

Gáspár Bencéné Vér Katalin*

ADATBÁNYÁSZAT A GAZDASÁGI ÉLETBEN

Az adatbányászat egy *döntéstámogatási módszer*, olyan üzleti intelligencia megoldás, amely új üzleti lehetőségeket segít megtalálni és kiaknázni a *nagytömegű adathalmazokban rejlő, nem ismert összefüggések feltárásával*. Az angol *Data mining* kifejezés tükörfordítása. Egyesíti az adatbázis-kezelés, a statisztika és a mesterséges intelligencia kutatások eredményeit. Az adatbányászat kifejezést a különböző informatikai cégek különböző értelemben használják, de a szigorúbb szakmai terminológia nem tekinti adatbányászatnak az adatokból egyszerű lekérdezésekkel, aggregálásokkal, illetve alap-statisztikai vizsgálatokkal történő információ-nyerést (<http://hu.wikipedia.org/wiki/Adatbányászat>).

Az elmúlt évek során az adatbányászat iránti érdeklődés rohamosan nőtt mind a korszerű informatikai szolgáltatásokat nyújtó cégek, mind az ilyen szolgáltatásokat igénybe vevő vállalatok körében. Az adatbányászati funkció hatékony infrastruktúrát biztosít az egész vállalatra kiterjedő, magas szintű üzleti adatelemző alkalmazások létrehozásához. Az új üzleti ismeretek kinyerésének és terjesztésének automatizálásával és integrálásával a vállalatok ki tudják használni az adataikban fekvő információkat, hatékonyabban működhetnek, és nagyobb versenyelőnyhöz juthatnak.

Az adatbányászat mint az adatelemzés eszköze és lehetősége – a statisztika hasonló kategóriáinak megfelelően – két nagyobb kategóriába sorolható:

- a *leíró adatbányászat* az adatok alap jellemzőinek meghatározását jelenti.
- a *következtetési adatbányászat* alapvetően összefüggések feltárásával foglalkozik.

A statisztikai eszközökkel „kis és közepes” adatmennyiségek esetén meg lehet találni bizonyos szabályszerűségeket és korrelációkat, de ezek az eszközök igazán nagy mennyiségű adattal már nem képesek megbirkózni. Az adatbányászatnak nincsenek ilyen korlátai. Az adatbányászat az adatok mélyére hatol. (Gáspár, 2006) Az adatbányászat alapadatai egyaránt lehetnek üzleti, kutatási, mérési adatok. Lényegében bármilyen nagytömegű adathalmaz elemei képezhetik az adatbányászat alapadat-állományait. Az adatok két köréhez kötődik az adatbányászat speciális területe: a szövegbányászat (textmining) és a webbányászat (webmining).

„A *szövegbányászat* a tudásmenedzsment egyik leghatékonyabb információtechnológiai eszköze. Tudásalapú technológia, amely képes a mesterséges intelligencia, a gépi tanulás, a természetes nyelvi feldolgozás, a nyelvtechnológia, a többváltozós matematikai statisztika, a valószínűségszámítás, a tartalomelmélet és még jónéhány tudományág legfrissebb eredményeinek produktív hasznosítására. A szövegbányászat segítségével olyan rejtett ismeretanyag nyerhető ki strukturálatlan szöveges dokumentum állományokból, amelyek kinyerése egyrészt emberi erővel kivitelezhetetlen lenne, másrészt pedig olyan tudásvagyont képvisel, amelyből rövid úton üzleti versenyelőny kovácsolható.” (Vázsonyi, 2006)

A *web-bányászat* olyan gépesített dokumentum-feldolgozó szakterület, amely az internethez kapcsolódóan a weboldalakon található képi, szöveges és egyéb alakú adatok feldolgozhatóvá való átalakításával foglalkozik. Az adatok átalakításának a célja többnyire megfelelő input-állományok kialakítása további, a felhasználó szempontjából értékes adatok kinyerése érdekében. Ezt a munkát speciális alakfelismerő programok végzik. A megrendelők többnyire a kereskedelmi és hírszerzés területéről való cégek. (<http://hu.wikipedia.org/wiki/Adatbányászat>) Az adatbányászat kialakulásának indítéka a napjainkban szinte az élet valamennyi területén, de különösen a gazdasági szférában, a vállalatoknál, a napi tevékenység során keletkező, *nagy tömegű adat*. Ez különösen igaz a *nagyvállalatokra* és a *pénzügyi szektorra*, ahol *minden tevékenységi lépés adatbázison alapul*.

