1	
2	The Data Avalanche is Here
3	Harvey J. Miller
4	Department of Geography
5	University of Utah
6	260 S. Central Campus Dr. Room 270
7	Salt Lake City, Utah 94112-9155
8	harvey.miller@geog.utah.edu
9	

10 1. Introduction

For most of its history, science has operated in a data-poor environment. Measurements of reality were difficult, expensive and cumbersome to obtain, store and manipulate. Consequently, much of the apparatus of science is designed to tease information from scarce observations. This has changed dramatically in recent decades. Science has moved from a data-poor to a data-rich environment. The costs of capturing, storing and manipulating digital data have collapsed to a stunning degree, and communications and information technologies are widely deployed in professional and personal settings.

18

19 In a recent paper in *Science*, Lazer et al. (2009) note that a computational social science is 20 emerging that is based on the capacity to collect and analyze massive amounts of data on 21 individual and group behavior. However, it is emerging in the private sector such as 22 Yahoo and Google and in government agencies such as the U.S. National Security 23 Agency. Little evidence of this approach appears in the major journals in the social and 24 economic sciences. The authors fear that computational social science will become the 25 exclusive domain of private companies, government agencies and a privileged set of 26 academics working with these entities on research that cannot be critiqued, published and 27 replicated. This will not facilitate the advancement of science or serve the broader public 28 interest in the accumulation and dissemination of knowledge.

29

In regional science, we have access to an unprecedented amount of fine-grained data on cities, transportation, economies and societies, much of these data referenced in geospace and time. There is a tremendous opportunity to discover new insights and knowledge about spatial economies that can inform theory and modeling in regional 1 science, as well as policy and infrastructure decisions. Yet, mirroring the larger trend 2 identified by Lazer et al. (2009), there is little presence in the mainstream literature in 3 regional science. While some activity is emerging, the level of apparent activity is far 4 short of the potential. Computational approaches to discovering patterns in spatio-5 temporal data mostly reside in the technical literatures in computer science, spatial 6 analysis and geographic information science. While this work is valuable, it does not 7 reflect the specific needs of researchers in regional science. It also does not reflect the 8 rich body of theory and models in regional science; valuable sources of background 9 knowledge that can help guide the exploration of massive spatio-temporal databases.

10

11 This paper addresses the potential for discovering new knowledge in regional science 12 through exploration of spatio-temporal databases. There are well-established methods for 13 knowledge discovery from databases, and a growing body of techniques tailored for 14 spatio-temporal data. A reason for the slow adoption in regional science may be due to a 15 lack of awareness about these techniques. But an equally formidable obstacle is 16 misunderstanding of the role of data mining and knowledge discovery in regional 17 science. Rather than being atheoretical or anti-theoretic, the knowledge discovery 18 process harmonizes well with traditional avenues to knowledge construction in science. 19 In fact, the knowledge discovery process benefits from domain expertise and theory to 20 focus searching through vast information spaces and distinguish between real and 21 spurious patterns discovered in these spaces.

22

23 The next section of this paper provides a brief history of knowledge discovery from 24 databases, with special reference to the parent disciplines of regional science, namely, 25 economics and geography. Section 3 discusses the new sources of data that are relevant 26 to regional science: fine-grained data on the individual people, objects and money 27 flowing through a spatial economic system. Section 4 reviews the general process of 28 discovering new knowledge from databases, while Section 5 make the case for a 29 specialized methods for geo-spatial data. Section 6 discusses the relationships between 30 knowledge discovery and regional science, with special emphasis on the role of theory. 31 This section also raises several challenges to closer integration between knowledge discovery and theory-building in regional science. Section 7 concludes the paper by
 identifying some additional challenges.

3

Throughout the remainder of this paper, I will use the term "knowledge discovery" to
describe the general process of exploring massive digital databases for novel information.
"Data mining" refers to the specific techniques applied in one stage of this process. This
is a critical distinction that will be more apparent later in the paper (Section 4).

8

9 2. <u>A Brief History of Knowledge Discovery from Databases</u>

10 To date, developments in knowledge discovery and data mining have mostly been driven 11 by computer scientists and others in related fields. The main contributors are the 12 subfields of machine learning and database research. The former involves the study of 13 how machines and humans can learn from data. Machine learning has its origins in 14 artificial intelligence, so early work in the 1950s and 1960s attempted to simulate human 15 learning in computers. This focus subsequently became more pragmatic and continued 16 research focused on developing algorithms and methods that could learn and perform 17 well on specific tasks. Database researchers, in contrast, became interested in knowledge 18 discovery as data warehousing began to grow in the 1990s. A *data warehouse* archives the records in the transactional databases of an organizational, often for liability purposes. 19 20 As data accumulated, interest grew in exploring these repositories for information to 21 support strategic and tactical decision making. Consequently, many of the exploratory 22 techniques developed from database research focus explicitly on the relational data model 23 (Smyth 2000).

24

A legacy of knowledge discovery's origins in machine learning and database research is a focus on algorithms and rules and a relative lack of traditional statistical concepts such as parameter estimation and testing. In addition, early applications of digital computers for data analysis in the 1960s lead to the conclusion that one can always search long enough and find a complex but often spurious model that will fit a dataset arbitrarily well. Consequently, "data mining" still has a negative connotation in many fields such as econometrics (Smyth 2000). 1

2 One of the more visible debates about data mining in the economics literature concerns 3 issues surrounding specification-search in statistical model-building. An early study by 4 Lovell (1983) analyzes specification-search methods with a simulated dataset and 5 suggests rules for deflating exaggerated claims of significance. Hoover and Perez (1999) extended this simulation design to evaluate general-to-specific modeling where one starts 6 7 with a complex model and reduces to a more elegant one. The theory is, given enough 8 data, only the true specification will survive a sufficiently stringent battery of statistical 9 tests designed to pare variables from the model. This "data mining" approach contrasts 10 with the traditional *specific-to-general* strategy where one starts with a spare model based 11 on theory and conservatively builds a more complex model. Hoover and Perez (1999) 12 report generally favorable results about the ability of the general-to-specific approach to 13 uncover the true model underlying the data. Campos and Ericsson (1999) and Hendry 14 and Krolzig (1999) also support this positive conclusion. However, Hand (1999) 15 dismisses this as irrelevant, claiming that given the size of the information space 16 involved, one would be extraordinarily lucky to configure an initial model that contains 17 the true one. Since one can only hope to discover an approximation of the true model, 18 the only relevant criterion is predictive performance not explanatory validity. This is a 19 controversial position, particularly in a theory-driven science such as economics 20 (Feelders 2002).

