

Contract nr.1 din 01/03/2024, etapa 1 - GenDeg

RAPORT ȘTIINȚIFIC ȘI TEHNIC 2024

Referință 7075/02.12.2024

Manager de proiect: ***Laura Andreica***

Istoricul versiunilor

Versiune	Autor	Modificări
0.1	Laura Andreica	Versiunea inițială
0.5	Laura Andreica	Versiunea intermediară
1.0	Laura Andreica	Versiunea finală

Cuprins

1	Introducere	5
2	Despre proiectul GenDeg	5
3	Activități planificate	6
4	Activități efectuate	7
4.1	Devieri de la planificare	7
5	Starea curentă a domeniului	7
5.1	Model Drift și Model Degradation	7
5.2	Tehnici de detectare a derivei modelelor	7
5.3	Strategii de mitigare a derivei modelelor	8
5.4	Rezultate și concluzii preliminare	8
6	Cerințele platformei	8
6.1	Cerințe funcționale	9
6.2	Cerințe nonfuncționale	9
6.3	Impactul așteptat al platformei	10
7	Calitatea datelor și metodologia propusă	10
7.1	Importanța Calității Datelor	10
7.2	Metodologia de Evaluare a Calității Datelor (DQAM)	10
7.2.1	Măsurile utilizate în DQAM	10
7.2.2	Etapele procesului de evaluare a calității datelor	11
7.3	Calculul scorului general	11
7.3.1	Calculul simetriei datelor	12
7.3.2	Importanța monitorizării calității datelor	12
7.4	Beneficiile metodologiei propuse	12
7.5	Validarea metodologiei	12
8	Rolul HOLISUN în cadrul proiectului	13
8.1	Contribuții tehnice	13
8.2	Impact asupra produselor și serviciilor HOLISUN	13
8.3	Abordare interdisciplinară	13
9	Extras din planul de riscuri	14
10	Rezultatele proiectului	15
10.1	Livrabile	15
10.2	Articole științifice	15

11 Diseminare și exploatare	16
11.1 Activități de diseminare	16
12 Concluzii	16

Parteneri



(a) AI Investments (Polonia) Coordonator



(b) InbestMe (Spania)



(c) Holisun SRL (Romania)

Figura 1: Partenerii proiectului *GenDeg*

1 Introducere

GenDeg își propune să inoveze evaluarea modelelor de prognoză bazate pe serii temporale prin monitorizarea generalizărilor și a degradării acestora în timp. Proiectul va dezvolta metode avansate care să funcționeze cu o varietate de algoritmi de inteligență artificială (AI) și să fie aplicabile în diverse domenii. Aceste progrese vor fi integrate în serviciile oferite de AI Investments, inbestMe și Holisun, adăugând caracteristici noi și inovatoare.

Soluțiile rezultate din proiect GenDeg vor avea un impact semnificativ în domeniul financiar și dincolo de acesta, utilizând cele mai recente progrese în învățarea automată (ML) și AI pentru a revoluționa procesele de investiții. Funcționalități precum selecția strategiilor viitoare optime (rezultate out-of-sample) și monitorizarea degradării modelelor (cum ar fi pierderea progresivă a "alphas") vor remodela industria investițiilor. Proiectul își propune, de asemenea, să demonstreze aplicabilitatea metodelor dezvoltate în alte sectoare prin intermediul Holisun care va folosi modele și rezultatele proiectului pentru îmbunătățirea aplicației de mentenanță predictivă, îmbunătățind astfel serviciile oferite de toți cei trei parteneri din consorțiu.

Prezentul raport oferă o imagine de ansamblu asupra cadrului operațional și a designului metodologiilor din cadrul proiectului GenDeg. Acesta subliniază contribuția proiectului la crearea unor soluții integrate pentru evaluarea și optimizarea modelelor de prognoză, oferind funcționalități inovatoare care pot aduce beneficii tangibile în domeniul financiar și în alte industrii.

2 Despre proiectul GenDeg

Proiectul **GenDeg** își propune să inoveze evaluarea modelelor de prognoză bazate pe serii temporale prin monitorizarea generalizabilității și degradării acestora în timp. Acest proiect vizează dezvoltarea unor algoritmi unici pentru identificarea degradării performanței modelelor predictive bazate pe inteligență artificială (AI) și evaluarea capacității acestora de a generaliza rezultatele pe perioade viitoare. Metodele rezultate vor putea fi aplicate într-o gamă largă de aplicații AI și vor contribui la îmbunătățirea produselor și serviciilor oferite de AI Investments, inbestMe și Holisun.

Figura 2 prezintă arhitectura logică de înalt nivel a proiectului **GenDeg**.

