

nye szerint is szükséges meghatározni. Lényegében katalogizálási kérdésről van szó, olyan igénnyel, hogy a leirási szabályok elég rugalmasak és redundánsak legyenek az egyedi gyűjtési helyek kívánalmaihoz képest.

Másfelől minden dokumentumtipust jellemez informatív értéke, azaz a benne foglalt információtipusok (pl. statisztikai adatok, műszaki paraméterek, cégnevek stb.) amelyek elemzése szükséges a tájékoztatás számára, illetve amelyek a tároló-visszakereső rendszer ismérrendszerét adják.

Általánosabb szinten egy-egy dokumentumtípus problémái – a szakirodalmi tájékoztatás szempontjából – ezekben a kérdésekben csucsosodnak ki. Bár Hegyi Nándor és Almásy László könyve ezekkel a kérdésekkel is foglalkozik, de valójában ezen a téren hiányérzetünk marad. Ennek leíró módszerük az oka. Szeretném megismételni, amit fentebb mondtam: az absztraktabb tárgyalásmódot számon lehet-e kérni a szerzőktől? Ez volt-e a feladatuk? Ha a könyvtári dokumentációs szakirodalomnak a feldolgozásra vonatkozó részét tekintjük, megállapíthatjuk, hogy ez a fajta tárgyalás még a klasszikus könyvtári dokumentumok esetében is hiányzik. A katalogizálásnak még a könyvek, periodikumok esetében sincs komolyan kimunkált elmélete. Olyan elmélete, amely az információfeldolgozás számára produkálná azokat a tipikus adatstruktúrákat, modelleket, amelyeket általánosított formában lehetne a továbbiakban kezelni és tárgyalni. Ugyanigy hiányosak azok a munkák, amelyek a visszakeresés kérdéseit szűkebb empiria szintjén túl tárgyalják. Miért lenne éppen a vállalati irodalom kivétel, s miért kellene éppen a legfiatalabb dokumentumtípus esetében ezt megkövetelni?

A szerzőktől talán tulzott követelés lenne mindezt elvárni. S hogy mégis szövétehető, nem ennek a konkrét könyvnek hibája, hanem az egész magyar könyvtári szakirodalomé.

HORVÁTH Tibor

The SMART Retrieval System. Experiment in Automatic Document Processing. E. Gerard Salton. New Jersey, 1971, Prentice-Hall, Inc.

Az automatizált SMART dokumentumvisszakereső rendszert a Harvard Egyetemen 1961 és 1964 közt tervezték meg. A rendszer a dokumentumokat ill. a keresési igényeket természetes nyelven fogadja el, a szövegek tartalmi elemzését teljesen automatikusan elvégzi, az elemzett dokumentumokat az elemzett kérdésekkel összeveti és kikeresi azokat a tárolt dokumentumokat, amelyek leginkább hasonlítanak a feltett kérdésre.

Más számítógépes visszakereső rendszerektől eltérően a SMART nem hagyatkozik a tárgyszavak vagy indexkifejezések manuális meghatározására, ugyanakkor nem alkalmazza kizárólagosan bizonyos szavak vagy kifejezések szövegbeni előfordulási gyakoriságát, hanem különféle segédeszközöket vesz igénybe a nyelvi elemzés elvégzése során.

Kiválóan alkalmas a SMART rendszer kísérletek elvégzésére. Ezt a lehetőséget számos alkalommal használták ki, s az eredmények jórésze az "Information Storage and Retrieval" report-sorozat köteteiben meg is jelent. Ezek (egybeszerkesztve, frissebb eredményekkel kiegészítve) alkotják e kötet gerincét.

