



# UNIVERSITATEA POLITEHNICA DIN BUCUREȘTI



**Școala Doctorală de Electronică, Telecomunicații și  
Tehnologia Informației**

**Decizie nr. 1092 din 24-07-2023**

## **REZUMAT TEZĂ DE DOCTORAT**

**Ing. Andrei-Mircea RACOVÎȚEANU**

---

**ÎNVĂȚAREA CU MARJĂ LARGĂ PENTRU ANALIZA  
IMAGINILOR**

**LARGE MARGIN LEARNING FOR IMAGE ANALYSIS**

---

### **COMISIA DE DOCTORAT**

<b>Prof. Dr. Ing. Mihai Ciuc</b> Univ. Politehnica din București	Președinte
<b>Prof. Dr. Ing. Corneliu Florea</b> Univ. Politehnica din București	Conducător de doctorat
<b>Prof. Dr. Ing. Cătălin Căleanu</b> Univ. Politehnica din Timișoara	Referent
<b>Conf. Dr. Ing. Ioan Buciu</b> Universitatea din Oradea	Referent
<b>Prof. Dr. Ing. Constantin VERTAN</b> Univ. Politehnica din București	Referent

**BUCUREȘTI 2023**

---

# Cuprins

<b>1</b>	<b>Introducere</b>	<b>1</b>
1.1	Prezentarea domeniului tezei de doctorat . . . . .	1
1.2	Scopul tezei de doctorat . . . . .	1
1.3	Conținutul tezei de doctorat . . . . .	2
<b>2</b>	<b>Rețele Convoluționale</b>	<b>3</b>
2.1	Straturile rețelelor convoluționale . . . . .	3
2.2	Procesul de învățare . . . . .	4
2.3	Arhitecturi de rețele convoluționale . . . . .	4
<b>3</b>	<b>Concepte de învățare automată</b>	<b>5</b>
3.1	Tipuri de învățare . . . . .	5
3.2	Învățare prin transfer/Adaptare de domeniu . . . . .	5
3.3	Algoritmi semi-supervizați . . . . .	6
3.3.1	Pseudo-etichetare . . . . .	6
3.3.2	Mean-Teacher . . . . .	6
3.3.3	MixMatch . . . . .	6
3.4	Metode de augmentare . . . . .	6
<b>4</b>	<b>Metode de structurare al spațiului descriptiv</b>	<b>7</b>
4.1	Funcția de cost centrică . . . . .	7
4.2	Funcția de cost insulară . . . . .	8
4.3	Funcția de cost circulară . . . . .	8
4.4	Funcția de cost marjă largă . . . . .	8
<b>5</b>	<b>Analiza expresiilor faciale</b>	<b>10</b>
5.1	Cuantificare expresiilor faciale . . . . .	10
5.2	Provocări în recunoașterea expresiilor faciale . . . . .	11
5.3	Legătura cu literatura de specialitate . . . . .	11
5.4	Soluții de detecție facială . . . . .	11
5.5	Baze de date . . . . .	12
5.5.1	Baze de date pentru expresii faciale . . . . .	12

5.5.2	Baze de date pentru detecția unităților de acțiune . . . . .	12
5.6	Recunoașterea expresiilor faciale . . . . .	13
5.6.1	Recunoașterea expresiilor faciale cu funcție de cost de marjă largă și pseudo-expresii . . . . .	13
5.6.2	Margin-Mix . . . . .	14
5.6.3	Injecția aleatoare în gradient pentru un transfer eficient în recunoașterea expresiilor faciale . . . . .	15
5.6.4	Detecția unităților de acțiune cu funcție de cost de marjă largă .	17
<b>6</b>	<b>Regăsirea de imagini similare</b>	<b>21</b>
6.1	Baza de date . . . . .	21
6.2	Legătura cu literatura de specialitate . . . . .	21
6.3	Metoda propusă . . . . .	22
<b>7</b>	<b>Concluzii</b>	<b>26</b>
7.1	Rezultate obținute . . . . .	26
7.2	Contribuții originale . . . . .	26
7.2.1	Publicații . . . . .	27
7.2.2	Dezvoltări ulterioare . . . . .	29
	<b>Bibliografie</b>	<b>30</b>

# Capitolul 1

## Introducere

### 1.1 Prezentarea domeniului tezei de doctorat

În ultimii ani, domeniul inteligenței artificiale și a vederii computerizate au cunoscut o evoluție foarte rapidă. Accesul mai facil la un volum mare de date împreună cu algoritmi de învățare automată au contribuit la dezvoltarea sistemelor care imită capacitățile umane. De progresul domeniilor amintite anterior beneficiază mai multe industrii cum ar fi: medicina, finanțele sau industria de divertisment.

Cele 2 probleme tratate în cadrul lucrării sunt recunoașterea expresiilor faciale și regăsirea imaginilor cu conținut similar. Detecția expresiilor faciale se concentrează pe recunoașterea automată a emoțiilor oamenilor. Aceasta implică dezvoltarea unor sisteme care pot interpreta și înțelege stările și reacțiile emoționale ale unei persoane bazate pe mișcările faciale. Regăsirea imaginilor cu conținut similar implică utilizarea unor tehnici de căutare și extragere a imaginilor dintr-un volum mare de date după conținutul vizual.

### 1.2 Scopul tezei de doctorat

Scopul principal al tezei este găsirea unor soluții mai eficiente pentru recunoașterea expresiilor faciale și regăsirea de imagini similare. Tehnicile propuse s-au bazat pe rețele convoluționale cât și pe introducerea unor noi funcții de cost și metode de augmentare.

Rețelele convoluționale au nevoie de cantități mari de date care de cele mai multe ori nu sunt ușor de achiziționat. Pentru a combate acest impediment, se pot folosi tehnici bazate pe învățarea semi-supervizată sau adaptarea de domeniu. Din acest motiv, soluțiile propuse depășesc categoria algoritmilor "supervizați" și au fost incluse în categoria "nesupervizată".

Prima problemă abordată este cea a expresiilor faciale. În acest caz, au fost folosite atât expresii faciale discrete cât și mișcări faciale cunoscute ca "unități de acțiune". Chiar dacă mișcările faciale sunt mai ușor de înțeles, unitățile de acțiune sunt o reprezentare mai obiectivă și mai greu de confundat. Pentru a îmbunătăți rezultatele a fost propusă o nouă

funcție de cost cu rol de grupare al spațiului descriptiv furnizat de rețeaua convoluțională. În plus, au fost testate și noi metode de augmentare și regularizare.

Cealaltă temă studiată a fost regăsirea imaginilor cu conținut similar. În această situație conținutul vizual al imaginilor poate fi foarte complex, motiv pentru care descriptorii obținuți cu o rețea convoluțională se pot confunda destul de ușor între ei. Aceeași funcție de cost a fost testată pentru a verifica dacă o organizare mai bună a spațiului descriptorilor îmbunătățește numărul de imagini similare returnate pentru o poză de referință.

### **1.3 Conținutul tezei de doctorat**

Lucrarea a fost împărțită în 7 capitole. Capitolele cu numărul 2,3 și 4 sunt prezentate conceptele teoretice care au reprezentat fundamentul rezultatelor experimentale. Acestea sunt aduse în prim plan în capitolele 5 și 6, urmând ca ultimul capitol să fie rezervat concluziilor.

Capitolul 2 redă informații despre rețelele convoluționale. Aici sunt prezentate principalele tipuri de straturi, funcțiile de cost, metodele de optimizare și principalele arhitecturi folosite. În capitolul 3 sunt aduse în discuție informații despre tipurile de învățare automată. Accentul este pus pe învățarea semi-supervizată/învățarea prin transfer și pe tehnicile de augmentare folosite. Capitolul 4 se concentrează pe funcțiile de cost care au rol de grupare al spațiului descriptiv. Sunt prezentate conceptele matematice, funcționalitatea, dar și o scurtă comparație a acestora. Capitolele 5 și 6 cuprind rezultatele obținute în cadrul tezei pentru cele 2 teme abordate, iar capitolul 7 este destinat concluziilor.

# Capitolul 2

## Rețele Convoluționale

Rețelele neuronale convoluționale sunt o clasă de modele de învățare profundă care au avut un succes deosebit în sarcini de vedere computerizată, cum ar fi clasificarea imaginilor, detectarea obiectelor sau segmentarea imaginilor. Acestea sunt concepute pentru a învăța și extrage automat caracteristici relevante din datele de intrare, făcându-le potrivite pentru sarcini care implică imagini.

### 2.1 Straturile rețelelor convoluționale

Ideea de bază din spatele rețelelor convoluționale este folosirea mai multor tipuri de straturi și procedee matematice astfel încât să fie surprinse trăsături cât mai relevante pentru datele de intrare. Primul strat relevant este evident cel convoluțional unde este aplicată efectiv operația de convoluție cu ajutorul unor filtre care conțin diferite ponderi. Aici apare și conceptul de conectivitate locală care înseamnă că nu există conexiuni între toți neuronii straturilor consecutive.

Următoarele straturi importante sunt cele de sub-eșantionare. Ele au rolul de a reduce dimensiunea hărților de trăsături produse de operațiile convoluționale. În acest mod efortul de calcul este mult mai mic, iar sistemul are o putere de generalizare mai bună. Cea mai cunoscută variantă de sub-eșantionare este *max-pooling* care ia în considerare valoarea maximă dintr-o vecinătate.

Fenomenul de supraînvățare poate apărea destul de des la rețelele convoluționale dacă nu sunt utilizate o serie de straturi care se numesc de regularizare. Un astfel de strat este cel de *dropout* [1] care implică eliminarea la întâmplare al unui anumit procent din conexiunile neuronale dintre straturi. Acestea sunt eliminate doar în procesul de antrenare, nu și în cel de testare. O altă variantă este normalizarea loturilor de date (*batch normalization*) [2] care constă în normalizarea activărilor fiecăruia strat. Astfel timpul de antrenare este redus, iar cazurile în care activările între straturi consecutive sunt complet diferite sunt eliminate. Ultimele straturi sunt cele dens conectate care sunt folosite adesea ca descriptori sau straturi decizionale.

## 2.2 Procesul de învățare

Un sistem de învățare automată folosește o serie de funcții matematice care se numesc funcții de cost pentru a măsura cât de departe se situează predicțiile acestuia de etichetele de referință. Există 2 tipuri de probleme în funcție de natura etichetelor: *clasificare* (etichete discrete) sau *regresie* (etichete continue). Pentru problemele de clasificare cea mai utilizată funcție de cost este *entropia încrucișată*, iar pentru problemele de regresie este *eroarea pătratică medie*.

Pe parcursul antrenării sistemul învață să reducă erorile printr-un algoritm de optimizare matematică (*optimizer*). Acesta ajustează ponderile după fiecare iterație în funcție de influența fiecărei ponderi la eroarea totală reprezentată de funcția de cost. Tot acest proces de ajustare al ponderilor care depinde de eroarea totală se numește propagarea înapoi a erorii (*backpropagation*).

