

MINISTERUL EDUCAȚIEI



UNIVERSITATEA TEHNICĂ
DIN CLUJ-NAPOCA

FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE
DEPARTAMENTUL CALCULATOARE

SISTEME INTELIGENTE

PROIECT

Student: **Crisan Mihai-George**

2024

Capitolul 1. Introducere

Proiectul are ca scop dezvoltarea și antrenarea unui sistem inteligent capabil de a recunoaște dacă o persoană fuge sau merge normal în funcție de datele obținute de 2 senzori (gyroskop și accelerometru) aflați în telefonul unei persoane.

Motivația alegerii proiectului și a setului de date este începerea regulată de a exersa, în principal prin alergare.

Acest sistem inteligent este gândit cu scopul utilitar de a detecta în timp real dacă o persoană este angajată în alergare, astfel eliminând nevoia de a porni un program de fitness manual și fără a mai avea grija că o perioadă de fugă nu a fost înregistrată.

Tehnologii folosite:

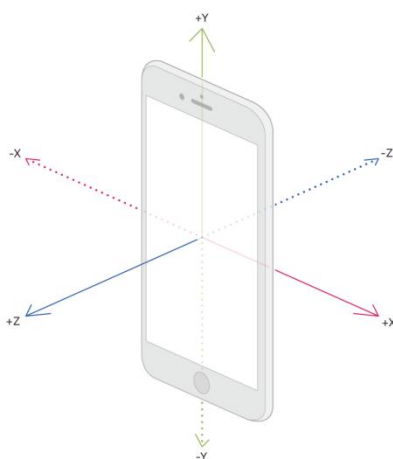
- Dataspell: IDE folosit în dezvoltarea proiectului
- Python: limbajul de programare folosit în dezvoltarea proiectului
- Jupyter: un tip de formatare a documentelor dezvoltate pe baza JSON
- Pandas: o librărie folosită în manipularea seturilor mari de date
- XGBoost: o librărie folosită pentru implementarea

Capitolul 2. Context

2.1. Context General

[Baza de date a fost aleasă de pe Kaggle](#), datele fiind colectate de utilizatorul Viktor Malyi utilizând un telefon iPhone 5c. În tabelă sunt înregistrate 88588 de valori ale datei, a timpului la care au fost înregistrate alături de valoarea fiecărei senzori.

Senzorii folosiți în colectarea datelor sunt accelerometrul și giroscopul integrat în telefon. Din cei doi senzori sunt înregistrate câte 3 valori pe axele înălțimii, lățimii și a lungimii, pentru accelerometru și pe 3 axe pentru giroscop. Datele înregistrate de senzori sunt numere reale, cu valori cuprinse între -5.35 și 5.6 pentru accelerometru și între -7.46 și 8.5 pentru giroscop.



(Fig 0)

În figura 0 sunt prezentate axele de colecție a datelor. Este important de menționat că datele au o variație mare deoarece la colectarea lor rotația telefonului are un rol important în determinarea semnului datelor.

Acest lucru oferă o acuratețe mai mare a algoritmilor antrenați pe aceste seturi de date și permit implementarea modelului în o gamă mai largă de aplicații.

2.2 Tema Proiectului

Tema proiectului este de a dezvolta un sistem de învățare automată folosind date reale și accesibile oricui prin intermediul internetului. Proiectul este împărțit în 7 etape principale, fiecare fiind o parte importantă din dezvoltarea sistemului.

2.3 Rezultate Propuse

La finalul proiectului dorim să obținem un sistem inteligent capabil de a detecta dacă o persoană se află în fugă sau în mers pe baza input-urilor senzorilor înregistrate în setul de date. Acuratețea propusă pe baza datelor este de peste 70%.

Capitolul 3. Aspecte Teoretice

3.1 Starea domeniului

Domeniul inteligenței artificiale și a sistemelor în starea lor curentă este unul vast. De la antrenarea sistemelor în a detecta celule canceroase, la strategiști politici ce folosesc sisteme de predicție pentru a calcula următoarele mișcări electoraleⁱ, sistemele inteligente devin din ce în ce mai ușor de încorporat în orice aplicație ce poate beneficia de analiza seturilor mari de date sau de predicții cu o acuratețe mare.

Utilizarea sistemelor inteligente implică o bună cunoștință a multor aspecte teoretice relevante. În următoarele subcapitole sunt descrise principalele concepte folosite în realizarea acestui proiect.