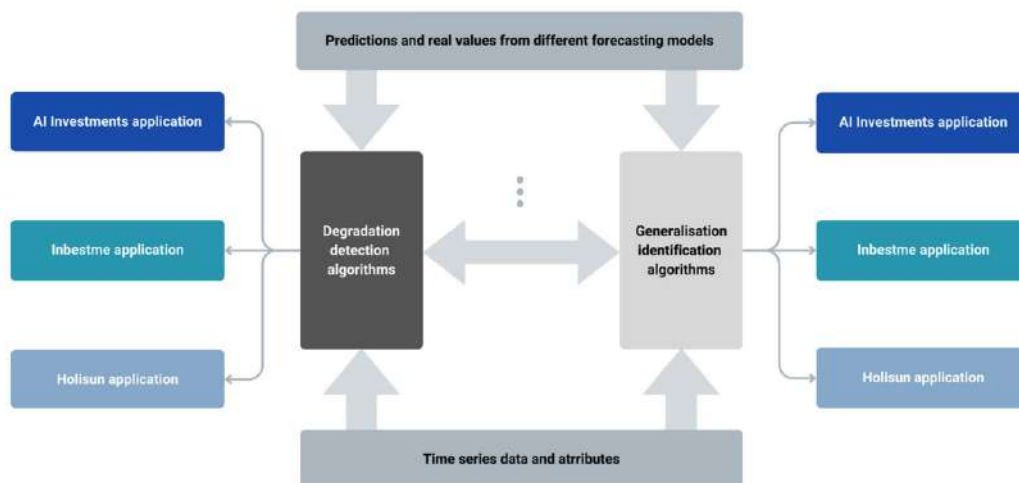


Figura 2: Arhitectura logică de înalt nivel a proiectului GenDeg

Aceasta evidențiază fluxurile de informații și componentele cheie implicate în procesele de detectare a degradării și identificare a generalizabilității modelelor predictive. Sistemul propus include următoarele componente majore:

- **Algoritmi de detectare a degradării:** Responsabili pentru identificarea scăderii performanței modelelor predictive în timp, utilizând date din serii temporale și valorile reale generate.
- **Algoritmi de identificare a generalizabilității:** Proiectați pentru a selecta modelele care au capacitatea de a oferi performanțe optime pentru date viitoare, maximizând acuratețea predicțiilor.

- **Integrarea aplicațiilor AI:** Soluțiile rezultate vor fi implementate în serviciile oferite de AI Investments, inbestMe și Holisun, demonstrând aplicabilitatea acestora atât în domeniul financiar, cât și în alte industrii.

Soluțiile propuse sunt construite pe baza datelor provenite din serii temporale și a caracteristicilor asociate acestora. Acestea includ atât valori generate de modele predictive, cât și valori reale. Algoritmii dezvoltăți vor îmbunătăți capacitatea de decizie a aplicațiilor AI, precum și performanța acestora în diverse contexte.

Datele colectate și rezultatele obținute vor fi puse la dispoziția comunității științifice pentru a încuraja reproducibilitatea și a stimula cercetările ulterioare. Proiectul **GenDeg** se prezintă ca o soluție inovatoare în domeniul prognozei bazate pe serii temporale, cu impact semnificativ asupra industriei financiare și dincolo de aceasta.

Proiectul este împărțit în 5 pachete de lucru (WP), fiecare cu sarcini clare care vor permite îndeplinirea obiectivelor proiectului printr-o structură coerentă care va facilita și managementul proiectului. Toate aceste WP-uri și sarcini detaliate sunt exemplificate în Diagrama GANTT de mai jos (Figura 3).

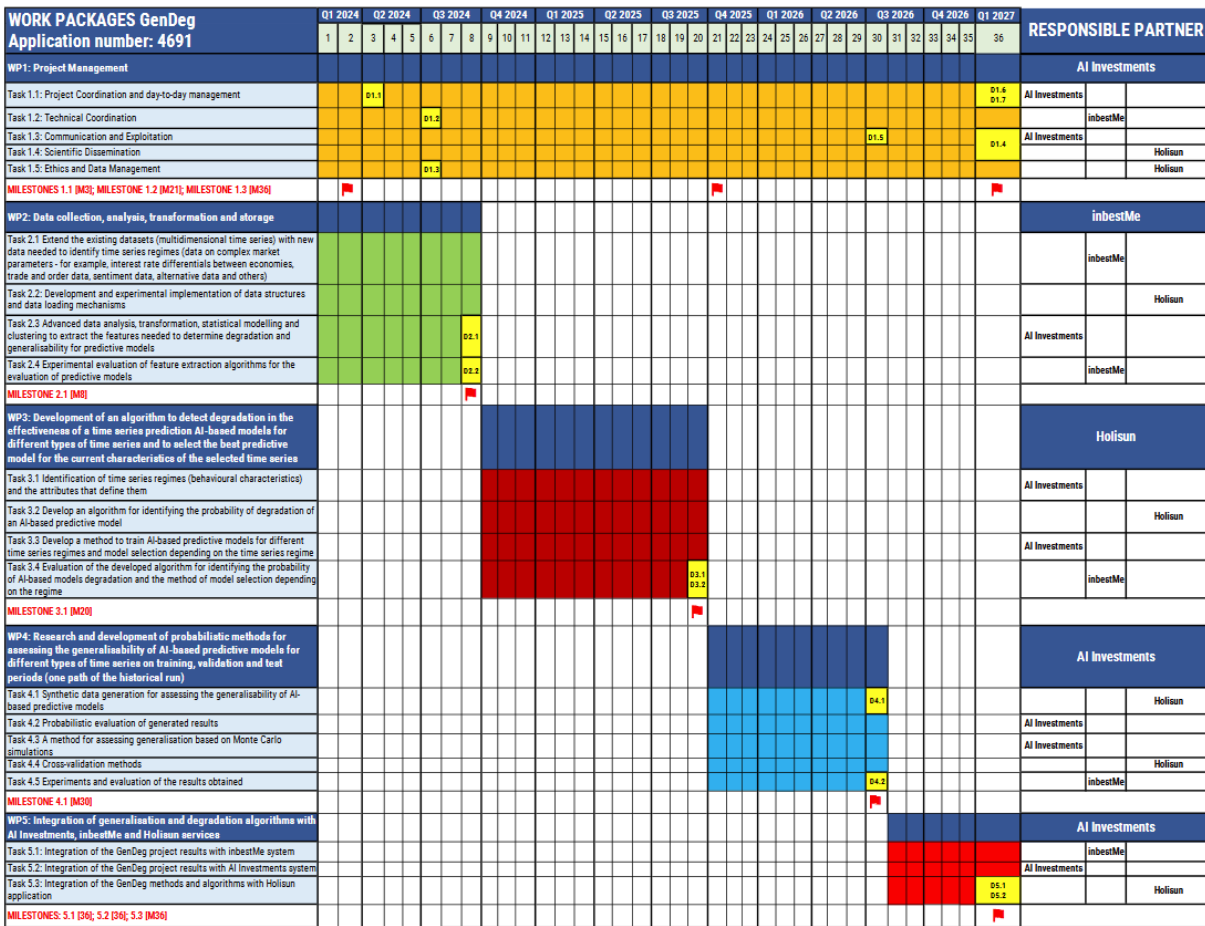


Figura 3: Diagrama Gantt care arată calendarul pachetelor de lucru și activitățile proiectului

