



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

Team Orpailleur

*Knowledge Discovery guided by Domain
Knowledge*

Lorraine

THEME SYM

Activity
R *eport*

2006

Table of contents

1. Team	1
2. Overall Objectives	1
2.1. Overall Objectives	1
3. Scientific Foundations	2
3.1. Knowledge Discovery in Databases	2
3.1.1. Symbolic Methods in Knowledge Discovery	2
3.1.2. KDD in biology, chemistry and medicine	3
3.1.3. Data Mining with Hidden Markov Models	7
3.1.4. The text mining process	9
3.2. Knowledge Representation, Knowledge Systems and Semantic Web	10
3.2.1. Classification-based Systems and Reasoning	10
3.2.2. The Semantic Web framework	11
3.2.3. Knowledge Management in Medicine: the Kasimir System	11
3.2.4. Spatial Knowledge Representation and Spatial Reasoning	13
3.2.5. Intelligent Access to Information	14
4. Software	14
4.1. A Data Mining Toolkit: the Coron Platform	14
4.2. Stochastic systems for knowledge discovery and simulation	15
4.2.1. CarottAge	15
4.2.2. GenExp	16
4.3. Softwares for the manipulation of documents for the Semantic Web	16
4.3.1. tamis: A software for text and rule mining	16
4.3.2. IntoWeb: Intelligent Access to Information	16
4.3.3. DefineCrawler: a Generic Crawler	17
4.4. Software for Spatial Reasoning	17
4.5. The Kasimir System	17
5. Other Grants and Activities	18
5.1. The European Network of Excellence Knowledge Web	18
5.2. The Eureka GenNet Project	19
5.3. National initiatives	19
5.3.1. ACI IMPBIO: the FouDAnGA Project	19
5.3.2. ACI IMPBIO: the ISIBIO Project	19
5.3.3. ACI “Masse de données en Astronomie”	19
5.3.4. cnrs tcan Project	20
5.3.5. Projects and Collaborations in Spatio-Temporal Reasoning	20
5.4. The “PRST IL” Programme	21
6. Dissemination	21
6.1. Scientific Animation	21
6.2. Teaching	21
6.3. Transfer activities	21
6.4. Awards	22
7. Bibliography	22

1. Team

Note on the organization of the report. Regarding the organization of this report, for convenience, applications and scientific results are not presented in specific sections, but, instead, follow the theoretical topics on which they are based.

Team Leader

Amedeo Napoli [Researcher (DR CNRS), HdR]

Administrative assistant

Antoinette Courrier [Secretary (Technicienne CNRS)]

Staff members

Marie-Dominique Devignes [Researcher (CR CNRS), HdR]

Florence Le Ber [Professor (ENGEES Strasbourg), HdR]

Jean Lieber [Associate Professor (MdC Université Henri Poincaré Nancy 1)]

Jean-François Mari [Professor (Université de Nancy 2), HdR]

Emmanuel Nauer [Associate Professor (MdC Université de Metz)]

Malika Smaïl [Associate Professor (MdC Université Henri Poincaré Nancy 1)]

Yannick Toussaint [Researcher (CR INRIA)]

Ph.D. Students

Mathieu d'Aquin [ATER at UHP Nancy 1 (until June 2006)]

Fadi Badra [PhD Student (MERT Grant)]

Rokia Bendaoud [PhD Student (Région and INRIA Grant)]

Adrien Coulet [PhD Student (CIFRE contract with Kika médical Nancy)]

Sébastien Hergalant [PhD Student and ATER at University of Nancy 2 (until September 2006)]

Nicolas Jay [PhD Student and lecturer (Faculté de Médecine, UHP Nancy 1)]

Mohamed Zied Maala [PhD Student (France Télécom Grant)]

Nizar Messai [PhD Student (Région and UHP Nancy 1 Grant)]

Frédéric Pennerath [PhD Student and lecturer (Supélec Metz)]

Laszlo Szathmary [PhD Student and ATER at UHP Nancy 1)]

Sylvain Tenier [PhD Student (CIFRE contract with INIST Diffusion Nancy)]

Post-doctoral fellows and visiting students

Charu Asthana [Engineer (from April 1st, 2006)]

Hacène Cherfi [Post-Doctoral fellow (Spring 2006)]

Bertrand Delecroix [Post-Doctoral fellow (from November 1st, 2006)]

Sandy Maumus [Post-Doctoral fellow (CPER Grant)]

Jesus Herrezuelo Santamaria [visiting student (Valladolid University ETSII, from October, 15th)]

Visiting scientists

Sergei Kuznetsov [Professor (High School of Economics, Moscow, Russia, February and August 2006)]

Petko Valtchev [Associate Professor (UQAM Montréal, Québec, July 2006)]

2. Overall Objectives

2.1. Overall Objectives

The “orpailleur” denotes in French a person who is searching for gold in the rivers. In the present case, gold nuggets correspond to knowledge units and may have two major different origins: explicit knowledge that can be given by domain experts, and implicit knowledge that must be extracted from data sources of different natures, e.g. rough data or textual documents. The main objective of the members of the Orpailleur team is to extract knowledge units from different data sources and to design structures for representing the extracted knowledge units. Knowledge-based systems may then be designed, to be used for problem-solving in a number of application domains such as agronomy, biology, chemistry, medicine, the Web...

The research work of the Orpailleur team may be considered mainly as KDDK, i.e. *Knowledge Discovery guided by Domain Knowledge* involving knowledge extraction and knowledge representation. First, the data sources are prepared to be processed, then they are mined, and finally, the extracted information units are interpreted for becoming knowledge units. These units are in turn embedded within a representation formalism to be used within a knowledge-based system. The mining processes are based on the *classification* operation, e.g. hidden Markov models, lattice-based classification, frequent itemset search, and association rule extraction. The mining process may be guided by a domain *ontology*, that is considered as a domain *model*, used for interpretation and reasoning, especially in the context of semantic Web.

The whole transformation process from rough data into knowledge units is based on the underlying idea of *classification*. Classification is a polymorphic process involved in a number of tasks within the data to knowledge transformation, e.g. mining operations, modeling of the domain for designing a domain ontology (or extending the ontology with extracted knowledge units), knowledge representation and reasoning. Finally, the knowledge extraction process and the associated knowledge base can be used for problem-solving activities within the framework of the Semantic Web, e.g. Web mining, intelligent information retrieval, content-based document mining...

3. Scientific Foundations

3.1. Knowledge Discovery in Databases

Keywords: *association rule extraction, bioinformatics, data mining methods, frequent itemset search, hidden Markov models for data mining, knowledge discovery in databases, lattice-based classification, text mining.*

Participants: Fadi Badra, Rokia Bendaoud, Adrien Coulet, Marie-Dominique Devignes, Sébastien Hergalant, Nicolas Jay, Florence Le Ber, Jean Lieber, Jean-François Mari, Sandy Maumus, Nizar Messai, Amedeo Napoli, Frédéric Pennerath, Malika Smail, Laszlo Szathmary, Sylvain Ténier, Yannick Toussaint.

Knowledge discovery is a process for extracting knowledge units from large databases, units that can be interpreted and reused within knowledge-based systems.

3.1.1. Symbolic Methods in Knowledge Discovery

Knowledge discovery in databases (KDD) consists in processing a huge volume of data in order to extract useful and reusable knowledge units from these data. An expert of the data domain, called hereafter the *analyst*, is in charge of guiding the extraction process, on the base of his objectives and of his domain knowledge. The extraction process is based on data mining methods returning information units from the data. The analyst selects and interprets a subset of the units for building “models” that may be further interpreted as knowledge units with a certain plausibility.

The KDD process is performed with a KDD system based on four main components: the databases (or the set of data), a domain ontology (and an associated knowledge-based system), data mining modules (either symbolic or numerical), and interfaces for interactions with the system, e.g. editing and visualization. For handling huge volume of data in a given domain, a KDD system may take advantage of domain knowledge, i.e. an ontology, and the problem-solving capabilities of a knowledge-based system working in the domain of data. In turn, closing the loop, the knowledge units extracted by the KDD system may be integrated within the ontology to be reused by the knowledge-based system for future problem-solving operations.

Lattice-based classification, frequent itemset search, and association rule extraction.

Symbolic methods for KDD mainly rely on lattice-based classification, frequent itemsets, and association rule extraction [39]. Lattice-based classification is used for extracting from a database (or a set of rough data) a set of concepts organized within a hierarchy i.e. a partial ordering. Lattice-based classification relies on the analysis of binary tables relating a set of individuals with a set of properties (or characteristics), where *true* stands for the individual *i* has the property *p*. The lattice may be built according to the so-called *Galois* correspondence, classifying within a formal concept a set of individuals, i.e. the extension of the concept,

sharing a common set of properties, i.e. the intension of the concept. In addition, lattice-based classification is the basic operation underlying the so-called *formal concept analysis*.

In parallel, the extraction of frequent itemsets consists in extracting from binary tables sets of properties occurring with a support or frequency, i.e. the number of individuals sharing the properties, greater than a given threshold. From the frequent itemsets, it is possible to generate association rules of the form $A \longrightarrow B$ relating the subset of properties A with the subset of properties B , that can be interpreted as follows: the individuals including A also include B with a certain support and a certain confidence. The number of rules that can be extracted is very large, and there is a need for pruning the sets of extracted rules for interpretation (most of the time, the analyst is in charge of interpreting the results of the rule extraction process). Different measures have been set on, mainly based on probability theory, that can be used for pruning the sets of extracted rules (i.e. rule mining). The Orpailleur team is also interested in the mining of rare itemsets and rare rules, an research work that is an originality of the team [25], [30], [38]. Accordingly, the team is currently developing a platform for knowledge extraction with symbolic methods, called CORON, that includes a collection of data filtering methods and symbolic data mining algorithms. The CORON platform is used in a number of KDD applications that are described in the following.

3.1.2. KDD in biology, chemistry and medicine

Data integration and knowledge extraction in bioinformatics.

Biological datasets have tremendously grown in size and complexity in the past few years. Genome sequences, biological structures, expression arrays, proteomics represent terabytes of data which are stored under variable formats in dispersed heterogeneous databases (DB). More than 800 such DB have been listed at the beginning of 2006. One of the major challenges in the post genomic era consists in exploiting the vast amounts of biological data stored in those DB. The extraction of knowledge from all these data is an increasingly challenging task which ultimately gives sense to the data production effort with respect to domains such as evolution and disease understanding, biotechnologies, systems biology, pharmacogenomics, etc. The knowledge discovery in biological databases process starts with two important steps: data selection from appropriate DB, and data integration. In the biological domain, these tasks are hampered by at least two distinct problems: (i) identifying the relevant DB, and (ii) managing the complexity and heterogeneity of biological data for their integration. Previous and present work within the Orpailleur group has been dealing with the first two aspects of the KDD process: selection of biological databases and heterogeneous biological data integration.