* Főiskolai tanár, Általános Vállalkozási Főiskola

Ahhoz, hogy olyan adathalmazunk/állományunk legyen, amellyel az adatbányászati munka megkezdődhet egy adatelemzési, logikai folyamatnak kell végbemennie. Ennek lépései:

- adatbázis audit – adatforrások és tárhelyek, kapcsolódásaik feltérképezése,
- adattisztítás,
- adatintegrálás,
- adatkiválasztás,
- adattranszformálás,
- egységes adatmodell kidolgozása.

Megfelelő adatállomány(ok) – ez sok esetben egy adattárházat jelent – birtokában kezdődhet maga az *adatbányászat*, amely ugyancsak több lépésből áll. Első lépés az adott üzleti probléma megértése. Ezt az adatok megismerése, és modellezés céljára történő előkészítése követi. Ezután kerül sor a statisztikai alapelemzések elvégzésére, majd a részletes adatelemzésre, amely magában foglalja az elemzési módszerek, szoftverek kiválasztását, és az elemzés elvégzését. A következő lépés a modellezés, az adott elemzés céljából érdekes összefüggések feltárása, valamint a kapott eredmények értékelése, majd a kapott eredmények ismertetése, bemutatása. Az eredmények felhasználása ezután kezdődhet el.

Mire keresi az adatbányászat a választ? A teljesség igénye nélkül csak néhány példa:

- személyek/szervezetek legvalószínűbb reakciói,
- együtt értékesített/vásárolt termékek/szolgáltatások,
- kitől várható hogy a konkurenciához pártol,
- csalások felderítése,
- azonos terméket vásárlók közös tulajdonságai,
- mi érdekli a látogatót a legjobban egy Web oldalon.

Az adatbányászat jellemzői

- Új távlatok (ahol a hagyományos módszerek nem elég jók, vagy nem elég gyorsak az igazán hatékony adatelemzéshez).
- Valós múltbeli összefüggések alapján segít a jövő döntéseinek megalapozásában.
- Általános módszer, amely minden üzleti területen hatékonyan alkalmazható.
- Hatékony eszköze az elektronikus kereskedelem bevezetésének és fejlesztésének.

Az adatbányászat módszertana a

- társításon,
- csoportosításon,
- osztályozáson,
- visszafelé haladáson (neurális hálók, regresszió, kapcsolatok felderítése),
- összegzésen alapul.

Vizsgáljuk meg, hogy *milyen információkat bányásznak* a különböző módszerek segítségével?

A *társítás (association)* egy eseményt, vagy tárgyat, például árucikket rendel egy másikhoz. Például: az asszociáció a marketingben megtalálja (felismeri), hogy mely termékeket vásárolnak együtt.

A *csoportosítás (clustering)* az adatok korábban nem ismert rendszereinek, együvé-tartozásának a felismerése, az adatokat leíró véges kategóriásorok azonosítása. Azonos tulajdonságok alapján csoportosítja a vizsgált adatokat. *Ezt az információt használják fel csoport-képzésre további elemzésekhez.*

Az *osztályozás (classification)* az adatok újfajta szerveződését eredményező minták észrevétele, amely alapján egy adat egy vagy több előre meghatározott osztályba rendezhető. (Például: drága sportkocsit vásárlók tipikusan fiatal, városi diplomások, magas jövedelemmel.)

A *visszafelé haladás*, amely neurális hálókat is alkalmaz, az adatokat valós-értékű előrejelzés változóhoz (*real-valued prediction variable*) társító funkció, illetve a változók közötti kapcsolatok felderítésére (*dependency modelling*) használatos módszer. Ilyen a hasonló idősorok keresése. Felfedi a hasonló sorozatokat egy adott időszakban, vagy felfedi a hasonló sorozat-párokat. Például: felfedi a hasonló árfolyam-mozgású részvényeket.

A *kivétel-keresés*, ahogy az elnevezése is mutatja, a „szokatlant” keresi. Például: felfedi a szokatlan hitelkártya tranzakciót, és ezáltal detektálja csalásokat.