3 Activități planificate

În perioada 01.03.2024 - 05.12.2024 au fost planificate următoarele activități:

- Cercetarea literaturii de specialitate;
- Întâlniri săptămânale de progres.

4 Activități efectuate

În perioada 01.03.2024 - 05.12.2024 au fost efectuate următoarele activități:

- Cercetare literaturii de specialitate cu privire la conceptele de Model Drift și Model Degradation;
- 1 întâlnire fizică de kick-off a proiectului cu toți partenerii proiectului 21 Martie 2024 Barcelona, Spania.
- Întâlniri săptămânale de progres, în fiecare zi de luni.

4.1 Devieri de la planificare

În perioada raportată nu au fost devieri de la planificare, sub nici un aspect.

5 Starea curentă a domeniului

În cadrul primului an al proiectului **GenDeg**, cercetările s-au concentrat pe înțelegerea și aprofundarea problemelor de *Model Degradation* (degradarea performanței modelelor predictive) și *Model Drift* (deriva modelelor predictive) în contexte complexe și dinamice. Studiile noastre au vizat sintetizarea literaturii de specialitate, identificarea celor mai relevante tehnici existente și formularea unor strategii viitoare pentru asigurarea performanței și generalizabilității modelelor predictive pe termen lung.

5.1 Model Drift și Model Degradation

Model Drift și Model Degradation reprezintă două dintre cele mai importante provocări în domeniul inteligenței artificiale (AI) și al învățării automate (ML). Aceste fenomene descriu pierderea treptată sau bruscă a performanței modelelor predictive în timp, cauzată de schimbări în mediul de date, distribuția caracteristicilor sau relațiile dintre variabilele independente și țintele predicției [5, 3].

Cercetările noastre au identificat două tipuri principale de Model Drift:

- **Concept Drift:** Apare atunci când relația dintre caracteristicile de intrare și ținta predicției se schimbă. Acest fenomen este frecvent întâlnit în aplicații financiare și de sănătate, unde contextul economic sau clinic evoluează constant [11].
- **Data Drift:** Reprezintă modificările în distribuția datelor de intrare. De exemplu, în aplicațiile de comerț electronic, preferințele utilizatorilor și modelele de cumpărare se schimbă frecvent, afectând modelele predictive [7].

Un exemplu practic îl reprezintă predicțiile financiare utilizate pentru investiții bursiere. Schimbările rapide în piață, declanșate de crize economice sau fluctuații valutare, pot reduce dramatic eficiența unui model de predicție [6]. În aceste cazuri, detectarea și adaptarea rapidă a modelelor devin esențiale.

Model Degradation, pe de altă parte, se referă la pierderea treptată a performanței modelului datorită modificării contextului operațional sau a faptului că modelele nu mai reflectă realitatea curentă a datelor. În primul an al proiectului, am explorat cauzele fundamentale ale degradării, incluzând supra-antrenarea modelelor, complexitatea excesivă a arhitecturii și inabilitatea de a încorpora dinamica schimbării datelor.

5.2 Tehnici de detectare a derivei modelelor

Detectarea derivei modelelor este o prioritate majoră în cercetările noastre. O abordare proactivă în acest sens poate preveni degradarea performanței prin detectarea timpurie a schimbărilor în datele de intrare sau în relațiile dintre acestea. Am analizat și evaluat mai multe tehnici existente, printre care:

- **Drift Detection Method (DDM):** Acest algoritm monitorizează rata de eroare a modelului și detectează derivatele bazându-se pe modificările statistice. Este eficient în detectarea schimbărilor bruște, dar mai puțin sensibil la schimbările graduale [10].

- **Early Drift Detection Method (EDDM):** O extensie a DDM, concepută pentru a fi mai sensibilă la schimbările treptate. Această metodă este recomandată pentru aplicații care implică date dinamice [9].
- **ADWIN (ADaptive WINdowing):** Utilizează o fereastră glisantă adaptativă pentru detectarea derivei în fluxurile de date în timp real. Este ideal pentru aplicații care necesită timpi de răspuns rapizi [8].
- **Page-Hinkley Test:** Această metodă detectează schimbările abrupte ale mediei datelor și este potrivită pentru fluxuri continue de date [13].
- **CUSUM:** Detectează modificările treptate prin monitorizarea sumei cumulative a deviațiilor față de o medie țintă. Este potrivit pentru detectarea derivei subtile și persistente [11].

Am constatat că metode precum ADWIN și CUSUM oferă performanțe superioare în detectarea derivei în fluxuri de date continue, fiind aplicabile în domenii precum analiza financiară și mentenanța predictivă.