Heterogeneous biological data integration: customized warehouses populated by a user-defined workflow of data retrieval from public databases (ACGR project)

Integration of heterogeneous biological data is addressed from a pragmatic and user-oriented point of view. In many concrete situations, data have to be collected from various public sources in order to answer complex queries. In previous work we have developed a generic solution for automated collection of biological data given a user-defined workflow. The Xcollect software described in the preceding reports provides users with structured documents containing retrieved data to answer his query. The need to combine various workflows and to store the retrieved data for further exploitation has led us to propose a concept of customized warehouse populated by user-defined workflows of data retrieval.

Given a specific need workflows of data retrieval from public databases (mostly web sources) are designed experimentally. A data model is then built to integrate the retrieved data in a database. Populating the database involves converting the format of retrieved data into an entry format for the database.

In practice retrieval workflows are Xcollect XML scenarios. The retrieved data are compliant with the Xcollect XML session DTD (Definition of Type of Document). The data model is implemented as a MySQL relational database. Mapping between the XML elements of a session document corresponding to a given scenario and the tables and attributes of the relational database model is expressed as a simple correspondence table. A generic python script takes this table as input and produces the XSL transformation file able to transform the XML session document into an SQL command. Execution of the resulting SQL command allows insertion of the retrieved data in the database.

This approach has been applied to the retrieval of candidate genes for a rare orphan disease : Aicardi Syndrome. Collected data concern human, mouse and fly genes related to the phenotype observed in diseased patients (for instance retina abnormalities or defects in neuron migration). Data such as chromosome localization, Gene Ontology (GO) annotations, homologous inter-species relationships, and interactions with other genes are retrieved thanks to various Xcollect scenarios and stored in the ACGR (Aicardi Candidate Gene Retrieval) database. For example users can retrieve from NCBI GENE database all human genes annotated by a given GO term (for example "neuron migration") and store in the ACGR database various pieces of information associated with these genes. Another scenario will retrieve from the same or another database all the genes interacting in some way with these genes and associated information such as chromosomal localization. Ultimately complex queries such as: "What are the genes located on chromosome X that interact with a gene annotated by the GO term 'neuron migration'" can be answered and give new insights on potential candidate genes.

The resulting ACGR database can be considered as a customized data warehouse integrating data from various public sources. Tracking of data origin is included in each Xcollect scenario and has been taken into account in the data model. Refreshing is performed by re-executing the Xcollect scenarios and producing a new release of the database. Data analysis and data mining methods can be plugged in the system according to user needs. Present work concerns prioritization of genes with respect to their probability of being the gene responsible for the disease since the best predictions must be tested now experimentally. Of course application to another disease is possible.

Organizing and Querying a metadata repository with Concept Lattices (The BioRegistry project).

The BioRegistry project aims at gathering and organizing knowledge about biological databases in order to facilitate and to optimize the selection of relevant databases with respect to a user query. In the "BioRegistry" repository the various metadata attached to biological databases are structured according to a model described in the previous report, and whenever possible expressed in terms of domain ontologies. Our model is compliant with the DCMI ("Dublin Core Metadata Initiative") recommendations and uses two main domain ontologies to valuate metadata fields describing the content of the databases. The *subjects* field for instance contains terms extracted from the biomedical thesaurus MeSH, maintained by NLM. This thesaurus was chosen because it is widely used to index scientific literature, it presents a broad coverage of many biological domains and is regularly updated to take into account changes in the topics addressed by scientific papers. Concerning the *organisms* field, the NCBI taxonomy of living organisms has been chosen since this taxonomy is also used to annotate biological sequences. In previous stage of the work, inclusion of several databases in the BioRegistry repository has been performed manually. To accelerate the process, an automatic procedure was designed to import the DBCAT metadata (see previous report). Since the DBCAT catalog is no more maintained, it was then decided to exploit the Nucleic Acids Research (NAR) 2006 catalog of molecular biology databases maintained at NCBI. Several text-mining programs were set up to translate the NAR information into controlled vocabulary terms. Terms found in the database short descriptions were cross-matched with the list of organism names available from the NCBI taxonomy. Retrieved terms were entered in the *organisms* subsection of the BioRegistry repository. In addition, we built a correspondence table between the NAR categories and sub-categories and MeSH terms to be included in the *subjects* field of the BioRegistry. Alert and survey mechanisms have to be designed to detect any change or new release in existing databases as well as new databases appearing on the Web.

Formal Concept Analysis (FCA) was set up for organizing the BioRegistry and visualizing the sharing of metadata across the DB. A formal context representing the relation between bioinformatics data sources and their metadata is provided, and the corresponding concept lattice is built (see previous report). The BR-explorer algorithm addresses the problem of retrieving the relevant data sources for a given query [6]. It starts by building the query concept representing the query, and then inserts the query concept in the concept lattice. Then, the BR-explorer algorithm fills a list of candidate concepts, according to a relevance criteria which was first set as follows : a relevant concept is a concept that shares at least one property with the query concept. The BR-explorer algorithm thus explores the ascendants of the query concept in the concept lattice, until the top concept is reached. Finally, the BR-explorer algorithm returns the set of relevant data sources ranked according

to their distance to the considered query. The distance measure was first set as follows : number of edges in the lattice between the query concept and the relevant concept. Various refinements of both the relevance criteria and the distance measure are currently under investigation for taking advantage of the semantics of queries and metadata. This work should be reusable for resource discovery and composition in the frame of the semantic web.

Knowledge extraction in pharmacogenomics.

Another ongoing research work consists in applying the whole KDD process to the pharmacogenomics context, i.e. from data selection and filtering to knowledge extraction guided by the domain knowledge. More precisely, the goal is to discover knowledge about interactions between clinical, genetic and therapeutic data. For example, a given genotype –set of selected gene versions–may explain adverse clinical reactions, e.g. hyperthermy, toxic reaction...to a given therapeutic treatment. Indeed more and more pharmaceutical firms are willing to include the exploration of particular genomic variants in their drug clinical trials in order to detect relevant relationships between the three vertices of the pharmacogenomics triangle, i.e., (i) drug (properties and administration), (ii) phenotype (biological and clinical data) and (iii) genotype (genome variations).

We first focused on the genotype vertex by building SNP-Ontology that formalizes available knowledge about genomic variations. This allowed us reconciling the various heterogeneous representations of both private data and data coming from public databases (dbSNP, UCSC, HapMap...). A UML class diagram was first designed as an intuitive description of the relevant knowledge and then transformed into an OWL formal model. A dedicated wrapper was developed for gathering data from several sources. These data (A-box) instantiated successfully the concepts of SNP-Ontology (T-box) and consistency checking was successfully conducted as well. This constituted a first validation of SNP-Ontology [16].

SNP-Ontology goes far beyond a simple controlled vocabulary or taxonomy which is more or less the state of most bio-ontologies today. We demonstrate that semantic relationships other than the classical "is-a" one have to be used for representing knowledge about genomic variations and for enabling reasoning over it. SNP-Ontology differs from PharmGKB ontology which is expressed as an XML schema. PharmGKB ontology has a wide scope and covers all three vertices of the pharmacogenomics triangle. However concerning representation of genomic variations, PharmGKB schema is less open than SNP-Ontology which is much more focused and complete on that aspect of pharmacogenomics. Indeed the SNP-Ontology makes it possible for all representations of a given variant to co-exist in the SNP-knowledge base and be declared equivalent [16].

In a second phase, we have been working on the construction of a modular and formal representation of domain knowledge in pharmacogenomics. The resulting ontology is called SO-Pharm for Suggested Ontology for Pharmacogenomics. We adapted some well-known methodologies for ontology construction to the case of pharmacogenomics, based on three steps: (i) specification, including definition of ontology domain and scope; (ii) conceptualization, involving definition of lists of terms and of concepts; (iii) formalization and implementation i.e. the translation of the conceptual model in a knowledge representation formalism (OWL in our case).

Domain and scope of SO-Pharm were primarily defined as follows. The domain considered should cover pharmacogenomics clinical trials. The ontology has to precisely represent groups of individuals involved in trials, their genotype, their treatment, their observed phenotype and the potential pharmacogenomics relations discovered between these concepts. SO-Pharm scope is to guide KDD in pharmacogenomics. Term lists were established either by domain experts or by extraction from existing data or knowledge resources in the domain. Relevant reusable resources were select at this stage. OBO (Open Biomedical Ontologies) ontologies (<http://obo.sourceforge.net>) were preferred and among them those involved in the OBO-Foundry project. As for the SNP-Ontology, a UML class diagram was used here for representing the conceptual model of SO-Pharm. Embedding and extension strategies were used to anchor existing ontologies to SO-Pharm concepts. Several highly specialized vocabularies such as Disease Ontology were embedded whereas formal ontologies, such as SNP-Ontology, extend definitions of more specific concepts pertaining from other ontologies. The consistency and the class hierarchy of SO-Pharm including reused ontologies have been validated with Racer 1.9 at each stage of the implementation thanks to standard reasoning mechanisms. As a preliminary validation of the ontology, several examples of published pharmacogenomics knowledge have been expressed with the

SO-Pharm concepts. The assertions of individuals and related information (clinical trial, treatment) lead us to enrich SO-Pharm concepts. The ontology construction method has been published [17].

In summary, SO-Pharm construction favors the reuse of concept definitions existing in other ontologies. This reuse mechanism will become more and more important since more and more autonomous ontologies are being produced in the biomedical domain, e.g. for representing phenotypes. SO-Pharm is available (in OWL format) at <http://www.loria.fr/~coulet/ontology/sopharm/version1.2/sopharm.owl>. We plan to submit SO-Pharm to OBO portal to gain in visibility and facilitate further improvements.

SO-Pharm is a crucial component for a future knowledge-based application dedicated to pharmacogenomics knowledge discovery. A complete validation has now to be conducted, aimed at evaluating how SO-Pharm is able to guide the KDD process. A significant issue will be to develop appropriate wrappers to achieve heterogeneous data integration as for SNP-Ontology. Then mining methods will have to be articulated with the ontology in order to extract new relevant knowledge units that will enrich the ontology.

Association rule extraction in a biological database.

Relying on the KDD principles, a research work is currently under investigation in the domain of biology for searching associations between biological parameters involving cardiovascular (CV) risk factors in a given population of individuals. The studies carried out here rely on a real-world individual database, the STANISLAS cohort. It is a ten-years study holding supposed healthy French families. Families are examined every five years. At the beginning of the study, in 1993, 1006 families (composed by two parents and at least two children) were recruited for medical examination at the “Centre de Médecine Préventive de Vandœuvre-lès-Nancy (France)”. Families have been examined further around 1998–1999, and 2003–2004.