## 2.3 Arhitecturi de rețele convoluționale

Odată ce rețelele convoluționale au devenit mai folosite s-au căutat o serie de arhitecturi standard care să poată fi folosite indiferent de natura problemei ce trebuie rezolvată. Prima arhitectură care a obținut rezultate impresionante a fost AlexNet [3]. Odată cu aceasta a fost introdusă și o nouă funcție de activare care anulează ponderile negative.

O altă familie cunoscută de arhitecturi este VGG [4] care a reușit să crească adâncimea rețelelor prin adăgarea de noi straturi. Spre deosebire de AlexNet, filtrele de convoluție au o dimensiune mai mică, iar la operațiile de sub-eșantionare ferestrele nu se mai suprapun. Totuși creșterea tot mai mare a numărului de straturi a favorizat apariția unui alt fenomen supărător care se numește "dispariția gradientilor" (*vanishing gradients*).

Arhitectura ResNet [5] a fost propusă pentru a combate fenomenul amintit mai sus. Principala componentă inovativă a fost blocul rezidual. Spre deosebire de celelalte arhitecturi amintite în arhitecturile reziduale intrarea unui strat este copiată la intrarea unor straturi superioare. În acest fel s-a putut crește din ce în ce mai mult adâncimea rețelelor fără a mai avea un impact negativ asupra performanței. După ce s-a constatat că creșterea excesivă a adâncimii nu mai aduce beneficii semnificative au fost propuse și alte variante de ResNet. Wideresnet [6] propune o creștere a hărților de trăsături mărinde astfel lățimea straturilor convoluționale.

# Capitolul 3

## Concepte de învățare automată

### 3.1 Tipuri de învățare

În general algoritmi de învățare automată pot fi împărțiți în 3 mari categorii după natura etichetelor datelor de intrare. Dacă datele conțin etichete algoritmi se numesc supervizați, altfel se numesc nesupervizați. Mai există și cea de-a treia categorie a algoritmilor semi-supervizați care folosesc de obicei o cantitate mai mică de date adnotate și o porțiune însemnată de date neadnotate. Aceștia din urmă au fost preponderent folosiți la partea experimentală.

### 3.2 Învățare prin transfer/Adaptare de domeniu

Datorită datelor insuficiente și incorect etichetate, tehnicile de învățare prin transfer și adaptare de domeniu au devenit din ce în ce mai ofertante. Învățarea prin transfer implică folosirea unui sistem pre-antrenat pentru o problemă similară sau diferită față de cea pe care a fost antrenat inițial. Cu alte cuvinte, această tehnică încearcă să crească performanțele pentru o sarcină nouă folosind caracteristici învățate anterior.

De cealaltă parte, adaptarea de domeniu presupune găsirea unei legături între două domenii care au o distribuție diferită a datelor. De obicei este antrenat un sistem pentru un domeniu inițial ca mai apoi să fie adaptat la un domeniu țintă prin diverse tehnici de ajustare cum ar fi: ponderarea funcțiilor de cost sau alinierea trăsăturilor. Rezultatul final este creșterea puterii de generalizare prin reducerea diferențelor structurale dintre cele două domenii.



## 3.3 Algoritmi semi-supervizați

### 3.3.1 Pseudo-etichetare

În momentul în care un sistem antrenat pe date care conțin etichete este aplicat pe o serie de eșantioane care nu dețin această informație se obține pseudo-etichetarea [7]. Practic, se obțin etichetele pentru setul de date neadnotat folosind un sistem pre-antrenat. Apoi, se adaugă datele cu etichetele generate la cele inițiale și se reantrenează tot ansamblul.

### 3.3.2 Mean-Teacher

Algoritmul *Mean-Teacher* [8] este un algoritm semi-supervizat care folosește 2 rețele similare. Prima este considerată "student" și are rol de clasificare. Cea de-a doua este rețeaua "profesor" care trebuie să reproducă cât mai bine ieșirea rețelei "student". Pe parcursul procesului de antrenare, unul dintre obiectivele principale este reducerea diferențelor de distribuție între rețeaua "student" și rețeaua "profesor".

### 3.3.3 MixMatch

*MixMatch* [9] este un algoritm care combină mai multe paradigme semi-supervizate pentru a maximiza performanța. Acesta folosește metoda *MixUp* [10] prin care se generează noi eșantioane și noi etichete. Se utilizează o funcție de cost care ia în calcul și componenta de informație furnizată de noile date generate.

## 3.4 Metode de augmentare

Ținând cont că rețelele convoluționale au nevoie de cantități mari de date s-au căutat variante care să rezolve această problemă. Metodele de augmentare se folosesc pentru a mări artificial numărul de eșantioane folosite la antrenare.

Cele mai uzuale metode de augmentare sunt legate de distribuția spațială a elementelor dintr-o imagine. Astfel se pot genera noi eșantioane prin operații de rotire, oglindire sau redimensionare. Se poate acționa și la nivel de pixel prin tehnici de filtrare sau contrastare.

Tehnica de augmentare *MixUp* [10] ajută la generarea de date noi printr-o combinație liniară a unei perechi aleatoare de eșantioane existente în baza de date folosită la antrenare. Pentru aceste date se creează etichete folosind aceeași metodă. Totuși, fiind dificil de interpretat ca valoare unica combinația liniară a 2 etichete discrete s-a folosit distribuția probabilistică obținută.

# Capitolul 4

## Metode de structurare al spațiului descriptiv

Pentru probleme complexe este posibil ca spațiul descriptiv furnizat de rețelele convoluționale să conțină exemple care se suprapun destul de mult. Din acest motiv, s-au realizat cercetări în direcția unor funcții de cost noi care să organizeze mai eficient caracteristicile obținute.

### 4.1 Funcția de cost centrică

Una din cele mai cunoscute astfel de funcții este cea centrică [11]. Principiul după care se ghidează este minimizarea distanțelor între punctele care sunt încadrate în aceeași categorie. Ecuația 4.1 descrie matematic funcția de cost centrică, unde  $e_i$  este descriptorul obținut de pe stratul dens conectat înainte de cel decizional, iar  $c_i$  este centroidul asociat eșantionului curent.

$$L_C = \frac{1}{2} \sum_{i=1}^N \|e_i - c_i\|_2 \quad (4.1)$$

Pentru a obține un spațiu mai aerisit, poziția centroizilor este recalculată după fiecare iterație conform 4.2.  $\Delta c_k^i$  este media datelor care aparțin clasei  $i$  în lotul de date curent, iar  $\alpha$  este un parametru subunitar care temperează o posibilă influență negativă la noua poziție a centroizilor a unor eșantioane etichetate eronat. Deoarece nu are rol de decizie/clasificare, funcția centrică se poate folosi doar cu o funcție de cost de acest tip.

$$c_k^{i+1} = c_k^i - \alpha \Delta c_k^i \quad (4.2)$$

## 4.2 Funcția de cost insulară

Spre deosebire de funcția de cost amintită în secțiunea precedentă, funcția insulară [12] aduce un nou plus. Pe lângă minimizarea varianței în aceeași clasă, aceasta încearcă să maximizeze distanța între centroizii asociați fiecărei etichete. Ecuația 4.3 descrie formula matematică, unde  $L_C$  este funcția de cost centrică, iar termenul al doilea cumulează distanțele unghiulare între centroidul curent și toate celelalte clase. La fel ca în cazul precedent se folosește cu o funcție de cost de clasificare.

$$L_{IL} = L_C + \lambda_1 \sum_{c_i \in C} \sum_{c_j \in C, c_i \neq c_j} \left( \frac{c_i \cdot c_j}{\|c_i\|_2 \|c_j\|_2} + 1 \right) \quad (4.3)$$

## 4.3 Funcția de cost circulară

Funcția de cost circulară [13] implică normalizarea eficientă a descriptorilor astfel încât aceștia să poată fi interpretați din punct de vedere geometric ca un cerc. În ecuația 4.4  $F_{x_i}$  este descriptorul asociat penultimului strat dens conectat, iar  $R$  este norma care este vizată pentru fiecare eșantion. Din punct de vedere geometric este asociată cu raza unui cerc.

$$L_R = \frac{\lambda}{2m} \sum_{i=1}^m (\|F_{x_i}\|_2 - R)^2 \quad (4.4)$$

## 4.4 Funcția de cost marjă largă

Funcția de cost centrică ia în calcul doar distanța în interiorul claselor, nu și distanța între punctele care se situează în clase distincte. Funcția insulară combate acest neajuns, dar bazându-se pe distanța unghiulară are o serie de limitări. Dacă unghiul dintre 2 grupuri de eșantioane este foarte mic acestea nu se pot separa suficient. Funcția circulară impune o reprezentare geometrică a descriptorilor sub forma unui cerc maximizând astfel distanțele unghiulare. Totuși, dacă 2 clase se suprapun în reprezentarea inițială există șanse mari ca acestea să coincidă și în reprezentarea circulară.

Luând în considerare limitările amintite în paragraful anterior, a fost propusă o nouă funcție de cost care să satisfacă nevoile actuale. Aceasta a fost intitulată funcția de marjă largă [14] și impune folosirea unei distanțe euclidiene între fiecare eșantion și centroizii asociați celorlalte clase. În acest mod sunt evitate limitările impuse de funcțiile de cost anterior prezentate.

Ecuația 4.5 redă definiția matematică a acestei funcții.  $e_i$  este descriptorul reprezentat de penultimul strat dens conectat,  $c_j$  reprezintă centroidul asociat descriptorului  $e_i$ , iar  $c_k$  reprezintă toți centroizii diferiți de centroidul lui  $e_i$ .

$$\mathbf{L}_{\text{LM}} = \sum_{i=1}^N \left( \left\| \frac{\mathbf{e}_i}{\|\mathbf{e}_i\|_2} - \frac{\mathbf{c}_j}{\|\mathbf{c}_j\|_2} \right\|_2 - \frac{1}{C-1} \sum_{k=1, k \neq j}^C \left\| \frac{\mathbf{e}_i}{\|\mathbf{e}_i\|_2} - \frac{\mathbf{c}_k}{\|\mathbf{c}_k\|_2} \right\|_2 \right) \quad (4.5)$$

În Figura 4.1 este prezentată pe scurt funcționalitatea funcției de marjă largă. Având la dispoziție 3 eșantioane ( $X_A; X_B; X_C$ ) cu etichetele A,B și C, acestea vor fi aduse mai aproape de centrozii claselor de care aparțin (minimizarea distanțelor marcate cu săgeți continue) și vor fi depărtate de centrozii celorlalte clase (maximizarea distanțelor marcate cu linie punctată).

