



UNIVERSITATEA DIN
BUCUREȘTI

FACULTATEA DE
MATEMATICĂ ȘI
INFORMATICĂ



SPECIALIZAREA INFORMATICĂ

Lucrare de licență

DETECTAREA TRANZACȚIILOR BANCARE FRAUDULOASE UTILIZÂND ÎNVĂȚARE AUTOMATĂ NESUPERVIZATĂ

Absolvent

Eftimie Petre-Laurențiu

Coordonator științific

Conf. dr. Cristian Rusu

București, iunie 2024

Rezumat

Detecția anomaliilor este un subiect ce apare frecvent în cercetarea din domeniul învățării automate. Acest fapt se datorează multitudinii de aplicații într-o gamă largă de industrii.

Totuși, găsirea unei tehnici bune de rezolvare a problemei este dificilă în cele mai multe cazuri, întrucât este costisitor sau aproape imposibil să obții date despre evenimentele rare ce merită semnalate.

În această lucrare, explorăm un set de date ce conține informații despre tranzacții bancare folosind 3 algoritmi de învățare automată nesupervizată: "One Class SVM", "Gaussian Mixture Model" și "Kernel Density Estimation". Performanța acestora din urmă este comparată prin metricile: "accuracy", "precision", "recall" și "f1 score".

Obiectivul este detectarea tranzacțiilor frauduloase.

Abstract

Anomaly detection is a topic appearing frequently in machine learning research. This is due to its many applications in various industrial settings.

However, finding a good algorithm to solve this problem is difficult in most cases, since obtaining data about rare events that are worth reporting is costly or almost impossible.

In this paper, we analyze a dataset containing information about bank transactions using 3 unsupervised machine learning algorithms: "One Class SVM", "Gaussian Mixture Model" and "Kernel Density Estimation". Their performance is compared using the following metrics: "accuracy", "precision", "recall" and "f1 score".

Our objective is detecting fraudulent transactions.

Cuprins

Capitolul 1

Introducere

1.1 Inteligența artificială

Este un domeniu de studiu din Informatică, discutat încă din prima jumătate a secolului al XX-lea, dar care a devenit popular abia la finalul secolului menționat anterior.

Aceasta are ca scop simularea inteligenței umane folosind diverși algoritmi și structuri de date. Există mai multe arii de cercetare în acest domeniu, precum: raționamentul, reprezentarea cunoștințelor, planificarea, procesarea limbajului natural, cât și învățarea automată care va fi prezentată detaliat în această lucrare. Totuși, obiectivul pe termen lung al cercetătorilor este de a simula inteligența generală, nu doar părți specifice din aceasta.

Popularitatea acesteia a început să crească odată cu mărirea spațiului de stocare, a puterii de procesare, cât și cu faptul că sistemele de calcul au devenit mai accesibile.

De asemenea, deși algoritmi existenți nu s-au mai îmbunătățit semnificativ, accesul la o cantitate enormă de date în ziua de astăzi, anume fenomenul "Big Data", a oferit un sprijin considerabil pentru performanța învățării inteligenței artificiale[**ai-history**].

În acest domeniu, datele colectate, cât și calitatea lor au un rol extrem de important. Oricât de bun ar fi algoritmul propus, în lipsa unor date reprezentative pentru problema ce trebuie rezolvată, performanța va fi îndoielnică. Putem spune că datele sunt cele care introduc particularitățile din lumea reală în spațiul teoretic al inteligenței artificiale.

Prin urmare, se observă că aplicabilitatea inteligenței artificiale a fost împiedicată în principal de considerente practice, nu neapărat teoretice.

Printre succesele recente din domeniu se remarcă AlphaFold, un model inventat de DeepMind care prezice structura proteinelor din secvențele de amino acizi, GPT-3 al OpenAI ce poate genera text care să răspundă la întrebările utilizatorilor aproape ca o ființă umană, dar și algoritmi pentru mașini autonome dezvoltați de cei de la Tesla.

Utilizarea inteligenței artificiale poate reduce costurile asociate sarcinilor ce ar necesita efortul unei persoane, dar de multe ori apare și problema eticii, precum în cazul înlocuirii

recrutorilor cu un program ce selectează candidații pentru un loc de muncă în funcție de anumite criterii care nu ar fi întocmai corespunzătoare.

Automatizarea detecției anomaliilor reprezintă una din posibilele aplicații ale inteligenței artificiale, subiect ce va fi abordat în această lucrare.

1.2 Definiția anomaliilor

O **anomalie** este o entitate ce diferă semnificativ de restul entităților din setul de date. Definiția lui Hawkins este următoarea [hawkins1980identification]: *"O anomalie este o observație ce deviază atât de mult față de restul observațiilor, încât să creeze suspiciunea că a fost generată de un mecanism diferit"*.