3.2 Entropia

În domeniul inteligenței artificiale entropia poate fi definită ca fiind o unitate de măsură pentru a determina nivelul de dezordine sau de incertitudine dintr-un sistem. Este folosită în a determina calitatea unui set de date, indicând dacă setul respectiv conține informații distribuite.

Presupunem următoarele cazuri:



(Fig. 1)



(Fig. 2)

Din figurile prezentate mai sus se observă că în figura 1 avem un set de date cu entropie mică (entropia este chiar 0, deoarece toate punctele sunt identice) iar în figura 2 avem un set de date cu entropie relativ ridicată.

Intr-un set de date cu C clase entropia este calculată utilizând formula următoare:

$$E = - \sum_i^C p_i \log_2 p_i$$

Unde p_i este probabilitatea de a alege un element din mulțimea setului de date (i.e. proporția setului de date care este i)

Originea entropiei se găsește în lucrarea “A Mathematical Theory of Communication”ⁱⁱⁱ, termenul de entropie fiind împrumutat din domeniul fizicii.

3.3 Corelații

Corelațiile dintre date ne arată relevanța unor anumite clase de date față de alte clase de aparținând aceluiaș set de date. Calcularea corelațiilor într-un set de date ne ajută în a dezvolta sisteme ce au acuratețe mai mareⁱⁱⁱ.

În lumea reală corelațiile nu promet că două evenimente sunt în directă legătură una față de cealaltă, însă pot indica înspre direcția potrivită. Evenimente precum detectarea realității atacurilor cibernetice^{iv} sau identificarea unei stări emoționale a unei persoane^v nu pot fi prezise în acuratețe maximă, însă prin analizarea datelor se pot obține linii principale care indică spre un anumit rezultat.

3.3 Modele

În proiect au fost alese inițial 4 tipuri de modele pentru a fi testate pentru implementarea finală: Regresia Liniară, Regresia Logistică, LSTM și XGBoost. În continuare sunt descrise fiecare model și o descriere simplificată a fiecărui tip de algoritm.

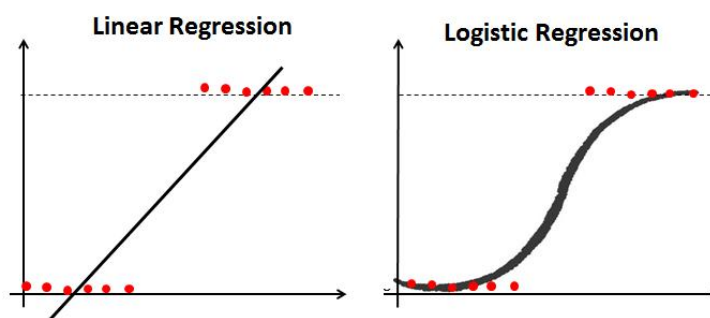
3.3.1 Regresie Liniară

Regresia liniară este unul dintre algoritmii de bază în dezvoltarea sistemelor inteligente. Regresia este folosită în a determina corelațiile dintre două sau mai multe variabile având relații de cauză-efect și pentru a oferi predicții folosind relațiile dintre obiecte^{vi}.

3.3.2 Regresia Logistică

Regresia logistică este asemănătoare în funcționalitate cu regresia liniară, singura diferență fiind modul de determinare a valorilor maxime.

Regresia logistică are valorile “blocate” între o valoare minimă și o valoare maximă, spre deosebire de regresia liniară ce explorează datele în continuitate



(Fig 3)

În figura alăturată se observă diferența dintre distribuția datelor în o regresie liniară față de o regresie logistică.

3.3.3 Long Short Term Memory

LSTM se încadrează în clasa sistemelor recurente, folosind reprezentări a input-uri recent introduse (short term), față de parametrii modificați în timp (long term)^{vii}.

Aceste sisteme se bazează pe truncarea în diferite serii a datelor, în proiect ele fiind împărțite în o serii temporale.

