

# Învățare Automată în Vederea Artificială și Procesarea Șirurilor de Caractere (Rezumat)

Conducător Științific:

Prof. Dr. Denis Enăchescu

Radu Tudor Ionescu

Departamentul de Informatică

Facultatea de Matematică și Informatică

Universitatea din București

București, Decembrie 2013

## Rezumat

Învățarea automată reprezintă o arie largă de cercetare ce are aplicații în multe domenii, precum vederea artificială, bioinformatica, regăsirea informației, procesarea limbajului natural, procesarea semnalelor, data mining, și multe altele. În varietatea metodelor de învățare automată propuse până în prezent, se numără metodele de învățare bazate pe similaritate. Învățarea bazată pe similaritate se referă la procesul de învățare prin folosirea similarității între perechi de exemple de antrenare. Procesul de învățare bazat pe similaritate poate fi atât supervizat, cât și nesupervizat, iar relația dintre perechi poate fi dată printr-o măsură de similaritate sau una de disimilaritate, sau chiar printr-o funcție de distanță.

Această teză studiază o serie de metode de învățare bazate pe similaritate, cum ar fi metoda celor mai apropiați vecini (Nearest Neighbor), metodele de tip nucleu (kernel), și algoritmi de clustering. În această teză este prezentată o metodă de tip Nearest Neighbor bazată pe o nouă măsură de disimilaritate pentru imagini. Metoda este aplicată pentru recunoașterea caracterelor scrise de mână, obținând rezultate foarte bune din punct de vedere al acurateții. Metodele de tip kernel sunt utilizate pentru mai multe probleme abordate în această teză. În primul rând, în teză este introdusă o nouă funcție kernel ce poate fi utilizată pentru histogramme de *cuvinte vizuale* (visual words). Folosind această funcție kernel, metodele de tip *sac de cuvinte vizuale* (bag of visual words) obțin o performanță foarte bună în cazul recunoașterii obiectelor în imagini. În al doilea rând, în teză sunt prezentate și câteva metode de tip nucleu bazate pe reprezentarea imaginilor sub formă de structură piramidală. Aceste metode sunt folosite pentru

recunoașterea expresiei faciale în imagini cu persoane care exprimă diferite stări de spirit. În al treilea rând, în teză este descrisă o abordare bazată pe *funcții de tip nucleu pentru șiruri de caractere* (string kernels) cu aplicații în identificarea limbii native în texte scrise în limba engleză de către persoane care nu sunt vorbitoare native de limba engleză. Abordarea propusă în teză obține cea mai bună performanță până în prezent pentru identificarea limbii native, fiind în același timp o abordare independentă de limbă. De asemenea, în această teză sunt studiați o serie de algoritmi de clustering. Algoritmii de clustering sunt aplicați pentru construirea arborilor filogenetici folosind secvențe de ADN mitocondrial provenit de la mamifere. Se poate observa imediat că problemele de învățare automată, prezentate în teza de față, se împart în două domenii, anume vederea artificială (computer vision) și procesarea de nivel înalt a șirurilor de caractere (string processing).

În ciuda faptului că vederea artificială și procesarea de nivel înalt a șirurilor de caractere par a fi domenii de studiu complet diferite, analiza imaginilor și a textelor este asemănătoare sub multe aspecte. Așa cum urmează a fi arătat în această teză, conceptul de a trata imaginile și șirurile de caractere (în special textele) într-un mod asemănător s-a dovedit a fi foarte fertil și productiv pentru anumite aplicații din vederea artificială. De fapt, una din metodele de ultimă oră folosite pentru clasificarea imaginilor este inspirată de reprezentarea sub formă de *sac de cuvinte* (bag of words) a documentelor text, o abordare extreme de răspândită și folosită în regăsirea informației și procesarea limbajului natural. Într-adevăr, modelul *sac de cuvinte vizuale*, ce are la bază contruirea unui vocabular de cuvinte vizuale prin folosirea unei metode de clustering pentru a grupa descriptorii locali de imagine, a demonstrat un nivel de acuratețe impresionant pentru clasificarea imaginilor, pentru regăsirea imaginilor, sau pentru alte sarcini asemănătoare. Prin adaptarea tehnicilor folosite în procesarea șirurilor de caractere pentru a analiza imagini, sau prin adaptarea tehnicilor folosite în vederea artificială pentru a analiza

