

# ASPECTE CHEIE ÎN GEOSPATIAL DATA MINING

dr. Angela Ioniță

Institutul Național pentru Cercetare-Dezvoltare în Informatică – ICI, București

**Rezumat:** Organizat în patru capitole, acest articol este un rezumat al Raportului Tehnic 630/2, elaborat în mai 2002, în cadrul Proiectului INTAS nr. 397: DATA MINING TECHNOLOGIES AND IMAGE PROCESSING: THEORY AND APPLICATIONS. În primul capitol, sunt prezentate aspecte legate de definiția termenilor *data mining* și *knowledge discovery* pentru ca, în capitolul al doilea, să se urmărească prezentarea caracteristicilor datelor geospațiale, care trebuie luate în considerare în contextul *data mining*, în general, și în contextul obiectivelor acestui proiect, în special. Următorul capitol este dedicat sistemelor software relevante pentru *geospatial data mining* (GeoMiner, ADaM, SPIN!). Ultimul capitol prezintă câteva concluzii privind impactul asupra științelor legate de informația geografică și aspectele critice al cercetărilor în domeniu. Menționăm, de asemenea, că acest articol, ca și întregul Raport Tehnic 630/2, este o sinteză bazată pe articolele, cărțile și prezentările menționate în capitolul dedicat referințelor bibliografice, pe care le-am considerat relevante pentru demersul din cadrul activităților Proiectului INTAS nr. 397.

Facem de la început precizarea că termenul "data mining" ar putea fi tradus ca "*mineritul datelor*". Deoarece până la data publicării acestui articol nu există consens asupra termenului în limba română, pentru a evita confuziile, îl vom utiliza pe cel în limba engleză. La fel vor fi utilizați și ceilalți termeni asupra cărora se mai lucrează în ceea ce privește completarea definițiilor (datorită dinamicii cercetării) și atribuirii unui termen românesc adecvat.

**Cuvinte cheie:** geographical information, geographical knowledge, Geographical Information Systems (GIS), data mining/knowledge discovery.

## Context

Acest articol este rezumatul Raportului tehnic 630/2, executat în cadrul activităților specifice Proiectului INTAS no. 39: DATA MINING TECHNOLOGIES AND IMAGE PROCESSING: THEORY AND APPLICATIONS.

Început în 1 noiembrie 2001, cu o durată de trei ani, obiectivele acestui proiect sunt:

- elaborarea metodelor pentru găsirea de șabloane, descrieri de date de volum mare și construirea de modele predictive;
- elaborarea și investigarea de metode de distilare șablon, bazate pe abordarea colectivă în clasificarea nesupervizată;
- dezvoltarea și investigarea de noi metode de căutare și analiză a dependențelor neliniare complexe;
- elaborarea și investigarea noilor metode pentru soluționarea activităților dinamice în cadrul *data mining*;
- efectuarea de prognoze agricole, bazate pe studierea factorilor meteorologici și pe geomonitorizarea regiunilor;
- modelarea curgerii râurilor folosind rețele neuronale;
- elaborarea de noi metode de discriminarea structurilor spațiale locale, în cadrul tipurilor de date GIS;
- elaborarea de noi instrumente pentru analiza texturii bazate pe caracteristici universale și morfologie matematică;
- elaborarea de modele de calcul adecvate pentru implementarea rețelelor neuronale ierarhice, pentru procese industriale.

Echipa acestui proiect este formată din:

- Lappeenranta University of Technology, Finland - coordonator;
- ICI - Romania;
- MSU - Moscow State University Computational Mathematics and Cybernetics Department, Russia;
- ISTC - Information Society Technologies Center, Armenia;
- GIC - V.M.Glushkov Institute of Cybernetics of the Ukrainian Academy of Sciences;
- CIT - Parallel Computing and Architectures Department, Research-Engineering Center of Informational Technologies (RECIT), National Academy of Sciences of Belarus.

Organizat în patru capitole, acest articol urmărește structura Raportului tehnic, prezentând, în primul capitol, aspecte legate de definiția termenilor *data mining* și *knowledge discovery* pentru ca, în capitolul al doilea, să se urmărească prezentarea caracteristicilor datelor geospațiale, care trebuie luate în considerare în contextul *data mining*, în general, și în contextul obiectivelor acestui proiect, în special. Următorul capitol este dedicat sistemelor software pentru *geospatial data mining* (GeoMiner, ADaM, SPIN!) relevante. Ultimul capitol prezintă câteva concluzii privind impactul asupra științelor legate de informația geografică și aspectele critice al cercetărilor în domeniu. Acest articol, ca de altfel întregul Raport Tehnic, se bazează pe o listă bogată și actualizată de referințe bibliografice.

Raportul Tehnic a fost executat ca suport pentru următoarele activități ale proiectului INTAS nr. 397:

- tehnici de distilare a datelor pentru masive de date și modelare predictivă (Task 1);

- activități dinamice în data mining (Task 4):
    - prognoze agricole pe baza factorilor meteorologici și a geomonitorizării regiunilor;
    - modelarea curgerii râurilor folosind rețele neuronale.
  - elaborarea de modele de calcul adecvate pentru implementarea clasificării deterministe (Task 7)
- și constituie contribuția ICI la task-urile menționate mai sus.

Sugerăm completarea conținutului acestui Raport Tehnic cu conținutul unui alt Raport Tehnic, efectuat tot de ICI în aprilie 2002, în cadrul activităților aceluiași proiect, dedicat "*Tehnicilor și algoritmilor specifici domeniului geospațial data mining*" [35] și cu rezultatele prezentate de drnd. mat. Laurențiu Leuștean în două articole:

- **Liquid Flow Time Series Prediction using Feed-Forward Neural Networks and Rprop Learning Algorithm**, publicat în "Studies in Informatics and Control", December 2001, pp. 287-300;
- **Liquid Flow Time Series Prediction using Feed-Forward Neural Networks and SuperSAB Learning Algorithm**, articol selectat pentru CONTI'2002 The 5<sup>th</sup> INTERNATIONAL CONFERENCE ON TECHNICAL INFORMATICS 18 - 19 October 2002, TIMISOARA, ROMANIA (în curs de publicare).

## 1. Introducere

Una dintre definițiile acceptate pe scară largă a termenilor data mining<sup>1</sup> și knowledge discovery este dată de Fayyad et al. [14]: "Data mining/knowledge discovery reprezintă procesul netrivial de identificare de șabloane de date inteligibile, valide, noi, potențial utilizabile".

Tehnologia Knowledge Discovery (KD) contribuie serios la dezvoltarea următoarei generații de sisteme informatice și sisteme de gestiune de baze de date, prin posibilitățile ei de a extrage informație nouă, înglobată în baze de date eterogene de volum mare, și de a construi cunoștințe. Un proces de knowledge discovery include "selectare de date din depozite mari de date, curățare, preprocesare, transformare și reducere, data mining, selecție (sau combinare) de model, evaluare și interpretare și consolidarea și utilizarea cunoștințelor extrase" [15]. În particular, *data mining* se ocupă cu dezvoltarea algoritmilor de extragere de noi șabloane din statistici, modelele cu rețele neuronale și vizualizare pentru clasificarea datelor și identificare de șabloane. *Knowledge discovery* are ca scop crearea mecanismelor prin care informația este transformată în cunoștințe prin intermediul ipotezelor de testare și al formalismelor teoretice.

Din definiția termenului *data mining*, pot fi observate următoarele aspecte:

1. *data mining* nu este o analiză simplă și nici nu este, în mod necesar, egală cu machine learning; nu este trivial de menționat că seturile de date considerate sunt mari; în cele mai multe cazuri, este posibilă o analiză statistică exhaustivă și este de dorit ca ea să fie cât mai riguroasă (multe metode din *data mining* conțin un grad de nedeterminare, care le permite scalarea la seturi de date masive);
2. unul dintre aspectele necunoscute la începutul procesului și care trebuie găsit, constă în faptul că *data mining* nu se aplică în cazul în care rezultatul este deja cunoscut, adică în cazul problemelor deterministe sau deductive. *data mining*, în sens larg, ar putea fi o activitate *abductivă* (numită *hypothesis* de filosoful și logicianul C.S. Peirce (1878)) care nu acoperă simultan o structură oarecare în cadrul datelor și o ipoteză de explicare a ei; aceasta va necesita structuri conceptuale sofisticate, prin care o ipoteză poate fi reprezentată într-o mașină; în cadrul *knowledge discovery*, accentul pare să se pună pe metode inductive de învățare, unde scopul este acela de a construi un model pentru intensitatea (intensitatea se referă în acest caz, mai degrabă, la descrierea generală a unei categorii decât la exemplele ei specifice) unei categorii oarecare din exemplele de instruire; deoarece structura este cunoscută în sens larg, aceasta nu se referă la acțiunea de *mining per se* ci, mai degrabă, la o formă de *knowledge discovery*; o singură excepție apare atunci când exemplele de instruire ele însele sunt doar o ipoteză, mai degrabă generată din date, decât *a priori*, într-un efort de stabilire a claselor cu care se reprezintă datele, așa cum este cazul unor instrumente ca AutoClass [10];
3. structura neacoperită trebuie să fie validă, adică trebuie arătat că poate fi o inferență semnificativă sau fezabilă cu un grad oarecare de confidențialitate; metricile de fiabilitate sunt cerute ca suport al ipotezelor prezentate și pentru a diferenția semnificativul din marginal sau irelevant;
4. constatările pot fi de domeniul noului; mașina nu are o imagine asupra a ceea ce este cunoscut sau nu de către experți, adică nu are metode de a mapa nouitatea pe domeniul discursului; prin urmare, este posibil să

<sup>1</sup> Acest termen ar putea fi tradus ca "*mineritul datelor*". Deoarece până la data publicării acestui articol nu există consens asupra termenului în limba română, pentru a evita confuziile, îl vom utiliza pe cel în limba engleză.

se postprocesează rezultatele astfel încât majoritatea inferențelor similare să fie grupate laolaltă într-o formă generalizată numită meta-instruire [4];

5. structura neacoperită trebuie să fie utilă, adică să fie explicabilă și aplicabilă într-o manieră care are sens în contextul domeniului de aplicare curent; seturile mari de date pot să conțină foarte multă structură care ea însăși nu este utilă și orientarea efortului pe aceste părți care sunt interesante este problematică deoarece este, prin definiție, necunoscută la început.

De obicei, *data mining* se referă la cazurile în care datele sunt prea mari sau extrem de complexe pentru a permite fie o analiză manuală, fie o analiză prin metode de interogare simple. *Data mining* constă din două etape principale:

- preprocesarea datelor, în timpul căreia caracteristicile relevante de nivel înalt sau atributele sunt extrase din date de nivel scăzut și
- *pattern recognition*, în care recunoașterea unui șablon de date se face prin intermediul acestor caracteristici (figura 1).

Preprocesarea datelor este, de cele mai multe, ori consumatoare de timp, prin urmare este critică. Pentru a asigura succesul în procesarea din *data mining* este important ca toate caracteristicile extrase din date să fie relevante pentru problemă și reprezentative ca date:

- în funcție de tipul datelor care trebuie "minerite", pasul de preprocesare constă în mai multe acțiuni; dacă dimensiunea datelor este foarte mare, va trebui să se recurgă la simplificarea și să se lucreze cu foarte puține instanțieri sau să se utilizeze tehnici multirezoluție și să se lucreze cu datele la o rezoluție macrogranulară. apoi zgomotul este înlăturat pe cât posibil și sunt extrase caracteristicile relevante; în unele cazuri, în care sunt disponibile date din diferite surse sau senzori, este necesară fuziunea datelor pentru a permite exploatarea tuturor datelor disponibile pentru problemă; la sfârșitul acestui prim pas, se obține un vector caracteristic pentru fiecare instanțiere; în funcție de problemă și de date, s-ar putea să fie necesară reducerea numărului de caracteristici, folosind tehnici de selecție sau de reducere a dimensiunii cum ar fi analiza componentei principale [37] sau a versiunilor sale nonlineare; după această preprocesare, datele sunt gata pentru detectarea de șabloane prin intermediul algoritmilor de clasificare, clustering, regresie etc.; aceste șabloane îi sunt afișate utilizatorului în vederea validării; *data mining* este un proces iterativ și interactiv; ieșirea din oricare dintre pași sau feedback-ul de la experți poate să conducă la o rafinare iterativă a unora sau a tuturor acestor activități.

