

MINISTERUL EDUCAȚIEI ȘI CERCETĂRII AL REPUBLICII MOLDOVA

Universitatea Tehnică a Moldovei

Facultatea Calculatoare, Informatică și Microelectronică

Departamentul Ingineria Software și Automatică

Admis la susținere

Șef departament:

I. Fiodorov, conf. univ. dr.

„\_\_\_\_\_” \_\_\_\_ 2022

# Analiza predicției personalității utilizând date din rețele sociale

Teză de masterat

Masteranda: \_\_\_\_\_ Sclifos Corina,  
gr. TI-201M

Conducător: \_\_\_\_\_ Duca Ludmila,  
asistent universitar

Consultant: \_\_\_\_\_ Cojocaru Svetlana,  
asistent universitar

Chișinău, 2022

## ADNOTARE

Teza de master are la bază următoarea structură de redare:

- introducere;
- analiza domeniului de studiu;
- metodologia;
- arhitectura sistemului;
- colectarea și procesarea datelor;
- concluzii;
- bibliografie.

Cuvintele cheie: Neural Network Language Model (NNLM), Long Short-Term Memory (LSTM), Convolutional Neural Network (CNN), dataset predicție, Big 5, model neuronal.

Scopul acestei teze de masterat constă în:

- analiza metodelor de predicție a personalității;
- analiza procesării limbajului natural;
- analiza modelului de limbaj de rețea neuronală;
- analiza tipurilor de personalitate și specificul lor;
- implimentarea unui model.

Personalitatea caracterizează în mod distinct un individ și influențează profund comportamentele acestuia. Rețelele de socializare oferă comunității virtuale o oportunitate fără precedent de a genera conținut și de a împărtăși aspecte din viața lor care reflectă adesea personalitatea lor. Interesul pentru utilizarea învățării profunde pentru a deduce trăsături din amprente digitale a crescut recent; cu toate acestea, au fost prezentate lucrări foarte limitate care explorează informațiile despre sentimente transmise. Prin urmare, prezentul studiu a folosit o abordare computațională pentru a clasifica personalitatea din rețelele sociale prin măsurarea percepțiilor publice care stau la baza factorilor care cuprind trăsăturile.

În cercetarea raportată în această teză, se planifică dezvoltarea unui sistem de detectare a personalității bazat pe sentimente pentru a deduce trăsături din texte scurte bazate pe dimensiunile personalității „Big 5”. Am exploatat spiritul Neural Network Language Model (NNLM) folosind un model unificat care combină o rețea neuronală recurentă numită Long Short-Term Memory (LSTM) cu o rețea neuronală convoluțională (CNN).

## ANNOTATION

The master's thesis has the following rendering structure:

- introduction;
- study analysis;
- methodology;
- system architecture;
- data collection and processing;
- conclusions;
- bibliography.

Keywords: Neural Network Language Model (NNLM), Long Short-Term Memory (LSTM), Convolutional Neural Network (CNN), dataset prediction, Big 5, neural model.

The purpose of this master's thesis is to:

- analysis of personality prediction methods;
- analysis of natural language processing;
- analysis of the neural network language model;
- analysis of personality types and their specificity;
- implementation of a neuronal model.

Personality distinctly characterizes an individual and profoundly influences behaviors. Social networks offer the virtual community or unprecedented opportunities to generate content and share aspects of their lives that often reflect their personality. The interest in using deep learning to deduce features from the fingerprint has recently increased; however, very limited papers have been presented that explore information about the feeling conveyed. Therefore, the present study used a computational approach to classify personality in social networks by measuring public perceptions that underlie traits.

In the research reported in this thesis it is planned to develop a personality detection system based on feelings to deduce features from short texts based on the dimensions of the "Big 5" personality. We exploited the spirit of the Neural Network Language Model (NNLM) using a unified model that combines or recurring neural network called Long Short-Term Memory (LSTM) with a convolutional neural network (CNN).

# CUPRINS

INTRODUCERE .....	11
1. ANALIZA DOMENIULUI .....	12
1.1 Prelucrarea limbajului natural .....	12
1.2 Model de limbaj de rețea neuronală .....	13
1.2.1 Rețele neuronale convoluționale.....	15
1.2.2 Rețele neuronale recurente .....	16
1.2.3 Reprezentarea cuvântului .....	18
1.3 Social Media Mining .....	20
1.4 Trăsături de personalitate .....	21
1.4.1 Tipuri de Personalitate Myers-Briggs .....	22
1.5 Rezumat .....	27
2. METODOLOGIE .....	28
2.1 Design de cercetare .....	28
2.1.1 Arhitectura sistemului .....	29
2.1.2 Instrumente și biblioteci .....	31
2.2 Colectarea datelor .....	31
2.2.1 Selectarea datelor.....	31
2.2.2 Recuperarea datelor .....	33
2.2.3 Analiza datelor.....	33
2.2.4 Pregătirea datelor .....	35
2.3 Clasificarea sentimentelor în Twitter .....	38
2.3.1 Resurse lexicale .....	39
2.4 Detectarea personalității .....	39
2.4.1 Înglobarea cuvintelor .....	40
2.4.2 Antrenarea rețelei .....	40
2.4.3 Evaluarea modelului .....	40
2.4.4 Rezumat .....	41
3. DEZVOLTAREA SISTEMULUI.....	42
3.1 Definierea sistemului .....	42
3.1.1 Modelul de dezvoltare software .....	43
3.1.2 Arhitectura software .....	44
3.2 Specificații de sistem .....	45
3.2.1 Cerințe preliminare .....	45
3.3 Implementarea .....	46