5.3 Strategii de mitigare a derivei modelelor

Pentru a atenua impactul derivei modelelor, este necesară integrarea unor strategii robuste care să minimizeze pierderile de performanță. În cadrul cercetărilor din primul an, ne-am concentrat pe următoarele metode:

- **Reantrenare periodică:** Modelele sunt reantrenate regulat pe seturi de date actualizate. Această abordare este eficientă pentru scenarii cu schimbări predictibile în date [3].
- **Învățare online (Online Learning):** Modelele sunt actualizate continuu, pe măsură ce noi date devin disponibile. Este o metodă eficientă pentru fluxurile de date în timp real [12].
- **Human-in-the-Loop:** Implicarea experților umani în validarea modelelor și ajustarea manuală a acestora în cazul detectării derivei. Este o metodă esențială în aplicațiile critice, precum sănătatea [1].
- **Metode de ansamblu (Ensemble Methods):** Combinarea mai multor modele predictive pentru a îmbunătăți robustețea și a reduce impactul derivei [2].

5.4 Rezultate și concluzii preliminare

Cercetările din primul an al proiectului **GenDeg** au oferit o perspectivă valoroasă asupra provocărilor asociate cu Model Drift și Model Degradation. Am identificat tehnici relevante de detectare și mitigare, adaptate pentru aplicațiile noastre financiare și industriale. De asemenea, am elaborat un set preliminar de indicatori pentru măsurarea performanței modelelor în medii dinamice, inclusiv:

- Rata de eroare cumulativă pe fluxuri continue.
- Stabilitatea predicțiilor în condiții de drift incremental.
- Timpul de răspuns al detectării derivei în aplicații critice.

Planurile pentru anul următor includ:

- Implementarea algoritmilor de detectare în fluxurile de date furnizate de partenerii proiectului.
- Dezvoltarea unui cadru scalabil pentru reantrenare automată și selecția modelelor generalizabile.
- Validarea soluțiilor dezvoltate pe seturi de date reale, cu focus pe aplicațiile financiare AI Investments și inbestMe, precum și în mentenanța predictivă cu soluțiile Holisun.

6 Cerințele platformei

Dezvoltarea unei platforme inovative pentru monitorizarea degradării performanței modelelor predictive și a derivei acestora implică cerințe clare și bine definite. Aceste cerințe sunt fundamentale pentru asigurarea performanței, scalabilității și robusteții sistemului. Pe baza analizei realizate, cerințele platformei **GenDeg** pot fi împărțite în două categorii principale: **cerințe funcționale** și **cerințe nonfuncționale**.

6.1 Cerințe funcționale

1. Detectarea degradării modelelor predictive

- Dezvoltarea unor algoritmi avansați capabili să monitorizeze și să detecteze degradarea performanței modelelor predictive în timp real.
- Integrarea acestor algoritmi în fluxurile operaționale ale platformei pentru actualizarea automată a modelelor.

2. Identificarea derivei modelelor (Model Drift)

- Detectarea schimbărilor în distribuția datelor sau în relațiile dintre caracteristicile acestora și variabilele țintă.
- Utilizarea de tehnici de analiză a fluxurilor de date pentru a preveni degradarea performanței.

3. Generalizabilitate ridicată

- Optimizarea modelelor predictive pentru a oferi rezultate relevante pe seturi de date neobservate anterior.
- Testarea continuă a modelelor pentru a garanta acuratețea predicțiilor în contexte variate.

4. Reantrenare automată

- Implementarea unui mecanism de reantrenare a modelelor pe baza datelor actualizate, asigurând adaptarea acestora la noile condiții.

5. Evaluarea performanței algoritmilor

- Dezvoltarea unor instrumente pentru evaluarea continuă a performanței algoritmilor pe baza KPI-urilor (Key Performance Indicators).

6. Integrarea în sisteme existente

- Asigurarea compatibilității cu infrastructurile IT ale partenerilor și a posibilității de extindere către alte industrii.

6.2 Cerințe nonfuncționale

1. Scalabilitate

- Platforma trebuie să poată susține creșterea dimensiunii și a complexității fluxurilor de date, menținând în același timp performanța.

2. Performanță ridicată

- Algoritmii implementați trebuie să aibă latență redusă și consum minim de resurse, chiar și în condiții de volum mare de date.

3. Reziliență și toleranță la defecțiuni

- Sistemul trebuie să fie capabil să facă față erorilor hardware și software, asigurând continuitatea operațiunilor.

4. Confidențialitate și securitate a datelor

- Implementarea unor măsuri stricte de protecție a datelor, inclusiv criptare, control al accesului și tehnici de anonimizare.

5. Consum redus de resurse

- Optimizarea algoritmilor și a fluxurilor de lucru pentru a reduce cerințele hardware și consumul energetic.

6. Interoperabilitate

- Sistemul trebuie să fie compatibil cu diverse tipuri de infrastructuri și tehnologii, facilitând integrarea cu alte platforme.

6.3 Impactul așteptat al platformei

Platforma **GenDeg** propune o soluție inovativă pentru monitorizarea performanței și generalizabilității modelelor predictive. Aceasta va contribui la:

- Reducerea costurilor operaționale prin detectarea timpurie a degradării modelelor și adaptarea automată la noile condiții.
- Creșterea preciziei predicțiilor și a fiabilității acestora în medii dinamice.
- Extinderea aplicabilității în diverse sectoare industriale, incluzând finanțele și mentenanța predictivă.

Această platformă va oferi un punct de referință pentru viitoarele soluții bazate pe inteligența artificială și va transforma modul în care modelele predictive sunt utilizate și gestionate.