21

22 Knowledge discovery methods in geography are less controversial than in economics due 23 to the rise of Geographic Information Systems (GIS) as a vehicle for spatial data 24 management and analysis. However, there is a growing concern about potential privacy 25 violations that may occurs from exploring geospatial data, especially when integrating 26 heterogeneous spatial databases. The concept of *locational privacy* results from these 27 concerns, and protocols are emerging to diminish or eliminate these risks without 28 destroying the utility of these data in basic and applied research (see Dobson and Fisher 29 2003; Duckham, Kulik and Birtley 2006).

1 As I will suggest below, controversy surrounding the scientific value of knowledge 2 discovery methods results from a misunderstanding of the relationship knowledge 3 discovery and statistics, as well as the relationship between knowledge discovery and 4 theory. Knowledge discovery complements rather than replaces traditional statistics by 5 providing an enhanced process for hypothesis generation; these hypotheses can (and 6 should!) be evaluated against theory as well as through the techniques of confirmatory 7 hypothesis testing. Rather than being antagonistic towards theory, knowledge discovery 8 benefits from a strong theoretical base: this can help guide the search process as well as 9 help evaluate the patterns that emerge from this process.

10

11 3. <u>New Data Sources</u>

12 The key opportunity for economic geographers and regional scientists at this juncture is 13 not just the increasing power of computers, the improving sophistication of knowledge 14 discovery methods, or even the avalanche of available data. Rather, it is the increasing 15 availability of fine-grained data on the location of individual people and objects densely 16 with respect to time, as well as data about online searches, transactions, social 17 connections and commentaries. These data can allow an unprecedented view of societies 18 and economies from the "bottom-up", as well as the aggregate structures, processes and 19 dynamics that emerge from the interactions of individuals in markets, institutions, and 20 regions. Although massive databases have been available for decades, the fine-grained 21 nature of these data is relatively new, and techniques are emerging for exploring these 22 data. This section discusses some of these new data sources.

23

24 **3.1. Data collection technologies**

Point-of-sale data. A motivation behind the wider diffusion of knowledge discovery techniques beyond the database and machine learning communities is the availability of point of sale (POS) data through product barcoding. POS data facilitates *market basket analysis*: the analysis and prediction of customer buying habits by finding associations among items that customers place in their shopping basket (Han and Kamber 2006). Indeed, one of the earliest data mining techniques used outside the database and machine learning communities was association rule mining (Smyth 2000). These data remain important for knowledge discovery in economics, marketing and related fields,
 particularly when combined with data associated with customer loyalty programs.
 However, new sources of data, especially data on individual behavior in space and time,
 is making knowledge discovery more relevant for geography and regional science.

5

6 Location-aware technologies. Location-aware technologies (LATs) are devices that can 7 report their geographic location densely with respect to time. Major strategies for 8 locational referencing include *radiolocation methods* based on the time, time difference, 9 or angle of the signals' arrivals at base stations from mobile clients, the global 10 *positioning system* (GPS) that exploits time differences of signals arriving from satellites 11 in Earth orbit and *interpolation methods* that use distances and directions from a known 12 location along a route to determine the current location (Greiner-Brzezinska 2004). 13 LATs enable location-based services (LBS). LBS provide information to individuals 14 based on their geographic location though client devices such as mobile phones (Benson 15 2001). Typical LBS include concierge services, navigation, social networking, 16 emergency response, fleet management, local news and tourist information (Spiekermann 17 2004). Worldwide deployment levels may reach 1 billion devices by 2010 (Bennahum 18 2001; Smyth 2001). LBS could be a very rich source of data on human activities in space 19 and time if issues regarding privacy and propriety can be resolved (Miller 2007).

20

21 An increasingly important LAT is radiofrequency identification (RFID) tags. RFID tags 22 attached to objects can transmit data to fixed readers using passive or active methods. 23 Passive tags are cheaper, smaller, and lighter, but have a very limited range. Also, 24 readers cannot track multiple passive tags simultaneously. Active tags are heavier and 25 more expensive since they contain a power source, but have a longer range and readers 26 than can track multiple tags simultaneously. In both systems, the RFID tags must self-27 identify because the reader conducts the location calculations. This means that RFID 28 systems have greater potential for surveillance and privacy violations than systems such 29 as the GPS where the client conducts the location referencing (Morville 2005). RFID 30 tags are becoming a central feature in the retailing industry due to their capabilities for 31 real-time supply chain management (Roberti 2003). Other RFID applications include automated toll collection, passports, airline baggage tracking, and VIP services at hotels,
 clubs and resorts. These applications can allow individual tracking through the products
 and services that individuals use, but have significant potential for privacy violations
 (Eckfeldt 2005; McGinty 2004; Shih et al. 2005).

- 5
- 6

7 Geosensor networks. Another technology that can capture data on activities in space 8 and time are geosensor networks. These are interconnected, communicating, and 9 georeferenced computing devices that monitor a geographic environment. The 10 geographic scales monitored can range from a single room to an entire city or ecosystem. 11 The devices are typically heterogeneous, ranging from temperature and humidity sensors 12 to video cameras and other imagery capture devices. Geosensor networks can also 13 capture the evolution of the phenomenon or environment over time. Geosensor networks 14 can provide fixed stations for tracking individual objects, identify traffic patterns and 15 determine possible stops for a vehicle, as it travels across a given domain in the absence 16 of mobile technologies such as GPS or RFID (Stefanidis 2006; Stefanidis and Nittel 17 2004).

18

19 The Internet. LATs and geosesnor networks only capture human activities in real 20 space. An increasing number of activities and interactions are occurring in *cyberspace*: 21 the domain implied by the world's collective information and communication 22 technologies. Phone calls, emails, texts and other forms of interpersonal communication 23 leave records that can be used to understand social networks, group interactions and 24 social dynamics. The World Wide Web provides a vast repository of what people are 25 saying, buying, searching and connecting with each other. Advances in natural language 26 processing as well as methods for analyzing media data are improving the ability to 27 evaluate social and economic behavior as mediated by these technologies (Lazer et al. 28 2009).

29

Brockmann, Hufnagel and Geisel (2006) provide a clever illustration of using web-based
data in knowledge discovery about human behavior. The website "Where's George?"

1 (www.wheresgeorge.com) allows users to input the serial number of a U.S. or Canadian 2 bill and the current location of the user. If the bill has been registered previously, a list of 3 all the locations and times where the bill has been appears; if the bill has not been 4 registered, a new list is created. This is clearly an incomplete, and likely biased, dataset. 5 Nevertheless, Brockmann, Hufnagel and Geisel (2006) analyzed the trajectories of 6 approximately 460,000 bills using this website and were able to conclude that human 7 travel as evidenced by the money flow patterns is consistent with a continuous-time random walk process. Although not a knowledge discovery exercise, this research 8 9 illustrates that the utility of the noisy, non-scientifically sampled surrogate data that is 10 readily available on the web.