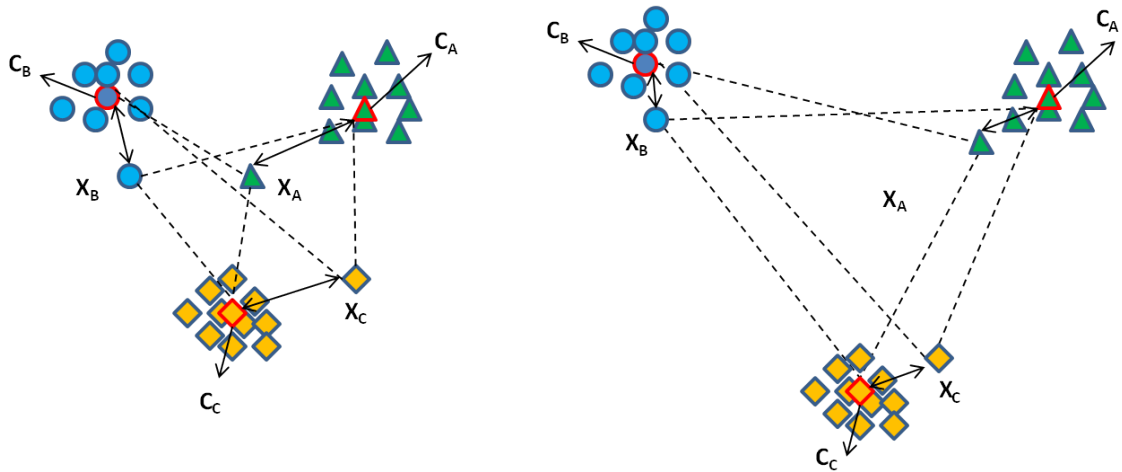


Figura 4.1 Funcționalitatea costului de marjă largă. Stânga: Înainte de marja largă. Dreapta: După marja largă

# Capitolul 5

## Analiza expresiilor faciale

Acest capitol este cel mai cuprinzător din întreaga teză și reprezintă efortul mai multor ani de cercetare. Printre sarcinile abordate se regăsesc recunoașterea de expresii discrete, dar și detecția de mișcări faciale. Pentru îmbunătățirea performanței s-au utilizat funcții noi de cost și metode inovative de augmentare și regularizare.

### 5.1 Cuantificare expresiilor faciale

Recunoașterea expresiilor faciale nu este o sarcină ușoară nici pentru oameni. Aceasta este cu atât mai dificilă pentru un sistem de învățare automată. De-a lungul timpului s-au căutat mai multe modele care să definească mai convenabil această problemă.

Primul model foarte folosit în acest domeniu este cel al expresiilor discrete definit de Ekman [15]. Acesta are la bază 6 expresii fundamentale : frică, dezgust, fericire, furie, surprindere, tristețe. De obicei, la acestea se mai adaugă și expresia neutră. Deși această variantă este simplu de folosit, s-au căutat și alte modele care să definească emoțiile într-o manieră mai obiectivă. Astfel, a fost introdus *sistemul de codare al unităților de acțiune* care presupune că fiecare expresie facială este formată dintr-o serie de activări ale unor mușchi faciali specifici (unitățile de acțiune).

Sistemul cuprinde un număr de 43 de unități de acțiune care au fost divizate în funcție de poziția la nivelul feței. Există mișcări faciale asociate celor mai importante elemente faciale cum ar fi: ochii, gura sau sprâncenele. Deși este un model care lasă mai puțin loc de interpretări, este necesar un personal destul de calificat care să recunoască cele mai fine mișcări faciale. Mai există și modele mai complexe care iau în calcul intensitatea expresiei sau dacă este pozitivă sau negativă. Un asemenea model a fost propus de Russel [16] și ia în considerare emoțiile compuse.

## 5.2 Provocări în recunoașterea expresiilor faciale

După cum s-a demonstrat și în secțiunea 5.1 problema expresiilor faciale este una foarte dificil de rezolvat. Mai mult, există și o multitudine de provocări care pot afecta performanța de recunoaștere. Un prim motiv este volumul de date inconsistent. Chiar dacă poate că este facil să obții multe imagini cu expresii faciale, procesul de etichetare poate fi lung și costisitor. Nu este o misiune ușoară să fie găsiți adnotatori umani profesioniști. Din această cauză există încă multe imagini cu etichete greșite.

Un alt motiv este timpul scurt de apariție al expresiilor faciale. De cele mai multe ori mișcăările faciale care alcătuiesc emoțiile sunt foarte fine chiar și pentru computer. De asemenea, există mai multe expresii care sunt ușor de confundat între ele cum ar fi frica și surprinderea. Ambele emoții implică deschiderea gurii și ridicarea sprâncenelor.

## 5.3 Legătura cu literatura de specialitate

Fiind o sarcină greu de rezolvat, recunoașterea expresiilor faciale a atras în ultima perioadă un număr mare de cercetători, motiv pentru care articolele pe această temă au crescut considerabil. În lucrări precum [17, 18] a fost abordată recunoașterea expresiilor fundamentale (secțiunea 5.1) cu rețele convoluționale .

Du [19] și Zhang [20] au observat inconsistența datelor de antrenare și au migrat spre soluții din zona învățării semi-supervizate. A fost abordată și recunoașterea unităților de acțiune. Corneanu [21], Zhao [22] sau Benitez [23] împreună cu colaboratorii sunt doar o parte din cei care au detectat mișcăările faciale și intensitatea acestora. Adaptarea de domeniu și învățarea prin transfer au fost abordate în [24–26].

## 5.4 Soluții de detecție facială

Cele mai multe seturi de date conțin imagini cu fețe sau expresii faciale care nu sunt încadrate convenabil. Informația legată de fundal care se regăsește în majoritatea imaginilor este redundantă și nu este folositoare rețelelor convoluționale. În plus, prin decuparea fețelor din pozele inițiale se micșorează dimensionalitatea care contribuie la un timp mai redus de antrenare.

Din multitudinea de soluții existente se remarcă algoritmi Viola-Jones [27] și MT-CNN [28]. Metoda lui Viola-Jones folosește o serie de trăsături la diferite scale pentru a surprinde fețele din imagini oricare ar fi dimensiunea acestora. Se mai folosește și imaginea integrală pentru accelerarea procesului de detecție care a permis identificarea facială în timp real.

De cealaltă parte MTCNN [28] este o tehnică bazată pe rețelele convoluționale. Aceasta poate fi organizat în 3 etape. Se detectează mai întâi fețele, se potrivesc ferestrele

de încadrare și la final se recunosc elementele faciale necesare pentru o aliniere cât mai corectă a feței.

## 5.5 Baze de date

Pentru a face față cerințelor impuse de metodele semi-supervizate testate în partea experimentală au fost folosite mai multe seturi de date diverse. S-a abordat problema recunoașterii de expresii fundamentale, dar și a detecției de mișcări faciale (unități de acțiune). Mai multe informații despre fiecare în parte se regăsesc în secțiunile ce urmează.

### 5.5.1 Baze de date pentru expresii faciale

**FER/FER+.** FER2013 [29] și FER+ [30] reprezintă 2 seturi de date care oferă imagini cu expresii faciale fundamentale. FER+ este o extindere a FER2013 în care s-au corectat o serie de etichete care au fost atribuite incorect în prima versiune. Dispune de un număr de aproximativ 35000 de imagini achiziționate în condiții naturale.

**Megaface.** Megaface [31] este o bază de date mult mai cuprinzătoare care conține aproximativ 1 milion de imagini achiziționate în condiții naturale. Nu dispune de etichete, motiv pentru care a fost folosită ca porțiune neetichetată de date pentru experimentele semi-supervizate.

**RAF-DB.** RAF-DB [32] seamănă cu FER deoarece conține expresii discrete. Cu toate acestea, dispune de un număr mai mic de fotografii obținute în condiții de laborator. Anotarea imaginilor a fost realizată de un număr de 40 de persoane.

**Facial expression in children.** În cadrul acestei lucrări a fost studiată și recunoașterea expresiilor faciale la copii. CAFE [33] este unul din cele mai cunoscute seturi de date pe această temă. LRIS [34] este un alt set care compensează limitările CAFE prin creșterea numărului de imagini și diversificarea emoțiilor.

### 5.5.2 Baze de date pentru detecția unităților de acțiune

**CK+.** CK+ [35] este un set de date care conține atât etichete discrete cât și unități de acțiune. Este prezentat în această secțiune pentru că a fost folosit doar pentru recunoașterea de unități de acțiune. Imaginile sunt organizate pe subiecți și fiecare secvență conține cadre începând de la expresia neutră până la intensitatea maximă a expresiei. Ultimele poze din fiecare secvență sunt anotate și la nivel de unități de acțiune.

**Emotionet.** Emotionet [23] este un set de date care conține aproximativ 1 milion de imagini. Spre deosebire de Megaface [31], dispune de 50000 de imagini etichetate la nivel de unități de acțiune. Imaginile sunt achiziționate în condiții naturale. Etichetele pentru mișcările faciale sunt binare.

**DISFA.** În comparație cu Emotionet [23], DISFA [36] oferă adnotări și pentru intensitatea unității de acțiune. Astfel, etichetele sunt cuprinse între 0 (unitatea de acțiune nu este activă) și 5 (unitatea de acțiune are intensitatea maximă). Conține 130000 de imagini captate în condiții de laborator împărțite în 27 de subiecți.

## **5.6 Recunoașterea expresiilor faciale**

Acest capitol s-a ocupat de îmbunătățirea rezultatelor pentru detecția expresiilor faciale folosind mai multe instrumente noi sau deja cunoscute în literatură. Printre acestea se număra învățarea semi-supervizată/învățarea prin transfer, funcțiile de cost cu rol de structurare al spațiului descriptorilor (centrică - Secțiunea 4.1, circulară - Secțiunea 4.3, marjă largă - Secțiunea 4.4), dar și metode de augmentare precum MixUp (Secțiunea 3.4).

### **5.6.1 Recunoașterea expresiilor faciale cu funcție de cost de marjă largă și pseudo-expresii**

Având ca motivație succesul funcțiilor de cost centrice și insulare, în această secțiune a fost explorat potențialul funcției de marjă largă care a fost propusă (Secțiunea 4.4).

Pentru acest scenariu a fost folosită o arhitectură Alexnet (Secțiunea 2.3) care a fost antrenată pentru detecția simultană a expresiilor faciale discrete și a unităților de acțiune. La nivel de arhitectură au fost folosite 2 straturi de decizie conectate în paralel din penultimul strat dens conectat.

S-a încercat folosirea la antrenare a mai multor seturi de date care nu conțin același tip de etichetă. De aceea a fost necesară o adaptare de domeniu plecând de la unități de acțiune la expresii faciale discrete. Legătura dintre cele 2 categorii s-a realizat cu un set de ecuații care descriu expresiile fundamentale ca fiind o sumă de unități de acțiune active simultan. Rezultatul a constat în obținerea unor pseudo-expresii.

A fost nevoie de aceste pseudo-expresii, deoarece funcția de cost de marjă largă are nevoie de conceptul de centroizi (clase discrete). Detecția de unități de acțiune este o problemă multi-clasă deoarece pot fi active mai multe mișcări faciale în același timp. Chiar dacă etichetele corespunzătoare unităților de acțiune sunt discrete, faptul amintit mai devreme poate conduce la un număr excesiv de clase cu reprezentare redusă.

Rezultatele pentru recunoașterea expresiilor faciale pot fi urmărite în tabelul 5.1, în timp ce performanța pentru unitățile de acțiune este redată în tabelul 5.2. Se poate observa că funcția de cost de marjă largă obține rezultate mai bune decât funcția centrică și insulară.