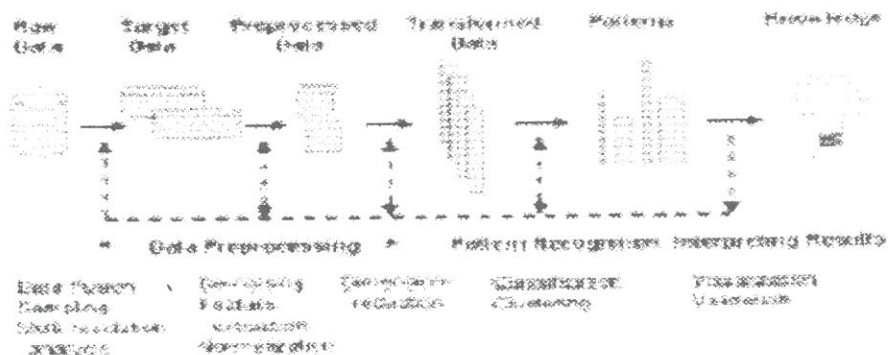


Figura 1: *Data mining* – un proces iterativ și interactiv

(adoptată de la Erick Cantu-Paz and Chandrika Kamath, "On the use of evolutionary algorithms in data mining")

Așa cum menționează Erick Cantu-Paz și Chandrika Kamath, în "On the use of evolutionary algorithms in data mining" există unele dezbateri legate de definirea termenului de *data mining*, în care majoritatea practicienilor și a celor care au propus definiții au căzut de acord că *data mining* este un domeniu multidisciplinar, împrumutând idei de la *machine learning* și inteligență artificială, statistică, calcul de înaltă performanță, procesare de imagine, optimizare, *pattern recognition* etc. Ca element de noutate, trebuie remarcată confluența ramurilor mature ale acestor tehnologii, la un moment dat, în care ele trebuie exploatate în analiza masivelor de date.

În [7] este prezentată în rezumat intersecția dintre comunitățile academice și activitățile de *knowledge discovery*. Acest rezumat (Tabelul 1) nu este unul exhaustiv, el prezentând doar câteva dintre inițiativele și direcțiile de cercetare cheie.

## 2. Cunoașterea geografică

Geografia este o disciplină integratoare: datele necesare abordării ei acoperă un spectru larg de domenii, de interese, de la aspectele sociale până la cele fizice. În continuare, sunt discutate problemele care provin dintr-un amalgam de perspective și în strânsă legătură cu o infrastructură adecvată dobândirii de informații.

• **Complexitatea asociată cu volumul datelor**

Ca toate disciplinele în care este implicat *data mining*, geografia este foarte bogată în date. Multe dintre bazele de date de tip consumator (de ex: care conțin aspecte medicale, tranzacții financiare), care sunt acum în construcție, conțin atribute spațiale și temporale, oferind, astfel, posibilitatea descoperirii sau confirmării cunoștințelor geografice [48]. Există acum seturi de date geografice de ordinul terabyte.

**Tabelul 1: Rezumatul punctelor de vedere ale diferitelor discipline asupra *data mining* și *knowledge discovery* (preluare din Buttenfield, B., Gahegan, M., Miller, H., Yuan, M.)**

	Baze de date	Statistică	Inteligență Artificială	Vizualizare
<b>Regăsire</b>	Reguli de asociere	Șabloane de analiză locală și teste de inferențiere globale	Rețele neuronale, arbori de decizie	Vizualizare exploratorie, Visual data mining
<b>Raportare</b>	Regula de listare	Semnificație și putere	Estimare probabilită, câștigare de informații	Un stimul în domeniul vizual
<b>Reprezentare</b>	Schema de actualizare, metadate	Modele de potrivire statist., locale sau globale	Grafuri conceptuale, metamodele	Partajare între scenă și observator
<b>Validare</b>	Testarea slabei semnificații	Teste de semnificație	Învățare urmată de verificare	Testarea subiecților umani
<b>Optimizare</b>	Reducerea complexității de calcul	Reducerea datelor și simplificare startificată	Căutare stocastică, metode de gradient	Metode ierarhice și adaptive

• **Complexități asociate domeniului**

Semnale interesante și relevante sunt, adesea, ascunse în întregime de șabloane puternice, care trebuie mai întâi să fie înlăturate. Multe dintre aspectele complexe își au originea în codependența spațială și temporală, care apare între diferitele varietăți de scări și din cauze extrem de diverse [57].

• **Complexități datorate variației locale**

Sistemele legate de Pământ sunt conectate intrinsec, ceea ce face dificilă analiza izolată, pe o parte, a unui sistem de la care să se modeleze alte aspecte. Rezultatul va apare în formă statistică.

• **Complexități cauzate de colectare și eșantionare**

Cu toate că datele sunt disponibile în a-și crește volumul, cazul cel mai frecvent este acela în care trebuie făcută mai degrabă o (re)sortare pentru a surprinde fenomenele de interes, decât să se treacă direct la măsurători.

• **Dificultăți în formalizarea domeniului geografic**

Una dintre cele mai mari dificultăți în ceea ce privește activitatea de *knowledge discovery* în cadrul domeniului geografic este însăși complexitatea domeniului. Există, deși neacceptat pe scară largă, un model conceptual al geografiei [27], iar modelele în GIS-urile comerciale variază, adesea, în mod fundamental, ceea ce conduce la trei probleme distincte:

- datele sunt, de cele mai multe ori, necomensurate ceea ce face ca ele să nu poată fi comparate sau combinate direct;
- este dificil de aplicat cunoașterea geografică formală procesului de *knowledge discovery*, deoarece astfel de cunoștințe nu sunt ușor disponibile;
- când noua cunoștință nu poate fi “acoperită”, este foarte dificil de reprezentat ceea ce vrea să spună ea.

Alte puncte de vedere, legate de aceste probleme, sunt prezentate de Yuan *et al.* [63] și Miller and Jiawei, [48]. Puse laolaltă, aceste puncte de vedere conduc la necesitatea ca activitățile de descoperire să fie situate în jurul expertului în domeniu: geograful care este, cel puțin în prezent, cea mai bună sursă de înțelegere și cunoaștere a domeniului.

În rezumat, exemplele de date geografice includ date ce descriu evoluția fenomenelor naturale în timp, date specifice științelor pământului, care descriu fenomene spațiotemporale în atmosferă, secvențe de imagini bi- și tridimensionale ale regiunilor geografice sau date care descriu localizarea indivizilor în spațiul geografic, în funcție de timp. De obicei, datele geospațiale conțin informații geometrice sau topologice, care se adaugă informațiilor despre factorii de mediu sau datelor legate de explorarea resurselor naturale [16], [40]. Astfel de date reprezintă cea mai mare parte din datele care au fost colectate de diverse organizații. Sistemele de

observare, cum sunt sateliții sau radarele, pot genera seturi de date foarte mari. De ex., se estimează că Sistemul de Observare a Pământului al NASA va transmite 50 GB pe oră.

Datele geospațiale (statice sau care implică evoluția în timp) se deosebesc, în mod esențial, de datele *on-line analytical processing* (OLAP), datele multidimensionale generale, datele relaționale clasice sau datele tranzacționale. Aplicând normele spațiului euclidian, datele geografice dau distanța informațională și topologică [30]. De obicei, domeniul atributelor este real. Seturile de date geospațiale pot să înregistreze informații suplimentare ca de ex., parametrii mediului, definind astfel spațiul dimensional al atributelor care pot fi corelate. Suplimentar, datele spațiotemporale pot să reprezinte schimbările intervenite în timp în geometria obiectelor. Activitățile specifice *data mining*, care pot fi aplicate datelor geospațiale, depind de date și de problemă. Problema generală poate fi formulată prin găsirea și descrierea structurii datelor care mai înainte nu au fost cunoscute, și nu este memorată explicit în baza de date [40].

Geographical Information Science a devenit un organism de cunoaștere, construit pe baza fundamentelor din domenii ca: vizualizare geografică (*geographic visualisation* = GVis), *geographic information systems* (GIS) și *knowledge discovery in databases* (KDD). Eforturile depuse în ultimii zece ani atât în cercetarea fundamentală, cât și în cea aplicativă, specifică acestor domenii, au căutat răspunsuri și la sporirea disponibilității datelor și la implicarea rapidă a tehnologiilor computerizate. În particular, în contextul dezvoltării de metode inovatoare și producerii de instrumente asociate, s-a căutat facilitarea de decizii politice, evaluări și planificări apelând, în acest fel, la o audiență mai mare din partea utilizatorilor de informație geografică.

Ținta principală în cercetarea legată de GVis a fost rolul instrumentelor de vizualizare interactivă care facilitează identificarea și interpretarea de șabloane și relații în cadrul datelor complexe [42]. GVis s-au construit pe o bază cartografică, cu GIS, analiza datelor spațiale și a imaginilor, eforturi care se corelează cu eforturile științifice, legate de vizualizarea informației, și cu eforturile exploratorii din statistică, legate de analiza datelor. Dezvoltarea KDD coincide cu creșterea exponențială a volumului de date generate de și disponibile pentru știință, guvernare și industrie, date generate, în particular, în format digital.

Termenul de "*knowledge discovery in databases*" a apărut în 1989, în contextul eforturilor de a face distincție între aplicarea de algoritmi specifici pentru extragerea de șabloane din date (subproces al *data mining*) și întregul proces, în care *data mining* este un pas în extragerea de cunoștințe din aceste șabloane [14]. Din literatura consultată, au reieșit numeroase metode KDD, care diferă în ceea ce privește dezvoltarea conceptuală, reflectând dezvoltarea separată din cadrul unor domenii ca sistemele de baze de date, machine learning, statistică și inteligență artificială [3], [11].

În ultimii zece ani, au fost acumulate mari cantități de date în GIS, cu scopul amplu de analiză exploratorie folosind metode GVis și KDD și instrumente asociate. Multe dintre datele legate de mediu, generate recent, așa cum s-a întâmplat cu cele de la Earth Observation System - EOS, cu sistemele de monitorizare, stațiile meteorologice, includ georeferențiere. Aspectele spațiale ale acestor date sunt, de fapt, unul dintre scopurile principale ale analizei – pentru studierea dispersiei poluanților, fragmentarea din silvicultură, monitorizarea amenajării teritoriului și alte aplicații. Observațiile repetate s-au dovedit critice, în ceea ce privește formularea de răspunsuri la cele mai importante întrebări ale științei mediului – cele legate de procesele specifice. Datele georeferențiate de mediu au, de obicei, componente spațiale și temporale.

În rezumat, comunitățile academice au căzut de acord că dezvoltarea tehnologiilor specifice *data mining* (DM) și *knowledge discovery* (KD) au deschis drumuri noi în cercetare, în general, și în cadrul cercetării informației geospațiale, în particular. Abilitatea de a "mineri" datele presupune că există mecanisme de furnizare și de acces la date. Chiar dacă serviciile de furnizare sunt pe cale să devină disponibile în medii locale și distribuite, rămân o mușumbe de probleme nerezolvate. O parte a activității academice trebuie să se dedice infrastructurii suport pentru *data mining* și *knowledge discovery*. Mecanismele existente nu sunt proiectate să manipuleze probleme specifice ale informației geospațiale.

În conformitate cu [7], trei caracteristici ale datelor geospațiale crează provocări în ceea ce privește dezvoltarea unui fundament robust de date. Dezvoltarea unei infrastructuri de date, necesară ca suport pentru GIScience, constituie o țintă a unei alte inițiative a University Consortium on Geographic Information Science (UCGIS) (spatial data infrastructure). Accentul nu este pus aici pe infrastructura de date spațiale în sine, ci pe dezvoltarea *data mining* în cadrul infrastructurii.