- 11
- 12

13 **3.2. Simulation**

14 The increasing availability of fine-grained empirical data is only part of the story. 15 Another part is the increasing ability to *generate* vast amounts of synthetic data by 16 simulating large and complex systems at the individual level. The avalanche of CPU 17 cycles bequeathed by improvements in computing engineering (and as described by the 18 now famous Moore's Law) is only one motivation. Equally important is a growing 19 recognition that complex systems such as cities and societies cannot be understood 20 through reductionist approaches. Rather, complex patterns and dynamics emerge from 21 the interactions of individual components of the systems: the whole is more than the sum 22 of the parts (Flake 1998).

23

Two major traditions for disaggregate modeling in regional science are microsimulation and agent-based modeling. *Microsimulation* is the older tradition: this refers to the modeling and analysis of phenomena at a disaggregate level to order to better understand its aggregate behavior. Microsimulation has a substantial history in social science, dating back to attempts to modeling the US economy in the 1950s (Clarke and Holm 1987). There are well-established standards and techniques for model estimation and validation (Boman and Holm 2004).

1 Agent-based modeling (ABM) simulates the dynamics of complex systems through the 2 behaviors and interactions of its individual units or *agents*. An agent is an independent, 3 goal-driven that are typically autonomous (it makes decisions based on its inputs and 4 goals without an external controlling mechanism) and *adaptive* (its behavior can improve 5 over time through a learning process). Agents interact by exchanging physical or virtual 6 (informational) resources (Maes 1995). Agents can represent people, households, 7 animals, firms, organizations, regions, countries, and so on, depending on the scale of the 8 analysis and the elemental units hypothesized for that scale. The increasing availability of 9 high-resolution data and GIS tools for handling these data facilitate ABM in simulating 10 human spatial systems such as cities and economies (Benenson and Torrens 2004). 11 Applications of ABM include economics (Epstein 1999; Tesfatsion 2009), land-use/land-12 cover change (Parker et al. 2003), social dynamics (Epstein and Axell 1996), 13 transportation (Balmer et al. 2004), and human movement at micoscales (Batty et al. 14 2003). ABM offers a rigorous but rich approach to simulating human phenomena from 15 the bottom-up, as well as the concepts of adaptation, self-organization, and emergence to 16 capture linkages between individual behavior and aggregate dynamics (Boman and Holm 17 2004).

18

19 Both microsimulation and ABM have potential for facilitating deeper understanding of 20 complex physical and human systems. Both techniques generate voluminous and 21 intricate results; essentially, massive spatio-temporal databases. Making sense of the 22 results of a large-scale simulation is the same challenge as understanding similarly 23 detailed and massive empirical data about the system in the real-world. This is 24 particularly challenging when conducting multiple simulation runs within an 25 experimental design, as is good practice.

26 27

28 4. What is Knowledge Discovery from Databases?

Knowledge discovery from databases (KDD) is based on a belief that information in the form of *interesting* patterns is hidden in massive databases. "Interesting" means that the information is *easily understood* by humans, *valid* for generalization, potentially *useful* and *novel* (Han and Kamber 2001). KDD is also predicated on the belief that traditional
 analysis methods are not capable of discovering the hidden and interesting information in
 massive, heterogeneous and nonscientific databases.

4

5 KDD accommodates data not normally amenable to statistical analysis. Statistical 6 techniques typically requires a clean (relatively noise free) numeric database 7 scientifically sampled from a large population with specific questions in mind. Many 8 statistical models require strict assumptions (such as independence, stationarity of 9 underlying processes, and normality). In contrast, the empirical and synthetic data being 10 generated and stored in many databases are noisy, non-numeric and possibly incomplete. 11 These data are also collected in an open-ended manner without specific questions in mind 12 or were generated as a byproduct of another activity (Hand 1998). KDD in its 13 contemporary form encompasses techniques from statistics, machine learning, pattern 14 recognition, numeric search and scientific visualization to accommodate the new data 15 types and data volumes being generated through information technologies.

16

17 **4.1. The Knowledge Discovery Process**

18 The KDD process usually consists of several stages, corresponding to major tasks 19 involved in preparing and exploring the data as well as interpreting results (Adriaans and 20 Zantinge 1996; Brachman and Anand 1996; Fayyad, Piatetsky-Shapiro and Smyth 1996; 21 Han and Kamber 2006; Matheus, Chan and Piatetsky-Shapiro 1993). Although the order 22 of the stages below represent a standard progression though the KDD process, in practice 23 they may not be executed in any particular order. Some steps may be skipped and others 24 repeated, depending on the judgment of the analyst guiding the process and the 25 intermediate results. This is a key point. KDD is not an automated, push-button process: 26 it demands intelligence decision-making, domain expertise and thoughtful reflection. 27 These are skills that human minds still do much better than computers.

- 28
- 29 30

• *Data selection* involves determining a subset of the records or variables in a database for knowledge discovery.

1 Data pre-processing involves "cleaning" the selected data to remove noise, ٠ 2 eliminating duplicate records, and handling missing data fields and domain 3 violations. The pre-processing step may also include *data enrichment* through 4 combining the selected data with other, external data (e.g., census data, market 5 data).

- 6 Data reduction and projection diminishes the dimensionality of the data through • 7 transformations to equivalent but more efficient representations of the information 8 space.
- 9 Data mining involves the application of low-level algorithms to uncover hidden • 10 patterns in the data.
- 11
- Interpreting and reporting stage involves evaluating, understanding and • 12 communicating the information extracted from the data.
- 13

14 4.2. Data Mining Techniques

Data mining involves the application of low-level functions or algorithms for revealing 15 hidden information in a database (Klösgen and Żytkow 1996). There are several major 16 17 classes of data mining functions: the type of information being sought determines the 18 particular type of data mining function to be applied (Han and Kamber 2006). Table 1 19 summarizes major data mining tasks. All of these techniques share a common 20 characteristic: *scalability*, or the ability to handle massive databases without unreasonable 21 increases in computational time.