Metodă	Învățare	Med. Ac.	Ac.
AlexNet - [32]	Superv	55.60	68.90
AlexNet + Feat.Sel.Net [17]	Superv	72.46	81.10
AlexNet + Island loss [12]	Superv	57.1	75.08
AlexNet + Center loss [11]	SSL	63.15	78.81
AlexNet + Island loss [12]	SSL	64.53	78.81
AlexNet + LM loss [14]	SSL	67.26	79.85

Tabel 5.1 Acuratețea obținută pentru setul de date RAF-DB. Tabel din [14]

Metodă	Înv	AU1	AU2	AU4	AU5	AU6	AU9	AU12	AU17	AU20	AU25	AU26	AU43	Med. redus	Med. total
AlexNet [3]	Sv	24.2	n/a	34.7	<b>39.5</b>	73.1	n/a	86.8	n/a	n/a	88.5	45.6	n/a	56.1	n/a
AlexNet cen. loss [11]	Sv	<b>34.4</b>	30.3	55.3	33.3	69.10	<b>46.1</b>	79.3	27.8	<b>32.3</b>	84.4	43.2	48.8	57.9	48.8
AlexNet +WSC [22]	SSL	25.3	n/a	34.5	39.3	<b>75.6</b>	n/a	<b>87.4</b>	n/a	n/a	<b>88.8</b>	47.4	n/a	57.0	n/a
AlexNet + Isl. loss [12]	T.	30.4	29.5	<b>56.7</b>	30.6	66.7	44.1	77.3	26.7	23.8	83.9	47.3	43.9	56.14	46.7
AlexNet + LM loss [14]	T.	34.1	<b>31.1</b>	56.6	33.9	71.0	45.1	78.1	<b>30.9</b>	25.3	83.8	<b>50.9</b>	47.2	58.33	49.0

Tabel 5.2 Scorul F1 (%) pentru detecția de unități de acțiune pentru baza de date Emotionet. Învățarea este fie supervizată, semi-supervizată sau prin transfer. “Med. redusă” este media pentru setul redus de AUs:1,4,5,6,12,25,26, Med. total este media pentru setul întreg. Tabel din [14]

## 5.6.2 Margin-Mix

Algoritmul *Margin-Mix* [37] combină funcția de cost de marjă largă cu o serie de tehnici de augmentare cum ar fi MixUp [10]. Așa cum a fost amintit și în secțiunea 3.4 tehnica MixUp este aplicată pentru învățarea pur supervizată. Deși combinația liniară între 2 imagini din setul de date este posibilă, rezultatul nu va putea fi atribuit unei clase în lipsa unor etichete inițiale.

Aici intervine conceptul de marjă largă care este folosit pentru etichetarea noilor exemple formate cu MixUp folosind descriptorii formați de o rețea convoluțională. Pentru a minimiza efectul suprapunerii datelor în spațiul descriptiv atribuirea etichetelor s-a realizat folosind o tehnică Fuzzy [38]. În acest fel unui eșantion i-a fost atribuită o distribuție de probabilitate pentru fiecare clasă, nu doar o etichetă unică. Mai departe, aceste noi eșantioane împreună cu etichetele precise vor putea fi folosite mai departe pe parcursul antrenării.

De-a lungul experimentelor a fost folosită o arhitectura Wide-ResNet [6]. Performanța obținută pentru seturile de date standard: CIFAR10, STL-10, SVHN se regăsește în Tabelul 5.3. Rezultatele pentru RAF-DB [32] au fost organizate în Tabelul 5.4. În

primul tabel rezultatele sunt asemănătoare cu cele din literatura de specialitate. Mai sugestive sunt rezultatele din cel de-al doilea tabel. Se remarcă o performanță considerabil mai bună pentru *Margin-Mix* în comparație cu alte metode pur supervizate. Este important de semnalat că pe măsură ce cantitatea de date etichetate este tot mai mică metoda propusă are cele mai bune rezultate. Totuși, atunci când este folosit tot setul de date (ultima coloană din tabelul 5.4) numerele sunt similare, ceea ce demonstrează într-adevăr că *Margin-Mix* este o soluție de luat în calcul atunci când setul de date nu conține suficiente adnotări.

Metodă/Etichete	SVHN		STL	
	1000	4000	1000	5000
Supervizat [6]	–	12.84	–	–
Model Π [39]	8.06	5.57	17.41	39.19
VAT [40]	5.63	18.68	11.05	–
MeanTeacher [8]	5.65	3.39	10.36	–
ICT [41, 6]	3.53	–	7.66	–
MixMatch [9]	3.27	2.89	10.18	5.59
<b>MarginMix [37]</b>	<b>3.35</b>	<b>8.33</b>	<b>9.85</b>	<b>5.80</b>

Tabel 5.3 Eroarea obținută pentru seturile de date STL și SVHN cu arhitectura Wide - ResNet.[37]

Metodă/Etichete	320	400	1000	4000	
Supervizat	nc	26.75	35.25	55.66	85.58
Supervizat [32]	–	–	–	–	84.13
MeanTeacher [8]	nc	28.83	36.53	60.36	–
MixMatch [9]	35.60	42.25	60.37	65.24	–
<b>MarginMix [37]</b>	<b>40.55</b>	<b>45.75</b>	<b>66.47</b>	<b>70.68</b>	<b>85.36</b>

Tabel 5.4 Acuratețea obținută pentru RAF-DB cu arhitectura Wide - ResNet. Linia de sus expune numărul de exemple etichetate luate în considerare. nc - nu s-a atins convergența [37]

### 5.6.3 Injecția aleatoare în gradient pentru un transfer eficient în recunoașterea expresiilor faciale

Tehnica propusă în [42, 43] are în componență învățarea semi-supervizată și tehnica de etichetare amintită în secțiunea 3.3.1 (Pseudo-etichetarea). Pseudo-etichetarea este simplu de utilizat, ea constând în folosirea unui sistem preantrenat pe o problemă supervizată pentru a eticheta o serie de eșantioane fără etichete.

Cu toate acestea, Pseudo-etichetarea pleacă de la premisa că datele etichetate, cât și cele neetichetate au o distribuție similară. Această idee nu este adevărată de cele mai

multe ori și poate avea cauze nefaste asupra performanței. În Figura 5.1 este expusă schematic modalitatea de funcționare a metodei discutate. Pentru cazul supervizat (a) lucrurile sunt clare, granița de separare fiind raportată la datele existente. Atunci când apar și date neadnotate (cercurile albastre) Pseudo-etichetarea obligă sistemul să devină prea sigur pe etichetele furnizate, chiar dacă ele nu sunt corecte. În acest caz, granița poate arăta ca la un sistem predispus supraînvățării. Dacă se folosește injecția aleatoare de gradient (IAG) se induce indirect un grad aleator de incertitudine la predicție, motiv pentru care granița riscă mai puțin să fie focusată prea mult asupra punctelor incerte.

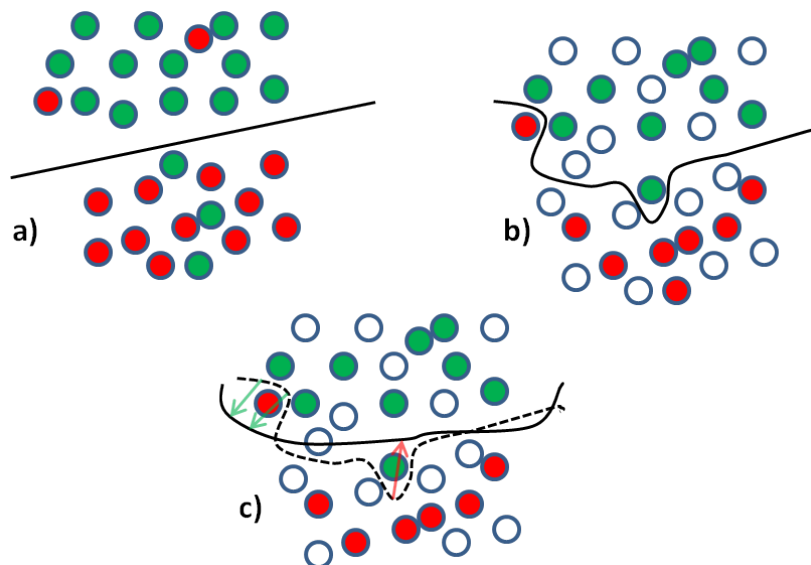


Figura 5.1 Granița de separație pentru a) Separare ideală supervizată b) Pseudo-etichetare - semi-supervizată c) Injecție aleatoare a gradientului (IAG) - Învățare prin transfer

Funcția matematică prin care s-au alterat gradientii este definită ca:

$$g(n, \lambda) = \begin{cases} \frac{\lambda n}{50}, & n < 50 \\ 0, & n \geq 50 \end{cases} \quad (5.1)$$

unde  $g$  este o variabilă aleatoare  $g : \{1, N_{epochs}\} \times [-1, 1] \rightarrow [-1, 1]$ ,  $\lambda$  este o valoare aleatoare subunitară din intervalul  $[-1; 1]$ , iar  $n$  este numărul de epoci. Această cantitate este adunată la costul curent, iar ponderile se modifică doar dacă câștigul de performanță este suficient de mare. În această manieră s-a combătut potențialul negativ care poate fi generat de distribuțiile diferite între seturile de date.

Rezultatele pentru baza de date RAF-DB sunt expuse în tabelul 5.5. Se observă că numerele obținute cu injecția aleatoare în gradient sunt preponderent mai mari decât metodele alese spre comparație. Merită semnalat că performanța este cu aproximativ 2-3% mai bună decât pseudo-etichetarea, ceea ce arată un transfer mai eficient de informație.

S-au efectuat o serie de teste și pe baza de date LRIS [34] care conține imagini cu expresii faciale la copii. Performanța este redată în tabelul 5.6. Și în această situație rezultatele au fost mult mai bune raportat la varianta supervizată aleasă ca reper.

Metodă / Metrică		Med. Ac.	Ac.
SUPERVIZAT	AlexNet [32]	55.60	68.90
	VGG-16 [32]	58.22	70.53
	DLP-CNN [32]	74.20	84.13
	ResNet-18 [44]	–	80.00
	RST [45]	72.46	81.10
	gCNN [46] - VGG16	–	85.07
	ensCNN [47]	75.14	86.31
TRANSFER	AlexNet + PE	69.5	78.5
	AlexNet + TE [42]	72.3	81.50
	AlexNet + IAG [43]	72.5	82.1
	VGG-16 + PE	74.6	83.25
	VGG-16 + TE [42]	76.50	84.5
	VGG-16 + IAG - [43]	76.82	85.15
	ResNet-50 + PE	77.12	84.8
	ResNet-50 + IAG - [43]	<b>78.22</b>	<b>86.67</b>

Tabel 5.5 Performanța pe setul RAF-DB. FSN - feature selection network, FSM - frame-to-sequence method, PE - Pseudo-etichetare, TE - Transfer etichete, RST - Rețea selecție trăsături, Med.Ac. - Acuratețe medie, Ac. - Acuratețe. Metoda propusă este notată cu IAG. Cu bold sunt marcate rezultatele cele mai bune. Tabel din [42]

Metodă	Acuratețe
VGG-16 [34] - supervizat	67.2
VGG-16 - +IAG [43]	68.5
ResNet-50 - supervizat	72.3
ResNet-50 + IAG [43]	76.6

Tabel 5.6 Performanța pentru setul de date LRIS care conține expresii faciale la copii. Tabel din [43]

#### 5.6.4 Detectia unităților de acțiune cu funcție de cost de marjă largă

Informațiile din secțiunea curentă folosesc principiile de marjă largă, pseudo-expresii și adaptare de domeniu exact ca în 5.6.1. Totuși, experimentele și scenariile testate au fost extinse. Aici, eforturile s-au concentrat doar pe detectia unităților de acțiune [48].