7 Calitatea datelor și metodologia propusă

Calitatea datelor (DQ) este esențială în dezvoltarea și implementarea modelelor predictive, deoarece datele de calitate scăzută pot duce la rezultate inexacte și decizii greșite. În cadrul proiectului **GenDeg**, ne concentrăm pe asigurarea unui set de date de înaltă calitate care să permită modelelor de machine learning (ML) să genereze predicții precise și generalizabile. Pentru a aborda această provocare, propunem utilizarea unei metodologii standardizate și automatizate pentru evaluarea calității datelor, denumită *Data Quality Assessment Methodology* (DQAM)[4].

7.1 Importanța Calității Datelor

În era Big Data, cantitatea și diversitatea informațiilor disponibile sunt în continuă expansiune. Această creștere pune presiune asupra sistemelor de învățare automată (ML) care trebuie să proceseze volume mari de date pentru a face predicții și a extrage informații valoroase. Cu toate acestea, calitatea datelor este adesea neglijată, iar datele incomplete, imprecise sau inconsistente pot afecta negativ performanța modelelor predictive. De asemenea, un set de date de calitate scăzută poate introduce erori în modelele ML, ceea ce poate duce la decizii eronate sau neîncredere în modelele generate.

Evaluarea calității datelor este, prin urmare, un pas esențial în orice proces de prelucrare a datelor, iar metodologia propusă de noi se bazează pe un set de indicatori relevanți care măsoară diferite aspecte ale calității datelor.

7.2 Metodologia de Evaluare a Calității Datelor (DQAM)

Metodologia de Evaluare a Calității Datelor (DQAM) propusă de noi este un sistem rapid, automatizat și flexibil, care poate fi integrat ușor în orice flux de lucru de machine learning, inclusiv în sisteme de învățare distribuită cum ar fi *Federated Learning*. DQAM este implementată în Python și este compatibilă cu Pandas DataFrames, ceea ce o face ușor de utilizat în cadrul proceselor de prelucrare a datelor.

7.2.1 Măsurile utilizate în DQAM

Pentru a evalua calitatea datelor, DQAM utilizează mai multe măsuri și metrici relevante, inclusiv:

- **Completitudinea datelor:** Măsoară procentajul valorilor lipsă din setul de date. Datele lipsă sunt o problemă majoră, deoarece nu oferă informații valoroase și reduc dimensiunea efectivă a setului de date.
- **Entropia informațională:** Aceasta măsoară diversitatea datelor și ajută la detectarea anomaliilor, absenței datelor sau prezenței unui număr mare de valori aberante. Calculată folosind Entropia Shannon, acest indicator ajută la identificarea datelor care sunt distribuite uniform și care nu conțin inconsistențe semnificative.
- **Simetria datelor:** Analizează dacă distribuția datelor este echilibrată sau distorsionată. O distribuție asimetrică poate afecta biasul modelului de ML, iar DQAM măsoară diferența dintre media și mediana valorilor, oferind o idee despre posibilele probleme de simetrie ale datelor.

- **Calitatea seriilor temporale:** În cazul datelor de tip serie temporală, DQAM analizează continuitatea acestora. Lipsa datelor într-o serie temporală poate indica probleme tehnice sau defecte ale senzorilor, care trebuie detectate și corectate.

7.2.2 Etapele procesului de evaluare a calității datelor

Metodologia DQAM se bazează pe patru etape principale:

1. **Înlocuirea valorilor lipsă cunoscute:** Prima etapă presupune identificarea și înlocuirea valorilor lipsă, un pas esențial pentru a pregăti datele pentru evaluare.
2. **Evaluarea calității pentru fiecare coloană:** Această etapă măsoară completitudinea, entropia și simetria datelor la nivelul fiecărei coloane din setul de date. Se calculează score-ul de calitate pentru fiecare coloană, care este apoi utilizat pentru a determina calitatea generală a setului de date.
3. **Evaluarea calității pentru datele de tip serie temporală:** În cazul în care setul de date include serii temporale, se măsoară continuitatea acestora. Discontinuitățile sunt detectate și cuantificate folosind indicatori precum indicele GINI pentru a evalua nivelul de inegalitate dintre intervalele de timp.
4. **Calculul scorului final:** După evaluarea fiecărei coloane și a seriei temporale, se calculează scorul final al calității datelor, care poate fi exprimat pe o scală de la 0 la 100.

7.3 Calculul scorului general

Calculul scorului general al calității datelor (DQ) implică agregarea mai multor metrici, fiecare ponderată în funcție de importanța sa. Acest lucru asigură că scorul final reflectă în mod fidel calitatea dataset-ului. În stabilirea scorului total, este logic să se aplice o penalizare mai mare pentru valorile lipsă (M_v) decât pentru simetria datelor (D_s). Pentru a aborda această problemă, propunem un set de ponderi derivate empiric, bazate pe contribuția metricii respective la cantitatea totală de informație a setului de date. Aceste ponderi pot fi ajustate pentru a corespunde cerințelor specifice ale aplicației.

Este important de menționat că, în Ecuțiile 1 și 2, suma ponderilor trebuie să fie întotdeauna egală cu 1.0.

$$DQ_c = 0.5 \cdot M_v + 0.35 \cdot V_e + 0.15 \cdot D_s \quad (1)$$

$$DQ_t = 0.7 \cdot \left(\frac{1}{n} \cdot \sum_{i=1}^n DQ_{c_i} \right) + 0.3 \cdot TSq \quad (2)$$

Unde:

- M_v - procentajul valorilor lipsă din coloana respectivă;
- V_e - entropia valorilor, exprimată ca procent din entropia maximă;
- D_s - simetria datelor, calculată conform Ecuției 3;
- DQ_c - scorul calității datelor pentru o coloană individuală;
- DQ_t - scorul total al calității datelor pentru întregul set de date;
- TSq - scorul calității datelor de tip serie temporală (*optional*);
- 0.7 și 0.3 - ponderi pentru contribuțiile scorurilor DQ_c și ale calității seriilor temporale, stabilite pe baza experimentelor preliminare.