Datele sunt colectate prin observarea unor procese, de preferat din viața reală, precum funcționarea unui motor sau traficul pe o rețea de calculatoare. În majoritatea timpului, procesele generează date ce corespund unei derulări normale, dar în anumite cazuri acestea se comportă anormal și astfel apar anomaliile.

"Normalitatea" observațiilor este definită de oameni. Algoritmii de învățare automată pot să semnaleze potențialele abateri, dar în final, decizia rămâne a omului.

De asemenea, caracteristicile anomaliilor depind de domeniul unde acestea apar. Ce este considerat anomalie în domeniul medical, va diferi aproape complet de domeniul finanțelor, spre exemplu.

Mai mult, unele observații nici măcar nu pot reprezenta comportament anormal în afara unui context bine definit sau fără a apărea în prezența altor observații care adunate ar forma o anomalie. Un exemplu de comportament ce poate fi observat în traficul de date pe o rețea al unei organizații ar fi conectarea de la distanță pe un anumit calculator. Acest lucru se întâmplă regulat și nu poate ridica nicio suspiciune de cele mai multe ori. Dar în combinație cu o oră târzie de accesare, când probabil angajații nu mai sunt la lucru, sau cu alte comenzi care sunt folosite succesiv pentru a comite un atac cibernetic, o observație ce în general este complet normală, acum prezintă caracteristici anormale.

Prin urmare, se observă că decizia de a eticheta o observație ca fiind anomalie sau nu este una dificilă chiar și pentru un om, iar în combinație cu faptul că obținerea de date ce ar putea fi anormale este de cele mai multe ori costisitoare, face ca evaluarea sistemelor de detecție să devină una problematică.

1.3 Avantajele învățării automate

Când sarcinile ce trebuie îndeplinite periodic, necesare bunei funcționări a unei companii spre exemplu, încep să ocupe o mare parte din timpul oamenilor, apare nevoia automatizării.

În general, soluțiile ce includ stabilirea unui set de reguli sunt de preferat, datorită simplității lor. Un exemplu banal ar fi verificarea siguranței unei parole introduse de utilizator la crearea unui cont pe o platformă web. Este suficient să verificăm dacă parola are un număr minim precizat de caractere sau dacă conține anumite simboluri ce reduc vulnerabilitatea în fața unui atac cibernetic.

Din păcate, multe din problemele cu care ne confruntăm nu au soluții simple. Aici intervine învățarea automată ce încearcă să aproximeze posibilul set de reguli, prea complex pentru a fi definit, ce ar putea rezolva problema. Am putea spune că această tehnică ”învăță” setul de reguli din datele puse la dispoziție, de unde și denumirea. Evident, o aproximare nu are cum să fie perfectă, dar este suficient de bună încât să scutească oamenii de o mare parte din treabă. În cazul tranzacțiilor bancare, ce probabil sunt efectuate în număr de sute de milioane pe zi, nu este greu de observat că verificarea manuală de către oameni este imposibilă, iar găsirea unui set simplu de reguli care să le caracterizeze din nou este o problemă dificilă. Totuși, învățarea automată ne poate ajuta să minimizăm fraudele ce trec nedetectate de la milioane la câteva zeci pe zi.

1.4 Caracterizarea anomaliilor nu este viabilă

Un exemplu ce ilustrează de ce obținerea de informații despre anomalii este costisitoare ar putea fi activitatea unui motor defect.

În acest caz, avem 2 variante:

- Una dintre ele este să **simulăm** într-un fel sau altul comportamentul unui motor defectuos, dar se observă că această operațiune este greu de realizat și chiar dacă am reuși să o ducem la capăt, datele colectate nu ar fi autentice și ar putea duce la o reprezentare incorectă a anomaliilor.
- A doua variantă implică **sabotarea intenționată** a motorului pentru a obține datele, dar acest lucru produce un cost prea mare de cele mai multe ori.

Este evident că definirea anomaliei în acest context nu este viabilă, iar cazul nu este singular. Tranzacțiile frauduloase ce constituie tema acestei lucrări suferă de aceleași probleme, ba mai mult, nici măcar nu avem opțiunea de a sabota intenționat derularea normală a tranzacțiilor. De asemenea, dacă anomaliile apar ca urmare a unor acțiuni cu caracter malițios, atunci adversarii responsabili pentru acestea vor încerca să evite mecanismele existente de detecție prin schimbarea frecvență a metodelor folosite. Prin urmare, ce era considerat comportament normal înainte, acum poate indica prezența unui atac.

În schimb, este mult mai ușor și cel mai probabil implică un cost aproape inexistent, să urmărim activitatea unui motor în stare bună de funcționare și să folosim datele respective pentru a defini ce înseamnă o observație normală.

