

1 International Journal of Pattern Recognition
 and Artificial Intelligence
 3 Vol. 21, No. 3 (2007) 1–16
 © World Scientific Publishing Company



5 **THE EVOLUTION OF KDD: TOWARDS DOMAIN-DRIVEN**
 DATA MINING*

7 LONGBING CAO[†] and CHENGQI ZHANG[‡]
Faculty of Information Technology
 9 *University of Technology, Sydney, Australia*
[†]lbcao@it.uts.edu.au
 11 [‡]chengqi@it.uts.edu.au

13 Traditionally, data mining is an autonomous data-driven trial-and-error process. Its typ-
 ical task is to let data tell a story disclosing hidden information, in which domain intel-
 15 ligence may not be necessary in targeting the demonstration of an algorithm. Often
 knowledge discovered is not generally interesting to business needs. Comparably, real-
 world applications rely on knowledge for taking effective actions. In retrospect of the
 17 evolution of KDD, this paper briefly introduces *domain-driven data mining* to comple-
 ment traditional KDD. *Domain intelligence* is highlighted towards actionable knowledge
 19 discovery, which involves aspects such as domain knowledge, people, environment and
 evaluation. We illustrate it through mining activity patterns in social security data.

21 *Keywords:* Data mining; knowledge actionability; domain-driven data mining.

23 **1. Introduction**

23 In the last decade, data mining, or KDD (knowledge discovery in database),¹³ has
 become an active research and development area in information technology fields. In
 25 particular, data mining is gaining rapid development in various aspects such as the
 data mined, the knowledge discovered, the techniques developed, and the appli-
 27 cations involved. Table 1 illustrates such key research and development progress
 in KDD.

29 A typical feature of traditional data mining is that KDD is presumed as an auto-
 mated process. It targets the production of automatic algorithms and tools. As a
 31 result, algorithms and tools developed have no capability to adapt to external envi-
 ronment constraints. Millions of patterns and algorithms are published in academia
 33 but unfortunately very few of them have been transferred into real business.

35 Many researchers and developers have realized the limitation of traditional
 data mining methodologies, and the gap between business and academic attention.

*This work was sponsored by Australian Research Council Discovery Grant (DP0667060,
 DP0773412), China Overseas Outstanding Talent Research Program of Chinese Academy of
 Sciences (06S3011S01), and UTS internal grants.

2 *L. Cao & C. Zhang*

Table 1. Data mining development.

Dimension	Key Research Progress
Data mined	<ul style="list-style-type: none"> — Relational, data warehouse, transactional, object-relational, active, spatial, time-series, heterogeneous, legacy, WWW — Stream, spatiotemporal, multimedia, ontology, event, activity, links, graph, text, etc.
Knowledge discovered	<ul style="list-style-type: none"> — Characters, associations, classes, clusters, discrimination, trend, deviation, outliers, etc. — Multiple and integrated functions, mining at multiple levels, mining exceptions, etc.
Techniques developed	<ul style="list-style-type: none"> — Database-oriented, association and frequent pattern analysis, multidimensional and OLAP analysis methods, classification, cluster analysis, outlier detection, machine learning, statistics, visualization, etc. — Scalable data mining, stream data mining, spatiotemporal data and multimedia data mining, biological data mining, text and Web mining, privacy-preserving data mining, event mining, link mining, ontology mining, etc.
Application involved	<ul style="list-style-type: none"> — engineering, retail market, telecommunication, banking, fraud detection, intrusion detection, stock market, etc. — Specific task-oriented mining — Biological, social network analysis, intelligence and security, etc. — Enterprise data mining, cross-organization mining, etc.

1 The research on challenges of KDD and innovative and workable KDD methodolo-
 3 gies and techniques has actually become a significant and productive direction of
 KDD. In the panel discussions of SIGKDD 2002 and 2003,^{2,9} a couple of grand
 5 challenges for extant and future data mining were identified. Among them, for
 instance, actionable knowledge discovery is one of key focuses, because it can not
 7 only afford important grounds to business decision makers for performing appro-
 priate actions, but also deliver expected outcomes to business. However, it is not a
 9 trivial task to extract actionable knowledge utilizing traditional KDD methodolo-
 gies. This situation partly results from the scenario that extant data mining is a
 data-driven trial-and-error process,² where data mining algorithms extract patterns
 11 from converted data through predefined models based on experts' hypothesis.

To bridge the gap between business and academia, it is important to understand
 13 the difference of objectives and goals of data mining in research and in real world.
 Real-world data mining presents extra constraints and expectation on mined results,
 15 for instance, financial data mining and crime pattern mining is highly constraint-
 based.^{3,9} The difference gets involved in key aspects such as the concerned problems,
 17 context mined KDD, interested patterns, the processes of mining, cared interests,
 and infrastructure supporting data mining.

19 To handle the above difference, real-world experience^{4,5} and lessons learned in
 data mining in capital markets¹⁶ show the significance of domain intelligence.

1 Domain intelligence consists of the involvement of domain knowledge²⁵ and
2 experts, the consideration of constraints, and the development of in-depth pat-
3 terns, which are essential for filtering subtle concerns while capturing incisive issues.
4 Combining these together, a sleek data mining methodology is necessary to find the
5 distilled core of a problem. They form the grounds of *domain-driven data mining*.

6 This paper provides a view of rethinking traditional data mining towards real-
7 world actionable knowledge discovery. The remainder of this paper is organized as
8 follows. Section 2 discusses the evolution of KDD. Section 3 presents major criteria
9 for measuring the actionability of knowledge. In Sec. 4, key components in domain-
10 driven data mining are stated. Section 5 briefly states domain-driven data mining
11 framework. A case study is demonstrated in Sec. 6. We conclude this research and
12 present future work in Sec. 7.

