

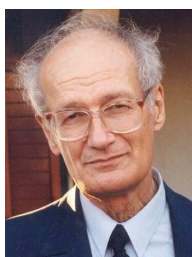
Bogdan Ionescu

Ionuț Mironică

**Conceptul de Indexare Automată după
Conținut în Contextul Datelor Multimedia**

București, 2013

Prefață



Ce vrea să zică asta - indexarea după conținut - cititorul va găsi în primul capitol, dar sunt tentat să zic și aici, în aceste rânduri, câteva cuvinte: problema nu e chiar nouă. Cu ceva zeci de ani în urmă am aflat că pe alte meleaguri oamenii se ocupau, pentru cuvinte, cu alcătuirea unor asemenea dicționare. Cele alfabetice, pe care le avem și noi, îți explică ce vrea să zică un cuvânt pe care îl ai dar al cărui sens nu îl știi; dar sunt și probleme de alt fel: acolo era un exemplu de întâmplare în academia spaniolă - un vorbitor nu-și aducea aminte cum se cheamă un om născut pe vapor (noi n-avem cuvânt pentru acest concept). Ne trebuie dicționare care să ne ducă de la concept la cuvânt. Despre unele popoare primitive se zice că aveau zeci de cuvinte pentru a denumi diferite tipuri de nori; noi n-avem, dar am putea eventual descrie formele lor, mișcarea lor, ca să precizăm la care ne referim când vrem să povestim o întâmplare concretă.

Într-o bibliotecă de un miliard de cărți, cu câte 500 de pagini fiecare și cu 2.000 de semne pe pagină avem nevoie doar de 50 de cifre binare pentru a identifica orice literă, ceea ce mi se pare extrem de puțin - la îndemâna umanului: le cuprindem cu ochiul dintr-o privire, pe un rând. Oare nu e posibil să avem căi/o cale de a ajunge la "obiectul" dorit dintr-o colecție vastă, cunoscându-l prin calitățile sale (făcute cumva măsurabile: da-nu, roșu-albastru-galben-verde, o valoare întreagă între 1 și 100, 17 grade de turtire a unui cerc în elipsă, etc.)? "Obiectele" de care vorbeam pot fi entități foarte complexe: o imagine, o secvență de film mut, entități "multimodale" (vorba, sunete, imagini, text, etc.). Parcă suntem tentați a zice da. Dar acum vine partea dificilă a problemei, și în același timp frumoasă prin efortul de

creație pe care ni-l cere (aspectul care ne provoacă, ne desfide, englezul ar zice "challenging"): pe de o parte, în cazul concret al unei colecții de un tip dat (de pietre, de găze, de filme), care sunt atributele, cum le definim ca să caracterizăm cât mai compact și mai corect, acea colecție; pe de altă parte, în fața unui obiect din colecție, cum măsurăm *automat*, adică nu prin intervenția omului (în cazul ăsta avem nevoie de un specialist în domeniu!), aceste atribute.

Fără acest mic amănunt aici, "automat", suntem pierduți fiindcă operația manuală de adnotare cu atribute a obiectelor este consumatoare de timp în așa măsură că ne face întreprinderea lipsită de sens.

În momentul de față al scurtei noastre istorii de câteva sute de ani, suntem în pericol de a fi "înecați în informații" care pe de o parte multe ne sunt vitale și pe de alta, în ansamblul lor ne copleșesc, fără a putea ajunge la cele de care avem nevoie suntem ca însetatul din pustiu peste care năvălește marea. Indexarea automată după conținut ne poate salva.

Extras din prefața cărții "Analiza și Prelucrarea Secvențelor Video: Indexarea Automată după Conținut", Editura Tehnică București, 2009.

Prof. univ. dr. ing. Vasile BUZULOIU (1938 - 2012)
București 17 Noiembrie 2008

Cuvântul autorului

Indexarea automată după conținut a datelor este un domeniu de actualitate ce câștigă din ce în ce mai mult teren datorită necesității crescânde de exploatare a volumelor mari de date multimedia.

Progresul tehnologic al dispozitivelor de achiziție și prelucrare a datelor (terminale mobile, sisteme de calcul, medii de stocare, dispozitive de redare și captură audio-video) cât și a infrastructurii de transmisie de date (protocoale de transmisie fără fir: WiFi, Bluetooth, rețele LAN de mare viteză, telefonie multimedia 3G și 4G) au condus practic la simplificarea stocării, transmisiunii și prelucrării volumului important de date specific multimedia (video, imagini, sunet, text).

Mărturie în acest sens este răspândirea Internet-ului în tot mai multe medii sociale și posibilitatea de accesare a acestuia de pe o categorie tot mai diversă de dispozitive electronice. La acestea se adaugă și succesul imens de care se bucură rețelele de socializare și platformele multimedia "on-line", Facebook, Twitter, LinkedIn, Google+, YouTube, Dailymotion, Picasa, Flickr sunt doar câteva exemple dintre acestea.

Dinamica partajării datelor pe Internet este una copleșitoare, aceasta realizându-se practic "în timp real" de pe orice terminal multimedia. Următoarele statistici sunt edificatoare în acest sens: în 2012 mai mult de 72 de ore video sunt încărcate în fiecare minut pe platforma YouTube, mai mult de 500 de ani de video de pe platforma YouTube sunt vizualizați zilnic de pe platforma de socializare Facebook, mai mult de 700 de înregistrări video de pe YouTube sunt partajate în fiecare minut pe rețeaua de socializare Twitter.

În societatea curentă, accesul la informația multimedia a devenit parte integrantă din viața noastră de zi cu zi. Problema cu care ne confruntăm nu este lipsa informației, ci imposibilitatea de a selecta dintr-un vast amalgam

de date, informațiile utile. Această problemă este cu atât mai dificilă cu cât conținutul acestor date a devenit din ce în ce mai complex.

Până nu demult, când făceam referire la informație multimedia ne adresam imaginilor, înregistrărilor audio sau eventual video. În prezent conceptul de multimedia vine să reunească toate aceste informații sub umbrela unei singure paradigme și anume aceea a reprezentării multimodale a informației. Datele multimedia sunt practic "metadate" ce reunesc orice tip de informație video, audio și textuală. Metodele de prelucrare trebuie să se adapteze acestor noi cerințe în care analiza de conținut este unitară și nu realizată independent pentru fiecare sursă de informații.

În acest context, lucrarea de față vine să realizeze o trecere în revistă a domeniului indexării automate după conținut a datelor multimedia și să discute soluțiile existente.

Lucrarea este structurată în felul următor. În primul rând este introdusă problematica indexării datelor și aplicațiile acesteia (Capitolul 1). Mai departe, este prezentat detaliat mecanismul de funcționare al unui sistem de indexare ce implică descrierea conținutului datelor, mecanismul de căutare a datelor și respectiv interacția cu utilizatorul (Capitolul 2). Capitolul 3 realizează o trecere în revistă a tehnicilor de descriere a conținutului datelor folosind informația vizuală, audio și respectiv textuală. Capitolul 4 se interesează de tehnicile de fuziune a informației specifice abordărilor multimodale ce exploatează date heterogene. Mai departe este adusă în discuție problema evaluării similarității datelor (Capitolul 5). Tehnicile de interacție cu utilizatorul în vederea îmbunătățirii performanțelor de indexare sunt prezentate în Capitolul 6. Capitolul 7 discută problematica vizualizării informației multimedia și în special a datelor video. În final, Capitolul 8 prezintă o serie de modalități de evaluare subiectivă și obiectivă a performanțelor unui sistem de indexare iar Capitolul 9 concluzionează lucrarea sintetizând paradigmele actuale ale sistemelor de indexare.

Lucrarea de față se dorește a fi un studiu introductiv al domeniului, furnizând cititorului o vedere de ansamblu asupra tehnicilor de prelucrare aferente sistemelor de indexare și a avantajelor și limitărilor acestora. Pentru o descriere detaliată, cititorul este îndrumat să consulte referințele bibliografice furnizate.

Ș.l. univ. dr. ing. Bogdan IONESCU
București 26 Aprilie 2013

Cuprins

1	Introducere	1
2	Mecanismul de indexare după conținut	7
2.1	Descrierea conținutului datelor	10
2.2	Formularea căutării	12
2.3	Căutarea datelor	14
2.4	Interacția cu utilizatorul	15
3	Descrierea conținutului multimodal	19
3.1	Informația vizuală	20
3.2	Informația audio	32
3.3	Informația textuală	34
3.4	Descriere semantică sau sintactică?	37
4	Fuziunea datelor	41
4.1	Metode de tip "early fusion"	41
4.2	Metode de tip "late fusion"	44
5	Conceptul de similaritate a datelor	49
5.1	Similaritatea descriptorilor	49
5.2	Similaritatea la nivel de structură	55
5.3	Similaritatea semantică	56
6	Tehnicile de tip "relevance feedback"	59
6.1	Algoritmul Rocchio	63
6.2	Estimarea importanței atributelor	64

<i>CUPRINS</i>	vi
6.3 Support Vector Machines	66
7 Vizualizarea conținutului multimedia	73
8 Evaluarea performanțelor indexării	79
8.1 Evaluarea subiectivă	79
8.2 Evaluarea obiectivă	85
8.2.1 Precision-Recall	86
8.2.2 F-measure	88
8.2.3 Curbă de precision-recall și ROC	89
8.2.4 Mean Average Precision	91
9 Paradigme ale indexării	93
Bibliografie	97

CAPITOLUL 1

Introducere

Dacă în urmă cu aproximativ un deceniu, cantitatea de informație multimedia disponibilă era una redusă, în zilele noastre putem vorbi despre o *explozie informațională*. Accesul la informația multimedia sau ”conținut”, fie că este vorba de imagini, sunet, text sau video, a devenit practic parte integrantă din viața noastră de zi cu zi. Evoluția tehnologică a dispozitivelor de achiziție și prelucrare a datelor (terminale mobile, sisteme de calcul, medii de stocare, dispozitive de redare și captură audio-video) cât și a infrastructurii de transmisie de date (protocoale de transmisie fără fir: WiFi, Bluetooth, rețele LAN de mare viteză, telefonia multimedia 3G și 4G) au dus la creșterea exponențială a volumului multimedia prin facilitarea stocării și prelucrării acestuia.

La acestea contribuie semnificativ și răspândirea Internet-ului în tot mai multe medii sociale precum și succesul imens de care se bucură rețelele de socializare ”on-line” (exemplu: Facebook¹, Twitter², LinkedIn³, Google+⁴) cât și platformele web multimedia (exemplu: YouTube⁵, Dailymotion⁶, Picasa⁷, Flickr⁸). Pe lângă producția de conținut multimedia să spunem comercial

¹<https://www.facebook.com>

²<https://twitter.com>

³<http://ro.linkedin.com>

⁴<https://plus.google.com>

⁵<https://www.youtube.com>

⁶<https://www.dailymotion.com>

⁷<http://picasa.google.com>

⁸<https://www.flickr.com>

(realizat de companii în vederea comercializării), accesul la rețele de socializare și platforme web a condus practic la facilitarea posibilității de a partaja și accesa date multimedia personale, generate de utilizatorii de rând, precum fotografii, filme din colecțiile personale, reportaje, "video blogging" și așa mai departe. Acestea reprezintă o sursă imensă de conținut multimedia, să luăm ca exemplu rețeaua de socializare Facebook care în 2012 însuma nu mai puțin de 1.2 miliarde de utilizatori ce partajează informații multimedia.

În prezent dinamica partajării datelor pe Internet este una copleșitoare aceasta realizându-se practic "în timp real" de pe orice terminal multimedia, atât mobil (de exemplu telefonul mobil) cât și fix. Prin simpla apăsare a unui buton, o înregistrare video sau imagine poate fi încărcată imediat "on-line". Următoarele statistici sunt edificatoare în acest sens: în 2012 mai mult de 72 de ore video sunt încărcate în fiecare minut pe platforma YouTube, mai mult de 500 de ani de video de pe platforma YouTube sunt vizualizați zilnic de pe platforma de socializare Facebook, mai mult de 700 de înregistrări video de pe YouTube sunt partajate în fiecare minut pe rețeaua de socializare Twitter. Dintre informațiile multimedia cel mai frecvent tranzacționate, conținutul video "on-line" reprezintă cea mai mare categorie de date vehiculate pe Internet, cuprinzând în 2012 26% din traficul total de date (sursa CISCO Systems⁹). Până în 2015 se estimează că mai mult de 1 milion de minute video (674 zile) vor traversa Internetul în fiecare secundă.

Astfel că problema cu care ne confruntăm acum nu este lipsa de informație, ci, dimpotrivă imposibilitatea de a selecționa din volumul informațional imens disponibil, *informația utilă* căutată. Am ajuns în punctul în care acest lucru nu mai poate fi realizat de operatori umani și este necesară preluarea acestei sarcini de către calculator.

Această problemă de cercetare se găsește la afluența unor domenii precum prelucrarea și analiza semnalelor ("signal processing"), vederii asistate de calculator ("computer vision") și al clasificării datelor ("data mining"). Importanța acestei direcții de cercetare a dat naștere unor domenii dedicate precum "multimedia" și al "căutării de informații" ("information retrieval"). Cercetările actuale vizează dezvoltarea de metode *automate* capabile să înțeleagă conținutul datelor și să îl pună la dispoziția utilizatorului într-un mod foarte apropiat de modul în care o persoană ar realiza acest lucru (apropiat de modul de percepție uman).

O potențială soluție la problema căutării informației multimedia a fost discutată cu mai mult timp în urmă în contextul căutării de imagini și consta în folosirea de *tehnici de indexare automată după conținut* [Smeulders 00]. Transpuse în contextul actual tehnologic, aceste tehnici trebuie acum să se

⁹<http://www.cisco.com>

adaptez, pe de-o parte unui volum imens de date, de exemplu în cazul video doar 1 minut este echivalentul a 1.500 de imagini statice și astfel o singură secvență video echivalează conținutul unei întregi colecții de imagini; cât și manipulării de conținut temporal, în mișcare (video) și multimodal (text-sunet-imagini). În ciuda unei disponibilități de putere de calcul în continuă creștere (în prezent un simplu telefon mobil folosește procesoare cu patru nuclee de prelucrare și frecvențe de 1.6 GHz) complexitatea acestei probleme necesită optimizarea și paralelizarea metodelor. Acestea trebuie să fie eficiente computațional pentru a putea fi aplicate la scară largă colecțiilor de pe Internet.

Cât de departe este tehnologia actuală pentru a realiza acest lucru? Să luăm ca exemplu cazul simplificat al căutării după conținut al imaginilor. În Figura 1.1 am prezentat rezultatele obținute pentru căutarea unor imagini ce conțin ”nuferi galbeni” folosind motorul de căutare propus de Google și anume Google Search by Image¹⁰ - considerat una dintre tehnologiile de vârf în prezent. Pentru a specifica datele dorite, am furnizat ca exemplu o imagine.



Figura 1.1: Exemplu de căutare după conținut pentru o imagine cu un nufăr galben (”water lily”, imagine stânga) folosind motorul de căutare Google Search by Image. Imaginile din dreapta reprezintă primele șase rezultate obținute în ordinea descrescătoare a similarității (ordine sus în jos și de la stânga la dreapta).

Se poate observa că în ciuda faptului că imaginile returnate au proprietăți vizuale similare cu imaginea dată drept referință, semnificația semantică a acestora poate fi complet diferită. De exemplu, primim ca rezultat alte tipuri de flori sau chiar o persoană cu un tricou având culori similare. Cu toate că sistemele de căutare după conținut a imaginilor au în acest moment aproape două decenii de existență, și aici ne referim nu la tehnologia în sine ci la sisteme funcționale, un exemplu în acest sens fiind sistemul Query By Image Content – QBIC propus de IBM în 1995 [Flickner 95], tehnologia actuală nu

¹⁰<http://images.google.com>

este încă capabilă să atingă un nivel apropiat de modul în care o persoană ar rezolva problema căutării, manual.

Tehnicile de căutare după conținut a informației video sunt și mai puțin dezvoltate în acest moment limitându-se în principal în a fi extensii temporale ale celor aplicate în cazul imaginilor statice (de exemplu pentru a lua în calcul dimensiunea temporală de mișcare). În prezent, nu există un sistem de căutare după conținut video disponibil public, încercările existente fiind doar experimentale, adaptate la baze video "off-line" de dimensiuni reduse (în cazul cel mai bun de sute de mii de secvențe) și limitate în a se adresa unor aplicații particulare (de exemplu căutarea de conținut de știri, sport, catalogarea colecțiilor de filme după gen, identificarea conținutului de animație și așa mai departe).

Platformele de căutare multimedia existente sunt limitate în a folosi doar informație textuală, precum descrierile asociate de către utilizatori datelor. De exemplu, o înregistrare cu turnul Eiffel poate fi însoțită de o descriere de genul "vizită turnul Eiffel, Paris 2013". Utilizatorul va căuta informația dorită furnizând tot o descriere textuală a acesteia, ca de exemplu caută toate înregistrările cu "turnul Eiffel", furnizând aceste cuvinte cheie. Informația furnizată va fi comparată cu cea asociată datelor obținând ca rezultat secvențele corespunzătoare, precum secvența etichetată anterior. Aparent problema pare a fi rezolvată. Totuși, informația textuală este limitată în a furniza doar o descriere globală și parțială a conținutului. În exemplul anterior, sistemul pe baza descrierilor existente nu va fi capabil de exemplu să identifice prezenta unei anumite persoane în acea înregistrare deoarece această informație lipsește din descriere. Mai mult, descrierile textuale nu pot fi determinate în mod automat, necesitând intervenția umană. Extrapolând această problemă la dimensiunea bazelor multimedia de pe Internet, asocierea de descrieri textuale care să detalieze conținutul datelor video devine practic imposibilă.

Soluția la problema căutării după conținut a datelor multimedia nu se găsește la nivel de modalitate individuală și anume la nivel de imagine, video, sunet sau chiar text. Soluția ține de o abordare globală interdisciplinară a acestei problematice prin interacționarea informațiilor multimodale extrase din toate sursele de informație disponibile, de la culoare, textură, forme, mișcare, informație temporală până la sunet, voce, text și așa mai departe. Aceasta constituie de fapt tendița actuală de cercetare. Folosirea independentă a surselor de informație se dovedește inefficientă pentru a rezolva o problemă atât de complexă precum înțelegerea automată a conținutului datelor multimedia. Ca referință în acest sens sunt campaniile TRECVID Video

Retrieval Evaluation Benchmarking Campaign¹¹, MediaEval Benchmarking Initiative for Multimedia Evaluation¹², ImageCLEF The CLEF Cross Language Image Retrieval Track¹³ sau PASCAL Challenge - Pattern Analysis, Statistical Modelling and Computational Learning¹⁴ ce anual prezintă tehnologiile și bunele practici curente din domeniu. Cititorul se poate raporta la acestea pentru o vedere de ansamblu a progresului tehnologic actual în acest domeniu.

În cele ce urmează vom face o trecere în revistă a tehnicilor ce stau la baza procesului de indexare după conținut, a tehnicilor de descriere a conținutului datelor și a surselor informaționale exploatare, a tehnicilor de fuziune a informațiilor multimodale, tehnicilor de integrare a opiniei utilizatorului în procesul de indexare, a problematicii vizualizării conținutului multimedia, a modului de evaluare al performanțelor unui sistem de indexare încheind cu prezentarea barierelor actuale ale sistemelor de indexare după conținut.

¹¹<http://trecvid.nist.gov>

¹²<http://www.multimediaeval.org>

¹³<http://www.imageclef.org>

¹⁴<http://pascallin.ecs.soton.ac.uk/challenges/VOC>

CAPITOLUL 2

Mecanismul de indexare după conținut

Conceptul de indexare folosit pentru căutarea datelor este definit ca fiind *procesul de adnotare* a informației existente într-o colecție de date, prin adăugarea de informații suplimentare relative la conținutul acesteia, informații numite și *indici* de conținut [Kyungpook 06]. Această etapă este necesară accesării colecției de date, deoarece permite catalogarea automată în funcție de conținut a datelor.

Într-o colecție de date suficient de vastă, putem spune că datele care nu au fost adnotate sunt practic inexistente pentru utilizator. Un exemplu simplu de sistem de indexare este însuși sistemul de fișiere al oricărui calculator personal. Acesta ne furnizează datele aflate pe diversele medii de stocare (disc dur, memorie externă, etc.) sub formă de fișiere ce sunt indexate după informații precum nume, extensie, dată, și așa mai departe. Să ne imaginăm situația în care un fișier a fost omis din această listă de indici, cu toate că el este prezent fizic pe suportul de stocare, acesta va fi invizibil și inaccesibil pentru utilizatorul de rând.

Procesul de adnotare a datelor este văzut din două perspective: pe de-o parte există *adnotarea manuală*, iar pe de altă parte *adnotarea automată*. Gradul de complexitate al adnotării este direct proporțional cu nivelul de detaliu dorit pentru accesarea datelor. Dacă se dorește ca utilizatorul să poată accesa datele folosind criterii mai complexe, ca de exemplu căutarea unei anumite secvențe video pentru care nu se cunoaște nici numele, nici extensia fișierului, dar totuși utilizatorul dispune de informații referitoare la conținutul vizual al acesteia, în această situație, procesul de indexare va fi mult mai complex, necesitând înțelegerea de către calculator a conținutului

datelor.

Astfel, în cazul unei indexări după criterii complexe de conținut, adnotarea manuală este foarte dificil de realizat, deoarece necesită un număr important de operatori umani. Aceștia ar trebui să ”răsfoiască” manual întregul conținut al bazei de date pentru definirea indicilor de conținut. Luând în calcul faptul că o astfel de colecție de date este în prezent practic nelimitată (exemplul sunt colecțiile de pe Internet), indexarea manuală devine imposibilă. În acest moment, cercetările existente în domeniu se focalizează pe dezvoltarea de algoritmi de adnotare automată a conținutului, mai ales în cazul datelor ce necesită un timp important pentru vizualizare, ca de exemplu documentele video.

Cu toate că adnotarea conținutului datelor este soluția optimală pentru a accesa informația utilă dintr-o vastă colecție de date, aceasta nu este și suficientă. Adnotarea în sine nu oferă decât o serie de date suplimentare, putem spune, de nivel semantic inferior (”low-level”), care deseori sunt inaccesibile utilizatorului neavizat. Pentru a accesa baza de date, utilizatorul trebuie să dispună de o modalitate prin care să poată accesa sau vizualiza ușor datele, fie pe baza indicilor, fie în mod direct. Aceasta trebuie să aibă o funcționalitate naturală și intuitivă. Sistemul care permite utilizatorului să vizualizeze conținutul bazei de date poartă numele de *sistem de navigare*.

Pe de altă parte, accesul la date presupune un proces de căutare. Utilizatorul trebuie să mai dispună, pe lângă sistemul de navigare, de un mecanism care să-i permită căutarea informațiilor dorite în baza de date. Căutarea se realizează prin formularea de cereri de căutare sau ”queries”. Pentru ușurință, o astfel de cerere trebuie să fie exprimată într-un limbaj natural, apropiat de limbajul uman, cum ar fi de exemplu ”caută filmele de acțiune” sau ”caută imaginile ce conțin peisaje”. Sistemul care răspunde acestor cerințe poartă numele de *sistem de căutare*. Figura 2.1 sintetizează aceste aspecte prezentând schematic modul de funcționare al unui sistem generic de indexare a datelor.

Astfel, pentru a sintetiza, mecanismul de indexare și căutare a datelor presupune realizarea următoarelor etape:

- **descrierea conținutului datelor:** într-o primă etapă, informația propriu-zisă din baza de date este reprezentată prin intermediul atributelor de conținut, informații pe baza cărora se realizează întregul proces de indexare (vezi Secțiunea 2.1);
- **formularea cererii de căutare:** utilizatorul furnizează o descriere a datelor pe care dorește să le găsească prin formularea unei cereri de căutare (”query”). Acest lucru poate fi realizat folosind un exemplu

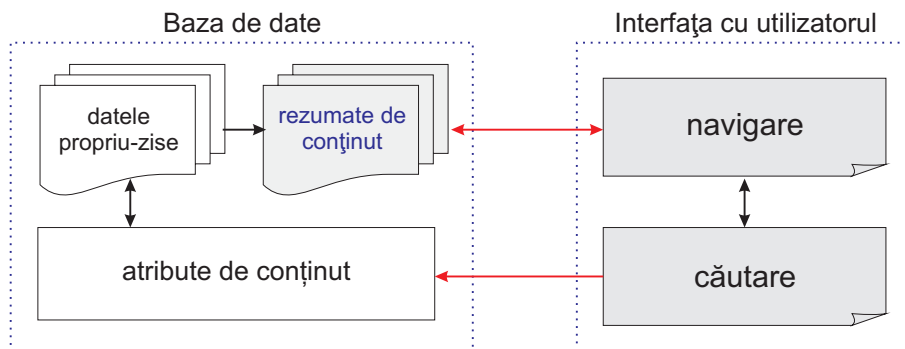


Figura 2.1: Principiul de funcționare al unui sistem de indexare după conținut.

a ceea ce caută, folosind o descriere textuală a conținutului datelor căutate, pe baza unei descrieri grafice schematice a proprietăților datelor căutate și așa mai departe (vezi Secțiunea 2.2);

- **conversia în descriptori:** sistemul de căutare traduce cererea utilizatorului în atribute de conținut folosind un mecanism similar cu cel folosit la adnotarea conținutului bazei de date. Acești descriptori pot fi proprietăți de culoare, forme, informație audio sau de mișcare (vezi Secțiunea 2.1);
- **căutarea propriu-zisă:** căutarea se realizează prin compararea atributelor cererii de căutare cu cele deja stocate în baza de date. Folosind diverse măsuri de distanță și similaritate între atribute, sistemul va căuta datele ce sunt cele mai apropiate (similare) de criteriile formulate (vezi Secțiunea 2.3);
- **interacția cu utilizatorul:** rezultatele căutării sunt furnizate utilizatorului de regulă folosind sistemul de navigare. Acesta presupune o interfață vizuală intuitivă în care utilizatorul poate vizualiza eficient conținutul datelor. În mod opțional, sistemul poate interacționa cu utilizatorul ("feedback") pentru a îmbunătății performanțele sistemului, de exemplu înregistrând opinia utilizatorului cu privire la relevanța datelor returnate de sistem (vezi Secțiunea 2.4).

În cele ce urmează vom detalia fiecare dintre aceste etape.

2.1 Descrierea conținutului datelor

Într-o primă etapă, informația propriu-zisă din baza de date este reprezentată prin intermediul atributelor de conținut. Sistemul va genera pentru fiecare document o colecție de atribute ce vor caracteriza proprietățile relevante ale conținutului acestuia (denumiți și *descriptori*). De exemplu, documentul X poate fi descris de atributele A_1, A_2, \dots, A_n unde valorile $\{a_1, a_2, \dots, a_n\}$ formează descriptorul de conținut. Atributele definesc ceea ce numim *spațiul de caracteristici* al datelor, de regulă un spațiu n -dimensional.

Atributele pot fi, fie date de *nivel semantic scăzut*, precum măsuri statistice, parametri numerici (de exemplu: histograme de culoare¹, câmpuri vectoriale de mișcare, histograme de orientare a conturilor din imagine), fie date simbolice de *nivel semantic superior* (de exemplu: nume obiecte de interes, percepția culorilor, recunoaștere text "încrustat" în imagine, identificarea prezenței umane). Cu alte cuvinte, informația inițială heterogenă și multimodală a fost convertită la o reprezentare uniformă într-un sistem unitar normalizat definit de spațiul de caracteristici. Fiecare document va fi caracterizat astfel de o anumită valoare a acestor atribute, definind un punct unic în spațiu.

Pentru a ilustra aceste aspecte, în Figura 2.2 am prezentat un exemplu concret de reprezentare a conținutului în cazul înregistrărilor audio (și în particular al sunetelor animalelor). Spațiul de caracteristici este definit în acest caz de trei atribute și anume: entropia Wiener² (A_1), amplitudine (A_2) și continuitate în timp (A_3) (spațiu tridimensional). Astfel, fiecare punct din spațiu, P_i (reprezentat grafic de un cerc) cu $i = 1, \dots, N$ unde N reprezintă numărul de înregistrări disponibile, reprezintă o înregistrare audio al cărei conținut a fost descris de valorile atributelor A_1, A_2, A_3 , și anume $P_i = \{a_{i1}, a_{i2}, a_{i3}\}$ (vezi și Secțiune 4 relativă la fuziunea descriptorilor). Dacă atributele sunt suficient de discriminatorii, înregistrările audio similare din punct de vedere al conținutului trebuie să conducă la puncte apropiate spațial (vezi cercurile de aceeași culoare) în timp ce înregistrările diferite trebuie să conducă la puncte distanțate spațial (vezi punctele de culori diferite).

Tot în această etapă a descrierii conținutului datelor, opțional, se pot

¹histograma unei imagini este o măsură a probabilităților discrete de apariție a culorilor (sau a intervalelor de culoare denumite și bini) în imagine, valorile acesteia reprezentând numărul de apariții al unei culori raportat la numărul total de pixeli. Astfel, histograma are sens de densitate de probabilitate a variabilei aleatoare determinată de valoarea unui pixel.

²entropia Wiener este definită ca fiind o măsură a lățimii și uniformității spectrului de putere audio. Ca referință, pe o scală de la 0 la 1, zgomotul alb (semnal aleator cu densitate spectrală de putere constantă) are o entropie 1 iar un ton pur are o entropie 0.

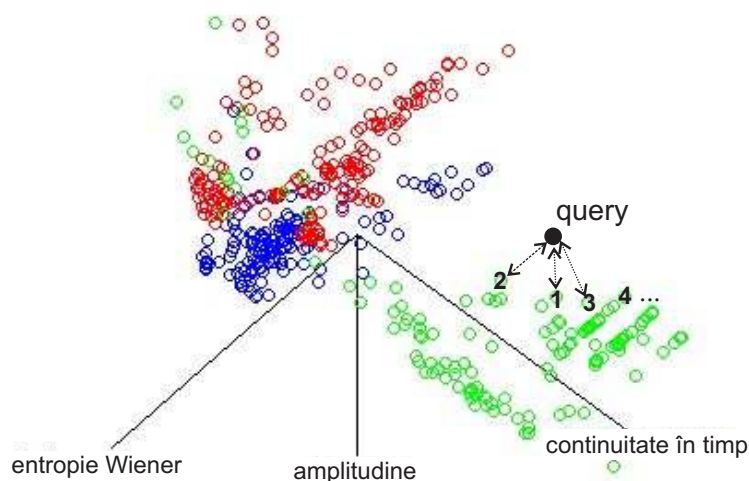


Figura 2.2: Exemplu de spațiu de caracteristici în cazul conținutului audio (sursă imagine programul de prelucrare audio "Sound Analysis Pro", <http://soundanalysispro.com/>).

genera *descrieri compacte*, precum scurte rezumate pentru secvențele video sau pasaje de text reprezentative pentru documentele textuale. Rolul acestor rezumate este acela de a eficientiza vizualizarea conținutului datelor. De exemplu, pentru o bază de documente video este practic imposibil ca utilizatorul să poată vizualiza rapid conținutul acesteia. În acest caz, sistemul poate furniza utilizatorului doar câteva imagini reprezentative sau un rezumat de câteva secunde (exemplu un "trailer") ce redă informația cheie din secvență.