șiruri de caractere, cunoștințele dintr-un domeniu pot fi transferate către celălalt domeniu. De fapt, multe dintre cele mai importante descoperiri științifice au fost realizate în urma transferului de cunoștințe între domenii de studiu diferite. Teza de față se încadrează în această direcție de cercetare, urmărind fie să prezinte noi abordări, fie să dezvolte abordările deja existente, prin transferarea și adaptarea metodelor între cele două domenii studiate (vederea artificială și procesarea șirurilor de caractere).

Mai întâi, este prezentată o nouă măsură de disimilaritate pentru compararea imaginilor. Această măsură de disimilaritate este inspirată de distanța rang (rank distance), folosită pentru compararea șirurilor de caractere. Principala sarcină care trebuie îndeplinită presupune extinderea și adaptarea distanței rang pentru a lucra cu intrări bidimensionale (imagini digitale), în loc de intrări unidimensionale (șiruri de caractere). În timp ce distanța rang este o măsură foarte precisă pentru compararea șirurilor de caractere, rezultatele empirice prezentate în această teză sugerează faptul că măsura de disimilaritate ce extinde distanța rang la imagini are o acuratețe foarte bună atât pentru recunoașterea cifrelor scrise de mână, cât și pentru clasificarea și analiza imaginilor cu texturi.

În al doilea rând, sunt prezentate câteva metode de îmbunătățire a modelului sac de cuvinte vizuale. Așa cum s-a menționat mai sus, acest model este inspirat de modelul sac de cuvinte folosit în procesarea limbajului natural și în regăsirea informației. Printre îmbunătățirile aduse acestui model se numără o nouă metodă de tip kernel (denumită PQ kernel), reprezentarea piramidală a imaginilor folosind vectori de prezență, dar și aplicarea *învățării locale* (local learning).

În al treilea rând, în teza de față este introdusă o nouă distanță pentru șiruri de caractere. Această distanță este inspirată de noua măsură de disimilaritate pentru imagini descrisă mai sus. Construită astfel încât să se conformeze unor principii mai generale, dar adaptată în același timp pentru secvențe de ADN, noua distanță obține rezultate

mai bune în comparație cu alte metode de ultimă oră folosite pentru analiza secvențelor ADN. Mai mult, această distanță este aplicată cu succes și în altă direcție, cea a comparării documentelor text. Mai exact, o metodă de tip nucleu, ce are la bază această nouă distanță, este folosită pentru identificarea limbii native în text. În concluzie, toate contribuțiile prezentate în această teză vin să susțină ideea că imaginile și șirurile de caractere pot fi tratate într-un mod asemănător, cu scopul de a îmbunătăți rezultatele metodelor de învățare.

Înainte de a încheia, trebuie atrasă atenția asupra faptului că metodele studiate și descrise în această teză obțin o performanță comparabilă sau uneori mai bună decât metodele de ultimă oră din cele două domenii abordate. În continuare, sunt prezentate câteva argumente care să confirme acest fapt. În primul rând, un model de tip sac de cuvinte vizuale, îmbunătățit din mai multe puncte de vedere, a obținut locul al patrulea la competiția Facial Expression Recognition (FER) Challenge organizată în cadrul workshop-ului ICML 2013 Workshop in Challenges in Representation Learning (WREPL). În al doilea rând, sistemul ce are la bază metode de tip nucleu pentru șiruri de caractere, prezentat în această teză, s-a clasat pe locul al treilea la competiția Native Language Identification Shared Task organizată în cadrul workshop-ului BEA-8 Workshop of NAACL 2013. În al treilea rând, articolul care introduce funcția de tip nucleu, denumită PQ kernel, folosită pentru a calcula similaritatea între histogramme de cuvinte vizuale, a obținut premiul Caianiello Best Young Paper Award la conferința ICIAP 2013.