

# Detectarea automată a genurilor muzicale

Adrian Simion<sup>1</sup>, Ștefan Trăușan-Matu<sup>1,2</sup>

<sup>1</sup>Universitatea Politehnică din București  
Splaiul Independenței nr. 313, sector 6, 060042, București,  
E-mail: [simion.adrian@gmail.com](mailto:simion.adrian@gmail.com)

<sup>2</sup>Institutul de Cercetări în Inteligența Artificială  
Calea 13 septembrie nr. 13, București

**Rezumat.** Această lucrare descrie și aplică diferite metode ce folosesc un calculator pentru a determina apartenența la un anumit gen muzical a unui fișier audio. Algoritmii au fost testați pe colecțiile de date MagnaTune și MARSYAS, dar instrumentele software implementate pot fi folosite pe o gamă variată de surse. Instrumentele vor face parte dintr-un sistem software mai larg numit ADAMS (Advanced Dynamic Analysis of Music Software) dezvoltat de autori. Acest sistem este bazat pe biblioteca open-source MARSYAS și conține un modul similar cu WEKA pentru sarcini de „data mining” și „machine learning”.

**Cuvinte cheie:** segmentare automată, clasificare audio, analiză muzicală bazată pe conținut, detectarea acordurilor, regiuni vocale și instrumentale, MIR.

## 1. Introducere

Cantitatea de informație digitală a cunoscut o creștere exponențială, acest fapt fiind susținut de progresul tehnologiei dar, în cazul muzicii, și de interesul utilizatorilor ce au la dispoziție din ce în ce mai multe surse de acces la acest tip de informație.

Trendul actual este explicat, în primul rând, de o componentă importantă a muzicii: "caracteristica existențială". Muzica din România poate fi apreciată în SUA, iar în Japonia putem găsi adepți ai muzicii clasice indiene. Această formă de exprimare poate fi înțeleasă și apreciată fără etape intermediare, cum ar fi traducerea în cazul textelor. Putem afirma că „muzica este o formă de exprimare ce poate fi împărtășită de oameni ce aparțin unor culturi diferite, deoarece aceasta depășește granițele limbilor naționale sau a fundalului cultural.” (Orio 2006)

În cele din urmă, muzica este o artă cultă și populară în același timp, ba chiar, câteodată, este imposibil de a separa cele două aspecte, un exemplu fiind jazz-ul și muzica tradițională.

Disponibilitatea imediată și cererea pentru conținutul muzical a indus noi cerințe pentru administrarea, publicitatea și distribuirea acestuia. Din această cauză a apărut nevoia de o analiză mai sofisticată și directă a conținutului decât cea făcută prin intermediul meta-datelor catalogate de oameni.

Noile tehnici au permis abordări ce fuseseră folosite până în acel moment în analiza muzicală teoretică. Una dintre problemele esențiale a fost prezentată de Frank Howes: „Există un fond vast de material muzical disponibil pentru studii comparative.” (Howes 1948) Ar fi foarte interesant de a descoperi și de a stabili o corelare între muzică și fenomenul social. Cu puterea de calcul actuală și cu progresul făcut în această direcție putem răspunde la întrebări de genul: Care este fondul etnic al unei anumite piese muzicale, din ce culturi face parte?

În lumina acestor posibilități și dezvoltări ale tehnologiei aveam nevoie de o nouă disciplină care să acopere și să răspundă diferitelor probleme. MIR (Music Information Retrieval) este o știință interdisciplinară ce își extrage informațiile din muzică. Originea MIR se află în domenii ca: muzicologie, psihologie cognitivă, știința calculatoarelor și lingvistică.

Un domeniu activ de cercetare este compus din noi metode și instrumente pentru regăsirea modelelor sau pentru compararea conținutului muzical. În acest sens a fost formată ISMIR (International Society for Music Information Retrieval) care este strâns legată de MIREX (Music Information Retrieval Evaluation eXchange). Sarcinile de evaluare includ detectarea automată a genurilor, detectarea acordurilor, segmentarea, extragerea liniilor melodice, interogarea prin fredonare, doar ca să amintim câteva. Această lucrare se va concentra mai mult pe detectarea automată a genurilor și pe segmentarea automată.

În cazul de față, segmentarea automată presupune împărțirea unei surse audio în regiuni omogene. Rezultatele obținute în urma unei detecții cât mai eficiente ar putea face clasificarea informației audio mult mai facilă, sau chiar o pot automatiza complet.

## 2. Studii anterioare și înrudite legate de segmentarea automată muzicală

Subiectul clasificării limbajului/muzicii a fost studiat de către mulți cercetători. Chiar dacă aplicabilitatea rezultatelor poate fi foarte variată, multe din aceste studii folosesc colecții similare de caracteristici, cum ar fi „energia de scurtă durată”, ZCR – Zero-Crossing-Rate (Kulkarni 2007), coeficienți cepstrali, desfășurarea spectrală, centroidul și intensitatea sonoră, alături de unele caracteristici unice, cum ar fi „dinamismul”. Un cepstrum, în general, se obține prin aplicarea transformatei Fourier inverse pe logaritmul spectrului unui semnal. Totuși, combinațiile de caracteristici ce au fost folosite pot diferi foarte mult, la fel ca și mărimea colecției.

În mod uzual câteva caracteristici de termen lung, cum ar fi media sau variația, și nu caracteristicile propriu-zise sunt folosite pentru discriminare.

Diferențele majore dintre diferitele studii constau în algoritmul pentru clasificarea exactă, chiar dacă în mod normal se folosesc unii clasificatori consacrați ca bază de pornire: KNN (K-Nearest Neighbor), rețele neuronale sau gaussieni multidimensionali.

Pentru studii, de obicei, sunt folosite diferite colecții de date pentru antrenarea și testarea algoritmului. Merită menționat că pentru aceste studii, în special pentru cele timpurii, colecțiile de date erau destul de mici. Următorul tabel descrie câteva din studiile precedente.