Dacă în urmă cu câțiva ani de zile extragerea de atribute putea fi considerată ca o etapă ce poate fi realizată "off-line", timpul de prelucrare nefiind critic, în prezent datorită dinamicii colecțiilor multimedia (să luăm ca exemplu YouTube ce raporta în 2012 o rată de încărcare de 72 de ore video pe minut) aceasta trebuie realizată mult mai rapid decât o prelucrare în timp real și trebuie să poată fi scalabilă (să poată fi aplicată unor colecții de date dinamice).

În acest punct al indexării, problema care apare este aceea a *relevanței atributelor folosite*. Diversitatea surselor de informație disponibile face dificilă eficientizarea reprezentării datelor. Cu cât crește dimensiunea spațiului de caracteristici și astfel numărul de atribute folosite la reprezentarea datelor cu atât tinde să crească redundanța informației și să scadă puterea dis-

criminatorie a acestora. Un descriptor eficient este acela care maximizează informația reprezentată și minimizează dimensionalitatea datelor. Mai multe informații relative la tehnicile existente de adnotare a conținutului sunt prezentate în Secțiunea 3.

2.2 Formularea căutării

Sistemul de căutare va permite utilizatorului să localizeze informațiile dorite pe baza formulării unei cereri de căutare, denumită și "query" (concept similar celui utilizat în contextul bazelor de date numit și interogare). În mod ideal, sistemul trebuie să poată permite ca aceasta să fie formulată într-un mod cât mai natural și cât mai apropiat de modul de percepție uman, pentru a putea fi la îndemâna oricărui utilizator.

Precizia rezultatelor căutării este în primul rând dependentă de modul de formulare a cererii de căutare a datelor sau cu alte cuvinte a modului de descriere a datelor care se doresc a fi găsite. Formularea adecvată a criteriilor de căutare nu este dependentă numai de sistemul de indexare aceasta depinzând în mare parte și de utilizator.

În primul rând, nivelul de cunoaștere de către utilizator a caracteristicilor datelor căutate este primul factor ce influențează căutarea. Se întâlnesc de regulă următoarele situații posibile [Maillet 03]:

- utilizatorul știe cu siguranță că datele căutate se află în baza de date. În acest caz, ținta este unică iar utilizatorul va fi capabil să formuleze eficient cererea de căutare. Utilizatorul va repeta căutarea până când va obține datele dorite;
- utilizatorul caută o anumită informație dar nu este sigur că aceasta este prezentă în baza de date. În acest caz, sistemul de indexare are rolul de a furniza algoritmi de căutare preciși și eficienți pentru ca utilizatorul să se decidă rapid dacă datele dorite sunt cu adevărat prezente în baza de date. Rafinarea ulterioară a căutării va permite identificarea mai precisă a datelor căutate;
- utilizatorul are informații vagi cu privire la ceea ce dorește să găsească în baza de date. În această situație, sistemul de navigare poate fi folosit pentru "răsfoirea" preliminară a conținutului și identificarea unor informații de interes repositionând utilizatorul în una dintre cele două situații enumerate anterior.

Odată identificată informația dorită este necesar un formalism care să permită enunțarea cererii de căutare în sistemul de căutare. Acesta face

practic legătura dintre modul de percepție uman și reprezentarea informației în sistemul respectiv. În funcție de natura datelor căutate, în literatură există o multitudine de abordări posibile:

- **folosirea vorbirii:** în cazul căutării textuale (informație sub formă de text) se poate folosi direct comanda vocală. Utilizatorul vorbește practic ceea ce dorește să caute, de exemplu: ”caută prognoza meteo pentru astăzi” sau ”caută informații despre posibilități de cazare în Paris”. Comanda este transformată folosind algoritmi de recunoaștere automată a vorbirii în text care este comparat mai departe cu datele din bază. Datorită limitărilor tehnologice a sistemelor de indexare multimedia, o astfel de abordare foarte generală rămâne viabilă doar în cazul căutării de text, ca de exemplu pe Internet (vezi sistemul Siri de pe dispozitivele iPhone³ sau sistemul Google Voice Search de pe dispozitivele cu sistem Android⁴);
- **folosirea de cuvinte cheie:** reprezintă o variantă intermediară a cazului anterior. Cererea de căutare este tot textuală dar este exprimată într-un mod mai restrictiv pe baza unor cuvinte cheie. Pentru ca acest mecanism să funcționeze, datele căutate trebuie să aibă asociate descrieri textuale similare, descrieri ce sunt generate de regulă de utilizatori (de exemplu în momentul în care datele sunt încărcate pe o platformă media on-line) sau în mod automat (metodele de adnotare textuală automată a conținutului multimedia - ”tagging” - sunt totuși încă destul de imprecise);
- **folosirea unui concept:** este de asemenea legată de specificarea unor cuvinte cheie. Diferența față de cazul anterior este dată de faptul că un concept este o noțiune destul de generală care face referire la o clasă de date și nu neapărat la un obiect particular. De exemplu, se dorește localizarea tuturor imaginilor ce conțin arbori, unde conceptul căutat este ”arbore”, sau a secvențelor în care apar case, conceptul căutat fiind acela de ”casă”. Noțiunea de căutare de concepte este asociată în prezent datelor video și constituie un pas intermediar în atingerea unui nivel de descriere textuală. La ora actuală sistemele de căutare după conținut video sunt limitate în a fi antrenate la a răspunde unui număr destul de limitat de concepte (de ordinul miilor - vezi campania TRECVID⁵);

³<http://www.apple.com/ios/siri>

⁴<http://www.google.com/mobile/voice-search>

⁵<http://trecvid.nist.gov>

- **folosirea unui exemplu:** în acest caz, cererea este formulată folosind un model al datelor. De exemplu, utilizatorul caută toate imaginile asemănătoare cu o anumită imagine de care dispune, imaginea fiind furnizată ca exemplu (vezi sistemul de căutare Google Image Search⁶). Tot în această categorie intră și cazul în care utilizatorul furnizează o descriere schematică a datelor căutate. De exemplu acesta nu dispune de o imagine de referință dar poate reprezenta schematic conținutul dorit generând o schiță a imaginii (poziționarea anumitor categorii de obiecte, prezența anumitor culori și așa mai departe - vezi sistemul QBIC al Hermitage Museum⁷);
- **folosirea gesturilor:** un mod interesant de formulare a cererii de căutare o reprezintă gesticularea obiectului care se dorește a fi căutat. Acest mod de căutare are totuși un interes mai mult științific deoarece limitările fiziologice fac imposibilă reprezentarea oricărui obiect prin intermediul gesturilor (vezi un exemplu în [Shirahama 11]);
- **fredonarea unui pasaj audio:** în cazul căutării înregistrărilor audio, de regulă muzicale, o modalitate inedită de formulare a cererii de căutare constă în fredonarea unui pasaj din melodia dorită (vezi de exemplu sistemul Midomi⁸).

2.3 Căutarea datelor

Pentru a fi înțelese de sistem, cererile de căutare trebuie mai întâi convertite în atribute de conținut folosind același mecanism ca și în cazul adnotării inițiale a bazei de date. În acest fel, cererea de căutare este reprezentată practic în spațiul de caracteristici definit în etapa anterioară, prin intermediul unui descriptor. Mai departe, căutarea propriu-zisă se efectuează prin compararea valorilor acestui descriptor cu valorile descriptorilor datelor din bază.

Rezultatele căutării vor fi acele date ale căror valori sunt cele mai apropiate din punct de vedere al unuia sau a mai multor *criterii de similaritate*, de exemplu valorile minime ale unei mărimi de distanță, folosirea unei baze de reguli de decizie și așa mai departe (vezi Secțiunea 5).

De exemplu, în cazul sistemului din Figura 2.2, cererea de căutare poate consta într-un exemplu de înregistrare audio. Utilizatorul dorește localizarea

⁶<http://images.google.com>

⁷<http://www.hermitagemuseum.org/cgi-bin/db2www/qbicSearch.mac/qbic?sellLang=English>

⁸<http://www.midomi.com>

tuturor înregistrărilor audio similare cu aceasta. Exemplul este convertit de sistem într-o serie de valori ale atributelor folosite la indexare, a_1, a_2, a_3 , definind descriptorul de căutare: $query = \{a_{q1}, a_{q2}, a_{q3}\}$. Rezultatele căutării vor fi acele înregistrări audio ce corespund punctelor cele mai apropiate de punctul definit de descriptorul de căutare (vezi Figura 2.2). Datorită subiectivității procesului de căutare, sistemul nu se limitează în a furniza un singur rezultat, ci va returna o clasificare ("ranking") a datelor în ordinea descrescătoare a similarității: poziția 1 - data cea mai similară, poziția 2 - următoarea dată cea mai similară, poziția 3, și așa mai departe.

În acest punct al procesului de indexare, problema principală este definierea conceptului de *similaritate* dintre date. Dacă în cazul datelor numerice soluția la această problemă se găsește în matematică (prin conceptul de metrică), lucrurile nu sunt așa de evidente în cazul datelor multimedia ce implică folosirea de descriptori de natură diferită (text-audio-vizuali). De exemplu, când două secvențe pot fi considerate similare? sau două pasaje de text? Aceasta este o problemă subiectivă chiar și pentru utilizator. Întreg procesul de indexare depinde de modul de definire al măsurii de distanță, schimbarea acesteia poate conduce la rezultate complet diferite. O prezentare detaliată a măsurilor de distanță folosite în contextul indexării datelor este realizată în Secțiunea 5.

2.4 Interacția cu utilizatorul

Ultima etapă a procesului de indexare constă în interacția cu utilizatorul. Aceasta este realizată de regulă prin intermediul *sistemului de navigare*. Sistemul de navigare este practic o interfață grafică ce deservește mai multe funcționalități.

O primă funcționalitate, independentă de procesul de căutare, este aceea de a furniza utilizatorului acces direct la datele din bază. În funcție de tipul datelor, poate fi necesară adoptarea unei strategii complexe. De exemplu, o bază de imagini poate fi vizualizată doar prin reprezentarea în miniatură a acestor imagini (folosind "thumbnails"). În cazul unei baze de secvențe video, acest lucru poate fi realizat prin prezentarea a câtorva imagini reprezentative pentru fiecare secvență. Totuși acest mod de prezentare nu este mereu suficient deoarece nu furnizează nici o informație relativă la conținutul de mișcare (de acțiune) specific. O serie de soluții de reprezentare a conținutului video sunt discutate în Secțiunea 7.

O a doua funcționalitate, și poate cea mai importantă, este aceea de a pune la dispoziția utilizatorului rezultatele obținute în urma etapei descrise anterior și anume a căutării după anumite criterii. Rezultatele sunt de re-

gulă vizualizate în ordinea descrescătoare a relevanței (similarității) față de cererea de căutare.

În final, o altă funcționalitate o constituie interacția cu utilizatorul. În ciuda progresului actual al tehnicilor de descriere a conținutului multimodal, procesul de indexare este inerent limitat de însăși natura datelor (vezi Secțiunea 9). Trebuie să ținem cont că practic imaginile și înregistrările sunt de fapt niște proiecții limitate, bidimensionale, ale lumii înconjurătoare. Astfel, după cum am prezentat și în introducerea acestui capitol, datorită puterii discriminatorii limitate a descriptorilor, rezultatele căutării nu sunt întotdeauna adaptate necesității utilizatorului. Pentru a ameliora acest aspect, de-a lungul timpului au fost studiate o serie de abordări ce tind să includă în procesul de indexare expertiza umană. Printre acestea, cea mai cunoscută poartă numele de "Relevance Feedback" (RF).

Un scenariu clasic de RF poate fi formulat în felul următor: pentru o anumită cerere de căutare rezultatele obținute sunt puse la dispoziția utilizatorului în ordinea descrescătoare a relevanței. Mai departe, utilizatorul este solicitat să marcheze un număr limitat dintre acestea (de regulă de ordinul zecilor) în funcție de relevanța lor. Utilizatorul va marca datele ca fiind relevante - datele corespund perfect cererii de căutare sau nerelevante - datele nu corespund. Pe baza acestor informații, sistemul de căutare calculează o nouă reprezentare a datelor căutate și returnează o rafinare a rezultatelor inițiale. Cu alte cuvinte, acest proces îmbunătățește răspunsul sistemului folosind informația de la utilizator pe post de "realitate" (sau "ground truth"⁹). Mai multe informații sunt prezentate în Secțiunea 6.

În acest punct al procesului de indexare avem la dispoziție un lanț complet de căutare ce pornește de la definirea cererii de căutare și se finalizează cu interacția cu utilizatorul relativ la rezultatele obținute. Problema care apare în acest punct este *evaluarea performanței sistemului*. Cum putem evalua performanțele unui sistem de indexare? Faptul că acesta furnizează rezultate bune pentru o serie de cereri de căutare (vezi exemplul din Figura 6.1) îi garantează performanța?

În realitate, performanțele sistemului variază în mod evident de la o cerere de căutare la alta (este posibil să avem date care sunt mai ușor de localizat datorită conținutului acestora). Avem nevoie de o modalitate generală care

⁹termenul de "ground truth" își are originea în domeniul cartografiei și implică procesul de colectare de informații despre un anumit fenomen, prin observarea practică pe teren a acestuia. Datele obținute constituie "realitatea de teren" folosită pentru calibrarea, validarea și interpretarea observațiilor sau a măsurărilor de la distanță a fenomenului în cauză sau a altor fenomene similare. În contextul indexării, "ground truth" reprezintă datele pentru care se cunoaște conținutul acestora, de exemplu faptul că o imagine reprezintă un anumit obiect sau că o secvență video este de un anumit gen.

CAPITOLUL 2. MECANISMUL DE INDEXARE DUPĂ CONȚINUT 17

să evalueze performanța sistemului, global, în orice situație. Acest lucru este realizat de regulă testând răspunsul acestuia la căutarea fiecărui document din baza de date considerată. Practic, fiecare document devine cerere de căutare.

Evaluarea performanței rezultatelor este mai departe realizată fie subiectiv, de exemplu pe baza opiniei utilizatorilor, fie obiectiv folosind măsuri numerice de performanță (exemplu numărul mediu de rezultate corecte, numărul mediu de rezultate eronate și așa mai departe). O trecere în revistă a abordărilor cel mai frecvent folosite în literatura de specialitate este prezentată în Secțiunea 8.

CAPITOLUL 3

Descrierea conținutului multimodal

După cum am menționat în secțiunile anterioare, procesul de adnotare al conținutului datelor constă în *crearea atributelor* sau a descriptorilor de conținut ce constituie baza sistemului de indexare. Practic căutarea datelor se realizează prin compararea valorilor acestor descriptori pentru cererea de căutare ("query") cu descriptorii informațiilor existente în baza de date.

În această secțiune vom face o trecere în revistă a tehnicilor existente și a surselor de informație folosite în cazul descrierii conținutului multimedia urmând ca acestea să fie detaliate în secțiunile următoare. În principal putem identifica trei surse majore de informație, și anume (vezi Figura 3.1):

- **informația vizuală:** acesta se referă la datele ce sunt percepute vizual, ca de exemplu culoare, formă, textură, mișcare, precum și derivate din acestea;
- **informația audio:** se referă la datele ce sunt percepute sub formă de semnale sonore, ca de exemplu voce, vorbire, muzică, sunete ambientale sau zgomot;
- **informația textuală:** se referă la datele reprezentate sub formă de text (caractere) ce pot provenii din textul atașat datelor (de exemplu textul ce înconjoară un obiect multimedia pe o pagina de web), textul obținut prin recunoașterea caracterelor ce apar "încrustate" în imagine (exemplu subtitrări), sau textul obținut prin recunoașterea vorbirii din informația audio.

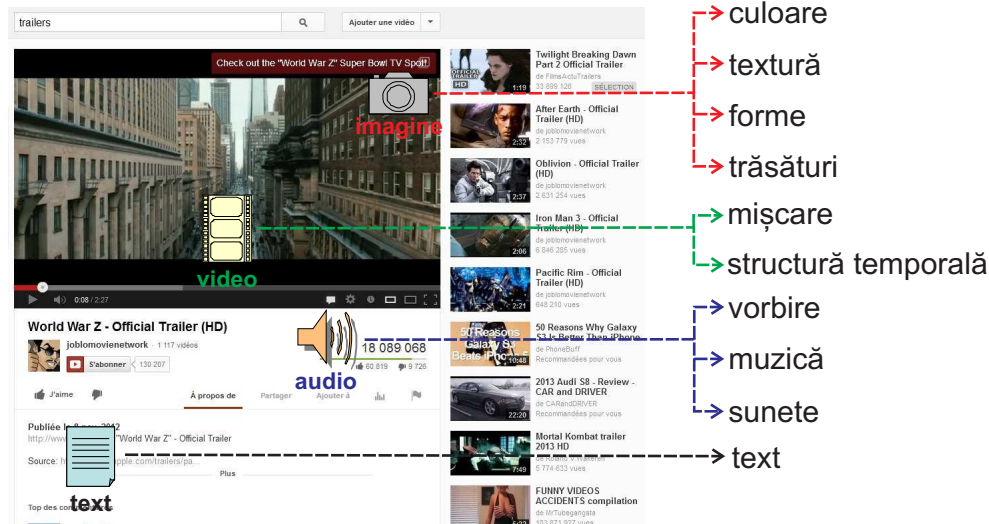


Figura 3.1: Surse de informație multimedia (sursă imagine platformă YouTube, <http://www.youtube.com>).

3.1 Informația vizuală

Informația de culoare reprezintă una dintre sursele de informație cel mai frecvent folosite în cazul descrierii conținutului imaginilor. Acest lucru se datorează în principal faptului că însuși sistemul vizual uman este bazat pe prelucrarea informației de culoare (unde luminoase de diverse frecvențe). Conținutul de culoare este analizat pe baza reprezentării acestuia folosind un anumit model de reprezentare a culorilor¹ sau spațiu de culoare.

Spațiile de culoare folosite variază de la cele clasice, precum sistemul RGB (Red - Roșu, Green - Verde, Blue - Albastru), sisteme ce separă componenta de intensitate de componentele cromatice, precum YCbCr (Y - luminozitate, Cb, Cr - diferențe cromatice), până la sisteme perceptuale în care culorile sunt structurate în așa fel încât să reflecte modul de percepție vizuală umană (culorile similare perceptual sunt alăturate în timp ce culorile opuse se găsesc separate), precum sistemul HSV (Hue - nuanță, Saturation - saturație, Value - măsură a intensității), $L^*a^*b^*$ (L - luminozitate, a,b - diferențe cromatice) în care distanța perceptuală dintre culori tinde să fie proporțională cu distanța matematică, sau HMMD (Hue - nuanță, Max - măsură a gradului de

¹un model de reprezentare a culorilor reprezintă un model matematic abstract ce descrie o culoare ca o combinație de numere, de regulă 3 sau 4 valori, ce corespund unor componente de culori primare (culorile primare sunt culori ce nu pot fi obținute din combinația altor culori).

întunecare sau "shade", Min - măsură a gradului de luminare sau "tint", D - măsură a tonalității sau "tone") sistem ce oferă o serie de avantaje în contextul indexării după conținut precum discretizarea mai eficientă a culorilor. Un studiu detaliat al spațiilor de culoare este prezentat în [Trémeau 04].

O etapă premergătoare descrierii conținutului de culoare constă în reducerea paletelor de culoare² [Orchard 91]. Să luăm exemplul spațiului RGB în care fiecare componentă de culoare este reprezentată pe 8 biți ceea ce conduce la un număr total de 16.777.216 de culori posibile. În practică gestionarea unui număr semnificativ de culori este atât ineficientă deoarece ochiul uman nu este sensibil la micile variații de culoare, cât și nerentabilă din punct de vedere computațional. Paleta de culoare este redusă la un număr semnificativ mai mic, de ordinul sutelor (exemplu de palete fixe: 256 de culori pentru paleta Windows pe 8 biți³; sau 216 culori pentru paleta Webmaster⁴) folosind tehnici de cuantizare a culorilor. De asemenea, analiza conținutului de culoare se poate realiza în urma segmentării imaginii în obiecte, proces de izolare a regiunilor din imagine ce corespund elementelor constitutive ale scenei. În acest fel descrierea culorilor este realizată la nivel de obiect și nu global la nivel de imagine.

Descrierea conținutului de culoare se realizează de regulă folosind descriptori de nivel semnificativ inferior precum histograme de culoare calculate în diverse spații de culoare, histograme ponderate, culori predominante, varianța de culoare, parametri de intensitate, descrierea repartiției spațiale a culorilor, cât și descriptori semantici precum prezența culorii pielii ("skin detection") ce indică prezența umană în scenă sau identificarea denumirii culorilor (asocierea de nume culorilor oferă informații asupra percepției acestora în imagine). Un studiu detaliat este prezentat în [Smeulders 00].

Un exemplu de descriere a conținutului de culoare este prezentat în Figura 3.2 unde sunt ilustrate histogramele de culoare pentru imagini de sport și respectiv animație (în fiecare caz sunt ilustrate câteva imagini reprezentative). Histograma de culoare este calculată folosind metoda propusă în [Ionescu 11] culorile fiind proiectate la paleta Webmaster de 216 culori. Se poate observa faptul că descriptorul de culoare astfel creat ilustrează particularitățile fiecărui tip de conținut, imaginile de sport au o tentă predominant verde în timp ce imaginile de animație sunt predominant galbene-portocaliu conform conținutului acestora.

²paleta de culoare a unei imagini reprezintă mulțimea tuturor culorilor prezente în această imagine. Aceasta reprezintă o sub-mulțime a spațiului de culoare în care este reprezentată imaginea.

³http://en.wikipedia.org/wiki/8-bit_color

⁴<http://www.visibone.com/colorlab>

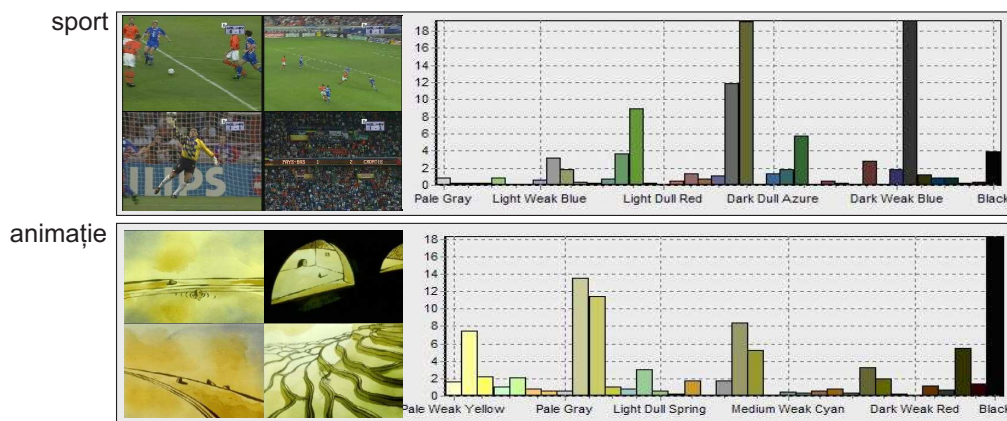


Figura 3.2: Exemplu de descriere a culorilor folosind histograme de culoare în cazul imaginilor de fotbal și respectiv de animație (pe axa orizontală sunt reprezentate culorile în timp ce valorile de pe axa verticală sunt proporționale cu procentul de apariție al acestora în imagini) [Ionescu 11].

Informația relativă la forme se referă la caracterizarea proprietăților obiectelor prezente în scenă din perspectiva proprietăților geometrice ale acestora, fiind specifică imaginilor. Analiza formelor presupune detecția în prealabil a obiectelor din scenă ce este realizată folosind tehnici de segmentare bazate pe contur sau pe regiuni de pixeli [Jain 89]. Succesul adnotării este astfel direct condiționat de calitatea segmentării imaginii.

Problema descrierii formelor nu este una simplă în principal datorită faptului că imaginea nu este altceva decât o proiecție bidimensională a lumii 3D, ceea ce înseamnă că una dintre dimensiunile obiectelor este pierdută. Astfel, formele extrase din imagine vor reprezenta numai parțial informația reală din scenă. Mai mult, imaginea este perturbată de zgomot⁵ și defecte de achiziție ceea ce cresc dificultatea obținerii unei reprezentări robuste.

Un descriptor de formă trebuie să fie eficient în sensul în care acesta trebuie să furnizeze suficientă putere discriminatorie pentru a identifica obiectele similare perceptual în contextul în care acestea pot fi reprezentate în contexte diferite (de exemplu scene diferite, momente temporale diferite), din diverse unghiuri, distorsionat, parțial sau suprapuse peste alte obiecte.

Acest lucru presupune atât o invarianță la zgomot, cât și la rotații, translații, modificări de scală sau în general la orice tip de transformare

⁵zgomotul în imagine se referă la acea informație perturbatoare ce alterează informația utilă.

afină⁶. Problema ocluziei obiectelor poate fi rezolvată prin integrarea de informații suplimentare, precum o evoluție temporală a imaginilor sau informație de adâncime (informație 3D).

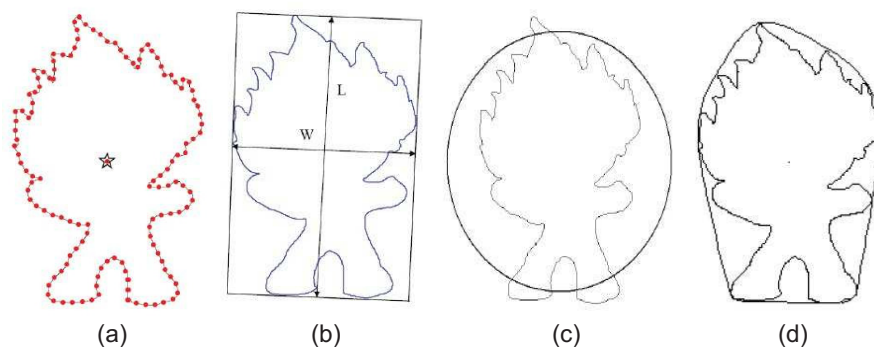


Figura 3.3: Exemplu de descriere a formelor: (a) reprezentarea centrului de greutate pe baza eșantionării uniforme a conturului, (b) determinarea parametrilor de elongație în funcție de rata de aspect a formei (W/L), (c) determinarea raportului de circularitate (arie obiect raportat la aria cercului de același perimetru), (d) convexitate (cea mai mică regiune convexă ce include obiectul). Sursă imagini [Mingqiang 08].

Descriptorii de formă sunt calculați fie folosind doar informația de contur exterior a obiectelor sau informația de contur în relație cu informația din interiorul obiectului (regiunea plină a obiectului). Abordările existente variază de la calculul unor parametri simpli precum suprafață, orientarea axelor principale ale obiectului, convexitate, curbură, lungime, la parametri mai complexi precum momente statistice invariante, parametri spectrali (Fourier sau wavelet), reprezentarea sub formă de coduri (descompunerea conturului în secvențe de segmente de dimensiune unitate și codarea acestora), descompunerea în poligoane, reprezentări de tip "scale-space" (conturul este caracterizat la mai multe niveluri de scală), reprezentare cu matrice de forme, și așa mai departe [Mingqiang 08]. Un studiu detaliat este prezentat în [Smeulders 00]. O serie de exemple sunt prezentate în Figura 3.3.

Informația de textură. Conceptul de textură este legat de caracteriza-

⁶o transformare afină (cuvântul "affinis" în Latină înseamnă "conectat cu") reprezintă o transformare geometrică ce are proprietatea de a păstra coliniaritatea punctelor precum și a rapoartelor de distanță dintre punctele ce se găsesc pe o aceeași dreaptă (de exemplu, punctul de mijloc al unei drepte în urma transformării își va păstra proprietatea). O transformare afină nu garantează totuși conservarea unghiurilor sau a lungimilor, dar are proprietatea de a păstra liniile paralele.

rea proprietăților materialelor prezente în imagini și presupune atât analiza informației de culoare cât și de contur. O textură este definită ca fiind o regiune din imagine ce prezintă caracteristici omogene, precum un motiv de bază ce se repetă în domeniul spațial sau frecvențial.

Un exemplu este ilustrat în Figura 3.4. Tehnicile de descriere a texturilor presupun cuantificarea acestor proprietăți pentru a caracteriza o serie de atribute specifice, precum asperitate, uniformitate, variabilitate, direcționalitate, regularitate, ca o funcție de variația spațială a intensității pixelilor din imagine (de regulă exprimată ca niveluri de gri). Metodele existente pot fi clasificate în abordări statistice, geometrice, pe bază de modele și pe bază de filtre [Tuceryan 93].

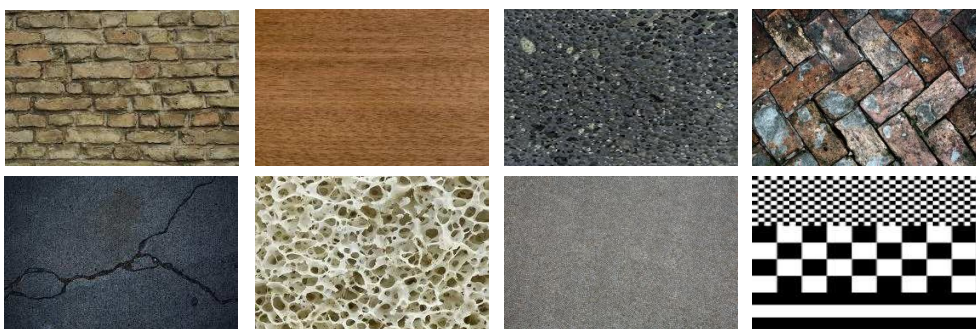


Figura 3.4: Exemplu de texturi (de la stânga la dreapta și de sus în jos): perețe de cărămidă, parchet lemn, ciment, pavaș piatră, zid de piatră, structură osoasă, pavaș de piatră radial și textură artificială (sursă imagini Wikipedia).

Una dintre cele mai utilizate abordări o constituie metodele statistice. Distribuția spațială a intensității pixelilor este caracterizată statistic, ca de exemplu prin calcularea probabilității de co-ocurență a unei anumite intensități în diverse direcții și distanțe față de un punct de referință. Statisticile pot fi calculate pentru valorile unui singur pixel (statistici de ordinul întâi) sau pentru perechi sau regiuni de pixeli (statistici de ordin superior). Astfel de exemple sunt parametrii extrași din matricele de co-ocurență (de exemplu: energie, contrast, corelație), parametrii de autocorelație sau histogramele de contur.

Abordările geometrice analizează textura din perspectiva proprietăților geometrice ale primitivelor acesteia (elementele texturii) precum arie, formă, lungime și a modului de distribuție al acestora într-o anumită rețea (sau "grid"). De exemplu, imaginea unui zid de cărămidă poate fi descrisă pe baza unei singure cărămizi (primitiva texturii în acest caz) și prin definirea rețelei de plasare a acesteia în spațiu.

O altă categorie de abordări sunt metodele bazate pe modele. Texturile sunt sintetizate pe baza unui model al cărui parametri descriu proprietățile esențiale ale acesteia. De exemplu, elementele texturii pot fi modelate ca puncte întunecate sau luminoase, ca tranziții verticale sau orizontale, ca linii. Exemple de astfel de modele sunt lanțurile Markov⁷ și modelarea fractală⁸.