1

Data mining task	Description	Techniques
Segmentation or clustering	Determine a finite set of implicit groups that describes the data.	Cluster analysis
Classification	Predict the class label that a set of data belongs to based on some training datasets	 Bayesian classification Decision tree induction Artificial neural networks Support vector machine
Association	Find relationships among data objects; predict the value of some attribute based on the value of other attributes	Association rulesBayesian networks
Deviations	Find data items that exhibit unusual deviations from expectations	Cluster analysisOutlier detectionEvolution analysis
Trends	Lines and curves summarizing the database, often over time	RegressionSequential pattern extraction
Generalizations	Compact descriptions of the data	Summary rulesAttribute-oriented induction

Table 1: Data mining tasks and techniques (Miller and Han 2009)

2

3 Segmentation or clustering involves partitioning a selected set of data into meaningful 4 groupings or classes. The commonly used data mining technique of *cluster analysis* 5 determines a set of classes and assignments to these classes based on the relative 6 proximity of data items in the information space. Cluster analysis methods for data 7 mining must accommodate the large data volumes and high dimensionalities of interest in 8 data mining; this usually requires statistical approximation or heuristics. *Classification* 9 refers to finding rules or methods to assign data items into pre-existing classes. There are 10 many classification methods developed in many years of research in statistics, pattern 11 recognition, machine learning and data mining, including decision tree induction, naïve 12 Bayesian classification, neural networks and support vector machines. Associations are 13 rules that predict the relationships between data objects or the value of some attribute 14 based on the value of other attributes. Deviations are data items that exhibit unexpected 15 deviations or differences from some norm. These cases are either errors that should be 16 corrected/ignored or represent unusual cases that are worthy of additional investigation. *Trends* are lines and curves fitted to the data; techniques include linear and logistic regression analysis that are very fast and easy to estimate. These methods are often combined with filtering techniques such as stepwise regression. *Generalization and characterization* are compact descriptions of the database. Techniques include *summary rules* that generate a relatively small set of logical statements that condense the information in the database.

7

8 4.3. Scientific Visualization and Knowledge Discovery

9 Visualization is a powerful strategy for integrating high-level human intelligence and 10 knowledge into the KDD process. The human visual system is extremely effective at 11 recognizing patterns, trends and anomalies. The visual acuity and pattern spotting 12 capabilities of humans can be exploited in many stages of the KDD process, including 13 OLAP, query formulation, technique selection and interpretation of results. These 14 capabilities have yet to be surpassed by machine-based approaches (Fayyad, Grinstein 15 and Wierse 2001).

16

17

18 5. Knowledge Discovery in Regional Science

19 **5.1. Geospatial Knowledge Discovery**

20 Geospatial knowledge discovery (GKD) is the process of extracting knowledge massive 21 georeferenced databases. GKD has emerged as a subdomain of KDD due to the unique 22 requirements of geospatial data, information and knowledge. There are several reasons 23 why GKD is unique.

24

Geospatial information is not only highly dimensioned, but also have the property that up to four dimensions (representing space and time) form a framework for the remaining dimensions. Fidelity with the real-world requires embedding georeferenced observations within a formal space that reflects empirical geographic and temporal relationships as faithfully as possible (or appropriate). Euclidean space is the most common framework, and often adopted implicitly by default, but it is only one of an infinite number of possibility. Techniques exist for estimating the geo-space implied by a matrix of distance, cost or flows between geographic locations and analyzing differences among
these spaces (see Ahmed and Miller 2007; Tobler 1994). Whatever space is adopted, it
carries with it a rich spectrum of implied relations among the embedded spatial objects,
including proximity, connectivity and direction (Miller and Wentz 2003). The
information implicit in the geographic measurement framework is ignored in many
knowledge discovery tools (Gahegan 2000a).

7

8 Geographic data usually exhibit the properties of *spatial dependency* and *spatial* 9 *heterogeneity.* Spatial dependency is the tendency of attributes to be more related with 10 proximity in geographic space due to the impeding effects of distance (Tobler 1970). 11 Spatial heterogeneity refers to the non-stationarity of most geographic processes. The 12 "friction" of distance combined with the relative uniqueness of each location means that geographic processes are local. Spatial dependency and spatial heterogeneity have 13 historically been regarded as nuisances confounding standard statistical techniques that 14 15 typically require independence and stationarity assumptions. However, these can also be 16 valuable sources of information about the geographic phenomena under investigation.

17

18 Spatial objects tend to be more complex than the non-spatial objects, particularly with 19 respect to their geometric footprint. Spatial objects often cannot be reduced to points in 20 some information without doing great harm to the representation of the real-world 21 phenomenon. Size, shape and boundary properties of geographic entities often affect 22 geospatial processes, sometimes due to measurement artifacts (e.g., recording flow only 23 when it crosses some fiat boundary), but often due to physical and human properties in 24 the real world (e.g., a mountain range forming a natural boundary between two nations, 25 or a divided highway creating an equally pervasive boundary between two urban 26 neighborhoods). Relationships such as distance, direction and connectivity are more 27 complex with dimensional objects (see Egenhofer and Herring 1994; Okabe and Miller 28 1996; Peuquet and Ci-Xiang 1987). The number and complexity of implied spatial 29 relationships increase dramatically with dimensional spatial objects such as lines, 30 polygons, surfaces and solids.

1 In addition to changes in their non-spatial attributes, spatial objects can also undergo 2 change over time with respect to their geometry and spatial identity. We often refer to 3 the former type of change as *motion*: this can include changes in morphology, location, or 4 both. Motion is deceptively complex concept: it can also occur in the whole object or its 5 composite parts, it can be continuous or discrete with respect to space and time and can 6 be conceptualized at the individual or collective scales (Galton 1995, 1997). Spatial 7 identity changes include events such as being created, destroyed, aggregated, 8 disaggregated, fusion with other spatial objects and fission into new spatial objects. New 9 spatial objects may also be spawned or cloned from existing objects (Frank 2001; 10 Hornsby and Egenhofer 2000; Medak 2001).

11

12 **5.2. Geospatial Data Mining Techniques**

13 Many of the traditional data mining tasks discussed above have analogous tasks in the 14 geographic data mining (Ester, Kriegel and Sander 1997; Han and Kamber 2006). 15 Spatial classification builds up classification models based on a relevant set of attributes 16 and attribute values that determine an effective mapping of spatial objects into predefined 17 target classes. Spatial clustering groups spatial objects such that objects in the same 18 group are similar and objects in different groups are unlike each other. Clustering can be 19 based on combinations of non-spatial attributes, spatial attributes (e.g., shape) and 20 proximity of the objects or events in space, time and space-time. Spatial trend detection 21 involves finding patterns of change with respect to the neighborhood of some spatial 22 object. Spatial characterization and generalization is therefore an important geographic 23 data mining task. Generalization-based data mining can follow one of two strategies in 24 the geographic case. Spatial dominant generalization first spatially aggregates the data 25 and then applies standard attribute-oriented induction method at each geographic 26 aggregation level. Non-spatial dominant generalization generates aggregated spatial 27 units that share the same high-level semantic description. Spatial association rules are 28 association rules that include spatial predicates in the precedent or antecedent.