The cohort is explored for searching for genotypes and intermediate phenotypes of cardiovascular diseases (CVD), which are multifactorial pathologies resulting from gene-gene and gene-environment interactions. There is a need for extracting implicit and new potential risk factors for CVD within an always increasing volume of data (mainly due to the development of technologies such as PCR multiplex or microarrays). In the STANISLAS cohort, information holds on environmental, clinical, biological and genetic data. The KDD experiments have given results in accordance with the domain knowledge, and as well, other results allowing new research insights for further investigations. Regarding statistical methods usually used in this context, the general idea of the present research work is to mine the cohort for extracting itemsets that may be in turn considered as hypotheses to be validated by statistical tests.

Experiments for extracting potential valuable information on the metabolic syndrome in the STANISLAS cohort have been carried on extracting frequent itemsets and association rules. A methodology for mining cohorts has been proposed [37], that can be applied with CORON, useful for various studies (not restricted to biology). The methodology is based on frequent itemsets and association rules extraction, and brings interesting results from a biological viewpoint, as it enables the expert of the domain to generate new research hypotheses validated by statistical tests or new lab experiments. However, based on the fact that the metabolic syndrome is a condition relatively rare in the STANISLAS cohort, which is composed of supposed healthy individuals, the mining work within the cohort has been oriented on the extraction of rare rather than frequent itemsets. In this way, different algorithms have been proposed for extracting rare itemsets and rules [25], [30], [38]. Sandy Maumus is in charge of the mining work on the STANISLAS cohort. She has defended her PhD thesis on November 15, 2005, and she won the PhD Thesis Award of the University Henri Poincaré (Faculty of Pharmacy) [1]. The work on rare itemset and rule extraction has to be followed, and different experiments are currently under investigation for evaluating the kind of information included in rare itemsets and rules, in comparison with frequent itemsets in the STANISLAS cohort.

Extracting knowledge in medico-economical databases.

Chronic diseases imply recurrent hospitalizations. In order to optimize healthcare resources, improve cooperation between hospitals treating chronic patients, it is very important to understand the factors that may determine the so-called *pathway* of a chronic patient. The patient pathway may be seen as a time-ordered sequence of events affecting the health of the patient. An event describes a set of information related to an hospitalization, such as, diagnoses, medical or surgical procedures, hospitalization locations, durations, and

costs... In France, the so-called PMSI (for “Programme de Médicalisation des Systèmes d’Informations”) is the name of the information system collecting for an hospital the information mentioned above.

At present, we are carrying out a research work on the data collected within the PMSI with the following objectives:

- The discovery of elements that may characterize the patient pathway.
- The classification of patients with respect to their pathway.
- The visualization of the patient pathway.

The first objective relies on the extraction of frequent patterns, sequential and not sequential, from the data of PMSI associated to the Lorraine Region. The database includes information on more than 800 000 hospitalizations per year. The two following objectives allow, based on the patterns that have been extracted, to build and to visualize a patient pathway classification, using concept lattices (or Galois lattices). More generally, this research work aims at investigating the relations that may exist between frequent itemsets, sequential itemsets, and knowledge representation and visualization with concept lattices. Various successful experiments have been carried out with data on cancer patients in the Lorraine Region [23], [22], [21].

Knowledge discovery in chemical reaction databases.

Chemical reactions are the main elements on which relies synthesis in organic chemistry, and this is why chemical reactions databases are of first importance. >From a problem-solving process point of view, synthesis in organic chemistry must be considered at several levels of abstraction: mainly a strategic level where general synthesis methods are involved, and a tactic level where actual chemical reactions are applied. The research work carried out in the present case is aimed at discovering general synthesis methods from chemical reaction databases in order to design generic and reusable synthesis plans.

A first research work based on levelwise frequent itemset search and association rule extraction, and on chemical knowledge, has been carried on, and has given substantial and promising results. At present, this first research work is extended, trying to adapt a graph-mining process for extraction knowledge from chemical reaction databases, but this time directly from the molecular structures and reactions themselves (both being represented as graphs in reaction databases). This research work is currently under investigation and should bring substantial results [29].

3.1.3. Data Mining with Hidden Markov Models

For designing a complete knowledge discovery system, we have developed stochastic models based on high-order hidden Markov models. These models are capable to map sequences of data into a Markov chain in which the transitions between the states depend on the n previous states according to the order of the model. The following experiments are based on second-order hidden Markov models (HMM2), i.e. the transitions between the states depend on the *two* preceding states, for discovering spatial and temporal dependencies in databases. The main advantage of HMM2 is the existence of a non-supervised training algorithm –the EM algorithm–, that allows the estimation of the parameters of the Markov model from a corpus of observations and an initial model. The resulting Markov model is able to segment each sequence of data into stationary and transient parts.

We focused our effort on two main points: (1) the elaboration of a process for mining spatial and temporal dependencies in order to extract knowledge units (for knowledge acquisition). This process involves an unsupervised classification of data. (2) The specification of adapted visualization tools giving a synthetic view of the classification results to the experts who have to interpret the classes and/or specify new experiments.

Several applications have been carried out during this last year, and two ANR projects in which the Orpailleur team is involved have been selected: the ADD-COPT project for “Agriculture et Développement Durable”, and the ECOGER project (for “Écologie pour la Gestion des Écosystèmes et de leurs Ressources”). In parallel, the research project called FONDANGA within the ACI IMPBIO (for “Informatique, Mathématiques, Physique en Biologie Moléculaire”) is running.

All these research works have taken advantage of the CAROTTAGE system, a generic data-mining system for spatio-temporal data, based on HMM2 (the CAROTTAGE system is a free software with a GPL license).

Two applications in agronomy.

The ANR project, called ADD-COPT for “Agriculture et Développement Durable”, aims at understanding the agriculture evolution for respecting the environment. An agriculture more respectful of the environment will modify the organization of the territory at several levels, i.e. spatial, economic and organizational levels. In this project, we work in collaboration with agronomists, but also with geographers, since there is a question on the representation of territories, with economists since bio-agriculture (organic agriculture) must remain economically viable, and with psychologists, who have to formalize how the different actors may share their knowledge for achieving this common objective of new agriculture. The goal of the ADD-COPT project is to specify an observatory of agricultural practices for supporting the different actors in the transformation process to this new agriculture: allowing these actors to confront and share their knowledge, to apprehend and analyze the observations made on the territory, and to assess the impacts of the changes in progress.

A second research project, called ECOGER, is still lying in the context of the mining of environmental data. It groups together various competences such as agronomy, zoology, and data mining. We are currently using the CAROTTAGE system to process at the same time temporal and space data, for allowing the agronomists to analyze data collected during several years on the ground occupation on a whole of points of the French territory. Preceding the ECOGER project, the CAROTTAGE system has been already used for understanding the risks that the bustard was facing while the disappearance of the meadows is clearly impacting its migration. In the ECOGER project, the CAROTTAGE system has to be used within a broader framework: environmental risks. The software has to be adapted to take into account the space organization of the successions of cultures. The challenge is double: whereas the software works with temporal data, it has to integrate spatial dimensions and, whereas it has been already tested on relatively homogeneous data, it has to be able to integrate data at different scales, e.g. satellite images, investigations with farmers...

An application in bioinformatics.

In the framework of the so-called *Contrat de plan État-Région*, we are carrying out a long term data mining project with the laboratory of genetics of the “Université Henri Poincaré Nancy 1”. The biological material is the soil-dwelling, filamentous bacteria belonging to the genus *Streptomyces*, that is the greatest source of antibiotics amongst microorganisms. In particular, the 8,7M bases of the *Streptomyces coelicolor* chromosome have been entirely sequenced and annotated. One objective in this research work is to detect genome heterogeneity islands, and inter sequences dependencies, using hidden Markov models without prior knowledge.

- The horizontal transfer understanding. Markovian models with “specie specific homogeneity” have been constructed, and coupled with transform filters. Their behavior generates regions with different statistical properties, allowing the user to separate “foreign DNA regions” with the own DNA regions of the studied specie itself. In *Streptomyces coelicolor*, the regions with such a statistical consensus have been detected, and correlated with potential events of horizontal transfer.
- The detection of promoters. A data mining method based on second order hidden Markov models (HMM2) that is able to process the whole genome sequences without prior hypotheses has been applied to the actinomycete genomes of *Streptomyces* and *Mycobacterium* species. The stochastic modeling of the genome with HMM2 allows the extraction, and the classification of short segments (5 to 12 bp) having a significant different structure and composition without any prior knowledge. These segments appear to be parts of binding sites for transcriptional factors.
- In order to confirm the applicability of this new method to the detection of transcriptional signals, the models have been applied to an experimentally determined co-regulated gene set (30 genes) dependent on SigR, a sigma factor of *Streptomyces coelicolor* involved in the oxidative stress response. A steady homogeneous second-order hidden state chain describes discrete heterogeneity visualized as peaks in the a posteriori observation of the hidden states.

The duration capabilities of the HMM2 shows very good performances in the modeling of short segments such as TFBS (Transcriptional factors binding sites), and RBS (Ribosome binding sites). On different genomes of the actinomycete family, the HMM2 reveal DNA heterogeneity, that are combined to predict known or potential TFBS or RBS. Based on these models, the present data mining method has proved to be efficient for the detection of DNA motifs involved in both transcriptional, i.e. sigma factor binding sites, and translational (RBS) regulation.

3.1.4. *The text mining process*

The objective of a text mining process is to extract new and useful knowledge units in a large set of texts. The text mining process relies on the principles of KDD, although it shows some specific characteristics due to the fact that texts are written in natural language. The information in a text is expressed in an informal way, following linguistic rules, making the mining process more complex. To avoid information dispersion, a text mining process has to take into account paraphrases, ambiguities, specialized vocabulary, and terminology. This is why the first steps of a text mining process are usually dedicated to linguistic knowledge acquisition: lexicon, terminology, markers of semantic relations, discourse markers, specific syntactic or semantic structures...

To carry out studies on text mining, the Orpailleur team is interested in linguistic resources, working on real-world texts in application domains such as biology, astronomy..., using robust information extraction tools. Language thus is considered as a way for accessing information, and not as an object to be studied for its own.

This year was mainly dedicated to building ontologies from texts. However, we kept a “background activity” on mining the web and started some investigations on information extraction tools.

Building ontologies from texts.

In astronomy, people are now spending much more time in studying texts for acquiring and synthesizing already known information instead of making new “physical” observations. Moreover, part of the categorization process assigning types (galaxy, star...) to celestial objects is done manually. These considerations lead to the following questions: could we (partly) automate the collect of celestial object properties from texts and could we structure this information in an ontology providing a more exhaustive knowledge and a more exhaustive description of celestial objects?

This year, we developed a prototype using Formal Concept Analysis to build ontologies from texts. This prototype relies on the idea that verbs may be used to characterize objects. For example, the sentence “We observed stars” enable us to say that stars are “observable”. We build a binary table (Objects \times Verbs) and then build the Galois lattice. Objects are then structured into classes following the properties they are associated with in the texts. A transformation function can convert the lattice into a hierarchy where objects of the domain are leaves. his approach has been tested over 72 scientific abstract texts, classifying 79 objects and 14 properties and is currently evaluated by astronomers. The lattice is composed of 16 formal concepts [15].