Az *előrejelzés (trendek)* az adatminták alapján becsli a jövőbeli értékeket. A gazdasági menedzserek számára legfontosabb információ az előrejelzés, de gyakran kombinálják az előzőekben ismertetett információkat is.

Az *összegzés (summarization)* az adatok egy meghatározott részhalmozának tömör leírására vonatkozik.

Az adatbányászat technikái a

- „felügyelt tanulás”
- „felügyelet nélküli tanulás”

Felügyelt tanulás

A *felügyelt tanulásnál* az adatelemzőnek ki kell jelölnie a célmezőt vagy a függő változót. A felügyelt tanulási technikával ezután átvizsgálja az adatokat, hogy *szabályszerűségeket és összefüggéseket fedezzen fel* a független változók és a függő változók között. *Felügyelt tanulási modellnél* az adatbányászati algoritmus aprólékosan átvizsgálja az adatokat, rejtett szabályszerűségeket tár fel, olyan *modellt* alkot, amely a lehető legjobban leírja a függőségeket.

Az adatokat általában három részre osztják:

- betanítási – kezdeti modell,
- tesztelési – modellfinomítás,
- kiértékelési adatok, előrejelzések készítése.

Felügyelet nélküli tanulás

A *felügyelet nélküli tanulásnál* a felhasználó nem adja meg a célt az adatbányászati algoritmus számára. Az *adattársítási és csoportosító algoritmusok* ilyenkor *semmit sem feltételeznek a célmezőről*. Ilyen esetben az adatbányászati algoritmusok minden előre definiált üzleti cél nélkül keresnek kapcsolatokat és csoportokat.

Az adatbányászat talán leggyakoribb gyakorlati felhasználási területei

- a direktmarketing,
- a kockázatelemzés,
- a keresztértékesítés,
- a visszaélések felderítése,
- az ügyfélmegtartás,
- az internetes oldalak vizsgálata.

Üzleti döntések adatbányászat segítségével

A cégek négy lehetőség közül választhatnak, amennyiben üzleti döntéseiket adatbányászat segítségével kívánják megtámogatni:

1. Megvásárolhatják az adott üzleti problémára létrehozott adatbányászati modell eredményeit, amely szabályok, ügyféllisták, pontszámok formáját öltheti. Ilyenkor nincs saját modell, lényegében információt vásárol a cég.

2. Megvásárolhatnak egy olyan, beágyazott adatbányászati alkalmazást tartalmazó szoftvert, amellyel kifejezetten a saját üzleti problémájukra vonatkozó modellezéseket lehet végezni, többek között csalásfelderítésre, lojalitásnövelésre vagy kampánymenedzselésre vonatkozóan.

3. Külső adatbányászati szakértőket bízhatnak meg előre definiált feladatokra egy-egy projekt keretei között.

4. A cégen belül saját adatbányászati csoportot hozhatnak létre.

Mind a négy verzióknak vannak pozitívumai és negatívumai. Mindig az adott körülményektől függ, hogy melyiket, vagy melyek kombinációit célszerű egy cégnek választania. A különböző felhasználók különböző területeken várnak választ és hasznot az adatbányászat alkalmazásából:

Kiskereskedők (retail szektor)

- A vásárlói szokások és preferenciák megismerése piaci vagy vásárlói kosárelemzéssel.
- A tisztességtelen fizetési magatartás felderítése.

Direktmarketing / Telemarketing

- Nagy megtakarítások pontosan megcélzott ügyfelek révén.

Gyártók

- A termelési folyamat ellenőrzése és ütemezése.

Légitársaságok

- Felfedezni az ügyfelek igényeit a stratégiai változtatásokhoz (pl. új útvonalak felkínálása vagy a szolgáltatások bővítése, stb.)

Telekommunikációs társaságok

- Annak „kitalálása”, milyen szolgáltatások lennének népszerűek az ügyfelek között
- A hívási csalások felfedezése szokatlan minták felderítésével.

A pénzügyi felhasználás területén kiemelhetők:**Bankok**

- Célzott marketing
- Kölcsön visszafizetés teljesítés előrejelzése

A nagy bankok (pl. Bank of American, Dresdner Bank) általában saját adatbányász munkatársakkal rendelkeznek. A kisebb bankok, amelyek korlátozott erőforrásokkal és technológiával rendelkeznek, többnyire megrendelik az adatbányászatot, illetve az adattárház-létrehozást illetve működtetést.