1.5 Aplicații ale detecției anomaliilor

În aproape toate domeniile apare problema semnalării unor fenomene sau evenimente ieșite din comun ce necesită atenția unei persoane pentru analiză. Mai jos prezentăm doar câteva din cazurile unde detecția anomaliilor este necesară, cât și ce ar reprezenta o anomalie pentru fiecare:

- **Fraude financiare:** ne dorim să semnalăm comportamentul ciudat observat într-o serie de tranzacții, precum activitatea generată de o persoană care folosește cardul de credit al altei persoane în mod neautorizat. Această aplicație constituie și obiectul lucrării noastre[**financial-fraud**].
- **Intruziune în rețele de calculatoare:** vrem să semnalăm activitatea neobișnuită ce poate indica un potențial atac cibernetic sau accesul neautorizat al unui terț malițios[**network-traffic**].
- **Controlul calității în manufactură:** suntem interesați să monitorizăm procesul de producție pentru a raporta eventualele defecțiuni ce ar afecta calitatea produsului[**quality-control**].
- **Domeniul medical:** vrem să identificăm anomalii în analizele de sânge, semnele vitale sau imaginile medicale ale pacientului pentru a depista și preveni anumite afecțiuni[**medical-images**].
- **Rețele sociale:** ne dorim să raportăm automat conținutul postat ce poate include remarci jignitoare, precum comentariile rasiste sau xenofobe, astfel fiind de ajutor în moderarea materialului încărcat pe platforme[**social-media**].

1.6 Anomalii în tranzacții bancare

Tranzacțiile fraudulente pot produce pagube financiare majore atât instituțiilor responsabile de operațiunile economice efectuate prin intermediul lor, cât și celor care apelează la respectivele servicii.

Pe lângă prejudiciul material adus, fraudele nedetectate deteriorează și reputația băncilor ce cad victimă acestui fenomen, iar pe termen lung acestea pot duce la pierderea potențialilor clienți.

De asemenea, în urma unor operațiuni frauduloase ce nu au fost prevenite, apar și tulburări ale activității economice desfășurate, deoarece trebuie investite resurse și timp pentru investigarea și soluționarea problemei. Alte costuri suplimentare ce mai pot apărea reprezintă procese juridice îndreptate împotriva instituției financiare afectate, cât și sancțiuni pentru nerespectarea normelor de siguranță.

Numai Marea Britanie a înregistrat pierderi de peste 1.2 miliarde de lire sterline în 2022 cu aproape 80% din fraude originând din mediul online, conform *UK Finance* [uk-finance].

Prin urmare, este nevoie de o soluție care să poată preveni măcar o mare parte din tranzacțiile frauduloase dacă nu chiar pe toate și care să nu necesite verificarea manuală a unui om decât pentru un număr mic de evenimente.

1.7 Detecția anomaliilor vs clasificare clasică

1.7.1 Greșeală comună

La prima vedere, aceste două probleme par să coincidă și ne determină să ne întrebăm de ce a mai apărut un nou domeniu de cercetare, anume detecția anomaliilor, când deja avem la dispoziție atâtea rezultate utile pentru problema clasificării. Din păcate, dacă nu cunoaștem numărul de clase ce pot apărea sau nu avem un eșantion reprezentativ pentru toate tipurile de clase, fapt ce se întâmplă deseori în practică, nu putem utiliza un clasificator clasic. De exemplu, în cazul tranzacțiilor bancare, nu putem ști cu certitudine care sunt toate tipurile de fraudă ce pot apărea și nici nu avem un număr suficient de exemple pentru a defini măcar tipurile de fraudă pe care le-am identificat până acum. Prin urmare, este mult mai ușor să modelăm problema în jurul unei singure clase de referință, decât să încercăm definirea unei clase adiționale ce de fapt este compusă dintr-o multitudine de subclase necunoscute.

De asemenea, problemele de clasificare de cele mai multe ori presupun accesul la etichetele asociate fiecărei observații, deci metodele supervizate sunt cele mai potrivite aici. Totuși, și dacă setul de date este etichetat, lucru ce se întâmplă în cel mai bun caz, un clasificator clasic ar avea o performanță slabă pentru problema detecției anomaliilor din cauza dificultății modelării claselor, după cum am menționat și mai sus.

De cele mai multe ori nu vom avea etichetele datelor colectate, sau cel puțin nu un număr suficient de mare, pentru detecția anomaliilor, așa că apare nevoia metodelor nesupervizate. În realitate, totuși, este posibil ca în timp să acumulăm măcar câteva etichete, chiar dacă la început soluția este pur nesupervizată, pe care le putem exploata folosind metode semi-supervizate care în general au o performanță mai bună decât cele menționate anterior.