Tabel 1. Câteva studii anterioare

Autor	Aplicabilitate	Caracteristici	Metoda de clasificare
Saunders, 1996	Monitorizare radio automatizată în timp real	STE (short-time energy), parametrii statistici ZCR	Clasificator multispacial gaussian
Scheirer și Slaney, 1997	Separarea discursului/muzicii pentru recunoașterea automată a limbajului	13 caracteristici temporale, spectrale și cepstrale (e.g., modulare la 4Hz, % din cadre cu energie joasă, derularea spectrală, spectroid central, fluxul spectral, ZCR, “ritmicitate”)	GMM (Gaussian Mixture Model), (KNN), arbori K-D, multidimensional, estimator gaussian MAP
Foote, 1997	Regăsirea documentelor audio pe baza asemănarilor acustice	12 MFCC, STE (short time energy)	Regăsirea modelelor cu ajutorul histogramelor, cuantificator de vectori bazat pe arbori antrenat pentru a maximiza informația reciprocă
Liu et al., 1997	Analiza semnalelor audio pentru clasificarea scenelor din programele	Coeficient pentru liniște, câmp dinamic pentru volum, frecvență 4Hz, media și diferența de înălțime,	O rețea neuronală ce folosește o structură de tip OCON (One-class-in-one

	televizate	raporturi pentru zgomot, centroidul frecvențelor, lățimea de bandă, energia în 4 sub-benzi	network)
Zhang și Kuo, 1999	Segmentarea/regăsirea audio pentru clasificarea scenelor video, indexarea înregistrărilor audio, parcurgerea bazelor de date	Caracteristici bazate pe energia de timp scurt, media ZCR, frecvența fundamentală pe durate scurte de timp	O procedură bazată pe reguli euristice pentru prima etapă, HMM pentru a doua etapă
Williams și Ellis, 1999	Segmentarea discursului vs. segmentarea semnalelor non verbale în sarcinile de recunoaștere automată a limbajului	Entropia medie per-cadru și probabilitatea medie „dinamică”, coeficient pt. energia de fundal, HMM	Raportul de probabilitate Gaussiană
El-Maleh al., 2000	Programarea automată și regăsirea vizuală/audio bazată pe conținut	LSF, LSF diferențial, filtru trece sus bazat pe măsurători ZCR	Clasificatori KNN și QCG (Quadratic Classifier Gaussian)
Bugatti et al., 2002	„Cuprins” pentru un document multimedia	Caracteristici bazate pe ZCR, flux spectral, energie pe termen scurt, coeficienți cepstrali, centroidul spectral, dimensiunea frecvenței silabice, raportul frecvenței înalte a spectrului puterii	Clasificator multispațial gaussian, rețele neuronale MLP
Lu, Zhang, și Jiang, 2002	Analiza conținutului audio în segmentarea video	Perechi liniare spectrale, periodicitatea benzilor, raport zgomot-cadru (NFR – noise-frame ratio), HZCRR (High zero-crossing rate ratio)	Clasificare în 3 pași: 1. KNN și perechi spectrale liniare – cuantificări vectoriale (LSP-VQ) pentru discriminarea discurs / non-discurs. 2. Reguli euristice pentru clasificarea non-discursului în muzică/zgomot de fond/liniște. 3. Segmentarea limbajului
Ajmera et al., 2003	Transcrierea automată a știrilor radiofonice	Măsurarea entropiei medii și „dinamismului” estimat la ieșirea unui perceptron multistrat (MLP – multilayer perceptron) antrenat pentru a emite probabilități posterioare.  Date de intrare MLP: primii 13 coeficienți cepstrali ai ordinului 12-a din filtrul de predicție perceptual liniar	HMM format din 2 stări cu limitări de durată (fără prag, supervizat, nesupervizat).
Burred și Lerch, 2004	Clasificarea audio (limbaj/ muzică /	Măsurători statistice a cadrelor pe timp scurt: ZCR, desfășurare/flux	

	zgomot de fond), clasificarea muzicală în genuri	centroid spectral, primii 5 MFCC, netezimea/centroidul spectrului audio, intensitatea ritmului, regularitatea ritmică, energia RMS, dimensiunea temporală, intensitatea sonoră	Clasificator KNN, clasificador GMM format din 3 componente
Barbedo și Lopes, 2006	Segmentarea automată pentru aplicații în timp real	Caracteristici bazate pe ZCR, derularea spectrală, intensitate și frecvențe fundamentale	KNN, SOM, rețele neuronale MLP, combinații liniare
Munoz-Expósito et al., 2006	Sistem inteligent pentru codarea audio automată	Centroid spectral LPC răsucit	GMM format din 3 componente cu sau fără sistem bazat pe reguli fuzzy
Alexandre et al, 2006	Clasificarea discursului / muzicii pentru clasificarea genurilor muzicale	Desfășurarea/centroidul spectral, ZCR, energie pe termen scurt, LSTER (low short time energy ratio), MFCC, VTW (voice to white)	Discriminant liniar Fisher, KNN

### 3. Semnale digitale audio

În momentul în care muzica este înregistrată, presiunea continuă a undei sonore este măsurată cu ajutorul unui microfon. Aceste măsurători sunt luate la intervale regulate de timp și fiecare măsurătoare este cuantificată.

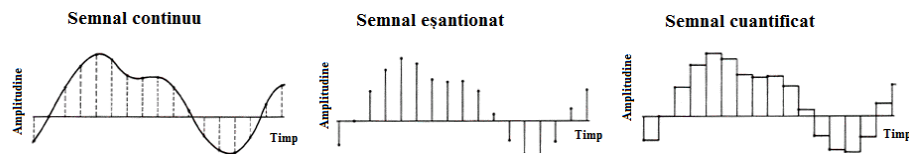


Figura 1. Reprezentare digitală a sunetului (în timp)

a. Muzica este un semnal    b. care este eșantionat    c. și cuantificat continuu;

Sunetele pot fi reprezentate ca o sumă de sinusoidale. Un semnal format din  $N$  eșantioane poate fi scris astfel:

$$x = \sum_{k=0}^{N/2} a_k^{(r)} \cos\left(2\pi\left(\frac{k}{N}\right)\right) + a_k^{(i)} \sin\left(2\pi\left(\frac{k}{N}\right)\right). \quad (1)$$

Semnalul poate fi reprezentat în domeniul *frecvență* folosind coeficienții  $\{(a_1^{(y)}, a_1^{(i)}), \dots, (a_{N/2}^{(y)}, a_{N/2}^{(i)})\}$ .

Amplitudinea și faza frecvenței  $k$  sunt date de:

$$X_M[k] = \sqrt{(a_k^{(r)})^2 + (a_k^{(i)})^2} \quad (2)$$

$$X_p[k] = \arctan\left(\frac{a_k^{(i)}}{a_k^{(r)}}\right) \quad (3)$$

Studii perceptuale bazate pe auzul uman demonstrează că informația dată de fază este relativ neimportantă atunci când aceasta este comparată cu informația preluată din amplitudine, astfel componenta fazică este de obicei ignorată în momentul extragerii caracteristicilor. [19]

*Spectroidul central* este o altă caracteristică spectrală ce este utilă în procesul de extragere și de analiză. În tabelul 1 se pot observa diferitele utilizări. Spectroidul central este centrul de gravitate al spectrului și este dat de:

$$C = \frac{\sum_{k=1}^{N/2} X_M[k] * k}{\sum_{k=1}^{N/2} X_M[k]} \quad (4)$$

Spectroidul central poate fi înțeles ca fiind o unitate de măsură a luminozității deoarece cântecele sunt considerate mai luminoase atunci când au componente ce conțin frecvențe mai înalte.