#### a) Depozitele de date geospațiale tind să fie foarte mari

Așa cum am menționat anterior, volumul datelor a fost un factor important în tranziția mai multor agenții de la furnizarea de date publice prin intermediul mecanismelor fizice (CD ROM, de ex.) la mecanisme electronice [50]. Mai mult, seturile de date GIS existente sunt, de cele mai multe ori, sparte în componente de tip caracteristică și de tip atribut, care sunt arhivate în mod convențional în sisteme hibride de management de date.

Cerințele algoritmice diferă substanțial de managementul de date relaționale (atribute) și pentru management de date topologice (caracteristice) [32]. Procedurile de calcul ale KD pot fi diversificate, dacă sunt pe cale să devină complet operaționale într-un mediu geospațial de calcul. Chiar în cazul apariției și dezvoltării de noi modele de date pentru GIS integrat, modelul de date hibrid (caracteristică/atribut) va fi păstrat. În practică, integrarea de cunoștințe va începe să acopere nu numai modele disparate de date dintr-o singură arhivă, ci și arhive disparate în sisteme disparate de management de date. În legătură cu acestea se află și gradul și diversitatea formatelor de date geografice, care prezintă provocări unice. Revoluția datelor geografice digitale a creat noi tipuri de formate de date printre care se află și cele tradiționale "vector" și "raster". Depozitele de date geografice includ date prost structurate cum sunt imagini și *geo-referenced multi-media* [8]. Descoperirea de cunoștințe geografice din date georeferențiate multimedia este un aspect mult mai complex al problemei KD pentru date multimedia [64].

#### **b) A doua caracteristică a datelor geospațiale se leagă de introducerea caracteristicilor datelor colectate ciclic**

Descoperirea de date trebuie să se acomodeze ciclurilor de colectare, care pot fi necunoscute (ca de ex: identificarea ciclurilor de schimbare în dislocările geologice majore) sau care se pot schimba de la un ciclu la altul atât în timp, cât și în spațiu (de ex: șabloane de dispersie a epidemiilor sau a substanțelor toxice). Integrarea informației din modele multiple de date este recunoscută ca fiind o țintă a unuia din temele de cercetare ale UCGIS (Achiziția și integrarea de date spațiale), și nu ne propunem să o dezvoltăm în acest material. În schimb, cercetătorii din domeniul KD pot să acorde atenție problemelor de raționament și modelare pe cicluri mai lungi sau mai scurte de timp. De exemplu, *geospatial knowledge discovery* (GKD) ar putea să suporte traseele unei furtuni în timp real sau predicția unei avalanșe sau alte evenimente localizate, legate de evoluția vremii. Aspectele legate de infrastructură, care urmează să fie cercetate includ, de ex., dezvoltarea de *real-time data mining* și utilizarea de instrumente KD pentru ghidarea corelării de șabloane de date descoperite în timp, validarea tendințelor datelor dealungul discontinuităților temporale ș.a. Deoarece înțelegem mai puțin despre natura timpului decât a spațiului, metodologiile de arhivare de date, în vederea facilitării de căutări spațiale ciclice, rămân nefinisate. Extensia la care se pot identifica șabloane de date va fi determinată parțial sau în întregime prin organizarea datelor într-o arhivă. Cercetarea privind cea mai bună structură sau pentru reordonarea datelor pentru activitățile specifice de KD nu este, încă, demarată în alte teme de cercetare ale UCGIS.

#### **c) A treia caracteristică de relevanță a DM/KD se aplică mai degrabă la o caracteristică a datelor fundamentale decât asupra datelor**

Peste trei milioane de Website-uri sunt online. În acest caz, cele mai bune motoare de căutare pot să localizeze cel mult o treime din paginile accesibile [49], [50]. Astfel de surse de date include, dar nu se limitează la site-uri publice de date de domeniu colectate în țările dezvoltate, site-uri de comunități de date extrem de localizate cum ar fi site-uri prezentând vecinătățile unor orașe sau ale comunității active dintr-un teritoriu și surse similare de date necunoscute sau cunoscute prin convenții în cadrul infrastructurii de date geospațiale. Acest tip de KD se ocupă de întregul Internet ca un depozit foarte mare, descentralizat de date și furnizează contribuții la o infrastructură globală. Este paradoxal faptul că o cantitate din ce în ce mai mare de date devine disponibilă prin Internet, în timp ce datele sunt din ce în ce mai greu de localizat, regăsit și analizat. Aceasta se întâmplă și datorită lipsei din Internet a unui catalog inteligibil sau index [6]. Fără o infrastructură coordonatoare, multe dintre sursele de date și serviciile disponibile astăzi rămân în esență inaccesibile.

Datele geografice au și alte proprietăți unice, care reclamă atenție și tehnici speciale, existente în domeniul măsurătorilor geografice. În timp ce alte aplicații KD implică spații superior dimensionate, datele geografice sunt unice întrucât până la patru dimensiuni ale spațiului informațional sunt interrelaționate și furnizează contextul de măsurare pentru dimensiunile rămase. Cadrul de măsurare cel mai des adoptat este topologia și geometria asociată cu spațiul euclidian. Cu toate acestea, unele fenomene geografice au proprietăți neeuclidiene. Exemple de acest fel sunt timpii de călătorie în zonele urbane, imaginile mentale ale spațiului geografic și propagarea bolilor în spațiu și timp [45]. Proiecția datelor spațiale, în contexte de măsurare potrivite, poate să ajute căutarea de șabloane în *geographic data mining*. Informația inerentă contextului de măsurare geografic este, adesea, ignorată de instrumentele de inducție și machine learning [18].

Adesea, atributele geografice măsurate arată proprietățile dependenței și eterogenității spațiale. Primele se referă la tendința atributelor de a fi legate de unele locații din spațiu (de obicei, acestea sunt vecinătățile). Ultimele se referă la nestaționaritatea majorității proceselor geografice, înțelegându-se că parametrii globali nu reflectă foarte bine procesul care apare la o locație particulară. În timp ce aceste proprietăți au fost tratate ca neplăceri, cercetarea contemporană, ajutată de progresele din domeniul tehnologiei informației geografice, furnizează instrumente care pot să explice aceste proprietăți pentru înțelegerea fenomenelor geografice [1], [5], [17], [25]. Unele cercetări în GKD sugerează că, ignorând aceste proprietăți, sunt afectate șabloanele derivate din tehnicile *data mining*. Este necesar un efort de cercetare mai mare, în ceea ce privește tehnicile scalabile de capturare a dependenței și eterogenității spațiale în GKD.

Un al treilea aspect unic al informației geografice în KD este complexitatea obiectelor și șabloanelor geospațiale. În majoritatea domeniilor negeografice, obiectele de date pot fi foarte bine reprezentate discret în cadrul spațiului informațional fără a pierde proprietăți importante. Nu este cazul obiectelor geografice: dimensiunea, forma și granițele pot să afecteze procesele geografice, înțelegând că obiectele spațiale nu pot fi

reduse la puncte sau caracteristici lineare simple, fără a se pierde informații. Relații cum sunt distanța, direcția și conectivitatea sunt mult mai complexe cu obiectele dimensionale [12], [51], [55]. Transformările în timp, între aceste obiecte, sunt complexe și purtătoare de informații [34]. Scara și granularitatea, în ceea ce privește măsurarea timpului, pot fi complexe, făcând ca, o simplă actualizare de dimensiune a spațiului, să includă timpul [35], [57]. Dezvoltarea de instrumente scalabile pentru extragerea de șabloane din colecții de obiecte spațiotemporale diferite este una dintre provocările majore. Astfel, deși complexitatea șabloanelor și a regulilor spațiotemporale poate fi descurajantă, capătă sens o altă provocare derivată probabil prin intermediul "meta-mining" [37].

### 3. Sisteme software pentru domeniul Geospatial Data Mining/Knowledge Discovery

În "A Data Miner's Tools", din BYTE/October 1995, Karen Watterson [60] explică trei categorii de software dedicate domeniului *data mining*. *Instrumentele interogare-și-raportare*, în forma lor simplificată și ușor de utilizat, cer asistare umană și legături în bazele de date sau alte formate speciale. *Instrumentele de analiză multidimensională* (multidimensional analysis = MDA) cer mai puțină intervenție umană, dar au nevoie de date în formate speciale. *Agenții inteligenți* sunt virtual autonomi, sunt capabili să-și facă propriile observații și să tragă concluzii pe care pot să le manipuleze ca forme libere în paragrafe de text.

*Data Mining* așa cum este definit în primul capitol, folosește *machine learning*, tehnici statistice și de vizualizare pentru a descoperi și prezenta cunoștințe într-o formă ușor inteligibilă de către factorii umani.

Explozia în datele georeferențiate, prilejuită de dezvoltarea tehnologiei informației (Information Technology = IT), cartografiei digitale, teledetecției și răspândirii pe scară largă a GIS, accentuează importanța dezvoltării metodelor inductive, conduse de date în analiză geografică și modelare pentru a facilita crearea de noi cunoștințe și pentru a ajuta procesul descoperirii științifice. În prezent, există un număr de instrumente pentru *data mining*, care sunt proiectate să asiste procesul de explorare a unor cantități mari de date în căutarea de șabloane recurente și relații (de ex: STATISTICA, Clementine, MineSet, Intelligent Data Miner etc.). Aceste pachete au fost dezvoltate, în principal, în scopul de a analiza baze de date comerciale foarte mari, pentru a modela și a face predicții în ceea ce privește comportamentul consumatorilor. Accentuarea pe predicție poate să le limiteze utilitatea în ceea ce privește GIS, unde *spatial pattern recognition* poate fi o activitate mult mai utilă. În contextul GIS, există trei argumente care pot fi invocate în favoarea dezvoltării și aplicațiilor instrumentelor pentru *data mining* pentru a exploata invazia geoinformațiilor. Mulți dintre specialiștii implicați în *data mining* cred că instrumentele lor vor lucra pe oricare și pe toate datele indiferent de subiectul de origine. *Geographical Data Mining* trebuie privit ca un tip special de *data mining* (așa cum este prezentat și în capitolul al doilea), care știe să proceseze funcții generice similare ca și instrumentele convenționale de *data mining*, dar modificate sau construite special pentru a ține cont de caracteristicile geoinformațiilor, de stilurile diferite și de cerințe de analiză și modelare, relevante pentru lumea GIS, dar și de natura specifică a explorării geografice.

#### 3.1. GeoMiner: Un sistem pentru descoperirea de cunoștințe pentru baze de date spațiale și Geographic Information Systems

##### Prezentare generală

GeoMiner este un sistem pentru descoperire de cunoștințe pentru baze de date spațiale, dezvoltat în Database Systems Research Laboratory, School of Computing Science, Simon Fraser University. Această secțiune se bazează pe prezentarea publică de la adresa: <http://db.cs.sfu.ca/DBMiner>.

*Spatial data mining* constă în a "mineri" informația spațială de nivel superior și cunoștințele din baze de date spațiale mari. Un prototip de sistem pentru *spatial data mining*, GeoMiner, a fost proiectat și dezvoltat pe baza experienței în cercetarea și dezvoltarea unui sistem relational *data mining*, DBMiner, și pe baza cercetărilor din domeniul *spatial data mining*. Puterea lui GeoMiner în *data mining* include "mineritul" a trei tipuri de reguli: reguli caracteristice, reguli de comparare și reguli de asociere în baze de date geospațiale cu o extensie planificată pentru a include "mineritul" regulilor de clasificare și clustering. GeoMiner include modulul de construcție *spatial data cube*, modulul *spatial on-line analytical processing (OLAP)*, precum și module de *spatial data mining*. De asemenea, a fost proiectat și implementat un limbaj de *spatial data mining language*, GMQL (Geo-Mining Query Language), ca o extensie a *Spatial SQL*, pentru *spatial data mining*. A fost construită și o interfață interactivă și prietenoasă și au fost implementate instrumente de vizualizare a cunoștințelor spațiale descoperite.

##### Descrierea proiectului

GeoMiner permite găsirea de cunoștințe interesante din baze de date spațiale mari. Metodele de *spatial data mining* au fost aplicate pentru a extrage astfel de cunoștințe din baze de date spațiale mari. O dată cu progresul făcut de cercetare în *data mining* și *data warehousing* din ultimii ani, au fost dezvoltate multe sisteme de *data mining* și *data warehousing* pentru "mineritul" de cunoștințe în baze de date relaționale și *data warehouses*.