Metodele bazate pe filtre sunt specifice domeniului prelucrării de semnal. Acestea se bazează practic pe filtrarea imaginii atât în domeniul spațial cât și frecvențial. Dintre filtrele cel mai des utilizate sunt operatorii de derivare (de exemplu Laplacian, Roberts) sau filtrele Gabor⁹. Un studiu detaliat al literaturii este prezentat în [Smeulders 00].

Informația de mișcare. Conceptul de mișcare este definit în contextul secvențelor de imagini, numite și imagini în mișcare. O secvență de imagini presupune o evoluție temporală a conținutului unei imagini (informație spațio-temporală; în cazul în care se adaugă și informație audio obținem ceea ce numim video - informație audio-vizuală). Dacă considerăm standardul de codare video PAL - Phase Alternating Line (unul dintre cele mai răspândite în Europa) o secundă dintr-o secvență video corespunde la o succesiune de nu mai puțin de 25 de imagini. Caracterizarea informației de mișcare presupune astfel caracterizarea schimbărilor (de regulă spațiale) ce au loc de la o imagine la alta. Aceste schimbări pot fi analizate local, doar pentru o anumită regiune din imagine (de exemplu mișcarea unui obiect în scenă), sau global pentru întreaga imagine (de exemplu mișcarea camerei video).

Pentru a putea descrie conținutul de mișcare este nevoie mai întâi de realizarea unei etape intermediare ce presupune identificarea acestuia în secvență. O abordare simplificată presupune detecția mișcării [Bovik 09]. Aceasta are ca scop localizarea acelor regiuni de pixeli din imagine în care survin schimbări în timp, de regulă de la o imagine la alta. Limitarea acestei abordări constă în faptul că nu se ține cont de natura acestor schimbări, acestea putând surveni, în special în cazul secvențelor editate în studio, independent de mișcare, de exemplu prin fluctuații de intensitate,

⁷un lanț Markov (denumit după Andrey Markov) reprezintă un sistem matematic caracterizat de tranziții succesive între un număr finit, măsurabil, de stări posibile. Acesta este un proces aleator fără memorie în sensul în care tranziția sistemului la o altă stare depinde doar de starea curentă și nu depinde de stările anterioare.

⁸un fractal (termen creat de Benoît Mandelbrot, din Latină "fractus" - neregulat) reprezintă o suprafață de formă neregulată sau fragmentată creată pe baza unor reguli deterministe sau stohastice ce implică un proces de omotetie internă (transformare geometrică în care punctele corespondente sunt coliniare cu un punct fix (centru), distanța față de el crescând sau reducându-se în raport constant - sursă "Marele Dicționar de Neologisme").

⁹un filtru Gabor (denumit astfel după Dennis Gabor) reprezintă un filtru liniar ce are proprietatea de a avea caracteristici similare filtrelor din sistemului vizual uman.

efecte speciale.

Exemple de astfel de metode includ detecția cu prag fix sau adaptiv (o regiune este declarată de mișcare dacă diferențele dintre pixeli pentru două imagini succesive sunt mai mari decât un anumit prag), tehnici de estimare a fundalului¹⁰ precum media alunecătoare, aproximare filtru median, metode statistice neparametrice, metode recursive și așa mai departe. Un exemplu de detecție este ilustrat în Figura 3.5.(b).

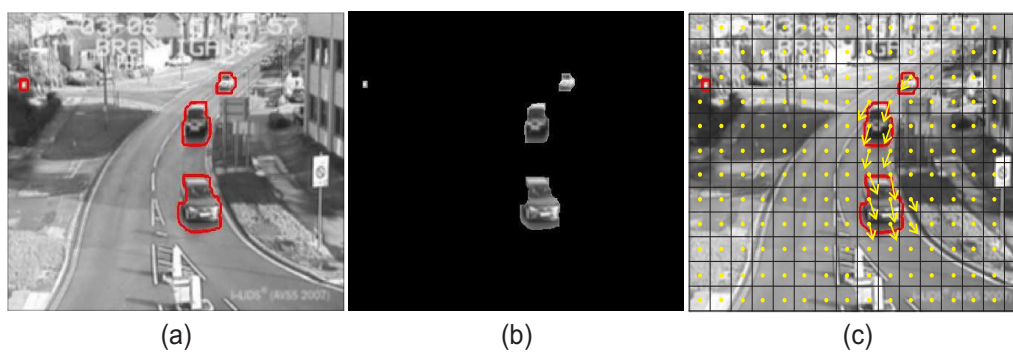


Figura 3.5: Exemple de determinare a conținutului de mișcare pentru o secvență de supraveghere video: (a) imagine din secvența originală (obiectele care se deplasează sunt încercuite cu roșu), (b) detecție de mișcare folosind aproximarea filtrului median (imaginea reprezintă regiunile care se schimbă), (c) câmpul vectorial de mișcare obținut cu o estimare pe blocuri de pixeli (vectorii de mișcare sunt ilustrați cu galben, punctele semnifică absența mișcării).

O a doua abordare o constituie tehnicile de estimare a mișcării [Bovik 09]. Acestea, spre deosebire de detecția mișcării, presupun estimarea deplasărilor pixelilor sau a regiunilor de pixeli de la o imagine la alta, estimare ce este cuantificată prin asocierea unui vector de mișcare. Acesta indică atât direcția deplasării pixelilor (orientare) cât și deplasarea spațială (amplitudine). În urma estimării, imaginea este practic reprezentată de un câmp de astfel de vectori de mișcare indicând modul de deplasare al fiecărui pixel sau bloc de pixeli. Un exemplu este prezentat în Figura 3.5.(c).

Tehnicile de estimare a mișcării variază de la abordări bazate pe metode diferențiale (bazate pe estimarea fluxului optic), parametrice, stohastice sau bazate pe blocuri ("block-based") acestea din urmă regăsindu-se în toate standardele de codare video precum cele dezvoltate de Moving Picture Ex-

¹⁰detecția fundalului sau "background" presupune localizarea acelor pixeli din imagine ce rămân aproximativ constanți de la o imagine la alta.

perts Group - MPEG¹¹ (informația relativă la deplasarea regiunilor permite reconstrucția imaginilor, ceea ce oferă un factor de compresie semnificativ).

Odată identificat conținutul de mișcare, descriptorii de conținut cuantifică o serie de proprietăți ale acestuia. Ca exemple de descriptorii putem enumera determinarea traiectoriei obiectelor din scenă, identificarea tipului de mișcare a camerei video ("zoom" - apropiere/depărtare, rotație, translație), determinarea activității de mișcare folosind cuantizarea varianței amplitudinii vectorilor de mișcare, determinarea distribuției spațiale și temporale a activității de mișcare, construcția de imagini MHI de "istorie a mișcării" (Motion History Images) formate prin acumularea informației de mișcare a fiecărui pixel într-o anumită fereastră temporală, determinarea de histograme de intensitate și așa mai departe. De menționat faptul că determinarea descriptorilor de mișcare depinde de succesul și calitatea detecției/estimării de mișcare folosită.

Informația de structură temporală. Descriptorii de conținut relativi la structura temporală se adresează secvențelor de imagini și în special secvențelor editate în studio, precum filme, reportaje, sport și așa mai departe (în general materiale destinate distribuției TV).

Descrierea structurii temporale video implică segmentarea temporală a acesteia prin descompunerea secvenței în unități structurale de bază numite și plane video [Lienhart 01]. Un plan video este practic o secvență de imagini întregită între pornirea și oprirea camerei video având proprietățile de unitate temporală și de loc (vezi Figura 3.6). Pentru a obține secvența finală, planele video sunt concatenate folosind diverse tranziții video. O tranziție video nu este altceva decât un efect vizual ce poate presupune fie o tranziție abruptă de tip "cut" (concatenarea directă a două plane succesive), fie tranziții graduale precum "fades" (aparitia sau disparitia imaginii dintr-o imagine constantă, de regulă neagră), "dissolves" (transformarea graduală a unei imaginii în alta), "mattes", "wipes" și așa mai departe [Bimbo 99]. Câteva exemple sunt ilustrate în Figura 3.6.

Practic segmentarea temporală implică localizarea în secvență a acestor tranziții. Ca frecvență de apariție, tranzițiile de tip "cut" sunt cele mai frecvente, de regulă 30 de minute video pot conține până la 300 de astfel de tranziții în timp ce frecvența tranzițiilor graduale este cu cel puțin un ordin

¹¹MPEG sau Moving Picture Experts Group, reprezintă o organizație internațională ce se ocupă cu dezvoltarea normelor pentru compresia, decompresia și analiza și codarea video. Aceasta este responsabilă pentru dezvoltarea standardelor clasice de codare, precum MPEG-1 folosit pentru formatul VideoCD, MPEG-2 folosit pentru stocarea pe DVD, MPEG-4 folosit la BD (Blu-Ray Disc), MPEG-7 standard de descriere a conținutului video pentru indexare multimedia sau MPEG-21 ce definește interoperabilitatea tuturor tipurilor de conținut multimedia.

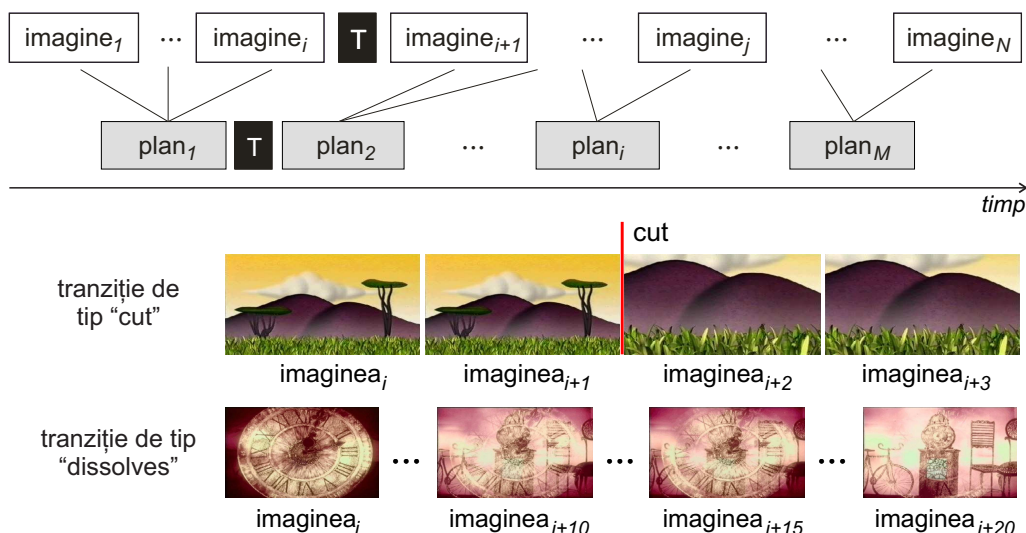


Figura 3.6: Structura temporală a unei secvențe video (T reprezintă o tranziție video, N este numărul de imagini al secvenței iar M numărul de plane video). În partea de jos a imaginii sunt ilustrate un exemplu de tranziție de tip "cut" (imagini film de animație "Gazoon" [CITIA 13]) și respectiv "dissolve" (imagini film de animație "Coeur de Secours" [CITIA 13]).

de mărime mai redusă.

Metodele de detecție a tranzițiilor video exploatează în general detecția discontinuității vizuale produse de acestea în fluxul video folosind abordări bazate pe analiza intensității pixelilor (de exemplu o tranziție de tip "cut" implică o diferență semnificativă a distribuției de culoare ce poate fi analizată folosind diferența dintre histograme, o tranziție de tip "fade" implică o variație graduală a intensității luminoase), analiza contururilor (de exemplu o tranziție de tip "dissolves" presupune un raport semnificativ de puncte de contur ce apar/dispar din imagine), analiza mișcării (de exemplu o tranziție de tip "cut" produce o discontinuitate a vectorilor de mișcare) sau analiza informației în domeniul comprimat (precum analiza coeficienților transformatei cosinus discrete din fluxul MPEG).

La un nivel de descriere superioară, segmentarea secvenței poate implica descompunerea acesteia în unități structurale de nivel semantic superior, precum gruparea planelor video în scene (grupuri de plane video ce sunt corelate din punct de vedere al conținutului semantic și presupun unitate de loc, de timp și de acțiune), în episoade (grupuri de scene ce sunt similare din punct de vedere al acțiunii globale, ca de exemplu episoadele unei serii TV) și așa mai departe.

Interesul în segmentarea temporală este dublu. Pe de-o parte, acesta constituie primul pas de analiză pentru marea parte a metodelor de analiză a conținutului video deoarece furnizează informații relative la structura semantică a acestuia. De exemplu, având la dispoziție structura de plane sau de scene video, analiza de conținut se poate realiza în interiorul acestora evitând astfel prelucrarea imaginilor de tranziție cât și asigurând unitatea semantică a conținutului.

Extragerea unui descriptor pentru un segment ales aleator din secvență riscă să amestece informații distincte. De exemplu, dacă considerăm cazul particular al unei secvențe de știri și segmentul ales conține atât înregistrarea prezentatorului cât și a unui reportaj extern, amestecarea informațiilor vizuale ale celor două subiecte complet diferite nu poate produce un descriptor reprezentativ.

Pe de altă parte, structura temporală furnizează ea însăși informații de conținut. Folosirea unui anumit tip de tranziții pentru a face legătura între planele video nu este aleatorie ci corespunde unor reguli cinematice de montaj bine definite [Reynertson 70]. De exemplu, folosirea frecventă a tranzițiilor de tip "cut" are ca efect creșterea dinamismului secvenței, tranzițiile de tip "dissolve" și "fade" sunt folosite frecvent pentru a schimba timpul sau locul acțiunii, o secvență de tip "fade-out - fade-in" introduce un moment de pauză în derularea acțiunii ca de exemplu pentru a trece la un alt capitol al narațiunii.

Descriptorii ce caracterizează informația de structură temporală exploatează în principal frecvența de apariție a schimbărilor de plan video în secvență, fie în mod direct prin măsuri precum determinarea duratei medii a planurilor video, ratei medii de schimbare de plan raportată la unitatea temporală (de regulă denumită ritm vizual), raportului tranzițiilor graduale din secvență, extragerea de imagini cheie la nivelul fiecărui plan și prelucrarea acestora folosind descriptori clasici de imagine (vezi secțiunile anterioare); fie derivând informații relative la activitatea vizuală a secvenței exploataând conceptul de acțiune (concept de regulă asociat frecvenței de tranziții de tip "cut", de exemplu o secvență de acțiune va avea o densitate ridicată de plane video de scurtă durată în timp ce o secvență a unui documentar este foarte probabil să conțină doar câteva plane video).

Trăsături. Informația legată de ceea ce numim trăsături ("features") este de fapt un caz particular de descriere a informației de contur în imagini și este strâns legată de noțiunea de puncte de interes ("interest points").

Un punct de interes în imagine reprezintă de regulă o regiune de pixeli (de dimensiuni reduse) a căror proprietăți o fac reprezentativă pentru înțelegerea conținutului structural al imaginii. Nu orice regiune de pixeli care conține

contururi este astfel un punct de interes.

De exemplu, dacă considerăm sistemul vizual uman, se știe faptul că ochiul este mai sensibil la percepția punctelor de inflexiune din imagini, precum unghiuri sau intersecții, decât la informațiile redundante, continue, precum liniile drepte. Aceste informații sunt acelea ce tind să fie percepute primele în imagine, fiind definatorii, și apoi pe baza lor să se realizeze o aproximare a scenei.

În Figura 3.7.(a) am ilustrat un exemplu în acest sens. În imagine sunt reprezentate patru treimi de cercuri dispuse simetric. La o primă vedere ochiul uman tinde să perceapă mai întâi cele patru colțuri contrastate (puncte de interes) și să extrapoleze informația imaginând un dreptunghi alb suprapus peste patru cercuri negre. Totuși în realitate, liniile ce definesc dreptunghiul nu există, fiind doar o iluzie.

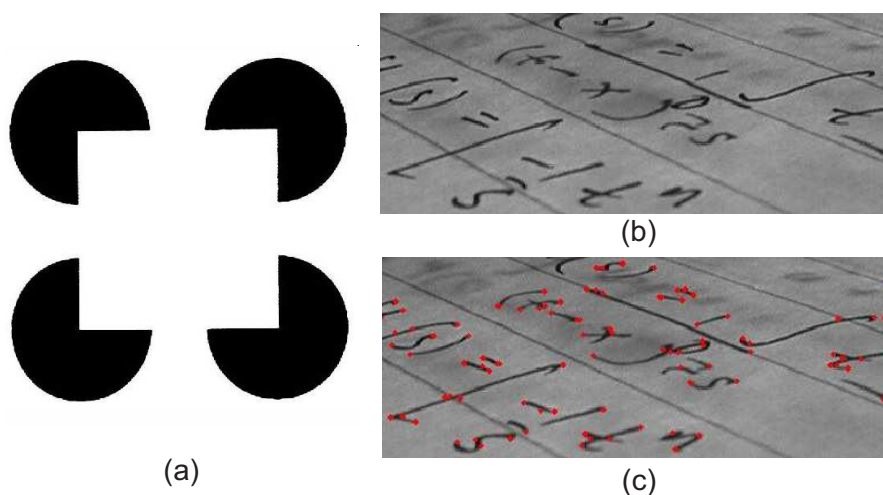


Figura 3.7: Exemplu de trăsături: (a) exemplu de iluzie optică în care cele patru treimi de cerc sunt percepute ca un dreptunghi suprapus peste patru cercuri negre (sursă <http://webvision.med.utah.edu/book/>), (b) și (c) ilustrează un exemplu de detector de colțuri, prima imagine fiind imaginea inițială iar în a doua imagine punctele roșii marchează trăsăturile detectate (sursă imagini Wikipedia).

În contextul imaginilor, aceste trăsături pot fi formalizate ca fiind acele puncte din imagine ce întrunesc următoarele proprietăți: au o definiție matematică bine precizată, au o poziție bine definită în imagine, informația locală din jurul punctului de interes este bogată informațional (cu alte cuvinte sunt definite de context), și cea mai importantă proprietate este aceea că acestea

trebuie să fie stabile la perturbații locale și globale precum deformări datorate transformărilor de perspectivă, schimbarea unghiului de vizualizare, schimbări de scală, rotații, translații cât și variații de iluminare (de exemplu: colțul unui dreptunghi își va păstra proprietatea indiferent dacă este întunecat, rotit, micșorat sau schimbată perspectiva). Datorită acestor proprietăți, punctele de interes sunt de departe cea mai eficientă modalitate de reprezentare a conținutului imaginilor în contextul indexării după conținut.

Tehnicile de detecție a punctelor de interes/trăsături au pornit inițial de la ideea detectării colțurilor în imagini, un astfel de detector fiind "Harris corner detector" ce folosește ipoteza conform căreia gradientii (diferențele) pe cele două direcții oX și respectiv oY trebuie să fie ambele semnificative pentru un colț. Un exemplu de astfel de detecție este prezentat în Figura 3.7.(c).

Alte metode mai elaborate sunt detectorul Harris Laplace (cunoscut ca detectorul Harris multi-scală) ce adaugă localizarea colțurilor folosind reprezentării ale imaginii pe diverse niveluri de scală, detectorul Hessian Laplace, abordări ce folosesc reprezentări de tip "scale-space" (similar informației de contur) precum Laplacian of Gaussian (LoG), Difference of Gaussian (DoG) sau Determinant of Hessian (DoH), detectorul Maximally Stable Extremum Regions (MSER) ce selectează anumite regiuni conexe din imagine dacă acestea sunt stabile în urma filtării repetate cu diverse praguri ("thresholding"¹²), până la binecunoscuții detectori Scale Invariant Feature Transform - SIFT (ce se bazează pe localizarea maximelor și minimelor obținute în urma aplicării unor funcții de diferențe de Gaussiene¹³) și respectiv Speeded Up Robust Features - SURF (ce folosește o descompunere wavelet de tip Haar și imagini integrale). O trecere în revistă a diferitelor tehnici de analiză a trăsăturilor în imagini poate fi consultată în [Gauglitz 11].

Având în vedere eficiența descriptorilor de trăsături în reprezentarea structurii imaginii, în special datorată invarianței acestora la o gamă largă de transformări, cercetările actuale în domeniu vizează extensia acestora pentru a putea exploata conținutul temporal specific secvențelor de imagini.

Dintre descriptorii de trăsături spațio-temporali putem menționa "Har-

¹²"thresholding" în prelucrarea de imagini reprezintă operația prin care valorile imaginii sunt transformate prin compararea cu un simplu prag de regulă obținând o imagine binară. Dacă valoarea din imagine este superioară pragului aceasta este schimbată într-o constantă (de regulă 1) și în caz contrar într-o altă constantă (de regulă 0).

¹³diferența de Gaussiene constă în realizarea diferenței dintre două variante încetșoșate ale imaginii inițiale, de regulă prima imagine fiind mai încetșoșată ("blurred"). Încetșoșarea unei imaginii presupune înlăturarea frecvențelor înalte (de exemplu zone ne-uniforme precum contururi). Prin realizarea diferenței între două astfel de imagini se obține un filtru trece bandă care conservă doar o gamă de frecvențe spațiale din imaginea inițială și astfel doar anumite informații din imagine.

ris3D corner detector” (extensie a detectorului de colțuri pentru a include pe lângă gradientii spațiali și gradienti temporali), detectorul Cuboid ce folosește filtre Gabor temporale (vezi explicația de la informația de textură) pentru a detecta acele trăsături cu proprietăți spațiale particulare și ce presupun o mișcare complexă, detectorul Hessian 3D ce se bazează pe estimarea determinantului matricei Hessiene¹⁴ în care derivatele parțiale sunt calculate și temporal, tehnici de eșantionare densă precum extragerea de regiuni 3D din secvență (de exemplu o porțiune de imagine pentru mai multe momente de timp) și descrierea acestora adaptând descriptori de trăsături precum SURF 3D (extensie a descriptorului SURF). Un studiu detaliat al descriptorilor spațio-temporali este prezentat în [Stöttinger 10].

3.2 Informația audio

Informația audio reprezintă o altă sursă importantă de informații relative la conținutul datelor multimedia. Aceasta se referă la caracterizarea sunetului, fie în contextul video unde acesta este sincronizat informației vizuale, fie independent (de exemplu fișiere audio de muzică, înregistrări, etc.). În general sunt vizate analiza și identificarea vorbirii, a zgomotului și a efectelor sonore sau analiza conținutului muzical.

Prelucrarea semnalului audio se realizează principial într-un mod similar prelucrării secvențelor de imagini fiind de asemenea o reprezentare temporală a datelor. Un semnal audio digital (discret) nu este altceva decât o secvență de eșantioane (valori de amplitudine ale undelor sonore) înregistrate în timp (vezi Figura 3.8.(a)). Acestea sunt prelucrate la nivel de cadre audio, un cadru audio fiind o secvență temporală ce conține un anumit număr de eșantioane (un exemplu de valoare uzuală este folosirea a 1024 de eșantioane). Important este faptul că aceste cadre nu sunt întotdeauna disjuncte, de regulă fiind suprapuse cu până la 50% din durată. Acest lucru asigură faptul că toate părțile semnalului audio vor fi bine reprezentate la nivel de cadre.

Metodele de descriere a conținutului audio se împart în două categorii. Metode ce analizează informația audio direct în domeniul temporal la nivel de cadru sau folosind o reprezentare statistică a distribuției acestora în documentul audio (descriptorii extrași la nivel de cadru sunt agregați pentru întreaga secvență prin statistici de medie, varianță, median și așa mai departe). Dintre descriptorii cei mai frecvent folosiți în acest caz putem

¹⁴matricea Hessiană (după numele matematicianului Ludwig Otto Hesse) reprezintă matricea pătratică a derivatelor parțiale de ordin doi ale unei anumite funcții de mai multe variabile. Definită în acest fel, aceasta are proprietatea de a descrie curbura locală a funcției.

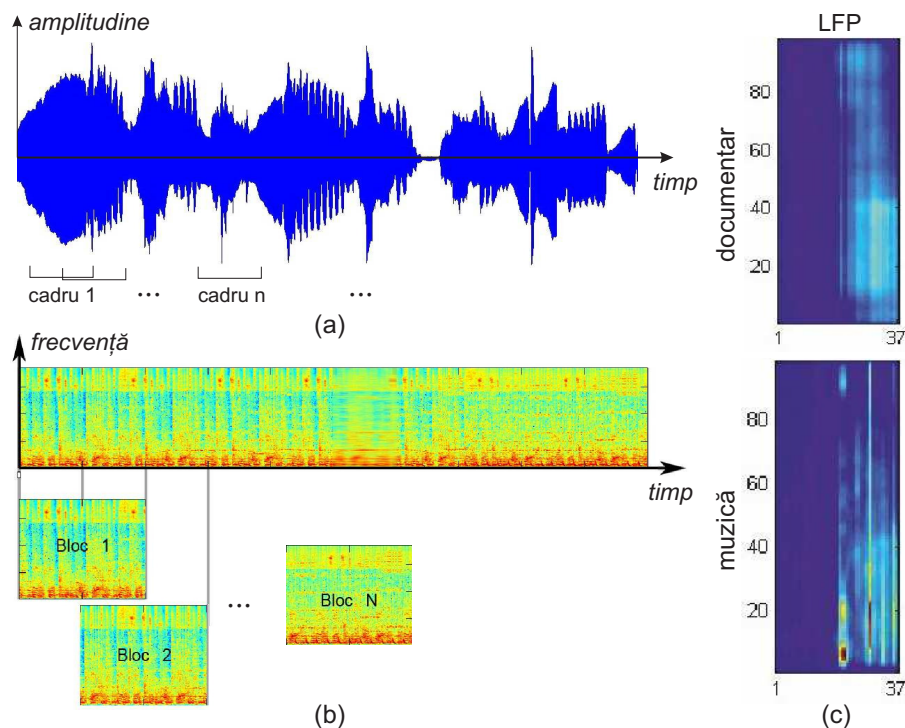


Figura 3.8: Exemplu de prelucrare a semnalului audio: (a) analiză în domeniul temporal, (b) analiză pe blocuri de cadre spectrale în domeniul frecvențial [Seyerlehner 10], (c) exemplu de descriptor Logarithmic Fluctuation Pattern (LFP) [Ionescu 12b] în cazul unui documentar și videoclip muzical (la acesta din urmă se observă aspectul ritmic prezent prin maximele LFP - reprezentate cu roșu și galben).

menționa: Zero-Crossing Rate (ZCR) ce reprezintă numărul de treceri prin zero ale semnalului raportat pe unitatea de timp, energia semnalului (Root Mean Square of Signal Energy sau RMS), rata de absență a sunetului sau coeficienții de autocorelație ai semnalului [Mathieu 10].

Totuși marea parte a metodelor analizează sunetul în domeniul frecvențial. Pentru aceasta, fiecare cadru audio este reprezentat în domeniul transformatei Fourier iar informația obținută este prelucrată într-o reprezentare frecvență (dată de reprezentarea Fourier a cadrelor audio) - timp (dată de succesiunea cadrelor audio în timp) ca de exemplu folosind spectrograma de amplitudine - reprezentarea temporală a amplitudinii transformatei Fourier a fiecărui cadru audio.

Dintre abordările folosite putem menționa distribuția energiei semnalului, centroizii frecvențelor, lărgimea de bandă, "pitch", "loudness" sau reprezen-

tarea coeficienților cepstrali Mel-Frequency Cepstral Coefficients (MFCC). La rândul ei, reprezentarea semnalului frecvență-timp poate fi prelucrată pe blocuri de cadre spectrale (valori uzuale sunt de ordinul a 10 până la 512 cadre per bloc) ceea ce are avantajul de a integra pe lângă informația de frecvență și informație temporală locală, ca de exemplu aspecte ritmice ale semnalului (vezi Figura 3.8.(b)).

Dintre descriptorii de acest gen putem menționa Spectral Pattern (informație relativă la timbru sonor), Logarithmic Fluctuation Pattern (informație relativă la aspectele ritmice ale semnalului), Correlation Pattern (informație relativă la schimbările de intensitate) sau Spectral Contrast Pattern (informație relativă la tonalitate) [Seyerlehner 10].

Un exemplu este prezentat în Figura 3.8.(c) în care am ilustrat descriptorul Logarithmic Fluctuation Pattern (LFP) în cazul coloanei sonore a unui documentar și respectiv a unui videoclip muzical. O caracteristică specifică muzicii este prezența de bătăi ritmice ("beats") ce sunt vizibile sub formă de maxime ale LFP (vezi zone colorate cu roșu și galben). În contrast, în cazul documentarului, structura LFP este plată ceea ce indică că nu există elemente repetitive percutante în fluxul audio.

O altă direcție de studiu importantă ce vizează analiza sunetului o reprezintă tehnicile de recunoaștere automată a vorbirii (Automatic Speech Recognition sau ASR [Lamel 08]). Acestea au ca obiectiv transformarea vorbirii din semnal audio în text ce poate fi prelucrat mai departe folosind tehnici specifice. Folosirea textului obținut în urma ASR furnizează informații prețioase relative la conținutul datelor. Avantajul descriptorilor textuali precum și limitările ASR sunt discutate în secțiunea următoare.

Raportat la descriptorii vizuali, descriptorii audio tind să furnizeze o putere discriminatorie mai bună în marea parte a aplicațiilor relative indexării după conținut, precum identificarea genului video sau detecția anumitor concepte video [Ionescu 12b] [Over 12].

3.3 Informația textuală

De departe cea mai eficientă sursă de informații pentru indexare o constituie textul. Aproape în totalitate, sistemele existente de căutare multimedia se bazează pe descriptorii textuali. Avantajul reprezentării textuale este acela că oferă un nivel de descriere semantică a conținutului, foarte apropiat de nivelul de percepție uman. Mai mult, exprimarea textuală este la îndemâna oricărui utilizator, ceea ce rezolvă problematica formulării cererilor de căutare.

Totuși dezavantajul principal al informației textuale este dat de posibilitatea limitată de automatizare a procesului de generare, aceasta necesitând

practic să fie furnizată de utilizator. De exemplu, la încărcarea on-line a unei imagini, de regulă utilizatorul va specifica o scurtă descriere textuală a conținutului acesteia, ca de exemplu "Turnul din Pisa". Aceasta va fi folosită ulterior drept descriptor de conținut pentru indexare. Totuși această informație este incompletă și nu descrie decât global conținutul, nu există informații relative la persoanele din scenă, la momentul zilei sau relativ la prezența altor obiecte. Acest lucru limitează această imagine să nu poată fi găsită decât în cazul căutării unor imagini cu turnul din Pisa, și nu pentru alte informații din scenă.

În cazul datelor multimedia, informația textuală poate fi obținută din mai multe surse. Conform celor menționate anterior, o primă sursă de descriptori textuali este însuși utilizatorul, aceste date fiind generate manual. În acest caz, informația textuală este de regulă reprezentată sub formă de mici rezumate de conținut referitoare la date ("synopsis", de exemplu în cazul unui film acestea pot fi rezumatul narațiunii acestuia), etichete de conținut ("user tags" ce reprezintă de regulă câteva cuvinte cheie ce descriu conținutul global), subtitrări în cazul filmelor, metadata¹⁵ (ce furnizează o serie de informații suplimentare de tipuri diferite, legături ("link") către alte surse de informații, proprietăți ale datelor), informații referitoare la localizarea geografică a datelor precum coordonatele GPS¹⁶ ale unei imagini (longitudine, latitudine), comentariile utilizatorilor relativ la conținutul datelor specifice de regulă rețelelor de socializare sau textul ce înconjoară elementul multimedia respectiv pe o pagină web.