## A SMART PROJEKT - HELYZETJELENTÉS ÉS TERVEK

A SMART kísérletek során különféle dokumentumgyűjteményeket használtak fel. Jónéhány kísérletet párhuzamosan valamennyire elvégeztek s a legfontosabb megállapítások a gyűjteménytől függetlennek bizonyultak. Ugy tűnik, hogy néhány eredmény minden műszaki jellegű dokumentumvisszakereső rendszerre igaz lehet. Ilyenek:

- a) Automatikus tartalom-elemzéshez a referátum mindig hatékonyabb, mint a cím; a teljes szöveg viszont a referátumnál nem annyival hatékonyabb, hogy javasolni lehetne feltétlen előnyben részesítését.
- b) A sulyozott tartalmi azonosítók hatékonyabbak a sulyozatlanoknál; a dokumentum és a keresőkép hasonlóságának mérése a korrelációs függvény jobb, mint az átfedés vizsgálata.
- c) A szókép normalizálása hatékonyabb, ha az alkalmazott szótár redundáns és nem túlzottan műszaki.
- d) A szinonimák felismerését biztosító szótár használata jelentősen emeli a visszakeresés hatékonyságát.
- e) A frázis-generáló módszerek (akár szótár, akár statisztikai asszociáció alapján) javítják a hatékonyságot, de a nyereség csak ritkán jelentős.
- f) A nagymértékben hierarchikus eljárások nem ígérnek jelentős eredményeket.
- g) Az automatikus szövegfeldolgozás nem lényegesen gyengébb a manuális indexelésnél; nagy és heterogén gyűjtemények esetén határozottan előnyösebb.
- h) A kipróbált eljárások minőségi sorrendje (nagyjából): referátumfeldolgozás frázis és szinonima felismeréssel; sulyozott szótó-összevetés és statisztikai asszociáció referátumok alapján; logikai szótó-összevetés sulyozás nélkül; címfeldolgozás és összevetés átfedés alapján.
- i) A felhasználók felőli visszacsatolás különösen ismétlődő keresések esetén igen hasznos.
- j) A dokumentumok szétválogatásán alapuló cluster-eljárás lényegesen csökkenti a keresési időt és alig rontja a hatékonyságot.

A SMART rendszer továbbfejlesztésére vonatkozó tervek három fő csoportra oszthatók:

A. Szövegfeldolgozási kísérletek, ezen belül: automatizált szótárszerkesztés; bibliográfiai hivatkozások felhasználása a tartalom azonosítására; az eredmények alkalmazhatósága nagyobb gyűjteményekre; összehasonlítás működő rendszerekkel; összehasonlítás klasszikus tárgyi indexekkel; idegennyelvű in-put problémája.

B. Visszacatolási és cluster-technikai kísérletek; optimális stratégia meghatározása a clusterek előállítására ill. a keresésre; visszacsatolási módszerek alkalmazása.

C. Real-time üzemmódu rendszerek; az on-line ("párbeszéd") működésű rendszerre való áttérés.

## A SMART RENDSZER ALKALMAZÁSA

A gyakorlati alkalmazás során a SMART visszakereső rendszert öt alapegységre bontották: 1. nyomtatott szövegek bevitele; 2. dokumentumok csoportosítása keresési célokra: cluster-képzés; 3. a keresendő dokumentumok kiválasztása; 4. keresés a dokumentumcsoportokon belül; 5. a keresés kiértékelése. Az egyes tevékenységekkel a kódot tanulmányai részletesen foglalkoznak.

## KIÉRTÉKELÉSI PARAMÉTEREK

Az automatikus indexelést alkalmazó dokumentumvisszakereső rendszerek a következő elemekkel jellemezhetők: természetes nyelvű dokumentumok halmaza; természetes nyelvű kérdések halmaza; indexnyelv; transzformáció a dokumentumok nyelvről az indexnyelvre; transzformáció a kérdések nyelvről az indexnyelvre; kereső függvény. Ezek alapján felépíthető a rendszer modellje. A kiértékelési mutatók elmélete feltételezi ezek mellett a visszakeresett dokumentumok relevancia szerinti – legalább részleges – rendezését.