13 **2. KDD: Data Driven versus Domain Driven**

14 One of the fundamental objectives of KDD is to discover knowledge of main interest
15 to real business needs and user preference. This forms a big challenge to extant and
16 future data mining research and applications. To better understand this conflict,
17 let us review traditional data-driven data mining methodologies and research, and
18 the expectation of real world KDD.

19 **2.1. Extant data mining: Data-driven interesting pattern discovery**

20 Conceptually, there is no problem with traditional data mining, which views data
21 mining as a process of data-driven interesting pattern discovery. After all, data
22 mining targets useful information hidden in data. However, attention there has been
23 simply or mainly paid to data itself. This may be evidenced by the research scope,
24 methodologies, and research interest of traditional data mining. We may generate
25 a picture of traditional data mining by summarizing its major characteristics from
26 the following aspects: (i) object mined: data is the object being mined, which is
27 expected to tell the whole story of a concern; (ii) aims of data mining are to develop
28 innovative approaches in this period. As a result of this motivation and trend,
29 almost all high-level papers talk about new approaches; (iii) datasets mined are
30 abstract or refined from real problems or data. Mining is not directly conducted
31 on raw data from business; (iv) correspondingly, the objective of data mining is
32 to develop or update and demonstrate new algorithms on a very nice data set;
33 (v) models and methods in data mining systems are usually predefined. It is the
34 data mining researcher rather than a user that can deploy an algorithm; (vi) the
35 process of data mining is packed as automated, in which a user is not necessary
36 and actually he/she cannot do much in the mining procedure; (vii) the evaluation
37 of mined results is basically based on technical metrics, if a threshold presumed
38 by data mining researchers is higher, then the algorithm is promising; (viii) among
39 (vii) the accuracy of an algorithm is taken as one of key criteria of quality judgment.

4 *L. Cao & C. Zhang*

1 In a summary, traditional KDD is a data-driven trial-and-error process targeting
2 automated hidden knowledge discovery.^{2,7} The goal of traditional data mining is to
3 let data create/verify research innovation, demonstrate and push the use of novel
4 algorithms discovering knowledge of interest to researchers.

5 **2.2. Real world KDD: Domain-driven actionable knowledge 6 discovery**

7 In the real world, discovering knowledge actionable in solving problems concerned
8 has been viewed as the essence of KDD. However, even up to now, it is still one
9 of the great challenges to extant and future KDD as pointed out by the panel
10 of SIGKDD 2002 and 2003^{2,9} and retrospective literature. This situation partly
11 results from the limitation of traditional data mining methodologies, which do not
12 take into much consideration the constrained and dynamic environment of KDD.
13 They naturally exclude human and problem domain in the loop of data mining. As
14 a result, very often data mining research mainly aims at developing, demonstrating
15 and pushing the use of specific algorithms. As a result, it runs off the rails in
16 producing actionable knowledge of main interest to specific user needs.

17 In the wave of rethinking original objectives of KDD, the following key points
18 have recently been highlighted: comprehensive constraints around a problem,³
19 domain knowledge and human role^{2,4,12} in KDD process and environment. A proper
20 consideration of these aspects in the KDD process has been reported to make KDD
21 promising to dig out actionable knowledge satisfying real life dynamics and requests
22 even though this is a very tough issue. This pushes us to think of what knowledge
23 actionability is, and how to support actionable knowledge discovery.

24 Aiming to complement the shortcoming of traditional data mining, in particular,
25 satisfying the real user needs in enterprise data mining, we study a practical method-
26 ology, called *domain-driven data mining*.⁷ The basic theory of domain-driven data
27 mining is as follows. On top of the data-driven framework, it aims to develop proper
28 methodologies and techniques for integrating domain knowledge, human role and
29 interaction, as well as actionability measures into KDD process. It targets to dis-
30 cover actionable knowledge in a practical constrained environment. This research
31 is very important for developing the next-generation data mining methodology and
32 infrastructure.^{2,7} It can assist in a paradigm shift from “*data-driven hidden pattern*
33 *mining*” to “*domain-driven actionable knowledge discovery*”, and provides supports
34 for KDD to be translated to the real business situations as widely expected.

35 In contrast with the traditional data mining, we also list the content of domain-
36 driven data mining research and development. Most importantly, in domain-driven
37 data mining, it is data and domain intelligence (including domain knowledge and
38 domain experts) that work together to tell a hidden story in business, which discov-
39 ers actionable knowledge to satisfy real user needs. It is the user who say “yes” or
40 “no” to mined results. Table 2 compares major aspects under research of traditional
41 data-driven and domain-driven data mining.

Table 2. Data-driven versus domain-driven data mining.