Câteva exemple de astfel de descrieri au fost prezentate în Figura 3.1 în care am ilustrat o pagină tipică de pe platforma YouTube. Se pot observa diferitele informații textuale, de la descrieri, tag-uri până la comentariile utilizatorilor asociate unei secvențe video. Toate aceste informații sunt informații relative la conținut.

O altă sursă de informație textuală o constituie textul conținut chiar de datele multimedia. Avantajul acestuia îl constituie faptul că poate fi extras folosind metode automate. O primă sursă este informația vizuală, ca de exemplu textul încrustat în imagine, scrisul de mână, subtitrările filmelor (în

¹⁵metadatale sunt definite uzual ca fiind "date despre date", sau altfel spus, date care descriu alte date, de orice fel și de orice tip. Cu alte cuvinte, metadatale oferă informații suplimentare la o serie de date. De exemplu, o pagină web, pe lângă textul propriu-zis poate conține metadata ce specifică limba în care este scrisă, modul de creare al paginii, diferite surse adiționale de informații și așa mai departe.

¹⁶sistemul GPS - Global Positioning System reprezintă un sistem de poziționare geografică bazat pe sateliți ce furnizează informație de localizare și timp independent de vreme și pentru oricare poziție de pe glob, atâta timp cât există posibilitatea de captare a semnalului de la cel puțin 3 sateliți GPS.

cazul în care nu sunt disponibile separat), textul grafic (de exemplu diverse indicatoare precum denumirea unei străzi, numele unui obiect, numărul de înmatriculare al unei mașini, scorului în secvențele sportive). Extragerea acestuia presupune folosirea de tehnici de recunoaștere automată a caracterelor sau OCR ("Optical Character Recognition"¹⁷).

O a doua sursă de text o constituie informația audio și în special vorbirea (de exemplu narațiuni, dialoguri, monologuri). Aceasta poate fi recunoscută și convertită în text folosind tehnicile de recunoaștere automată a vorbirii sau Automatic Speech Recognition (ASR) [Lamel 08]. Textul obținut în acest fel oferă informații prețioase de conținut, totuși tehnicile de ASR sunt limitate pe de-o parte de diversitatea limbilor existente cât și de imposibilitatea de a furniza o transcriere eficientă în condiții de zgomot de fundal (cum este cazul filmelor).

Odată obținută informația textuală aceasta poate fi folosită direct ca descriptor de conținut. Totuși în marea parte a cazurilor informația textuală tinde să fie redundantă și de dimensiune semnificativă (de exemplu sute de mii de cuvinte) necesitând o reprezentare mai eficientă. Dintre metodele cel mai frecvent folosite putem enumera reprezentarea de tip Term Frequency–Inverse Document Frequency (TF–IDF) [Knees 09] și Bag-of-Words (B-o-W) [Wallach 06].

TF–IDF este un model statistic ce se bazează pe determinarea gradului de importanță al unui termen pentru un anumit document dintr-un corpus de date. Valoarea TF-IDF va crește proporțional cu numărul de apariții al termenului în document (term frequency) dar în același timp este compensată de frecvența de apariție a cuvântului în corpus (inverse document frequency) ceea ce ajută la verificarea a cât de comun sau rar este termenul pentru toate documentele din corpus. Informația textuală poate fi astfel sintetizată prin valorile TF-IDF pentru un set de termeni cheie prestabiliți în funcție de aplicație sau extrași chiar din document.

Modelul B-o-W este un model similar ce ține cont de frecvența de apariție a cuvintelor. În acest model textul este reprezentat sub forma unei colecții, ne-ordonate, de cuvinte ("bag of words"), ignorând astfel orice reguli gramaticale. Pe baza acestei reprezentări se alcătuieste mai întâi un dicționar de cuvinte eliminând cuvintele care se repetă. Descriptorul textual pentru un anumit document va consta astfel în reprezentarea sub formă de histogramă a numărului de apariții ale fiecărui cuvânt din dicționar în documentul respectiv. Documentele ce descriu date similare vor avea o frecvență comparabilă

¹⁷recunoașterea automată a caracterelor reprezintă procesul mecanic sau electronic de traducere a imaginilor ce conțin scris de mână, scris de mașină sau text imprimat (de regulă rezultate în urma procesului de scanare) în text editabil de către calculator.

a anumitor termeni.

3.4 Descriere semantică sau sintactică?

În general descriptorii de conținut obținuți în urma adnotării conținutului datelor multimedia pot fi clasificați în funcție de nivelul semantic al informațiilor furnizate în trei categorii:

Descriptori sintactici ("low-level"), constau de regulă în adnotarea datelor cu descrieri numerice. Acest mod de descriere corespunde în general primelor sisteme de indexare (cu toate acestea multe dintre metode sunt folosite și în sistemele existente - vezi secțiunile anterioare). Adnotarea sintactică este definită generic ca fiind adnotarea ce se referă la *relațiile dintre unitățile de nivel scăzut constituente ale datelor multimedia și modul de constituire a structurii acestora*. Aceasta se poate realiza pe baza atributelor numerice, de nivel semantic redus, ca de exemplu parametri statistici calculați la nivel de pixel sau regiuni de pixeli, proprietăți geometrice ale obiectelor, structura temporală a unei secvențe sau vectori de mișcare. De regulă, descriptorii obținuți în urma procesului de adnotare sunt valori numerice ce descriu atribute de tipul celor enumerate mai sus dar și relațiile sintactice ce pot exista între acestea. Extrași la acest nivel de percepție, descriptorii sintactici sunt dificil accesibili utilizatorului de rând. De exemplu, căutarea unei imagini în funcție de procentul de apariție al unei culori sau a unei secvențe de imagini care să conțină 30% mișcare de translație și 20% mișcare de rotație, nu constituie o descriere prea relevantă pentru utilizator.

Descriptori simbolici ("mid-level"), aceștia corespund unui nivel de descriere intermediar, ce se găsește între cele două extreme: numeric și semantic, ca de exemplu denumirea culorilor dintr-o imagine, detectarea unei scene de dialog sau a prezenței umane în scenă, identificarea unui anumit tip de conținut. De regulă descriptorii de nivel semantic intermediar sunt determinați, indirect, pe baza descrierilor sintactice.

Descriptori semantici ("high-level"), în contrast cu adnotarea sintactică, adnotarea semantică a conținutului presupune o descriere perceptuală ce tinde să atingă un nivel similar cu nivelul de percepție uman. Informațiile numerice obținute în urma analizei sintactice pot fi convertite în concepte semantice precum conceptele lingvistice folosind informații "a priori" despre conținutul datelor, și/sau trecând printr-o etapă intermediară de descriere simbolică. Un sistem semantic este definit generic ca fiind *orice sistem ce implică o colecție de simboluri (vocabularul sistemului), reguli ce permit con-*

stituirea de propoziții, reguli de desemnare și reguli de validare. În cazul sistemelor de indexare, termenul de "semantic" își conservă acest sens. Acesta se traduce prin *codarea interpretării datelor pentru a servi unei aplicații specifice* [Smeulders 00]. Astfel, descrierea semantică implică existența unui *set de simboluri și reguli* ce permit interpretarea lingvistică a anumitor evenimente sau proprietăți ale datelor multimedia.

Acest mod de descriere presupune dezvoltarea de tehnici capabile să furnizeze o înțelegere completă a conținutului necesitând de cele mai multe ori o abordare multimodală (imagine-sunet-text). De exemplu, dacă ne limităm în a folosi doar informația furnizată de o imagine, să luăm cazul unei imagini ce surprinde un jucător de fotbal, singurele caracteristici ce reies din analiza imaginii sunt fizionomia acestuia și prezența sa în scenă. Pe de altă parte, dacă dispunem de secvența ce îl surprinde pe jucător, putem determina dacă acesta va marca golul, modul în care acesta joacă, contextul înregistrării, cum ar fi meciul despre care este vorba și așa mai departe, informații semantice esențiale pentru înțelegerea conținutului secvenței. În ciuda dificultății sporite de generare automată, acest mod de reprezentare al datelor este unul dintre cele mai eficiente și constituie direcția actuală de cercetare în domeniu.

Pentru a înțelege mai bine diferența dintre cele trei categorii de adnotări de conținut, în Figura 3.9 am ilustrat un exemplu concret de adnotare sintactică, simbolică și respectiv semantică în cazul unei secvențe de fotbal (din motive de vizualizare, secvența este reprezentată prin ilustrarea a câtorva imagini reprezentative).

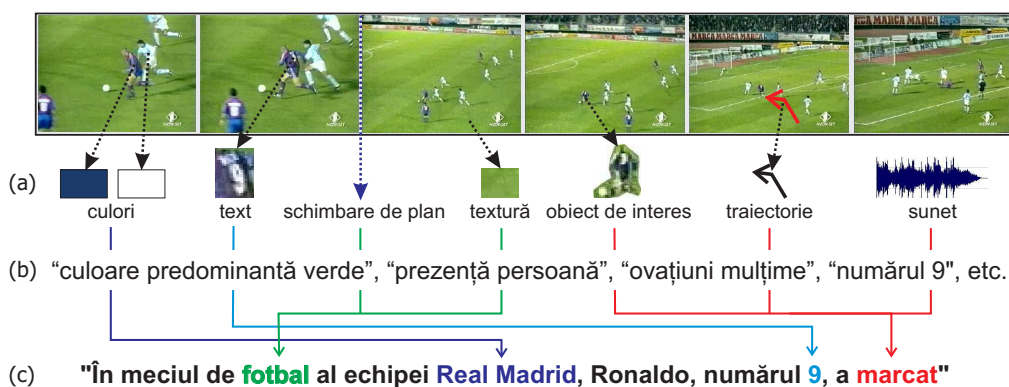


Figura 3.9: Exemplu de descriere sintactică (a), simbolică (b) și semantică (c) în cazul unei secvențe de imagini (axa orizontală reprezintă axa temporală).

Astfel, pornind de la informația video (imagine-sunet), adnotarea sintactică va fi capabilă să furnizeze doar informații relative la scenă și la pro-

prietățile acesteia, precum culoare, prezență text, textură, traiectoria obiectelor în mișcare, ritmul de desfășurare al acțiunii sau detecție zgomot audio specific mulțime.

Folosind aceste informații se poate obține o descriere simbolică de nivel semantic intermediar al conținutului video precum detecția culorii predominante ce corespunde gazonului, detecția unei persoane în mișcare, detecția ovațiilor mulțimii specifice unui gol, detecția tricoului cu numărul 9 și așa mai departe. Aceste informații, nu sunt simple date numerice dar totuși nu furnizează o înțelegere semantică a conținutului secvenței.

O adnotare semantică va da sens acestor informații într-un mod unitar, de exemplu textura verde va indica că este vorba despre un meci de fotbal, culorile jucătorilor (obiecte în mișcare) vor dezvălui echipele, recunoașterea numerelor de pe tricou va identifica jucătorii, segmentarea obiectului de interes, urmărirea acestuia și prezența zgomotului specific mulțimii vor indica marcarea golului. Astfel că sistemul va ”înțelege” sensul acțiunii secvenței și anume că este vorba despre un meci de fotbal al echipei Real Madrid în care jucătorul cu numărul 9, Ronaldo, marchează.

CAPITOLUL 4

Fuziunea datelor

În cele mai multe dintre cazuri, pentru reprezentarea conținutului multimedial este necesară combinarea mai multor tipuri de descriptori. De exemplu, conținutul unei secvențe de imagini poate fi reprezentat atât pe baza structurii temporale, cât și folosind descriptori de mișcare, descriptori audio și așa mai departe. Strategiile de fuziune a datelor se bazează pe ipoteza conform căreia o decizie obținută pe baza mai multor descriptori poate oferi performanțe superioare unei decizii bazate pe un singur tip de descriptor.

Astfel, se pune problema găsirii unei modalități de agregare (fuziune) a acestor date, formând în general un nou descriptor ce sintetizează cât mai bine puterea discriminatorie a descriptorilor individuali.

Cu alte cuvinte, ideal, noul descriptor trebuie să păstreze acele proprietăți distincte ale descriptorilor individuali (de exemplu informația audio descrie proprietăți diferite față de informația structurală) și să elimine informațiile redundante (similare), exploatând cât mai bine complementaritatea acestora în reprezentarea informației. În general există două tipuri de abordări ale problemei fuziunii datelor, tehnici de tip ”*early fusion*” și respectiv ”*late fusion*” [Snoek 05].

4.1 Metode de tip ”early fusion”

Tehnicile de tip ”early fusion” realizează agregarea datelor ”timpuriu” în lanțul de prelucrare, înainte de a fi folosite la indexare sau în alte procese de analiză. Fuziunea datelor are loc în spațiul de caracteristici (vezi Secțiunea

2.1) și constă practic în concatenarea propriu-zisă a tuturor descriptorilor fără a ține cont de redundanța acestora.

De exemplu, dacă obiectul multimedia X este descris de descriptorii de conținut $desc_1 = \{a_1, a_2, \dots, a_n\}$, $desc_2 = \{b_1, b_2, \dots, b_m\}$ și respectiv $desc_3 = \{c_1, c_2, \dots, c_l\}$, unde a , b și c reprezintă valorile atributelor acestora, descriptorul agregat este dat de concatenarea valorilor și anume $desc_{e-f} = \{a_1, \dots, a_n, b_1, \dots, b_m, c_1, \dots, c_l\}$. Acesta definește astfel un nou spațiu de caracteristici $(n + m + l)$ -dimensional.

O problemă care apare o reprezintă necesitatea normalizării valorilor datelor într-un anumit interval comun. Descriptorii diferiți tind să aibă intervale de variație diferite ale valorilor, de la normalizări diferite, de exemplu valori între $[0; 1]$ sau $[a; b]$ (unde a și b sunt două valori cunoscute) până la intervale de valori variabile și care depind de tipul datelor.

Dintre tehnicile de normalizare cel mai frecvent folosite putem enumera normalizarea min-max:

$$a_i = \frac{a_i - \min\{a_i\}}{\max\{a_i\} - \min\{a_i\}} \quad (4.1)$$

unde a_i sunt attributele descriptorului, $i = 1, \dots, n$ cu n numărul de valori ale acestuia, $\min\{a_i\}$ și $\max\{a_i\}$ reprezintă operatorii ce returnează valoarea minimă și respectiv maximă a tuturor valorilor descriptorilor (pentru toate obiectele multimedia considerate) pentru atributul a_i . Calculată în acest fel, normalizarea min-max asigură o normalizare a valorilor în intervalul $[0; 1]$.

Normalizarea z-score se folosește de calculul abaterii pătratice medii:

$$a_i = \frac{a_i - \text{medie}\{a_i\}}{\sigma\{a_i\}} \quad (4.2)$$

unde ca și în cazul anterior, operatorii $\text{medie}\{a_i\}$ și $\sigma\{a_i\}$ returnează valoarea medie și respectiv abaterea pătratică medie a tuturor valorilor descriptorilor pentru atributul a_i . În acest caz normalizarea se realizează pe o distribuție de medie zero și dispersie unu.

O altă abordare constă în calculul statisticii mediane:

$$a_i = \frac{a_i - \text{median}\{a_i\}}{\text{median}\{|a_i - \text{median}\{a_i\}|\}} \quad (4.3)$$

unde operatorul $\text{median}\{a_i\}$ returnează statistica mediană¹ a mulțimii tuturor valorilor descriptorilor pentru atributul a_i iar operatorul $|\cdot|$ returnează valoarea absolută.

¹valoarea mediană a unei mulțimi se obține prin ordonarea valorilor acesteia în ordine crescătoare și alegerea valorii de mijloc.

Dacă ordinul intervalului de variație al valorilor descriptorului diferă foarte mult, ca de exemplu printr-un ordin de mărime logaritm, $[0; 1]$ comparativ cu $[0; 1000]$, normalizarea se poate realiza folosind scalarea zecimală:

$$a_i = \frac{a_i}{10^n}, \quad n = \log_{10}(\max\{a_i\}) \quad (4.4)$$

În cazul în care nu se cunoaște intervalul de variație al valorilor descriptorului se poate opta pentru o normalizare folosind funcții duble sigmoide:

$$a_i = \left[1 + \exp\left(-2 \cdot \frac{a_i - t}{r}\right) \right]^{-1} \quad (4.5)$$

unde t este de regulă valoarea medie a distribuției valorilor descriptorului iar $r = r_1$ dacă $a_i < t$ și $r = r_2$ în caz contrar. Constantele r_1 și r_2 reprezintă valorile unor intervale alese la dreapta și respectiv stânga valorii lui t . Aceste aspecte sunt ilustrate în Figura 4.1.

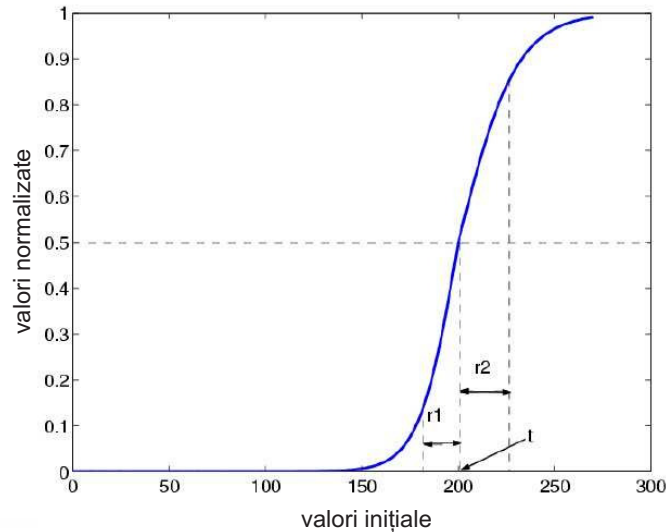


Figura 4.1: Exemplu de normalizare folosind funcții dublu sigmoide (axa oX corespunde valorilor inițiale iar axa oY valorilor normalizate).

Principalul dezavantaj al tehnicilor de tip "early fusion" este dat de dimensionalitatea datelor, descriptorul obținut prin fuziune având ca dimensiune suma dimensiunilor descriptorilor individuali, ceea ce conduce la un număr semnificativ de valori (un astfel de descriptor agregat în cazul video poate avea uzual zeci de mii de componente).

Cu cât dimensiunea datelor este mai ridicată cu atât este mai probabil ca puterea discriminatorie să scadă deoarece datele similare tind să se disperseze în spațiul de caracteristici ceea ce face dificilă separarea acestora (și astfel indexarea). De asemenea, folosind concatenarea descriptorilor nu se poate controla contribuția pe care o are fiecare descriptor individual asupra sistemului. Descriptorii de dimensiune mai mare vor tinde să aibă pondere principală în reprezentarea conținutului raportat la descriptorii cu un număr redus de valori (de exemplu descriptorii care conțin o singură valoare).

4.2 Metode de tip "late fusion"

Pe de altă parte, tehnicile de tip "late fusion" realizează fuziunea "târziu" în lanțul de prelucrare bazându-se pe exploatarea individuală a puterii discriminatorii a fiecărui descriptor sau modalități în parte.

Să considerăm pentru exemplificare un sistem de indexare după conținut bazat pe clasificarea datelor. Tehnicile de clasificare sunt tehnici de învățare asistată de calculator ("machine learning"). Problema pe care o rezolvă poate fi formulată în felul următor: având la dispoziție un set necunoscut de date se dorește realizarea unei partiționări a acestora în clase de similaritate (etichetarea acestora ca aparținând unei anumite categorii). Pentru aceasta, sistemul poate dispune de o serie de exemple de partiții, numite și date de antrenare - date pentru care se cunoaște apartenența la clase și pentru care problema clasificării este deja soluționată (de regulă de către un expert). Pe baza datelor de antrenare, clasificatorul învață mecanismul de asociere în clase urmând să-l aplice ulterior datelor noi necunoscute (ne-etichetate) [Witten 05]. Principiul este ilustrat schematic în Figura 4.2.

În contextul indexării după conținut, tehnicile de clasificare transpun problema căutării într-o problemă inversă de partiționare a bazei de date în funcție de conținutul căutat. Problema indexării se transpune astfel într-o problemă de partiționare adecvată a datelor în categoriile căutate de utilizator. De exemplu, dacă se dorește căutarea unui anumit gen video, baza de date va fi clasificată după diferite clase de gen (de exemplu film, muzică, știri), dacă se dorește găsirea unui anumit obiect într-o bază de imagini, acestea vor fi clasificate în clase de obiecte (de exemplu "minge", "mașină", "casă"). Procesul de clasificare se realizează pe baza reprezentării datelor cu descriptorii de conținut.

În momentul căutării datelor, cererea de căutare a utilizatorului ("query") va fi asociată uneia dintre clasele determinate anterior, rezultatele căutării fiind acele documente ce au fost etichetate ca aparținând acestei clase. Ca și în cazul mecanismului de indexare clasic (vezi Secțiunea 2) datele vor fi

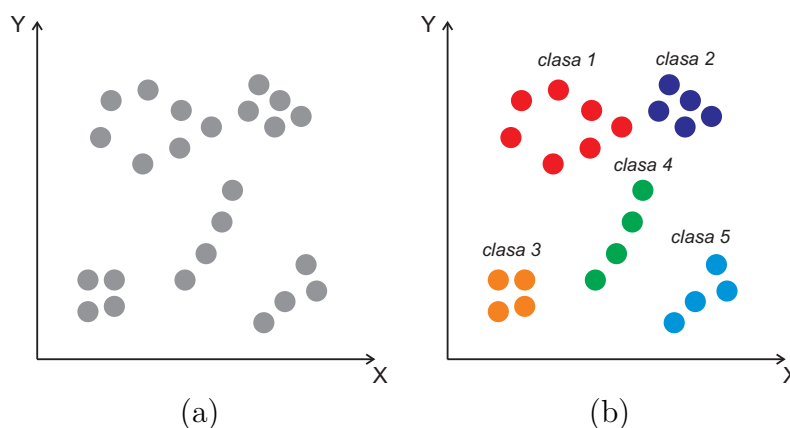


Figura 4.2: Principiul clasificării datelor: (a) datele de intrare reprezentate în spațiul de caracteristici, (b) repartizarea în clase obținută în urma clasificării (obiectele din aceeași clasă sunt reprezentate cu aceeași culoare).

returnate în ordinea descrescătoare a relevanței. Pentru ca acest lucru să fie posibil, clasificatorul în locul unei decizii binare de apartenență sau non-apartenență va furniza un grad de relevanță, de regulă o valoare reală în intervalul $[0; 1]$, unde 1 reprezintă apartenența sigură la clasă, iar 0 cazul contrar. Astfel, rezultatele sunt returnate utilizatorului în ordinea descrescătoare gradului de relevanță furnizat de clasificator pentru clasa ce aparține căutării. Acest mecanism este exemplificat în Figura 4.3 în contextul căutării după gen a documentelor video.

Revenind la problematica fuziunii datelor, fuziunea de tip "late fusion" se realizează în acest caz prin fuzionarea rezultatelor clasificărilor obținute independent pentru fiecare tip de descriptor sau modalitate, cât și pentru tipuri de clasificatori diferiți. În acest fel, agregarea datelor nu este realizată la nivel de descriptor ci la nivelul gradului de relevanță atribuit de fiecare clasificator descriptorilor, beneficiind de puterea discriminatorie a fiecărui descriptor în parte. Dintre tehnicile de tip "late fusion" cel mai frecvent folosite putem enumera:

- **fuziunea paralelă:** presupune rularea aceluiași sistem în paralel pentru descriptori și tipuri de clasificatori diferiți. Agregarea finală a rezultatelor se realizează pe baza agregării rezultatelor obținute individual (vezi Figura 4.4);
- **fuziunea serială:** presupune rularea în cascadă a sistemelor, fiecare ieșire a unui clasificator fiind folosită la intrarea unui alt clasificator ca de exemplu pentru clasificarea datelor ce au fost clasificate eronat

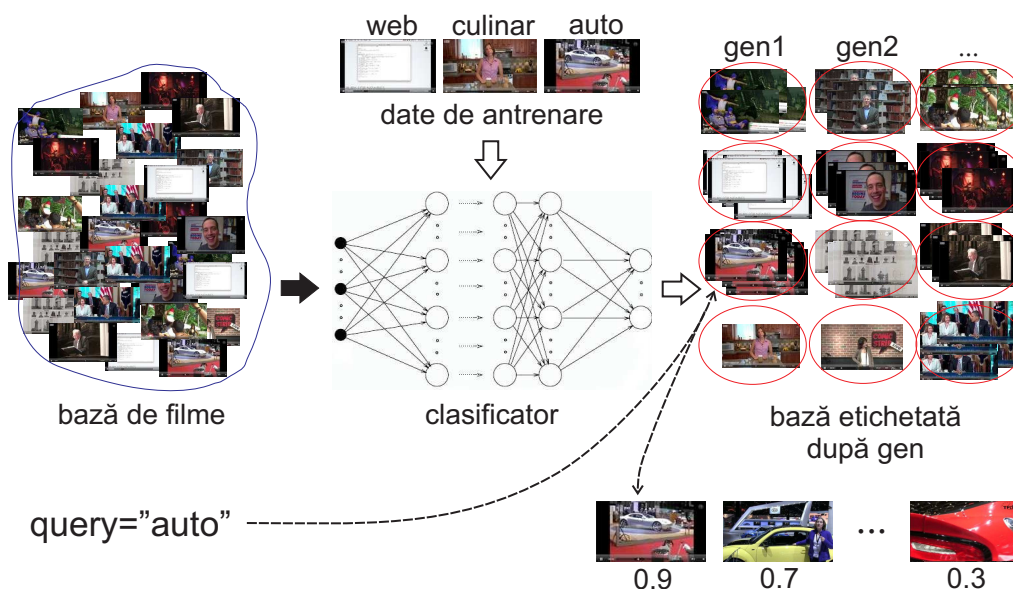


Figura 4.3: Exemplu de sistem de clasificare folosit pentru căutarea după gen a secvențelor video. Clasificatorul este mai întâi antrenat folosind un set redus de exemple și un set predefinit de genuri urmând să catalogheze automat baza video necunoscută. Cererea de căutare primește ca rezultat secvențele ce au fost atribuite clasei căutată în ordinea descrescătoare a gradului de relevanță furnizat de clasificator (sursă imagini blip.tv).

de sistemul anterior. Fiecare dintre sisteme rulează pentru descriptori și clasificatori diferiți. Principiul este inspirat de tehnicile de tip "boosting" în care mai mulți clasificatori "slabi" (cu performanțe reduse) sunt combinați pentru a obține un clasificator cu performanțe ridicate (vezi AdaBoost [Witten 05]);

- **fuziunea ierarhică:** sistemele sunt organizate ierarhic, fie de tip "bottom-up" în care mai mulți clasificatori converg către un clasificator final, sau de tip "top-down" unde în funcție de rezultatele unui clasificator inițial, deciziile se separă ierarhic pe mai multe niveluri de clasificatori. Acest mod de reprezentare este similar arborilor de decizie (vezi Random Forest sau Random Tree [Witten 05]).
- **fuziunea mixtă:** constă în combinarea mai multor modalități de fuzionare din categoriile enumerate anterior.

În continuare vom detalia modul de luare al deciziei în cazul fuzionării paralele. Acesta este ilustrat în Figura 4.4. Având la dispoziție N clasi-

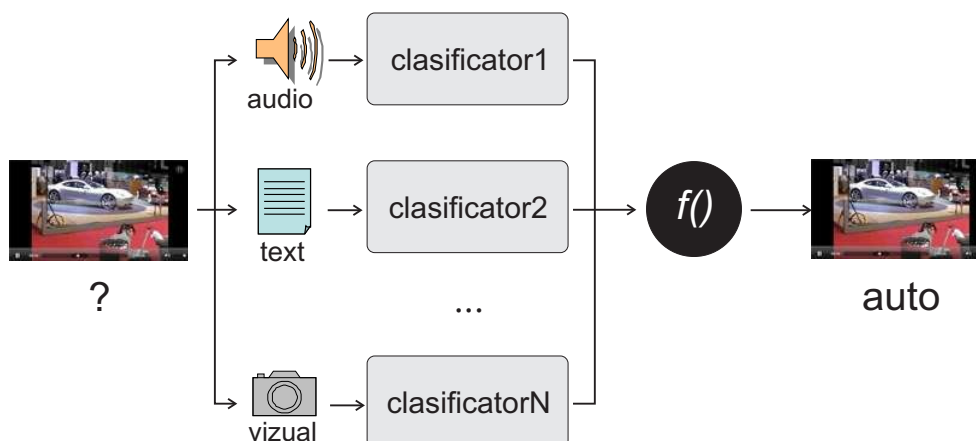


Figura 4.4: Principiul fuzionării de tip "late fusion" paralel. Catalogarea datelor de intrare se realizează pe baza unei funcții de agregare, $f(\cdot)$, a ieșirilor mai multor tipuri de clasificatori antrenați folosind descriptori diferiți.

ficatori ce sunt antrenați folosind descriptori de conținut diferiți, fuzionarea de tip "late fusion" a descriptorilor presupune determinarea unei funcții care combină gradele de relevanță furnizate de fiecare clasificator în parte, $f(x_1, \dots, x_N)$, unde x_i reprezintă gradul de relevanță atribuit de clasificatorul i datelor de intrare. Acestea reprezintă probabilitățile de apartenență la clasele considerate, $x_i = \{p_{i,c_1}, p_{i,c_2}, \dots, p_{i,c_M}\}$ unde c_1, \dots, c_M reprezintă clasele considerate (de exemplu genurile video în Figura 4.3) iar $p_{i,c}$ reprezintă probabilitatea ca datele să fie atribuite ca aparținând clasei c .

În mod natural, fiecare clasificator va tinde să furnizeze grade de apartenență diferite fiind antrenat pentru descriptori diferiți. Funcția $f(\cdot)$ trebuie determinată în așa fel încât rezultatele obținute de clasificatorul agregat să fie cât mai bune și superioare fiecărui clasificator individual. Agregarea se va realiza pentru gradele de relevanță ale fiecărei clase în parte.

O modalitate de definire a lui $f(\cdot)$ o reprezintă combinația liniară a gradelor de relevanță și anume:

$$f_{CombMean}(d, c_j) = \sum_{i=1}^N \alpha_i \cdot p_{i,c_j} \quad (4.6)$$

unde d reprezintă documentul curent, p_{i,c_j} reprezintă probabilitatea de apartenență la clasa c_j , $j = 1, \dots, M$ cu M numărul de clase considerate, atribuită de clasificatorul i iar α_i reprezintă un set de ponderi. Un caz particular îl reprezintă considerarea de ponderi egale ceea ce conduce la însumarea gradelor de relevanță pentru fiecare clasă.