The main limitation of this approach is that properties are Boolean attributes and n-ary relations cannot be modeled. Our current work is based on Relational Analysis where concepts are related to other concepts.

Knowledge extraction from Web pages.

This research is concerned with the design of a system for extracting knowledge from Web pages. Knowledge is encoded as a “semantic annotations” for manipulating the documents by their content [11].

Most of the current works consider the annotation process as an Information Extraction task. They rely on patterns which aims at identifying in the documents concepts of the ontology. We propose a new approach which relies on classification. Then the annotation process integrates both the syntactic structure of a web page as well as semantic constraints coming from the ontology [32]. The ontology is implemented within the Web Ontology Language (OWL) (reasoning mechanisms such as classification and subsumption are available).

More precisely, the semantic annotation of an element in a Web page relies on two main operations: (i) identification of the syntactic structure of a specific element in the Web page using the DOM Structure

(Representation of the page as trees and subtrees), (ii) identification in the ontology of the most specific concept subsuming the extracted element, that will be used for building the annotation.

The global context of the present research work is the study of research themes within the European Research Community. The objective is to use the information provided by research teams on their website to generate knowledge about the European Research Community, for technological watch, analysis of research themes, or detection of new research directions. Preliminary results have been described within an article presented at the Web intelligence Conference [31].

3.2. Knowledge Representation, Knowledge Systems and Semantic Web

Keywords: *case-based reasoning, classification-based reasoning, description logics, knowledge representation, knowledge-based information retrieval and extraction, lattice-based classification, object-based representation systems, qualitative spatial reasoning Semantic Web.*

Participants: Fadi Badra, Florence Le Ber, Jean Lieber, Amedeo Napoli, Emmanuel Nauer, Laszlo Szathmary, Yannick Toussaint.

Knowledge representation is a process for representing knowledge within a knowledge representation formalism, giving knowledge units a syntax and a semantics. The **Semantic Web** is a framework for building knowledge-based systems for manipulating documents on the Web by their contents, i.e. taking into account the semantics of the elements included in the documents.

3.2.1. Classification-based Systems and Reasoning

A knowledge system relies on a knowledge base and a reasoning module for problem solving and knowledge management in a given domain. Knowledge units are represented within a knowledge representation formalism where they have a syntax and a semantics. Inference can be drawn from already known knowledge units for deriving new units, that are useful for solving the current problem. Moreover, the units extracted from data by data mining procedures have to be represented within a knowledge representation formalism to be taken into account in the framework of a knowledge system.

In the team Orpailleur, two kinds of formalisms are particularly studied, namely description logics (DL) and object-based knowledge representation (OBKR) formalisms. Knowledge units are represented within concepts (also called classes), with attributes (properties of concepts, or relations, also called roles in DL), and individuals. The hierarchical organization of concepts relies on a subsumption relation that is a partial ordering. These formalisms provide inference services such as subsumption, concept and individual classifications. Concept classification is used to insert a concept at the right location in the concept hierarchy (searching for its most specific subsumers and its most general subsumees). Individual classification is used for recognizing the concepts an individual may be instance of. In both cases, subsumption and classification are the main operations: this is why these systems are denoted here by “classification-based systems”.

Classification-based reasoning may be extended into case-based reasoning (CBR), that relies on three main operations: retrieval, adaptation, and memorization. A source case $(srce, Sol(srce))$ lies in a case base, and can be seen as a problem statement $srce$ together with its solution $Sol(srce)$. Then, given a new target problem, say tgt , retrieval consists in the search for a memorized case whose problem statement $srce$ is similar to the target problem tgt . Then, when $srce$ exists, its solution $Sol(srce)$ is adapted to fulfill the constraints attached to tgt . When there is enough interest, the new pair $(tgt, Sol(tgt))$ can be memorized as a new case for further problem solving. In the context of a concept hierarchy, retrieval and adaptation may be both based on classification. Moreover, a number of studies within the Orpailleur team has been carried out on CBR, especially on “adaptation-guided retrieval”, that consists in searching for a source case whose solution will be adaptable for the target problem, giving a kind of guarantee regarding the building of the solution of the source case.

In parallel with knowledge representation (and knowledge extraction), knowledge management is oriented toward the management of what could be called the “cycle” of knowledge, including acquisition, memorization, retrieval, maintenance, dissemination (or exchange) of knowledge. There is also a need for coupling knowledge with data, with respect to representation and management. This means in particular that, besides knowledge extraction from databases, there are some other needs such as knowledge-based information retrieval, content-based manipulation of documents, and knowledge mining. These new directions of investigation are particularly important in the framework of the semantic Web.

3.2.2. *The Semantic Web framework*

Today people try to take advantage of the Web by searching for information (navigation, exploration), and by querying documents using search engines (information retrieval). Then people try to analyze the obtained results, a task that may be very difficult and tedious. Tomorrow, the Web will be “semantic” in the sense that people will search for information with the help of machines, that will be in charge of posing questions, searching for answers, classifying and interpreting the answers. The Web will become a space for exchange of information between machines, allowing an “intelligent access” and “management” of information. However, a machine may be able to read, understand, and manipulate information on the Web, if and only if the knowledge necessary for achieving those tasks is available. This is why ontologies are of main importance with respect to the task setting up a semantic Web. Thus, there is a need for representation languages for annotating documents, i.e. describing the content of documents, and giving a semantics to this content. Knowledge representation languages are (the?) good candidates for achieving the task: they have a syntax with an associated semantics, and they can be used for retrieving information, answering queries, and reasoning.

The semantic Web has gained a great interest in the research work of the Orpailleur team. Indeed, it constitutes a good platform for experimenting a number of ideas on knowledge representation, reasoning, knowledge management, and knowledge discovery (and especially text mining) as well. Investigations mainly hold on the content-based manipulation of textual documents using annotation, ontologies, and a knowledge representation language. The idea is to build an XML-based “bridge” between documents, and the knowledge units of the domain of documents, lying in domain ontology. The annotations attached to documents, and the queries, are built with the help of the concepts of the domain ontology. Then, the manipulation of annotations, e.g. information retrieval, query answering, reasoning on the content of documents, is left to the reasoning module associated with the knowledge representation formalism.

3.2.3. *Knowledge Management in Medicine: the Kasimir System*

The objective of the KASIMIR research project is decision support and knowledge management for the treatment of cancer. This is a multidisciplinary research project in which participate researchers in computer science (Orpailleur), in ergonomics (*Laboratoire d'ergonomie du CNAM*, Paris), experts in oncology (*Centre Alexis Vautrin* or CAV, Vandœuvre-lès-Nancy), Oncolor (a healthcare network in Lorraine involved in oncology), and Hermès (an association for the sharing of resources in informatics and medicine).

For a cancer localization, e.g. the breast, the treatment is based on a protocol similar to a medical guideline. This protocol is built according to evidence-based medicine principles. For most of the cases (about 70%), a straightforward application of the protocol is sufficient, and provides a solution, i.e. a treatment, that can be directly reused.

A case out of the 30% remaining cases is “out of the protocol”, meaning that either the protocol does not provide a treatment for this case, or the proposed solution raises difficulties, e.g. contraindication, treatment impossibility, etc. For such an out of the protocol case, oncologists try to *adapt* the protocol (actually they discuss such a case during the so-called “breast therapeutic decision meetings”, including experts of all domains in breast oncology, e.g. chemotherapy, radiotherapy and surgery). In addition, protocol adaptations are studied from the ergonomics and computer science viewpoints. These adaptations can be used to propose *evolutions* of the protocol based on a confrontation with actual cases. The idea is then to make suggestions for protocol evolutions based on frequently performed adaptations.

Adaptation knowledge acquisition.

The adaptation in KASIMIR, as well as in many CBR systems, requires knowledge. The adaptation knowledge acquisition (AKA) is a current research work, that takes two directions: AKA from experts and semi-automatic AKA.

AKA from experts consists in analyzing adaptations performed by experts. Interviews of experts confronted to decision problems requiring adaptation have been recorded to be afterward analyzed, and modeled within adaptation patterns [12].

Semi-automatic AKA is based on the “mining of the protocols”. A protocol can be seen as a set of rules “situation→decision”. Knowing how the decisions change when the situations change from one rule to another rule provides a specific adaptation rule. By generalizing these specific rules, general adaptation rules may be obtained. This generalization process has been implemented thanks to a frequent close itemset extraction module of the CORON platform (see § 4.1). This requires a formatting of the situations and decisions of the protocol following the itemset mode. A system, called CABAMA, realizes this case base mining for adaptation knowledge acquisition, and provides pieces of information that can be used for building adaptation rules [33]. This AKA is not fully automated: an analyst pilots CABAMA, following the principles of knowledge discovery. More precisely, the analyst uses filters to orientate the mining process, and interprets the extracted pieces of information in adaptation rules.

AKA from experts and semi-automatic AKA are not completely satisfying: the former provides generic adaptation patterns that are intelligible, but cannot be directly operational, while the latter provides adaptation rules that can be directly implemented, but are difficult to understand (and thus, to validate). A future research work will combine the two kinds of AKA for producing operational *and* intelligible adaptation knowledge units.

Knowledge representation for decision support tools.

Two versions of KASIMIR are currently used: one based on an *ad hoc* object-representation formalism (OBRF), the other one based on semantic Web principles, in a semantic portal (as explained below). A number of knowledge bases corresponding to specific cancers (decision protocols) has been developed. Moreover, the inference engine has been extended for taking into account a fuzzy representation of concepts and fuzzy hierarchical classification. The system tries to detect and to propose more than one treatment for “borderline cases”: this has been implemented for the OBRF version of KASIMIR, and its implementation in the semantic portal is under development [13], [34]. Another study is about “multiple viewpoint representation and reasoning”, that may be useful for modeling the reasoning of the breast therapeutic decision committee, i.e. each viewpoint represents a domain in breast oncology. In [35], the formalism C-OWL for the representation of multiple contextualized ontologies in the semantic Web is adapted for the purpose of multiple viewpoint representation and decentralized CBR.

A semantic portal for oncology.

The current version of the KASIMIR system is embedded within a semantic portal for oncology, i.e., a Web server relying on the principles and technologies of the semantic Web for providing an intelligent access to knowledge and services in oncology.

One of the main issues of the semantic Web relies on interoperability between applications and knowledge modules (e.g. ontologies). Thus, building a semantic portal implies a standardization of knowledge and software components of the KASIMIR system. For the knowledge bases, standardization relies on a sharable domain model, and leads to the definition of general ontologies in oncology. This kind of “knowledge base re-engineering” requires to replace the *ad hoc* knowledge representation formalism of KASIMIR with OWL, the knowledge representation formalism of the semantic Web. The representation of protocols is also re-engineered in order to take a better advantage of the expressiveness of the OWL formalism.