29

30 **5.3. Geovisualization**

1 Geographic visualization (GVis) is the integration of cartography, GIS, and scientific 2 visualization to explore geographic data and communicate geographic information to 3 private or public audiences (see MacEachren and Kraak 1997). GVis and GKD are 4 highly complementary, perhaps even more than KDD and scientific visualization due to 5 importance of geometry and position in spatial data. Integration between GVis and 6 GKD can occur at several levels, including determining high-level goals for the GKD 7 process, specification of appropriate geographic data mining tasks for achieving the high-8 level goals and choices include specific tools and algorithms to achieve the data mining 9 task specified (MacEachren et al. 1999).

10

11 6. <u>Regional Science and Knowledge Discovery: The Role of Theory</u>

12 KDD is sometimes dismissed by those outside its community as a black-box technique, 13 or as a "fishing expedition." It should be clear by this point in the paper that this is not 14 the case: KDD is complex and require thoughtful guidance by an interested and 15 knowledgeable expert. Nevertheless, too often the "expert" is a computer scientist who 16 enjoys building and testing tools for the sake of the tools themselves. KDD in general 17 has not spread widely into the domain sciences beyond the fields of marketing and some 18 of the physical sciences such as chemistry and medicine. GKD is also not well known or 19 widely used outside of the GIS community: it is almost invisible in economic geography 20 and regional science. This is unfortunate since regional scientists are missing 21 opportunities to advance their science by exploiting newly available data on economic 22 systems and increasingly sophisticated tools for exploring these data. They are also 23 missing an opportunity to influence the development of these tools for the greater benefit 24 of regional scientists.

25

GKD has a supporting role in modeling and theory-building: it does not attempt to replace or diminish these processes. Also, there is a crucial role for theory in supporting the GKD process.

29

30 6.1. Knowledge Discovery to Support Theory

1 KDD is essentially a hypothesis generation process; much more so than traditional 2 inferential statistics. Although statistical analysis is appropriately viewed as an inductive 3 process, it is often embedded within a broader deductive process. Statistical models are 4 confirmatory, requiring the analyst to specify a model based on some theory, test the 5 hypotheses generated by that model, and revise the theory depending on the results. In 6 contrast, the deeply hidden patterns sought through KDD are difficult or impossible to 7 specify a priori, at least with any reasonable degree of completeness. The number of 8 potential hypothesis implied by a massive database is often too large to test exhaustively, 9 even if many can be dismissed a priori as trivial or absurd. KDD is more concerned 10 about prompting investigators to formulate *new* predictions and hypotheses from data: 11 KDD seek novel information which (by definition) is surprising and would not have 12 otherwise come to mind (Elder and Pregibon 1996; Hand 1998). In this sense, KDD is 13 similar to a telescope or microscope: it is a way for researchers to look at the data in 14 different ways, discover something new, and formulate theories, models or hypothesis 15 based on this novel information within the context of existing theory and knowledge.

16

17 In addition to supporting the hypothesis formulation stage of scientific knowledge 18 construction, Roddick and Lees (2009) argue that the knowledge discovery process can 19 work directly in concert with traditional knowledge construction in science. Echoing the 20 concerns expressed by Hand (1999), they note that exponentially large information space 21 implied by a massive database may lead to an induction fallacy where hypothesis 22 developed from the data are consistent but do not reflect the true model. Roddick and 23 Less (2009) suggest a strategy that not only focuses the search for patterns through data 24 selection and reduction, but also can support conceptual model building. The investigator 25 supplies a conceptual model as a starting point for generating experimental hypotheses 26 and a set of associated null hypotheses. Data mining techniques focus on the subregions 27 of the information space indicated by these hypotheses. In accordance with the usual 28 notion of scientific induction, confidence in the conceptual model increases if discovered 29 patterns support its hypotheses. Unsupported hypotheses suggest a modification to the 30 conceptual model or the investigator's reasoning about the phenomenon. This process 31 can also accept competing conceptual models and provide a ranking based on support for the data. However, this is a modeling building and data focusing strategy, not a modeling
 testing procedure: final acceptance of rejection of a model should be based on standard
 confirmatory procedures and tests.

4

5 6.2. Theory to Support Knowledge Discovery

6 Despite data reduction techniques and the computational efficiency of data mining tools, 7 the information space implied by a massive, heterogeneous database may be so large or 8 so complex that it cannot be exhaustively explored. A good data mining system should 9 be able to generate all of the interesting patterns in a database, but only the interesting 10 ones. Unfortunately, while some dimensions of interestingness can be such as validity 11 can be assessed automatically, the qualities of comprehensible, useful and novel are more 12 difficult to assess using automated methods. Consequently, it is typical for a data mining 13 system to generate an unacceptably large number of patterns, most of which could be 14 dismissed as difficult to interpret, useless or trivial. This is particularly true in GKD: the 15 number of spatial predicates and transformations over time is so large that the number of 16 candidate patterns to evaluate can be overwhelming (Roddick and Lees 2009).

17

18 *Background knowledge* is strategy to maximize the completeness of a data mining 19 system. It is a computational representation of domain knowledge to focus data mining 20 techniques on subregions of the information space that are likely to contain interesting 21 Background knowledge reflects known facts about the domain being patterns. 22 investigated in a way that can be exploited by data mining algorithms. For example, a 23 common type of background knowledge is a *concept hierarchy*: this is a sequence of 24 mappings from low-level to high-level semantics. It is typically organized as a tree 25 whose leaves correspond to measured attributes in the database and parents represent 26 higher level semantics that synoptically summarize the lower-level semantics of its 27 children. Concept hierarchies provide a logical map for aggregating or drilling-down 28 attributes automatically during the data mining process (Han and Kamber 2006). 29 Background knowledge can be derived from assumptions about the system, observable 30 facts, expert knowledge, or theory.

1 Interestingness measures quantitatively separate interesting patterns from uninteresting 2 ones by assessing the simplicity, certainty, utility and novelty of the generated patterns 3 (Silberschatz and Tuzhilin 1996; Tan, Kumar, and Srivastava 2002). These can be used 4 to guide data mining algorithms using feedback, or evaluate the discovered patterns after 5 discovery. There are different types of interestingness measures defined across the 6 dimensions mentioned above, and some are specific to the type of data mining technique. 7 For example, the resulting *rule length* when expressing the pattern in conjunctive normal 8 form is a measure of simplicity (and therefore legibility). A simple *confidence* measure 9 for association rules is the number of times two facts appear together relative to the 10 appearance of one fact alone. *Support* refers to the number of times in the database for 11 which a pattern is true; this measures the general utility of the pattern (Han and Kamber 12 2006).