Un alt exemplu este atribuirea unei ponderi superioare acelor date care sunt mai probabile să fie relevante pentru o clasă, astfel:

$$f_{CombMNZ}(d, c_j) = F(d)^\gamma \cdot \sum_{i=1}^N \alpha_i \cdot p_{i,c_j} \quad (4.7)$$

unde $F(d)$ reprezintă numărul de clasificatori pentru care documentul d apare în primele k documente din punct de vedere al valorii de relevanță (k este o constantă stabilită a priori) iar $\gamma \in [0, 1]$ este un parametru de control.

Noile valori de relevanță obținute în urma agregării sunt folosite mai departe pentru indexarea datelor în mod similar în care acestea erau folosite în cazul considerării unui singur clasificator.

Comarate cu abordările de tip "early fusion", tehnicile de tip "late fusion" sunt mai avantajoase din punct de vedere computațional deoarece agregarea se face folosind dimensiunea inițială a descriptorilor. Este mai eficientă clasificarea unor descriptori de dimensiuni reduse și agregarea rezultatelor decât clasificarea unui descriptor agregat de dimensiuni semnificativ mai mari. Principalul dezavantaj al acestor metode este totuși dat de pierderea eventualei corelații dintre descriptori ce se obține în cazul concatenării acestora și care poate furniza un nivel de discriminare superior folosirii individuale a acestora.

În ciuda diferențelor dintre cele două abordări, "early fusion" și respectiv "late fusion", nu există o metodă preferențială în defavoarea celeilalte, ambele abordări dovedindu-se eficiente în contexte diferite. Astfel că tehnica de fuziune a datelor rămâne dependentă de aplicație [Lan 12].

CAPITOLUL 5

Conceptul de similaritate a datelor

Așa cum am prezentat în Secțiunea 2.3, în procesul de căutare după conținut a datelor, descrierea eficientă a conținutului nu este suficientă pentru a asigura indexarea acestora în baza de date. La fel de importantă este definirea conceptului de similaritate (sau opus, disimilaritate) dintre date sau dintre descriptorii acestora.

Practic identificarea rezultatelor căutării se realizează prin localizarea datelor ce sunt "similare" până la un anumit nivel cu cererea de căutare ("query"). Cu alte cuvinte este necesară definirea unei funcții, $S(O_1, O_2)$, capabilă să evalueze în ce măsură două obiecte multimedia, O_1 și O_2 , arată sau sună în mod similar, în ce măsură au o structură similară sau în ce măsură conduc la aceeași percepție sau interpretare a conținutului [Worring 03].

În general, evaluarea similarității dintre date se poate realiza fie la nivel de *descriptori*, la nivel de *structură* ("layout") sau la nivel *semantic*, fie folosind combinații ale acestora.

5.1 Similaritatea descriptorilor

În acest caz, similaritatea datelor este evaluată numeric folosind valorile descriptorilor de conținut aferente acestora iar funcția $S()$ este de regulă o măsură de distanță (metrică). Datele vor fi considerate similare în măsura în care valoarea distanței dintre descriptorii acestora este minimă.

În cele ce urmează vom face o trecere în revistă a diverselor metrici folosite în domeniul căutării informației. Marea parte dintre acestea sunt în mod

natural inspirate din matematică [Deza 06].

Una dintre abordările clasice este folosirea distanței Minkovski, ce este definită ca:

$$S_{Mink}(A_{O_1}, A_{O_2}) = \sqrt[r]{\sum_{i=1}^n [A_{O_1}(i) - A_{O_2}(i)]^r} \quad (5.1)$$

unde $A_O(i)$ reprezintă valoarea de indice i a descriptorului aferent obiectului multimedia O , cu $i = 1, \dots, n$ elemente (de regulă descriptorii de conținut sunt vectori n -dimensionali de valori, vezi și Secțiunea 2.1).

În cazul în care considerăm parametrul $r = 1$ obținem norma L1 sau distanța Manhattan:

$$S_{Manh}(A_{O_1}, A_{O_2}) = \sum_{i=1}^n |A_{O_1}(i) - A_{O_2}(i)| \quad (5.2)$$

unde operatorul $|\cdot|$ reprezintă valoarea absolută.

Pentru $r = 2$ obținem mai departe norma L2 cunoscută sub numele de distanța Euclidiană:

$$S_{Euclid}(A_{O_1}, A_{O_2}) = \sqrt{\sum_{i=1}^n [A_{O_1}(i) - A_{O_2}(i)]^2} \quad (5.3)$$

În cazul în care nu toate elementele descriptorului au aceeași importanță, distanța dintre fiecare pereche de valori poate fi ponderată diferit obținând astfel distanța Euclidiană ponderată:

$$S_{wEuclid}(A_{O_1}, A_{O_2}) = \sqrt{\sum_{i=1}^n w_i \cdot [A_{O_1}(i) - A_{O_2}(i)]^2} \quad (5.4)$$

unde w_i , cu $i = 1, \dots, n$ reprezintă ponderile fiecărei valori.

O altă măsură de distanță ce este folosită de regulă când descriptorii de conținut sunt reprezentați sub formă de histograme (de exemplu histograma color a unei imagini) o constituie intersecția histogramei. Aceasta este de fapt o măsură a disimilarității și este definită ca suma minimelor valorilor histogramelor:

$$S_{inter}(h_{O_1}, h_{O_2}) = \sum_{i=1}^n \min\{h_{O_1}(i), h_{O_2}(i)\} \quad (5.5)$$

unde $h_O(i)$ cu $i = 1, \dots, n$ reprezintă histograma color a obiectului multimedia O iar operatorul $\min\{\cdot\}$ returnează valoarea minimă a unui set de elemente.

Tot în cazul evaluării diferențelor dintre histogramme și în special dintre histogrammele color ale imaginilor, în cazul folosirii distanțelor clasice, este foarte probabil ca pentru distribuții ale unei aceleiași nuanțe (de exemplu roșu deschis și roșu) să obținem valori semnificative ale distanței, de ordin de măsură similar ca pentru distanța față de o distribuție a unei nuanțe complet diferite (de exemplu albastru), în ciuda faptului că diferențele în primul caz ar trebui să fie reduse, culorile fiind asemănătoare. O distanță care tinde să contracareze acest efect este distanța pătratică dintre histogramme:

$$S_{hist2}(h_{O_1}, h_{O_2}) = \sqrt{(h_{O_1} - h_{O_2})^T \cdot A \cdot (h_{O_1} - h_{O_2})} \quad (5.6)$$

unde h_O reprezintă vectorul histogramă cu n elemente, T reprezintă transpusa unei matrice iar $A = [a_{i,j}]$, $i, j = 1, \dots, n$, reprezintă o matrice pătratică de valori ce indică corelația dintre elementele histogramelor de indici i cu cele de indice j (de regulă A este simetrică și are elementele de pe diagonala principală egale cu 1).

Alte măsuri de distanță frecvent folosite sunt distanța Canberra:

$$S_{Candb}(A_{O_1}, A_{O_2}) = \sum_{i=1}^n \frac{|A_{O_1}(i) - A_{O_2}(i)|}{|A_{O_1}(i)| + |A_{O_2}(i)|} \quad (5.7)$$

distanța Bray-Curtis:

$$S_{B-C}(A_{O_1}, A_{O_2}) = \frac{\sum_{i=1}^n |A_{O_1}(i) - A_{O_2}(i)|}{\sum_{i=1}^n [A_{O_1}(i) + A_{O_2}(i)]} \quad (5.8)$$

distanța SquaredChord:

$$S_{SChord}(A_{O_1}, A_{O_2}) = \sum_{i=1}^n \left[\sqrt{A_{O_1}(i)} - \sqrt{A_{O_2}(i)} \right]^2 \quad (5.9)$$

distanța Lorentzian, Soergel, Czekanowski, Motyka, Ruzicka, Tanimoto, Wave-Hadges, Clark, Person și așa mai departe. Pentru mai multe detalii cititorul se poate raporta la [Deza 06].

O abordare diferită este distanța Bhattacharyya ce măsoară similaritatea a două distribuții de probabilitate. În cazul în care descriptorii sunt considerați a avea o distribuție normală Gaussiană, distanța poate fi scrisă ca fiind:

$$S_{Bhatta}(A_{O_1}, A_{O_2}) = \frac{1}{8} \cdot (\mu_{A_{O_1}} - \mu_{A_{O_2}})^T \cdot (\Sigma_{O_1, O_2})^{-1} \cdot (\mu_{A_{O_1}} - \mu_{A_{O_2}}) + \frac{1}{2} \cdot \ln \left(\frac{\det(\Sigma_{O_1, O_2})}{\sqrt{\det(\Sigma_{O_1}) \cdot \det(\Sigma_{O_2})}} \right) \quad (5.10)$$

unde μ_{A_O} reprezintă vectorul medie al distribuției de probabilitate a descriptorului A_O , Σ_O reprezintă matricea de covarianță a distribuției lui A_O , Σ_{O_1, O_2} reprezintă media aritmetică a matricelor de covarianță pentru distribuțiile lui A_{O_1} și A_{O_2} (vezi și [Ciuc 05]), T reprezintă transpusa unei matrice iar operatorul $\det(\cdot)$ returnează determinantul unei matrice.

O altă perspectivă o constituie reprezentarea datelor sub formă de mulțimi. Distanța Hausdorff evaluează gradul de apropiere a două submulțimi într-un anumit spațiu și folosind o anumită metrică, astfel:

$$S_{Haus}(A_{O_1}, A_{O_2}) = \max\left\{ \sup_i \inf_j d(A_{O_1}(i), A_{O_2}(j)), \sup_j \inf_i d(A_{O_1}(i), A_{O_2}(j)) \right\} \quad (5.11)$$

unde $i, j = 1, \dots, n$, \inf și \sup reprezintă infimum și respectiv supremum al unei mulțimi (de regulă valoarea minimă și respectiv maximă), $d(\cdot)$ reprezintă o anumită metrică (de exemplu norma L1) iar $\max\{\cdot\}$ returnează valoarea maximă a unei mulțimi. În acest caz, valorile descriptorilor pot fi văzute din perspectiva elementelor unei mulțimi.

Un alt caz interesant este distanța cosinus. Să presupunem că descriptorii de conținut sunt vectori de caractere iar datele ce trebuie comparate sunt documente textuale, atunci similaritatea dintre acestea poate fi evaluată folosind produsul scalar:

$$S(A_{O_1}, A_{O_2}) = \sum_{i=1}^n A_{O_1}(i) \cdot A_{O_2}(i) \quad (5.12)$$

Acum dacă descriptorii textuali sunt reprezentați sub formă de histograme ale căror valori indică numărul de apariții al unui anumit cuvânt în document (eventual ponderat de un factor de importanță - cuvintele sunt alese pentru un dicționar predefinit; vezi TF-IDF în Secțiunea 3.3) atunci similaritatea se reduce la o înmulțire a valorilor histogramei pentru cele două documente. Astfel, atunci când un cuvânt apare frecvent în cele două documente, contribuția acestuia la produs va fi semnificativă.

Problema care apare este faptul că documentele mari vor conține mai multe cuvinte și vor tinde să devină mai similare decât documentele ce conțin mai puțin text. Astfel că în practică descriptorii sunt normalizați la dimensiunea acestora $\|A_O\|^2 = \sum_{i=1}^n A_O^2(i)$ ceea ce conduce la formularea distanței cosinus astfel:

$$S_{cos}(A_{O_1}, A_{O_2}) = \frac{A_{O_1} \cdot A_{O_2}}{\|A_{O_1}\| \cdot \|A_{O_2}\|} \quad (5.13)$$

unde \cdot reprezintă produsul scalar (denumirea de cosinus vine de la faptul că distanța este practic cosinusul unghiului celor doi vectori normalizați).

În cazul comparării de obiecte, de exemplu prin intermediul a două imagini binare (în care obiectul are valoarea 1 și fundalul 0) se poate folosi distanța Baddeley definită în felul următor:

$$S_{Badd}(I_{O_1}, I_{O_2}) = \left[\frac{1}{M \cdot N} \sum_{p \in S} |d_{I_{O_1}}(p) - d_{I_{O_2}}(p)|^q \right]^{1/q} \quad (5.14)$$

unde I_O reprezintă o imagine binară, $M \cdot N$ reprezintă numărul total de pixeli din setul S , $d_{I_O}(p)$ reprezintă o anumită metrică de distanță de la punctul p la cel mai apropiat punct al obiectului conținut în imaginea I_O iar q este exponentul (de regulă se consideră $q = 2$). Definită în acest fel, distanța Baddeley oferă un anumit grad de invarianță la translația obiectelor și modificarea factorului de scală.

O problemă aparte o ridică compararea descriptorilor de dimensiuni diferite, ca de exemplu histogramele color a două imagini cu palete de culoare diferite (binii histogramei și numărul acestora sunt diferite). O soluție în acest sens este propusă de distanța Earth Mover's Distance (EMD). Aceasta se bazează pe evaluarea costului minim aferent transformării unuia dintre descriptori în celălalt și este formulată ca o problemă de optimizare. EMD este definită în felul următor:

$$S_{EMD}(A_{O_1}, A_{O_2}) = \frac{\sum_{i=1}^m \sum_{j=1}^n d_{i,j} \cdot f_{i,j}}{\sum_{i=1}^m \sum_{j=1}^n f_{i,j}} \quad (5.15)$$

unde cei doi descriptori A_{O_1} și respectiv A_{O_2} au dimensiuni diferite, m și respectiv n , $d_{i,j}$ reprezintă distanța dintre valorile $A_{O_1}(i)$ și respectiv $A_{O_2}(j)$ iar $f_{i,j}$ este o funcție de cost ce reprezintă deplasarea între $A_{O_1}(i)$ și $A_{O_2}(j)$ determinată ca minimizând valoarea costului total $\sum_{i=1}^m \sum_{j=1}^n d_{i,j} \cdot f_{i,j}$ cu o serie de constrângeri [Rubner 00].

O altă categorie de distanțe sunt cele inspirate din teoria informației a lui Shannon, precum divergențele Kullback–Leibler:

$$S_{K-L}(A_{O_1}, A_{O_2}) = \sum_{i=1}^n A_{O_1}(i) \cdot \ln \frac{A_{O_1}(i)}{A_{O_2}(i)} \quad (5.16)$$

sau divergența Jeffrey:

$$S_{Jeff}(A_{O_1}, A_{O_2}) = \sum_{i=1}^n [A_{O_1}(i) - A_{O_2}(i)] \cdot [\ln(A_{O_1}(i)) - \ln(A_{O_2}(i))] \quad (5.17)$$

Acestea sunt aplicate cu precădere la compararea descriptorilor specifici datelor audio, unde este relevantă distribuția statistică a valorilor acestora.

Pentru a ilustra importanța alegerii adecvate a măsurii de distanță, în Figura 5.1 am prezentat rezultatele obținute pentru o căutare de imagini cu ”relevance feedback” (vezi Secțiunea 2.4) și folosind metrici și descriptori de conținut diferiți [Mironică 12b]. Graficele ilustrează performanța căutării pe baza valorii MAP (Mean Average Precision, vezi Secțiune 8; reprezentată pe axa oY - valoarea maximă este 1 ce indică o performanță de 100%) raportată la metrica folosită (axa oX). Pentru descrierea conținutului imaginilor au fost folosiți descriptori de trăsături de tip SIFT și SURF (vezi Secțiune 3.1). Testele au fost efectuate pe două baze de imagini, baza Microsoft Object Class Recognition¹ (puncte roșii) și respectiv Caltech-101² (puncte albastre).

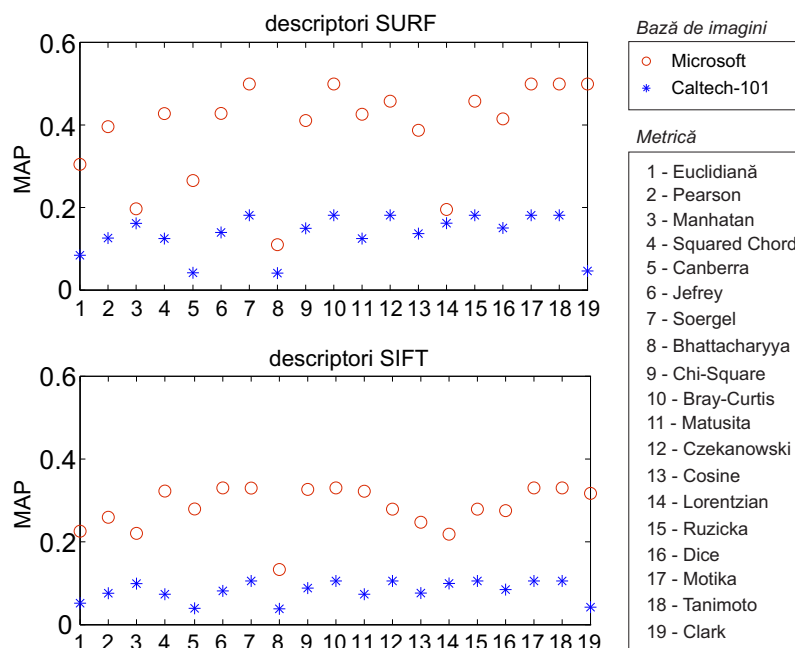


Figura 5.1: Exemplu de influență a metricii asupra performanțelor căutării de imagini [Mironică 12b] (MAP reprezintă Mean Average Precision - valoarea maximă 1 corespunde unei performanțe de 100%).

Se poate observa faptul că în funcție de metrică, performanțele sistemului variază semnificativ, de exemplu pentru baza Microsoft valorile MAP variază de la 10% la 50% pentru descriptorii SURF și de la 10% la 38% pentru SIFT. Pe lângă alegerea adecvată a descriptorilor (se observă și în figură faptul că

¹Microsoft Object Class Recognition <http://research.microsoft.com/en-us/projects/objectclassrecognition>.

²Caltech-101 http://www.vision.caltech.edu/Image_Datasets/Caltech101.

descriptorul SURF este mai performant în contextul sistemului prezentat), alegerea adecvată a metricii joacă un rol cel puțin la fel de important.

5.2 Similaritatea la nivel de structură

Aceasta presupune evaluarea gradului de similaritate a două obiecte multimedia, O_1 și O_2 , din punct de vedere al structurii acestora, ca de exemplu modul de aranjare spațială a obiectelor în imagini, modul de structurare al unei paginii de text, structura temporală a unui document video. O modalitate eficientă de caracterizare a structurii este prin intermediul descrierii acesteia cu șiruri de caractere [Worring 03].

Să considerăm în continuare exemplul datelor video. Un document video, din punct de vedere structural, este constituit ca o înșiruire de plane video separate de tranziții (vezi Secțiunea 3.1). Informația structurală poate consta în descrierea acestei structuri. Documentul video poate fi reprezentat ca un șir de caractere de genul "scswsdcscs", unde s reprezintă un plan video ("shot"), c reprezintă o tranziție de tip "cut", w reprezintă o tranziție graduală de tip "wipe" iar d reprezintă un "dissolves". Informația temporală este dată de ordinea simbolurilor în șir, astfel acest document video începe cu un plan urmat de un "cut", un plan video, o tranziție "dissolves" și așa mai departe.

Pentru a compara similaritatea descriptorilor astfel obținuți o variantă eficientă o reprezintă folosirea distanței de editare ("edit distance"), ce folosește un concept similar distanței Earth Mover's Distance (EMD) descrisă anterior. Având la dispoziție descriptorii structurali de conținut ai celor două obiecte multimedia, $A_{O_1} = \{a_{1,1}, a_{1,2}, \dots, a_{1,n}\}$ și respectiv $A_{O_2} = \{a_{2,1}, a_{2,2}, \dots, a_{2,m}\}$, unde n și m reprezintă numărul de caractere, un alfabet Σ ce descrie simbolurile posibile (valorile lui a), un set E de operații de editare și costurile aferente acestora, distanța de editare dintre A_{O_1} și A_{O_2} reprezintă costul minim de transformare a șirului A_{O_1} în șirul A_{O_2} pe baza operațiilor din E .

În cazul a două secvențe video, să presupunem că descriptorii acestora sunt $V_{O_1} = \{scswsdcscs\}$ și $V_{O_2} = \{sdsWSCscscs\}$, mulțimea $\Sigma = \{c, w, d, s\}$ iar operațiile de editare posibile sunt $E = \{inserare, ștergere, înlocuire\}$ iar costurile aferente acestora sunt egale. Operațiile necesare pentru a transforma pe V_{O_1} în V_{O_2} constau în două "înlocuiri" ale lui c cu d și respectiv două operații de "inserare" pentru adăugarea lui c și s la sfârșit. Astfel distanța de editare dintre cei doi descriptori este în acest caz 4.

5.3 Similaritatea semantică

Așa cum am prezentat și în Secțiunea 3.4, tendința actuală a sistemelor de indexare este aceea de a determina descriptorii de conținut ce oferă un nivel de înțelegere al conținutului cât mai apropiat de nivelul de percepție uman. Acest lucru se realizează în principal prin identificarea a ceea ce numim "concepte". Un concept este practic o reprezentare textuală a entităților reprezentate de obiectul multimedia, exemple de concepte în cazul imaginilor sau a documentelor video fiind "cer", "mașină", "persoană", "casă", și așa mai departe [Over 12].

Reprezentarea conceptelor poate fi realizată fie prin adnotarea manuală a acestora de către utilizator, fie folosind tehnici automate de adnotare sau folosind informații derivate din ontologii. Ontologia constituie un mod formal de reprezentare a cunoașterii sub forma unui set de concepte dintr-un domeniu și a relațiilor dintre acestea folosind următoarele componente: obiecte sau instanțe, clase (mulțimi, colecții, concepte), atribute (proprietăți, trăsături, parametri ai obiectelor și claselor), relații (descriu modul în care clasele și instanțele sunt relaționate), restricții, reguli (afirmații de tip dacă-atunci ("if-then") sau antecedent-consecvent ce descriu o serie de implicații logice), axiome și evenimente (modul de schimbare al atributelor sau al relațiilor).

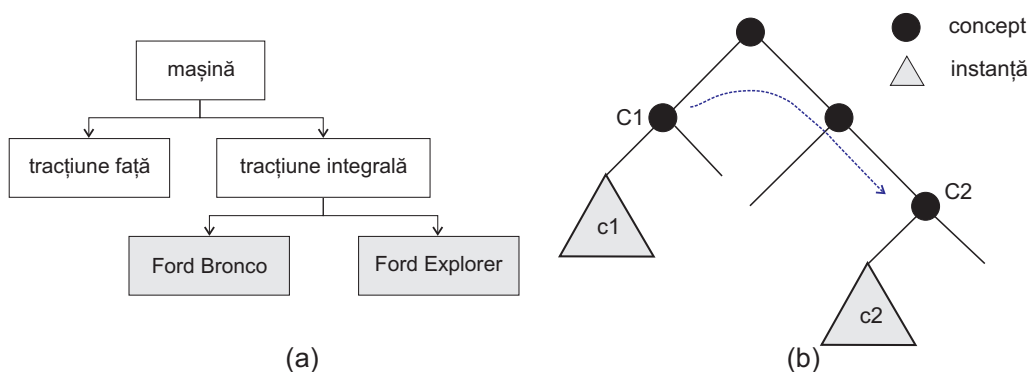


Figura 5.2: Exemple de ontologii: (a) definirea clasei "mașină" și a obiectelor "Ford Bronco" și "Ford Explorer" (exemplu din Wikipedia), (b) calculul distanței dintre două concepte, C_1 și C_2 (exemplu din [Worring 03]).

Un exemplu este ilustrat în Figura 5.2.(a) unde este prezentată o ontologie simplificată pentru clasa "mașină". O clasă poate fi subordonată unei alte clase, de exemplu clasa "mașină" poate fi considerată subclasă pentru clasa "autovehicul", deoarece toți membrii acesteia sunt implicit și membrii clasei "autovehicul", sau la rândul ei poate să conțină alte clase subordonate,

de exemplu clasele "tracțiune față" și "tracțiune integrală" în exemplul din figură.

Acest mod de reprezentare crează o structură ierarhică în care la nivelul ierarhic superior se găsesc clasele cele mai generale iar la cel inferior clasele cele mai specifice. Relațiile de subordonare implică mostenirea proprietăților claselor superioare ("părinți") către clasele inferioare ("copii"). O partiție a ontologiei reprezintă un set de clase și regulile asociate acestora ce asigură faptul că obiectele pot fi clasificate în subclasa cea mai apropiată.

De exemplu, Figura 5.2.(a) conține de fapt diagrama parțială a unei ontologii ce corespunde unei partiții a clasei "mașină" în clasele "tracțiune față" și "tracțiune integrală". Regula de partiționare determină dacă o anumită mașină poate fi clasificată în una dintre cele două subclase. În acest mod de reprezentare, obiectele sunt descrise de atribute. Tipul unui obiect și tipul atributelor determină modul de relaționare între acestea. O relație dintre un obiect și un atribut reflectă faptul că acesta este specific obiectului de care este relaționat.

În exemplul din Figura 5.2.(a), obiectul "Ford Explorer" poate conține atribute de tipul:

- <se numește> Ford Explorer,
- <are drept componentă> ușă (număr minim și maxim 4),
- <are drept componentă unul dintre> {motor 4.0 litrii, motor 4.6 litrii},
- <are drept componentă> transmisie cu 6-viteze,

Mai multe informații relative la ontologii pot fi găsite în [Gómez-Pérez 04].

În cazul comparării descrierilor semantice reprezentate sub formă de concepte, o metodă simplă constă în evaluarea distanței ce trebuie parcursă în arborele unei ontologii pentru a ajunge de la un concept la altul. Un exemplu este prezentat în Figura 5.2.(b) în care am ilustrat conceptul C_1 și C_2 în contextul unei ontologii. Având la dispoziție instanțele c_1 și respectiv c_2 ale acestor concepte, obținute de exemplu din datele multimedia ce trebuie comparate, o măsură a similarității acestora poate fi determinată ca numărul de pași necesari în arbore pentru a ajunge de la conceptul C_1 la C_2 , și anume 3 pentru acest exemplu (ce corespunde numărului de laturi ale arborelui ce trebuie parcurse, vezi săgeată în figură).

CAPITOLUL 6

Tehnicile de tip "relevance feedback"

Așa cum a fost prezentat și în Secțiunea 2.4, conceptul de "relevance feedback" în contextul sistemelor de indexare după conținut se referă la interacția cu utilizatorul în vederea îmbunătățirii rezultatelor inițiale ale căutării. În general, mecanismul de "relevance feedback" funcționează după următorul algoritm [Manning 08]:

1. **căutarea datelor dorite:** utilizatorul realizează o anumită căutare specificând datele dorite prin formularea unei cereri de căutare ("query"). Sistemul, pe baza mecanismului implementat, returnează rezultatele ce au caracteristicile cele mai apropiate de "query" folosind un anumit criteriu de similaritate. Până în acest punct, procesul este identic procesului de căutare al unui sistem de indexare (vezi Secțiunea 2);
2. **evaluare rezultate** de către utilizator: în funcție de performanța sistemului, rezultatele obținute pot fi mai mult sau mai puțin relevante pentru datele căutate. În acest punct, utilizatorul analizează rezultatele și le clasifică manual ca fiind, fie relevante pentru căutare (rezultat corect), fie ne-relevante (rezultat eronat). De regulă acest proces are loc pentru un număr limitat de rezultate, de ordinul zecilor;
3. **rafinare rezultate:** informațiile furnizate de utilizator sunt folosite drept referință ("ground truth"). Pe baza acestora, sistemul va recalcula o reprezentare mai bună a rezultatelor căutării furnizând utilizatorului o rafinare a rezultatelor în funcție de asemănarea cu datele

indicate drept relevante pentru căutare. Acest pas constituie practic algoritmul efectiv de "relevance feedback";

4. **re-iterare algoritm:** în funcție de calitatea noilor rezultate obținute, întreg procesul poate fi repetat prin reluarea punctului 2 până când rezultatele obținute sunt satisfăcătoare pentru utilizator sau îndeplinesc un anumit criteriu de performanță.

Un exemplu este prezentat în Figura 6.1 pentru căutarea de imagini în baza de date Microsoft Object Class Recognition¹.

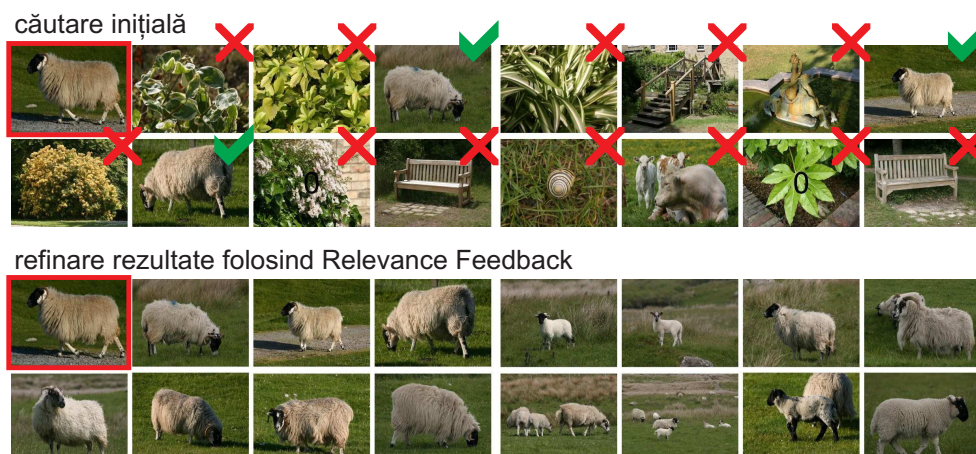


Figura 6.1: Exemplu de "relevance feedback" în cazul căutării de imagini (metoda propusă în [Mironică 12a]). Imaginile de mai sus reprezintă rezultatele sistemului de indexare pentru căutarea imaginilor similare cu imaginea marcată de chenarul roșu (\checkmark și \times reprezintă rezultatele corecte și respectiv eronate marcate de utilizator) în timp ce imaginile din partea de jos reprezintă rafinarea rezultatelor cu "relevance feedback".

Cererea de căutare a fost formulată prin furnizarea unei imagini exemplu (vezi imagine marcată de dreptunghiul roșu). În prima parte a figurii sunt ilustrate rezultatele obținute inițial de sistemul de indexare (din motive de spațiu sunt prezentate doar primele 15 rezultate; rezultatele sunt afișate în ordine descrescătoare a similarității, de la stânga la dreapta și de sus în jos). Pentru aceste rezultate utilizatorul a marcat imaginile corecte (simbol \checkmark) și respectiv cele eronate (simbol \times).

În partea de jos a figurii sunt prezentate rezultatele rafinate în urma aplicării metodei de "relevance feedback" ierarhic propusă în [Mironică 12a].

¹vezi notă de subsol 1.

Se poate observa o îmbunătățire semnificativă a performanțelor sistemului, imaginile returnate în acest caz corespunzând în totalitate cererii de căutare.