A visszakereső rendszer teljesítménye kiértékelésének elvi alapja egy összehasonlítás: a teljes gyűjteményből, adott kérdésre, mely tételeket keres ki a rendszer és melyeket keres(ne) ki a kérdést feltevő használó. A relevancia különféle fokai közti megkülönböztetés problémája a visszakeresett tételek rangsorolásával kerülhető meg. Még így is két különböző, realizálható értékelési szempont érvényesülhet: vagy a visszakereső kérdéseket tekintjük elemeiknek (makroértékelés), vagy a releváns dokumentum és az irreleváns dokumentumok halmaza kapcsolatát (mikroértékelés). Az első kérdés-orientált, a második dokumentum-orientált szempont. A kétféle megközelítés összehasonlítására a visszahíváspontosság görbét alkalmazták.

A teljesítménymutatók céljaik szerint három csoportba oszthatók: adott kísérleti szituáción belüli összehasonlítást szolgáló "belső" mutató; a kísérleti eredmények működő rendszerekre való kiterjesztésének mérőeszközei. Emellett alkalmazhatók az értékelés eltérő szempontjai, s figyelembe veendők a megkívánt tulajdonságok. A rangsoroló rendszerekhez kialakíthatók egy-adatos mutatók (log-pontosság, rang-visszahívás, normalizált visszahívás, normalizált pontosság), melyek függetlenek a rangsoron belül a releváns-irreleváns küszöb meghatározásától. Megadhatók a küszöbtől függő teljesítménygörbék is. A kétféle (makro és mikro) átlagolási eljárás eredményei a pontosságvisszahívás görbével vethetők össze, s az eredmények kiterjeszthetők eltérő általánosságú gyűjteményekre ill. a keresés során a clustertechnika alkalmazására is.

A nyelvi elemzés célja a természetes nyelvű dokumentum vagy kereső kérdés számára tartalmi azonosítók meghatározása automatikus módszerekkel. Ennek során számos alapprobléma merül fel: a csak szintaktikai funkciót hordozó szavakat törölni kell; a szinonimákat fel kell ismerni; a hononimák eltérő jelentését azonosítani kell; a szintaktikai parafrázisokat szintén; lehetőleg figyelembe kell venni a szövegen belüli indirekt utalásokat; fel kell tárnai a szavak közti implicit kapcsolatokat; vigyázni kell a szavak jelentésének esetleges időbeli változására. Az elemzés elvégzéséhez négyféle szótár látszik szükségesnek: negatív szótár; tézaurusz; frázis-szótár; a kifejezések hierarchikus (fa-típusu) elrendezése. A tézauruszokon belül is több típust különböztetünk meg: szinonima-szótár és végződés-jegyzék; frázis-szótárak; fogalmi hierarchia. A tézaurusz szerkesztésében az automatizálás különféle mértékben érvényesíthető.

A SMART rendszer a következő, a dokumentum-elemzés szempontjából fontos lehetőségekkel rendelkezik: szótövek meghatározása; fogalmak azonosítása számmal; fogalmak hierarchikus elrendezése; statisztikai asszociációs módszerek; szintaktikai elemzés; statisztikai frázis-felismerés; kérdés-dokumentum összevetés. A rendszer biztosítja a különféle feldolgozási eljárások értékelését és összehasonlítását. Ennek során elsősorban a következő problémákat vizsgálták: a rendelkezésre álló dokumentum (szöveg) hossza; összevetés-függvények és súlyozott kifejezések; hierarchikus kiterjesztés; manuális indexelés.

Néhány éve megindultak a kísérletek idegen (nem-angol) nyelvű dokumentumok automatikus feldolgozására is. A legtöbb eredmény eddig a német nyelvvel kapcsolatban született. Próbálkoztak többnyelvű tézaurusz kialakításával is. Egy vegyes nyelvű dokumentumgyűjteményen lefolytatott kísérletsozrot alapján úgy tűnik, hogy a SMART módszerei más nyelvekre is alkalmazhatók.