13

Computer scientists have become very clever at developing ways to focus the knowledge discovery process. However, they are starving for content: the domain-specific facts that can guide discovery. Regional science, on the other hand, have a very rich body of theory and models that can be used as background knowledge to guide the discovery process and as interestingness measure to filter spurious patterns from potentially interesting ones. However, there are three challenges to translating domain knowledge from regional science to the discovery process

21

22 First, similar to many human sciences, regional science concepts can be abstract, vague, 23 fluid and multi-level. For example, the core concept of a "region" can very depending on 24 the context. Similarly, definitions of entities such as "household," "firm" and "city" are 25 not always clear or consistent across different research inestigations. This is 26 understandable given the complex nature of the phenomena of interest in regional 27 science, but it does point to a need for more consistent definitions, even if these are 28 flexible. In brief, regional science needs a clear ontology that provides formal definitions 29 of its basic core concepts and the relationships among these concepts. Jackson (1994) 30 advocates the use of object-oriented modeling languages in regional science, arguing that 31 this can improve not only the modeling process itself through facilitating incremental

model building, but also the communication and collaboration among regional scientists.
The latter point is the one that is particularly relevant for the discussion here: the development of standards and definitions through object-oriented modeling languages can allow the creation and maintenance of common and extensible libraries of core regional science concepts that can be shared across researchers and applications.

6

7 A second challenge concerns differences in knowledge representation between regional 8 science versus computer science. In regional science, most knowledge is implicit: it is 9 embedded in formal theories, models and equations. However, knowledge discovery 10 techniques require *explicit* representations in the form of rules, hierarchies and concept 11 networks. Translating knowledge embedded in regional science theories and models to a 12 form that can be exploited by automated data mining techniques is not trivial: these explicit representations must not only be effective (can capture all of the implicit 13 14 knowledge, but not more) but *efficient* (can be applied without undue computational 15 burden).

16

17 A third challenge is a need for effective methods for spatial pattern evaluation when 18 comparing reality to theory. If an observed spatial pattern matches a theoretical 19 prediction, this is interesting but not novel (unexpected). At the other extreme, if an 20 observed spatial pattern matches a spatial null (the trivial pattern that would be expected 21 by chance) is neither interesting nor novel. The spatial patterns that fall between theory 22 and null (in other words, those that are not trivial and at variance with theory), may be 23 both interesting and novel, and therefore worthy of further investigation. However, there 24 is no clear metric for comparing observed spatial patterns to null and theoretical spatial 25 patterns. One difficulty is the lack of a good spatial null. Michael Goodchild has argued 26 persuasively that complete spatial randomness (CSR) is a "straw target:" the likelihood of 27 any spatial pattern following CSR is so remote that this null model is likely to be rejected 28 in all cases. But, what is the alternative? It is unclear that there is a general spatial null 29 model across all domains and applications, and the appropriate spatial null will depend on 30 the behavioral and other processes assumed or postulated in the specific theory or model. 31 Also, it is unclear how to measure deviations from null or expected spatial patterns.

1 Differences between observed and null or expected patterns can be expressed 2 geometrically, but it is not clear how a geometric difference can be translated into one-3 dimensional, unambiguous quantitative measure of departure from a norm.

4 5

6 **7. Other Challenges**

In additional to the challenges mention above, there are challenges facing the wider use
of knowledge discovery techniques in regional science.

9

10 Semantically poor data. LATs and ICTs can generate a vast amount of fine-grained 11 data on peoples' activities in geo-space and cyberspace. However, much of these data are 12 semantically poor. For example, while it is easy to use vehicle-based GPS receivers or 13 location-enabled mobile phones to capture individual trajectories in space, it is more 14 difficult to get the attributes associated with these trajectories such as the characteristics 15 of the person, the activities conducted, or the planned activities that could not be 16 conducted. Some of the information can be recovered by integrating heterogeneous data 17 (for example, map-matching trajectories to detailed land-use data). However, this can be 18 cumbersome, raise privacy concerns, and introduce error. Although existing efforts to 19 exploit surrogate and auxiliary data to recover or infer data semantics should continue, 20 working with voluminous but semantically poor data may require asking some research 21 questions differently, or asking the same questions in a different way.

22

23 Algorithms and infrastructures. Spatial models and techniques can be computationally 24 complex. In theory, a spatial model requires evaluation of pairwise distances among all 25 locations in space. This implies quadratic computational times in the worse-case: this is 26 unacceptably slow for knowledge discovery and data mining. The worse-case can be 27 avoided in most cases since we know that most spatial processes attenuate with distance 28 (indeed, this is why Waldo Tobler originally coined the often-cited "Tobler's First of 29 Law of Geography" in Tobler (1970): it was a justification for skipping calculations 30 between entities that are spatially remote). However, in general this is a heuristic, and the 31 implications of these shortcuts for the knowledge discovery process must be understood.

1

There is also a need for more effective use of parallel and high-performance computing environments in spatial knowledge discovery. The use of parallel, grid and cloudcomputing environments is seriously under-researched in regional science, geography and even Geographic Information Science despite the fact that many spatial models are easily adapted to these environments due to the inherent parallelism of spatio-temporal data, as well as parallelism in the computational tasks required (see Armstrong, Wang and Cowles 2005).

9

10 Education. Finally, educating the next generation of economists, regional scientists, and 11 quantitative geographers is a major challenge. It is already difficult to provide students 12 with the necessary background in economic and geographic theory to conduct cutting-13 edge research in regional science: adding computer science to this mix will be 14 challenging, especially since the language and ways of thinking in these disparate fields 15 are in many ways orthogonal to each other. Also, given the dynamic nature of digital 16 technologies, it will be necessary to continue education throughout one's career in order 17 to keep abreast of developments in the field. I am not saying that regional scientists 18 should become computer scientists: in fact, I think this would be a very bad thing. But, 19 regional scientists need the appropriate expertise to choose, implement and interpret 20 knowledge discovery and data mining techniques effectively.