Aspects	Traditional Data-Driven	Domain-Driven
Object mined	Data tells the story	Data and domain (business rules, factors etc.) tell the story
Aim	Developing innovative approaches	Generating business impacts
Objective	Algorithms are the focus	Systems are the target
Dataset	Mining abstract and refined data set	Mining constrained real life data
Extendibility	Predefined models and methods	Ad-hoc, running-time and personalized model customization
Process	Data mining is an automated process	Human is in the circle of data mining process
Evaluation	Evaluation based on technical metrics	Business say “yes” or “no”
Accuracy	Accurate and solid theoretical computation	Data mining is a kind of artwork
Goal	Let data create/verify research innovation; Demonstrate and push the use of novel algorithms discovering knowledge of interest to research	Let data and domain knowledge tell hidden story in business; discovering actionable knowledge to satisfy real user needs

3. What Makes KDD of Interest to Business

In traditional data mining, often mined patterns are nonactionable to real needs due to gaps of interests between academia and business.¹¹ Therefore, it is critical to get a clear answer to the problem “what makes KDD of interest to business”.²⁰ Answers to it may be quite varying. Basically, traditional data mining focuses on developing and refining *technical objective* measures. A typical example is those metrics developed for associations.²² Recently, *subjective* metrics are also paid attention by researchers. On the other hand, domain-driven data mining verifies and validates the usability of a pattern based not only on *technical* measures but also on *business* concerns. A more likely scenario is to integrate technical concerns with business ones, and generate an integrative measurement system to justify the quality of mined results. To this end, the concept of knowledge actionability is essential for recognizing interesting links permitting users to react to them to better service business objectives. The measurement of *knowledge actionability* should be from both *objective* and *subjective* perspectives. Table 3 summarizes the measurement of interest of data-driven versus domain-driven data mining.

Table 3. Measurement of interest of data-driven versus domain-driven data mining.

	Interest	Traditional Data-Driven	Domain-Driven
Technical	Objective	Technical objective <i>tech_obj()</i>	Technical objective <i>tech_obj()</i>
	Subjective	Technical subjective <i>tech_subj()</i>	Technical subjective <i>tech_subj()</i>
Business	Objective	—	Business objective <i>biz_obj()</i>
	Subjective	—	Business subjective <i>biz_subj()</i>
	Integrative	—	Actionability <i>act()</i>

1 3.1. Technical significance versus business expectation

2 The development of actionability is a progressive process in data mining. In the
 3 framework of traditional data mining, the so-called actionability is mainly embodied
 4 in terms of technical significance. *Technical interesting* $tech_int()$ measures whether
 5 a pattern is of interest or not in terms of specific statistical significance corre-
 6 sponding to a particular data mining method. There are two steps in technical
 7 interest evolution. The original focus basically was on *technical objective interest*
 8 $tech_obj()$,^{10,14} which aims to capture the complexities of pattern structure and sta-
 9 tistical significance. For instance, *coefficient* is developed for measuring objective
 10 interest of correlated stocks. Recent work appreciated *technical subjective* measures
 11 $tech_sub()$,^{17,19,21} which also recognize to what extent a pattern is of interest to
 12 a particular user. For example, probability-based belief is used to describe user
 13 confidence of unexpected rules.¹⁹

14 Let $X = \{x_1, x_2, \dots, x_m\}$ be a set of items, DB be a database consisting of
 15 transactions, x is an itemset in DB . Let P be interesting evidence discovered in
 16 DB through a modeling method M . For the above two procedures, we have the
 17 follows.

18 Phase 1: $\forall x \in X, \exists P: x.tech_obj(P) \longrightarrow x.act(P)$

19 Phase 2: $\forall x \in X, \exists P: x.tech_obj(P) \wedge x.tech_subj(P) \longrightarrow x.act(P)$

20 Gradually, data miners realize that the actionability of a discovered pattern must
 21 be assessed by and satisfies domain user needs. To achieve business expectations,
 22 *business interestingness* $biz_int()$ measures to what degree a pattern is of interest
 23 to a business person from social, economic, personal and psychoanalytic factors.
 24 Similar to $tech_int()$, recently *business objective interest* $biz_obj()$ is recognized by
 25 some researchers, say profit mining²⁴ and domain-driven data mining,⁷ involving
 26 $biz_int()$. At this stage, we get Phase 3 as:

27 Phase 3: $\forall x \in X, \exists P: x.tech_obj(P) \wedge x.tech_subj(P) \wedge x.biz_obj(P) \longrightarrow$
 28 $x.act(P)$

29 Moreover, *business subjective interest* $biz_sub()$ also plays essential roles in
 30 assessing $biz_int()$. This leads to a comprehensive cognition of actionability as indi-
 31 cated by Phase 4 advocated in domain-driven data mining.

32 Phase 4: $\forall x \in X, \exists e: x.tech_obj(P) \wedge x.tech_subj(P) \wedge x.biz_obj(P) \wedge$
 33 $x.biz_subj(P) \longrightarrow x.act(P)$

3.2. Knowledge actionability

34 Based on the above assessment, knowledge actionability should highlight both aca-
 35 demic and business concerns.⁷ *Actionability* recognizes technical significance of an
 36 extracted pattern that also permits users to specifically react to it to better service
 37 their business objectives. Since the satisfaction of technical interest is the antecedent
 38 of actionability, we view actionable knowledge as what satisfies not only technical
 39 interestingness $tech_int()$ but also user-specified business interest $biz_int()$. We have
 40 the following definition.
 41

1 **Definition 1.** (Knowledge Actionability) Given a mined pattern, its actionable
 2 capability $act(P)$ is described as the degree of its satisfaction with both technical
 3 and business interests.

$$\forall x \in X, \exists P : x.tech_int(P) \wedge biz_int(P) \longrightarrow act(P) \quad (1)$$

Further, it is instantiated in terms of objective and subjective factors from both technical and business sides.

$$\forall x \in X, \exists P : x.tech_obj(P) \wedge x.tech_subj(P) \wedge x.biz_obj(P) \wedge biz_subj(P) \longrightarrow act(P) \quad (2)$$

5 In this case, there are two sets of interest measures needed to be calculated
 6 when a pattern is extracted. For instance, we say a mined association trading rule
 7 is (technically) interesting because it satisfies requests on *support* and *confidence*.
 8 Moreover, if it also beats the expectation of user-specified market index return *IR*
 9 then it is a generally actionable rule.