*Spatial data mining* este un subdomeniu care se ocupă cu extragerea de cunoștințe implicite, relații spațiale sau de alte șabloane interesante, memorate neexplicit în baze de date spațiale. Cu o mare cantitate de date spațiale, colectate de sistemele satelitare, sistemele de teledetecție, sistemele de vânzări regionale și de alte instrumente de colectare de date, este inevitabilă dezvoltarea de instrumente pentru descoperirea de cunoștințe interesante din bazele de date spațiale mari. În plus, multe dintre bazele de date relaționale conțin și ele informații spațiale cum sunt adrese ale reședințelor clienților sau localizarea unui magazin/depozit. Aceste adrese pot fi geocodificate prin extragerea coordonatelor spațiale, care pot fi memorate într-o bază de date spațială împreună cu alte date spațiale. Este important să se "mineze" cunoștințe legate atât de obiecte spațiale, cât și de obiecte nespațiale în bazele mari de date. În studiile făcute până acum, din păcate, nu sunt cunoscute multe sisteme de *spatial data mining*.

Progresele recente în ceea ce privește structurile datelor spațiale și ale bazelor de date spațiale fac posibilă crearea de baze de date spațiale mari, care pot fi efectiv interogate. Aceste progrese în combinație cu cercetările din cadrul raționamentului spațial, precum și cu cele din domeniul *data mining* în baze de date relaționale, promovează cercetările din cadrul *spatial data mining*.

Grupul de cercetare GeoMiner din cadrul Intelligent Database Systems Research Laboratory a lucrat mulți ani în domeniul *data mining* și, în special, în *spatial data mining*. Sistemul GeoMiner include modulul de construcție *spatial data cube*, modulul *spatial on-line analytical processing* (OLAP), precum și module de *spatial data mining* pentru "mineritul" de reguli caracteristice, reguli de comparare, reguli de clasificare, reguli de asociere și clustering.

Arhitectura SAND (Spatial And Nonspatial Data) se aplică la modelarea bazelor de date spațiale, iar modulele de *spatial data mining* includ atât "mineritul" cunoștințelor spațiale, cât și relațiile între componentele spațiale și nespațiale. De asemenea, a fost proiectat și implementat un limbaj de *spatial data mining*, GMQL (Geo-Mining Query Language), ca o extensie a Spatial SQL. A fost construită și o interfață interactivă și prietenoasă și au fost implementate instrumente de vizualizare a cunoștințelor spațiale descoperite.

#### Arhitectura

Sistemul GeoMiner este o extensie și o evoluție de la un sistem relațional pentru *data mining* DBMiner, cercetat și dezvoltat de același laborator. Sistemul DBMiner (a se vedea referințele și sistemul prin intermediul <http://db.cs.sfu.ca/DBMiner>) conține următoarele cinci module funcționale: caracterizator, comparator, asociator, predictor și clasificator. Numeroase alte module se află, încă, în faza de cercetare și dezvoltare. DBMiner este implementat prin integrarea tehnicilor de *data mining* și *data warehousing*, incluzând *data cube construction and manipulation*, *attribute-oriented induction*, *multi-level association analysis*, *statistical data analysis*, *machine learning* etc. pentru "mineritul" datelor relaționale.

GeoMiner este construit pe baza sistemului DBMiner. Funcțiile pentru operațiile OLAP nespațiale și pentru "mineritul" datelor nespațiale sunt direcționate către sistemul DBMiner, în timp ce funcțiile pentru "mineritul" datelor spațiale și pentru relațiile între datele spațiale și nespațiale sunt procesate de funcții dedicate ale GeoMiner.

Principalele funcții ale sistemului includ "mineritul" a cinci tipuri de reguli de cunoștințe în baze de date spațiale, integrarea de tehnologii pentru *data mining* și *data warehousing*, "mineritul" interactiv de reguli multinivel, integrarea cu baze de date relaționale comerciale și GIS, precum și multe forme de ieșire, incluzând hărți generalizate, relații generalizate, reguli multinivel etc.

Arhitectura generală a GeoMiner constă din:

1. o interfață grafică cu utilizatorul pentru "mineritul" interactiv și afișare de rezultate ale *data mining* sub formă de tabele, hărți etc.;
2. un set de module de descoperire, incluzând cele cinci module existente: geo-caracterizatorul, geo-comparatorul, geo-clasificatorul, analizorul de geo-cluster și geo-asociatorul, precum și alte două module în dezvoltare: geo-predictorul și analizorul de geo-pattern;
3. un server de baze de date spațiale, care include MapInfo Professional 4.1;
4. *data cube mining engine* care se bazează pe nucleul de descoperire DBMiner pentru manipulare de date multidimensionale;
5. baza de date și baza de cunoștințe care memorează datele spațiale și nespațiale și
6. ierarhiile conceptuale.



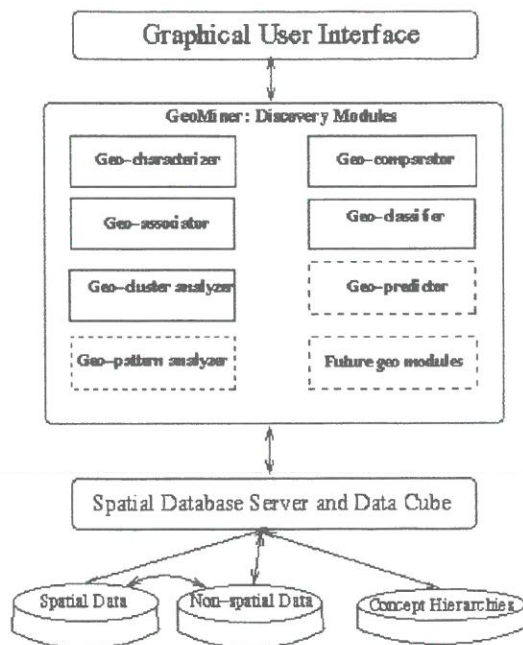


Figura 2: Arhitectura GeoMiner (preluată de la <http://db.cs.sfu.ca/DBMiner>.)

## Modulele funcționale importante

### Geo-caracterizatorul

Acest modul “minerește” un set de reguli caracteristice la niveluri multiple de abstractizare, de la un set de date relevant dintr-o bază de date spațiale. El furnizează utilizatorilor un punct de vedere multinivel, multiunghiular asupra datelor dintr-o bază de date spațială.

### Geo-comparatorul

Acest modul “minerește” un set de reguli de comparare, care sunt în contrast cu caracteristicile generale ale diferitelor clase ale diferitelor seturi de date, dintr-o bază de date. Se compară un set de date, cunoscut ca și clasă țintă, cu unul sau mai multe seturi de date, cunoscute ca și clase contrastante.

### Geo-asociatorul

Acest modul găsește un set de reguli de asociere spațial relaționate din unul sau mai multe seturi de date dintr-o bază relațională. O regulă de asociere arată apariția curentă de șabloane (sau relații) dintr-un set de articole de date dintr-o bază de date. O regulă de asociere spațială tipică este de forma “ $A \rightarrow B (s\%, c\%)$ ” unde A și B sunt seturi de predicate spațiale sau nonspațiale, s este suportul regulii (probabilitatea ca A și B să fie păstrate amândouă în toate cazurile posibile), iar c este gradul de confidențialitate al regulii (probabilitatea condițională ca B să fie adevărat sub condiția A).

### Analizorul de Geo-cluster

Acest modul găsește clustere de puncte cu descrieri nespațiale relevante. El folosește un algoritm eficient, CLARANS, pentru a procesa spatial clustering.

Prin folosirea inducției orientate-atribut se găsesc descrierile nespațiale ale clusterelor.

### Geo-classicatorul

Acest modul adoptă o metodă de inducție generalizată, bazată pe arbori de decizie, pentru a construi un arbore de decizie care clasifică setul datelor relevante în conformitate cu unul dintre atributele nespațiale. Se afișează arborele de clasificare și, prin poziționarea pe unul dintre noduri, utilizatorul poate să pună în evidență regiunile de hartă corespunzătoare.

## Evoluția GeoMiner

Echipa de la Database Systems Research Laboratory, School of Computing Science, Simon Fraser University își propune să continue cercetarea și dezvoltarea în următoarele direcții:

- sporirea puterii și eficienței mecanismelor de descoperire, incluzând îmbunătățirea calității regulilor și performanței sistemului pentru modulele funcționale existente;
- nu cu mult timp în urmă, echipa s-a concentrat pe descoperirea de cunoștințe doar dintr-o singură hartă tematică; intenția sa este de a aborda algoritmi eficienți, care să permită manipularea de hărți tematice multiple;
- versiunea actuală de clasificator construiește arborele de decizie, bazându-se pe proprietăți nespățiale ale obiectelor, iar, la sfârșit, este procesată vizualizarea spațială datelor; echipa și-a propus să dezvolte noi algoritmi care pot să țină cont și de proprietățile spațiale ale obiectelor clasificate;
- în plus față de modulele menționate anterior, echipa va implementa un geo-predictor; câteva dintre aspectele raționale ale predictorului au fost deja implementate în sistemul DBMiner; de asemenea, urmează să fie încorporată și o parte de analiză de șablon în cadrul sistemului GeoMiner;
- este planificată și proiectarea și îmbunătățirea interfețelor de nivel înalt, prietenoase pentru *interactiv spatial data mining* și *visual presentation* pentru cunoștințele descoperite.

### 3.2. Data Mining System Toolkit for Earth Science Data

Information Technology and Systems Center (ITSC) de la Universitatea din Huntsville, Alabama, a dezvoltat un sistem pentru *data mining*, care permite cercetătorilor și altor utilizatori să culeagă informații și cunoștințe din munții de date din Științele Pământului. Seturile de date din cadrul Științelor Pământului și, în general, datele spațiale sunt variabile ca formate, scări și rezoluție. Prin cercetările desfășurate în domeniile *data mining*, *coincidence data searching*, *object - relational databases*, *Open GIS*, *data integration* și *interoperable distributed data systems*, ITSC a produs o varietate de sisteme software, care ajută la a face seturile de date universal accesibile, accesibile și utile. Prin angajarea construcțiilor orientate-obiect și a instrumentelor de dezvoltare ITSC s-a creat situația de a utiliza structura sistemului pentru a furniza și alte funcționalități cum ar fi subsetare de date generice și generarea de date la cerere.

Această secțiune se bazează pe articolul: “*Data Mining System Toolkit for Earth Science Data*”, scris de Ken Keiser, John Rushing, Helen Conover, Sara Graves, care este disponibil la adresa: [http://webtech.ceos.org/eogeo99/Papaers/Keiser/Adam\\_EOGEO.html](http://webtech.ceos.org/eogeo99/Papaers/Keiser/Adam_EOGEO.html).

Sistemul Algorithm Development and Mining (ADaM) dezvoltat la ITSC a fost instrumentul în detectarea fenomenelor, extragerea de caracteristici și furnizarea de instrumente pentru analiză de date și procesare, care integrează variabile spațiale și temporale din seturile de date specifice Științelor Pământului. Acest sistem constă dintr-un cadru de bază pentru interschimbul de date și un set de module *plug in*, care lucrează în acest context. Modulele includ filtre de intrare, module de analiză și filtre de ieșire. Filtrele de intrare/ieșire translatează seturile de date specifice într-o zonă de tipuri de date și formate la și de la o reprezentare internă puternică și flexibilă. În cadrul sistemului de mining, datele sunt administrate într-o structură standard, care permite aliniere spațială și temporală pentru integrare și analiză. Modulele de analiză încărcate dinamic, dezvoltate și personalizate de către cercetători, manipulează structura internă, sunt reutilizabile, și nu sunt constrânse de diferențele din formatele externe. Sistemul ADaM are multe module de *data mining*, *pattern recognition*, *image processing*, *subsetting*, *gridding*. Acest sistem are, de asemenea, filtre pentru o varietate de formate specifice, precum și filtre noi și pot fi adăugate cu ușurință module de analiză. Cercetătorii de la ITSC împreună cu alți colaboratori utilizează sistemul pentru *data mining* în studiile despre clasificarea de textură, procesare de imagini și analiză statistică pentru aplicații cum sunt managementul în silvicultură și modelarea atmosferică.