Arhitecturile folosite în acest caz au fost preponderent cele din familia ResNet [5]. La fel ca în secțiunea 5.6.1 s-a realizat detectia de unități de acțiune folosind și informația furnizată de stratul decizional asociat pseudo-expresiilor. Astfel au fost

folosite 3 funcții de cost diferite asociate pentru fiecare problemă în parte : entropia încrucișată binară pentru predicția mișcărilor faciale, entropia încrucișată pentru predicția pseudo-expresiilor și funcția de marjă largă pentru gruparea descriptorilor. Funcția de cost totală este redată în ecuația 5.2. Constantele  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$  sunt folosite pentru echilibrul numeric dintre cei 3 termeni ai costului final.

$$L_T = \lambda_1 L_{BCE} + \lambda_2 L_{SE} + \lambda_3 L_{LM} \quad (5.2)$$

Una dintre problemele apărute în bazele de date folosite la experimente a fost faptul că anumite unități de acțiune au o frecvență de apariție mult mai redusă în comparație cu altele. Acest aspect contribuie la o detecție mult mai redusă a acestora, deoarece la nivel descriptiv unitățile de acțiune care apar mai puțin se vor suprapune cu celelalte așa cum se poate remarca în Figura 5.2. În figură sunt reprezentate în partea stângă pseudo-expresiile reprezentate de unitățile de acțiune care apar mai rar, iar în partea dreaptă pseudo-expresiile reprezentate de toate unitățile de acțiune existente.

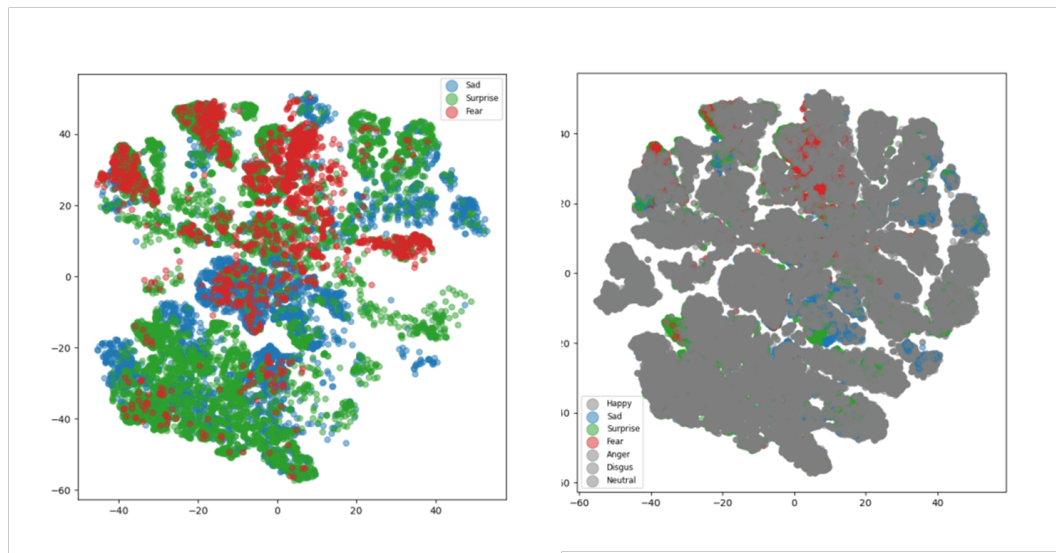


Figura 5.2 Reprezentarea la nivel descriptiv a pseudo=expresiilor. Stânga – Descriptorii asociați pseudo-expresiilor reprezentate de unitățile de acțiune care apar mai rar. Dreapta - Descriptorii asociați pseudo-expresiilor reprezentate de toate unitățile de acțiune existente. Figura din [48]

În tabelele 5.7 și 5.8 pot fi urmărite rezultatele obținute pe setul de date DISFA [36] și Emotionet [23] în comparație cu alte tehnici similare din literatura de specialitate. Pentru setul DISFA se poate observa că rezultatele obținute cu funcția de cost de marjă largă pentru unitățile de acțiune cu frecvență mai mică de apariție surclasează toate celelalte tehnici din literatura de specialitate chiar dacă media totală nu este neaparat cea mai bună. Pe setul de date Emotionet, diferențele se păstrează.

Un posibil motiv pentru care performanța este mai crescută în cazul unităților de acțiune care apar mai puțin se datorează conceptului de marjă largă care reușește să

separe mai bine la nivel descriptiv pseudo-expresiile care sunt reprezentate de unitățile de acțiune mai rare. Acest lucru se poate vedea în Figura 5.3 unde expresiile formate din mișcări faciale care apar mai rar (punctele colorate) sunt mult mai compacte și mai puțin suprapuse cu celelalte (punctele gri - unități de acțiune care apar mai des) atunci cand este utilizat conceptul de marjă largă.

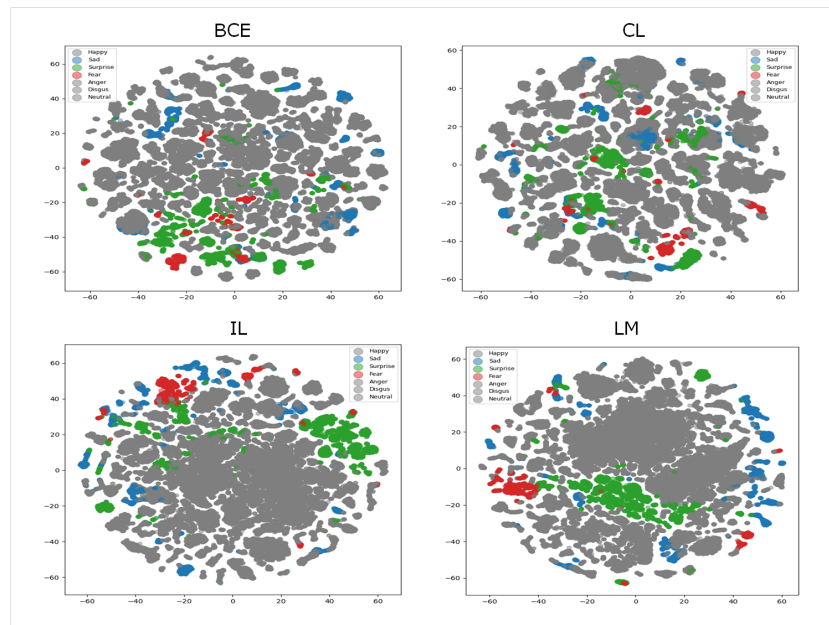


Figura 5.3 Evoluția spațiului descriptiv la epoca 40 pentru toate pseudo-expresiile. BCE (Entropie Binară Încrucișată) , CL(Funcție centrică, IL (Funcție insulară) , LM (Marjă largă) Figură din [48]

Tabel 5.7 Comparație cu literatura de specialitate pentru DISFA. AD reprezintă adaptare de comeniu. ML reprezintă marjă largă. F1-tot este media tuturor unităților de acțiune. F1-rare este media unităților de acțiune care apar mai rar .Cele mai bune rezultate sunt marcate cu bold. Rezultatele obținute cu ML sunt colorate cu gri. Tabel din [48]

Metodă	F1-tot	F1-rare
AD (DISFA neetichetat) PreActRes18 - - preantrenat Funcție de marjă largă	55.4	<b>51.2</b>
AD (DISFA neetichetat) PreActRes18 - - fără preantrenare Funcție de marjă largă	50.4	44.7
AD (DISFA neetichetat) PreActRes18 Funcție centrică [11]	46.6	38.9
AD (DISFA neetichetat) PreActRes18 Funcție insulară [12]	47.3	40.0
AD (DISFA unlabeled) PreActRes18 Funcție circulară [13]	46.0	39.4
DRML [49]	26.7	14.2
ROI [50]	48.4	23
DSIN [21]	53.6	42.6
JAA [51]	56.0	48.2
SRERL [52]	55.9	45.6
MLT-RM [53]	<b>60.1</b>	46.6
UGN-B [54]	60.0	49.65

Tabel 5.8 Rezultate pentru Emotionet. Comparație cu literatura de specialitate. F1-tot este media tuturor unităților de acțiune. F1-rare este media unităților de acțiune care apar mai rar. PsExpr = Pseudo-expresii SV = Supervizat ; AD = Adaptare Domeniu; EBI = Entropie Binară încrucișată; ML = Marjă largă; FC = Funcție centrică; FI = Funcție insulară; FCR = Funcție circulară. Rezultatele cele mai bune sunt marcate cu bold. Cele mai bune rezultate obținute cu ML sunt colorate cu gri. Tabel din [48]

Metodă	F1-tot	F1-rare
SV-BCE +PrActRes18	47.08	35.19
SV-BCE +PrActRes18 - CE(PSEExpr)	48.48	36.75
SV-NM	50.16	38.90
AD - ML -Alexnet	49.04	35.96
AD - ML +PrActRes18	52.12	40.74
AD - ML +PrActRes18 Imagenet preantrenat	54.31	43.25
AD - ML+PrActRes18 DISFA preantrenat	<b>55.89</b>	<b>45.58</b>
AD - FC[11] +PrActRes18	48.92	36.83
AD - FCR[13] +PrActRes18	49.14	36.34
AD - FI[12] +PrActRes18	50.38	38.20

# Capitolul 6

## Regăsirea de imagini similare

Regăsirea de imagini similare este o problemă tot mai întâlnită în diverse domenii cum ar fi: motoarele de căutare de imagini sau aplicații medicale. De-a lungul timpului au fost folosiți o multitudine de descriptori care au avut rolul să descrie informația vizuală din imagini cât mai eficient. În acest capitol s-a studiat posibilitatea folosirii funcției de cost de marjă largă pentru a obține o serie de descriptori mai buni folosind rețele convoluționale.

### 6.1 Baza de date

Baza de date utilizată în experimente se numește Places365 [55] și este alcătuită din aproximativ 1.8 milioane de imagini distribuite în 365 de clase. Imaginile conțin scene diverse. Setul de antrenare dispune de un număr situat între 3000 și 5000 de imagini. Setul de testare are 900 de poze pentru fiecare clasă.