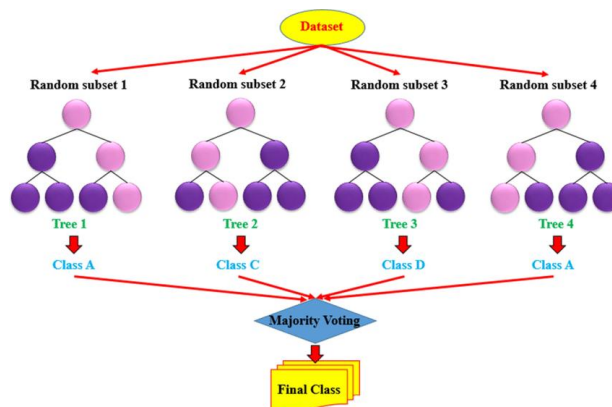
Acest algoritm nu a prezentat rezultate satisfăcătoare din perspectiva acurateței, astfel nu a fost dezvoltat în detaliu.

3.3.4 XGBoost

XGBoost (eXtrem Gradient Boosting) face parte din categoria “Gradient Boosting”, ce s-a dovedit în timp ca fiind una dintre cele mai rapide și acurate sisteme inteligente.^{viii}

Algoritmii de tip gradient boost încearcă să gasească o funcție a căror rezultate sunt apropiate de valoarea label-urilor oferite în setul de date, ele fiind împărțite în diferite subseturi de date.

XGBoost este foarte util în rezolvarea problemelor de clasificare, ceea ce este exact ce proiectul necesită: clasificarea între două acțiuni pe baza datelor.



(Fig 4)

În figura 4 este prezentat modul de funcționare al algoritmului într-o manieră abstractă.

Capitolul 4. Implementarea Aspectelor Teoretice

În proiect au fost implementate aspectele teoretice relevante, ele având ca scop înțelegerea datelor sau realizarea modelelor în sine.

4.1 Implementarea și Analiza Corelației

Corelația a fost calculată utilizând implementările standard utilizate prezente în librăriile utilizate.

Rezultatul este graficul prezent în figura 5:

	wrist	activity	acceleration_x	acceleration_y	acceleration_z	gyro_x	gyro_y	gyro_z
wrist	1.000	-0.113	-0.610	0.087	0.324	-0.022	-0.068	0.009
activity	-0.113	1.000	-0.018	0.640	-0.192	0.041	0.012	-0.008
acceleration_x	-0.610	-0.018	1.000	-0.265	-0.552	-0.022	-0.004	-0.061
acceleration_y	0.087	0.640	-0.265	1.000	0.106	0.011	0.072	-0.023
acceleration_z	0.324	-0.192	-0.552	0.106	1.000	0.035	-0.021	0.050
gyro_x	-0.022	0.041	-0.022	0.011	0.035	1.000	0.094	0.318
gyro_y	-0.068	0.012	-0.004	0.072	-0.021	0.094	1.000	0.287
gyro_z	0.009	-0.008	-0.061	-0.023	0.050	0.318	0.287	1.000

(Fig 5)

Din graficul de corelație se pot trage următoarele axe utile pentru înțelegerea datelor:

- Din [figura 0](#) se poate explica de ce există o corelație mare între "wrist" și direcțiile de accelerare pe X și Y.
- "Acceleration_y" și "activity" prezintă o corelație ridicată, implicându-se creșterea accelerației atunci când "activity" este 1 (1 indicând alergat)
- Din rotația telefonului pe încheietură cât și din încheietura pe care a fost ținută se explică și corelațiile între accelerațiile prezente în tabel. ("acceleration_x" cu "acceleration_z")
- Se observă că datele giroscopice nu au corelația ridicată în legătura cu restul datelor.

4.2 Implementare și Analiza Entropiei

În proiect au fost folosite diferite metode de a calcula entropia (în baza e), însă a fost aleasă o metodă generală pentru a calcula entropia.

În setul de date ales, s-a calculat o entropie maximă posibilă de 11.39. Pentru fiecare dintre seturile de coloane cu o importanță mărită rezultând în următoarele rezultate;

- Acceleration_x -> 9.93
- Acceleration_y -> 9.79
- Acceleration_z -> 9.34

Aceste date indică că setul de date are o multitudine de date diferite, ceea ce poate duce spre un model cu o rată de detecție mai ridicată.

4.3 Implementarea Modelelor

Modele au fost implementate utilizând diferite tehnici. Modele de regresie au fost implementate utilizând librăria “sklearn”, modelul de LSTM a fost implementat utilizând librăria “tensorflow” iar modelul XGBoost utilizând librăria “xgboost”.