Toate metricile sunt normalizate în intervalul $[0, 1]$, bazat pe numărul total de rânduri, ceea ce le face ușor de interpretat ca procente. Ponderile pot fi ajustate în funcție de aplicație.

7.3.1 Calculul simetriei datelor

Simetria unui set de date poate influența semnificativ calitatea informațiilor extrase. Ecuația 3 este utilizată pentru a cuantifica simetria unei coloane:

$$Ds = 1 - \frac{|Mean - Median|}{Range} \quad (3)$$

Unde:

- *Mean* - media aritmetică a valorilor din coloană;
- *Median* - mediana valorilor din coloană;
- *Range* - intervalul valorilor (*Max* – *Min*).

7.3.2 Importanța monitorizării calității datelor

Asigurarea calității datelor este un proces continuu, care necesită monitorizare constantă, curățare și îmbunătățire. Datele de înaltă calitate nu doar că îmbunătățesc performanța și fiabilitatea modelelor de învățare automată, dar facilitează și înțelegerea sistemului monitorizat. Mai mult, indicatorii calității datelor pot acționa ca un mecanism secundar de alarmare, detectând anomalii sau disfuncționalități în procesul de colectare a datelor.

Proiectanții de sisteme ar trebui să ia în considerare fluctuațiile temporale ale metricilor calității datelor și să le compare cu mecanismele tradiționale de alarmare pentru a identifica modele sau inconsistențe. Această abordare integrată asigură robustețea și fiabilitatea întregului sistem, facilitând soluții avansate bazate pe date.

7.4 Beneficiile metodologiei propuse

Implementarea DQAM în proiectul GenDeg adresează unele dintre cele mai importante provocări legate de calitatea datelor. Prin automatizarea procesului de evaluare și prin oferirea unor măsuri de calitate precise și cuantificabile, metodologia contribuie semnificativ la îmbunătățirea fiabilității și performanței modelelor ML. Aceasta permite:

- Detectarea timpurie a problemelor de calitate a datelor, ceea ce duce la o mai bună gestionare a datelor și reducerea erorilor în modelele predictive.
- Creșterea preciziei și generalizabilității modelelor de învățare automată, prin asigurarea că datele utilizate sunt de înaltă calitate.
- Reducerea timpului și resurselor necesare pentru curățarea și preprocesarea datelor.

7.5 Validarea metodologiei

Metodologia DQAM a fost validată într-un caz de utilizare din domeniul agriculturii, unde a fost aplicată pentru evaluarea calității datelor provenite de la senzori care monitorizează condițiile de mediu. Rezultatele obținute au demonstrat eficiența metodei în identificarea și corectarea problemelor de calitate ale datelor, cum ar fi valorile lipsă și erorile de măsurare, iar scorurile obținute au permis o îmbunătățire semnificativă a predicțiilor și recomandărilor pentru fermieri.

În concluzie, metodologia DQAM propusă oferă o abordare inovatoare și eficientă pentru evaluarea calității datelor în contextul modern al Big Data și al Machine Learning. Aplicarea acesteia în proiectul GenDeg va contribui semnificativ la creșterea fiabilității și performanței sistemului, asigurând în același timp o gestionare mai bună a datelor.

8 Rolul HOLISUN în cadrul proiectului

Rolul HOLISUN în cadrul proiectului **GenDeg** este esențial, contribuind atât la dezvoltarea algoritmilor inovativi de identificare a degradării performanței modelelor predictive, cât și la integrarea acestora în soluții software scalabile. HOLISUN aduce expertiza sa în cercetarea și dezvoltarea inteligenței artificiale și în construirea de sisteme software performante, având ca obiectiv îmbunătățirea preciziei predicțiilor și extinderea durabilității modelelor AI.

8.1 Contribuții tehnice

HOLISUN se implică activ în:

- **Dezvoltarea algoritmilor de detecție a degradării modelelor predictive:** HOLISUN explorează metode avansate, inclusiv abordări bazate pe învățare automată și analize probabilistice, pentru a monitoriza și detecta scăderile de performanță ale modelelor predictive în timp. Acest proces este esențial pentru menținerea eficienței și a fiabilității modelelor implementate.
- **Implementarea metodelor pentru generalizarea performanței modelelor:** HOLISUN lucrează la evaluarea capacității modelelor predictive de a performa pe seturi de date neobservate anterior. Prin aceasta, se urmărește reducerea riscurilor de suprapotrivire și creșterea aplicabilității modelelor pe perioade viitoare.
- **Integrarea rezultatelor în soluții software:** HOLISUN asigură integrarea algoritmilor dezvoltați în platforme software adaptabile, optimizând infrastructura de procesare a datelor și fluxurile de lucru existente.

8.2 Impact asupra produselor și serviciilor HOLISUN

Rezultatele proiectului vor fi incluse în portofoliul de soluții inteligente HOLISUN, oferind clienților:

- Sisteme software capabile să monitorizeze și să gestioneze performanța modelelor predictive în timp real.
- Reducerea costurilor de mentenanță prin detectarea timpurie a degradării și prin reantrenarea eficientă a modelelor.
- Extinderea duratei de viață a soluțiilor software, prin reducerea necesității de intervenții manuale frecvente.

8.3 Abordare interdisciplinară

HOLISUN colaborează îndeaproape cu ceilalți parteneri din proiect, incluzând AI Investments și inbestMe, pentru a asigura o integrare fluidă a algoritmilor dezvoltați în domenii diverse precum finanțele și mentenanța predictivă. Această abordare interdisciplinară sporește potențialul de aplicare a rezultatelor proiectului, oferind soluții inovative adaptate mai multor sectoare.