### 3.1 Transformări în domeniile timp-frecvență

În MIR și în analiza sunetelor în general este foarte uzual să facem o transformare între domeniile timp și frecvență. Pentru acest lucru aparatul matematic ne dă transformata Fourier discretă, transformata Fourier rapidă, transformata Fourier, transformata cosinus discretă, transformata undină discretă și transformata gammaton.

Analiza muzicală nu se ocupă cu transformarea în mulțimea complexă, deoarece muzica este întotdeauna o serie de valori reale în timp și are doar frecvențe pozitive.

Fiind dat un semnal  $x$  cu  $N$  eșantioane, funcțiile de bază pentru transformata Fourier discretă va fi  $N/2$  unde sinus și  $N/2$  unde cosinus ce corespund coeficienților anteriori.

Operatorul de proiecție este corelarea, o unitate a similarității între două serii temporale. Coeficienții pot fi aflați folosind:

$$a_k^{(r)} = \frac{2}{N} \sum_{i=0}^{N-1} x[i] \cos(2\pi \frac{k}{N} i) \quad (5)$$

$$a_k^{(i)} = -\frac{2}{N} \sum_{i=0}^{N-1} x[i] \sin(2\pi \frac{k}{N} i) \quad (6)$$

Transformata Fourier discretă este calculată într-un mod eficient cu ajutorul transformatei Fourier rapide. O limitare a celor două reprezentări, cea temporală și cea spectrală, este dată de faptul ca niciuna dintre reprezentări nu prezintă simultan informația obținută din analiza frecvenței și din analiza timpului. O reprezentare timp-frecvență este găsită folosind transformata Fourier de timp-scurt: Mai întâi, semnalul audio este divizat în secvențe (suprapuse) de segmente. Fiecare segment este multiplicat de o *funcție fereastră*.

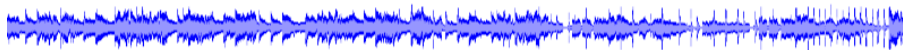


Figura 2. Magnatune apa\_ya-apa\_ya-14-maani-59-88.wav (domeniu temporal)

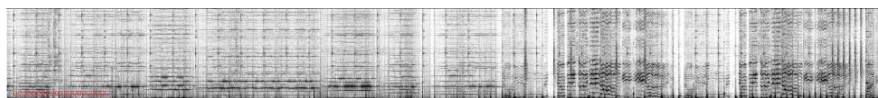


Figura 3. Magnatune apa\_ya-apa\_ya-14-maani-59-88.wav (spectrogramă)

Figurile 2 și 3 au fost obținute cu ajutorul unei versiuni modificate a executabilului MARSAYS sound2png cu ajutorul următoarelor comenzi:

- `./sound2png -m waveform ../audio/magnatune/0/apa_ya-apa_ya-14-maani-59-88.wav ../saveres/magnatunewav.png -ff Adventure.ttf`
- `./sound2png -m spectrogram ../audio/magnatune/0/apa_ya-apa_ya-14-maani-59-88.wav ../saveres/magnatunespec.png -ff Adventure.ttf`

O altă transformare utilă este cea bazată pe transformata undină.

### 3.2 MFCC (Mel-Frequency Cepstral Coefficients)

Cel mai uzual set de caracteristici folosit în recunoașterea limbajului și în sistemele de adnotare muzicală sunt MFCC. Aceștia sunt de fapt trăsături de scurtă durată ce caracterizează amplitudinea spectrală a semnalului audio. Pentru fiecare interval scurt de timp (25 ms), este găsit vectorul de caracteristici folosind un algoritm explicat în Alg1. Primul pas este de a obține amplitudinea fiecărei frecvențe în domeniul frecvență folosind transformata Fourier discretă. După acest pas se logaritmează amplitudinea deoarece intensitatea perceptuală s-a dovedit a fi aproximativ logaritmică. În pasul următor componentele ce formează frecvența sunt unite în 40 de cadre ce au fost determinate conform scării Mel. Scara Mel este maparea dintre frecvența reală și modelul perceput de frecvență ca fiind aproximativ logaritmic. Deoarece o secvență temporală a acestor vectori 40-dimensionali Mel-frecvență ar fi fost foarte redundantă, am putea reduce aceste dimensiuni folosind PCA. În locul acestei abordări, comunitatea a adoptat transformata cosinus discretă, ce aproximează PCA dar nu are nevoie de date de antrenament, pentru a reduce numărul de dimensiuni la un număr de 13 MFCC. (Turnbull 2005)

- |   |
|---|
| 1: Calcularea spectrului folosind transformata Fourier discretă |
| 2: Logaritizarea spectrului                                     |
| 3: Aplicarea scării Mel netezirea                               |
| 4: Decorelare folosind transformata cosinus discretă            |

*Algoritm 1. Calcularea vectorului de caracteristici MFCC*

## 4. Descrierea problemei

O caracteristică obișnuită ce-i ajută pe producătorii de discuri să răspundă cerințelor celor din publicul țintă, pe muzicologi să studieze influențele muzicale și pe entuziaști să-și organizeze colecțiile muzicale este identificarea genurilor.

Conceptul de gen este în mod natural subiectiv deoarece influențele, ierarhia și intersecția unui cântec cu un anumit gen nu este un lucru stabilit de comun acord. Acest punct de vedere este susținut de o comparație a trei furnizori de servicii muzicale online ce a arătat că există mari diferențe în



numărul genurilor, cuvintele ce le descriu și structura ierarhiei acestora. (Pachet 2000)

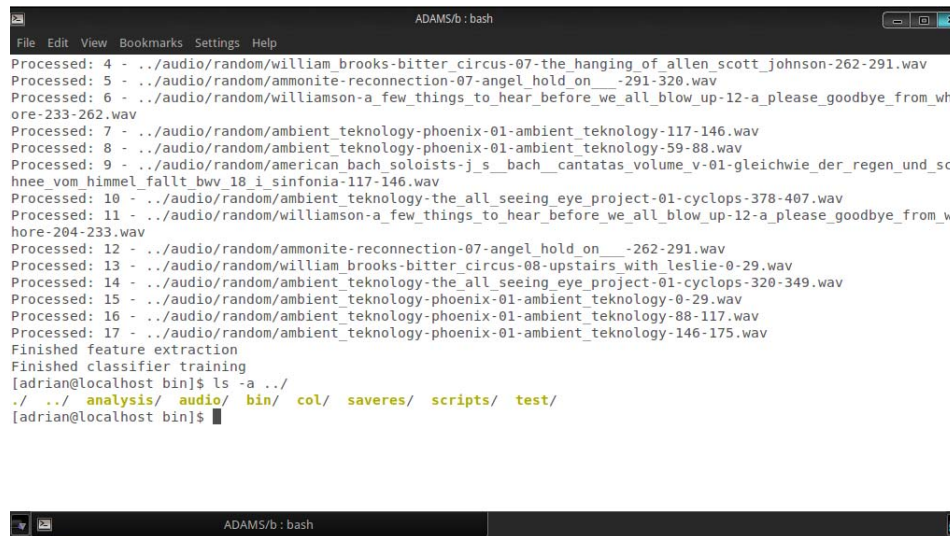
Chiar dacă există multe contradicții cauzate de natura subiectivă, conceptul de gen a stârnit interesul comunității MIR. Diferite lucrări legate de acest subiect reflectă presupunerile autorilor în legătură cu genurile. Legile drepturilor de autor au împiedicat cercetătorii să stabilească o bază de date comună de cântece, făcând foarte dificilă compararea rezultatelor.