This work also implies a new software architecture, for the KASIMIR reasoner and the editing, visualization and maintenance modules. This architecture must take into account constraints related to the distributed and dynamic environment of the semantic Web. In order to query the protocols represented within OWL, an instance editor called EDHIBOU has been developed. Another interface, called NAVHIBOU, has been developed for navigating in the class hierarchies built by a reasoner based on OWL. Moreover, since the KASIMIR inference

engine is based on subsumption, a study on the integration of an extended inference engine taking into account inferences based on CBR, and the integration within the semantic Web, has to be carried out. A service of CBR based on an OWL representation has been developed for this purpose (see the thesis of Mathieu d'Aquin, defended at the end of 2005 [2], [36]).

Going further: knowledge discovery for the semantic Web.

The semantic portal of KASIMIR is operational in the sense that, given a decision protocol represented in OWL and an adaptation knowledge base, it can be used to apply or to adapt the protocol to specific situations. Besides, some ongoing research in the KASIMIR project aims at acquiring knowledge, especially adaptation knowledge, as explained above.

This is the goal of the thesis of Fadi Badra, initiated in October 2005, to combine these two research issues, i.e. how knowledge discovery techniques can be used to feed a semantic portal, and how the knowledge server embedded in this portal can be used to assist the knowledge discovery processes.

>From a longer term perspective, the goal is the following: having a clear distinction between the notions of data and knowledge, try to build a distributed system, with knowledge bases, heterogeneous data bases, inference engines, knowledge discovery modules, allowing communications with human beings, such as experts and end-users.

3.2.4. Spatial Knowledge Representation and Spatial Reasoning

In this framework, we work on two major themes, the representation of spatial structures in knowledge-based systems, and the design of reasoning models on these structures e.g. hierarchical classification and CBR. This research work is applied to answer agronomic questions regarding the recognition and the analysis of farmland spatial structures. Besides, we have been involved in the organization of the workshop RTE 2006 on spatial and temporal reasoning.

Lattice-based classification of spatial relations.

This work has been initiated during the thesis of Ludmila Mangelinck (1995–1998), in collaboration with the INRA BIA laboratory in Nancy. It has been carried out in the context of the design of a knowledge-based system for agricultural landscape analysis.

In this framework, we have designed a hierarchical representation of topological relations based on a *Galois lattice* –or *concept lattice structure*– relying on the Galois lattice theory. A Galois lattice is a multi-faceted tool for designing hierarchies of concepts: it allows the construction of a hierarchical structure both for representing knowledge and for reasoning. In a concept lattice structure, a concept may be defined by an *extension*, i.e. the set of individuals being instances of the concept, and by an *intension*, i.e. the set of properties shared by all individuals. In our framework, the extension of concepts corresponds to topological relations between regions of an image, and the intension of concepts corresponds to properties computed on that image regions (*computational operations*). Thus, a concept lattice structure emphasizes the correspondence between qualitative models, e.g. topological relations, and quantitative data, e.g. vector or raster data.

Currently, this work is continuing with a deeper study of Galois lattices for linking qualitative topological relations, and computational operations on numerical (raster or vector) data. In particular, we focus on the comparison of lattices built on different sets of relations, or computational operations [40], [10].

CBR on spatial organization graphs.

This work has been undertaken in the framework of Jean-Luc Metzger thesis (2000 – 2005), in collaboration with INRA SAD. The objective was to develop a knowledge-based system, called ROSA, for comparing and analyzing farm spatial structures. The reasoning in the ROSA system follows the principles of case-based reasoning (CBR). In our research work, CBR relies on the agronomic assumption that there exists a strong relation between the spatial and the functional organizations of farms, and thus, that similar spatial organizations correspond to similar functional organizations. According to this assumption, and given a set of previously studied farm cases, the ROSA system has to help agronomists to analyze new problems holding on land use and land management in farms.

The development of the system is achieved and tests have been done [7]. This part of the project is stopped since J.-L. Metzger left the team.

Besides, the analysis of the knowledge acquisition and modeling processes, undertaken with the help of researchers in socio-psychology and linguistics (CODISANT, LPI-GRC, Université Nancy 2 and GRIC UMR 5612 CNRS, Lyon) is continuing [3], [4].

3.2.5. *Intelligent Access to Information*

The availability and retrieval of information is of main importance in scientific and technical domains, e.g. for research and technological watch purposes. Nowadays, there is a large quantity of data available, and this requires to implement adapted tools for exploiting this mass of data. A research work holds on the definition and implementation of a toolbox allowing an “intelligent” access to information, by combining information retrieval, hypertext navigation, and data-mining. This toolbox can be used for document retrieval on the Web, bibliographical search or domain analysis.

In this framework, the design of a semantic-based algorithm for comparing and classifying documents is under investigation [27]. The annotations of documents are represented as labeled trees, where nodes and edges are composed of concepts lying in a domain ontology associated with the topics of the considered documents. A reasoning process based on classification is carried out for comparing the labeled trees representing documents, i.e. the annotations, and thus for comparing the documents. This comparison process allows to compute a semantic similarity measure between documents, and then to classify documents according to their content.

Another important idea underlying the toolbox is that data-mining and information retrieval are complementary tasks for accessing and analyzing data. Data-mining allows the guiding of information retrieval by taking advantage of the knowledge units extracted from the data, for example the extraction of a lattice from the data may provide an organization on which the information retrieval process may rely. Conversely, information retrieval allows the guiding of the data-mining process by making available information on data that can be used for example for pruning a set of extracted rules, or for providing a focus for a classification process.

>From a practical point of view, the toolbox, called “IntoWeb”, provides a set of tools for implementing the core tasks of the knowledge extraction process (see 4.3.2). For building a generic knowledge-based information retrieval system, it is needed to precisely define the kinds of objects to be manipulated, and the manipulation operations. The objects may be, among others, URL (reference to web documents), hypertextual documents (web documents), full-text documents, XML documents, vectors (sets of valued properties), etc. Operations may be applied to these objects for producing new objects, containing characteristic information, e.g. objective of a research, constraints for guiding a data mining or information retrieval process, etc. These operations may be based on information retrieval or knowledge discovery, e.g. finding all hypertextual documents identified by a set of URL, computation of the vector representation of a full-text or an XML document, extraction of an annotation tree from a textual document according to an ontology, extraction of a set of association rules from a set of XML documents, classification of web documents according to an ontology, etc.

An experiment is currently under study in the field of astronomy (MDA PROJECT, see § 5.3.3). The focus is on the building of a prototype ontology in the field of astronomy (more precisely units of measurement for celestial objects) [28]. A work has also been carried out on scientific bibliographic data for improving retrieval and navigation services on this kind of data [14]. In this way, knowledge about publications and their domain is stored in an ontology. The ontology is then used for representing concepts and their relationships and for reasoning on documents. Reasoning may help researchers by providing more efficient (focused) document retrieval and navigation.

4. Software

4.1. A Data Mining Toolkit: the Coron Platform

Keywords: *association rule extraction, data mining, frequent generators, frequent closed itemsets, frequent itemsets, rare itemsets.*

Participants: Sandy Maumus, Amedeo Napoli, Laszlo Szathmary [contact person].

One of the goals of data mining is to extract hidden relations among objects and properties in databases. Usually frequent itemsets are used to find association rules, but the process produces a large number of rules, leading to the associated problem of “mining the set of extracted rules”. Studies have shown that it can be more interesting to find only a subset of frequent itemsets, namely *frequent closed itemsets* (FCIs) and *frequent generators* (FGs). In turn, FCIs and FGs can be used for finding “minimal non-redundant” association rules.

We have developed a collection of programs for data mining that are grouped in the so-called CORON platform. The platform contains a rich set of well-known algorithms in the data mining community, such as APRIORI, APRIORI-CLOSE, CLOSE, PASCAL, ECLAT, CHARM, and, as well, several original algorithms such as PASCAL+, ZART, CARPATHIA, ECLAT-Z, and CHARM-MFI. The toolkit is composed of three main parts: (i) CORON-base, (ii) ASSRULEX, (iii) pre- and post-processing modules.

With CORON-base, it is possible to extract different kinds of itemsets, e.g. frequent itemsets, frequent closed itemsets, frequent generators, etc. Each of the algorithms has advantages and disadvantages with respect to the form of the data that are mined. Since there is no best universal algorithm for any arbitrary dataset, CORON-base offers the possibility for users to choose the algorithm that best suits their dataset and needs.

Finding association rules is one of the most important tasks in data mining. The second part of the system, ASSRULEX (Association Rule eXtractor) can generate different sets of association rules. This can lead to another data mining problem: which rules are the most useful? Beside all possible rules, some useful rule subsets can be extracted, e.g. minimal non-redundant association rules, generic basis, informative basis.

The CORON toolkit supports the whole life-cycle of a data mining task. We have modules for cleaning the input dataset, and reduce its size if necessary. The module RULEMINER facilitates the interpretation and the filtering of the extracted rules. The association rules can be filtered by (i) attribute, (ii) support, and/or (iii) confidence. It is also possible to color the most important attributes in the list of rules, for finding the most interesting rules from a given viewpoint.

Until now, studies in data mining have mainly concentrated on frequent itemsets and generation of association rules from them. Recently, we started to investigate the complement of frequent itemsets, namely the rare (or non-frequent) itemsets. In the literature, the problem of rare itemset mining and the generation of rare association rules has not yet been studied in detail, though such itemsets also contain important information just as frequent itemsets do. A particularly relevant field for rare itemsets is medical diagnosis. CORON already contains some algorithms that are designed to extract rare itemsets and rare association rules, e.g. APRIORI-RARE, ARIMA and BTB.

The CORON toolkit is developed entirely in Java, which provides a maximal portability. The system is operational, and it has already been tested within several research projects, e.g. for mining the STANISLAS cohort, or in the CABAMAKA project (which is part of the KASIMIR system, see § 3.2.3). Moreover, the CORON implementation of the TITANIC algorithm has been integrated into the GALICIA platform, that is developed at the University of Montréal, Canada.

4.2. Stochastic systems for knowledge discovery and simulation

Keywords: *Hidden Markov models, stochastic process.*

Participants: Florence Le Ber, Jean-François Mari [contact person].

4.2.1. CarottAge

One aspect of data-mining is to provide a synthetic representation of data that a domain analyst can interpret. The purpose of the CAROTTAGE system is to build a partition –called the hidden partition– in which the inherent noise of the data is withdrawn as much as possible. Then spatio-temporal data are explored for extracting homogeneous classes both in temporal and spatial dimensions, giving also a clear view of the transitions between the classes.

CAROTTAGE is a free software, under a GPL license, taking as input an array of discrete data where the rows represent the spatial sites and the columns the time slots, and building a partition with the associated *a posteriori* probability. This probability may be plotted as a function of time, and is a meaningful feature for the analyst searching for stationary and transient behaviors of data. This software is currently used by INRA researchers interested in mining the successions of land use processes, e.g. in order to build models simulating the contamination of cave and surface waters.