### 6.2 Legătura cu literatura de specialitate

Înainte ca rețelele convoluționale să devină populare, pentru problema regăsirii de imagini similare s-au căutat variante foarte diverse de descriptori. Printre cei apăruiți înainte de perioada cu învățare adâncă se regăsesc modelele binare locale [56], histograma gradientilor orientați [57] sau histograme de culoare. Apoi, s-a trecut la descriptori care detectează puncte similare într-o imagine cum ar fi: SIFT [58] sau SURF [59].

În ultimii ani au prins contur și descriptori care evidențiază mai bine informația vizuală cum ar fi histograma de cuvinte vizuale [60]. Evident că odată cu dezvoltarea rețelelor convoluționale tot mai multă lume a dorit să folosească straturile dens conectate pe post de descriptori [61, 62].



Scenariu/Metrică [%]	Prec.med.-5-interogări	Prec.med.-8-interogări	Er. Top-1	Er. Top-5	Acc.	ASG PR
<b>EI-pre-antrenat-reper</b>	29.93	28.35	66.74	37.53	53.10	9.17
<b>ML-fine-tunat</b>	29.89	28.56	66.99	37.95	54.69	9.23
<b>ML-antrenare de la zero</b>	<b>31.35</b>	<b>29.98</b>	66.18	37.79	53.51	11.18
<b>Resnet-152 [19]</b>	-	-	-	-	54.74	-

Tabel 6.1 Rezultate experimentale Places 365 (Prec.med.- precizie medie, Er. - Eroare, EI- entropie încrucișată , Acc. - acuratețe, ASG PR - aria de sub graficul precizie-reamintire, ML- marjă largă).

Tabel din [63]

### 6.3 Metoda propusă

Pentru problema de regăsire de imagini similare a fost folosită o rețea convoluțională ResNet care a fost antrenată pe setul de date Places 365 [55]. După această etapă s-a extras penultimul strat dens conectat și s-a folosit pe post de descriptor la partea de regăsire.

Rețeaua a fost antrenată cu conceptul de marjă largă (Secțiunea 4.4) pentru a crește densitatea spațială a spațiului descriptiv. Rezultatele sunt în tabelul 6.1. Ca metrici de performanță s-au folosit precizia medie, eroarea și aria de sub graficul curbei precizie-reamintire. Se poate observa că rezultatele sunt doar ușor în favoarea metodei prezentate.

Pentru a stabili utilitatea funcției de marjă largă s-au creat 3 noi scenarii care să conțină la nivel de descriptori date separabile, date neseperabile și date din ambele categorii. O figură sugestivă cu spațiul datelor înainte de antrenare se regăsește mai jos.

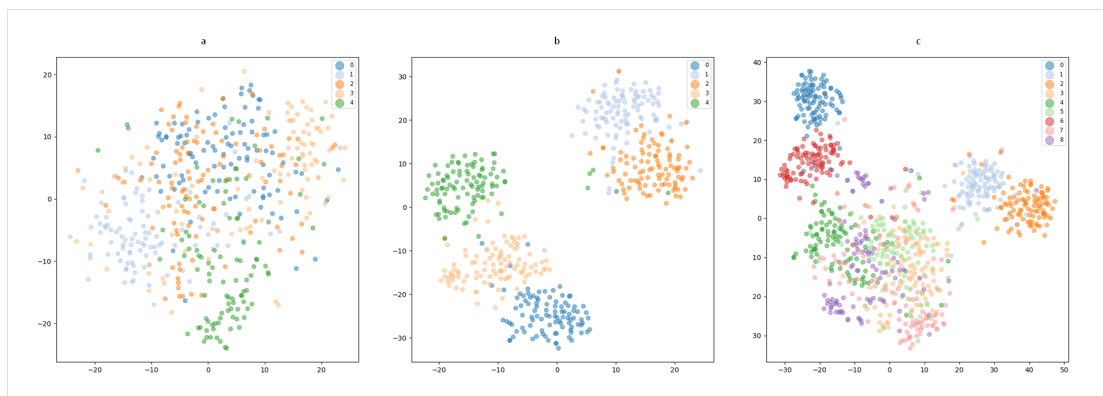


Figura 6.1 Reprezentarea spațiului descriptiv pentru cele 3 scenarii. a – date neseperabile; b- date separabile; c- date separabile cu date neseperabile

Figură din [63]

În tabelul 6.2 sunt redade rezultatele pentru cele 3 scenarii analizate. Ce trebuie remarcat este că funcția de marjă largă conduce la rezultate mai bune atunci când spațiul datelor este mai aglomerat. Pentru cazul cel mai defavorabil cu date foarte

Scenariu/Metrică [%]	Pr.med. 5-int.	Pr.med. 10-int.	Er. Top-1	Er. Top-5	Acc	ASG PR
<b>EI-date separabile</b>	92.82	92.93	6.85	2.46	94.82	82.81
<b>ML-date separabile</b>	92.82	92.93	6.85	2.46	94.82	<b>85.16</b>
<b>EI-date neseperabile</b>	55.83	56.00	43.61	12.69	67.20	37.25
<b>ML-date neseperabile</b>	59.75	58.97	40.60	13.45	70.08	<b>44.75</b>
<b>EI-date separabile și neseperabile</b>	72.78	71.50	25.82	7.13	78.30	52.12
<b>ML-date separabile și neseperabile</b>	72.11	72.43	22.15	9.32	79.85	<b>58.96</b>

Tabel 6.2 Rezultate experimentale Places 365 pentru noile scenarii (Pr.med.- precizie medie, Er. - Eroare, EI- entropie încrucișată , Acc. - acuratețe, ASG PR - aria de sub graficul precizie-reamintire, ML- marjă largă).

Tabel din [63]

greu separabile metricile de performanță sunt clar mai bune decât în cazul utilizării descriptorilor proveniți din antrenarea cu entropia încrucișată.

Pentru a demonstra valoarea conceptului de marjă largă poate fi urmărită figura 6.2. Aici este expusă evoluția spațiului descriptiv pentru entropia încrucișată și marja largă. Așa cum se poate observa spațiul este mult mai compact cu marja largă, ceea ce înseamnă că rețeaua va furniza descriptorii mult mai utili pentru problema de regăsire de imagini.

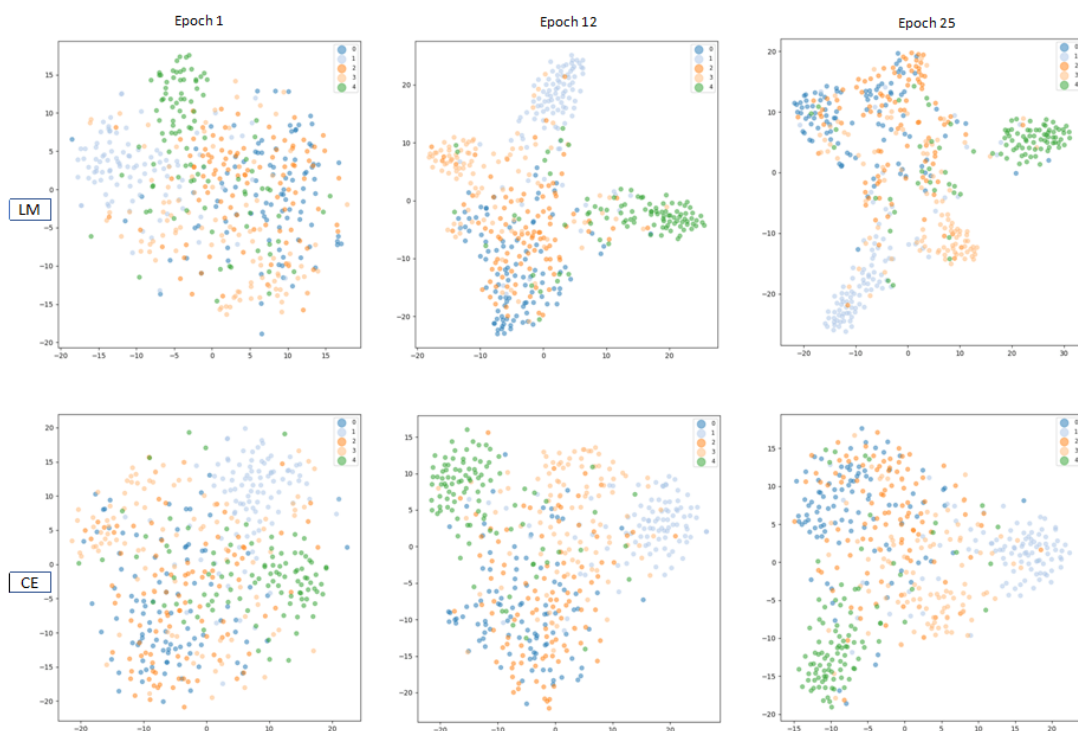


Figura 6.2 Evoluția spațiului descriptiv pe parcursul antrenării pentru cazul cu date suprapuse (ML –sus, EI- jos)

Figură din [63]

Ideea amintită anterior poate fi verificată în figura 6.3. Aici se pot vizualiza primele 5 imagini returnate pentru o serie de imagini de interogare. Se poate vedea că atunci când se folosesc descriptorii furnizați de rețeaua antrenată cu funcția de marjă largă numărul de imagini returnate care se află în clasa corespunzătoare este mai mare. Cu toate acestea, există și cazuri când entropia încrucișată oferă rezultate mai bune. Acest aspect demonstrează că pentru date foarte suprapuse nici măcar conceptul de marjă largă nu e suficient.

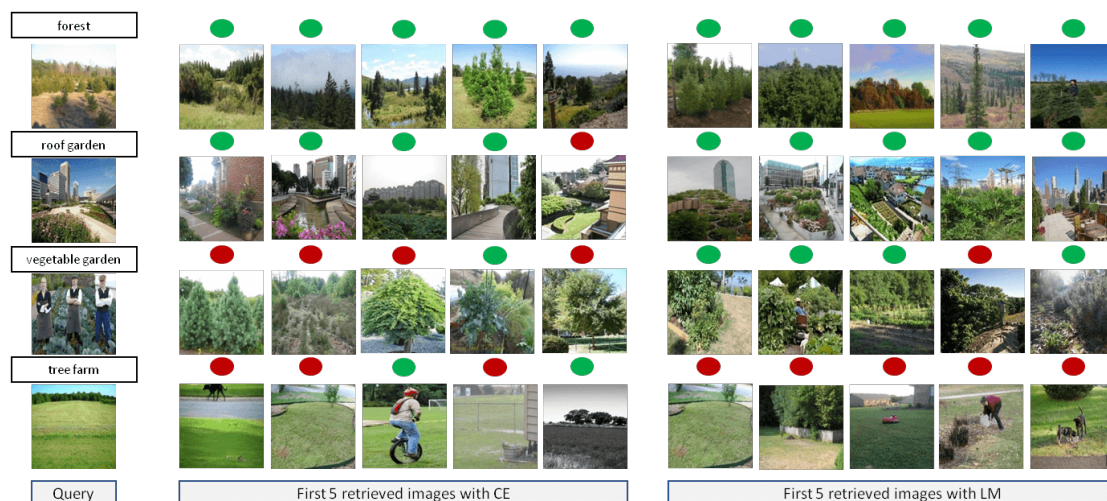


Figura 6.3 Exemple de primele 5 imagini returnate cu entropia încrucișată și marja largă. Cercurile roșii reprezintă imagini returnate care au clasa diferită față de imaginea de interogare. Cercurile verzi reprezintă imagini returnate care au aceeași clasa ca imaginea de interogare  
 Figură din [63]

Chiar dacă marja largă are limitările sale, obține o compactare mai bună a descriptorilor. În figura 6.4 este reprezentat la nivel descriptiv setul de referință din care sunt extrase imaginile pentru problema de regăsire. După cum se poate constata, atunci când se folosește rețeaua antrenată cu marjă largă, datele sunt mai compacte și mai puțin suprapuse. Acest fapt ușurează deciziile sistemului în ceea ce privește regăsirea de poze cu conținut similar. Creșterea ariei de sub graficul curbei precizie-reamintire din figura 6.5 confirmă utilitatea acestei tehnici pentru probleme cu date aglomerate.