## 5. Descrierea experimentului

Colecțiile de date folosite pentru antrenament și testare au fost MAGNATUNE (2012) și două colecții ce au fost construite în etapele timpurii ale MARSYAS.

Autorii au creat un sistem software plecând de la biblioteca open-source MARSYAS pentru a facilita acest experiment și pentru alte sarcini de analiză audio ulterioare. Deoarece sistemul ADAMS (Advanced Dynamic Analysis of Music Software) este construit într-un design modular, diferitele sarcini (descrise mai jos) pot fi automatizate, iar sunetul poate „curge” prin aceste module până când este făcută întreaga analiză.

Structura directorului ADAMS poate fi observată în figura următoare:



```
ADAMS/b : bash
File Edit View Bookmarks Settings Help
Processed: 4 - ../audio/random/william_brooks-bitter_circus-07-the_hanging_of_allen_scott_johnson-262-291.wav
Processed: 5 - ../audio/random/ammonite-reconnection-07-angel_hold_on___-291-320.wav
Processed: 6 - ../audio/random/williamson-a_few_things_to_hear_before_we_all_blow_up-12-a_please_goodbye_from_w
ore-233-262.wav
Processed: 7 - ../audio/random/ambient_teknology-phoenix-01-ambient_teknology-117-146.wav
Processed: 8 - ../audio/random/ambient_teknology-phoenix-01-ambient_teknology-59-88.wav
Processed: 9 - ../audio/random/american_bach_soloists-j_s_bach_cantatas_volume_v-01-gleichwie_der_regen_und_sc
hnee_vom_himmel_fallt_bwv_18_i_sinfonia-117-146.wav
Processed: 10 - ../audio/random/ambient_teknology-the_all_seeing_eye_project-01-cyclops-378-407.wav
Processed: 11 - ../audio/random/williamson-a_few_things_to_hear_before_we_all_blow_up-12-a_please_goodbye_from_w
hore-204-233.wav
Processed: 12 - ../audio/random/ammonite-reconnection-07-angel_hold_on___-262-291.wav
Processed: 13 - ../audio/random/william_brooks-bitter_circus-08-upstairs_with_leslie-0-29.wav
Processed: 14 - ../audio/random/ambient_teknology-the_all_seeing_eye_project-01-cyclops-320-349.wav
Processed: 15 - ../audio/random/ambient_teknology-phoenix-01-ambient_teknology-0-29.wav
Processed: 16 - ../audio/random/ambient_teknology-phoenix-01-ambient_teknology-88-117.wav
Processed: 17 - ../audio/random/ambient_teknology-phoenix-01-ambient_teknology-146-175.wav
Finished feature extraction
Finished classifier training
[adrian@localhost bin]$ ls -a ../
./ ../ analysis/ audio/ bin/ col/ saveres/ scripts/ test/
[adrian@localhost bin]$
```

Figura 4. Structura principală a sistemului ADAMS

Sarcinile de „învățare automată” (machine learning) au fost realizate prin intermediul WEKA (2012), încărcând fișierele compatibile arff create cu ajutorul MARSYAS.

Sistemul de operare ales pentru aceste experimente a fost Mandriva Linux 2011, versiunea compilatorului fiind „gcc (GCC) 4.6.1 20110627 (Mandriva)”.

Au fost folosiți următorii extractori:

- BEAT: Caracteristici histogramice legate de ritm (Beat histogram features)
- LPCC: Coeficienți cepstrali LPC derivați (LPC derived Cepstral coefficients)
- LSP: Perechi spectrale liniare (Linear Spectral Pairs)
- MFCC: Mel-Frequency Cepstral Coefficients
- SCF: Spectral Crest Factor (MPEG-7)
- SFM: Unitate de măsură spectrală a netezimii (Spectral Flatness Measure MPEG-7)
- SFMSCF: caracteristici SCF și SFM
- STFT: Centroid, Desfășurare spectrală, Flux, ZC – Valori pozitive (Zero-Crossings)
- STFTMFCC: Centroid, Flux-Desfășurare spectrală, ZC, MFCC

La fiecare experiment pentru extractorii specificați sunt prezentate de asemenea matricele de confuzie pentru a avea o idee despre clasificarea reală și cea prezisă realizată cu ajutorul sistemului de clasificare.

## 5.1 Experimentul 1: Clasificarea folosind „Caracteristici timbrale”

Acest experiment folosește următorii extractori: TZC (Time Zero-Crossings), Centroidul Spectral, Flux-Desfășurare spectrală și MFCC. Extragem aceste caracteristici cu opțiunea –timbral și creăm de asemenea un fișier ce va fi încărcat cu ajutorul mediului WEKA pentru analiză cu ajutorul comenzii:

```
./adamsfeature -sv -timbral ../col/all.mf -w
../analysis/alltimbral.arff
```

Pentru celelalte etape s-au folosit comenzi similare. În urma experimentelor s-au ales următorii clasificatori: BN (Bayes Network), NB (Naive Bayes), DT (Decision Table), FC (Filtered Classifier) și NNGE.

Rezultatele sunt evidențiate în următorul tabel:

Tabel 2. Caracteristici Timbrale – Rezultatele clasificatorilor

Clasificator	Timp Cons model	Clasif. corect	Clasif. incorect	Er. Med. Abs.	Er. Med. Pătr.	Er. Rel. Abs.	Er. Rel. Pătr.
Bayes Network	1.78	62.5%	37.5%	0.0753	0.2648	41.82%	88.28%
Naive Bayes	0.04	55%	45%	0.0902	0.2925	50.09%	97.51%
Decision Table	15.49	51.6%	48.4%	0.1467	0.2599	81.53%	86.64%
Filtered Classifier	4.55	87.8%	12.2%	0.0348	0.1318	19.31%	43.94%
NNGE	10.69	100%	0%	0	0	0	0

Tabelul 2 a fost construit încărcând fișierul alltimbral.arff în WEKA și antrenând clasificatorii existenți.