Biztosítótársaságok

- Az ügyfelek jobb megismerése.
- A biztosítási csalások hatékonyabb felderítése.

Pénzügyi szolgáltatások

Értékpapír-elemzők intenzíven használják nagytömegű pénzügyi adat elemzésére kereskedési és kockázati modellek létrehozása céljából, befektetési stratégiák kialakításához az adatbányászatot. Ez a felhasználási terület nagyon széles, és olyan alkalmazásokat foglal magába, mint:

- a devizakereskedelem,
- a részvény-kiválasztás,
- a jelzálog kiválasztás.

Hitelkártya kibocsátó társaságok (pl.: American Express, Citibank)

- Hitelkártya igénylések jóváhagyása.
- Vásárlásokat jóváhagyó döntések.
- A kártyatulajdonosok vásárlási szokásainak elemzése.
- Csalásfelderítés.

Pénzügyi felügyelet és szabályozó hatóságok (pl: NASD, Internal Revenue Service, FBI)

- Pénzmosás felderítése.
- Belső kereskedelemre utaló minták felfedezése.
- A piac szabályainak megsértése.

Az adatbányászat eredményeinek *vizuális* megjelenítése elsősorban a real-time (valós idejű) pénzügyi alkalmazásokban, például a pénzügyi, tőzsdei adatbányászatnál különösen fontos lehet. Ennek oka a

nagy mennyiségű információ

- Fundamentális adatok, jövedelmek, becslések, vetítések, osztalékok, könyvszerinti érték stb.
- Technikai jellemzők, pl. részvények teljesítménye, kvantitatív elemzések az ázsiai derivatív piacokra.
- Különböző elemzők és brókerek véleménye stb.

a döntési elemek nagy száma, amelyeket az egyes értékpapírok esetében figyelembe kell venni, ami azt eredményezi, hogy nehéz az *anomáliák azonosítása és kiküszöbölése* az egyes piacokon.

A különböző információk egyszerű integrálása nem jelent segítséget a nagytömegű adat miatt. Az adatok tradicionális ábrázolása (táblázatkezelő oldalak ipari csoportok és országok szerint, oszlop-, kör- vagy akár 3D-s diagramok) nem tudják az adatokat megfelelően integrálni és prezentálni. A megoldás a *valós idejű vizuális adatbányászat*: különböző színek, formák, méretek, prezentációs stílus (pl. villogás és forgás) segítségével hajtják végre az adatelemek vizuális megjelenítését oly módon, hogy ezeket a megjelenítési módokat rendelik az egyedi adatelemekhez azok fontos numerikus értékeinek vagy nem-numerikus tulajdonságainak leírásához. (Például egy nyíl megjelenítése jelölheti minden egyedi adatelemhez a jelenlegi feltételek egy halmazát.) Ilyen vizuális adatbányászati megoldás például a NASD Advanced Detection System (ADS), amelynek célja a piacok integritásának védelme a belterjes kereskedelemtől, és más szabályok megsértésétől. Az eszköz, amellyel ezt az adatbányászatot és a vizualizációt végzik a Metaphor Mixer (MM) a Maxus Systems International Inc. terméke. (Metaphor Mixer)

Összefoglalásként elmondhatjuk, hogy az adatbányászati funkció hatékony infrastruktúrát biztosít az egész vállalatra kiterjedő, magas szintű üzleti adatelemző alkalmazások létrehozásához. Az új üzleti ismeretek kinyerésének és terjesztésének automatizálásával és integrálásával a vállalatok ki tudják használni nagytömegű adataikban fekvő információkat, hatékonyabban működhetnek, és nagyobb versenyelőnyhöz juthatnak (Gáspár, 2006).

IRODALOM

Gáspár Bencéné dr. Vér Katalin (2006): *Üzleti intelligencia rendszerek*. Budapest. ÁVF.

Metaphor Mixer (MM) by Maxus Systems International Inc.: <http://www.maxusystems.com/>

Vázsonyi Miklós (2006): <http://www.szovegbanyaszat.hu/Szovegbanyaszat>

<http://hu.wikipedia.org/wiki/Adatbanyaszat>

<http://www.datamining.hu/>

http://www.eflow.hu/Tudasmenedzsment/Intelligens_portalok/Adatbanyaszat/adatbanyaszat.html