## CLUSTER-GENERÁLÁS ÉS KERESÉS

A visszakeresési hatékonyság maximalizálása helyett, a felmerülő költségek miatt, az optimális megoldás látszik célszerűnek. Ez pedig, főleg nagy és viszonylag heterogén adattár esetén, nem a teljes adattár, hanem csak annak bizonyos részei keresését teszi indokolttá. A SMART rendszer ezt a következőképpen oldja meg:

- a) Minden egyes dokumentum tartalmi azonosítóját az összes többivel összevetve a hasonló azonosítóju (tartalmu) dokumentumokat csoportokba, clusterekbe vonja össze.
- b) Minden clusterhez meghatároz egy reprezentatív elemet, az un. centroid vektort.
- c) A keresést két lépésben végzi: először a keresőképet összeveti valamennyi centroid vektorral, majd a keresést már csak azokban a clusterekben végzi el, amelyek centroidja hasonlít a keresőképhez.

A cluster-képzéshez megfelelő algoritmus áll rendelkezésre, amely az eljárást automatizálhatóvá teszi. Az algoritmus kiértékeléséhez szükséges muta-

tók figyelembe veszik az "átnézett" clusterok számát, ez ugyanis befolyásolja a rendszer teljesítményét.

Egy másik (gyors) algoritmus a gyűjtemény tételeinek automatikus osztályozását célozza. Ezt két gyűjteményen is kipróbálták. Az eredmények kiértékelésére statisztikai szignifikancia-próbát alkalmaztak.

A cluster-képzés egy további módszere volna a keresőképből és az ehhez más módon visszakeresett dokumentumokból való kiindulás. Ezt is alkalmazták kísérletileg, és így összehasonlítást tehettek a clusteres ill. a teljeskörű keresés teljesítménye között.

## VISSZACSATOLÁS

Dokumentum-visszakereső rendszer hatékonyságának értékeléséhez el kell különíteni azokat a változókat, amelyek a rendszer viselkedését befolyásolják. Az e célból felállított modell a visszakereső rendszer három alapfunkcióját határozta meg: indexelés, a keresőkép megfogalmazása és a keresőkép-dokumentum összevetés. Ezek közül a legnagyobb változatosság a keresőkép megfogalmazásában lehet (ezt ugyanis a használó végzi). Éppen ez adja a visszakeresés eredményére vonatkozó használói véleményt a visszacsatolás jelentőségét a keresőkép folyamatos javításának folyamatában. E módosítási eljárás algoritmizálása és automatizálása szintén megtörtént. A relevancia-visszacsatolás egyike a leghatékonyabb módszereknek, amelyeknek célja, hogy a felhasználót bevonja az információkeresés folyamatába.

A megfelelően felhasznált visszacsatolás egyuttal alapja lehet a felhasználó és a gép közti párbeszédnek, az interaktív üzem módnak, amelynek során a visszakeresést a használó folyamatosan ellenőrzi és irányítja. A visszacsatolást alkalmazó kereső eljárások hatékonyságának mérésére szolgáló mutatórendszer kidolgozták és kipróbálták.

Ez a lehetőség annál is inkább jelentős, mert az automatikus vagy fél-automatikus információkereső rendszerek jelenlegi vagy a közeli jövőben várható teljesítménye lényegesen kisebb, mint azt sok potenciális használó képzeletében. Úgy tűnik, hogy a nyelvi elemzés vagy az adatszerkeztetés terén várható változások sem idéznek majd elő gyökeres javulást. A legjelentősebb mérvű fejlődés az interaktív kereső technika alkalmazásától várható, ennek során ugyanis a gép szinte "kikényszeríti" a használóból igényeinek egyre pontosabb megfogalmazását. Az interaktív keresés jól párosítható a jelenlegi, viszonylag merev file-szerkeztetés helyett egy dinamikus file alkalmazásával.