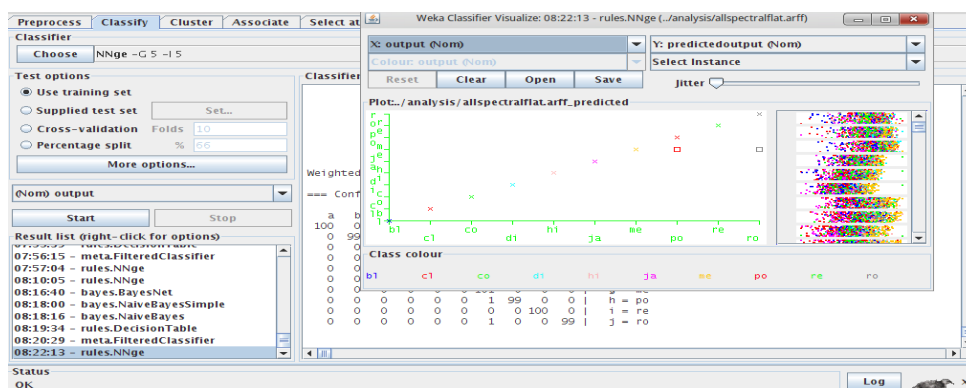


Figura 5. Grafic de predicție a erorilor în mediul WEKA

```

=== Confusion Matrix === Bayes Network
  a b c d e f g h i j <-- classified as
63 0 5 7 2 1 6 8 2 6 | a = bl
4 82 0 0 0 10 0 0 0 3 | b = cl
6 2 66 7 0 7 4 1 1 6 | c = co
1 1 6 64 4 1 2 6 9 6 | d = di
0 0 0 17 45 1 2 16 19 0 | e = hi
19 14 2 1 0 58 2 1 1 2 | f = ja
2 2 2 2 2 4 77 2 1 7 | g = me
7 0 1 10 7 2 1 66 3 3 | h = po
1 1 3 11 9 4 0 5 61 5 | i = re
5 1 8 14 0 9 14 1 5 43 | j = ro

=== Confusion Matrix === Naive Bayes
  a b c d e f g h i j <-- classified as
40 0 14 8 3 13 6 4 0 12 | a = bl
0 90 0 0 0 7 0 0 0 2 | b = cl
9 3 59 3 0 2 5 1 0 18 | c = co
1 0 5 45 5 2 3 2 6 31 | d = di
3 0 1 14 51 0 2 7 12 10 | e = hi
0 0 1 1 3 5 67 2 0 22 | g = me
2 0 4 10 11 2 2 51 9 9 | h = po
2 0 19 8 8 3 0 2 44 14 | i = re
4 3 11 3 0 4 12 1 5 57 | j = ro

=== Confusion Matrix === Decision table
  a b c d e f g h i j <-- classified as
26 0 14 12 1 5 5 9 23 5 | a = bl
2 73 0 0 0 15 4 0 0 5 | b = cl
2 1 66 11 0 4 3 2 1 10 | c = co
9 1 6 44 4 3 2 7 23 1 | d = di
3 0 4 15 45 0 0 9 24 0 | e = hi
9 10 4 5 0 57 1 0 6 8 | f = ja
5 1 9 11 0 1 60 4 0 10 | g = me
6 0 2 15 17 0 2 47 7 4 | h = po
7 0 4 7 10 2 0 2 67 1 | i = re
2 2 11 19 1 5 16 3 10 31 | j = ro

=== Confusion Matrix === Filtered classifier
  a b c d e f g h i j <-- classified as
95 0 3 0 0 0 0 0 1 1 0 | a = bl
3 91 0 0 0 4 0 0 0 1 | b = cl
7 0 87 0 0 1 2 1 0 2 | c = co
2 0 6 91 0 0 0 0 1 0 | d = di
2 0 1 4 86 0 1 2 4 0 | e = hi
2 3 0 1 1 93 0 0 0 0 | f = ja
2 1 0 3 1 1 90 1 0 2 | g = me
2 1 3 2 3 1 1 86 0 1 | h = po
1 0 2 4 4 1 2 0 84 2 | i = re
1 0 7 2 2 6 3 2 2 75 | j = ro

=== Confusion Matrix === NNGE
  a b c d e f g h i j <-- classified as
100 0 0 0 0 0 0 0 0 0 | a = bl
0 99 0 0 0 0 0 0 0 0 | b = cl
0 0 100 0 0 0 0 0 0 0 | c = co
0 0 0 100 0 0 0 0 0 0 | d = di
0 0 0 0 100 0 0 0 0 0 | e = hi
0 0 0 0 0 100 0 0 0 0 | f = ja
0 0 0 0 0 0 100 0 0 0 | g = me
0 0 0 0 0 0 0 100 0 0 | h = po
0 0 0 0 0 0 0 0 100 0 | i = re
0 0 0 0 0 0 0 0 0 100 | j = ro

```

Figura 6. Matricele de confuzie pentru clasificarea cu ajutorul caracteristicilor timbrale

## 5.2 Experimentul 2: Clasificare folosind „Caracteristici spectrale”

Acest experiment folosește următorii extractori: Centroidul Spectral și Flux-Desfășurare spectrală.

Folosind aceiași clasificatori rezultatele sunt:

Tabel 3. Caracteristici Spectrale – Rezultatele clasificatorilor

Clasificator	Timp Cons model	Clasif. corect	Clasif. incorect	Er. Med. Abs.	Er. Med. Pătr.	Er. Rel. Abs.	Er. Rel. Pătr.
Bayes Network	1.78	46.5%	53.5%	0.1192	0.2742	66.21%	91.41%
Naive Bayes	0.23	42.5%	57.5%	0.1205	0.2924	66.92%	97.47%
Decision Table	0.72	46.1%	53.9%	0.1491	0.2655	82.82%	88.49%
Filtered Classifier	0.41	63.6%	36.4%	0.099	0.2225	54.98%	74.15%
NNGE	20.2	100%	0%	0	0	0	0

```

=== Confusion Matrix === Bayes Network
  a b c d e f g h i j <-- classified as
41 2 11 6 0 9 16 2 7 6 | a = bl
 1 76 2 0 0 11 3 0 0 6 | b = cl
16 6 36 3 0 13 12 2 3 9 | c = co
 4 0 4 41 11 0 10 18 5 1 | d = di
 4 0 2 15 40 0 1 19 18 1 | e = hi
10 23 12 3 1 34 5 6 0 6 | f = ja
 4 1 4 4 1 3 80 2 1 1 | g = me
 5 0 3 3 21 1 2 53 7 5 | h = po
 9 0 4 5 21 0 2 7 51 1 | i = re
14 8 8 10 1 5 23 9 9 13 | j = ro

=== Confusion Matrix === Naive Bayes
  a b c d e f g h i j <-- classified as
45 2 16 3 1 0 25 1 5 2 | a = bl
 2 81 2 0 2 5 5 0 0 2 | b = cl
19 13 37 2 1 3 15 2 0 8 | c = co
11 0 6 41 2 0 21 12 6 1 | d = di
 5 0 1 20 16 0 4 24 26 4 | e = hi
16 34 20 0 2 12 6 5 1 4 | f = ja
 5 0 7 8 0 0 77 2 0 2 | g = me
 3 2 6 14 3 1 2 57 5 7 | h = po
11 0 5 10 2 1 3 15 51 2 | i = re
11 10 13 11 4 2 33 4 4 8 | j = ro

=== Confusion Matrix === Decision table
  a b c d e f g h i j <-- classified as
31 4 21 10 1 12 16 2 2 1 | a = bl
 1 79 2 0 0 12 4 0 0 1 | b = cl
11 9 44 9 1 11 9 2 1 3 | c = co
 5 0 10 47 4 1 11 15 3 4 | d = di
 2 0 4 10 39 0 3 22 20 0 | e = hi
14 23 8 5 2 35 5 5 0 3 | f = ja
 0 2 6 11 1 5 62 3 0 11 | g = me
 7 4 1 8 18 1 0 57 3 1 | h = po
 7 0 6 9 16 2 2 9 49 0 | i = re
11 4 13 14 3 9 18 5 5 18 | j = ro

=== Confusion Matrix === Filtered classifier
  a b c d e f g h i j <-- classified as
74 1 6 2 0 8 4 0 1 4 | a = bl
 1 81 2 0 0 13 1 0 0 1 | b = cl
10 8 52 5 0 13 6 1 1 4 | c = co
 3 0 4 66 3 1 8 0 4 3 | d = di
 5 0 0 11 52 0 1 13 16 2 | e = hi
 9 15 4 0 1 63 4 1 0 3 | f = ja
 5 1 3 5 0 3 77 1 1 5 | g = me
 5 1 5 5 12 0 0 66 3 3 | h = po
 7 0 5 6 12 0 1 7 61 1 | i = re
14 4 8 4 2 3 14 4 3 44 | j = ro

=== Confusion Matrix === NNGE
  a b c d e f g h i j <-- classified as
100 0 0 0 0 0 0 0 0 0 | a = bl
 0 99 0 0 0 0 0 0 0 0 | b = cl
 0 0 100 0 0 0 0 0 0 0 | c = co
 0 0 0 100 0 0 0 0 0 0 | d = di
 0 0 0 0 100 0 0 0 0 0 | e = hi
 0 0 0 0 0 100 0 0 0 0 | f = ja
 0 0 0 0 0 0 101 0 0 0 | g = me
 0 0 0 0 0 0 0 100 0 0 | h = po
 0 0 0 0 0 0 0 0 100 0 | i = re
 0 0 0 0 0 0 0 0 0 100 | j = ro
    
```

Figura 7. Matricele de confuzie pentru clasificare folosind caracteristici spectrale

### 5.3 Experimentul 3: Clasificare folosind „MFCC”

Acest experiment folosește extractorii bazați pe MFCC (Mel-Frequency Cepstral Coefficients).

Table 4. Caracteristici MFCC – Rezultatele clasificării

Clasificator	Timp Cons model	Clasif. corect	Clasif. incorect	Er. Med. Abs.	Er. Med. Pătr.	Er. Rel. Abs.	Er. Rel. Pătr.
Bayes Network	1.23	63.3%	36.7%	0.0764	0.2475	42.42%	82.50%
Naive Bayes	0.22	58.5%	41.5%	0.0847	0.2694	47.07%	89.80%
Decision Table	6.4	49.1%	50.9%	0.1481	0.2638	82.27%	87.94%
Filtered Classifier	0.81	87.1%	12.9%	0.0363	0.1348	20.18%	44.92%
NNGE	3.74	99.8%	0.2%	0.0004	0.02	0.22%	6.66%

```

=== Confusion Matrix === Bayes Network
a b c d e f g h i j <-- classified as
47 1 6 13 0 1 6 11 8 7 | a = bl
0 92 0 0 0 3 1 0 0 3 | b = cl
5 3 69 3 1 0 4 5 1 9 | c = co
3 0 5 48 1 1 9 10 10 13 | d = di
2 0 0 11 58 0 5 1 22 1 | e = hi
3 10 1 1 0 77 2 0 1 5 | f = ja
3 0 1 0 0 0 84 6 0 7 | g = ne
1 0 9 13 9 3 3 45 12 5 | h = po
6 0 2 9 6 0 0 2 68 7 | i = re
4 0 13 8 3 3 18 5 1 45 | j = ro

=== Confusion Matrix === Naive Bayes
a b c d e f g h i j <-- classified as
45 0 8 12 0 4 8 4 7 12 | a = bl
0 93 0 0 0 2 1 0 0 3 | b = cl
6 3 55 14 1 0 3 2 2 14 | c = co
2 0 5 47 1 1 16 6 10 12 | d = di
1 0 0 13 64 0 5 5 9 3 | e = hi
6 9 0 3 1 57 6 0 0 18 | f = ja
5 0 0 2 0 0 87 3 1 3 | g = ne
2 0 6 22 10 5 3 31 9 12 | h = po
6 0 3 12 7 1 0 2 61 8 | i = re
2 0 9 11 4 2 23 2 2 45 | j = ro

=== Confusion Matrix === Decision table
a b c d e f g h i j <-- classified as
36 1 14 8 2 1 7 22 9 0 | a = bl
0 74 0 0 0 21 1 0 0 3 | b = cl
18 2 29 8 1 4 3 19 14 2 | c = co
4 0 9 47 7 1 5 17 7 3 | d = di
5 0 3 13 45 1 1 19 13 0 | e = hi
2 25 2 0 0 60 4 1 3 3 | f = ja
4 0 1 16 0 1 70 7 0 2 | g = ne
6 0 4 19 6 1 3 58 3 0 | h = po
13 0 6 2 11 2 1 8 57 0 | i = re
11 0 11 17 2 4 18 15 7 15 | j = ro

=== Confusion Matrix === Filtered classifier
a b c d e f g h i j <-- classified as
91 0 1 2 0 1 2 1 1 1 | a = bl
0 95 1 0 0 1 2 0 0 0 | b = cl
5 0 86 1 0 1 0 4 0 3 | c = co
3 0 1 85 2 0 2 2 1 4 | d = di
3 0 0 2 88 1 3 2 1 0 | e = hi
0 3 1 2 1 92 0 0 0 1 | f = ja
1 0 0 2 0 0 93 1 1 3 | g = ne
3 0 2 3 2 0 1 87 1 1 | h = po
6 0 1 1 7 0 0 3 79 3 | i = re
4 0 5 4 1 1 4 5 1 75 | j = ro

=== Confusion Matrix === NNGE
a b c d e f g h i j <-- classified as
100 0 0 0 0 0 0 0 0 0 | a = bl
0 99 0 0 0 0 0 0 0 0 | b = cl
0 0 100 0 0 0 0 0 0 0 | c = co
0 0 0 100 0 0 0 0 0 0 | d = di
0 0 0 0 100 0 0 0 0 0 | e = hi
0 0 0 0 0 100 0 0 0 0 | f = ja
0 0 0 0 0 0 100 0 0 0 | g = ne
0 0 0 0 0 0 0 100 0 0 | h = po
0 0 0 0 0 0 0 0 100 0 | i = re
0 0 0 0 0 0 0 0 0 100 | j = ro

```

Figura 8. Matricele de confuzie pentru clasificarea bazată pe caracteristici MFCC

## 5.4 Experimentul 4: Clasificare folosind „ZC”

Acest experiment folosește extractorii bazați pe ZC.