10 In the real-world mining, business interests $biz_int()$ may differ or conflict tech-
 11 nical significance $tech_int()$. The relationship between them may present as one of
 12 four scenarios as listed in Table 4.

13 Clearly, actionable knowledge mining targets patterns confirming the relation-
 14 ship $tech_int() \Leftrightarrow biz_int()$. However, it is a kind of artwork to tune thresholds and
 15 balance significance and difference between $tech_int()$ and $biz_int()$. Quite often a
 16 pattern with high $tech_int()$ creates bad $biz_int()$. Contrarily, it is not a rare case
 17 that a pattern with low $tech_int()$ generates good $biz_int()$. In this case, it is domain
 18 users who can better tune thresholds and difference. Besides the above-discussed
 19 work on developing useful technical and business interest measures, there are some
 20 other things to do to reach and enhance knowledge actionability such as efforts on
 21 selecting actionability measures, testing actionability, enhancing actionability and
 22 assessing actionability in domain-driven data mining process.⁷

23 4. Towards Domain Driven Data Mining

24 Data mining research and development is boosted by challenges from the real world.
 25 For instance, some typical recent progress made in data mining includes stream

Table 4. Relationship between technical significance and business expectation.

Relationship Type	Explanation
$tech_int() \Leftarrow biz_int()$	The pattern e does not satisfy business expectation but technical significance
$tech_int() \Rightarrow biz_int()$	The pattern e does not satisfy technical significance but business expectation
$tech_int() \Leftrightarrow biz_int()$	The pattern e satisfies business expectation as well as technical significance
$tech_int() \not\Leftarrow biz_int()$	The pattern e satisfies neither business expectation nor technical significance

8 *L. Cao & C. Zhang*

1 data mining handling stream data, link mining studying linkage across entities.
2 Challenges and prospects coming from the real world force us to rethink some key
3 points in data mining. This includes problem understanding and definition, KDD
4 context, patterns mined, mining process, interest system, and infrastructure sup-
5 ports. The outcome of this retrospection and rethinking is a paradigm shift from
6 traditional data-driven-focused research towards domain-driven-oriented research
7 and development. The domain-driven data mining has potential for making KDD
8 available for satisfying real user needs rather than demonstrating algorithms if
9 relevant points can be appropriately considered and supported from technical, pro-
cedural and business perspectives.

11 **4.1. Problem: Domain-free versus domain-specific**

12 In traditional data mining, researchers pay a large amount of time in constructing
13 research problems, which in real-world data mining comes from real challenges. As
14 a typical phenomenon, even though a problem may come from a real scenario, it
15 is always abstracted and pruned into a very general and brilliant research issue
16 to fill in innovation and significance requirements. Such research issue is usually
17 domain-free, which means it does not necessarily involve specific domain intelli-
gence. Undoubtedly, this is important for developing the science of KDD.

18 On the other hand, in real-world scenarios, challenges always come from specific
19 domain problems. Therefore, objectives and goals of applying KDD are basically
20 problem-solving to satisfy real user needs. Problem-solving and satisfying real user
21 needs present strongly usable requirements. Requirements mainly come from a spe-
22 cific domain involving concrete functional and nonfunctional concerns. The anal-
23 ysis and modeling of these requirements request domain intelligence, in particular
24 domain background knowledge and involvement of domain experts. Therefore, real-
25 world data mining is more likely domain-specific. However, domain-specific data
26 mining is not necessarily specific domain-problem oriented. Here *domain* can refer
27 to either a big industrial sector, for instance, telecom or banking, or a categorical
28 business such as customer relationship management.

29 Domain intelligence can play significant roles in real-world data mining. Domain
30 knowledge in business field often takes forms of precise knowledge, concepts, beliefs,
31 relations, or vague preference and bias. For instance, in cross-market mining, traders
32 often take “beating market” as a personal preference to judge an identified rule’s
33 actionability. The key to taking advantage of domain knowledge in the KDD pro-
34 cess is knowledge and intelligence integration, which involves how it can be rep-
35 resented and filled into the knowledge discovery process. Ontology-based domain
36 knowledge representation, transformation and mapping between business and data
37 mining system is a proper approach to model domain knowledge. Ontology-based
38 specifications build a *business ontological domain* to represent domain knowledge in
39 terms of ontological items and semantic relationships. Ontological representation⁶
40 can be developed to manage the above items and relationships.

1 Through ontology-based representation and transformation, business terms are
mapped to data mining system's internal ontologies. We build an internal *data*
3 *mining ontological domain* for KDD system collecting standard domain-specific
terms and discovered knowledge. To match items and relationships between two
5 domains and reduce and aggregate synonymous concepts and relationships in each
domain, ontological rules, logical connectors and cardinality constraints are studied
7 to support ontological transformation from one domain to another, and semantic
aggregations of semantic relationships and ontological items intra or inter domains.

9 **4.2. KDD context: Unconstrained versus constrained**

10 Law, business rule and regulation are common forms of constraints in human soci-
11 ety. Similarly, actionable knowledge discovery can only be well conducted in a con-
strained rather than unconstrained context. Constraints involve technical, economic
13 and social aspects in the process of developing and deploying actionable knowledge.
For instance, constraints can be something involving aspects such as environmental
15 reality and expectations on data format, knowledge representation, and outcome
delivery in the mining process. Other aspects of domain constraints include domain
17 and characteristics of a problem, domain terminology, specific business process, poli-
cies and regulations, particular user profiling and favorite deliverables. In particular,
19 we highlight following types of constraints — *domain constraint*, *data constraint*,
interest constraint and *deployment constraint*.