#### Contextul de lucru ADaM

Inițial, ITSC a dezvoltat sistemul ADaM în cadrul unui grant de la NASA pentru a investiga noi metode de procesare a volumelor mari de date de la Earth Observing System (EOS) și teledetecție. Scopul grantului a fost acela de a face căutări pe valori de date, precum și pe metadate, și de a îmbunătăți metadatele limitate, disponibile, de obicei, pentru Științele Pământului prin catalogarea informațiilor bazate pe conținutul datelor. Ca subobiectiv, s-a urmărit adăugarea de algoritmi care să poată suporta detectarea unei varietăți de fenomene geofizice în datele “minerite”. ITSC a dezvoltat și procesele de generare și de memorare de metadate bazate pe conținut, care pot fi regăsite de cercetători prin intermediul interfețelor bazate pe web, care îi conduc la seturile actuale de date cerute în numeroase analize și studii [33]. Acest sistem pentru *data mining* a fost utilizat și pentru alte studii care se refereau la clasificare de textură, procesare de imagini și analiză statistică.

## Mediul ADaM

Motorul ADaM este proiectat pentru a extrage conținut bazat pe metadata din arhivele dedicate Științelor Pământului [33]. Acesta poate detecta fenomene sau evenimente care sunt de interes pentru oamenii de știință și poate memora informația într-un mod care facilitează procesul de căutare și de ordonare a datelor. Unele rezultate de acest tip sunt memorate în Eureka, un motor de căutare de date spațiale, coincident folosit pentru a găsi coincidențe între fenomenele *mining-generated*, evenimentele climatologice și informații statistice, cum ar fi granițele de județ sau ale bazinelor hidrologice. Acest motor furnizează și alte capacități de ordonare ca de ex: generarea de date prin intermediul aplicațiilor client. Procesarea client poate să includă *gridding*, *resampling*, filtrare, conversie de format sau alte analize depinzând de cerințele clienților. De ex., ADaM poate să genereze lunar o imagine a totalului pluviometric de la datele provenind de la radar. Atât instrumentele de interogare spațială, cât și aplicațiile client sunt aplicații web și, în felul acesta, clienții sunt capabili să ruleze aproape în orice mediu. Figura 3 prezintă o arhitectură generală a mediului pentru *data mining*, ADaM.

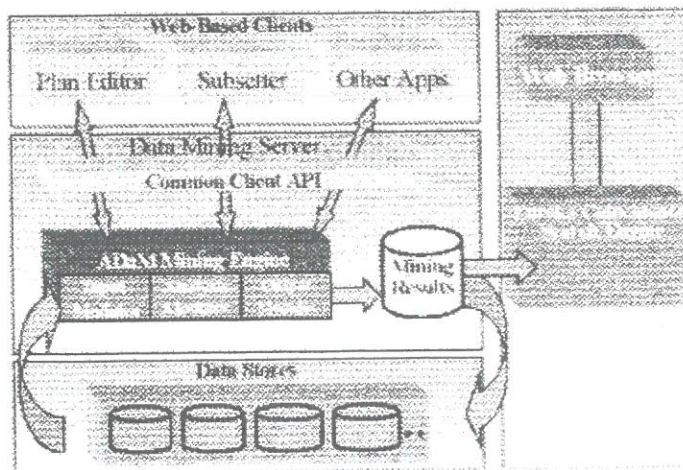


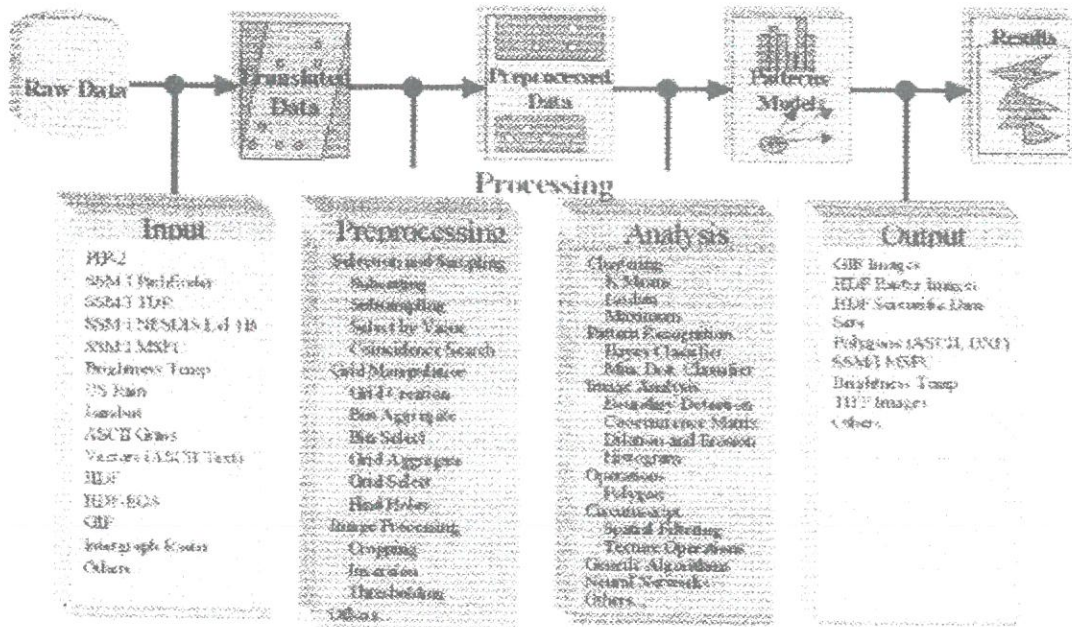
Figura 3: Mediul ADaM pentru data mining (preluată din prezentarea "*Data Mining System Toolkit for Earth Science Data*", făcută de Ken Keiser, John Rushing, Helen Conover, Sara Graves și disponibilă la adresa: [http://webtech.ceos.org/eogeo99/Papers/Keiser/Adam\\_EOGEO.html](http://webtech.ceos.org/eogeo99/Papers/Keiser/Adam_EOGEO.html))

### Arhitectura de procesare ADaM

Specificul datelor spațiale constă în faptul că tind să înțeleagă mai multe forme și dimensiuni în termeni de formate, scări și rezoluții. Acestea s-au dovedit a fi probleme consumatoare de timp atunci când cercetătorii se așteptau să integreze seturi noi de date în studiile și modelele lor. În ideea depășirii acestor probleme, ADaM a fost proiectat pentru a manipula intern toate datele într-un mediu comun fără a pierde formatele de intrare și de ieșire actuale. Această metodă izolează efectiv modulele de analiză internă de setul de date specific. Se întâmplă, în mod frecvent, ca să fie necesari mai mulți pași distincți pentru îndeplinirea unei activități. Este adeseori de dorit ca aceasta să se producă fără scrierea de rezultate intermediare pe disc. Pentru atingerea acestui scop, proiectanții lui ADaM au adoptat metoda *data pipeline*. Rezultatele provenite dintr-o operație de procesare constituie intrare pentru următoarea operație. Mai multe operații pot fi înlănțuite într-un singur plan de procesare. Figura 4 ilustrează arhitectura funcțională a motorului ADaM.

### 3.3. Integrarea GVis, GIS și KDD pentru explorarea datelor spațio-temporale

Această secțiune prezintă o metodă de abordare a unui proces pe trei niveluri ca o strategie pe termen lung pentru a integra metode GVis și KDD cu un GIS temporal. Scopul principal este acela de a dezvolta tehnici inovatoare și instrumente asociate pentru explorarea de date (care includ informație cantitativă și calitativă) pentru, a ajunge la șabloane spațiotemporale valide, noi, potențial utile și inteligibile și pentru a procesa date. Rezultatul este un proiect de vizualizare bazată pe cunoștințe pentru GIS temporal, care implică procesul de percepție cognitivă și procesarea automată de informații cu multe decizii luate de utilizatorii de informație geografică (GI) în ceea ce privește modul în care se potrivesc modelele sau modul în care se determină șabloanele spațiotemporale din date.



**Figure 4: Fluxul de date în cadrul AdaM (preluată de la “Data Mining System Toolkit for Earth Science Data”, autori: Ken Keiser, John Rushing, Helen Conover, Sara Graves [http://webtech.ceos.org/eogeo99/Papaers/Keiser/Adam\\_EOGEO.html](http://webtech.ceos.org/eogeo99/Papaers/Keiser/Adam_EOGEO.html))**

Bazându-se pe articolul scris de Monica Wachowicz, “Integrating GVis, GIS and KDD for Exploring Spatio-Temporal Data”, această secțiune prezintă o discuție preliminară asupra modului în care integrarea GVis, GIS și KDD poate naște noi provocări rezultând din complexitatea procesului de explorare de date și capacitățile suplimentare ale sistemului integrat țintă. Din punct de vedere sistemic, rezultatele sunt legate de suportul efectiv al interacțiunilor utilizator-date atât în modele de date, cât și în interfețele de nivel înalt. Din punctul de vedere al utilizatorilor de informație geografică, principalul rezultat constă în a face procesul de construire de cunoștințe cât mai flexibil și de a facilita explorarea interactivă a datelor spațiotemporale multivariate. În secțiunile următoare, sunt trecute în revistă principiile și dezvoltările cheie din ultimi zece ani din domeniile GVis, GIS și KDD. Această trecere în revistă furnizează o bază de la care se poate explora integrarea potențială a GVis, GIS și KDD.

Cunoștințele câștigate din cercetările desfășurate în vederea integrării GVis, GIS și KDD vor fi fructificate în proiectarea viitoarei generații de GIS, care va furniza un singur mediu pentru explorarea bazei de date, analizare și vizualizare. Cu alte cuvinte, această integrare va furniza instrumente pentru regăsirea informației, explorare și analizare și va permite utilizatorilor să controleze interactiv vizualizări animate ale conținutului bazei de date, interogări și rezultate ale interogărilor. Instrumentele de succes vor fi: *puternice* pentru a furniza imediat tehnologie cu valoare-adăugată, *flexibile* în procesarea de date și prezentarea ieșirilor pentru a evita conflictele cu clienții, *sensibile* în a furniza un mediu superior interactiv și *ușor de utilizat* pentru a încuraja experimentarea spontană.

Analiza exploratorie de date este un proces iterativ, în care interogările de nivel înalt (conceptual) conduc la interogări specifice, ale căror răspunsuri sunt examinate de către utilizatorii de informație geografică în vederea găsirii de șabloane interesante care, la rândul lor, pot sugera noi interogări. Provocarea implicată de acest tip de explorare este, în principal, legată de furnizare rapidă de operații de analiză, incrementale și reversibile, care generează continuu bucle de reacție. Este necesară o interfață utilizator grafică pentru a reorganiza instantaneu datele, pentru a juxtapune caracteristici în conformitate cu diverse criterii și pentru a vizualiza relațiile între diversele proprietăți ale acestor caracteristici. O caracteristică poate fi un obiect fizic, obiect abstract sau un eveniment. Integrarea GVis, KDD și GIS va combina ușurința în utilizare a sistemelor de manipulare directă/interactivă cu puterea sistemelor de interogare de baze de date.

Pentru atingerea acestui deziderat, este necesar un proces de management pentru a combina metodele și instrumentele asociate, care facilitează înțelegere științifică a seturilor mari de date, folosind un mediu de explorare singular. În această secțiune, este propusă și discutată o metodă de procesare în trei niveluri pentru executarea acestei integrări la trei niveluri:

- nivelul conceptual pentru înlănțuirea etapelor găsite în GIS, GVis și KDD,
- nivelul operațional pentru integrarea metodelor dezvoltate independent în fiecare dintre domenii,
- nivelul de implementare pentru combinarea diferitelor instrumente într-un singur mediu de lucru (sistem).

Principalul rezultat al acestei abordări constă în proiectarea unei vizualizări bazate pe cunoștințe pentru GIS temporal, în care accentul se mută de la metoda deductivă la cea condusă de date. În loc de a se căuta ipoteze pentru un model care pare că se potrivește datelor disponibile, accentul este pus acum pe date și pe descoperirea inerentă de relații în cadrul lor. Mai mult, această abordare va suporta metode noi de a interacționa cu seturi mari de date, va avea flexibilitate în a se ocupa de scări de variație a spațiului și timpului, mecanisme noi de identificare și trasare de incertitudine, precum și versatilitate în ceea ce privește manipularea de formate multiple de date.