Prin implicarea sa, HOLISUN consolidează capacitatea proiectului **GenDeg** de a livra soluții de ultimă generație care să adreseze provocările legate de degradarea modelelor predictive și să susțină adoptarea inteligenței artificiale în medii operaționale complexe.

9 Extras din planul de riscuri

În Tabelul 1 este reprezentat planul de riscuri ce ține de partea de implementare a proiectului **GenDeg**.

Tabela 1: Tabel de analiză a riscurilor și metode de mitigare pentru proiectul **GenDeg**.

Risc	Probabilitate	Impact	Valoare	Mitigare
Erori în identificarea degradării modelelor predictive	Mediu	Mare	12	Dezvoltarea unei funcționalități de monitorizare continuă pentru verificarea corectitudinii rezultatelor și a KPI-urilor algoritmilor.
Vulnerabilitatea algoritmilor GenDeg la amenințări externe	Mic	Mare	10	Implementarea unor configurații de disponibilitate ridicată, teste riguroase de asigurare a calității și verificări de securitate pentru fiecare versiune de software.
Preocupări privind confidențialitatea datelor în prognozele demografice	Mare	Mediu	15	Implementarea criptării robuste a datelor, a controlului accesului, tehnici de anonimizare și respectarea reglementărilor relevante privind protecția datelor.
Proiecții inexacte ale infrastructurii de mediu	Mic	Mare	10	Validarea regulată a modelelor cu date istorice și scenarii multiple pentru a gestiona incertitudinea.
Dependența excesivă de Gen-Deg pentru intervenții climatice	Mediu	Mediu	12	Refinarea continuă a modelului pentru a include cele mai recente avansuri științifice privind clima și colaborarea cu experți în domeniu.
Capacitate insuficientă de procesare a datelor	Mic	Mediu	8	Achiziția de servicii de cloud și hardware corespunzătoare; extinderea capacităților de procesare la nevoie.
Probleme cu scalabilitatea sistemului în scenarii comerciale	Mediu	Mare	15	Testarea sarcinii pentru fluxuri de date ridicate și adaptarea parametrilor de cercetare la cerințele de producție.

10 Rezultatele proiectului

10.1 Livrabile

În perioada raportată am furnizat livrabile din tabela 2 și am început lucrul intens asupra celorlalte livrabile.

Tabela 2: Tabel cu livrabile și statusul acestora.

Nr. livrabil	Termen	Livrabil	Status livrabil
D1.3	M6	Planul de management al datelor	Livrat M6
D3.1	M20	Descrierea algoritmilor de detectare a degradării, inclusiv abordarea științifică și starea curentă, precum și rezultatele evaluării - raport	În lucru

10.2 Articole științifice

În perioada de raportare s-a lucrat intens la mai multe articole științifice, unele dintre acestea au fost prezentate la conferințe, sau publicate în jurnale, iar unele urmează să fie prezentate/publicate.

În Tabelul sunt listate toate articolele din cadrul proiectului:

Tabela 3: Lista de articole

Articolul	Detalii Conferința / Jurnal	Link-ul pentru Open-Access
Daniela Delinschi, Rudolf Erdei, Emil Pașca, Oliviu Matei, " Data Quality Assessment Methodology "	19th International Conference on Soft Computing Models in Industrial and Environmental Applications (SOCO 2024)	În curs de publicare
Emil Pașca, Rudolf Erdei, Daniela Delinschi, Oliviu Matei, " Augmenting API Security Testing with Automated LLM-Driven Test Generation "	17th International Conference on Computational Intelligence in Security for Information Systems (CISIS 2024)	Link articol
Oliviu Matei, Rudolf Erdei, Daniela Delinschi, Iulia Bărbăian, Jose Barata, Sanaz Nikghadam-Hojjati " Collaborative networks in orchestration-based software architectures "	The 51st International Conference on Computers and Industrial Engineering (CIE51)	În curs de publicare
Rudolf Erdei, Daniela Delinschi, Emil Pașca, Laura Andreica, Oliviu Matei " Selective Survey of Distributed Learning Methodologies for Agricultural Applications: Challenges and Strategies for Ensuring Privacy and Resilience "	SCIENTIFIC BULLETIN, Seria C: Inginerie Electrică și Știința Calculatoarelor	În curs de publicare

11 Diseminare și exploatare

11.1 Activități de diseminare

Proiectul a fost diseminat în următoarele moduri:

- pe pagina web: <https://research.holisun.com/ro/proiecte/predictive-analysis/gendeg-ro>, având un număr de 180 de vizitatori lunari
- pe contul de LinkedIn: <https://www.linkedin.com/company/holisun>, cu 400 de adepți
- pe pagina de Facebook: <https://www.facebook.com/Holisun.IT/>, având 1881 de urmăritori

Au fost desfășurate o serie de activități de diseminare în cadrul unor evenimente de afaceri, expoziții și evenimente de brokeraj sau networking, listate în Tabelul 4.