4.2.2. *GenExp*

In the framework of the project “Impact des OGM” initiated by the French ministry of research, we have developed a software called *GenExp* for simulating bidimensional random landscapes, and then studying the dissemination of vegetable transgenes. The *GenExp* system is based on the **CAROTTAGE** system, and on computational geometry. The simulated landscapes are given as input for programs such as *Mapod-Maïs* or *GeneSys-Colza* for studying the transgene diffusion. This year, we have released a new version of *GenExp* allowing an interaction with R subroutines. This version is on the way to receive a GPL License.

4.3. Softwares for the manipulation of documents for the Semantic Web

Keywords: *Semantic Web, association rule extraction, frequent itemset search, information retrieval, knowledge discovery from databases, navigation, text mining.*

Participants: Hacène Cherfi, Amedeo Napoli, Emmanuel Nauer [contact person], Yannick Toussaint [contact person].

4.3.1. *tamis: A software for text and rule mining*

The system, called **TAMIS** for “Text Analysis by Mining Interesting rulesS” is currently under development. This system allows the navigation through a large set of association rules, such as those produced by a text mining experiment. The **TAMIS** system is based on a user-friendly interface, and it can be easily used by non-computer scientists, e.g. analysts, experts in the domain of the analyzed data. The association rules are extracted by a mining algorithm, e.g. using the **CORON** platform in the present case, encoded in a predefined XML format. The **TAMIS** system stores the rules in a database, and proposes eight different statistical measures for sorting the rules, e.g. support, confidence, interest, conviction, dependence... In this way, the analyst may focus on smaller sets of interesting rules satisfying a given set of constraints. These constraints may be expressed by means of operations on the values of the statistical measures, and on the content of the left/right hand side of a rule.

4.3.2. *IntoWeb: Intelligent Access to Information*

Two systems are under development. A first system, called “**IntoBib**”, is a generic system designed for the exploitation of bibliographical data. Two kinds of objects are manipulated within the **IntoBib** system, namely bibliographical references and properties –or points of view– about these references, e.g. authors, keywords... The available operations on these specific objects are references filtering using one or more points of view, conceptual clustering of similar references with respect to a given point of view, and extraction of correlation between references. Accordingly, the **IntoBib** system is based on a toolbox providing a number of modules, among which, hypertext navigation, retrieval of bibliographical references, extraction of correlation between references, search for equivalent references (duplicates), conceptual clustering of similar references, normalization of fields e.g. author name, keywords...

The second system, called “**IntoWeb**”, extends the **IntoBib** system. The objective is to provide a more generic environment for an intelligent access to information, by combining information retrieval, hypertext navigation, and data-mining. The **IntoWeb** system contains a set of tools implementing the core tasks of a knowledge extraction process, i.e. collecting, filtering, and mining data. Solving a given problem of information retrieval, or data mining, is performed by a well chosen sequence of operations that are available in the system.

4.3.3. DefineCrawler: a Generic Crawler

The “DefineCrawler” system can be seen as an information retrieval “meta-system”, in the sense that it can be parameterized for satisfying different information retrieval tasks. The DefineCrawler system is based, on a classical information retrieval architecture, and on search engines available on the Web. A number of parameters have been retained, to be adjusted within an XML file for implementing and controlling different information retrieval system behaviors.

- Initialization parameters (*Start*) include the maximum depth of the crawl (*Depth*), a set of starting points for navigation (*URL*, possibly making reference to the URL of a search engine), the directory where have to be stored the data collected by the crawler (*Directory*), the number of parallel processes crawling the Web (*NbThread*), a halting condition (*Stop*) making possible the specification of a maximal crawling time, and thus ensuring a termination of the information retrieval process.
- Validation parameters (*Validation*) include a set of conditions (connected by Boolean operators) that must be satisfied by the documents, for eliminating documents without interest with respect to the query, e.g. documents that do not satisfy some criteria, that are not in a fixed language...
- Evaluation parameters within which additional conditions can be set, in order to evaluate the returned documents. The evaluation and validation conditions can be combined to calculate a score for a returned document. This score is then used to rank the returned documents.

Every validation and evaluation condition is defined by an external instruction, allowing the use of various commands or tools, e.g. for checking the presence of an element, for counting the occurrences of some elements, for calculating a similarity between documents...

4.4. Software for Spatial Reasoning

Keywords: *land organization, qualitative spatial reasoning, typological relations.*

Participant: Florence Le Ber [contact person].

Rosa, for “Reasoning on Organization of Space in Agriculture”, is a system developed in collaboration with agronomists, whose objective is to record and to maintain an agronomic knowledge base on farms, and to solve problems in agronomy, based on this knowledge base. Two kinds of knowledge elements are considered: domain knowledge, and knowledge on spatial organization and functioning of specific farms. The domain knowledge is described by a hierarchy of spatial concepts and relations (spatial occupation and relations). The spatial organization of farms is described by the so-called “space organization graphs” (SOGs) linking spatial entities through spatial relations. A vertex of a SOG (either a spatial entity or a relation) is labeled and linked to a concept of the domain knowledge hierarchy. The functioning of farms is described within “explanations” attached to SOGs. An explanation holds on a particular function of the considered farm organization and functioning. The association of a particular SOG with an explanation composes a case, to be used within a case-based reasoning process. The Rosa system is under development, and is implemented within the RACER description logic system.

4.5. The Kasimir System

Keywords: *case-based reasoning, classification-based reasoning, edition and maintenance of knowledge, semantic portal.*

Participants: Fadi Badra, Jean Lieber [contact person], Amedeo Napoli.

The objective of the KASIMIR system is decision support and knowledge management for the treatment of cancer. A number of tools have been developed within the KASIMIR system: mainly modules for the editing of treatment protocols, visualization, and maintenance. The ontology editor PROTÉGÉ has been customized for editing the KASIMIR protocols, and it has been connected with the KASIMIR inference engine. The use of the PROTÉGÉ editor involves a simplification of the protocol editing, and the detection of errors during the editing, thanks to the inference engine.

Two visualization modules have been integrated in PROTÉGÉ, allowing the display of the KASIMIR hierarchy of concepts representing the protocol being edited: PALÉTUVIER and HYPERTREE (HYPERTREE has been initially developed in the ECOO team at LORIA). The combined use of these two visualization modules, and of the classical tree widget of PROTÉGÉ, provides several useful features for hierarchy visualization, navigation, and global or focused views.

Finally, a maintenance module has been developed and integrated into PROTÉGÉ, that compares two versions of a protocol in order to separate changed and unchanged elements. This module can be used in particular during an editing session, to visualize the modifications since the beginning of the session.

Actually, two versions of KASIMIR are currently used: one version is based on an *ad hoc* object-based representation formalism, and the other version is developed within the semantic portal, as introduced in the section 3.2.3. This latter is based on OWL and on some extensions of OWL, and has motivated the development of the two user interfaces, namely EDHIBOU and NAVHIBOU, presented above. The software CABAMAKA (see also section 3.2.3) for case base mining for adaptation knowledge acquisition is part of the KASIMIR system.

5. Other Grants and Activities

5.1. The European Network of Excellence Knowledge Web

“Knowledge Web” is the name of a European network of excellence initiated in 2004. Three INRIA teams are involved in Knowledge Web, namely ACACIA at INRIA-SOPHIA, EXMO at INRIA-RHÔNE-ALPES and Orpailleur. The current World Wide Web (www) is the syntactic Web, where the structure of the content of documents is presented, while the content of documents itself is inaccessible to computers. The next generation of the Web, the Semantic Web, aims at alleviating such problem, and provide specific solutions targeted to concrete problems. The Web resources will be much easier and more readily accessible by both human and computers, with an additional semantic information in a machine-understandable and machine-processible form. The Semantic Web will have much higher impact on eWork and eCommerce than the current version of the Web already had. Still, there is a long way to go transferring the Semantic Web from an academic adventure into a technology provided by software industry. Supporting this transition process of Ontology technology from Academia to Industry is the main and major goal of the “Knowledge Web” project. This main goal naturally translates into three main objectives, given the nature of such a transformation:

- Industry requires immediate support in taking up this complex and new technology. Languages and interfaces need to be standardized to reduce the effort and provide scalability to solutions. Methods and use-cases need to be provided to convince and to provide guidelines for how to work with this technology.
- Important support to industry is provided by developing high-class education in the area of Semantic Web, Web services, and Ontologies.
- Research on Ontologies and the Semantic Web has not yet reached its goals. New areas such as the combination of Semantic Web with Web services realizing intelligent Web services require serious new research efforts.

More briefly, it is the mission of Knowledge Web to strengthen the European software industry in one of the most important areas of current computer technology: Semantic Web enabling eWork and eCommerce. Naturally, this includes education and research efforts to ensure the durability of impact and support of industry.

5.2. The Eureka GenNet Project

The research and development GenNet project is a European EUREKA-labeled project, involving two industrial societies, namely the French *KIKA medical* society, and the Belgian *Phenosystems* society. Two members of the Orpailleur group drive a so-called “thèse Cifre” on the integration of clinical and genetic data for mining and pharmacogenomics knowledge extraction. This research work is in progress, and more developments are needed before substantial results may be obtained.

5.3. National initiatives

5.3.1. ACI IMPBIO: the FouDAnGA Project

The FouDAnGA proposal, for “Fouille de données pour l’annotation de génomes d’actinomycètes” has been selected in June 2004 as an ACI IMPBIO project in bioinformatics. This project involves two research teams from LORIA (namely ADAGE and Orpailleur), and the Laboratory of Genetics and Microbiology of the University UHP Nancy 1. Since a number of years, these three teams have been collaborating within the PRST “Intelligence logicielle – Bioinformatique et applications à la génomique” (see hereafter). Being selected as an ACI IMPBIO project has reinforced and structured the initial project, allowing two students to complete their thesis.

The scientific motivation of this project is to extract subsequences from DNA with informative and significant values in molecular genetics. In particular, the signals implied in the gene regulation are under investigation. The models used correspond to the bacteria of the group of the actinomycetes –in particular to *Streptomyces*– that is the main producer of antibiotics and of metabolites with therapeutic interest, and with *Mycobacteries* –for example *M. tuberculosis*– that is responsible for tuberculosis.

A steady homogeneous second-order hidden state chain describes discrete heterogeneities distributed with a strong bias in the intergenic regions. The a posteriori observation of the hidden states specifies short DNA loci (5 to 12 pb) corresponding mostly to targets for DNA binding proteins, including transcriptional regulators. The analysis of the *Streptomyces coelicolor* genome allows the detection of the exact location of all 30 SigR promoters, as well as 92 other known or putative relevant regulatory sequences described so far. These DNA motifs represent about 7,8% of the 3000 extracted from a database corresponding to 1,15 Mb of chromosomal DNA.