### Nivelul conceptual

La nivelul conceptual, sunt identificate cerințele utilizatorului translatate în scopuri de nivel înalt, care urmează să fie atinse în construcția cunoștințelor folosite în progresul științific, sporirea profitului în afaceri, administrarea resurselor naturale și aplicații specifice. Aspectele critice ale procesului de construcție de cunoștințe de la acest nivel sunt:

- ce tip de date spațiotemporale se estimează a fi explorate (adică specifice mediului, socio-economice, discrete sau continue),
- ce rezultate particulare sunt așteptate de la acest proces (adică generarea de ipoteze, predicția unei schimbări ulterioare),
- cine sunt utilizatorii cunoștințelor obținute (adică, oameni de știință, analiști politici).

Deciziile privind acest nivel conceptual acționează ca și constrângeri de integrare ale etapelor GVis, GIS și KDD (a se vedea tabelul 2 pentru o trecere în revistă a etapelor actuale ale GVis, GIS și KDD). Abordările particulare pot să identifice cerințe mai degrabă pentru explorare vizuală de date, decât pentru un algoritm automatizat pentru data mining, conducând la o amalgamare a etapelor GVis (cunoașterea percepției vizuale) și a etapelor GIS (cunoașterea datelor). În acest caz, amalgamarea va fi definită în conformitate cu constrângerile din domeniile GVis și GIS. O constrângere comună constă în faptul că, înainte ca datele să poată fi vizualizate efectiv, pot fi cunoscute diverse metadate referitoare la structură și tip. Datele sunt organizate în baze de date (fișiere flat, tabele) în care fiecare bază de date se leagă de o metodă particulară de capturare de date.

**Tabelul 2: Etapele principale în GVis, GIS și KDD (preluare din Monica Wachowicz, "Integrating Gvis, GIS and KDD for Exploring Spatio-Temporal Data")**

<b>GVis</b> <i>Explorare condusă de percepție vizuală</i>	<b>GIS</b> <i>Explorare condusă de instrucțiuni de interogare</i>	<b>KDD</b> <i>Explorare condusă de procesarea automată de informații</i>
Explorare	Colectare de date	Selecția de date
Confirmare	Modelare de date	Pre-Procesare
Sinteză	Manipulare de date	Data Mining
Prezentare	Ieșire/Prezentare	Interpretare/Evaluare

### Nivelul operațional

Nivelul operațional se ocupă cu specificarea și combinarea celor mai potrivite metode, în vederea atingerii scopurilor conceptuale. Această integrare a metodelor GVis, GIS și KDD este foarte importantă în obținerea de avantaje de la experții umani, specialiștii în analiză (expertiza de domeniu și abilități de visual pattern recognition) și de la computere (putere mare de procesare). În continuare, se propune o perspectivă de integrare bazată pe metodele GVis, GIS și KDD pentru explorarea științifică a datelor spațio-temporale (tabelul 3), bazată pe o metodologie *task analysis* (Kirwan and Ainsworth, 1992). *Task analysis* a fost utilizată, la început, pentru definirea operațiunilor utilizator, precum și pentru proiectarea interfețelor utilizator și de sistem [19], [38].

*Task analysis* "...este un set de tehnici care pot fi utilizate pentru a determina care activitate a utilizatorului se dorește a fi îndeplinită, cum se planifică îndeplinirea ei și cum este ea îndeplinită – informații care plasează datele în context și plasează utilizatorul pe primul loc în procesul de formulare de cerințe" [38].

Ca prim pas către integrarea inteligibilă a GVis, GIS și KDD, este sugerat procesul de construcție și înlănțuire de instrumente care suportă un concept *metodă-activitate-operație*. Metodele sunt de forma celor descrise în secțiunile anterioare (a se vedea tabelul 3 pentru o trecere în revistă) și stabilesc "cum" se poate face explorarea în vederea atingerii scopului/scopurilor conceptual/conceptuale. O *activitate* este o frază conținând "ce" trebuie îndeplinit prin structurarea unei organizări ierarhice sau secvențiale de sarcini. O *operație* este acțiunea elementară perceptuală, motorie sau cognitivă, a cărei execuție este necesar să fie îndeplinită de activități.

**Tabelul 3: Trecere în revistă a taxonomiei dezvoltate în GVis, GIS și KDD  
(preluare de la Monica Wachowicz, "Integrating GVis, GIS and KDD for Exploring Spatio-Temporal Data")**

<b>GVis</b> <i>Pattern Identification Model</i> <i>MacEachren and Ganter (1990)</i>	<b>GIS</b> <i>Virtual GeoData Model, Albrecht</i> <i>(1996)</i>	<b>KDD</b> <i>Process Model [14]</i>
<b>Feature Identification</b> focusing, sequencing, multivariate glyphs, space- time cubes, small multiples, animation	<b>Search</b> interpolation, search-by- region, search-by-attribute	<b>Classification</b> symbolic methods, statistical methods
<b>Feature Comparison</b> scatterplot matrices, parallel coordinate plots, small multiples, map overlay, multivariate colour schemes, linking brushing	<b>Location Analysis</b> buffer, corridor, overlay, Voronoi/Thiessen	<b>Clustering</b> rule-based, set functions optimisation
<b>Feature Interpretation</b> cone tress, spider diagrams, spatialisation of information	<b>Terrain Analysis</b> slope, catchment areas	<b>Summarisation</b> data cube, attribute-oriented induction
	<b>Distribution/Neighbourhood</b> proximity, nearest neighbour	<b>Predictive Modelling</b> neural networks, induction trees
	<b>Spatial Analysis</b> Multivariate analysis, pattern/dispersion	<b>Change and Deviation Detection</b> Bayesian change detection
	<b>Measurements</b> distance, area, volume, fractal dimension	

#### Nivelul de implementare

La nivelul de implementare, opțiunile se fac în legătură cu algoritmi potriviți, care susțin activitățile, și în legătură cu mediile software/hardware specifice, în care se realizează operațiile. În rezumat, implementarea care suportă o metodă-activitate-operație la nivel operațional ar trebui să furnizeze instrumente interactive cum sunt:

- operații "drag" și "drop" pentru a crea interogări de la vizualizări și vizualizări de la interogări;
- limbaje vizuale de interogare pentru manipularea directă a datelor;
- sesiuni de explorare constând din interogări și vizualizări care pot fi salvate independent de orice date și pot fi reutilizate într-un set diferit de date;
- navigare printre obiecte multiple.

Scopul principal este acela de a furniza comunicație transparentă între mediile software, în care s-a propus ca metodele GIS, GVis și KDD să fuzioneze cu instrumente inteligibile. Obținerea acestei integrări de la nivelul de implementare va cere regândirea procesului de dezvoltare GIS în ideea de a se ocupa de aspecte critice cum sunt cele legate de modul în care utilizatorii de informație geografică vor explora seturi mari de date, căutând șabloane de date spațiotemporale printr-un singur mediu de explorare de date. La un moment dat, s-a obținut implementarea de bază a mișcării formatelor de date către instrumente GVis pentru a permite utilizatorilor să acceseze direct bazele de date cu informație geografică. Din păcate, o dată ce instrumentul GVis a remis o vizualizare tridimensională sau animație, este, în general, imposibil să se activeze interogări sau algoritmi pentru *data mining* de la dispozitivul de vizualizare. Prin urmare, eforturile au fost în zona de expansiune și integrare a diferitelor sisteme într-un mediu singular, care va cere standarde de date GIS deschise (open GIS data standards), legături software, funcții/operații și rețele de comunicații.

#### 3.4. SPIN! Spatial Mining for Data of Public Interest

SPIN! a fost lansat și finanțat de către Comisia Europeană în cadrul Programului Cadru 5 prin contractul IST-1999-10536 SPIN!

Conținutul acestei secțiuni este preluat din materialul disponibil la adresa: <http://www.ccg.leeds.ac.uk/spin/overview.html>.

## Contextul de lucru

Pătrunderea rapidă pe piață a tehnologiilor pentru *Data Mining* și *Geographic Information Systems* (GIS) este dirijată de presiunea venită din partea sectorului public, agențiilor de mediu și din partea sectorului industrial pentru furnizarea de soluții inovatoare pentru un spectru larg de probleme cum sunt de ex: servicii pentru asigurarea sănătății publice, agenții de mediu, care evaluează impactul schimbării folosinței terenului datorită schimbării climei, companii de geo-marketing care studiază segmentarea clienților/consumatorilor pe baza localizării spațiale.

Pentru suportarea acestor tipuri de analize, majoritatea GIS-urilor au doar funcționalitate de bază pentru analiză spațială. Multe se limitează la analize care implică afișare statistică descriptivă cum ar fi histogramme sau *pie charts*. *Data mining* care este căutare parțial automatizată, de șabloane ascunse în baze de date mari, oferă potențiale beneficii pentru aplicarea GIS ca bază pentru luarea deciziei în organizațiile din sectorul public și privat.

SPIN! este unul dintre cele mai inteligibile și mai ambițioase proiecte care și-a propus să pună laolaltă unele dintre cele mai interesante abordări din domeniul *data mining* și cartografie interactivă.

## Scopuri

Obiectivul principal al proiectului SPIN! este acela de a oferi posibilități pentru analizarea datelor georeferențiate. În acest scop, a fost dezvoltat sistemul Spatial Data Mining, care integrează nivelul de funcționalitate la care s-a ajuns până acum în GIS și *Data Mining*, într-o arhitectură deschisă, extensibilă, *internet-enabled plug-in*. Domeniul *Data Mining* va progresa prin adaptarea metodelor din *Machine Learning* și *Bayesian Statistics* la analiza spațială. Domeniul GIS va progresa prin dezvoltarea de noi metode pentru vizualizarea informației temporale și spațiale. Sistemul SPIN! pentru *spatial mining* va fi testat și evaluat în cadrul unor aplicații dedicate analizei seismice și vulcanologice și pentru diseminarea bazată pe web a datelor de recensământ.

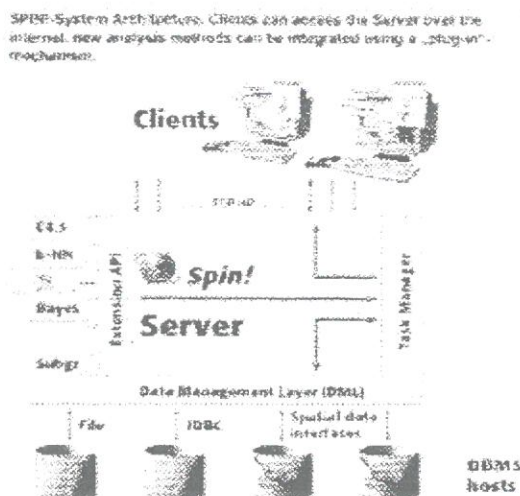


Figura 5: Arhitectura SPIN! (preluată din materialul disponibil la adresa: <http://www.ccg.leeds.ac.uk/spin/overview.html>)

În ultimii ani, câțiva dintre partenerii proiectului au dezvoltat componente tehnologice și instrumente științifice cerute de nucleul unui astfel de sistem. Pe parcursul acestui proiect aceste eforturi individuale și expertiza asociată se vor reuni la nivel european. În această idee, vor fi dezvoltate pisele care lipsesc și va fi construită o platformă integrată GIS-Data-Mining. Partenerii din sectorul industrial vor dezvolta un model de business pentru brokering bazat pe web, cu date statistice georeferențiate, și vor face estimări privind impactul economic al tehnologiei.

Mai multe detalii despre cele prezentate în acest capitol se găsesc în Raportul tehnic 630/2, elaborat în mai 2002: "Some key aspects in geospatial data mining" [36], iar aspecte legate de tehnicile și algoritmi folosiți în *geospatial data mining* sunt prezentate în Raportul tehnic 630/1 elaborat în aprilie 2002: "Techniques and algorithms in geospatial data mining" [35].