Érdekes kísérletet folytattak a visszacsatolás egy másik alkalmazási lehetőségére vonatkozólag is: ez nem egyszerűen módosítja a keresőképet a visszajelzés eredménye alapján, hanem attól függően két vagy több új keresőképet ad meg, melyek kissé különböznek az eredetitől, s a használó ezek közül választ. A kísérletek bizonyos fokú teljesítmény-növekedésre mutatnak.

Egy további kísérleti eljárás a negatív relevancia visszacsatoláson alapul, amely akkor is működik, ha az eredeti keresőkép egyetlen releváns dokumentumot sem eredményezett.

Biztató eredményeket mutat az a kísérlet is, amely a visszacsatolás alapján elhatárolt releváns ill. irreleváns dokumentumok közti szignifikáns hasonlóságból indul ki, s a megismételt-módosított keresőképet ennek alapján állítja elő. Az eljárás egyelőre további kipróbálás alatt áll.

A visszacsatolás alkalmazása nem szükségképpen merül ki a keresőkép iteratív módosításában. Alternatív lehetőség az indexelt dokumentumbázis (un. dokumentumtér) ill. a dokumentumok tartalmi leírása (un. dokumentumvektor) módosítása, transzformációja is. A SMART rendszerben mindkettőre fejlesztettek ki algoritmust s alkalmazták kísérleti körülmények közt. Mindkét esetben a folyó, ismételt keresés mellett biztosították a korábbi keresési eredmények felhasználásának lehetőségét is.

## MŰKÖDÉSI JELLEMZŐK

A SMART-kísérletsorozat – elsősorban távlati alkalmazás szempontjából – jelentős része az üzemszerűen működő rendszerek tanulmányozása és összevetése a SMART eredményeivel, tapasztalataival. Ilyen vizsgálat volt az interaktív keresés és az automatikus információ-kijelzést tartalmazó módszerek vizsgálata, a relevancia-ítéletek és az információkereső rendszerek teljesítmény-értékelése közti kapcsolat elemzése, végül a manuális ill. az automatikus indexelési módszerek összehasonlítása. A SMART-rendszer alapelveinek, módszereinek gyakorlati alkalmazását azonban még jónéhány vizsgálatnak kell megelőznie, főleg a felmerülő költségek és a teljesítmény viszonyát illetően.

SÁRDY Péter

Fejér Megyei Könyvtáros. Fejér megye könyvtárosainak tájékoztatója. Kiadja a Vörösmarty Mihály Megyei Könyvtár, Székesfehérvár.

Könyvtári szakirodalmunk egyik mostohagyereke a (ma már pontatlanul) "hálózati híradó"-nak nevezett műfaj. Bevallottan "helyi érdekűnek" kétszül, s ezt a könyvtáros-kollegák annyira komolyan veszik, hogy tudomásom szerint egyetlen intézményben sem található meg ezeknek a kiadványoknak többé-kevésbé teljes sorozata. Természetesen azt, hogy e szigorú elbírálás igazságos-e, csak a híradók elemzése döntheti el.

Az un. hálózati híradók nagyjából két csoportra oszthatók: a területi és a szakterületi hatókörűekre. Az előbbinek klasszikus példája a megyei híradó, az utóbbié a szakmai (műszaki, mezőgazdasági, orvostudományi, stb.) hálózat híradói. Miután pedig az élet gyakran produkál nehezen beskatulyázható jelenségeket, itt is találunk átmeneti kategóriákat. Példaként elég a hálózatnyi létszámmal működő nemzeti könyvtárunk "házi" (de nemcsak házon belül érdekes) híradójára vagy egyes egyetemi könyvtáraink hasonló kiadványaira utalni (pl. "Műszaki Egyetemi Könyvtáros"), mely utóbbiak a központi könyvtáron kívül a tanszéki, kari és egyéb hálózati könyvtárak problémáit is műsorra tűzik. A szakmai hálózati híradók a nagy szakkönyvtári-információs szervezetek fejlő-