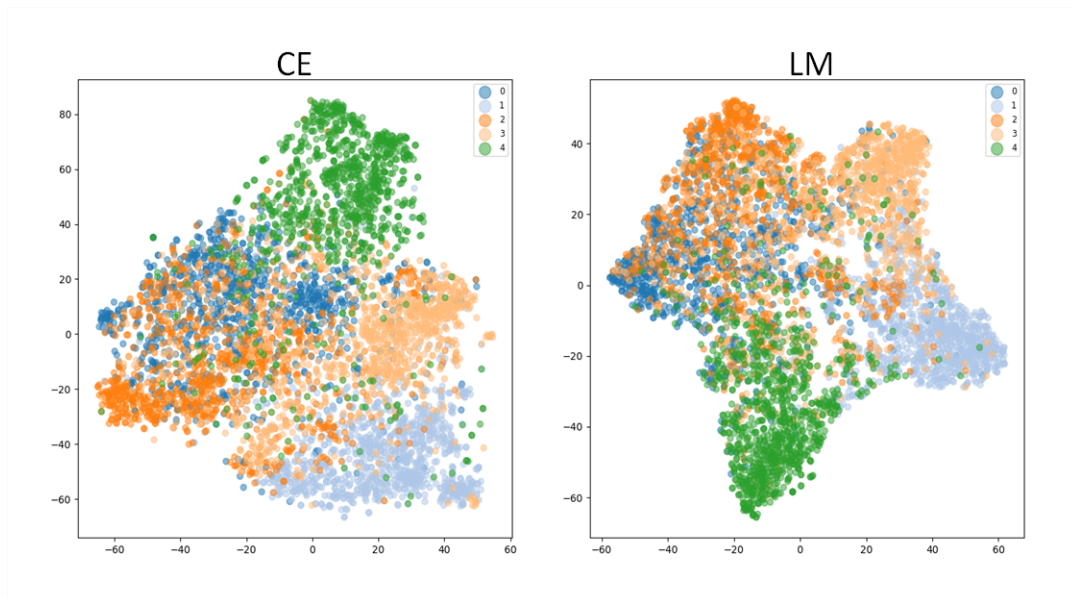


Figura 6.4 Reprezentarea descriptorilor pentru setul de referință pentru EI (stânga) și ML (dreapta) în cazul cu date neseparabile  
Figură din [63]

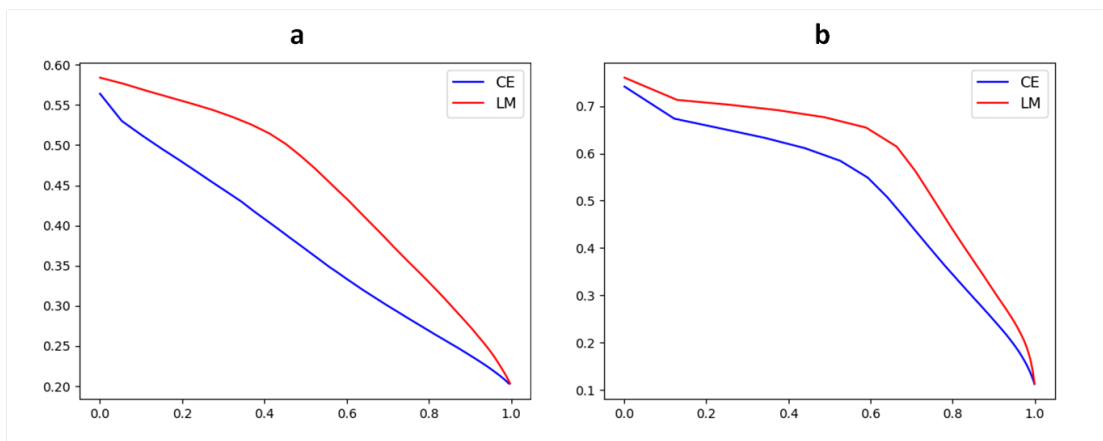


Figura 6.5 Curbele precizie-reamintire pentru EI și ML. a- scenariul cu date neseparabile; b- scenariul cu date neseparabile și date separabile)  
Figură din from [63]

# Capitolul 7

## Concluzii

### 7.1 Rezultate obținute

Pe parcursul lucrării au fost abordate 2 teme importante de cercetare: recunoașterea expresiilor faciale și regăsirea de imagini similare. Prima dintre ele s-a concentrat și pe detecția de mișcări faciale intitulate unități de acțiune. Au fost elaborate mai multe metode bazate pe învățarea semi-supervizată/adaptare de domeniu, funcții de cost cu rol de structurare al spațiului descriptiv și metode noi de augmentare/regularizare. Rezultatele satisfăcătoare au confirmat potențialul metodelor propuse.

Celălalt domeniu de interes a fost regăsirea de imagini simale. În această situație s-a testat dacă o funcție nouă de cost (marjă largă) poate grupa mai eficient descriptorii folosiți mai departe pentru găsirea de imagini asemănătoare. S-a demonstrat că acest concept este util în special pentru cazurile în care datele se suprapun foarte mult la nivel descriptiv.

### 7.2 Contribuții originale

- A fost abordată o nouă tehnică de recunoaștere a expresiilor faciale și a unităților de acțiune, care are mai multe componente originale. A fost folosită o soluție de adaptare a domeniului pentru a face legătura între mișcările faciale și expresiile discrete pentru a beneficia de potențialul învățării semi-supervizate. În plus, a fost utilizată o funcție de cost cu rol de grupare a spațiului descriptiv pentru o performanță crescută. [14]
- A fost propusă o soluție pentru recunoașterea expresiilor faciale, concentrându-se pe utilizarea unei noi metode de regularizare bazată pe injectarea aleatoare în gradient. [42, 43]

- A fost propusă o nouă modalitate care combină datele adnotate și datele neadnotate cu o tehnică de creștere a numărului de eșantioane. Algoritmul a fost testat atât în cazul expresiilor faciale, cât și pe baze de date standard [37]
- Eficacitatea unei funcții de cost cu rol de grupare a spațiului descriptiv a fost testată pentru sarcina de regăsire de imagini. În acest caz, o grupare mai eficientă a descriptorilor a crescut semnificativ performanța de regăsire pentru scenariile cu date mai suprapuse. [63]
- Metoda din [14] a fost extinsă pentru mai multe seturi de date care conțin unități de acțiune. Aici, funcția de pierdere propusă a fost studiată în raport cu alte funcții similare din literatura de specialitate și s-a dovedit a fi mai eficientă. De asemenea, s-a demonstrat că pierderea propusă contribuie la o mai bună recunoaștere a unităților de acțiune care apar mai rar în seturile de date, confirmând capacitatea acesteia de a grupa corespunzător datele. [48] [64]
- Pentru toți algoritmi propuși, a fost efectuată o comparație extinsă cu literatura de specialitate. Au fost discutate metode similare cu abordările studiate pentru a obține o idee mai obiectivă despre eficacitatea acestora.

### 7.2.1 Publicații

- Andrei Racoviteanu, Corneliu Florea, Mihai Badea, and Constantin Vertan. Spontaneous emotion detection by combined learned and fixed descriptors. In 2019 International Symposium on Signals, Circuits and Systems (ISSCS), pages 1–4. IEEE, 2019
- Andrei Racoviteanu, Iulian Felea, Laura Florea, Mihai Badea, and Corneliu Florea. Clustering based reference normal pose for improved expression recognition. In International Conference on Advanced Concepts for Intelligent Vision Systems, pages 51–61. Springer, 2018
- Mihai Badea, Constantin Vertan, Corneliu Florea, Laura Florea, and Andrei Racoviteanu. Improving small convolutional neural networks with semi-supervised learning. UPB Scientific Bulletin, Series C: Electrical Engineering, pg Series C, Vol. 84, Iss. 3, 2022, pp 107-119
- Andrei Racoviteanu, Mihai Badea, Corneliu Florea, Laura Florea, and Constantin Vertan. Dual task training for face expression recognition. In 2020 12th International Conference on Electronics, Computers and Artificial Intelligence (ECAI), pages 1–4. IEEE, 2020

- M. Boeru, A. Racovițeanu and C. Florea, "Facial Expressions Recognition by Structuring the Embeddings Space," 2021 International Conference on e-Health and Bioengineering (EHB), Iasi, Romania, 2021, pp. 1-4
- B. Stoica, L. Florea, A. Bădeanu, A. Racovițeanu, I. Felea and C. Florea, "Visual saliency analysis in paintings," 2017 International Symposium on Signals, Circuits and Systems (ISSCS), Iasi, Romania, 2017, pp. 1-4
- Badea, M., Florea, C., Racovițeanu, A., Florea, L., Vertan, C. (2023). Timid semi-supervised learning for face expression analysis. *Pattern Recognition*, 138, 109417.
- Florea, Corneliu, et al. "Automatic Real-Estate Image Analysis for Retrieval and Classification." *Bulletin of the Polytechnic Institute of Iași. Electrical Engineering, Power Engineering, Electronics Section* 68.2 (2022): 35-45.
- Corneliu Florea, Mihai Badea, Laura Florea, Andrei Racoviteanu, and Constantin Vertan. Margin-mix: Semi-supervised learning for face expression recognition. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII* 16, pages 1–17. Springer, 2020
- Corneliu Florea, Laura Florea, Mihai-Sorin Badea, Constantin Vertan, and Andrei Racoviteanu. Annealed label transfer for face expression recognition. In *BMVC*, page 104, 2019
- Andrei Racoviteanu, Mihai-Sorin Badea, Corneliu Florea, Laura Florea, and Constantin Vertan. Large margin loss for learning facial movements from pseudo-emotions. In *BMVC*, page 108, 2019
- Andrei Racoviteanu, Corneliu Florea, Mihai-Sorin Badea. Large margin loss for Image Retrieval. Accepted to *UPB Scientific Bulletin, Series C: Electrical Engineering*
- Andrei Racoviteanu, Corneliu Florea, Laura Florea, and Constantin Vertan. Normalized Margin Loss for Action Unit Detection. Submitted to *MVAP*
- Project "Technologies and innovative video/audio systems for the recognition/identification of people and simulated behavior" - SPIA-VA, PN-III-P2-2.1-SOL-2016-02-0002
- Project "TRANSLATE" , TE 66/2020, PN-III-P1-1.1-TE-2019-0543.
- Project "Innovative Artificial Intelligence systems in the field of real estate portals" - online number 137-221-A2, MySMIS number: 129132
- Project "OPTIMizarea rezultatelor cercetării aplicative a doctoranzilor și cercetătorilor postdoctorat" - nr. 62461/03.06.2022, SMIS code 153735.