În funcție de modul în care sunt preluate informațiile de la utilizator ("feedback"), întâlnim trei tipuri de algoritmi:

- **"feedback" explicit:** corespunde algoritmului descris anterior în care utilizatorul el însuși specifică care dintre rezultate sunt corecte și care sunt eronate;
- **"feedback blind"** sau pseudo-feedback: presupune simularea interacției cu utilizatorul și se bazează pe ipoteza conform căreia sistemul de indexare este suficient de performant (prin prisma descriptorilor de conținut folosiți și al mecanismului de căutare) astfel încât este foarte probabil ca primele rezultate returnate să conțină un număr semnificativ de rezultate corecte. În acest caz, interacția cu utilizatorul este substituită prin considerarea implicită a primelor k rezultate drept relevante [Larson 10]. Pe măsură ce datele căutate devin din ce în ce mai complexe (exemplu multimodale), ipoteza de relevanță a primelor rezultate devine din ce în ce mai dificil realizabilă ceea ce conduce la performanțe limitate pentru această abordare;
- **"feedback" indirect:** interacția cu utilizatorul se realizează în acest caz în mod indirect, pe baza observării "comportamentului" de căutare a diverși utilizatori în situații diferite. De exemplu, sistemul poate utiliza informații despre datele pe care utilizatori diferiți le-au accesat în urma căutării unor documente cu conținut asemănător (faptul că documentele respective au fost accesate conferă un grad de încredere ridicat privind relevanța conținutului acestora) [Kelly 03]. Aceste informații pot fi stocate cu ușurință de motoarele de căutare actuale și în special de cele "on-line" bazate pe text, ca de exemplu istoricul căutării pe Internet ce implică accesarea de documente web, mesagerie electronică, articole de știri, filme, cărți, programe TV și așa mai departe.

În funcție de durata relativă a procesului de antrenare a sistemului, algoritmi de "relevance feedback" se împart în algoritmi cu *antrenare cu termen scurt de învățare* ("short-term relevance feedback") și respectiv cu *antrenare cu termen lung* ("long-term relevance feedback").

Antrenarea cu termen scurt de învățare presupune interacția cu utilizatorul doar în sesiunea curentă fiind și categoria cea mai studiată de metode în contextul sistemelor de indexare multimedia. Dintre abordările cele mai frecvent folosite putem enumera: algoritmi de schimbare a punctului de interes, algoritmi de determinare a importanței descriptorilor de conținut,

algoritmi statistici sau algoritmi ce implementează procesul de "relevance feedback" ca o problemă de clasificare binară a datelor în două clase de relevanță: date relevante și date ne-relevante pentru utilizator. O parte dintre aceste metode sunt detaliate ulterior în secțiunile următoare.

Principalele provocări ale acestui tip de abordare pot fi sintetizate cu următoarele:

- numărul rezultatelor căutării pentru care utilizatorul furnizează relevanța acestora este de regulă mult mai mic decât dimensiunea descriptorilor de conținut (dimensiunea spațiului de caracteristici) folosiți pentru reprezentarea datelor (vezi Secțiunea 2.1), oferind astfel o capacitate de selecție limitată din punct de vedere statistic;
- realizarea interacției cu utilizatori diferiți va conduce în general la rezultate diferite și uneori chiar contradictorii. Persoane diferite au moduri de percepție diferită cu privire la proprietățile aceluiași concepte, de exemplu un expert va percepe într-un mod diferit conținutului imaginii unei opere de artă față de o persoană neavizată. Acest lucru va conduce la varierea performanțelor sistemului de "relevance feedback" în funcție de utilizator;
- discrepanța dintre numărul de rezultate relevante și cele nerelevante. De cele mai multe ori numărul de rezultate relevante returnate de sistem tinde să fie foarte mic nefiind suficiente pentru ca sistemul să se poată adapta la acestea. Aceeași problemă apare și în situația opusă, când nu există practic nici un rezultat nerelevant, situație în care sistemul nu poate face diferența dintre cele două cazuri;
- rafinarea rezultatelor în timp real este de asemenea un punct critic. Având în vedere interacția directă cu utilizatorul, pentru a fi rentabil, sistemul trebuie să poată furniza noile rezultate cât mai rapid, implicând un timp de așteptare minim din partea utilizatorului.

Învățarea de lungă durată se folosește nu numai de informațiile obținute de la utilizator în sesiunea curentă, ci de toate informațiile furnizate de-a lungul timpului de utilizatori diferiți și în sesiuni diferite. Acestea sunt de regulă stocate de cele mai multe ori sub forma unei reprezentări matriceale a relațiilor descoperite ca existând între informațiile din baza de date, relații ce sunt actualizate pe măsură ce se obțin noi informații de la utilizatori.

Ca și în cazul anterior, există o serie de limitări ale acestui mod de abordare, cele mai semnificative fiind:

- acești algoritmi sunt mai dificil de implementat în cazul bazelor de date ce presupun frecvent eliminarea și adăugarea de date noi;

- gradul de succes depinde foarte mult de cantitatea de informații de "feedback" stocate anterior, de cele mai multe ori în realitate preferându-se utilizarea unei combinații între strategii de învățare de scurtă și lungă durată;
- datorită utilizării mai multor surse de "feedback" informația stocată tinde să fie neomogenă și foarte probabil să nu acopere toate datele;
- ca și în cazul anterior, procesul de rafinare trebuie să poată fi implementat în timp real. Suplimentar complexității datelor de prelucrat, sistemul trebuie să fie capabil să analizeze și un volum semnificativ de date de "feedback" de la utilizatori. De regulă pentru a soluționa această problemă, se preferă împărțirea bazei de date pe diverse niveluri de relevanță folosind ierarhii arborescente de conținut.

6.1 Algoritmii Rocchio

Algoritmii de schimbare a punctului de interogare constituie una dintre primele abordări de tip "relevance feedback" ale problemei rafinării rezultatelor căutării, dezvoltate inițial în contextul căutării de documente textuale, exemplul fiind algoritmul propus de Rocchio [Rocchio 71]. Pornind de la modul de reprezentare al datelor într-un sistem clasic de indexare după conținut în care fiecare document este reprezentat ca un punct în spațiul de caracteristici definit de descriptorii de conținut asociați (vezi și Figura 2.2), o anumită cerere de căutare a utilizatorului ("query") este descrisă la rândul ei în același spațiu sub forma unui punct numit și punct de interogare.

Acest lucru este ilustrat schematic în Figura 6.2. Axele a_1, a_2, \dots, a_n reprezintă valorile atributelor de conținut ce definesc spațiul de caracteristici n -dimensional. Fiecare punct reprezintă valorile descriptorilor unui document din baza de date. Cererea de căutare este reprezentată în acest caz de dreptunghiul verde (punctul de interogare). În urma procesului de căutare, sistemul returnează ca rezultat datele cele mai apropiate punctului de interogare marcate în Figura 6.2 de cercul punctat (punctele care se află la o anumită distanță de "query", de regulă în interiorul unei sfere). Aceste rezultate sunt prezentate utilizatorului de regulă în ordinea descrescătoare a distanței față de punctul de interogare.

Conform algoritmului de "relevance feedback", utilizatorul marchează mai departe rezultatele ca fiind, fie relevante, fie nerelevante pentru datele căutate; de exemplu punctele marcate în figură cu cercuri verzi și respectiv punctele marcate cu "+" de culoare roșie.

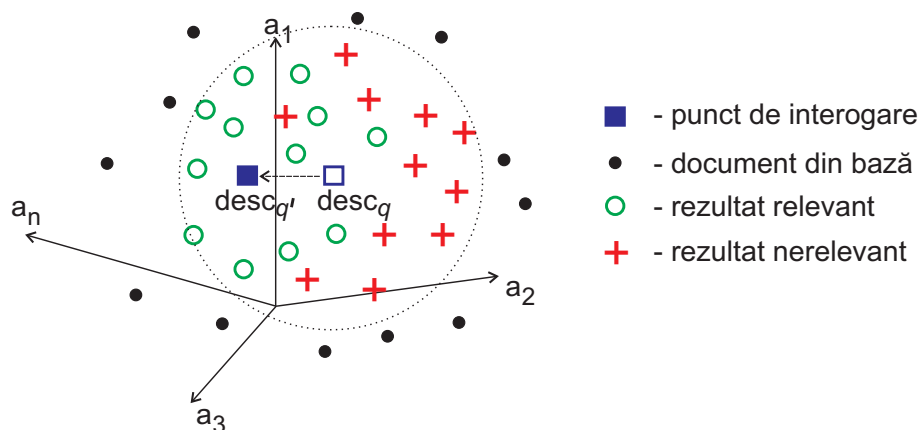


Figura 6.2: Modul de schimbare a punctului de interogare în cazul metodei propuse de Rocchio [Rocchio 71] (reprezentarea obiectelor din bază în spațiul de caracteristici; $desc_q$ reprezintă punctul de interogare inițial iar $desc_{q'}$ noul punct de interogare calculat - notațiile sunt explicate în text).

Algoritmul lui Rocchio utilizează mulțimea de documente relevante, R , și respectiv de documente nerelevante, N , pentru a redefini un nou punct de interogare folosind următoarea relație:

$$desc_{q'} = \alpha \cdot desc_q + \beta \cdot \frac{1}{\|R\|} \sum_{desc_i \in R} desc_i - \gamma \cdot \frac{1}{\|N\|} \sum_{desc_j \in N} desc_j \quad (6.1)$$

unde $desc_{q'}$ reprezintă noul punct de interogare, $desc_q$ reprezintă punctul de interogare inițial, α (ponderea punctului inițial de interogare), β (factorul de importanță al rezultatelor relevante) și γ (factorul de importanță al rezultatelor nerelevante) sunt o serie de ponderi alese empiric (valorile acestora sunt cuprinse în intervalul $[0; 1]$), $\|\cdot\|$ este operatorul ce returnează numărul de elemente ale unei mulțimi iar $desc = \{a_1, \dots, a_n\}$ reprezintă descriptorii de conținut ai rezultatelor căutării.

Definit în acest fel, noul punct de interogare tinde să se deplaseze spre centroidul mulțimii R a rezultatelor marcate ca fiind relevante, ceea ce în urma reluării mecanismului de căutare va conduce la rezultate mai relevante.

6.2 Estimarea importanței atributelor

Algoritmii de estimare a importanței atributelor ("Feature Relevance Estimation") [Rui 99] pleacă de la ipoteza conform căreia pentru o anumită căutare

ponderea atributelor descriptorilor de conținut poate influența relevanța rezultatelor. În mod implicit, atributele descriptorilor sunt considerate a avea o contribuție identică la localizarea datelor celor mai similare, acest lucru fiind realizat pe baza calculului unei măsuri de distanță (de exemplu distanța Euclidiană, vezi și Secțiunea 5). Pe baza interacției cu utilizatorul, ponderile atributelor pot fi modificate astfel încât să se îmbunătățească rezultatele căutării.

Folosind notațiile anterioare, dacă $desc = \{a_1, \dots, a_n\}$ reprezintă descriptorul de conținut al datelor, unde a_i cu $i = 1, \dots, n$ reprezintă valorile atributelor acestuia, atunci se va considera un anumit vector de ponderi, $W = \{w_1, \dots, w_n\}$, unde w_i reprezintă ponderea atributului a_i . Aceste valori sunt inițial considerate egale cu 1 (cu alte cuvinte nu există ponderare).

Sistemul de indexare realizează căutarea datelor pe baza comparării descriptorilor și returnează rezultatele în ordinea descrescătoare a similarității. Ca și în cazul anterior, utilizatorul marchează rezultatele relevante și respectiv nerelevante. Pe baza acestor informații se va modifica ponderea individuală a fiecărui atribut.

O variantă o reprezintă calculul lui w_i în funcție de abaterea pătratică medie a valorilor atributelor σ_i , și anume:

$$w_i = \frac{1}{\sigma_i} \quad (6.2)$$

unde σ_i reprezintă abaterea pătratică medie a valorilor atributului a_i pentru documentele marcate drept relevante de utilizator. Definit în acest fel, un atribut cu grad de importanță ridicat va tinde să aibă o valoare relativ constantă pentru fiecare document în timp ce un atribut mai puțin discriminant pentru datele căutate va tinde să aibă o gamă mult mai mare de valori, ponderea acestuia fiind redusă proporțional.

O altă abordare constă în folosirea de ponderi ce depind de rezultatele căutării individuale după fiecare atribut în parte:

$$w_i = \frac{2 \cdot ||R_i||}{T} \quad (6.3)$$

unde R_i reprezintă mulțimea documentelor relevante în cazul unei căutări folosind drept descriptor doar atributul a_i , $||\cdot||$ este operatorul ce returnează numărul de elemente ale unei mulțimi iar T reprezintă numărul total de documente relevante din bază.

Odată determinate ponderile atributelor, acestea sunt folosite la rafinarea rezultatelor căutării prin calcularea similarității documentelor pe baza unei

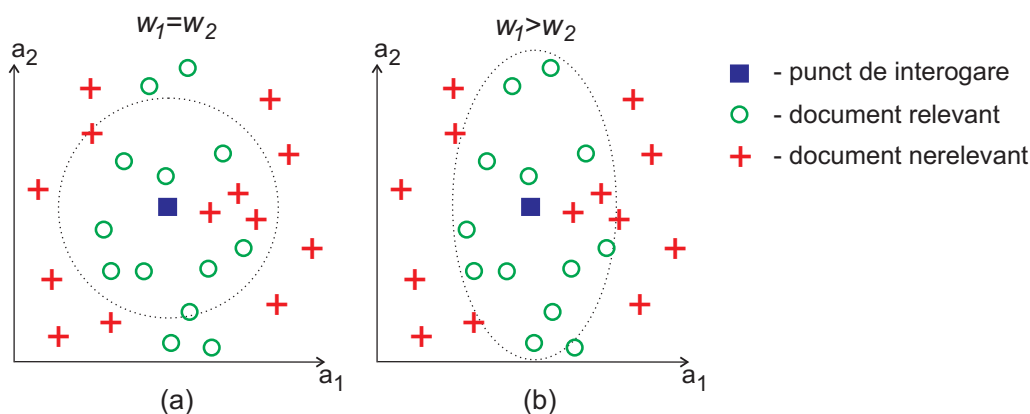


Figura 6.3: Estimarea importanței atributelor cu metoda [Rui 99] (reprezentarea obiectelor din bază în spațiul de caracteristici - pentru exemplificare s-au ales doar două atribute): (a) reprezentarea rezultatelor căutării (delimitate de cercul punctat), (b) modificarea rezultatelor în funcție de noua pondere a atributelor (delimitate de elipsa punctată).

măsuri de distanțe ponderate:

$$d_{FRE}(desc_x, desc_y, W) = \sqrt{\frac{\sum_{i=1}^n w_i \cdot (a_{xi} - a_{yi})^2}{\sum_{i=1}^n w_i}} \quad (6.4)$$

unde $desc_x$ și $desc_y$ reprezintă descriptorii de conținut a două documente iar a_{xi} și a_{yi} cu $i = 1, \dots, n$ atributele acestora.

Modificarea ponderilor asociate fiecărui atribut individual al descriptorului în funcție de rezultatele relevante se traduce în spațiul de caracteristici prin modificarea regiunii de selecție a rezultatelor de la o sferă la un elipsoid, adaptându-se mulțimii de documente relevante. Acest lucru este ilustrat schematic în Figura 6.3.

6.3 Support Vector Machines

Motivați de succesul implementării tehnicilor de învățare asistată de calculator ("machine learning") în contextul sistemelor de indexare după conținut, algoritmi de clasificare și-au găsit aplicabilitate și în cazul tehnicilor de "relevance feedback". Astfel, problema îmbunătățirii performanțelor sistemului de căutare pe baza utilizării informației furnizate de utilizator este transformată într-o problemă clasică de clasificare.

Documentele marcate ca fiind relevante și respectiv nerelevante sunt folosite pentru antrenarea unui anumit clasificator care să permită catalogarea datelor în una dintre cele două clase: documente relevante și respectiv documente nerelevante. Mai departe, documentele din bază sunt trecute prin clasificator și vor fi astfel re-locate uneia dintre cele două clase. Practic, informația de la utilizator este folosită pe post de "ground truth"² pentru determinarea unui set de reguli care să permită partiționarea datelor în cele două clase de relevanță. În urma clasificării, datele vor primi un nou rang calculat în funcție de un grad de relevanță atribuit de clasificator, ceea ce conduce global la rafinarea rezultatelor inițiale.

Dintre tehnicile de clasificare a datelor cel mai frecvent întâlnite în contextul de "relevance feedback" putem menționa: Support Vector Machines (SVM), k-Nearest Neighbors (kNN) sau arborii de decizie (ca de exemplu Random Forests). Pentru mai multe detalii relativ la tehnicile de clasificare a datelor cititorul se poate raporta la [Ionescu 09] [Witten 05] (vezi și explicația din Secțiunea 4.2).

În cele ce urmează ne vom limita la prezentarea unuia dintre algoritmi de clasificare foarte populari care s-a dovedit eficient în rezolvarea diferitelor probleme de indexare a conținutului multimedia și anume Support Vector Machines.

Support Vector Machines (SVM) realizează clasificarea datelor prin construcția unui hiperplan³ ce separă în mod optimal datele de intrare în două categorii [Welling 05]. Aceasta este o problemă de clasificare liniară. Având în vedere că există o multitudine de hiperplane ce pot separa datele, SVM restricționează căutarea la acele hiperplane ce permit o separare maximă între cele două clase (maximizarea "marginii" dintre date).

Cu alte cuvinte, se caută hiperplanul cu proprietatea ca acesta să maximizeze distanța față de cel mai apropiat punct din spațiul de caracteristici. Acesta este denumit "hiperplanul marginii maxime" ("maximum-margin hyperplane"). Un exemplu este prezentat în Figura 6.4.(a) unde spațiul de caracteristici este separat de hiperplanul H_1 ce nu permite separarea datelor, hiperplanul H_2 care are o margine redusă și respectiv hiperplanul H_3 ce maximizează separarea dintre clase (vezi distanțele față de cele mai apropiate puncte).

Formalizarea problemei de clasificare abordată de SVM este următoarea:

²vezi notă de subsol 9.

³un hiperplan este un concept folosit în domeniul algebrei liniare pentru a generaliza noțiunea de linie - folosită în geometria Euclidiană a planului, sau de plan - folosită în geometria Euclidiană tridimensională, pentru cazul n -dimensional, cu $n > 3$.

având la dispoziție un set de date de antrenare, D , constituit ca fiind:

$$D = \{(X_i, c_i) | X_i \in R^n, c_i \in \{-1, 1\}\} \quad (6.5)$$

unde X_i este un vector n -dimensional (de exemplu descriptorul datelor), c_i indică clasa din care face parte vectorul X_i (valori etichete -1 și 1), i reprezintă indicele vectorului curent, cu $i = 1, \dots, p$, iar p reprezintă numărul de vectori considerați; se caută hiperplanul marginii maxime ce permite separarea punctelor din clasa $c_i = 1$ de cele din clasa $c_i = -1$ (vezi Figura 6.4.(b)).

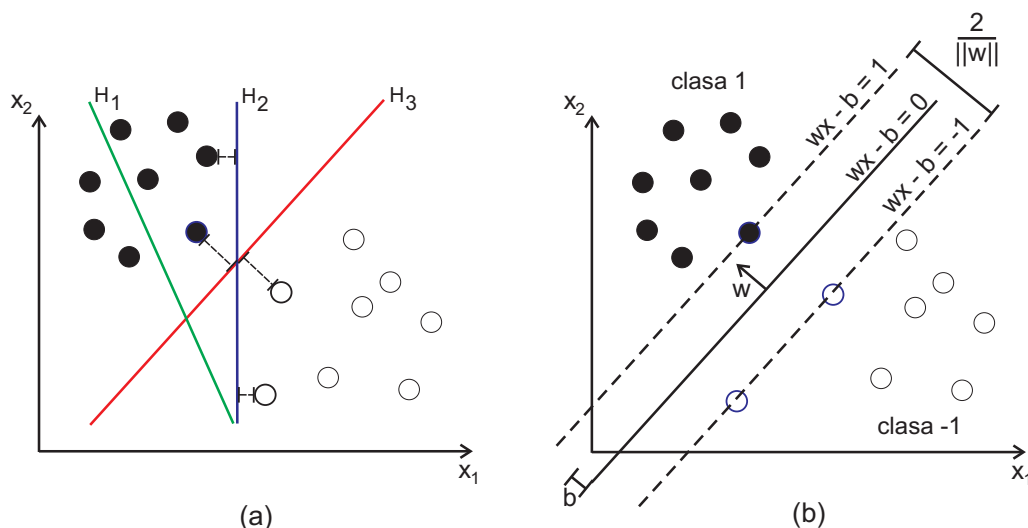


Figura 6.4: Principiul SVM (cercurile reprezintă vectorii de caracteristici, X_1 și X_2 formează spațiul de caracteristici): (a) exemple de hiperplane de separare a datelor, (b) hiperplanul marginii maxime în cazul a două clase (sursă exemple Wikipedia).

Un hiperplan oarecare poate fi definit ca un set de puncte X ce satisfac următoarea relație:

$$W \cdot X - b = 0 \quad (6.6)$$

unde W reprezintă un vector normal (perpendicular pe hiperplan), \cdot reprezintă produsul scalar iar parametrul $\frac{b}{\|W\|}$ va defini decalajul hiperplanului față de originea axei de coordonate, de-a lungul vectorului W (vezi Figura 6.4.(b)).

În scopul definirii marginii maxime, căutăm valorile lui W și b astfel încât acestea să maximizeze distanța dintre hiperplanele paralele, cele mai

depărtate, dar care încă separă datele. Acestea sunt date de ecuațiile:

$$W \cdot X - b = 1 \quad (6.7)$$

$$W \cdot X - b = -1 \quad (6.8)$$

Distanța dintre acestea este $\frac{2}{\|W\|}$, astfel că problema maximizării se transformă într-o problemă de minimizare a valorii $\|W\|$. De asemenea, pentru a preveni ca punctele să se găsească pe margini, se folosesc o serie de constrângeri suplimentare, astfel marginea maximală este determinată de condițiile următoare:

$$W \cdot X_i - b \geq 1, \quad X_i \in c_1 \quad (6.9)$$

$$W \cdot X_i - b \leq -1, \quad X_i \in c_{-1} \quad (6.10)$$

sau

$$c_i \cdot (W \cdot X_i - b) \geq 1 \quad (6.11)$$

pentru oricare $i \in [1; p]$.

Transformată într-o problemă de optimizare, clasificarea SVM poate fi enunțată astfel: alege parametrii W și b astfel încât să minimizeze valoarea $\|W\|$ cu constrângerea ca: $c_i \cdot (W \cdot X_i - b) \geq 1$, pentru oricare i . Această clasificare este valabilă totuși doar în cazul în care datele sunt liniar separabile. În realitate, mai ales în contextul descriptorilor de conținut multimodali, este puțin probabil ca separarea acestora să se poată realiza liniar.

Pentru a crea un clasificator SVM neliniar, la maximizarea marginii dintre clase se folosesc ceea ce numim funcții nucleu sau "kernel functions". Operațiile de înmulțire scalară sunt înlocuite acum de nuclee de funcții neliniare, $k(X, X')$ unde X și X' sunt doi vectori. În acest fel, hiperplanul marginii maximale va fi potrivit datelor într-un spațiu de caracteristici transformat neliniar. Dintre nucleele cel mai frecvent folosite putem menționa:

- nucleu polinomial omogen:

$$k(X, X') = (X \cdot X')^d \quad (6.12)$$

unde d este un număr întreg;

- nucleu polinomial neomogen:

$$k(X, X') = (X \cdot X' + 1)^d \quad (6.13)$$

- funcție radială:

$$k(X, X') = \exp(-\gamma \|X - X'\|^2) \quad (6.14)$$

unde $\gamma > 0$;

- funcție radială Gaussiană:

$$k(X, X') = \exp\left(-\frac{\|X - X'\|^2}{2\sigma^2}\right) \quad (6.15)$$

unde σ^2 reprezintă varianța statistică;

- funcție sigmoidă:

$$k(X, X') = \tanh(\kappa \cdot X \cdot X' + c) \quad (6.16)$$

unde $\tanh(\cdot)$ reprezintă tangenta hiperbolică, $\kappa > 0$ iar $c < 0$.

Cu toate că SVM este un clasificator binar, acesta poate fi folosit cu succes pentru a rezolva probleme de clasificare multi-clasă specifice indexării după conținut. Una dintre metodele cele mai uzuale constă în transformarea clasificării multi-clasă într-o succesiune de clasificări binare [Kotsiantis 07] (de exemplu folosind clasificatori binari ce clasifică o clasă față de toate celelalte - "one-versus-all", sau care clasifică fiecare pereche de clase - "one-versus-one").

În cele ce urmează vom prezenta un studiu comparativ al performanțelor a o serie de algoritmi de "relevance feedback" în contextul unei căutări de imagini folosind descriptori de conținut. Algoritmii vizați sunt: Rocchio, Relevance Feature Estimation, Support Vector Machines (SVM), arbori de decizie (TREE), AdaBoost (BOOST), Random Forests și clasificare ierarhică [Mironică 12b]. Testele sunt efectuate folosind o bază "off-line" și anume Microsoft Object Class Recognition⁴ ce conține imagini cu 23 de categorii de obiecte (de exemplu animale, persoane, avioane, mașini și așa mai departe). Căutarea presupune identificarea tuturor imaginilor ce conțin un anumit obiect furnizat de utilizator.

Rezultatele sunt prezentate în Figura 6.5. Graficele ilustrează performanța căutării pe baza valorii MAP (Mean Average Precision, vezi Secțiune 8; reprezentată pe axa oY - valoarea maximă este 100 ce indică o performanță de 100%) raportată la numărul de sesiuni de "feedback" ale utilizatorului (axa oX ; vezi explicația de la începutul Secțiunii 6). Rezultatele prezentate sunt obținute folosind descriptori de culoare clasici (vezi Secțiune 3.1).

Ceea ce se observă imediat este faptul că performanța căutării crește semnificativ cu numărul sesiunilor de "feedback", de exemplu cu până la 20% în cazul metodei de clasificare ierarhică (comparat cu rezultatele din prima sesiune). De asemenea, comparativ cu rezultatele obținute fără a aplica

⁴vezi notă de subsol 1.

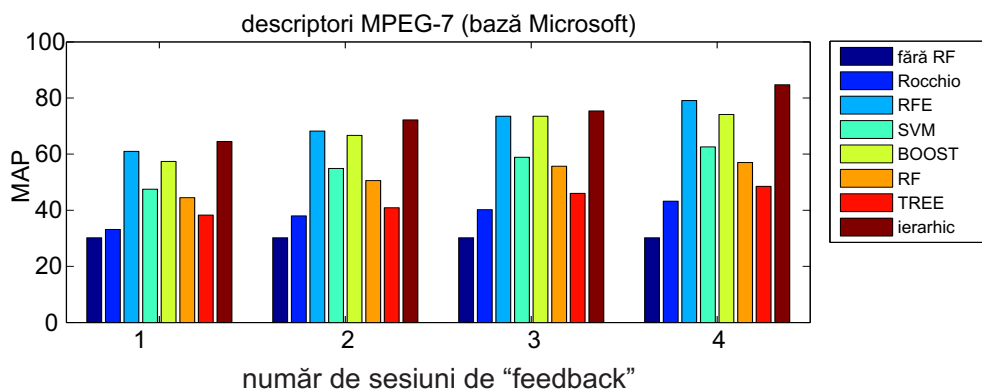


Figura 6.5: Compararea performanțelor tehnicilor de "relevance feedback" în contextul căutării de imagini [Mironică 12b] (notațiile sunt explicate în text).

"relevance feedback", performanța sistemului se poate dubla pentru o singură sesiune sau chiar ajunge la o performanță de peste 80% după mai multe sesiuni, ceea ce este într-adevăr un rezultat relevant. Creșterea semnificativă a performanței se realizează de regulă pentru primele sesiuni de "feedback", în general după prima sesiune, urmând să se diminueze progresiv cu creșterea numărului de sesiuni. De exemplu, clasificarea ierarhică furnizează o creștere a performanței cu 31% în prima sesiune (față de căutarea fără "relevance feedback") și apoi doar de 7% (față de prima sesiune) și de 4% față de a doua sesiune.

Din punct de vedere al metodelor, în ciuda superiorității clare a unor abordări față de altele, nu putem trage o concluzie generală, rezultatele fiind de regulă dependente de baza de test și de sistemul de indexare folosit. În exemplul ilustrat, metoda de clasificare ierarhică urmată de estimarea importanței atributelor (RFE) furnizează performanțele cele mai ridicate.

CAPITOLUL 7

Vizualizarea conținutului multimedia

Vizualizarea conținutului datelor multimedia este la rândul ei o problemă ce trebuie luată în calcul. Aceasta este integrată sistemului de navigare (vezi Secțiunea 2).

În contextul imaginilor, dificultatea vizualizării datelor ține în cea mai mare parte de volumul de date ridicat ce trebuie accesat, conținutul unei imagini putând fi reprezentat simplu prin reprezentarea acestuia la o rezoluție scăzută (de exemplu pe bază de miniaturi). Astfel, o bază de imagini poate fi vizualizată eficient prin vizualizarea miniaturilor imaginilor din aceasta sub formă de planșe. În Figura 7.1 am prezentat ca exemplu modul de vizualizare folosit de platforma de căutare Flickr¹. Se observă faptul că informația furnizată poate fi analizată foarte rapid de utilizator, timpul necesar fiind de ordinul zecilor de secunde.

În contextul secvențelor de imagini, pe lângă volumul mare de date se mai adaugă și problema vizualizării conținutului video dinamic. Este evident faptul că vizualizarea în parte a fiecărei secvențe este aproape imposibilă iar reprezentarea acestora cu o singură imagine este nerelevantă deoarece nu surprinde informația definitorie care ține de conținutul de mișcare și de evoluția în timp. O soluție la această problemă constă în folosirea de rezumate de conținut ce reprezintă practic modalități de reprezentare compactă a conținutului, atât vizual cât și temporal.