5.3.2. ACI IMPBIO: the ISIBIO Project

The ISIBIO project for “Information Systems Integration in Biology” is a research project, supported since July 2004 by the Ministry of Research in the framework of the ACI IMPBIO initiative. In this interdisciplinary project, the interest is on the exploration of the role of metadata and ontologies in the integration of information systems in biology. The ISIBIO project reinforces the existing collaborations between people from different disciplines, and stimulate new interactions at both the national and the international levels, by organizing twice a year an international seminar.

The second ISIBio seminar took place in Paris (Institut Pasteur) on December 12–13, 2006 (<http://bioinfo.loria.fr/projects/isibio/isibio-presentation>). The ISIBio group also co-organized the second OGSB workshop “Ontologies, Grille et Intégration Sémantique pour la Biologie” held in conjunction with the JOBIM conference in Bordeaux, on July 4th, 2006. The third ISIBio seminar has been hold on November 21st, 2006 in Nancy.

5.3.3. ACI “Masse de données en Astronomie”

This research project “Knowledge Discovery and Ontology Design in Astronomy” is carried out in collaboration with the CDS in Strasbourg (“Centre de données astronomiques de Strasbourg”), and the IRIT computer science laboratory in Toulouse. Researchers in astronomy use every day an information network made of journal articles available under an electronic form, and a number of databases, such as the SIMBAD database recording bibliographical entries and measure sets on about three millions of astronomical objects, and the catalog server VizieR recording astronomical catalogs and measure tables published in the astronomical journals.

Interested researchers should have access to the content of documents, e.g. journal articles, astronomical object catalogs, or measure tables. For facilitating this access, researchers in astronomy have at their disposal a base of the so-called UCD for “Unified Content Descriptors”, i.e. a hierarchical database that has been extracted and designed at the CDS from the content of astronomical catalogs and tables.

The research work currently carried out in collaboration with the CDS concerns the study and the design of an ontology for representing the UCD and astronomical objects as well, starting from a collection of articles –thus involving text mining– and for extending the UCD base. This ontology will be used for a number of important and different tasks for researchers in astronomy, such as intelligent information retrieval based on the content of documents, information manipulation for matching and comparing the content of the astronomical documents. This research work can be seen as a contribution to the research works on the Semantic Web, where the purpose is to attach semantics to astronomical documents, for defining an annotation method of astronomical documents, and for a knowledge-based information retrieval method in heterogeneous astronomical sources.

A methodology for building an OWL ontology of the UCDS is currently under study [28]. The specific task in which this ontology has been used is for retrieving the UCDS representing at the best the description of an astronomical object given by a set of properties. An approximate 2-step classification process is performed by exploiting the metadata linking lexical items used in the descriptions of astronomical objects, and concept properties defining the UCDS in the ontology. The recognition of composed UCDS depending on several concepts has to be studied further. The classification of simple and composed UCDS presents similarities to the works on disjunctive classification, where concepts are defined by union of properties: in this case, owning a subset of properties for an object is sufficient to be classified as an instance of the concept.

5.3.4. *cnrs tcان Project*

A research work on Adaptation Knowledge Acquisition (AKA) for the KASIMIR system (see section 3.2.3) is carried out in the framework of the CNRS interdisciplinary project **TCAN** for “Traitement des connaissances, apprentissage et NTIC”. The objective of AKA is to provide knowledge in the form of *adaptation meta-rules*:

- Automated AKA is based on the mining of the protocols. A protocol can be seen as a set of rules *situation* \rightarrow *decision*. Knowing how the decisions change when the situations change from one rule to another rule may provide a specific adaptation rule. Clustering and generalizing these specific adaptation rules produce general adaptation rules, that have to be validated by experts.
- Supervised AKA is based on the analysis of adaptations performed by experts. Interviews of experts confronted to decision problems requiring adaptation have been recorded to be afterwards analyzed and modeled within adaptation rules.

Orpailleur is involved in this TCAN project, together with the “laboratoire d’ergonomie du CNAM in Paris”, and the Centre Alexis Vautrin in Nancy. Beyond the application framework, this research work will involve progress in the AKA methodology and techniques, that is an original research area in CBR (at its beginning, despite its importance for knowledge-intensive approaches in CBR).

5.3.5. *Projects and Collaborations in Spatio-Temporal Reasoning*

- Géomatique (CNRS–STIC): “Modélisation, comparaison et interprétation d’organisations territoriales agricoles” (in charge of Florence Le Ber).
- Impact des OGM (MENRT): “Modélisation de la dispersion de transgènes à l’échelle de paysages agricole” (in charge of Florence Le Ber).
- Eau, environnement, sociétés Ressources – Usages – Risques Gestion (CNRS–SHS): RIBAVAL project “Conception d’un outil pour la simulation du fonctionnement d’un bassin versant et définition des conditions d’utilisation pour la co-gestion” (in charge of Florence Le Ber).
- Programme fédérateur “Agriculture et Développement Durable”: Conception d’Observatoires de Pratiques Territorialisées de la Durabilité de l’Agriculture (COPTDA) (in charge of Jean-François Mari).

- Collaborations: ENGEES Strasbourg, INRA in Nancy-Mirecourt, Paris-Grignon, Dijon, and Toulouse, Laboratoire ESE UPRESA 8079 CNRS/Paris-Sud, Équipe Codisant, LPI GRC, Université de Nancy 2, GRIC UMR 5612 CNRS Lyon, and ENGREF Clermont-Ferrand.

5.4. The “PRST IL” Programme

The acronym PRST-IL stands for “Programme Régional Scientifique et Technique Intelligence Logicielle” in which is involved the LORIA Laboratory.

- The PRST IL project ILD-ISTC for “Ingénierie des langues et du document, information scientifique, technique et culturelle”.

The Orpailleur team is involved within the regional research project ILD-ISTC. In this context, research work is carried out in association with the URI team at INIST CNRS on the design of an operational text mining platform (mainly for technological watch with respect to scientific texts).

- The PRST IL project BIOINFO for “Bioinformatique et applications à la génomique” (<http://bioinfo.loria.fr/Bioinfo-Loria/>).

The Orpailleur team is involved in three main collaborations with biology laboratories, namely “Fouille de données pour l’annotation de génomes d’actinomycètes” (with the Laboratory for Microbial Genetics LGM UHP-INRA), “Vers une exploitation sémantique des sources de données biologiques du Web” (with EA 3446, CRB-INSERM-U724, EA 4002), and “Combinaison de méthodes symboliques-numériques de fouilles de données pour l’étude et l’analyse de la cohorte Stanislas” (with INSERM U525 (Équipe 4)).

6. Dissemination

6.1. Scientific Animation

- The members of the Orpailleur team are involved, as members or as head persons, in a number of national research groups.
- The members of the Orpailleur team are involved in the organization of conferences, as members of conference program committees, as members of editorial boards, and finally in the organization of journal special issues.

6.2. Teaching

- The members of the Orpailleur team are involved in teaching at all levels of teaching in the universities of Nancy (especially “Université Henri Poincaré Nancy-1” and “Université de Nancy 2”; actually, it must be noticed that most of the members of the Orpailleur team are employed on university positions).
- The members of the Orpailleur team are also involved in student supervision, again at all university levels, from under-graduate until post-graduate students.
- Finally, the members of the Orpailleur team are involved in HDR and thesis defenses, being thesis referees or thesis committee members.

6.3. Transfer activities

A delegation from the Orpailleur team was present on the INRIA booth at EUROBIO 2006 (<http://www.eurobio2006.com/>). The project “Transfer and services for Genomics and Biomolecular Modelling in Lorraine” (G-BioModel) was presented and contacts were established with five companies present on the booth: AUREUS PHARMA, GENE-IT, GENOMINING, GENOSTAR, and HELIOS BioScience. Two demonstrations, SNP-CONVERTER and VSM-G, have illustrated the activities in bioinformatics within the Orpailleur team. A flyer is available on the bioinfo web site <http://bioinfo.loria.fr/Members/pierronl/eurobio-2006/>.

6.4. Awards

- L. Szathmary, S. Maumus, P. Pétronin, Y. Toussaint and A. Napoli: Best paper award, EGC 2006 for the article “Vers l’extraction de motifs rares”. EGC 2006, 6ièmes Journées d’Extraction et Gestion des Connaissances, 18–21 janvier, Lille, France, pages 499–510, 2006.
- Sandy Maumus: PhD Thesis Award of Université Henri Poincaré (Faculty of Pharmacy), “Approche de la complexité du syndrome métabolique et de ses indicateurs de risque par la mise en oeuvre de méthodes numériques et symboliques de fouille de données. Thèse de l’Université Henri Poincaré Nancy 1, Novembre 2005.

7. Bibliography

Year Publications

Doctoral dissertations and Habilitation theses

- [1] S. MAUMUS. *Approche de la complexité du syndrome métabolique et de ses indicateurs de risque par la mise en oeuvre de méthodes numériques et symboliques de fouille de données*, PhD Thesis Award of Université Henri Poincaré Nancy 1 (Faculté de Pharmacie), Thèse de biologie, Université Henri Poincaré Nancy 1, 2005.
- [2] M. D’AQUIN. *Un portail sémantique pour la gestion des connaissances en cancérologie*, Thèse d’université, Université Henri Poincaré Nancy 1, soutenue le 15 décembre 2005, 2005.

Articles in refereed journals and book chapters

- [3] C. BRASSAC, F. LE BER. *Inscription spatiale d’une activité cognitive collective de représentation de l’espace*, in "Intellectica", vol. 2-3, n^o 41-41, 2006.
- [4] S. LARDON, F. LE BER, C. BRASSAC, P. CARON, M. MAINGUENAUD, J.-M. PRÉAU. *Conception collaborative d’objets géo-graphiques. Application aux jeux de territoire*, in "Revue internationale de Géomatique", vol. 16, n^o 2, 2006, p. 269–284.
- [5] F. LE BER, M. BENOÎT, C. SCHOTT, J.-F. MARI, C. MIGNOLET. *Studying crop sequences with CARROTAGE, a HMM-based data mining software*, in "Ecological Modelling", vol. 191, n^o 1, 2006, p. 170–185.
- [6] N. MESSAI, M.-D. DEVIGNES, M. SMAÏL-TABBONE, A. NAPOLI. *Trellis de concepts et ontologies pour interroger l’annuaire de sources de données biologiques BioRegistry*, in "Ingénierie des Systèmes d’Information", vol. 11, n^o 1, 2006, p. 39–60.
- [7] J.-L. METZGER, S. LARDON, F. LE BER. *Comparaison d’organisations spatiales agricoles : le système ROSA*, in "Revue internationale de Géomatique", vol. 16, n^o 2, 2006, p. 195–210.
- [8] H. MURAD, P. COLLET, E. BRUNNER, H. SCHOHN, P. BECUWE, M.-D. DEVIGNES, M. DAUÇA, L. DOMENJOUR. *Immunoselection and characterization of a human genomic PPAR binding fragment located within POTE genes*, in "Biochimie", 2006.
- [9] H. MURAD, P. COLLET, C. HUIN-SCHOHN, N. AL-MAKDISSY, G. KERJAN, A. CHEDOTAL, M. DONNER, M.-D. DEVIGNES, P. BECUWE, H. SCHOHN, L. DOMENJOUR, M. DAUÇA. *Effects of PPAR and RXR ligands in semaphorin 6B gene expression of human MCF-7 breast cancer cells*, in "International Journal of Oncology", vol. 28, 2006, p. 977–984.