Tabela 4: Lista de activități de diseminare

Nume	Data	Link	Participanți	Rezultate
Cluj Innovation Days 2024	21.03.2024-22.03.2024	https://clujinnovationdays.com/	Oliviu Matei	Prezentare <i>GenDeg</i>
BOOSTing European collaboration among Industry 4.0 stakeholders	16.04.2024-26.04.2024	https://boosting-european-collaboration-among-industry.b2match.io/	Rudolf Erdei	Prezentare <i>GenDeg</i>
Clean Energy Transition Partnership	12.09.2024	https://www.b2match.com/e/clean-energy-transition-partnership-2024	Rudolf Erdei Daniela Delinschi	Prezentare <i>GenDeg</i>

12 Concluzii

În acest prim an al proiectului **GenDeg**, HOLISUN și partenerii săi au concentrat eforturile pe studiul și cercetarea modalităților de identificare a degradării performanței modelelor predictive și a derivei acestora. Activitățile au inclus analiza stadiului actual al domeniului, cu accent pe identificarea limitărilor algoritmilor existenți, precum și pe cerințele pentru crearea unor soluții robuste și generalizabile. Scopul principal al acestor analize a fost construirea unei înțelegeri aprofundate a contextului tehnologic și a provocărilor asociate cu utilizarea inteligenței artificiale în medii dinamice și complexe.

În anul următor, HOLISUN își va concentra eforturile pe finalizarea dezvoltării arhitecturii platformei și integrarea algoritmilor avansați de detecție și mitigare a derivei modelelor în fluxuri de date reale. Aceasta va include implementarea unui cadru software pentru monitorizarea continuă a performanței modelelor predictive și generarea de soluții bazate pe machine learning (ML). Validarea soluțiilor tehnice propuse va avea loc în colaborare cu partenerii consorțiului, utilizând seturi de date reale din domeniul financiar și al mentenanței predictive.

Sistemul rezultat va integra mai multe module software și va include funcționalități pentru detectarea timpurie a degradării modelelor, optimizarea performanței predictive și generarea de estimări generalizabile. Planurile de extindere ale proiectului includ adăugarea unor componente inovatoare, cum ar fi algoritmi de ensemble learning și soluții de învățare online, care vor contribui la creșterea eficienței și a rezilienței modelelor predictive în diverse aplicații.

De asemenea, primul an al proiectului a marcat începutul activităților de diseminare și exploatare a rezultatelor. În acest sens, HOLISUN a creat și întreținut website-ul oficial al proiectului, care include informații detaliate despre obiective și progres. În plus, au fost administrate profilele de social media asociate proiectului, pe platforme precum LinkedIn, pentru a facilita colaborarea și schimbul de cunoștințe cu experți din domeniu. Aceste activități au asigurat o diseminare eficientă a rezultatelor, promovând impactul pozitiv al proiectului în domeniul analizei predictive și al rezilienței modelelor AI.

Pe partea științifică, rezultatele obținute în acest prim an au fost diseminate prin prezentarea a două articole la conferințe internaționale de prestigiu. De asemenea, alte două articole au fost trimise spre publicare în jurnale academice de înaltă calitate, contribuind astfel la creșterea vizibilității proiectului în comunitatea științifică.

Prin aceste realizări, proiectul **GenDeg** a făcut progrese semnificative în direcția dezvoltării unor soluții avansate pentru monitorizarea și îmbunătățirea performanței modelelor predictive, demonstrând potențialul său de a adresa provocările complexe ale inteligenței artificiale moderne.

Referințe

- [1] Barbon, S., Baumer, T., Meding, T., Benavides, A.: A framework for human-in-the-loop interpretation of machine learning models. *Artificial Intelligence Review* **51**(2), 103–122 (2018)
- [2] de Barros, R., de Carvalho Santos, S.: An overview and comprehensive comparison of ensembles for concept drift. *Information Fusion* **52**, 213–244 (2019)
- [3] Bayram, F., Ahmed, B., Kassler, A.: From concept drift to model degradation: An overview on performance-aware drift detectors. *Knowledge-Based Systems* **245**, 108632 (2022)
- [4] Delinschi, D., Erdei, R., Pasca, E., Matei, O.: Data quality assessment methodology. In: *The 19th International Conference on Soft Computing Models in Industrial and Environmental Applications SOCO 2024: Salamanca, Spain, October 9–11, 2024 Proceedings, Volume 2*. p. 199. Springer Nature (2024)
- [5] Gama, J., Žliobaite, I., Bifet, A., Pechenizkiy, M., Bouchachia, A.: A survey on concept drift adaptation. *ACM computing surveys (CSUR)* **46**(4), 1–37 (2014)
- [6] Hoens, T., Chawla, N.: Learning from class-imbalanced data: Review of methods and applications. *Expert Systems with Applications* **39**(5), 5712–5721 (2012)
- [7] Lu, J., Liu, A., Dong, F., Gu, F., Gama, J., Zhang, G.: Learning under concept drift: A review. *IEEE Transactions on Knowledge and Data Engineering* **31**(12), 2346–2363 (2019)
- [8] Moharram, H., Awad, A., El-Kafrawy, P.: Optimizing adwin for steady streams. In: *Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing*. pp. 450–459 (2022)
- [9] Pesaranghader, A., Viktor, H.: Fast hoeffding drift detection method for evolving data streams pp. 96–111 (2016)
- [10] Pinagé, F., dos Santos, E., Gama, J.: A drift detection method based on dynamic classifier selection. *Data Mining and Knowledge Discovery* **34**(1), 50–74 (2020)
- [11] Ross, G., Adams, N., Tasoulis, D., Hand, D.: Exponentially weighted moving average charts for detecting concept drift. *Pattern Recognition Letters* **33**(2), 191–198 (2012)
- [12] Shalev-Shwartz, S.: Online learning and online convex optimization. *Foundations and Trends® in Machine Learning* **4**(2), 107–194 (2012)
- [13] Xie, L., Moustakides, G., Xie, Y.: Window-limited cusum for sequential change detection. *IEEE Transactions on Information Theory* **69**(9), 5990–6005 (2023)