Modelele au fost antrenate pe aceleași date, împărțite între setul de antrenament și setul de validare (împărțirea s-a făcut în proporție de 80/20).

Rezultatele obținute sunt următoarele:

- Regresie liniară : valoarea erori absolute medie : ~ 0.29

Valoarea erori absolute medie este mică, indicând o performanță ridicată.

- Regresie logistică : acuratețea modelului este de $\sim 85\%$

Acuratețea este una ridicată, indicând o performanță ridicată a modelului

- LSTM: valoarea erori absolute medie este de 0.19,

Valoarea erori absolute medie este foarte mică, indicând o performanță foarte ridicată.

- XGBoost: acuratețea modelului este de 99%

Din rezultatele obținute se observă cum modelul clasifică datele cu o acuratețe foarte ridicată, indicând că acest model este aproape perfect în a prezice dacă o persoană se află în fugă sau în mers.

Acest model este cel ales pentru implementare.

Capitolul 5. Testare și validare

5.1 Împărțirea datelor și antrenarea modelului final

Înainte de a se antrena modele este important de menționat că nu toate datele din setul de date inițial nu au fost necesare pentru antrenarea modelului.

Datele ce nu au fost luate în considerare sunt:

- "Username" : datele au fost colectate doar de un utilizator.
- "Time" & "Date" : ora și data pot oferi o direcție greșită modelului, inducând posibilitatea de a prezice o rutină în loc de a prezice o activitate^{ix}.

5.2 Antrenarea modelului final

Modelul final a fost antrenat pe setul de antrenament (80% din totalul setului de date) și a fost valid pe setul de validare (20% din totalul setului de date).

Modelul a fost inițiat folosind "XGBClassifier" inclus în librăria "xgboost", inițiat cu variabilele "random_state" = 42, "verbosity" = 0, "silent" = 0.

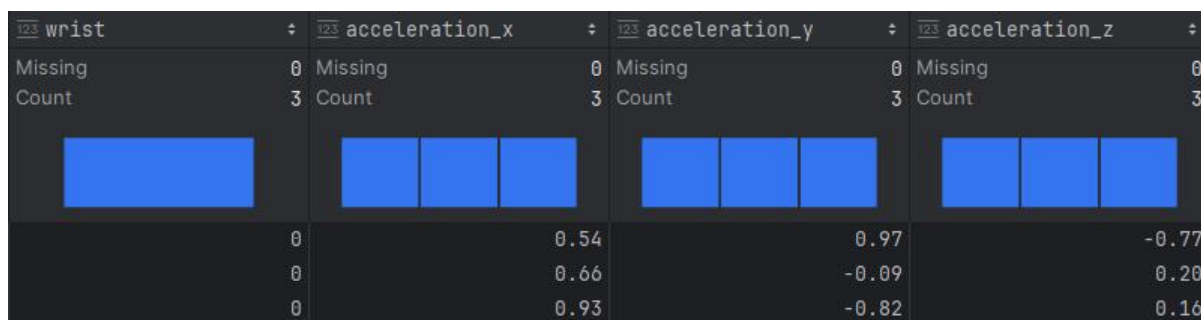
După antrenarea modelului a rezultat o acuratețe de 99%.

Modelul a fost reantrenat de mai multe ori pentru a asigura că nu a fost un incident unic. De asemenea se recomandă reantrenarea pe termen lung regulată a modelului pentru a ne asigura că datele coincid cu realitatea, deoarece datele colectate pot fi influențate de factori externi precum vârsta sau sănătatea persoanei responsabile de colectarea datelor.^x