## 4. Concluzii și tendințe

Progresele făcute, în ultimii zece ani, în domeniul teledetecției și capturării de date spațiale/geografice au sporit în mod dramatic posibilitățile de colectare zilnică de informație geografică la nivel de terabytes. Cu toate acestea, bunăstarea datelor geografice nu poate fi complet realizată atunci când informația implicită din date este dificil de separat, ceea ce duce la confruntarea specialiștilor în informație geografică cu o cerere urgentă de metode și de instrumente noi, care pot să transforme în mod inteligent și automat *datele* geografice în *informații* și să sintetizeze *cunoștințe* geografice. Trebuie să se apeleze la noi abordări în reprezentarea geografică, procesarea de interogări, analiza spațială și vizualizarea de date [18], [46], [62]. Specialiștii în informatică sunt confrunțați cu aceeași provocare ca rezultat al revoluției digitale ce expediază date de ordinul terabytes de la carduri de credit tranzacționale, examinări medicale, apeluri telefonice și de la alte numeroase activități umane. Eforturile comunităților din domeniul inteligenței artificiale, statistică și baze de date, au sprijinit tehnologiile KDD în extragerea de informații din masivele de date pentru asigurarea suportului în luarea deciziei [2], [24], [31].

În prima parte a acestui capitol, a fost identificat potențialul impact asupra *geographic information science* și asupra cercetării în sens mai larg. În partea a doua, este prezentată o listă de subiecte de cercetare, oferită de Miller and Han [46].

### 4.1. Impactul potențial asupra *geographic information science* și asupra cercetării

Există cerințe unice și provocări în ceea ce privește descoperirea de cunoștințe geografice, în cadrul *geographic information science*. Majoritatea bazelor de date digitale sunt, în cel mai bun caz, o reprezentare foarte simplă a cunoștințelor geografice la nivelul geometriei elementare, constrângerilor topologice și provenite de la măsurători. Se estimează că GIS-urile bazate pe cunoștințe vor construi cunoștințe geografice de nivel înalt, în cadrul bazelor de date geografice digitale, pentru analizarea fenomenelor complexe [58], [61]. Descoperirea de cunoștințe geografice este o sursă potențială bogată, pentru GIS bazat pe cunoștințe și pentru analiză spațială inteligentă. Dezvoltarea reprezentărilor cunoștințelor geografice descoperite, care sunt efective în GIS bazate pe cunoștințe și analiză spațială, este una din abordările critice provocatoare.

#### Descoperirea de cunoștințe geografice în cercetarea geografică

Informația geografică a fost întotdeauna un obiect de uz curent central în cercetarea geografică. Din punct de vedere istoric, cercetarea geografică a apărut în medii sărace în date. Multe dintre progresele din cercetarea geografică vor sta la baza îmbunătățirilor tehnologiilor pentru georeferențiere, capturare, memorare și procesare de date geografice. Revoluția produsă în domeniul datelor geografice digitale poate fi considerată ca fiind cea mai dramatică schimbare de atitudine în mediul dedicat cercetării geografice din istoria științei. Aceasta generează, poate, cea mai importantă "meta-chestiune" pentru cercetarea geografică a acestui secol și anume "care sunt problemele la care nu s-a găsit un răspuns până acum?"

Ne aflăm încă la începuturile istoriei descoperirii de cunoștințe geografice. În acest moment, în opinia noastră, poate fi oferită doar o listă a aplicațiilor *geographic knowledge discovery* în *geographic information science* și, în general, în cercetarea geografică.

- **Interpretarea hărților și extragerea de informații**

Malebra ș.a. [43] au demonstrat utilizarea de instrumente *inductive machine learning* în cadrul mediilor GIS. Sistemul lor poate să extragă și să interpreteze caracteristici umane și fizice complexe din hărți topografice pentru introducerea lor într-un GIS și pentru analiză.

- **Extragerea de informație de la imaginile provenite din teledetecție**

Creșterea rezoluțiilor spațiale, temporale și spectrale, furnizată de progresele făcute de tehnologiile specifice teledetecției, a condus la crearea de masive de baze de date de imagini. Aceste baze de date dau cercetătorilor posibilitatea de a analiza și înțelege informația din aceste date. Gopal ș.a. [29] au folosit rețelele neuronale artificiale, combinate cu tehnici de vizualizare, pentru interpretarea și înțelegerea șabloanelor extrase din imaginile provenite de la teledetecție.

- **Carcateristicile de mediu ale hărților**

Multe dintre fenomenele geografice au atribute complexe, multidimensionale, care sunt dificil de rezumat și de integrat folosind metodele analitice tradiționale. Eklund ș.a. [13] au folosit tehnici de *inductive learning* și rețele neuronale artificiale pentru a clasifica și reprezenta pe hartă tipuri de sol. Lees ș.a. [41] au folosit metode de inducție, bazate pe arbori de decizie, pentru reprezentarea pe hartă a tipurilor de vegetație, în zonele în care metodele de clasificare din teledetecție confundă terenul cu tulburările neobișnuite (de ex. incendii).



- **Extragerea de șabloane spațio-temporale**

Identificarea șabloanelor neobișnuite din masivele de baze de date spațiotemporale poate fi dificilă, iar numărul de șabloane posibile poate fi foarte mare. Mesrobian ș. a. [44] au dezvoltat Open Architecture Scientific Information Sistem (OASIS) pentru interogarea, explorarea și vizualizarea fenomenelor geofizice din baze de date mari, eterogene, distribuite. Componenta Conquest Scientific Query Processing Sistem a lui OASIS, identifică activitatea ciclonică din datele despre vreme și climă prin extragerea de șabloane neobișnuite din presiunea aerului și vânt în timp. În alte domenii, Openshaw și colegii lui [52], [53] au dezvoltat tehnici exploratorii, bazate pe metode simple de interogare, pentru afișarea de clustere spațiotemporale în datele criminalistice.

- **Interacțiune, flux și mișcare**

Interacțiunea spațială, fluxul și mișcarea în spațiul geografic pot să furnizeze noi puncte de cercetare în structura spațială a sistemelor geografice fizice și umane. Structura spațială și interacțiunea spațială sunt strâns legate: locația influențează șabloanele de interacțiune, în timp ce șabloanele de interacțiune influențează locația entităților și a activităților. Din rațiuni de maleabilitate, analiza spațială și rețeaua analitică formulează ipoteze puternice în legătură cu influențele asupra fluxului, interacțiunii, mișcării și localizării, doar prin capturare directă și efecte de aproximare în spațiu și timp. Influențele de ordin  $n$  pot fi îngropate în masive de baze de date de interacțiune, fluxul și mișcarea fiind capturate prin sisteme de monitorizare în timp real, sisteme de transport inteligente și dispozitive "position-aware" cum sunt telefoanele celulare și clienții *wireless Internet*. Marble et al. (1997) au descris metode de vizualizare pentru a explora matrici de interacțiune masive. Smyth (2001) a explorat posibilitățile pentru descoperire de cunoștințe geografice din traiectoriile spațiu-timp ale dispozitivelor mobile.

## 4.2. Aspecte critice în provocările cercetării

Există numeroase aspecte critice în descoperirea de cunoștințe geografice și data mining. Miller și Han [46] oferă următoarea listă de subiecte de cercetare în acest domeniu:

- **Dezvoltarea și suportul pentru geographic data warehouses**

Trebuie spus că nu există o adevărată *geographic data warehouse* (GDW). Proprietățile spațiale sunt deseori reduse la atribute spațiale simple în cadrul *data warehouses*. Crearea de GDW integrate cere rezolvarea unor probleme de interoperabilitate spațială și temporală a datelor, incluzând diferențele din semantică, sistemele de referință, geometrie, acuratețe și poziție.

- **Cele mai bune reprezentări în geographic knowledge discovery**

Tehnicile actuale de *geographic knowledge discovery* (GKD) utilizează, în general, reprezentări foarte simple ale obiectelor geografice și ale relațiilor spațiale. Tehnicile pentru *geographic data mining* vor recunoaște obiecte geografice mult mai complexe (linii și poligoane) și relații (distanțe neeuclidiene, direcție, conectivitate și interacțiune în cadrul spațiului geografic de atribute). În aceste reprezentări geografice și relații trebuie să fie complet integrat timpul.

- **Geographic knowledge discovery folosind diferite tipuri de date**

Pot fi dezvoltate tehnici de *geographic knowledge discovery* GKD astfel încât să poată fi manipulate diverse tipuri de date începând cu modelele tradiționale raster și vector, până la imagini multimedia și georeferențiate precum și tipuri de date dinamice (video streams, animație și realitate virtuală).

- **Interfețe utilizator pentru geographic knowledge discovery**

GKD necesită schimbarea de optică de la cercetătorii orientați tehnic către comunitățile de cercetare GIScience și către alte comunități de cercetare. Această schimbare de optică necesită *interfețe* și instrumente care pot să ajute cercetătorii în aplicarea acestor tehnici la chestiuni specifice.

- **Verificarea conceptelor și benchmarking**

Ca și în alte domenii, și în KDD și DM trebuie să fie unele cazuri de test, definitive sau *benchmarks*, care să ilustreze puterea și utilitatea GKD în descoperirea de cunoștințe geografice neașteptate. O abordare conexasă este și cea legată de *benchmarking* în vederea determinării efectelor datelor de diferite calități asupra cunoștințelor geografice descoperite.

- **Construirea cunoștințelor geografice descoperite în cadrul GIS și analiza spațială**

Sunt absolut necesare reprezentări ale cunoștințelor geografice descoperite, care sunt potrivite pentru GIS și analiză spațială. Aceasta poate să includă interfețe GIS bazate pe *on line analytical processing* (OLAP) și instrumente inteligente pentru ghidarea analizei spațiale.

În cartea "Advances in Knowledge Discovery and data Mining", publicată de MIT Press, Usama M. Fayyad și editorii fac următoarea remarcă: "...din combinarea celor doi termeni "data mining" și "data warehousing",

este de așteptat să se construiască un pod între comunitățile de specialiști din statistică, baze de date și machine learning, care să atragă foarte mult interesul dezvoltatorilor de sisteme informatice”.

Bazată pe studiul literaturii de specialitate și pe experiența proprie, predicția formulată de specialiști este aceea că modul de familiarizare cu termenii “data mining” și “data warehousing” va fi similar cu ce s-a întâmplat cu termenul “GIS” în ultimii cinci ani.

## Mențiuni

Acest articol este un rezumat al Raportului Tehnic, elaborat de dr. Angela Ioniță (Institutul Național pentru Cercetare-Dezvoltare în Informatică – ICI București) în mai 2002, în cadrul proiectului INTAS nr. 397: **DATA MINING TECHNOLOGIES AND IMAGE PROCESSING: THEORY AND APPLICATIONS** coordonat de Lappeenranta University of Technology din Finlanda.