### **7.2.2 Dezvoltări ulterioare**

Tehnica de adaptare a domeniului folosită pentru expresiile faciale poate fi folosită și în alte contexte. Localizarea punctelor faciale și clasificarea poziției capului sunt exemple practice, deoarece cele trei unghiuri ale capului pot fi corelate în raport cu punctele faciale; un alt exemplu este descrierea imaginilor și detectarea obiectelor, unde descrierile sunt derivate dintr-un anumit set de obiecte.

În cazul regăsirii imaginilor, funcția de cost cu rol de grupare a spațiului descriptiv s-a dovedit a fi foarte eficientă. În acest context, ar putea fi extinsă și la bazele de date cu expresii faciale. Având în vedere că multe expresii sunt similare la nivel descriptiv, această soluție are un anumit potențial.



# Bibliografie

- [1] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [2] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [3] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [4] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [6] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [7] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 896, 2013.
- [8] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017.
- [9] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019.
- [10] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- [11] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *European conference on computer vision*, pages 499–515. Springer, 2016.
- [12] Jie Cai, Zibo Meng, Ahmed Shehab Khan, Zhiyuan Li, James O’Reilly, and Yan Tong. Island loss for learning discriminative features in facial expression recognition. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 302–309. IEEE, 2018.

- [13] Y. Zheng, D. Pal, and M Savvides. Ring loss: Convex feature normalization for face recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5089–5097, 2018.
- [14] Andrei Racoviteanu, Mihai-Sorin Badea, Corneliu Florea, Laura Florea, and Constantin Vertan. Large margin loss for learning facial movements from pseudo-emotions. In *BMVC*, page 108, 2019.
- [15] Paul Ekman and Wallace V Friesen. Facial action coding system. *Environmental Psychology & Nonverbal Behavior*, 1978.
- [16] James A Russell. A circumplex model of affect. *Journal of personality and social psychology*, 39(6):1161, 1980.
- [17] Shuwen Zhao, Haibin Cai, Honghai Liu, Jianhua Zhang, and Shengyong Chen. Feature selection mechanism in cnns for facial expression recognition. In *BMVC*, page 317, 2018.
- [18] Zhiding Yu and Cha Zhang. Image based static facial expression recognition with multiple deep network learning. In *Proceedings of the 2015 ACM on international conference on multimodal interaction*, pages 435–442, 2015.
- [19] Changde Du, Changying Du, Hao Wang, Jinpeng Li, Wei-Long Zheng, Bao-Liang Lu, and Huiguang He. Semi-supervised deep generative modelling of incomplete multi-modality emotional data. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 108–116, 2018.
- [20] Zixing Zhang, Jing Han, Jun Deng, Xinzhou Xu, Fabien Ringeval, and Björn Schuller. Leveraging unlabeled data for emotion recognition with enhanced collaborative semi-supervised learning. *IEEE Access*, 6:22196–22209, 2018.
- [21] Ciprian Corneanu, Meysam Madadi, and Sergio Escalera. Deep structure inference network for facial action unit recognition. In *Proceedings of the european conference on computer vision (ECCV)*, pages 298–313, 2018.
- [22] Kaili Zhao, Wen-Sheng Chu, and Honggang Zhang. Deep region and multi-label learning for facial action unit detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3391–3399, 2016.
- [23] C Fabian Benitez-Quiroz, Ramprakash Srinivasan, Qianli Feng, Yan Wang, and Aleix M Martinez. Emotionet challenge: Recognition of facial expressions of emotion in the wild. *arXiv preprint arXiv:1703.01210*, 2017.
- [24] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-supervised learning with ladder networks. *Advances in neural information processing systems*, 28, 2015.
- [25] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. *Advances in neural information processing systems*, 27, 2014.
- [26] Augustus Odena. Semi-supervised learning with generative adversarial networks. *arXiv preprint arXiv:1606.01583*, 2016.
- [27] Paul Viola and Michael Jones. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the 2001 IEEE computer society conference on computer vision and pattern recognition. CVPR 2001*, volume 1, pages I–I. Ieee, 2001.

- [28] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters*, 23(10):1499–1503, 2016.
- [29] Ian J Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, et al. Challenges in representation learning: A report on three machine learning contests. In *International conference on neural information processing*, pages 117–124. Springer, 2013.
- [30] Emad Barsoum, Cha Zhang, Cristian Canton Ferrer, and Zhengyou Zhang. Training deep networks for facial expression recognition with crowd-sourced label distribution. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pages 279–283, 2016.
- [31] Ira Kemelmacher-Shlizerman, Steven M Seitz, Daniel Miller, and Evan Brossard. The megaface benchmark: 1 million faces for recognition at scale. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4873–4882, 2016.
- [32] Shan Li and Weihong Deng. Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition. *IEEE Transactions on Image Processing*, 28(1):356–370, 2018.
- [33] Vanessa LoBue and Cat Thrasher. The child affective facial expression (cafe) set: Validity and reliability from untrained adults. *Frontiers in psychology*, 5:1532, 2015.
- [34] Rizwan Ahmed Khan, Arthur Crenn, Alexandre Meyer, and Saida Bouakaz. A novel database of children’s spontaneous facial expressions (liris-cse). *Image and Vision Computing*, 83:61–69, 2019.
- [35] Patrick Lucey, Jeffrey F Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar, and Iain Matthews. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *2010 IEEE computer society conference on computer vision and pattern recognition-workshops*, pages 94–101. IEEE, 2010.
- [36] S Mohammad Mavadati, Mohammad H Mahoor, Kevin Bartlett, Philip Trinh, and Jeffrey F Cohn. Disfa: A spontaneous facial action intensity database. *IEEE Transactions on Affective Computing*, 4(2):151–160, 2013.
- [37] Corneliu Florea, Mihai Badea, Laura Florea, Andrei Racoviteanu, and Constantin Vertan. Margin-mix: Semi-supervised learning for face expression recognition. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*, pages 1–17. Springer, 2020.
- [38] James C Bezdek, Robert Ehrlich, and William Full. Fcm: The fuzzy c-means clustering algorithm. *Computers & geosciences*, 10(2-3):191–203, 1984.
- [39] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016.
- [40] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.

- [41] Vikas Verma, Kenji Kawaguchi, Alex Lamb, Juho Kannala, Arno Solin, Yoshua Bengio, and David Lopez-Paz. Interpolation consistency training for semi-supervised learning. *Neural Networks*, 145:90–106, 2022.
- [42] Corneliu Florea, Laura Florea, Mihai-Sorin Badea, Constantin Vertan, and Andrei Racoviteanu. Annealed label transfer for face expression recognition. In *BMVC*, page 104, 2019.
- [43] Andrei Racoviteanu, Corneliu Florea, Laura Florea, and Constantin Vertan. Randomization injection for efficient transfer in face expression recognition. In *Submitted to Applied Intelligence*, 2023.
- [44] Valentin Vielzeuf, Corentin Kervadec, Stéphane Pateux, Alexis Lechervy, and Frédéric Jurie. An occam’s razor view on learning audiovisual emotion recognition with small training sets. In *Proceedings of the 20th ACM International Conference on Multimodal Interaction*, pages 589–593, 2018.
- [45] Shuwen Zhao, Haibin Cai, Honghai Liu, Jianhua Zhang, and Shengyong Chen. Feature selection mechanism in cnns for facial expression recognition. In *BMVC*, page 317, 2018.
- [46] Yong Li, Jiabei Zeng, Shiguang Shan, and Xilin Chen. Occlusion aware facial expression recognition using cnn with attention mechanism. *IEEE Transactions on Image Processing*, 28(5):2439–2450, 2018.
- [47] Yanling Gan, Jingying Chen, and Luhui Xu. Facial expression recognition boosted by soft label with a diverse ensemble. *Pattern Recognition Letters*, 125:105–112, 2019.
- [48] Andrei Racoviteanu, Corneliu Florea, Laura Florea, and Constantin Vertan. Normalize margin loss for action units detection. In *Submitted to MVAP*, 2023.
- [49] Kaili Zhao, Wen-Sheng Chu, and Honggang Zhang. Deep region and multi-label learning for facial action unit detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3391–3399, 2016.
- [50] W. Li, F. Abtahi, and Z. Zhu. Action unit detection with region adaptation, multi-labeling learning and optimal temporal fusing. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 6766–6775, 2017.
- [51] Zhiwen Shao, Zhilei Liu, Jianfei Cai, and Lizhuang Ma. Deep adaptive attention for joint facial action unit detection and face alignment. In *Proceedings of the European conference on computer vision (ECCV)*, pages 705–720, 2018.
- [52] Guanbin Li, Xin Zhu, Yirui Zeng, Qing Wang, and Liang Lin. Semantic relationships guided representation learning for facial action unit recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8594–8601, 2019.
- [53] Jiyuan Cao, Zhilei Liu, and Yong Zhang. Cross-subject action unit detection with meta learning and transformer-based relation modeling. In *2022 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2022.
- [54] Tengfei Song, Lisha Chen, Wenming Zheng, and Qiang Ji. Uncertain graph neural networks for facial action unit detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 5993–6001, 2021.

- [55] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017.
- [56] Timo Ojala, Matti Pietikainen, and Topi Maenpaa. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on pattern analysis and machine intelligence*, 24(7):971–987, 2002.
- [57] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, volume 1, pages 886–893. Ieee, 2005.
- [58] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004.
- [59] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. *Lecture notes in computer science*, 3951:404–417, 2006.
- [60] Josef Sivic and Andrew Zisserman. Video google: A text retrieval approach to object matching in videos. In *Computer Vision, IEEE International Conference on*, volume 3, pages 1470–1470. IEEE Computer Society, 2003.
- [61] Weixun Zhou, Shawn Newsam, Congmin Li, and Zhenfeng Shao. Learning low dimensional convolutional neural networks for high-resolution remote sensing image retrieval. *Remote Sensing*, 9(5):489, 2017.
- [62] Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez. Aggregating local descriptors into a compact image representation. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 3304–3311. IEEE, 2010.
- [63] Andrei Racoviteanu, Corneliu Florea, and Mihai Badea. Large margin loss for image retrieval. In *Accepted to UPB Scientific Bulletin, Series C: Electrical Engineering*, 2023.
- [64] Andrei Racovițeanu, Mihai Badea, Corneliu Florea, Laura Florea, and Constantin Vertan. Dual task training for face expression recognition. In *2020 12th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)*, pages 1–4. IEEE, 2020.