Tehnicile de rezumare automată a conținutului video [Truong 07] vizează două categorii de rezumat, și anume *rezumatul în imagini* (static) ce re-

¹<http://www.flickr.com/search/?q=Tour+Eiffel&z=t>

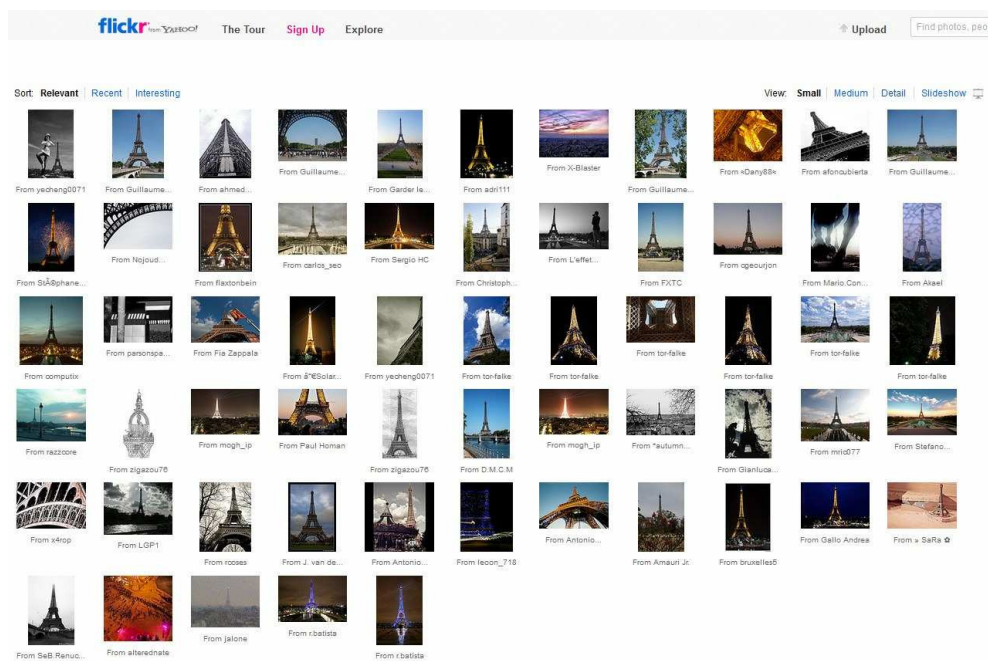


Figura 7.1: Exemplu de vizualizare a conținutului unei colecții de imagini pe platforma Flickr (rezultate obținute în urma căutării de imagini cu ”Tour Eiffel”).

prezintă o colecție de imagini reprezentative și *rezumatul în mișcare* (dinamic), ce reprezintă o colecție de pasaje reprezentative ale secvenței. Rezumatele de conținut permit utilizatorului să-și facă rapid o idee globală asupra conținutului secvenței. Astfel, rezumatul static permite reprezentarea conținutului vizual al secvenței în doar câteva imagini (de exemplu câte o imagine pentru fiecare scenă reprezentativă), ce sunt ușor accesibile utilizatorului prin sistemul de navigare, timpul de vizualizare fiind neglijabil. Pe de altă parte, rezumatul dinamic aduce un plus de informație la nivelul acțiunii prezente în secvență, informație ce nu este disponibilă în rezumatul static. Totuși, fiind el însuși o secvență, în funcție de nivelul de detaliu furnizat, timpul necesar vizualizării acestuia este mai ridicat decât în cazul rezumatului static, dar net inferior timpului de vizualizare integrală a secvenței (un exemplu sunt rezumatele de tip ”trailer” care tind să surprindă doar conținutul de acțiune).

Pe lângă aspectul vizualizării propriu-zise a datelor, așa cum am enunțat și anterior, principala problemă a vizualizării colecțiilor multimedia este dată de necesitatea parcurgerii unui volum semnificativ de date, indiferent dacă

este vorba de imagini sau video. În cele ce urmează vom trece în revistă câteva sisteme de navigare multimedia ce întegrează tehnici inteligente de reprezentare a conținutului datelor:

- MediaTable [Rooij 10] (vezi Figura 7.2): permite categorizarea imaginilor și secvențelor video. Sistemul folosește o vizualizare tabulară ce permite o vedere de ansamblu asupra colecției multimedia și a descrierilor textuale atașate cât și o serie de interfețe grafice ce permit sortarea, filtrarea, selectarea și vizualizarea documentelor.

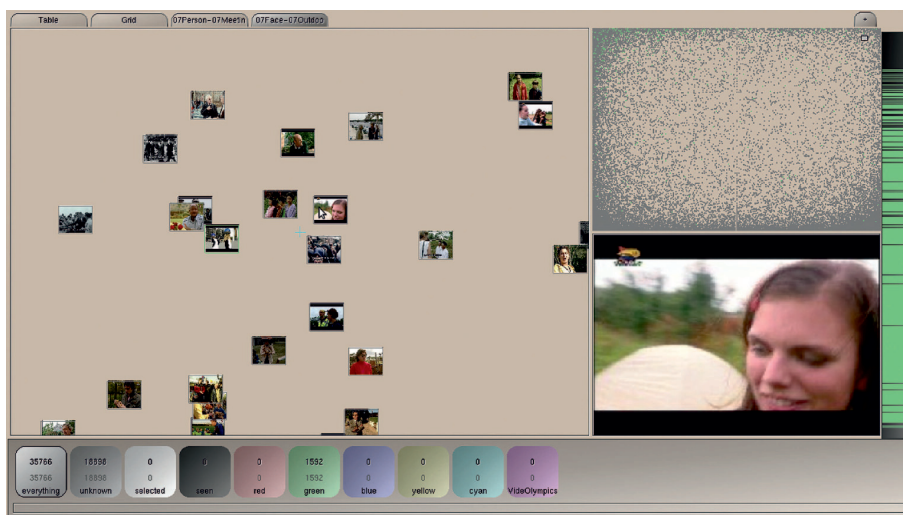


Figura 7.2: Sistemul MediaTable [Rooij 10].



Figura 7.3: Sistemul 3D MARS [Nakazato 01].

În Figura 7.2 este ilustrat un exemplu de vizualizare a datelor după conținut. Graficul din colțul din dreapta sus reprezintă o hartă a distribuției tuturor documentelor din bază în timp ce imaginea din colțul dreapta jos detaliază conținutul documentului selectat curent.

- 3D MARS [Nakazato 01] (vezi Figura 7.3): permite vizualizarea colecțiilor de imagini folosind un sistem de reprezentare 3D de tip realitate virtuală. Imaginile sunt reprezentate în funcție de conținutul de culoare, textură și respectiv structural (un exemplu este prezentat în imaginile din Figura 7.3).
- MediaMill Forkbrowser [Rooij 08] (vezi Figura 7.4): folosește un sistem de vizualizare intercalată atât a rezultatelor căutării video cât și a conținutului temporal. Pe axa de adâncime spre partea superioară sunt reprezentate rezultatele unei anumite căutări, pe axa orizontală este prezentat conținutul temporal al unui segment al secvenței curente ("timeline"), pe axele diagonale sunt ilustrate succesiuni de plane video al căror conținut este similar cu imaginea vizualizată curent în centru ("similarity threads") iar pe axa de adâncime în partea de jos este prezentat istoricul căutărilor ("history").

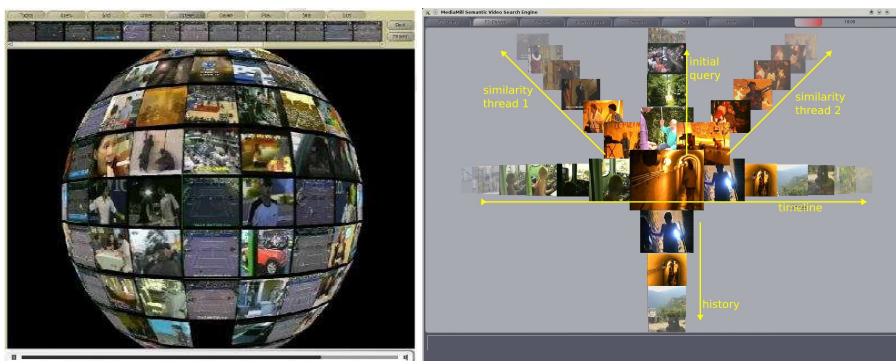


Figura 7.4: Sistemul MediaMill: Forkbrowser [Rooij 08].

- Reprezentare 3D cilindrică [Schoeffmann 11] (vezi Figura 7.5): permite reprezentarea colecțiilor de imagini sub forma unor reprezentări de tip "storyboard" ilustrate folosind o reprezentare cilindrică 3D. Diferite categorii de imagini sunt reprezentate folosind cilindrii diferiți, utilizatorul putând selecta categoria dorită. Pentru vizualizarea curentă, imaginile prezentate în prim plan sunt reprezentate detaliat în timp ce imaginile din fundal sunt reprezentate schematic. Folosind interfața

grafică, utilizatorul poate derula imaginile rulate pe cilindru cât și deta-
lia o anumită regiune a cilindrului.



Figura 7.5: Sistem de reprezentare 3D cilindrică [Schoeffmann 11].

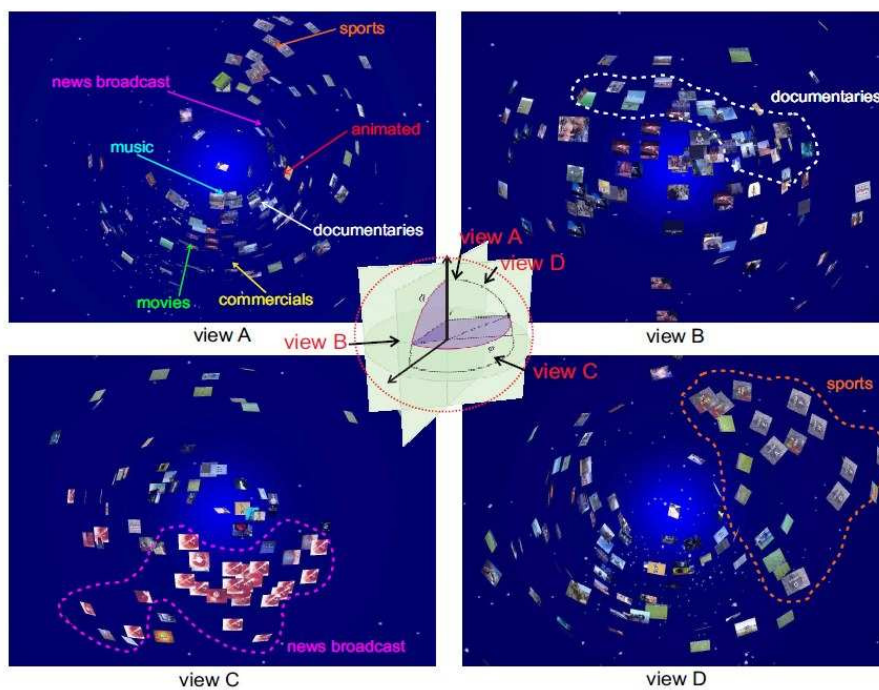


Figura 7.6: Sistemul MovieGlobe [Ionescu 12a].



Figura 7.7: Sistemul nepTunes [Knees 07].

- MovieGlobe [Ionescu 12a] (vezi Figura 7.6²): permite reprezentarea colecțiilor multimedia de imagini și filme într-un spațiu 3D virtual. Fiecare obiect multimedia este reprezentat ca un punct în acest spațiu. Distribuția obiectelor în spațiul 3D este realizată în funcție de similaritatea conținutului acestora. Utilizatorul se poate deplasa virtual și vizualiza conținutul obiectelor întâlnite. În Figura 7.6 este prezentat un exemplu de reprezentare a filmelor în funcție de gen (animație, sport, film, etc.).
- nepTune [Knees 07] (vezi Figura 7.7): permite vizualizarea conținutului colecțiilor de muzică sub forma unor peisaje 3D virtuale pe care utilizatorul le poate explora. Peisajele "muzicale" sunt adaptate automat pe baza analizei conținutului audio preferințelor fiecărui utilizator.

²o demonstrație este disponibilă la http://imag.pub.ro/~bionescu/index_files/MovieGlobe.avi

CAPITOLUL 8

Evaluarea performanțelor indexării

Așa cum am menționat și în Capitolul 2.4, alături de problematica descrierii eficiente a conținutului datelor cât și a conceptului de similaritate între date de regulă heterogene, un aspect cel puțin la fel de important îl constituie evaluarea performanțelor. Cu toate că un sistem de indexare poate funcționa corect din punct de vedere al algoritmilor implementați și al tehnicilor de reprezentare a datelor, acest lucru nu implică și faptul că rezultatele obținute sunt relevante pentru utilizator. Pentru validarea sistemului este necesară evaluarea globală a performanțelor acestuia, atât pentru seturi de date cât mai diverse cât și pentru utilizatori diferiți.

Metodele existente se împart în două categorii: metode de *evaluare subiectivă* ce au la bază utilizatorul și respectiv metode de *evaluare obiectivă* ce se bazează pe calculul unor măsuri matematice. Acestea sunt descrise în cele ce urmează.

8.1 Evaluarea subiectivă

Campanii de evaluare. Evaluarea subiectivă a performanțelor implică însuși utilizatorul. Practic calitatea rezultatelor obținute de sistem este evaluată pe baza opiniei utilizatorilor (care până la urmă este chiar "consumatorul produsului"), ca de exemplu prin realizarea a ceea ce numim "user studies" (sau campanii de evaluare).

Utilizatorului i se pun la dispoziție rezultatele obținute de sistem și acesta va completa un chestionar cu privire la gradul de satisfacție și relevanța

acestora relativ la datele căutate. Procesul se repetă în general pentru un număr cât mai semnificativ de rezultate precum și pentru cât mai mulți utilizatori. De regulă, experimentele respectă un protocol bine definit și sunt realizate în aceleași condiții pentru toți utilizatorii pentru a nu exista factori externi diferiți care să influențeze răspunsurile la întrebări. În final, răspunsurile obținute relativ la performanța sistemului sunt analizate din punct de vedere statistic și se concluzionează asupra performanțelor medii globale ale sistemului.

Prezentăm pentru exemplificare o astfel de campanie de evaluare realizată în cazul tehnicilor de rezumare automată de conținut. Sistemul evaluat este un sistem de generare automată a unui rezumat în imagini a unui document video [Ionescu 10] (o colecție de imagini considerate ca fiind reprezentative pentru conținutul secvenței respective). Având în vedere subiectivitatea unui astfel de proces, se dorește validarea acestuia de către utilizatori. Primul pas al campaniei constă în definirea protocolului de evaluare, și anume acel algoritm pe care îl vor urma utilizatorii. Definirea precisă a unui protocol asigură în primul rând standardizarea testului prin realizarea acestuia în același mod de către toți participanții la evaluare.

În cazul exemplului considerat protocolul folosit este unul simplu și constă în următoarele etape: 1. vizualizarea într-o cameră de proiecție a secvenței video originale (izolarea utilizatorului de alte surse de informație și focalizarea asupra datelor evaluate), 2. prezentarea succesivă a imaginilor rezumatului propus (câte o imagine pe secundă), 3. completarea unui chestionar de către utilizator, 4. repetarea procesului pentru diverse secvențe video. Chestionarul folosit cuprinde următoarele întrebări:

- întrebarea 1 - "În ce măsură estimați că rezumatul propus este relevant pentru conținutul secvenței?". Evaluarea acestei întrebări se realizează pe o scară de valori de la 0 la 10 cu următoarea semnificație: 0 nu știu, 1-2 deloc, 3-4 foarte puțin, 5-6 parțial, 7-8 în mare parte, 9-10 în totalitate. Pentru fiecare grad de apreciere sunt furnizate două niveluri;
- întrebarea 2 - "Cum estimați durata rezumatului din punct de vedere al numărului de imagini furnizate?". Evaluarea pentru această întrebare se realizează tot pe o scară de la 0 la 10 cu următoarea semnificație: 0 nu știu, 1-2 prea scurtă, 3-4 scurtă, 5-6 suficientă, 7-8 ridicată, 9-10 prea lungă.

În Figura 8.1 sunt prezentate rezultatele obținute în urma testării rezumatelor pentru 10 secvențe de animație (sursă [CITIA 13]) de către un număr de 27 de utilizatori. Graficele ilustrează scorul mediu obținut pentru fiecare secvență și întrebare în parte cât și abaterea standard a acestor rezultate (un

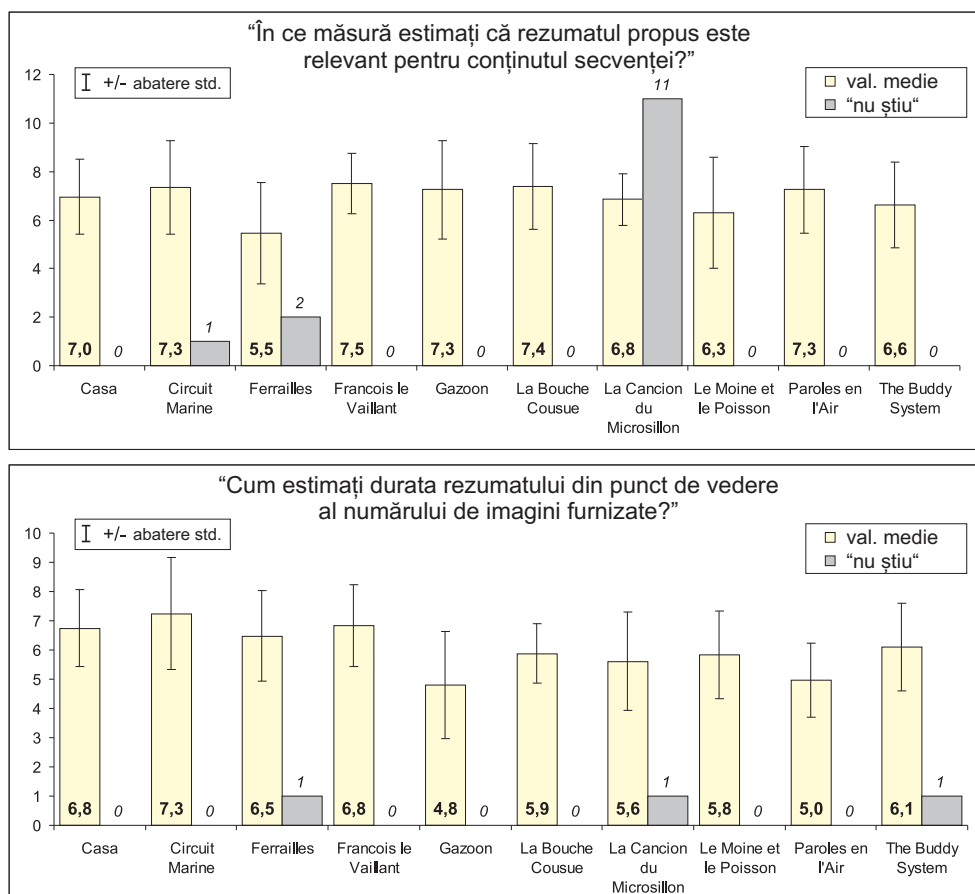


Figura 8.1: Exemplu de rezultate ale unei campanii de evaluare a performanței în cazul metodei propuse în [Ionescu 10] (axa oX corespunde secvențelor testate, axa oY corespunde scorului mediu furnizat de utilizatori, segmentele verticale ilustrează abaterea standard, bărele gri reprezintă numărul de răspunsuri ”nu știu” furnizate de utilizatori).

indicator al gradului de dispersie al răspunsurilor pentru utilizatori diferiți și implicit al subiectivității - cu cât această valoare este mai mare cu atât răspunsurile furnizate de utilizatori au fost mai diferite).

Ceea ce se observă imediat este faptul că rezultatele sunt dependente atât de utilizator cât și de date. De exemplu, există situații în care utilizatorii nu pot furniza un răspuns relevant, de exemplu pentru secvența ”La Cancion du Microsillon” numărul de răspunsuri ”nu știu” este semnificativ (11 din 27); sau dispersia răspunsurilor este foarte ridicată ceea ce atestă un nivel ridicat de subiectivitate, de exemplu pentru secvența ”Le Moine et le Poisson” unde

abaterea standard este de 2.3.

Totuși, pe baza acestor date se poate concluziona la nivel global relativ la calitatea rezultatelor sistemului, în acest exemplu tehnica de rezumare propusă obținând la întrebarea 1 un scor mediu global de 6.9, ceea ce corespunde faptului că este capabilă să reprezinte ”în mare parte” conținutul original; cât și un scor mediu global de 6.1 la întrebarea 2 ceea ce corespunde faptului că durata rezumatului propus tinde să fie adecvată.

Crowd-sourcing. O alternativă actuală la realizarea fizică de campanii de evaluare o constituie folosirea mediului on-line și anume a Internetului. Una dintre dificultățile principale ale unei campanii de evaluare o constituie dificultatea de a dispune de un număr semnificativ de utilizatori la un anumit moment de timp într-o anumită locație. Astfel că o soluție mai eficientă o constituie organizarea campaniei on-line, utilizatorii nefiind restricționați a fi prezenți fizic și putând realiza evaluarea la momentul dorit în funcție de disponibilitatea lor de timp. Mai mult, participarea on-line permite accesarea unui număr semnificativ de utilizatori din toată lumea.

Un domeniu aparte își găsește în prezent aplicație în contextul sistemelor de evaluare a performanțelor algoritmilor multimedia și anume acela de ”crowd-sourcing”. Cu toate că dezvoltarea ”crowd-sourcing” nu este legată de acest context, fiind dezvoltată în principal pentru realizarea unei structuri de prestare de servicii la distanță - conceptul de ”crowd-sourcing” fiind definit ca ”procesul de formulare a unei anumite sarcini de lucru, divizarea acesteia în micro-sarcini ce pot fi realizate foarte ușor și rapid de personal necalificat și distribuirea acestora spre rezolvare către un grup necunoscut de utilizatori de pe Internet” - posibilitatea de a accesa un număr practic nelimitat de utilizatori face din aceasta un candidat ideal pentru evaluarea subiectivă.

În prezent domeniul de ”crowd-sourcing” se stabilește ca domeniu de sine stătător asociat metodelor de analiză multimedia. Tot mai multe studii dovedesc faptul că rezultatele obținute în urma ”crowd-sourcing” pot fi comparabile cu cele obținute de utilizatori experți [Nowak 10]. Totuși sistemul de ”crowd-sourcing” nu este perfect și nu orice evaluare poate fi proiectată prin intermediul ”crowd-sourcing”.

Principala problemă este dată de controlul calității rezultatelor. Dacă în cazul campaniilor de evaluare utilizatorii sunt aleși astfel încât să fie familiarizați cu domeniul precum și să fie motivați în a furniza o evaluare de calitate (voluntar, în interes de cercetare, eventual remunerat), în cazul ”crowd-sourcing” nu există un control direct asupra alegerii utilizatorilor iar calitatea rezultatelor nu poate fi controlată în mod direct, participanții la studiu fiind motivați în principal de un câștig financiar asociat fiecărei sarcini de

rezolvat care este extrem de redus (exemplu 4\$ pe oră). Din perspectiva organizării evaluării, singurul mecanism de creștere a calității evaluării este dat de modul de concepere al evaluării care trebuie să fie unul intuitiv, simplu, rapid și atractiv pentru utilizator.

Dintre platformele de "crowd-sourcing" existente una dintre cele mai populare este Amazon Mechanical Turk¹. Aceasta este totuși limitată în a fi accesibilă doar pentru persoane ("requesters" - persoanele care formulează sarcinile ce trebuiesc rezolvate de utilizatori) care au coordonate bancare în Statele Unite. O alternativă la aceasta este platformă este Crowdfunder². Cererile de lucru create în Crowdfunder pot fi publicate pe diverse canale de "crowd-sourcing" ce includ și platforma Amazon Mechanical Turk.

În ceea ce privește controlul calității, există o serie de facilități care țin mai mult de modul de alegere al utilizatorilor decât de evaluarea acestora. De exemplu, în cazul platformei Amazon Mechanical Turk se poate opta pentru a alege utilizatori din anumite locații geografice, alege utilizatori în funcție de performanța acestora dovedită în alte sarcini efectuate anterior (cel mai probabil în domenii complet diferite) sau pe baza numărului de sarcini realizate anterior. Există și posibilitatea de refuzare a rezultatelor considerate nesatisfăcătoare fără a implica costuri suplimentare. În cazul platformei Crowdfunder aceasta introduce conceptul de "gold units" prin care încearcă să elimine utilizatorii cu performanțe slabe precum și posibilitatea de generare de răspunsuri automate sau aleatorii. Practic, utilizatorilor li se cere să răspundă la cel puțin 4 întrebări al căror răspuns este deja cunoscut de sistem și doar în cazul în care obțin o precizie de minim 70% răspunsurile acestora la sarcina curentă de rezolvat sunt luate în calcul. În cazul platformei Crowdfunder nu există posibilitatea de a refuza răspunsurile considerate ca fiind nerelevante.

Indiferent de modul de implicare al utilizatorilor în procesul de evaluare, acest mod de abordare presupune un anumit grad de subiectivitate. Persoane diferite pot percepe diferit anumite informații (vezi și exemplul din Figura 8.1). Astfel, se pune problema găsirii unei modalități de evaluare a gradului de subiectivitate dintre evaluările furnizate de utilizatori, informație ce este de regulă furnizată împreună cu rezultatele obținute.

Una dintre abordările cele mai frecvent folosite constă în evaluarea gradului de concordanță dintre evaluările realizate de utilizatori diferiți pentru aceleași date, ceea ce se numește "inter-annotator agreement". Prezentăm în continuare modul de calcul al coeficientului Kappa [Carletta 96] ce reprezintă

¹<https://www.mturk.com/mturk>

²<http://crowdfunder.com>

o măsură statistică a concordanței dintre răspunsurile furnizate de utilizatori diferiți. Spre deosebire de alte mărimi similare, coeficientul Kappa ia în calcul și concordanța rezultatelor obținută din întâmplare (aleator).

Să considerăm cazul a doi utilizatori care evaluează un număr de N entități ca aparținând a C categorii (categoriile considerate sunt complementare). De exemplu poate fi vorba de etichetarea a N imagini ca fiind relevante sau nerelevante ($C = 2$ în acest caz). În acest caz coeficientul Kappa este dat de relația următoare:

$$\kappa = \frac{Pr(a) - Pr(e)}{1 - Pr(e)} \quad (8.1)$$

unde $Pr(a)$ reprezintă probabilitatea observată relativă de concordanță între utilizatori iar $Pr(e)$ reprezintă probabilitatea ipotetică de concordanță datorată întâmplării. Dacă răspunsurile utilizatorilor sunt în concordanță completă atunci valoarea lui κ este 1 iar similar, dacă există o disconcordanță totală între răspunsuri κ este 0. În realitate o valoare a lui κ superioară a 0.6 este considerată ca fiind perfectă.

Pentru exemplificare să considerăm următoarele date (sursă Wikipedia): avem la dispoziție 50 de propuneri de proiecte de cercetare ce sunt evaluate fiecare de câte doi evaluatori (notați A și respectiv B). Aceștia atribuie propunerilor categoria "da" sau "nu" (semnificând acceptarea acestora pentru finanțare sau nu). Presupunând că datele obținute sunt cele prezentate în Tabelul 8.1 (numerele corespund numărului de proiecte pentru care evaluatorii au furnizat răspunsul da sau nu) atunci probabilitățile $Pr(a)$ și $Pr(e)$ sunt estimate în felul următor:

- $Pr(a)$: evaluatorii A și B au acordat împreună calificativul "da" pentru 20 de proiecte și respectiv "nu" pentru 15 proiecte astfel că probabilitatea de concordanță a răspunsurilor este $Pr(a) = (20 + 15)/50 = 0.7$;

Tabelul 8.1: Exemplu de calcul al coeficientului Kappa (sursă Wikipedia).

		B	B
		"da"	"nu"
A	"da"	20	5
A	"nu"	10	15

- $Pr(e)$: în acest caz se observă următoarele: evaluatorul A a răspuns "da" pentru 25 de proiecte și "nu" tot pentru 25 ceea ce înseamnă că evaluatorul A răspunde cu "da" pentru 50% din cazuri. Similar, evaluatorul B a răspuns "da" pentru 30 de proiecte și "nu" pentru 20

ceea ce înseamnă că evaluatorul B răspunde cu "da" pentru 60% din cazuri.

Probabilitatea ca cei doi evaluatori să răspundă cu "da" în mod aleator este $0.5 \cdot 0.6 = 0.3$ iar probabilitatea ca ambii să răspundă cu "nu" este $0.5 \cdot 0.4 = 0.2$. Astfel, per total probabilitatea de concordanță aleatoare este $0.3 + 0.2 = 0.5$.

Aplicând relația anterioară obținem în acest caz un coeficient Kappa de 0.4 care indică o concordanță relativ scăzută a rezultatelor.

8.2 Evaluarea obiectivă

O altă abordare a problemei evaluării performanței sistemelor de indexare după conținut o constituie metodele de evaluare așa zisă obiectivă. Acestea se bazează pe evaluarea performanțelor cuantificând erorile de căutare cu diverse măsuri statistice matematice. Pentru a putea evalua o măsură de eroare este necesară cunoașterea apartenenței datelor la clasele căutate (datele să fie etichetate) sau cu alte cuvinte "ground truth".

Având în vedere faptul că este practic imposibil să dispunem de "ground truth" în cazul unei baze de date dinamice (de exemplu de pe Internet) sau chiar de dimensiune semnificativă, lucru ce ar face procesul de căutare inutil atâta timp cât datele sunt deja cunoscute, validarea obiectivă se realizează preliminar folosind seturi de date de test. Sistemul se calibrează astfel pentru performanță optimă folosind aceste baze de test urmând a fi implementat practic ulterior în contextul real. Pentru ca rezultatele unui astfel de proces de evaluare să fie relevante la scară reală, seturile de date folosite trebuie să fie reprezentative și cât mai diverse.

Ca ordin de măsură, în contextul actual, bazele de test pentru sistemele de căutare după conținut a imaginilor tind să conțină până la milioane de imagini în timp ce în contextul video acestea sunt de ordinul sutelor de mii. Principala limitare este dată de efortul necesar etichetării acestora ce presupune analiza lor manuală de către experți umani. De exemplu, dacă dorim validarea unui sistem de căutare a secvențelor de gol într-o bază video de înregistrări de fotbal, fiecare dintre secvențe trebuie parcursă manual și etichetate momentele de timp în care apar secvențele căutate. Pe baza acestor date, rezultatele obținute de sistemul de căutare automată pot fi comparate cu rezultatele ideale obținute manual.

În literatura de specialitate există o multitudine de abordări propuse pentru evaluarea obiectivă a performanțelor, pentru o descriere exhaustivă a

acestora cititorul se poate raporta la [Manning 08]. În cele ce urmează vom detalia unele dintre abordările cele mai frecvent întâlnite.

8.2.1 Precision-Recall

Dacă analizăm problema căutării datelor din perspectiva unui sistem de clasificare (vezi exemplu Secțiune 4.2) și anume, rezultatele obținute în urma căutării corespund de fapt unei clasificări binare a datelor existente, acestea fiind etichetate fie ca aparținând clasei obiectului căutat ("query", clasa A), fie ca aparținând celorlalte clase existente (clasa B), atunci erorile de căutare pot fi sintetizate în modul următor (vezi Tabel 8.2):

- tp sau "true positive": reprezintă obținerea unui rezultat corect și anume obiectul returnat de sistem a fost prezis ca aparținând clasei A (clasa căutată) acesta corespunzând și în realitate clasei A ;
- fp sau "false positive": reprezintă obținerea unui rezultat fals și anume obiectul returnat de sistem a fost prezis ca aparținând clasei A dar în realitate acesta corespunde unui obiect din clasa B ceea ce conduce la o predicție falsă;

Tabelul 8.2: Erori statistice în cazul clasificării datelor.

		clasa reală	
		clasa A	clasa B
clasa prezisă	clasa A	tp (true positive)	fp (false positive)
	clasa B	fn (false negative)	tn (true negative)

- fn sau "false negative": reprezintă obținerea tot a unui rezultat fals și anume sistemul a prezis că obiectul returnat aparține clasei B în realitate acesta fiind din clasa A fapt ce conduce la o non-detectie, obiectul A (din clasa căutată) fiind pierdut;
- tn sau "true negative": reprezintă prezicerea rezultatului ca fiind un obiect din clasa B în măsura în care acesta este în realitate tot din clasa B această situație fiind o confirmare a absenței obiectului căutat de tip A .

Cu alte cuvinte, în urma căutării se pot obține două situații de eroare: obiectul căutat este estimat eronat ca fiind un obiect din altă clasă, eroare cuantizată de raportul fp ; și respectiv obiectul căutat nu este găsit, situație cuantizată de raportul fn .