3.3.1 Preprocesare .....	46
3.3.2 Modelul .....	47
3.4 Antrenarea .....	48
3.5 Rezumat .....	50
4 REZULTATE ȘI EVALUARE .....	51
4.1 Evaluare .....	51
4.2 Analiză .....	52
4.2.1 Eficacitatea modelului și implicații .....	52
4.2.2 Word Clouds .....	53
4.2.3 Test pe rețele sociale .....	55
4.3 Rezumat .....	55
CONCLUZII .....	57
BIBLIOGRAFIE .....	58

## INTRODUCERE

În ultimii ani, creșterea și colectarea informației a început într-un ritm accelerat odată cu apariția rețelelor sociale, în special sub formă de tipuri de date textuale. Conform raportului Social Media Trend publicat în [1], există 3,8 miliarde de utilizatori activi de social media din lume începând cu ianuarie 2020, cu o creștere estimată de 9,2% a utilizatorilor în fiecare an. Adesea, oamenii folosesc rețelele sociale pentru a se exprima asupra anumitor probleme legate de viața lor și de ființele familiale, psihologie, probleme financiare, interacțiunea cu societatea și mediul, precum și politica. În unele cazuri, aceste expresii pot fi folosite pentru a caracteriza comportamentul individual și personalitatea. Studiile anterioare (de exemplu [2]) demonstrează că există o corelație puternică între personalitățile utilizatorilor și comportamentul lor online pe rețelele de socializare. Exemple de aplicații care pot profita de informațiile despre personalitatea utilizatorilor sunt sistemele de recrutare, sistemele de consiliere personală, marketing online, sistemele de recomandări personale și sistemele de notare a creditelor bancare [3].

Datorită ambiguităților inerente ale limbajelor naturale, dezvoltarea unui model eficient de predicție a personalității bazat pe mesajul pe care utilizatorul îl partajează pe social media poate fi o sarcină extrem de dificilă. Având în vedere aceste ambiguități, s-au făcut multe progrese în domeniul procesării limbajului natural (NLP). Până în prezent NLP a permis computerelor să înțeleagă cuvinte sau propoziții scrise în limbajul uman [4]. Concepte lingvistice precum partea de vorbire (substantive, verbe, adjective) și structurile gramaticale sunt de obicei utilizate în NLP.

Predicția automată a personalității a devenit un subiect larg discutat. Referințele menționate mai sus arată că personalitatea poate fi definită ca un model de influență sau personalitate folosit pentru a caracteriza indivizii unici. Predicția personalității existente a exploatat algoritmul de învățare profundă și învățare automată pentru a îmbunătăți acuratețea clasificării. Cu toate acestea, această abordare are limitări pentru a extrage caracteristici contextuale din propoziție datorită limitării algoritmului de calcul și a problemei vocabularului în afara utilizării setului de date predefinit.

În această lucrare, procesarea limbajului natural, extragerea textului, gruparea vor fi folosite pentru a face analiza personalității din datele rețelelor sociale.

## **BIBLIOGRAFIE**

1. Violino B. Social media trends. Association for Computing Machinery. Commun ACM.
2. Riccardi G. Alam F Stepanov EA. Personality traits recognition on social network—Facebook. AAAI Workshop—Technical Report.
3. Myint PH. Aung ZMM. Personality prediction based on content of facebook users: a literature re- view. Proceedings - 20th IEEE/ACIS International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing.
4. Kuchy L. Ben-Porat O Hirsch S. Predicting strategic behavior from free text.
5. F. J. Damerau N. Indurkha. Handbook of natural language processing (2nd ed. 2010).
6. Schutze H. Manning C. Foundations of statistical natural language processing. (1999).
7. I. Zitouni D. Bikel. Multilingual natural language processing applications.
8. Rath S. K. Tripathy A. Agrawal A. Classification of sentiment reviews using n-gram machine learning approach. (2016).
9. Y. Kim. Con-volutional Neural Networks for Sentence Classification 2014.
10. Bengio Y. LeCun Y. The HANDBOOK of BRAIN theory AND NEURAL networks.
11. Collobert R. Natural language processing (almost) from scratch.
12. S. Haykin. NEURAL Networks: A Comprehensive FOUNDATION.
13. D. Shi. A study on neural network language modeling.
14. Danihelka I. Graves A. Kalchbrenner N. Grid long short-term memory. (2015).
15. Shalev A. Mandelbaum A. Word embeddings and their use in sentence classification tasks. (2016).
16. A. Black W. Ling C. Dyer. Two/Too Simple Adaptations of Word2Vec for Syntax Problems.
17. Craswell N. Diaz F. Mitra B. Query expansion with locally- trained word embeddings. (2016).
18. Alaedini Z Dalvi-Esfahani M Niknafs A. Social Media Addiction and Empathy: Moderating impact of personality traits among high school students. Telematics Inform. 2020.
19. Tang Y. Han S Huang H. Knowledge of words: An interpretable approach for personality recognition from social media
20. Cuzzocrea A Howlader P Pal KK. Predicting facebook-users' personality based on status and linguistic features via flexible regression analysis techniques. Proc ACM Symposium Appl Comput