Tabel 5. Caracteristici ZC – Rezultatele clasificării

Clasificator	Timp Cons model	Clasif. corect	Clasif. incorect	Er. Med. Abs.	Er. Med. Pătr.	Er. Rel. Abs.	Er. Rel. Pătr.
Bayes Network	0.09	34.7%	65.3%	0.1437	0.2789	79.83%	82.50%
Naive Bayes	0.01	34.5%	65.5%	0.1441	0.2869	80.06%	89.80%
Decision Table	0.22	42.4%	57.6%	0.1511	0.2691	83.95%	87.94%
Filtered Classifier	0.15	44%	56%	0.1403	0.2649	77.94%	44.92%
NNGE	0.52	99.8%	0.2%	0.0004	0.02	0.22%	6.66%

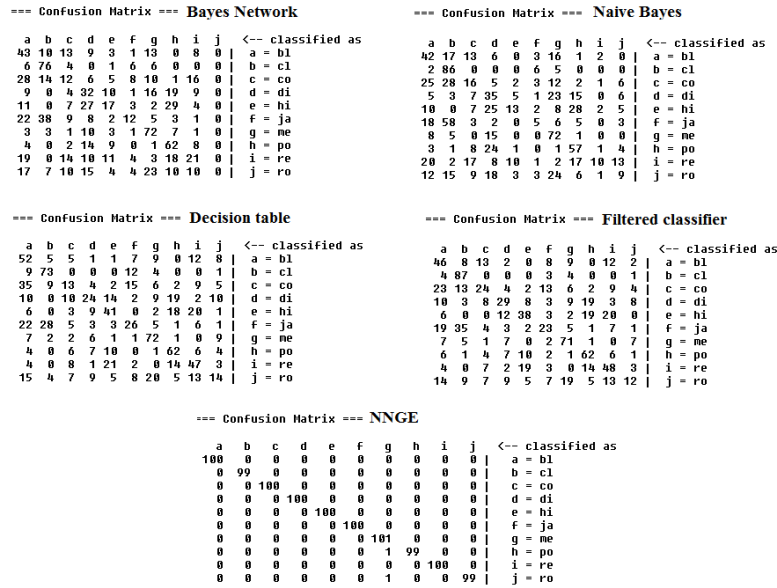


Figura 9. Matricele de confuzie pentru clasificarea bazată pe caracteristici ZC

### 5.5 Experimentul 5: Clasificare folosind „SFM”

Tabel 6. Caracteristici SFM – Rezultatele clasificării

Clasificator	Timp Cons model	Clasif. corect	Clasif. incorect	Er. Med. Abs.	Er. Med. Pătr.	Er. Rel. Abs.	Er. Rel. Pătr.
Bayes Network	1.78	58.4%	41.6%	0.0838	0.2738	46.53%	91.28%
Naive Bayes	0.15	53.2%	46.8%	0.0935	0.294	51.96%	97.99%
Decision Table	12.35	50.4%	49.6%	0.1472	0.2621	81.78%	87.37%
Filtered Classifier	2.1	83.8%	16.2%	0.045	0.15	25.01%	50.12%
NNGE	9.24	99.8%	0.2%	0.0004	0.02	0.22%	6.66%

Acest experiment folosește extractorii bazați pe SFM (Spectral Flatness Measure).

```

=== Confusion Matrix === Bayes Network
  a b c d e f g h i j <-- classified as
39 0 9 13 4 12 0 1 13 9 | a = bl
 0 78 12 0 0 8 1 0 0 0 | b = cl
2 10 50 13 0 5 0 5 3 12 | c = co
2 1 5 63 9 0 2 5 4 9 | d = di
3 0 0 7 65 6 10 6 1 2 | e = hi
5 13 4 1 2 61 2 6 0 6 | f = ja
0 0 0 8 1 0 82 1 1 8 | g = me
3 1 5 10 4 4 3 55 5 10 | h = po
2 0 3 10 6 2 6 10 52 9 | i = re
2 1 12 18 3 4 14 5 2 39 | j = ro

=== Confusion Matrix === Naive Bayes
  a b c d e f g h i j <-- classified as
34 0 4 23 3 11 0 0 16 9 | a = bl
17 0 17 0 0 10 1 0 0 0 | b = cl
2 8 39 20 0 6 0 7 1 17 | c = co
0 1 1 62 5 1 9 3 6 12 | d = di
1 0 0 9 61 7 10 3 5 4 | e = hi
17 9 4 5 2 49 3 4 1 6 | f = ja
0 0 0 7 0 0 83 1 1 9 | g = me
6 1 5 15 5 2 4 52 1 9 | h = po
6 0 2 15 6 4 7 15 39 6 | i = re
3 0 13 16 4 2 15 1 3 43 | j = ro

=== Confusion Matrix === Decision table
  a b c d e f g h i j <-- classified as
43 0 9 14 1 3 1 4 23 2 | a = bl
10 78 3 0 0 6 1 0 1 0 | b = cl
10 13 45 14 0 2 1 2 1 12 | c = co
8 1 8 43 1 3 1 10 14 11 | d = di
3 0 0 10 44 4 8 15 14 2 | e = hi
13 9 7 6 2 51 1 4 4 3 | f = ja
2 0 0 9 0 0 81 1 3 5 | g = me
9 3 3 12 11 8 4 32 13 5 | h = po
11 1 1 12 8 4 1 8 52 2 | i = re
8 2 8 19 0 5 13 6 4 35 | j = ro

=== Confusion Matrix === Filtered classifier
  a b c d e f g h i j <-- classified as
88 1 1 1 1 4 1 2 1 0 | a = bl
2 92 3 0 0 2 0 0 0 0 | b = cl
4 2 93 0 0 0 0 1 0 0 | c = co
2 0 9 82 1 1 1 1 2 1 | d = di
3 1 1 5 86 0 0 2 1 1 | e = hi
3 3 1 1 2 88 0 1 1 0 | f = ja
0 0 2 3 1 0 94 1 0 0 | g = me
2 2 3 3 3 2 1 82 2 0 | h = po
4 1 3 2 4 3 0 9 73 1 | i = re
3 2 12 3 4 2 8 2 4 60 | j = ro

=== Confusion Matrix === NNGE
  a b c d e f g h i j <-- classified as
100 0 0 0 0 0 0 0 0 0 | a = bl
 0 99 0 0 0 0 0 0 0 0 | b = cl
 0 0 100 0 0 0 0 0 0 0 | c = co
 0 0 0 100 0 0 0 0 0 0 | d = di
 0 0 0 0 100 0 0 0 0 0 | e = hi
 0 0 0 0 0 100 0 0 0 0 | f = ja
 0 0 0 0 0 0 101 0 0 0 | g = me
 0 0 0 0 0 0 0 1 99 0 | h = po
 0 0 0 0 0 0 0 0 0 100 | i = re
 0 0 0 0 0 0 0 0 1 0 99 | j = ro

