 Text representation 0000000 Learning 00000 Conclusion

Challenges in text mining

- Text are unstructured data
- Polysemy
 - "Mining expert comments"
 - "Mining expert comments"
- Synonymy
 - "Trade union", "Workers' organisation"
- Inflected forms, stop words, ...
- ⇒ **Requires pre-processing**

20/8/2008gr 14/45

Introduction 0000 The Text Mining Process 000000 Text representation 0000000 Learning 00000 Conclusion

The retained approach

- Want a tool that can be used by any non text mining expert.
- ⇒ **No pre-processing** (grammatical tagging, lemmatisation, stemming) **in the application stage.**

20/8/2008gr 16/45

Introduction 0000 The Text Mining Process 000000 Text representation 0000000 Learning 00000 Conclusion

Two main steps

- 1 Representing texts by a set of quantitative variables (**whole corpus**):
 - Extracting useful terminology.
 - Grouping terms into reduced number of descriptor concepts.
 - Quantifying descriptor concepts (tf×idf).
- 2 Learning prediction rules (**learning sample**):
 - Classification trees

20/8/2008gr 17/45

Introduction 0000 The Text Mining Process 000000 Text representation 0000000 Learning 00000 Conclusion

Retained key concepts, i.e. types of violation

Whose presence is the attribute to predict

- Original list of 27 key concepts (types of violation)
- Merged into the following 9 key concepts for Convention 87

v ₁	Right to life and physical integrity (not observed)
v ₂	Right to liberty and security of person / Right to a fair trial (not observed)
v ₃	Right to establish and join workers' organizations
v ₄	Trade union pluralism
v ₅	Dissolution or suspension of workers' organizations (not observed)
v ₆	Election of representatives / Eligibility criteria
v ₇	Organization of activities / Protection of property / Financial independence
v ₈	Approval and registration of workers' organizations
v ₉	Restrictions on the right to industrial action

20/8/2008gr 18/45

Introduction 0000 The Text Mining Process 000000 Text representation 0000000 Learning 00000 Conclusion

From text to quantitative representation

Defining the descriptor concepts

- Two major approaches
 - *n*-grams (hard to account for semantic)
 - **Bag of words**

```

    graph LR
      A[Bag of words] -- "Domain expert  
dedicated software" --> B[Useful terms]
      B -- "grouped into" --> C[Descriptor concepts]
  
```

trade union action	trade union pluralism	17 for Conv. 87
...	trade union activity	9 for Conv. 98

20/8/2008gr 21/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion 0000

Defining descriptor concepts

With help of legal experts

- Whole corpus (1 file), Grammatically tagged (TreeTagger)
- Process based on both :
 - Statistical characteristics (discriminating terms, co-occurrences)
 - Similarity of meaning (according to domain expert)
- 3 steps :
 - Preliminary set of concepts (by grouping new terms with already extracted most frequent ones, EXIT)
 - Extensional induction with legal expert (preliminary list is adapted according to knowledge base)
 - Amended list is again compared with corpus (find infrequent terms with meaning similar to defined concepts)
- **Result** : Concepts, each defined by a list of terms.

20/8/2008gr 22/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion 0000

Retained descriptor concepts for Convention 87

c1	Life and physical integrity	c10	Industrial action
c2	Liberty and security of persons	c11	Essential service
c3	Property and financial independence	c12	Arbitration
		c13	Strike action
c4	Service	c14	Union establishment limitations
c5	Pluralism		
c6	Election	c15	Specific workers
c7	Opinion and expression freedom	c16	Number of workers
c8	Restrictions on trade union activities	c17	Supervision
c9	Trade union approval		

20/8/2008gr 23/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion 0000

Examples of descriptor concept lists

Election (c6)	Restrictions on trade union activities (c8)
electoral procedure/ representative member/ trade union office/ eligibility/ re-election/ representatives of organization/ representatives of trade union/ vote/ elect member/ union leader/ elect their representative/ elect their own representative/ to elect representative/ to elect representatives/ elect its representatives/ elect their trade union representative/ election of trade union representative/ elected representative/ elected workers' representative/ election/ elect freely their representative/ elected their representative/ union officer/ union office/ representatives of association/ representatives of union/ representatives of workers' organization/ ...	trade union activities/ right to organize/ right of trade unions to organize/ right to publication/ right to assembly/ right to disseminate information/ freedom of opinion/ freedom of expression/ political opinion/ political activity/ hold meeting/ right to organize/ right of association/ right of workers organization/ right of workers to organize/ right to hold trade union/ holding office/ formulate their programmes/ right of organizations to organize/ right of first-level unions to organize/ right of unions to organize/ right of workers' organizations to organize/ right of trade union organizations to organize/ right of these employees to organize/ rights of workers' organizations to organize/ right of workers' trade unions to organize/ political activities/ ...