21 The real-world business problems and requirements are often tightly embed-
ded in domain-specific business process and business rules (*domain constraint*).
23 Potential matters to satisfy or react on domain constraints may consist of building
domain model, domain metadata, semantics and ontologies,⁶ supporting human
25 involvement, human-machine interaction, qualitative and quantitative hypotheses
and conditions, merging with business processes and enterprise information infras-
27 tructure, fitting regulatory measures, conducting user profile analysis and model-
ing, etc. Relevant hot research areas include interactive mining, guided mining, and
29 knowledge and human involvement etc.

31 Patterns that are actionable to business are often hidden in large quantities of
data with complex data structures, dynamics and source distribution (*data con-*
33 *straint*). Constraints on particular data may be embodied in terms of aspects such
as very large volume, ill-structure, multimedia, diversity, high dimensions, high fre-
quency and density, distribution and privacy, etc. Data constraints seriously affect
35 the development and performance requirements of mining algorithms and systems,
and constitute some grand challenges to data mining. As a result, some popular
37 researches on data constraints-oriented issues are emerging such as stream data
mining, link mining, multirelational mining, structure-based mining, privacy min-
39 ing, multimedia mining and temporal mining.

41 Often mined patterns are not actionable to business even though they are sen-
sible to research. There may be huge conflicts of interest or gaps between academia

1 and business (*interest constraint*). What makes this rule, pattern and finding more
 2 interesting than the other? In the real world, simply emphasizing technical inter-
 3 est such as objective statistical measures of validity and surprise is not adequate.
 4 Social and economic interests (we refer to Business Interests) such as user prefer-
 5 ences and domain knowledge should be considered in assessing whether a pattern
 6 is actionable or not. Business interests may be instantiated into specific social and
 7 economic measures in terms of a problem domain. For instance, *profit*, *return* and
 8 *roi* are usually used by traders to judge whether a trading rule is interesting enough
 9 or not.

10 Furthermore, often interesting patterns cannot be deployed to real life if they are
 11 not integrated with business rules and processes (*deployment constraint*). The deliv-
 12 ery of an interesting pattern must be integrated with the domain environment such
 13 as business rules, process, information flow, presentation, etc. In addition, many
 14 other realistic issues must be considered. For instance, a software infrastructure
 15 may be established to support the full lifecycle of data mining; the infrastructure
 16 needs to integrate with the existing enterprise information systems and workflow;
 17 parallel KDD²³ may be involved with parallel supports on multiple sources, par-
 18 allel I/O, parallel algorithms, memory storage; visualization, privacy and security
 19 should receive much-deserved attention; false alarms should be minimized.

20 Some other types of constraints include knowledge type constraint, dimen-
 21 sion/level constraint and rule constraint.¹² Several types of constraints play sig-
 22 nificant roles in effectively discovering knowledge actionable to business world. In
 23 practice, many other aspects such as data stream and scalability and efficiency
 24 of algorithms may be enumerated. They consist of domain-specific, functional,
 25 nonfunctional and environmental constraints. These ubiquitous constraints form
 26 a constraint-based context for actionable knowledge discovery. All the above con-
 27 straints must, to varying degrees, be considered in relevant phases of real-world
 28 data mining. In this case, it is even called *constraint-based data mining*.^{3,12}

29 4.3. Pattern: Generic versus actionable patterns

30 Many mined patterns are more useful to data miners than to business persons.
 31 Generally interesting patterns are useful because they satisfy technical interest mea-
 32 surement. These rules are *generic patterns* or technically interest rules.

33 However, they are not necessarily useful for solving business problems. To
 34 improve this situation, we advocate in-depth pattern mining which aims to develop
 35 patterns actionable in business world. It targets the discovery of actionable patterns
 36 to support smart and effective decision-making, namely a pattern P must satisfy

$$37 \quad \forall P : x.tech_int(P) \wedge x.biz_int(P) \longrightarrow x.act(P). \quad (3)$$

38 Therefore, in-depth patterns can be delivered through improving either technical
 39 interests *tech_int()* or business interests *biz_int()*. As discussed in Sec. 3 on pattern
 40 interests, both technical and business interest measures must be satisfied from both
 41 objective and subjective perspectives.

1 Technically, it could be through enhancing or generating more effective interest
measures.¹⁸ For instance, a series of research have been done on designing right
3 interest measures for association rule mining.²² It may also be through developing
alternative models for discovering deeper patterns. Some other solutions include
5 further mining actionable patterns on a discovered pattern set. Additionally, tech-
niques can be developed to deeply understand, analyze, select and refine the target
7 data set in order to find in-depth patterns. Actionable patterns in most cases can be
created through rule reduction, model refinement or parameter tuning by optimiz-
9 ing generic patterns. In this case, actionable patterns are a revised optimal version
of generic patterns, which capture deeper characteristics and understanding of the
11 business. Of course, such patterns can also be directly discovered from data set with
sufficient consideration of business constraints.

13 On the other hand, for those generic patterns identified based on technical mea-
sures, their business interest needs to be checked so that business requirements and
15 user preference can be put into proper consideration. Domain intelligence, including
business requirements, objectives, domain knowledge and qualitative intelligence of
17 domain experts, can play roles in enhancing pattern actionability. This can be
achieved through selecting and adding business features, involving domain knowl-
19 edge, supporting interaction with users, tuning parameters and data set by domain
experts, optimizing models and parameters, adding factors into technical interest
21 measures or building business measures, improving result evaluation mechanism
through embedding domain knowledge and human involvement.