Pe baza acestor erori sunt definite măsurile de performanță numite *precision* și *recall* astfel:

$$precision = \frac{tp}{tp + fp}, \quad recall = \frac{tp}{tp + fn} \quad (8.2)$$

Definite în acest fel, *precision* este o măsură a falselor detecții iar *recall* o măsură a non-detețiilor. Plaja de valori a acestora se găsește în intervalul $[0; 1]$ unde 1 reprezintă cazul ideal în care nu există nici o falsă detecție ($fp = 0$) și respectiv toate documentele existente în bază au fost găsite ($fn = 0$). Se poate observa faptul că valoarea $tp + fn$ este o constantă și reprezintă numărul total de obiecte de tip A existente în baza de date (numărul celor identificate corect + numărul celor care nu au fost returnate).

Dacă analizăm problema căutării datelor din perspectiva unui sistem de indexare clasic în care rezultatele căutării sunt reprezentate în ordinea descrescătoare a relevanței acestora relativ la obiectul căutat (vezi exemplu Secțiune 6) atunci modul de calcul al *precision* și *recall* este un pic diferit. Diferența provine din faptul că evaluarea performanței se realizează de această dată pe un set de rezultate ordonate și care nu reprezintă neapărat toate documentele disponibile din baza de date (se pot returna doar o parte din acestea în urma căutării - de exemplu în cazul bazelor de date de pe Internet rezultatele căutării sunt limitate la un număr ce poate fi gestionat de utilizator).

În acest context, *precision* este o măsură a procentului din documentele returnate ce sunt relevante pentru obiectul căutat ("query"):

$$precision = \frac{|\{\text{documente relevante}\} \cap \{\text{documente returnate}\}|}{|\{\text{documente returnate}\}|} \quad (8.3)$$

unde operatorul $|\cdot|$ returnează numărul de elemente ale unei mulțimi.

Similar, *recall* este o măsură a procentului de documentele relevante pentru obiectul căutat ce au fost returnate în urma căutării și anume:

$$recall = \frac{|\{\text{documente relevante}\} \cap \{\text{documente returnate}\}|}{|\{\text{documente relevante}\}|} \quad (8.4)$$

Dat fiind faptul că aceste măsuri sunt evaluate pentru o anumită căutare particulară, pentru a obține o măsură globală de performanță de regulă se calculează valorile medii ale acestora pentru un anumit număr de căutări. Dacă baza de date este cunoscută, atunci se poate realiza o evaluare exhaustivă în care fiecare document din bază este folosit pentru a specifica cererea de căutare iar performanța sistemului este estimată ca valoare medie pentru toate căutările efectuate.

8.2.2 F-measure

Având în vedere cele două situații de eroare ce trebuie luate în calcul pentru evaluarea performanțelor indexării și anume numărul de false detecții și respectiv numărul de non-detcții, se pune problema care dintre acestea este mai importantă. Astfel, de exemplu un sistem de indexare care furnizează *precision* de 95% și *recall* de 80% este preferabil unui sistem ce furnizează *recall* de 95% și respectiv *precision* de 80%? Cu alte cuvinte, care dintre cele două situații sunt mai dezavantajoase, un sistem în care rata de documente relevante returnate este mai mare (număr de false detecții redus) iar numărul total de documente relevante returnate din numărul total existent în bază este mai mic (numărul de non-detcții mai mare), sau situația inversă?

În realitate, răspunsul depinde strict de domeniul de aplicație. Figura 8.2 prezintă estimativ importanța celor două măsuri pentru o serie de domenii de aplicație [Worring 12]. Astfel, dacă considerăm ca domeniu de aplicație căutarea datelor pe Internet atunci cel mai important parametru este *precision* deoarece se dorește ca rezultatele căutării să fie cât mai precise. În același timp nu este la fel de important faptul că în urma căutării nu obținem toate rezultatele relevante existente, în cazul Internetului acesta fiind un număr practic nelimitat, ci este suficientă obținerea a unei submulțimi a acestora. Este cunoscut faptul că în practică în urma căutării într-un sistem "on-line" ne limităm de regulă în a analiza doar primele câteva zeci de rezultate.

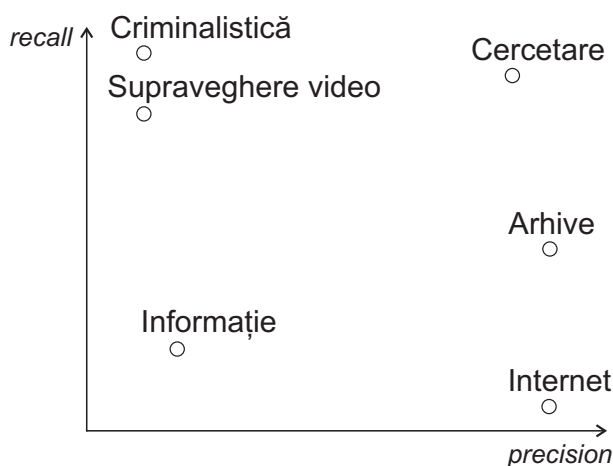


Figura 8.2: Gradul de importanță al *precision* și *recall* în funcție de domeniul de aplicație (bazat pe informațiile prezentate în [Worring 12]).

Pe de altă parte, dacă considerăm ca domeniu de aplicație un sistem spe-

cific expertizei criminalistice, de exemplu un sistem de identificare a amprentelor, în acest caz este mai important parametrul de *recall*. Cu alte cuvinte, este mai important ca sistemul să fie capabil să returneze toate documentele relevante existente în baza de date chiar dacă numărul de detecții false este ridicat. Acestea pot fi reduse ulterior printr-o analiză manuală a rezultatelor dar absența unor documente relevante pentru căutare din rezultate nu mai poate fi corectată.

În acest context, în literatura de specialitate există un parametru care combină contribuția celor două măsuri și anume *F – measure*. Acesta este definit astfel:

$$F - measure = (1 + \beta^2) \cdot \frac{precision \cdot recall}{\beta^2 \cdot precision + recall} \quad (8.5)$$

unde β reprezintă un parametru de reglaj al contribuției celor două măsuri. În funcție de valoarea lui β , *F – measure* poate evidenția mai mult contribuția uneia dintre cele două măsuri permițând adaptarea evaluării la domeniul de aplicație.

Dacă $\beta = 1$ atunci *precision* și *recall* au ponderi egale ceea ce conduce la mărimea *F1 – score* definită ca fiind media armonică dintre *precision* și *recall*, astfel:

$$F1 - score = 2 \cdot \frac{precision \cdot recall}{precision + recall} \quad (8.6)$$

8.2.3 Curbă de precision-recall și ROC

Având în vedere faptul că de regulă sistemele de indexare returnează rezultatele în ordinea descrescătoare a relevanței față de cererea de căutare ("ranking"), valorile estimate pentru *precision* și *recall* sunt dependente de dimensiunea ferestrei de analiză a rezultatelor returnate. De exemplu, nu este același lucru dacă evaluăm performanța pentru 100 de rezultate returnate sau pentru 200, în cazul din urmă fiind mai probabil ca numărul de rezultate corecte să fie mai mare.

Se pune astfel problema evaluării performanței pentru puncte de operare ("operating points") diferite. Una dintre modalitățile cele mai frecvent folosite este aceea de a reprezenta grafic *precision* în funcție de *recall* pentru toată plaja de dimensiuni a ferestrei de rezultate până în punctul în care în aceasta se regăsesc toate datele căutate existente în baza de date.

Algoritmul de generare este următorul: pentru o anumită căutare în baza de date se consideră doar primele N_i rezultate obținute, valoare astfel aleasă încât între acestea să se găsească exact i rezultate corecte. Valoarea lui i va varia de la 1 la $tp + fn$ (vezi ecuația 8.2), și anume până în momentul în care

regăsim în fereastra considerată toate rezultatele corecte existente în baza de date.

În aceste condiții, *precision* și *recall* se evaluează în felul următor:

$$\begin{aligned} precision &= \frac{1}{N_1}, \frac{2}{N_2}, \dots, \frac{i}{N_i}, \dots, \frac{tp + fn}{N_{tp+fn}} \\ recall &= \frac{1}{tp + fn}, \frac{2}{tp + fn}, \dots, \frac{i}{tp + fn}, \dots, 1 \end{aligned} \quad (8.7)$$

unde N_{tp+fn} reprezintă acea dimensiune a ferestrei pentru care obținem toate rezultatele corecte existente în bază. Reprezentat în acest fel, graficul *precision – recall* oferă o imagine asupra performanței sistemului pentru toată plaja de puncte de operare, punctându-se stabili performanța punctuală în oricare dintre acestea.

Figura 8.3 prezintă câteva exemple de grafice *precision – recall* pentru un sistem perfect în care *precision* și *recall* sunt 100%, un sistem complet ineficient în care *precision* și *recall* sunt 0% și un sistem real, sistemul propus în [Ionescu 13]. Primele două variante sunt variantele extreme, de performanță maximă și relativ minimă, în realitate performanțele sistemelor existente regăsindu-se între aceste două curbe (vezi Figura 8.3.(c)).

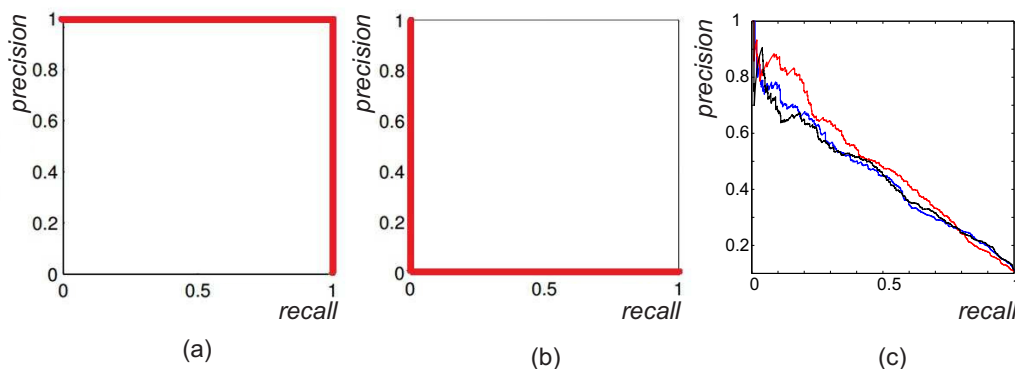


Figura 8.3: Exemple de grafice de tip *precision – recall* pentru: (a) un sistem perfect, (b) un sistem complet ineficient, (c) un sistem real de căutare automată a segmentelor de violență din filme [Ionescu 13] (curbele sunt obținute pentru diferite valori ale parametrilor sistemului).

O altă interpretare a graficului *precision – recall* este aceea din perspectiva raportului de documente găsite corect (*tpr*) raportat la raportul de documente returnate eronat (*fpr*), ceea ce se numește curbă de tip Receiver Operational Characteristic sau ROC. Cele două rapoarte sunt definite în

modul următor:

$$tpr = \frac{tp}{tp + fn}, \quad fpr = \frac{fp}{fp + tn} \quad (8.8)$$

unde tp reprezintă numărul de documente returnate corect (vezi ecuația 8.2), fn reprezintă numărul de documente căutate care nu sunt returnate (non-dectție), fp reprezintă numărul de documente fals detectate iar tn reprezintă numărul de documente ignorate (documente care sunt prezise corect ca neaparținând clasei căutate). Definite în acest fel, tpr este o măsură a numărului de documente returnate corect iar fpr o măsură a numărului de documente returnate eronat.

Figura 8.4 prezintă două exemple de curbe ROC, în cazul unui sistem perfect în care tpr este 100% iar fpr este 0% cât și în cazul unui sistem complet ineficient în care numărul de rezultate corecte este egal cu numărul de rezultate false, un astfel de sistem neputând fi practic utilizat. În realitate, pentru ca un sistem de indexare să ofere performanțe bune, curba ROC asociată trebuie să se situeze între cele două grafice, cât mai apropiată de sistemul ideal.

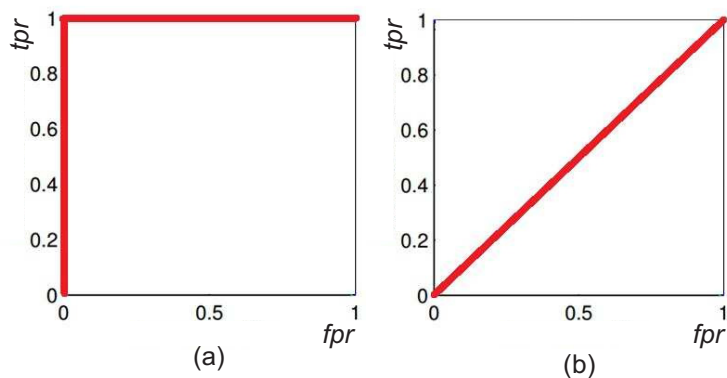


Figura 8.4: Exemple de grafice de tip ROC pentru: (a) un sistem perfect, (b) un sistem complet ineficient în care numărul de documente returnate eronat este egal cu numărul de documente returnate corect.

8.2.4 Mean Average Precision

În ultimii ani, pornind din contextul sistemelor de indexare video, s-a impus ca standard de evaluare a performanțelor sistemelor de indexare ceea ce numim Mean Average Precision sau MAP³. MAP furnizează o măsură a ca-

³vezi utilitar http://trec.nist.gov/trec_eval

lității sistemului pentru diferite valori ale *recall* (vezi ecuație 8.2), totul prin intermediul unei singure mărimi. Acesta se dovedește în practică a furniza o bună stabilitate și discriminanță în evaluarea diferitelor sisteme.

MAP este estimat în modul următor: dacă pentru o anumită cerere de căutare q_j , unde $j = 1, \dots, |Q|$ cu Q reprezentând mulțimea căutărilor posibile pentru sistemul considerat (de exemplu, dacă sistemul permite indexarea în funcție de obiecte atunci Q reprezintă mulțimea tuturor obiectelor din bază) iar operatorul $||\cdot||$ returnează numărul de elemente ale unei mulțimi; mulțimea documentelor relevante din bază este $\{d_1, \dots, d_{m_j}\}$ (numărul de documente relevante pentru q_j este m_j) iar R_{jk} reprezintă mulțimea primelor documente returnate până la documentul d_k (fereastra de rezultate care include și documentul d_k), atunci MAP este definit ca:

$$MAP(Q) = \frac{1}{||Q||} \cdot \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} precision(R_{jk}) \quad (8.9)$$

unde *precision* este calculat așa cum a fost definit în ecuația 8.2. Cu alte cuvinte, MAP reprezintă media *precision* pentru fereastra de rezultate ce include toate documentele relevante pentru o căutare (termenul *Average*), valoare ce este la rândul ei mediată pentru toate căutărilor posibile (termenul *Mean*). În cazul în care sistemul nu returnează nici un document relevant atunci MAP este 0%.

Pentru o singură cerere de căutare ("query") MAP poate fi aproximat ca fiind aria dintre graficului *precision-recall* și axa orizontală (vezi Figura 8.3) și astfel pentru un set de căutări acesta va reprezenta aria medie a graficelor de *precision-recall*.

CAPITOLUL 9

Paradigme ale indexării

În capitolele anterioare am trecut în revistă punctual marea parte a problemelor de prelucrare aferente sistemelor de indexare automată după conținut a datelor multimedia.

În acest ultim capitol vom face o trecere în revistă a barierelor tehnologice, principale, ce trebuie depășite pentru a putea soluționa eficient problema căutării informației. Acestea sunt enunțate în literatură sub denumirea de paradigme:

- **paradigma senzorială** ("sensor gap") reprezintă discrepanța care există între informațiile prezente în lumea reală 3D și informațiile înregistrate de senzori (de exemplu camere foto, video, microfoane, etc.), informații ce sunt folosite pentru analiza conținutului datelor. De exemplu, în cazul imaginilor acestea nu sunt decât proiecții plane 2D al lumii 3D. Mai mult, același obiect de interes poate conduce la un număr nelimitat de reprezentări diferite datorate perturbației senzorilor sau a factorilor externi (vezi exemple din Figura 9.1). Astfel, o primă paradigmă ce trebuie depășită este aceea a modelării informației incomplete de care dispunem și a variabilității acesteia. Practic metodele de analiză de conținut încearcă să estimeze informațiile lipsă, fie pe baza unor modele, sau prin compensarea cu informații suplimentare obținute din alte surse;
- **paradigma semantică** ("semantic gap") reprezintă discrepanța care există între informațiile extrase în mod automat din date și semnificația semantică pe care le-o putem atribui acestora. Cu alte cuvinte, în ciuda



Figura 9.1: Un anumit obiect poate fi înregistrat sub o multitudine de reprezentări diferite datorate schimbării unghiului din care este reprezentat, schimbării de iluminare, schimbării fundalului sau ocluziei cu alte obiecte (sursă imagini [Snoek 10]).

faptului că un sistem poate funcționa corect din punct de vedere al algoritmilor, și chiar mai mult, poate fi antrenat să răspundă optimal pentru un anumit domeniu de aplicație sau set de date, în realitate rezultatele obținute pot să nu corespundă așteptărilor și a modului de percepție uman;

- **paradigma modelării** ("model gap") reprezintă imposibilitatea de a determina un model general pentru toate obiectele sau entitățile informaționale existente în lume fiind limitați în a modela cazuri particulare, precum obiecte, concepte, evenimente și așa mai departe. Diversitatea informațională existentă face imposibilă acoperirea tuturor cazurilor posibile;



Figura 9.2: Există o multitudine de obiecte și concepte ce trebuie modelate pentru a putea fi accesate la nivel de informație.

- **paradigma intenției** ("intention/query gap") reprezintă discrepanța dintre informațiile pe care utilizatorul dorește să le găsească și modul de exprimare a criteriilor de căutare într-un sistem de indexare (vezi Figura 9.3). Cele mai performante metode existente permit specificarea criteriilor de căutare sub formă textuală. Acest mod de reprezentare este limitat la un număr redus de informații ce pot fi furnizate (de regulă cel mult o propoziție) nerefectând în totalitate informația reală dorită;



Figura 9.3: Există o multitudine de "întrebuințări" ale aceluiași concept, de exemplu "kiwi" poate reprezenta atât o companie aeriană, un fruct sau o pasăre, "bear" (urs) este foarte similar cu "beer" (bere) sau "grid" (caroiaj) cu "greed" (lacom) (exemplu din cursul Indexarea Conținutului Vizual, Constantin Vertan, Universitatea Politehnica din București).

- **paradigma utilității** ("utility gap") reprezintă discrepanța care există între rezultatele furnizate de sistem și utilitatea reală practică a acestora pentru utilizator. Ca și în cazul paradigmei semantice, sistemul poate fi performant și să returneze utilizatorului o multitudine de informații relevante relativ la datele căutate, dar câte dintre aceste informații vor servi în mod real util utilizatorului.

Bibliografie

- [Bimbo 99] A. Del Bimbo. *Visual Information Retrieval*. Morgan Kaufmann Publishers, San Francisco, USA 1999.
- [Bovik 09] Alan C. Bovik. *The Essential Guide to Video Processing*. Academic Press, ISBN: 978-0-12-374456-2, 2009.
- [Carletta 96] J. Carletta. *Assessing agreement on classification tasks: The kappa statistic*. Computational Linguistics, vol. 22, nr. 2, pag. 249–254, 1996.
- [CITIA 13] CITIA. *City of Moving Images, International Animated Film Festival of Annecy, France*. <http://www.citia.info>, 2013.
- [Ciuc 05] M. Ciuc & C. Vertan. *Prelucrarea Statistică a Semnalelor*. Editura MatrixRom, http://www.miv.ro/books/MCiuc_CVertan_PSS.pdf, 2005.
- [Deza 06] E. Deza & M.M. Deza. *Dictionary of Distances*. Elsevier Science, 1st edition, ISBN-10:0444520872, 2006.
- [Flickner 95] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Patkovic, D. Steele & P. Yanker. *Query by Image and Video Content: The QBIC System*. IEEE Computer, vol. 28, nr. 9, pag. 23–32, septembrie 1995.

- [Gauglitz 11] S. Gauglitz, T. Höllerer & M. Turk. *Evaluation of Interest Point Detectors and Feature Descriptors for Visual Tracking*. Int J. Comput Vis, vol. DOI 10.1007/s11263-011-0431-5, 2011.
- [Gómez-Pérez 04] A. Gómez-Pérez, M. Fernández-López & O. Corcho. *Lecture Notes: Multimedia Information Systems*. Ontological Engineering: With Examples from the Areas of Knowledge Management, E-commerce and the Semantic Web, Springer. ISBN 978-1-85233-551-9., 2004.
- [Ionescu 09] B. Ionescu. *Analiza și Prelucrarea Secvențelor Video: Indexarea Automată după Conținut*. Editura Tehnică București, ISBN 978-973-31-2354-5, 2009.
- [Ionescu 10] B. Ionescu, L. Ott, P. Lambert, D. Coquin, A. Pacureanu & V. Buzuloiu. *Tackling Action - Based Video Abstraction of Animated Movies for Video Browsing*. SPIE - Journal of Electronic Imaging, vol. 19, nr. 3, 2010.
- [Ionescu 11] B. Ionescu, C. Rasche, C. Vertan & P. Lambert. *A Contour-Color-Action Approach to Automatic Classification of Several Common Video Genres*. Springer-Verlag LNCS - Lecture Notes in Computer Science, Eds. M. Detyniecki, P. Knees, A. Nurnberger, M. Schedl and S. Stober, vol. 6817, pag. 74–88, 2011.
- [Ionescu 12a] B. Ionescu, K. Seyerlehner, C. Rasche, C. Vertan & P. Lambert. *Content-based Video Description for Automatic Video Genre Categorization*. International Conference on MultiMedia Modeling, 2012.
- [Ionescu 12b] B. Ionescu, K. Seyerlehner, C. Rasche, C. Vertan & P. Lambert. *Video Genre Categorization and Representation using Audio-Visual Information*. SPIE - Journal of Electronic Imaging, vol. 21, nr. 2, 2012.
- [Ionescu 13] B. Ionescu, J. Schlüter, I. Mironică & M. Schedl. *A Naive Mid-level Concept-based Fusion Approach to Violence Detection in Hollywood Movies*. ACM International Conference on Multimedia Retrieval, 2013.
- [Jain 89] Anil K. Jain. *Fundamentals of digital image processing*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1989.

- [Kelly 03] D. Kelly & J. Teevan. *Implicit Feedback for Inferring User Preference: a Bibliography*. International Conference on Research and Development in Information Retrieval, vol. 37, nr. 2, pag. 18–28, 2003.
- [Knees 07] P. Knees, M. Schedl, T. Pohle & G. Widmer. *Exploring Music Collections in Virtual Landscapes*. IEEE MultiMedia, vol. 14, nr. 3, pag. 46–54, 2007.
- [Knees 09] P. Knees, T. Pohle, M. Schedl, D. Schnitzer, K. Seyerlehner & G. Widmer. *Augmenting Text-Based Music Retrieval with Audio Similarity*. International Society for Music Information Retrieval, 2009.
- [Kotsiantis 07] S.B. Kotsiantis. *Supervised Machine Learning: A Review of Classification Techniques*. Informatica, vol. 31, pag. 249–268, 2007.
- [Kyungpook 06] National University Kyungpook. *Artificial Intelligence Laboratory*. <http://ailab.kyungpook.ac.kr/vindex/video-view.html>, 2006.
- [Lamel 08] L. Lamel & J.-L. Gauvain. *Speech Processing for Audio Indexing*. Int. Conf. on Natural Language Processing, vol. LNCS, 5221, pag. 4–15, 2008.
- [Lan 12] Z. Lan, L. Bao, S.-I. Yu, W. Liu & A.G. Hauptmann. *Double Fusion for Multimedia Event Detection*. International Conference on Multimedia Modeling, Klagenfurt, Austria, 2012.
- [Larson 10] Ray R. Larson. *Blind Relevance Feedback for the Image-CLEF Wikipedia Retrieval Task*. CLEF 2010 LABs and Workshops, Notebook Papers, pag. 22–23, 2010.
- [Lienhart 01] R. Lienhart. *Reliable Transition Detection in Videos: A Survey and Practitioner's Guide*. MRL, Intel Corporation, http://www.lienhart.de/Publications/IJIG_AUG2001.pdf, august, Santa Clara, USA 2001.
- [Maillet 03] S.M. Maillet. *Content-Based Video Retrieval: An Overview*. <http://vipser.unige.ch/~marchand/CBVR/>, 2003.

- [Manning 08] C.D. Manning, P. Raghavan & H. Schütze. *Introduction to Information Retrieval*. Cambridge University Press, <http://nlp.stanford.edu/IR-book/>, 2008.
- [Mathieu 10] B. Mathieu, S. Essid, T. Fillon, J. Prado & G. Richard. *YAAFE an Easy to Use and Efficient Audio Feature Extraction Software*. 11th ISMIR conference, Utrecht, Netherlands, 2010.
- [Mingqiang 08] Y. Mingqiang, K. Kidiyo & R. Joseph. *A Survey of Shape Feature Extraction Techniques*. Pattern Recognition, pag. 43–90, 2008.
- [Mironică 12a] I. Mironică, B. Ionescu & C. Vertan. *Hierarchical Clustering Relevance Feedback for Content-Based Image Retrieval*. 10th International Workshop on Content-Based Multimedia Indexing, Annecy, France 2012.
- [Mironică 12b] I. Mironică, B. Ionescu & C. Vertan. *The Influence of the Similarity Measure to Relevance Feedback*. 20th European Signal Processing Conference EUSIPCO, 2012.
- [Nakazato 01] M. Nakazato & S. T. Huang. *3D MARS: Immersive virtual reality for content based image retrieval*. IEEE International Conference on Multimedia and Exposition, pag. 45–48, 2001.
- [Nowak 10] S. Nowak & S. Rüger. *How reliable are annotations via crowdsourcing? a study about inter-annotator agreement for multi-label image annotation*. Int. Conf. on Multimedia Information Retrieval, pag. 557, 2010.
- [Orchard 91] M. Orchard & C. Bouman. *Color Quantization of Images*. IEEE Trans. on Sig. Proc., vol. 39, nr. 12, pag. 2677–2690, 1991.
- [Over 12] Paul Over, George Awad, Martial Michel, Jonathan Fiscus, Greg Sanders, Barbara Shaw, Wessel Kraaij, Alan F. Smeaton & Georges Quénot. *TRECVID 2012 – An Overview of the Goals, Tasks, Data, Evaluation Mechanisms and Metrics*. In Proceedings of TRECVID 2012. NIST, USA, 2012.

- [Reynertson 70] A. J. Reynertson. *The Work of the Film Director*. Hastings House, 1970.
- [Rocchio 71] J. Rocchio. *Relevance Feedback in Information Retrieval*. The Smart Retrieval System – Experiments in Automatic Document Processing, Prentice Hall, Englewood Cliffs NJ, pag. 313–323, 1971.
- [Rooij 08] O. Rooij, C. G. M. Snoek, & M. Worring. *Mediamill: Fast and effective video search using the ForkBrowser*. ACM International Conference on Image and Video Retrieval, 2008.
- [Rooij 10] O. Rooij, M. Worring & J. J. van Wijk. *MediaTable: Interactive Categorization of Multimedia Collections*. IEEE Computer Graphics and Applications, vol. 30, nr. 5, pag. 42–51, 2010.
- [Rubner 00] Y. Rubner, C. Tomasi & L.J. Guibas. *The Earth Mover’s Distance as a Metric for Image Retrieval*. International Journal of Computer Vision, vol. 40, nr. 2, pag. 99–121, 2000.
- [Rui 99] Y. Rui, T. Huang & S.-F. Chang. *Image Retrieval: Current Techniques, Promising Directions and Open Issues*. Journal of Visual Communication and Image Representation, vol. 10, nr. 1, pag. 39–62, 1999.
- [Schoeffmann 11] K. Schoeffmann & L. Boeszoermyeni. *Image and Video Browsing with a Cylindrical 3D Storyboard*. ACM International Conference on Multimedia Retrieval, 2011.
- [Seyerlehner 10] K. Seyerlehner, M. Schedl, T. Pohle & P. Knees. *Using Block-Level Features for Genre Classification, Tag Classification and Music Similarity Estimation*. 6th Annual Music Information Retrieval Evaluation eXchange (MIREX-10), Utrecht, Netherlands, 2010.
- [Shirahama 11] K. Shirahama & K. Uehara. *Query by Virtual Example: Video Retrieval Using Example Shots Created by Virtual Reality Techniques*. Sixth International Conference on Image and Graphics, pag. 829–834, 2011.

- [Smeulders 00] A.W.M. Smeulders, M. Worring, S. Santini, A. Gupta & R. Jain. *Content-Based Image Retrieval at the End of the Early Years*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 22, nr. 12, pag. 1349–1380, decembrie 2000.
- [Snoek 05] C. G. M. Snoek, M. Worring & A. W. M. Smeulders. *Early versus Late Fusion in Semantic Video Analysis*. ACM Multimedia, 2005.
- [Snoek 10] C.G.M. Snoek & A.W.M. Smeulders. *Video Search Engines*. IEEE Conference on Computer Vision and Pattern Recognition, <http://staff.science.uva.nl/~cgmsnoek/videosearch2010/>, 2010.
- [Stöttinger 10] Julian Stöttinger, Bogdan Tudor Goras, Nicu Sebe & Allan Hanbury. *Behavior and properties of spatio-temporal local features under visual transformations*. 2010.
- [Trémeau 04] A. Trémeau, C. Fernandez-Maloigne & P. Bonton. *Image Numérique Couleur: De l'Acquisition au Traitement*. DUNOD ISBN 2-10-006843-1, 2004.
- [Truong 07] B.T. Truong & S. Venkatesh. *Video Abstraction: A Systematic Review and Classification*. ACM Transactions on Multimedia Computing, Communications and Applications, vol. 3, nr. 1, 2007.
- [Tuceryan 93] M. Tuceryan & A. K. Jain. *Texture analysis*. The Handbook of Pattern Recognition and Computer Vision (2nd Edition), pag. 235–276, 1993.
- [Wallach 06] Hanna M. Wallach. *Topic Modeling: Beyond Bag-of-Words*. University of Cambridge, https://people.cs.umass.edu/~wallach/talks/beyond_bag-of-words.pdf, 2006.
- [Welling 05] M. Welling. *Support Vector Machines*. Note de curs, University of Toronto, Department of Computer Science, Canada, http://www.ics.uci.edu/~welling/classnotes/papers_class/SVM.pdf, 2005.

- [Witten 05] I.H. Witten & E. Frank. *Data Mining - Practical Machine Learning Tools and Techniques*. Elsevier, Morgan Kaufman Publishers, second edition, pag. 265–270, 2005.
- [Worring 03] M. Worring. *Lecture Notes: Multimedia Information Systems*. Intelligent Sensory Information Systems, University of Amsterdam, 2003.
- [Worring 12] M. Worring. *Multimedia Analytics: Exploration of Large Multimedia Collections*. keynote la International Workshop on Content-Based Multimedia Indexing, http://www.polytech.univ-savoie.fr/fileadmin/polytech_autres_sites/sites/cbmi2012/templates/fichiers/cbmi2012-worring.pdf, 2012.