```

Figura 10. Matricele de confuzie pentru clasificarea bazată pe caracteristici SFM

## 6. Concluzii

Au fost făcute cinci experimente pentru determinarea genului muzical al fișierelor audio. Caracteristicile extrase au variat de la un experiment la altul pentru a determina care se pretează mai bine pentru colecția de date folosită. Cei cinci clasificatori au dat rezultate diferite în funcție de caracteristicile extrase, iar acestea au fost evaluate cu unelte software de „machine-learning” bine-cunoscute cum ar fi WEKA și de asemenea cu un sistem de analiză dezvoltat pe baza bibliotecii MARSYAS.

Rezultatele arată că se pot obține concluzii satisfăcătoare chiar și prin abordări simple cum ar fi clasificarea NB (Naive Bayes), dar rezultate cu o relevanță mai mare au fost obținute folosind tehnici mai avansate. Faptul că NN (cel mai apropiat vecin – Nearest Neighbor) a produs rezultate foarte bune nu înseamnă că va avea același comportament pe o altă colecție de date.



Îmbunătățiri ale metodelor prezentate pot fi obținute testându-le pe o colecție extinsă de date și determinând influențele intrinsece ale fiecărui gen asupra altuia.

Concluziile acestor influențe pot avea un înțeles mai mare din punct de vedere social, cum ar fi analiza blues-ului și modul în care au apărut derivatele sale. În unele cazuri putem găsi rezultate improbabile cum ar fi că genul muzical death-metal își are rădăcinile în muzica jazz.

## Referințe

- Ajmera J., McCowan I., și Bourlard H., *Speech/music segmentation using entropy and dynamism features in a HMM classification framework*, Speech Communication, vol. 40, nr. 3, pp. 351-363, 2003
- Alexandre E., Rosa M., L. Caudra, și Gil-Pita R., *Application of Fisher linear discriminant analysis to speech/music classification*, Proceedings of the 120th Audio Engineering Society Convention (AES '06), Paris, France, Mai 2006, nr 6678
- Barbedo J. G. A. și Lopes A., *A robust and computationally efficient speech/music discriminator*, Journal of the Audio Engineering Society, vol. 54, nr. 7-8, pp. 571–588, 2006
- Bugatti A., Flammini A., și Migliorati P., *Audio classification in speech and music: a comparison between a statistical and a neural approach*, EURASIP Journal on Applied Signal Processing, vol. 2002, nr. 4, pp. 372–378, 2002.
- Burred J. J. și Lerch A., *Hierarchical automatic audio signal classification*, Journal of the Audio Engineering Society, vol. 52, no. 7-8, pp. 724–739, 2004
- CONFUSION MATRIX  
[http://www2.cs.uregina.ca/~hamilton/courses/831/notes/confusion\\_matrix/confusion\\_matrix.html](http://www2.cs.uregina.ca/~hamilton/courses/831/notes/confusion_matrix/confusion_matrix.html) (Vizitat pe 2012/01/23)
- El-Maleh K., Klein M., Petrucci G., și Kabal P., *Speech/music discrimination for multimedia applications*, Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '00), vol. 6, pp. 2445–2448, Istanbul, Turcia, Iunie 2000.
- Foote J. T., *A similarity measure for automatic audio classification*, in Proceedings of the AAAI Spring Symposium on Intelligent Integration and Use of Text, Image, Video, and Audio Corpora, Stanford, Calif, USA, Martie 1997.
- Howes, F. *Man Mind and Music*. Marin Secker & Warbug LTD., 1948
- ISMIR. <http://www.ismir.net/> (Vizitat pe 2012/01/23)
- Kulkarni K., Iyer D. și Sridharan S.R., *Audio Segmentation*, Stanford University, Stanford, 2007
- Liu Z., Huang J., Wang Y., și Chen I. T., *Audio feature extraction and analysis for scene classification*, Proceedings of the 1st IEEE Workshop on Multimedia Signal Processing

- (MMSP '97), pp. 343–348, Princeton, NJ, USA, Iunie 1997.
- [19] Logan B., *Mel-Frequency Cepstral Coefficients for music modeling*, ISMIR '00: International Symposium on Music Information Retrieval, 2000
- Lu L., Zhang H.-J., și Jiang H., *Content analysis for audio classification and segmentation*, IEEE Transactions on Speech and Audio Processing, vol. 10, no. 7, pp. 504–516, 2002.
- Mangaturne. <http://tagatune.org/Magnatagatune.html> (Vizitat pe 2012/01/23)
- MARSYAS. <http://marsyas.info/> (Vizitat pe 2012/01/23)
- MIREX [http://www.music-ir.org/mirex/wiki/MIREX\\_HOME](http://www.music-ir.org/mirex/wiki/MIREX_HOME) (Vizitat pe 2012/01/23)
- Munoz-Exposito J. E., Galan S. G., Reyes N. R., P. V. Candeaș și F. R. Pena, *A fuzzy rules-based speech/music discrimination approach for intelligent audio coding over the Internet*, Proceedings of the 120th Audio Engineering Society Convention (AES '06), Paris, Franța, Mai 2006, nr. 6676
- Orio N., *Music Retrieval: A Tutorial and Review*, Department of Information Engineering, University of Padova, Via Gradenigo, 6/b, Padova 35131, Italy, 2006
- Pachet F. și Cazaly D., *A taxonomy of musical genres*, RIAO '00: Content-Based Multimedia Information Access, 2000.
- Saunders J., *Real-time discrimination of broadcast speech/music*, Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '96), vol. 2, pp. 993–996, Atlanta, Ga, USA, Mai 1996.
- Scheirer E. și Slaney M., *Construction and evaluation of a robust multifeature speech/music discriminator*, Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '97), vol. 2, pp. 1331–1334, Munich, Germany, Aprilie 1997.
- Turnbull D., *Automatic music annotation*, Department of Computer Science, UC San Diego, 2005
- WEKA. <http://www.cs.waikato.ac.nz/ml/weka/> (Vizitat pe 2012/01/23)
- Williams G. și Ellis D. P. W., *Speech/music discrimination based on posterior probability features*, Proceedings of the 6th European Conference on Speech Communication and Technology (EUROSPEECH '99), pp. 687–690, Budapest, Ungaria, Septembrie 1999.
- Zhang T. și Kuo C.-C. J., *Hierarchical classification of audio data for archiving and retrieving*, Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '99), vol. 6, pp. 3001–3004, Phoenix, Ariz, USA, Martie 1999.