21

1 Literature Cited

23	Adriaans P. and Zantinge, D. (1996) Data Mining; Harlow, U.K.: Addison-Wesley
4 5	Ahmed, N. and Miller, H. J. (2007), "Time-space transformations of geographic space for exploring, analyzing, and visualizing transportation systems," <i>Journal of</i>
6 7	Transport Geography, 15, 2-17.
8 9	Armstrong, M.P., Wang, S., and Cowles, M. (2005) "Using a computational grid for geographic information analysis," <i>Professional Geographer</i> , 57, 365-375.
10	
11	Balmer, M., Nagel, K., and Raney, B. (2004) "Large-scale multi-agent simulations for
12 13	transportation applications," <i>Journal of Intelligent Transportation Systems: Technology, Planning, and Operations,</i> 8, 205–221.
14	
15 16	Batty, M., Desyllas, J., and Duxbury, E. (2003) "The discrete dynamics of small-scale spatial events: agent-based models of mobility in carnivals and street parades."
10 17 18	International Journal of Geographical Information Science 17, 673–697.
10	Penerson I and Tomana D. M. (2004) Cassimulations Automata haved Modeling of
20	Urban Phenomena. Chichester, UK: John Wiley & Sons.
21	
22 23	Bennahum, D. S. (2001) "Be here now," <i>Wired</i> , 9 (11), 159–163
24 25 26	Benson, J. (2001) "LBS technology delivers information where and when its needed," Business Geographics, 9 (2), pp. 20–22.
27 28 29 30	Boman, M., and Holm, E. (2004) "Multi-agent systems, time geography and microsimulations," in Olsson, MO. and Sjöstedt, G. (eds) <i>Systems approaches and their applications</i> . Dordrecht, The Netherlands: Kluwer Academic, 95–118.
31 32 33 34 35	 Brachman, R. J. and Anand, T. (1996) "The process of knowledge-discovery in databases: A human-centered approach," in U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth and R. Ulthurusamy (eds.) <i>Advances in Knowledge Discovery and Data Mining</i>, Cambridge, MA: MIT Press 37-57.
36 37 38	Brockmann, D., Hufnagel, L. and Geisel, T. (2006) "The scaling laws of human travel," <i>Nature</i> , 439, 462-465.
39 40 41	Campos, J. and Ericsson, N. R. (1999) "Constructive data mining: Modeling consumers' expenditures in Venezuela," <i>Econometrics Journal</i> , 2, 226-240.
42 43 44	Clarke, M., and Holm, E. (1987) "Microsimulation methods in spatial analysis and planning," <i>Geografiska Annaler</i> 69B, 145–164.
45 46	Dobson, J. E. and Fisher, P. F. (2003) "Geoslavery," <i>IEEE Technology and Society</i> <i>Magazine</i> 22(1): 47-52

1	
23	Duckham, M., Kulik, L. and Birtley, A. (2006) "A spatiotemporal model of strategies and counter-strategies for location privacy protection" in M. Raubal, H. I. Miller, A.
4	U Frank and M F Goodchild (eds.) Geographic Information Science: 4th
5	International Conference, GIScience 2006, Münster, Germany, September 2006
6	Proceedings, Springer, 47-64.
7	
8	EckTeldt, B. (2005) what does RFID do for the consumer? Communications of the $ACM 48$ (0) 77 70
9	ACM 48 (9), 77-79.
10	Egenhofer M I and Herring I P (1004) "Categorizing binary topological relations
12	between regions lines and points in geographic databases " in M Egenhofer D
13	M Mark and I R Herring (eds.) The 9-intersection: Formalism and its Use for
14	Natural-language Spatial Predicates. National Center for Geographic Information
15	and Analysis Technical Report 94-1, 1-28.
16	
17	Elder, J. and Pregibon, D. (1996) "A statistical perspective on knowledge discovery," pp.
18	83-113 in U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth and R. Ulthurusamy
19	(eds.) Advances in Knowledge Discovery and Data Mining, Cambridge, MA: MIT
20	Press, 83-113
21	
22	Epstein, J. M. (1999). "Agent-based computational models and generative social
23	science," Complexity 4 (5), 41–60.
24	
25	Epstein, J. M., and Axtell, R. (1996). Growing Artificial Societies: Social Science from
20	the Bottom Up. Cambridge, MA: MIT Press
21	Ester M Kriegel H-P and Sander I (1997) "Spatial data mining: A database
20	approach " M Scholl and A Voisard (eds.) Advances in Spatial Databases
30	Lecture Notes in Computer Science 1262 Berlin: Springer 47-66
31	Lecture rotes in computer Science 1202, Derini, Springer, 17 00.
32	Favvad, U., Grinstein, G. and Wierse, A. (2001) Information Visualization in Data
33	Mining and Knowledge Discovery, San Matel, CA: Morgan Kaufmann.
34	
35	Fayyad, U. M., Piatetsky-Shapiro, G. and Smyth, P. (1996) "From data mining to
36	knowledge discovery: An overview" in U.M. Fayyad, G. Piatetsky-Shapiro, P.
37	Smyth and R. Ulthurusamy (eds.) Advances in Knowledge Discovery and Data
38	Mining, Cambridge, MA: MIT Press, 1-34.
39	
40	Feelders, A. (2002) "Data mining in economic science," in J. Meij (ed.) Dealing with the
41	Data Flood: Mining Data, Text and Multimedia, The Haag, The Netherlands:
42	STT/Beweton, 165-1/5.
43	Elaka C W (1009) The Commutational Density of M (C) (E) (
44 15	Flake, G. w. (1998) The Computational Beauty of Nature: Computer Explorations of Exact als Chaos Complex Systems and Adaptation Combined MAXMIT Dress
45 46	Tracials, Chuos, Complex systems, and Adaptation, Camonage, MA: MIT Pless.
то	

1 2 3	Frank, A. (2001) ""Socio-economic units: Their life and motion," in A. Frank, J. Raper, JP. Cheylan (eds.) <i>Life and Motion of Socio-economics Units</i> , London: Taylor and Francis, 21-34.
4 5 6 7	Galton, A. (1995) "Towards a qualitative theory of movement," in A. Frank and W. Kuhn (eds.) Spatial Information Theory: Proceedings of the European Conference on Spatial Information Theory (COSIT), Springer-Verlag, 377-396.
o 9 10	Galton, A. (1997) "Space, time and movement," in O. Stock (ed.) Spatial and temporal Reasoning, Kluwer, 321-353
11 12 13	Gahegan, M. (2000a) "On the application of inductive machine learning tools to geographical analysis," <i>Geographical Analysis</i> , 32, 113-139
14 15 16 17	Grejner-Brzezinska, D. (2004). "Positioning and tracking approaches and technologies," in Karimi, H. A. and Hammad, A. (eds) <i>Telegeoinformatics: Location-based</i> <i>Computing and Services</i> , Boca Raton, FL: CRC Press, pp. 69–110.
19 20 21	Han, J. and Kamber, M. (2006), <i>Data Mining: Concepts and Techniques</i> , 2 nd ed., San Matel, CA: Morgan Kaufmann.
21 22 23 24	Hand, D. J. (1998) "Data mining: Statistics and more?" American Statistician, 52, 112- 118
25 26 27 28	Hand, D. J. (1999) "Discussion contribution on 'Data mining reconsidered: Encompassing and the general-to-specific approach to specification search' by Hoover and Perez," <i>Econometrics Journal</i> , 2, 241-243.
29 30 31	Hendry, D. F. and Krolzig, HM. (1999) "Improving on 'Data mining reconsidered,' by K. D. Hoover and S. J. Perez," <i>Econometrics Journal</i>, 2, 202-219.
32 33 34 35	Hoover, K. D. and Perez, S. J. (1999) "Data mining reconsidered: Encompassing and the general-to-specific approach to specification search," <i>Econometrics Journal</i> , 2, 167-191.
36 37 38 39	Hornsby, K. and Egenhofer, M. J. (2000) "Identity-based change: A foundation for spatio-temporal knowledge representation," <i>International Journal of Geographical Information Science</i> , 14, 207-224
40 41 42 43	Jackson, R. W. (1994) "Object-oriented modeling in regional science: An advocacy view," <i>Papers in Regional Science</i> , 73, 347-367.
44 45	Klösgen, W. and Żytkow, J. M. (1996) "Knowledge discovery in databases terminology," in U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth and R. Ulthurusamy (eds.)