## Bibliografie

1. **ANSELIN, L.:** Local Indicators of Spatial Association. LISA. În: *Geographical Analysis*, 27, 1995.
2. **BHANDARI, E. E. COLET, J. PARKER, Z. PINES, R. PRATAP, R. PRATAP, K. RAMANUJAM:** Advanced Scout: Data Mining and Knowledge Discovery in NBA Data. În: *Data Mining and Knowledge Discovery*, 1, 1997.
3. **BRACHMAN, R.J., T. ANAND:** The Process of Knowledge Discovery in Databases in *Advances in Knowledge Discovery and data Mining*, U.M. Fayyad et al. (eds.) AAAI Press/ The MIT Press, 1996.
4. **BRADSIL, P. B. K. KONOLIGE (Eds.):** Meta-Learning, Meta-Reasoning and Logics, Kluwer Academic Press., Boston, MA, USA, 1990.
5. **BRUNSDON, C., A. S. FOTHERINGHAM, M. E. CHARLTON:** Geographically Weighted Regression: A Method for Exploring Spatial Nonstationarity. În: *Geographical Analysis*, 28, 1996.
6. **BUTTENFIELD, B.P.:** Looking Forward: Geographic Information Services and Libraries in the Future. În: *Cartography and GIS*, 25(3), 1998.
7. **BUTTENFIELD, B.P., M. GAHEGAN, H. MILLER, M. YUAN:** Geospatial Data Mining and Knowledge Discovery.
8. **CÂMARA, A. S. J. RAPER (Eds.):** Spatial Multimedia and Virtual Reality, London: Taylor and Francis, 1999.
9. **CANTU-PAZ, ERIK and CHANDRIKA KAMATH:** On the Use of Evolutionary Algorithms in Data Mining.
10. **CHEESEMAN, P. J. STUTZ:** Bayesian Classification: Theory and Results. În: Fayyad, U., Piatetsky-Shapiro, G, Smyth, P. and Uthurusamy, R. (Eds.) *Advances in Knowledge Discovery and Data Mining*, Cambridge, MA: AAAI/MIT Press, 1996.
11. **CHEN, M., J. HAN, P.S. YU:** Data Mining – An Overview from a Database Perspective: Focussing Techniques for Efficient Class Identification. În: *Proc. of Int. Symposium on large databases (SSD'95)*, Maine, 1996.
12. **EGENHOFER, M. J. J.R. HERRING:** Categorizing Binary Topological Relations Between Regions, Lines and Points in Geographic Databases. În: M. Egenhofer, D. M. Mark and J. R. Herring (Eds.), *The 9-intersection: Formalism and its Use for Natural-language Spatial Predicates*, Santa Barbara, CA: National Center for Geographic Information and Analysis Technical Report 94-1), 1994.
13. **EKLUND, P. W., S. D. KIRKBY, S. D., A. SALIM:** Data mining and soil salinity analysis. În: *International Journal of Geographical Information Science*, 12, 1998.
14. **FAYYAD, U., G. PIATETSKY-SHAPIRO, P. SMYTH:** From Data Mining to Knowledge Discovery in Databases. În: *AI Magazine*, Fall, 1996.
15. **FAYYAD, U.:** Editorial. *Data Mining and Knowledge Discovery*, 1997.
16. **FLEWELLING, D. M., M. J. EGENHOFER:** Using Digital Spatial Archives Effectively. În: *The Int. Journal of Geographical Information Science*, 1999.
17. **FOTHERINGHAM, A. S., M. CHARLTON, C. BRUNSDON:** Two Techniques for Exploring Non-Stationarity in Geographical Data. În: *Geographical Systems*, 1997.

18. **GAHEGAN, M.:** On the Application of Inductive Machine Learning Tools to Geographical Analysis. În: *Geographical Analysis*, 2000.
19. **GAHEGAN, M.** Gahegan, M. and O'Brien, D. (1997). *A Strategy and Architecture for the Visualisation of Complex Geographical Datasets*. *International Journal of Pattern Recognition and Artificial Intelligence*, 11(2);
20. **GAHEGAN, M., M. TAKATSUKA, M. WHEELER, F. HARDISTY:** GeoVISTA Studio: A Geocomputational Workbench. În: Proc. 4th Annual Conference on GeoComputation, UK, August 2000. URL: <http://www.ashville.demon.co.uk/gc2000/>.
21. **GAHEGAN, M.:** National Academies White Paper. În: *Intersection of Geospatial Information and Information Technology*. September, 2001.
22. **GAHEGAN, M.:** Data Mining and Knowledge Discovery in the Geographical Domain.
23. **GAHEGAN, M., M. WACHOWICZ, M. HARROWER, T. M. RHYNE:** The Integration of Geographic Visualization with Knowledge Discovery in Databases and Geocomputation. În: *Cartography and Geographic Information Systems*, special issue on the ICA research agenda, 2001.
24. **GARDNER, C.:** IBM Data Mining Technology, Stamford, Connecticut: IBM Cooperation, 1996.
25. **GETIS, A., J.K. ORD:** The Analysis of Spatial Association by Use of Distance Statistics. În: *Geographical Analysis*; 1992.
26. **GETIS, A., J.K. ORD:** Local spatial statistics: An overview. În: P. Longley and M. Batty (Eds.) *Spatial Analysis: Modelling in a GIS Environment*, Cambridge, UK: GeoInformation International, 1996.
27. **GOODCHILD, M. F.:** Geographical Data Modeling. În: *Computers and Geosciences*, Vol. 18, No. 4, 1992.
28. **GOODCHILD, M. F., B.P. BUTTENFIELD, P. ADLER, A. KRYGIEL, H. ONSRUD, R. KAHN:** Distributed Geolibraries. National Research Council Monograph., Washington, D.C.: National Academy Press, 1999.
29. **GOPAL, S., W. LIU, C. WOODCOCK:** Visualization Based on Fuzzy ARTMAP Neural Network for Mining Remotely Sensed Data. În H. J. Miller and J. Han (Eds.) *Geographic Data Mining and Knowledge Discovery*, London: Taylor and Francis, 2001 (in press).
30. **GUNOPULOS, D.:** Data mining techniques for geospatial applications.
31. **HEDBERG, S. R.:** Search for the Mother Lode: Tales of the First Data Miners. În: *IEEE Expert*, 11(5), 1996.
32. **HEALEY, R.:** Database Management Systems. În: Maguire, D., Goodchild, M.F., and Rhind, D., (Eds.), *Geographic Information Systems: Principles and Applications*, London: Longman, 1991.
33. **HINKE, T., J. S. RUSHING, S. KANSAL, GRAVES, AND H. RANGANATH:** For Scientific Data Discovery: Why Can't the Archive be More Like the Web. În: Proc. Of the 9th Int. Conf. on Scientific Database Management, Olympia, WA, Aug. 11-13, 1997.
34. **HORNSBY, K., M. JEGENHOFER:** Identity-based Change: A Foundation for Spatio-Temporal Knowledge Representation. În: *International Journal of Geographical Information Science*, 14, 2000.
35. **IONIȚĂ, A.:** Techniques and Algorithms in Geospatial Data Mining. Technical Report, TR ICI 630/1, April 2002.
36. **IONIȚĂ, A.:** Some Key Aspects in Geospatial Data Mining. Technical report, TR ICI 630/2, May, 2002.
37. **JACKSON, J.E.:** An User's Guide to Principial Components, New York, NY:John Wiley, 1991.
38. **KNAPP, L.:** A Task Analysis Approach to the Visualisation of Geographic Data. În: *Cognitive Aspects of Human-Computer Interaction for Geographic Information Systems*, T.L. Nyerges et al. (Eds.), Kluwer Academic Publishers, 1995.
39. **KOPERSKI, K., J. HAN, J.:** Discovery of Spatial Association Rules in Geographic Information Databases. În: Proc. of the 4th International Symposium on Large Spatial Databases, SSD95, Maine, 1995.
40. **KOPERSKI, K., J. HAN, J. ADHIKARY:** Mining Knowledge in Geographic Data. În: *Comm. ACM*, 1999. Available at URL: <http://db.cs.sfu.ca/sections/publication/kdd/kdd.html>
41. **LEES, B. G. K. RITMAN:** Decision-Tree and Rule-Induction Approach to Integration of Remotely Sensed and GIS Data in Mapping Vegetation in Disturbed or Hilly Environments. În: *Environmental Management*, 15, 1991.
42. **MACEACHREN, A.M., M. J. KRAAK:** Exploratory Cartographic Visualization: Advancing the Agenda. În: *Computers and Geoscience* (4); 1997.

43. **MALEBRA, D., F. ESPOSITO, A. LANZA, F. LISI:** Machine Learning for Information Extraction from Topographic Maps. În: H. J. Miller and J. Han (Eds.), *Geographic Data Mining and Knowledge Discovery*, London: Taylor and Francis, 2001 (in press).
44. **MESROBIAN, E, R. MUNTZ, E. SHEK, S. NITTEL, M. LA ROUCHE, M. KRIGUER, C. MECHOSO, J. FARRARA, P. STOLORZ, H. NAKAMURA:** Mining Geophysical Data for Knowledge. În: *IEEE Expert*, 1996.
45. **MILLER, H. J.:** Geographic Representation in Spatial Analysis. În: *Journal of Geographical Systems*, 2, 2000.
46. **MILLER, H. J., J. HAN:** Discovering Geographic Knowledge in Data Rich Environments: A report on a specialist meeting, SIGKDD Explorations: Newsletter of the, Association for Computing Machinery, Special Interest Group on Knowledge Discovery and Data Mining, 1(2), 2000, available at <http://www.acm.org/sigs/sigkdd/explorations>.
47. **MILLER, H. J., J. HAN (Eds):** *Geographic Data Mining and Knowledge Discovery*, London: Taylor & Francis, 2001.
48. **MILLER, H. J., H. JIAWEI:** *Geographic Data Mining and Knowledge Discovery: An Overview*. În: H. J. Miller and J. Han (Eds.) *Geographic Data Mining and Knowledge Discovery*, London: Taylor and Francis, 2001 (in press).
49. \* \* \*: NPR National Public Radio Morning Edition. Report on Emergence of Commercial Internet Search Engines such as Yahoo and Other .dot-coms., 3 April, 1998.
50. \* \* \*: NRC Data Foundation for the National Spatial Data Infrastructure. National Research Council Mapping Science Committee, Sugarbaker, L.A., Chair, Washington, D.C.: National Academy Press, 1995.
51. **OKABE, A, H. J. MILLER:** Exact Computational Methods for Calculating Distances Between Objects in a Cartographic Database. În: *Cartography and Geographic Information Systems*, 1996.
52. **OPENSHAW, S.:** The Modifiable Areal Unit Problem. *CATMOG 38*, Norwich: Geo Abstracts, 1984.
53. **OPENSHAW, S.:** Two Exploratory Space-time-attribute Pattern Analysers Relevant to GIS. În: A. S. Fotheringham and P. A. Rogerson (Eds.), *Spatial Analysis and GIS*, London: Taylor and Francis, 1994.
54. **OPENSHAW, S.:** *Geographical Data Mining: Key Design Issues*, 1999, available at the address: [http://www.geovista.psu.edu/sites/geocomp99/Gc99/051/gc\\_051.htm](http://www.geovista.psu.edu/sites/geocomp99/Gc99/051/gc_051.htm)
55. **PEUQUET, D. J., C. X. ZHANG:** An Algorithm to Determine the Directional Relationship Between Arbitrarily-shaped polygons in the Plane. În: *Pattern Recognition*; 1987.
56. **PEIRCE, C. S.:** Deduction, Induction and Hypothesis. În: *Popular Science Monthly*, 13, 1878.
57. **RODDICK, J. F., B. LEES:** Paradigms for Spatial and Spatio-temporal Data Mining. În: H. J. Miller and J. Han (Eds.), *Geographic Data Mining and Knowledge Discovery*, London: Taylor and Francis, 2001 (in press).
58. **SRINIVASAN, A., J. A. RICHARDS:** Analysis of GIS Spatial Data Using Knowledge-based Methods. În: *International Journal of Geographical Information Systems*, 7, 1993.
59. **WACHOWICZ, M.:** Integrating Gvis, GIS and KDD for Exploring Spatio-Temporal Data.
60. **WATTERSON, K.:** A Data Miner's Tools. În: *BYTE*, October 1995.
61. **YUAN, M.:** Use of Knowledge Acquisition to Build Wildfire Representation in Geographic Information Systems. În: *International Journal of Geographical Information Systems*, 11, 1997.
62. **YUAN, M.:** Representing Spatiotemporal Processes to Support Knowledge Discovery in GIS databases. În: T. K. Poiker and N. Chrisman (Eds.), *Proc. of the 8th International Symposium on Spatial Data Handling Spatial Data Handling*; 1998.
63. **YUAN, M., B. BUTTENFIELD, M. GAHEGAN, H. MILLER:** *Geospatial Data Mining and Knowledge Discovery. A UCGIS White Paper on Emergent Research Themes*, 2001. URL: <http://www.ucgis.org/emerging/>.
64. **ZAIANE, O. R., J. HAN, Z-N. LI, J. HOU:** Mining Multimedia Data. Proceedings, CASCON'98: Meeting of Minds, Toronto, Canada, November 1998; available at: <http://db.cs.sfu.ca/sections/publication/smmdb/smmdb.html>
65. <http://db.cs.sfu.ca/DBMiner>
66. <http://db.cs.sfu.ca/DBMiner>