20/8/2008gr 24/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion 0000

Quantifying descriptor concepts

The $tf \times idf$ measure

$tf_{ij} \times idf_j$ weight of concept j in text i

- tf_{ij} frequency of concept j in text i
- idf_j inverse document frequency of concept j
(= $\ln(d/d_j)$ with d_j #docs with concept j , d total #docs)

20/8/2008gr 26/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion 0000

Excerpt of data representing comments

CEACR Comment	Descriptor Concepts								
	c1	c2	c3	c4	c5	c6	c7	c8	...
Algeria 1991	0	0	0	0	2.75	0	0.8	0	...
Antigua and Barbuda 1991	3.0	1.53	0	0	0	0	0	0	...
Argentina 1991	0	0	0	0	20.59	2.39	0.8	0	...
Bangladesh 1991	1.0	0.77	2.35	1.24	0	1.59	5.59	0	...
Belgium 1991	0	0	0	0	0	0	0.8	0	...
Bolivia 1991	0	0.77	0	1.24	1.37	4.77	0	0	...
Bulgaria 1991	0	0	0	0	2.75	0	0	0	...
Burkina Faso 1991	0	0	0	0	0	0	0.8	2.19	...
...									

20/8/2008gr 27/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion 0000

Limits and possibilities of improvement

- Design of descriptor concept is crucial stage of the process.
- Semi-automatized process that strongly relies on the domain expert.
- Time-consuming
- Clever tuning through individual intervention of
 - domain expert
 - text mining expert
- Resulting descriptor concept somewhat subjective (not strictly reproducible)
- **Need for improvement and Systematization**

20/8/2008gr 28/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

Improvement and systematization

- Through access to **lexicon** and **ontology** of the concerned legal domain (which do not yet exist)
- Lexicon : list of words (terms) with synonyms.
- Ontology
 - puts together the characteristic terminology
 - organizes it in terms of concepts and sub-concepts
 - describes interrelation between concepts.

20/8/2008gr 29/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

Learning Categorization Rules

- Learning sample
 - 78 texts out of 671 for C87
 - 60 texts out of 509 for C98
- Assign label (violation) to texts in sample
 - Learn a classifier, for each violation k
 - $v_k = f_k(c_1, c_2, \dots)$, $k = 1, 2, \dots, 9$
- We used **classification trees**

20/8/2008gr 32/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

Grown tree for restriction on organization of trade union activities v_7

20/8/2008gr 33/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

Error rates, Convention 87

Key concept (violation)	Learning error rate	Cross-validation error rate	std err	Test sample (size 21) number of errors
v_3	14.10%	n.a.*	n.a.*	3
v_4	5.13%	5.13%	2.50%	0
v_6	12.82%	14.1%	3.94%	4
v_7	15.38%	n.a.*	n.a.*	7
v_8	7.69%	7.69%	3.01%	4
v_9	2.56%	2.56%	1.79%	2

*Cross-validation is not available for v_3 and v_7 , because first split is enforced.

Example of prediction error

20/8/2008gr 35/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

False Positives, False Negatives, Recall and Precision Convention 87

Key Concept	Positives		Negatives		% with key concept		Recall	Precision
	true	predic.	true	predic.	reported	predic.		
v_3	30	32	37	46	50.0%	41.0%	76.9%	93.8%
v_4	29	31	45	47	39.7%	39.7%	93.5%	93.5%
v_6	35	38	33	40	53.8%	48.7%	83.3%	92.1%
v_7	50	59	16	19	67.9%	75.6%	94.3%	84.7%
v_8	29	30	43	48	43.6%	38.5%	85.3%	96.7%
v_9	57	59	19	19	73.1%	75.6%	100.0%	96.6%

20/8/2008gr 36/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

Degrees of freedom in assignment rules

- Two class problem (violation yes/no)
- Assignment rule : Select "yes" if $p(\text{yes}) > s$
- Assignment controlled through threshold s :
 - large : favors precision (few false positives)
 - small : favors recall (few false negatives)
- Choice depends on what we want to do with predictions
 - Help legal expert identifying issues with given countries : **recall** more important (easier to find out false positives)
 - Synthetic analysis : **precision** more important.