- [10] A. NAPOLI, F. LE BER. *The Galois lattice as a hierarchical structure for topological relations*, in "Discrete Applied Mathematics", Accepted and to be published, 2006.
- [11] Y. TOUSSAINT. *Des outils informatiques pour accéder au contenu des textes : l'apport des outils de traitement de la langue et de fouille de textes*, in "Pérenniser le document numérique", L. CALDERAN, B. HIDOINE, J. MILLET (editors). , ADBS Éditions, 2006, p. 83–100.
- [12] M. D'AQUIN, J. LIEBER, A. NAPOLI. *Adaptation Knowledge Acquisition: a Case Study for Case-Based Decision Support in Oncology*, in "Computational Intelligence (an International Journal)", vol. 22, n^o 3/4, 2006, p. 161–176.
- [13] M. D'AQUIN, J. LIEBER, A. NAPOLI. *Towards a Semantic Portal for Oncology using a Description Logic with Fuzzy Concrete Domains*, in "Fuzzy Logic and the Semantic Web", E. SANCHEZ (editor). , chap. 19, Elsevier, 2006, p. 379–393.

Publications in Conferences and Workshops

- [14] S. AL-SUDANI, R. ALHULO, A. NAPOLI, E. NAUER. *OntoBib: an Ontology-Based System for the Management of a bibliography*, in "Workshop on Knowledge Management and Organizational Memories at ECAI'2006, Riva del Garda, Italy", 2006.
- [15] R. BENDAOU, Y. TOUSSAINT, A. NAPOLI. *Construction et enrichissement d'une ontologie à partir de base de textes*, in "CORIA 2006, Lyon, France", ARIA (editor). , 2006, p. 353–358.
- [16] A. COULET, M. SMAIL-TABBONE, P. BENLIAN, A. NAPOLI, M.-D. DEVIGNES. *SNP-Converter: an Ontology-Based solution to Reconcile Heterogeneous SNP Descriptions for Pharmacogenomic Studies*, in "Proceedings of the third International Workshop on Data Integration in the Life Sciences 2006 (DILS'06), Hinxton, UK", 2006.
- [17] A. COULET, M. SMAIL-TABBONE, A. NAPOLI, M.-D. DEVIGNES. *Suggested Ontology for Pharmacogenomics (SO-Pharm): Modular Construction and Preliminary Testing*, in "Proceedings of the International Workshop on Knowledge Systems in Bioinformatics (KsinBIT'06), in conjunction with OTM'06, Montpellier, France", 2006.
- [18] B. FUCHS, J. LIEBER, A. MILLE, A. NAPOLI. *Réflexions sur la place du RàPC dans trois domaines de recherche actuels*, in "Actes du quatorzième atelier raisonnement à partir de cas RàPC'06, Besançon", B. CHEBEL-MORELLO (editor). , 2006.
- [19] B. FUCHS, J. LIEBER, A. MILLE, A. NAPOLI. *Une première formalisation de la phase d'élaboration du raisonnement à partir de cas*, in "Actes du quatorzième atelier raisonnement à partir de cas RàPC'06, Besançon", B. CHEBEL-MORELLO (editor). , 2006.
- [20] C. GRAC, A. HERRMANN, F. LE BER, M. TRÉMOLIÈRES, A. BRAUD, A. HANDJA, N. LACHICHE. *Mining a database on Alsatian rivers*, in "Proceedings of the 7th International Conference on Hydroinformatics, HIC 2006, Nice, France", P. GOURBESVILLE, J. CUNGE, V. GUINOT, S.-Y. LIONG (editors). , vol. III, Research Publishing, 2006, p. 2263–2270.
- [21] N. JAY, F. KOHLER, A. NAPOLI. *Using Formal Concept Analysis for mining and interpreting patient flows within a healthcare network*, in "Proceedings of CLA 2006", 2006.

- [22] N. JAY, A. NAPOLI, F. KOHLER. *Cancer Patient Flows Discovery in DRG Databases*, in "Ubiquity: Technologies for Better Health in Aging Societies. Proceedings of MIE 2006.", A. HASMAN, R. HAUX, J. V. D. LEI, E. D. CLERCQ, F. ROGER (editors)., Studies in Health Technology and Informatics, IOS Press, 2006, p. 725–730.
- [23] N. JAY, A. NAPOLI, F. KOHLER. *Mise en évidence de filières de soins dans le traitement du cancer à l'aide des motifs fréquents et treillis de Galois.*, in "XIXièmes Journées EMOIS", 2006.
- [24] F. LE BER, C. LAVIGNE, J.-F. MARI, K. ADAMCZYK, F. ANGEVIN. *GenExp, un logiciel pour simuler des paysages agricoles, en vue de l'étude de la diffusion de transgènes*, in "Actes du Colloque International de Géomatique et d'Analyse Spatiale (SAGEO 2006), Strasbourg", C. WEBER, P. GANÇARSKI (editors)., Actes sur CD, 2006.
- [25] S. MAUMUS, A. NAPOLI, L. SZATHMARY, Y. TOUSSAINT. *Réflexions sur l'extraction de motifs rares*, in "Comptes-rendus des 13 ièmes rencontres de la Société Francophone de Classification (SFC-06), Metz", M. NADIF, F.-X. JOLLOIS (editors)., Presses Universitaires de Montréal, 2006, p. 157–162.
- [26] N. MESSAI, M.-D. DEVIGNES, A. NAPOLI, M. SMAIL-TABBONE. *BR-Explorer: An FCA-based algorithm for Information Retrieval*, in "Fourth International Conference on Concept Lattices and their Applications (CLA'06), Yasmine Hammamet, Tunisia", 2006, p. 285–290.
- [27] E. NAUER, A. NAPOLI. *A Proposal for Annotation, Semantic Similarity and Classification of Textual Documents*, in "Proceedings of the 12th International Conference on Artificial Intelligence: Methodology, Systems, Applications (AIMSA 2006), Varna, Bulgaria, 13-15th September, 2006", J. EUZENAT, J. DOMINGUE (editors)., LNAI 4183, Springer, Berlin, 2006, p. 201–212.
- [28] E. NAUER, A. RICHARD, S. DERRIÈRE, F. GENOVA, A. NAPOLI, Y. TOUSSAINT. *Construction d'une ontologie de descripteurs UCD en astronomie*, in "17ièmes journées francophones d'Ingénierie des connaissances (IC 2006), Nantes", 2006.
- [29] F. PENNERATH, A. NAPOLI. *La fouille de graphes dans les bases de données réactionnelles au service de la synthèse en chimie organique*, in "Extraction et gestion des connaissances (EGC'2006), Lille", G. RITSCHARD, C. DJERABA (editors)., RNTI-E-6, Cépaduès-Éditions Toulouse, 2006, p. 517–528.
- [30] L. SZATHMARY, S. MAUMUS, P. PETRONIN, Y. TOUSSAINT, A. NAPOLI. *Vers l'extraction de motifs rares*, in "Extraction et gestion des connaissances (EGC'2006), Lille", G. RITSCHARD, C. DJERABA (editors)., RNTI-E-6, Cépaduès-Éditions Toulouse, 2006, p. 499–510.
- [31] S. TENIER, Y. TOUSSAINT, A. NAPOLI, X. POLANCO. *Instantiation of relations for semantic annotation*, in "The 2006 IEEE/WIC/ACM International Conference on Web Intelligence - WI 2006, Hong Kong", IEEE Computer Society Press, 2006.
- [32] S. TÉNIER, A. NAPOLI, X. POLANCO, Y. TOUSSAINT. *Annotation sémantique de pages Web*, in "Extraction et gestion des connaissances (EGC'2006), Lille", G. RITSCHARD, C. DJERABA (editors)., RNTI-E-6, Cépaduès-Éditions Toulouse, 2006, p. 305–310.
- [33] M. D'AQUIN, F. BADRA, S. LAFROGNE, J. LIEBER, A. NAPOLI, L. SZATHMARY. *Adaptation Knowledge Discovery from a Case Base*, in "Proceedings of the 17th European Conference on Artificial Intelligence

(ECAI-06), Trento", G. BREWKA, S. CORADESCHI, A. PERINI, P. TRAVERSO (editors). , IOS Press, 2006, p. 795–796.

[34] M. D'AQUIN, J. COJAN, J. LIEBER, A. NAPOLI. *Vers l'implantation d'un moteur d'inférences pour une logique de descriptions avec domaine concret flou*, in "Actes des rencontres francophones sur la logique floue et ses applications (LFA-06)", 2006, p. 129–135.

[35] M. D'AQUIN. *Raisonnement à partir de cas décentralisé pour le Web sémantique*, in "Actes du quatorzième atelier raisonnement à partir de cas RàPC'06, Besançon", B. CHEBEL-MORELLO (editor). , 2006.

[36] M. D'AQUIN, J. LIEBER, A. NAPOLI. *Case-Based Reasoning within Semantic Web Technologies*, in "Twelfth International Conference on Artificial Intelligence: Methodology, Systems, Applications (AIMSA-06)", 2006, p. 190–200.

Internal Reports

[37] S. MAUMUS, L. SZATHMARY, A. NAPOLI, S. VISVIKIS-SIEST. *Towards a global methodology for mining cohorts with biological data*, Research Report, LORIA, 2006, <http://hal.inria.fr/inria-00000640>.

[38] L. SZATHMARY, S. MAUMUS, A. NAPOLI. *Mining Rare Association Rules*, Research Report, LORIA, 2006, <http://hal.inria.fr/inria-00102909>.

[39] L. SZATHMARY, A. NAPOLI, S. O. KUZNETSOV. *ZART: A Multifunctional Itemset Mining Algorithm*, Research Report, LORIA, 2006, <http://hal.inria.fr/inria-00001212>.

Miscellaneous

[40] F. LE BER, M. HUCHARD. *Variations sur l'utilisation des treillis de Galois pour la classification de connaissances et la modélisation par objets*, Conférence invitée, 2ièmes Rencontres Inter-Associations, RIAS 2006, Lyon, 2006.

[41] N. MESSAI, M.-D. DEVIGNES, A. NAPOLI, M. SMAIL-TABBONE. *BR-Explorer: A sound and complete FCA-based retrieval algorithm*, Poster at the Fourth International Conference on Formal Concept Analysis (ICFCA 2006), Dresden, Germany, 2006.