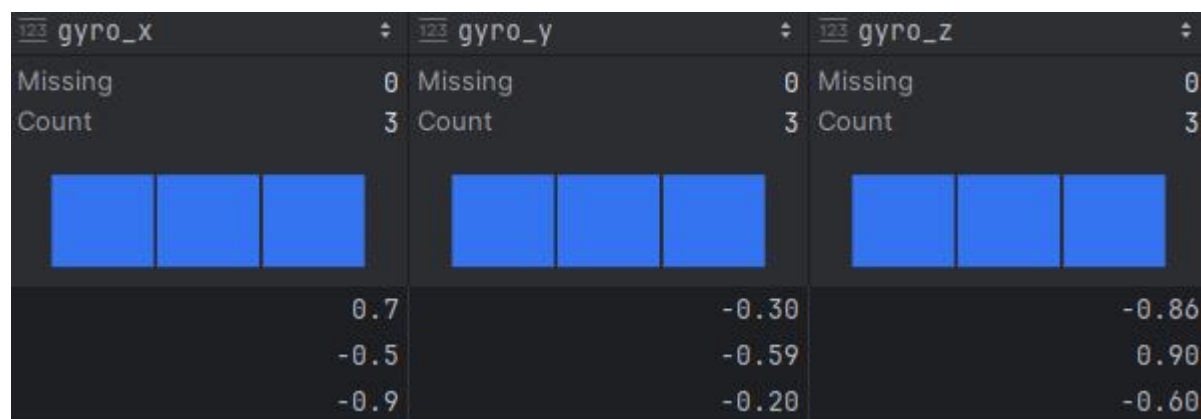
Capitolul 6. Testare manuală și Rezultate

6.1 Colectarea datelor pentru testare

Datele de testare au fost colectate în mod asemănător cu metoda de colectare a setului de date inițial. Aplicația utilizată pentru colectare nu a fost făcută publică însă folosind aplicația “Sensor Recorder” publicată de Nils Ackermann au fost colectate aceleași date care ar fi fost colectate de aplicația utilizată de autorul setului de date.



(Fig 6)



(Fig 7)

Datele ce au fost colectate sunt alese special pentru a testa dacă sistemul deosbește alergatul (primele două rânduri) de mers (ultimul rând)

Au fost testate aproximativ 50 de instanțe de date în diferite ipostaze, însă aceste 3 au fost alese, ele fiind cele mai apropiate de media datelor.

6.2 Testarea pe date culese manual și validarea modelului

Pe datele culese manual, modelul a prezentat o acuratețe de 100%. Acest lucru se explică prin setul restrâns de date colectate manual.

În cele 3 instanțe prezente în figurile 6 și 7 rezultatele au fost următoarele : alergat, alergat, mers.

Rezultatele coincid cu activitatea în care ele au fost recoltate.

Capitolul 7. Concluzie

În acest proiect s-a încercat crearea unui model capabil să determine activitatea unei persoane pe baza datelor colectate de senzorul giroscopic cât și a senzorului de accelerație prezenți într-un telefon.

Au fost antrenate diferite modele, fiecare generând un rezultat bun, însă pentru modelul final a fost ales XGBoost, având o rată de acuratețe extrem de ridicată.

Datele colectate au fost curățate în mare parte deja de utilizatorul care le-a înregistrat, însă în acest proiect nu au fost utilizate toate câmpurile, ele având ocazia de a influența negativ modelul final.

În viitor se recomandă implementarea modelului în o aplicație IOS, pentru ca el să poată fi testat în timp real. Acest lucru ar duce la o validare totală a modelului, ce ar putea să influențeze modul în care aplicațiile de urmărire al exercițiilor abordează înregistrarea alergării.

ⁱ Min Yin, Jennifer Wortman Vaughan, Hannah Wallach - Understanding the Effect of Accuracy on Trust in Machine Learning Models

ⁱⁱ [Claude E. Shannon - A Mathematical Theory of Communication](#)

ⁱⁱⁱ [Mark A. Hall - Correlation-Based Feature Selection for Machine Learning](#)

^{iv} [Ibrahim Ghafir, Mohammad Hammoudeh, Vaclav Prenosil, Liangxiu Han, Robert Hegarty, Khaled Rabie, Francisco J. Aparicio-Navarro - Detection of advanced persistent threat using machine-learning correlation analysis](#)

^v [Sunil Kumar, Ilyoung Chong - Correlation Analysis to Identify the Effective Data in Machine Learning: Prediction of Depressive Disorder and Emotion States](#)

^{vi} [Gulden Kaya Uyanik, Neşe Guler - A study on multiple regression analysis](#)

^{vii} [Sepp Hochreiter, Jurgen Schmidhuber - LONG SHORT- TERM MEMORY](#)

^{viii} [Didrik Nielsen - Tree Boosting with XGBoost](#)

^{ix} [Alexander Martin Musgnug - Thre predictive reframing of machine learning applications: good predictions and bad measurements](#)

^x [Diego Klabjan, Xiaofeng Zhu - Neural Network Retraining for Model Serving](#)