20/8/2008gr 37/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

Conclusion 1

- Ad hoc text mining process for identifying issues reported by ILO legal texts.
- **Rule construction is long** and time-consuming.
 - Semi-automatic building of the prediction system based on strong interaction with domain expert.
- **Resulting rules are fully automatic.**
 - Can fast and straightforwardly be applied to any new text.

20/8/2008gr 39/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

Conclusion 2

- Improvement and systematization
- **Characterization of relevant descriptor concepts.**
 - Global criterion measuring simultaneously discrimination of all targets (violations).
 - Lexicon and ontology of the domain.
- Improvement at **learning level.**
 - Account for preference for precision or recall during tree growing (not only for class assignment).

20/8/2008gr 40/45

Introduction 0000 The Text Mining Process 000000 Text representation 00000000 Learning 000000 Conclusion

THANK YOU!

20/8/2008gr 41/45

References Example of report : Gabon 1992

References

Heitz, T., M. Roche, and Y. Kodratoff (2005). Extraction de termes centrée autour de l'expert. *Revue des nouvelles technologies de l'information RNTI E-5*, 685–690.

Ritschard, G., D. A. Zighed, L. Baccaro, I. Georgiou, V. Pisetta, and M. Studer (2007). Mining expert comments on the application of ILO Conventions on freedom of association and collective bargaining. Working Papers 2007.02, Department of Econometrics of the University of Geneva.

Salton, G. and C. Buckley (1988). Term-weighting approaches in automatic text retrieval. *Information Processing and Management* 24(5), 513–523.

Schmid, H. (1994). Probabilistic part-of-speech tagging using decision trees. In *International Conference on New Methods in Language Processing, Manchester*.

20/8/2008gr 42/45

References Example of report : Gabon 1992

Convention 87 : Gabon 1992

(back)

Actual and predicted violations in CEACR 1992 report for Gabon

Violation	3	4	6	7	8	9
Actual	0	1	0	0	0	1
Prediction	9%	94%	92%	53%	10%	97%

- Report contains two references to the "Election" descriptor concept ⇒ 94% probability to mention problems regarding election of representatives (v_6).
- However, first mention of "election" is within a sentence about new rules adopted for **insuring representative election**, second one in a sentence for asking the Government to **indicate the election results**.
- Positive reference to "election" confused with a violation.

20/8/2008gr 44/45

References Example of report : Gabon 1992

CEACR 1992 report for Gabon

(back)

The Committee notes, in particular, the Government representative's statement that the recognition of individual liberties in the new Constitution of Gabon, which came into force on 26 March 1991, has a corollary in the overall social plan, which is the abolition of trade union **monopoly**, that is to say the establishment of genuine and complete freedom of association. It notes that a draft new Labour Code which was discussed during a tripartite meeting from January to April 1991, attended both by the unitary employers' and workers' central organisations and by other organisations of workers and employees, has already been examined by the Government and was to be presented before the end of 1991. According to the Government, the amendment envisaged includes the repeal of section 174 of the present Labour Code which obliges all workers' or employers' organisations to affiliate with the Trade Union Confederation of Gabon (COSYGA) or the Employers' Confederation of Gabon (CPG). The Government also states that Act No. 13/80 of 2 June 1980, establishing a trade union solidarity tax deducted for the COSYGA, is no longer applied and that the tax has not been deducted since March 1990. Legislation is to be adopted for its formal repeal.

With regard to the provisions on compulsory arbitration restricting workers' right to **strike** (sections 239, 240, 245 and 249 of the Labour Code), the Government representative stated that a draft law specifically on the right to **strike**, which takes into account the requirements of the Convention, has been prepared and may be incorporated into the revised Labour Code.

The Committee notes the Government's reply in its last report to the effect that : (1) COSYGA, whose members wish the organisation to continue under the same name, has complied with the laws of the Republic of Gabon and adopted new rules under which it is now protected from any influence on the part of political parties and religions; (2) the new rules of COSYGA settle clearly the problem of the social assets of COSYGA vis-à-vis the new unions; (3) the sole object of occupational organisations is to examine and defend members' economic, industrial, commercial, agricultural and artisanal interests and there are no longer any restrictions on the establishment of these organisations; and (4) future **elections** of staff delegates and members of the Economic and Social Cooperation Committees will demonstrate that the various unions in establishments and enterprises are representative.

In the light of this information, the Committee asks the Government to provide a copy of the new COSYGA rules with its next report and to indicate the results of the above-mentioned **elections**.

20/8/2008gr 45/45