1 Advances in Knowledge Discovery and Data Mining, Cambridge, MA: MIT Press 2 573-592. 3 4 Lazer, D., Pentland, A., Adamic, L., Aral, S., Barabási, A.-L., Brewer, D. Christakis, 5 N., Contractor, N., Fowler, J., Gutmann, M., Jebara, T., King, G., Macy, M., Roy, D. and Van Alstyne, M. (2009) "Computational Social Science," Science, 6 7 323, 721 – 723. 8 9 Maes, P. (1995) "Modeling adaptive autonomous agents," in Langton, C. (ed.) Artificial 10 Life: An Overview. Cambridge, MA: MIT Press, 135–162. 11 12 MacEachren, A. M. and Kraak, M.-J. (1997) "Exploratory cartographic visualization: 13 Advancing the agenda," Computers and Geosciences, 23, 335-343. 14 MacEachren, A. M., Wachowicz, M., Edsall, R., Haug, D. and Masters, R. (1999) 15 "Constructing knowledge from multivariate spatiotemporal data: Integrating 16 geographic visualization with knowledge discovery in database methods," 17 International Journal of Geographical Information Science, 13, 311-334 18 Matheus, C. J., Chan, P. K. and Piatetsky-Shapiro, G. (1993) "Systems for knowledge 19 discovery in databases," IEEE Transactions on Knowledge and Data Engineering, 20 5,903-913. 21 22 McGinty, M. (2004) "RFID: is this game of tag fair play?" Communications of the ACM 23 47 (1), pp. 15–18. 24 25 Medak, D. (2001) "Lifestyles," in A. Frank, J. Raper, J.-P. Cheylan (eds.) Life and 26 Motion of Socio-economics Units, London: Taylor and Francis, 139-153. 27 28 Miller, H. J. (2007) "Place-based versus people-based geographic information science," 29 Geography Compass, 1, 503-535. 30 31 Miller, H. J. and Han, J. (2009) "Geographic data mining and knowledge discovery: An 32 overview," in Miller, H. J. and Han, J. Geographic Knowledge Discovery and 33 Data Mining, second edition. CRC Press. 34 35 Miller, H. J. and Wentz, E. A. (2003) "Representation and spatial analysis in geographic 36 information systems," Annals of the Association of American Geographers, 93, 37 574-594. 38 39 Morville, P. (2005) Ambient findability: What we Find Changes Who we Become. 40 Sebastopol, CA: O'Reilly Media. 41 42 Okabe, A. and Miller, H. J. (1996) "Exact computational methods for calculating 43 distances between objects in a cartographic database," Cartography and 44 Geographic Information Systems, 23, 180-195. 45

1 2 3 4 5	Parker, D. C., Manson, S. M., Janssen, M. A., Hoffmann, M. J. and Deadman, P. (2003). "Multi-agent systems for the simulation of land-use and land-cover change: a review," Annals of the Association of American Geographers 93, 314– 337.
6 7 8 9	Peuquet, D. J. and Ci-Xiang, Z. (1987) "An algorithm to determine the directional relationship between arbitrarily-shaped polygons in the plane," <i>Pattern Recognition</i> , 20, 65-74.
10 11 12	Roberti, M. (2003) "Wal-Mart spells out RFID vision," <i>RFID Journal</i> [online]. Retrieved on 16 June 2003 from http://www.rfidjournal.com/
13 14 15 16	Roddick, J. F. and Lees, B. G. (2009) "Spatio-temporal data mining: Paradigms and methodologies," in H. J. Miller and J. Han (eds.) <i>Geographic Data Mining and Knowledge Discovery</i> , 2ed., CRC Press, in press.
17 18 19	Shih, DH., Lin, CY., and Lin, B. (2005) "RFID tags: Privacy and security aspects," <i>International Journal of Mobile Communications</i> , 3, 214–230.
20 21 22 23	Silberschatz, A. and Tuzhilin, A. (1996), "What makes patterns interesting in knowledge discovery systems", <i>IEEE Transactions on Knowledge and Data Engineering</i> , 8, 970-974.
24 25 26 27	Smyth, C. S. (2001). "Mining mobile trajectories," in Miller, H. J. and Han, J. (eds) Geographic Data Mining and Knowledge Discovery. London: Taylor and Francis, 337–361.
28 29 30 31	Spiekerman, S. (2004) "General aspects of location-based services," in Schiller, J. and Voisard, A. (eds) <i>Location-based services</i> . San Francisco, CA: Morgan- Kaufmann, pp. 9–26.
32 33 34	Stefanidis, A. (2006). The emergence of geosensor networks. <i>Location Intelligence</i> [online]. Retrieved on 27 February 2006 from http://locationintelligence.net
35 36 37	Stefanidis, A., and Nittel, S. (eds) (2004). Geosensor networks. Boca Raton, FL: CRC Press.
38 39 40 41	Tan, PN., Kumar, V. and Srivastava, J. (2002) "Selecting the right interestingness measure for association patterns", <i>Proceedings</i> , 2002 ACM SIGKDD International Conference on Knowledge Discovery in Databases (KDD'02), Edmonton, Canada, 32-41
42 43 44 45 46	Tesfatsion, L. (2009) <i>Agent-Based Computational Economics: Growing Economies from</i> <i>the Bottom Up.</i> Website; <u>http://www.econ.iastate.edu/tesfatsi/ace.htm</u> . Last accessed 2 April 2009.

- Tobler, W. R. (1970) "A computer movie simulating urban growth in the Detroit region.
 Economic Geography, 46, 234–40.
- 3
- 4 Tobler, W. R. (1994) "Bidimensional regression," *Geographical Analysis*, 26, 187–212.