23 **4.4. Infrastructure: Automated versus human-mining-cooperated**

Traditional data mining is an automated trial and error process. Deliverables are
25 presumed as automated predefined algorithms and tools. It is arguable that such
automated methodology has both strengths and weaknesses. The good side is to
27 make user life easy. However, it meets with challenges such as a lack of capability
in involving domain intelligence and adapting to dynamic situations in the business
29 world. In particular, automated data mining has trouble in handling enterprise data
mining applications.

31 Actionable knowledge discovery in constrained context determine that real-
world data mining is more likely to be human involved rather than automated.
33 Human involvement is embodied through the cooperation between human (includ-
ing users and business analysts, mainly domain experts) and data mining sys-
35 tem. This is achieved through the complementation between human qualitative
intelligence such as domain knowledge and field supervision, and mining quantita-
37 tive intelligence like computational capability. Therefore, real-world data mining is
likely to be present as a human-machine-cooperated interactive knowledge discovery
39 process.

41 Human role can be embodied in the full period of data mining from business
and data understanding, problem definition, data integration and sampling, feature

12 *L. Cao & C. Zhang*

1 selection, hypothesis proposal, business modeling and learning to evaluation, refine-
2 ment and interpretation of algorithms and resulting outcomes. For instance, expe-
3 rience, metaknowledge and imaginary thinking of domain experts can guide or
4 assist with selection of features and models, adding business factors into modeling,
5 creating high quality hypotheses, designing interest measures by injecting business
6 concerns, and quickly evaluating mining results. This assistance can largely improve
7 the effectiveness and efficiency of mining actionable knowledge.

8 Usually, humans serve on feature selection and result evaluation. Humans
9 can play roles in a specific stage or full stages of data mining. Humans can be
10 an essential constituent or the centre of data mining system. The complexity
11 of discovering actionable knowledge in constraint-based context decides to what
12 extent and how humans must be involved. As a result, human-mining cooperation
13 presents to varying degrees, human-centered, guided mining,²⁹ or human-assisted
14 mining.

15 To support human involvement, human mining interaction, or perhaps presented
16 as interactive mining,^{1,2} is absolutely necessary. Interaction often takes explicit
17 forms, for instance, setting up direct interaction interfaces to fine tune param-
18 eters. Interaction interfaces may take various forms as well, such as visual interfaces,
19 virtual reality technique, multimodal, agents,¹⁵ etc. On the other hand, it could
20 also go through implicit mechanisms, for example accessing a knowledge base or
21 communicating with a user assistant agent. Interaction communication may be
22 message-based, model-based, or event-based. Interaction quality relies on perfor-
23 mance such as user-friendliness, flexibility, run-time capability, presentable capabil-
24 ity and understandability.

25 **5. Domain-Driven KDD Framework**

26 We have presented a domain-driven data mining framework.⁷ Domain-driven data
27 mining consists of the following key components (i) problem understanding and
28 the definition is domain-specific and must involve domain intelligence, (ii) data
29 mining is in a constraint-based context, (iii) pattern discovery targets mining in-
30 depth patterns, (iv) data mining presented as a loop-closed iterative refinement
31 process, (v) the mined results must be actionable in business, and (vi) building
32 a human-machine-cooperated infrastructure supporting domain-driven data min-
33 ing. In domain-driven framework, data mining and domain experts complement
34 each other with regards to in-depth granularity through interactive interfaces. The
35 involvement of domain experts and their knowledge can assist in developing highly
36 effective domain-specific data mining techniques and reduce the complexity of the
37 knowledge producing process in the real world. In-depth pattern mining discov-
38 ers more interesting and actionable patterns from a domain-specific perspective.
39 A system following this framework can embed effective supports for domain knowl-
40 edge and experts' feedback, and refine the lifecycle of data mining in an iterative
41 manner.

1 6. Case Study

3 Here we briefly illustrate the development of actionable activity patterns in social
 4 security data⁸ using domain-driven data mining. Taking frequent activity sequence
 5 mining as an instance, we identify those i -itemset ($i = 2, 3, 4, \dots$) frequent activity
 6 sequences likely associated with the occurrence of government customer debt using
 7 sequential association mining. Due to the imbalance of class and item distribution
 8 of debt-related activities, we split activities into two classes with domain super-
 9 vision: debt-related activity set and nondebt related activity set. To handle such
 10 unbalanced data, we develop both technical and business metrics for measuring
 11 the actionability of a pattern. The following technical metrics are defined: *global
 support, local support, class difference rate, relative risk ratio.*

12 **Definition 2.** The global support of a pattern $\{P \rightarrow \$\}$ in activity set A is
 13 defined as $Supp_A(P, \$) = |P, A|/|A|$.

14 If $Supp_A(P, \$)$ is larger than a given threshold, then P is a frequent activity
 15 sequence in A leading to debt. $Supp_A(P, \$)$ reflects the global statistical significance
 of the rule $\{P \rightarrow \$\}$ in activity set A .

16 **Definition 3.** The local support (L_SUPP) of a rule $\{P \rightarrow \$\}$ in target activity
 17 set D is defined as $Supp_A(P, \$) = |P, D|/|D|$. On the other hand, the local support
 18 of rule $\{P \rightarrow \bar{\$}\}$ in activity set $A - D$ (i.e. nondebt activity set) is defined as
 19 $Supp_{A-D}(P, \bar{\$}) = |P, A - D|/|A - D|$. The class difference rate $Cdr(P, \frac{D}{A-D})$ of P
 20 in two independent classes D and $A - D$ is defined as:

$$21 \quad Cdr(P, \frac{D}{A-D}) = Supp_D(P, \$) / Supp_{A-D}(P, \bar{\$}). \quad (4)$$

22 If $Cdr(P, \frac{D}{A-D})$ is larger than a given threshold, then P far more frequently
 23 leads to debt than nondebt. This measure indicates the difference between targeted
 24 class and untargeted class. An obvious difference between them is expected for
 25 positive frequent impact-targeted activity patterns.

26 **Definition 4.** Given local support ($SUPP$) $Supp_D(P, \$)$ and $SUPP_{A-D}(P, \bar{\$})$, the
 27 relative risk ratio $Rrr(P, \frac{\$}{\bar{\$}})$ of P leading to target activity classes D and nontarget
 28 class $A - D$ is defined as:

$$29 \quad Rrr(P, \frac{\$}{\bar{\$}}) = Prob(\$|P) / Prob(\bar{\$}|P) = Prob(P, \$) / Prob(P, \bar{\$}) \quad (5)$$

$$30 \quad Rrr(P, \frac{\$}{\bar{\$}}) = Supp_A(P, \$) / Supp_A(P, \bar{\$}) \quad (6)$$

31 If $Rrr(P, \frac{\$}{\bar{\$}})$ is larger than a given threshold, then P far more frequently leads
 32 to debt than results in nondebt. This indicates statistical difference of a sequence P
 33 leading to debt or nondebt in a global manner. An obvious difference between them
 34 is expected to distinguish frequent impact-targeted activity patterns. In addition,
 35 if the statistical significance of P leading to $\$$ and $\bar{\$}$ are compared in terms of local

Table 5. Technical interest metrics in activity sequence mining in social security area.

PATTERN	LSUPP	SUPP	CONF	LIFT	ZSCORE	$Cdr(P, _{A-D}^D)$	$Rrr(P, _{seq}^{\$})$
$A, E \longrightarrow DET$	0.0186	0.0157	0.845	1.69	3.73		

1 classes, then relative risk ratio $Rrr(P, |_{seq}^{\$})$ indicates the difference of a pattern's
 2 significance between targeted class and untargeted class as defined in Definition 4.

3 A number of sequential activity patterns are mined based on the above and
 4 traditional measures such as left side support (LSUPP), confidence (CONF), lift
 5 (LIFT) and z score (ZSCORE). For instance, the following Table 5 illustrates one
 6 sequential activity pattern ($A, E \longrightarrow DET$) likely associated with debt in balanced
 7 mix data (where A and E are activity labels).

8 We then prune this pattern set by developing business interest metrics, for
 9 instance, the following specify the impact of a mined activity sequence on averaged
 10 debt amount and debt duration: *pattern average debt amount*, and *pattern average*
 11 *debt duration*.

12 **Definition 5.** The total debt amount $d_amt()$ is the sum of all individual debt
 13 amounts $d_amt_i (i = 1, \dots, f)$ in f itemsets holding the pattern ACB . Then we get
 pattern average debt amount for the pattern ACB :

$$14 \quad \overline{d_amt()} = \sum_1^f d_amt()_i / f \quad (7)$$

15 **Definition 6.** Debt duration $d_dur()$ for pattern ACB is the average duration of
 16 all individual debt durations in f itemsets holding ACB . Debt duration $d_dur()$ of an
 17 activity is the number of days a debt keeps valid, $d_dur() = d_end_date - d_start_date$
 18 $+ 1$, where d_end_date is the day a debt is completed, d_start_date is the day a debt
 19 is activated. Pattern average debt duration $\overline{d_dur()}_i$ is defined as:

$$20 \quad \overline{d_dur()} = \sum_1^f d_dur()_i / f \quad (8)$$

21 For instance, the following lists technical and business interest measures of activ-
 22 ity sequence rule " $L, O \longrightarrow DET$ " for Australian social security benefit recipients.
 23 If the activity "O" follows "L" in customer contacts, then the customer is likely to
 24 be in government customer debt. The technical interest tells the statistical signifi-
 25 cance of this rule, while business interest shows governmental officers how important
 26 this rule leads to debt cost to the Government.

- Technical interest:

- 27
- 28
- 29 - support = 0.01251
- 30 - confidence = 0.60935
- 31 - lift = 1.2187

- 1 • Business interest:
 - 2 – $\overline{d_amt}()$ = 29,526, the averaged debt amount in cents of those debt-related
 - 3 activity sequences supporting the rule
 - 4 – $\overline{d_dur}()$ = 15.5, the averaged debt duration in days of those debt-related
 - 5 activity sequences supporting the rule

7. Conclusions and Future Work

7 The retrospection of traditional data mining has disclosed the significance of develop-
 8 ing KDD methodologies and supports targeting actionable knowledge discovery.
 9 Domain-driven data mining provides complementary supports and ideas on tradi-
 10 tional data-driven data mining. It adequately utilizes domain intelligence includ-
 11 ing domain expertise, knowledge, constraints, environment, human cooperation for
 12 deep and actionable pattern mining satisfying business expectation fitting in busi-
 13 ness rules and processes.

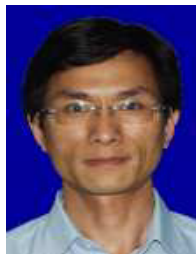
14 Domain-driven data mining has been used in telecom data mining, financial data
 15 mining and government service mining. They have shown that it has a potential to
 16 strengthen traditional KDD where a great number of rules are mined while few of
 17 them are interesting to business, and promote a wide deployment of data mining
 18 into business. Our further work is performed on qualitative analysis of the impact
 19 of domain intelligence on KDD, as well as the representation and integration of
 domain knowledge into KDD systems.

21 References

- 22 1. C. Aggarwal, Towards effective and interpretable data mining by visual interaction,
 23 *ACM SIGKDD Explor. Newslett.* **3**(2) (2002) 11–22.
- 24 2. M. Ankerst, Report on the SIGKDD-2002 panel the perfect data mining tool: inter-
 25 active or automated? *ACM SIGKDD Explor. Newslett.* **4**(2) (2002) 110–111.
- 26 3. J. F. Boulicaut and Jeudy, B. Constraint-based data mining, *The Data Mining and*
 27 *Knowledge Discovery Handbook* (Springer, 2005), pp. 399–416.
- 28 4. L. Cao and R. Dai, Human-computer cooperated intelligent information system based
 29 on multi-agents, *ACTA AUTOMATICA SINICA* **29**(1) (2003) 86–94.
- 30 5. L. Cao and R. Dai, Agent-oriented metasynthetic engineering for decision making,
 31 *Int. J. Inform. Technol. Dec. Mak.* **2**(2) (2003) 197–215.
- 32 6. L. Cao *et al.*, Ontology-based integration of business intelligence, *Web Intell. Age.*
 33 *Syst.: an Int. J.* **4**(4) (2006) (to appear).
- 34 7. L. Cao *et al.*, Domain-driven data mining: a practical methodology, *Int. J. Data*
 35 *Warehousing and Mining* **2**(4) (2006) 49–65.
- 36 8. L. Cao *et al.*, Mining impact-targeted activity patterns in unbalanced data, Technical
 37 Report, University of Technology Sydney, 2006.
- 38 9. U. Fayyad, G. Shapiro and R. Uthurusamy, Summary from the KDD-03 panel — Data
 39 mining: the next 10 years, *ACM SIGKDD Explor. Newslett.* **5**(2) (2003) 191–196.
- 40 10. A. A. Freitas, On objective measures of rule surprisingness, *PKDD98*, 1998, pp. 1–9.
- 41 11. O. F. Gur Ali and W. A. Wallace, Bridging the gap between business objectives and
 parameters of data mining algorithms, *Decision Support Syst.* **21** (1997) 3–15.

16 L. Cao & C. Zhang

- 1 12. J. Han, Towards human-centered, constraint-based, multi-dimensional data mining, An invited talk at Univ. Minnesota, Minneapolis, Minnesota, 1999.
- 3 13. J. Han and M. Kamber, *Data Mining: Concepts and Techniques*, 2nd edition (Morgan Kaufmann, 2006).
- 5 14. R. J. Hilderman and H. J. Hamilton, Applying objective interestingness measures in data mining systems, *PKDD00*, 2000, pp. 432–439.
- 7 15. M. Klusch *et al.*, The role of agents in distributed data mining: issues and benefits, *Proc. IAT03*, 2003, pp. 211–217.
- 9 16. L. Lin and L. Cao, Mining in-depth patterns in stock market, *Int. J. Intell. Syst. Technol. Appl.* 2006 (to appear).
- 11 17. B. Liu, W. Hsu, S. Chen and Y. Ma, Analyzing subjective interestingness of association rules, *IEEE Intell. Syst.* **15**(5) (2000) 47–55.
- 13 18. E. Omiecinski, Alternative interest measures for mining associations, *IEEE Trans. Knowl. Data Engin.* **15** (2003) 57–69.
- 15 19. B. Padmanabhan and A. Tuzhilin, A belief-driven method for discovering unexpected patterns, *KDD-98*, 1998, pp. 94–100.
- 17 20. A. Silberschatz and A. Tuzhilin, What makes patterns interesting in knowledge discovery systems, *IEEE Trans. Knowl. Data Engin.* **8**(6) (1996) 970–974.
- 19 21. A. Silberschatz and A. Tuzhilin, On subjective measures of interestingness in knowledge discovery, *Knowl. Discov. Data Min.* 1995, pp. 275–281.
- 21 22. P. Tan, V. Kumar and J. Srivastava, Selecting the right interestingness measure for association patterns, *SIGKDD*, 2002, pp. 32–41.
- 23 23. D. Taniar and J. W. Rahayu, Chapter 13: Parallel data mining, *Data Mining: A Heuristic Approach*, eds. H. A. Abbass, R. Sarker and C. Newton (Idea Group Publishing, 2002), pp. 261–289.
- 25 24. K. Wang, S. Zhou and J. Han, Profit mining: From patterns to actions, *EBDT*, 2002.
- 27 25. S. Yoon, L. Henschen, E. Park and S. Makki, Using domain knowledge in knowledge discovery, *Proc. Eighth Int. Conf. Information and Knowledge Management (ACM Press, 1999)*.
- 29



Longbing Cao is a IEEE Senior Member, received the Ph.D. in complex systems and intelligence Sciences from Chinese Academy of Sciences, and the Ph.D. in computer science from University of Technology, Sydney.

Currently, he is with Faculty of Information Technology, University of Technology, Sydney, Australia.



Chengqi Zhang is a IEEE Senior Member, received the Ph.D. in computer science from the University of Queensland, and D.Sc. degree in computer science from Deakin University. Currently, he is an research professor

with Faculty of Information Technology, University of Technology, Sydney, Australia.