1.7.2 Scopul detecției anomaliilor

Diferența constă în faptul că detecția anomaliilor are la bază **o singură clasă de referință**. Observația fie aparține acestei clase, fie aparține oricărei alte clase care este diferită de aceasta din urmă. Multitudinea de "clase diferite" ne indică faptul că este greu

să definim ce înseamnă ”*diferit de normal*”, dar este relativ ușor să definim ce înseamnă ”*normal*”.

Un exemplu trivial ar putea fi să detectăm dacă animalul din imagine este un câine sau nu. În acest caz, normal înseamnă câine, iar anormal înseamnă orice alt animal care nu este câine. Se observă că este relativ ușor să caracterizăm conceptul de câine, în timp ce definirea conceptului de ”diferit de câine” este complexă, fapt ce ar pune în dificultate un clasificator clasic care prin definiție are nevoie de un număr cunoscut de clase cu observațiile reprezentative aferente.

1.7.3 Scopul clasificării clasice

La clasificarea clasică, clasele ce trebuie identificate sunt bine definite și de cele mai multe ori, trăsăturile lor se suprapun în mai multe locuri. De asemenea, această problemă acordă o importanță **egală** tuturor categoriilor implicate, pe când eșecul de a semnala o anomalie este în general mult mai dăunător față de detectarea unei observații normale ca fiind anormală.

Pentru a continua exemplul precedent cu imaginile cu animale, putem reformula problema ca de această dată să diferențiem între pisici și câini. Deja se vede că ambele clase sunt mult mai bine definite și putem găsi atât trăsături comune, precum nasul, cât și trăsături definitorii, precum mustățile. Vrem să detectăm la fel de bine ambele animale, pe când într-o problemă clasică de detecție a anomaliilor, precum identificarea tranzacțiilor frauduloase, vrem preponderent să nu ratăm evenimentele cu caracter malițios.

1.8 Probleme conexe pentru detecția anomaliilor

1.8.1 Outlier detection

Outlier detection are scopul de a identifica observații fie noi, fie previzibile, din setul de antrenare, care deviază de la comportamentul normal. **Outliers** pot fi valori extreme, erori cauzate de zgomot, posibil din instrumentele de măsură sau senzorii folosiți pentru colectarea datelor, sau comportament anormal.

Această abordare este utilă atunci când avem un set de date ”**poluat**” cu observații anormale și ne dorim să extragem din el doar porțiunea ce conține observații deviante. De asemenea, se poate folosi atât pe seturi adnotate, cât și pe seturi fără etichete.

1.8.2 Novelty detection

Novelty detection are scopul de a identifica observații noi ce diferă semnificativ de datele de antrenare care au rolul de a caracteriza cât mai bine clasa normală. Din acest motiv, este foarte important ca setul de antrenare să nu conțină decât date normale pentru

a putea identifica caracteristicile clasei de referință, astfel modelând comportamentul normal.

Prin urmare, această abordare este utilă atunci când avem un set de date ”**curat**” (fără anomalii) și ne așteptăm ca observațiile noi să fie rare și să aibă trăsături distincte față de setul de antrenare. De asemenea, se presupune că datele sunt adnotate pentru a putea extrage punctele cu eticheta normală și să le folosim la antrenare.

Novelty detection este problema pe care o vom studia detaliat în această lucrare.

1.8.3 Change detection

Change detection este folosit atunci când dorim să găsim anomalii într-o serie de timp, anume să analizăm dacă comportamentul seriei se schimbă semnificativ. De obicei, suntem interesați ori doar să aflăm dacă o schimbare a apărut, ori să aflăm pozițiile în timp unde schimbările respective au avut loc.

Aici se remarcă 2 categorii, **offline**, unde presupunem că avem la dispoziție o serie de timp de lungime finită pe care o analizăm cu scopul de a detecta puncte anormale, și **online**, unde avem un flux de date din care vrem să extragem, în timp real, punctele problematice. Deși este probabil impropriu, putem considera pentru seriile de timp, intuitiv, că outlier detection este similar cu cazul offline din change detection, iar novelty detection este asemănător cu cel online din această problemă.

În această lucrare, nu vom aborda problema change detection.

Capitolul 2

Preliminarii

2.1 Tipuri de învățare automată

Există numeroase tipuri de învățare automată, așa că le vom prezenta doar pe cele relevante lucrării noastre. De asemenea, notațiile pentru setul de date introduse în această secțiune vor fi folosite și în restul lucrării.

2.1.1 Supervizată

Formal, avem un set de date cu N elemente, sub formă de perechi:

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \quad (2.1)$$

unde observațiile $x_i \in \mathbb{R}^d$ și etichetele $y_i \in \mathbb{R}$ sau $y_i \in \{0, 1\}$. Problema se reduce la aproximarea cât mai bună a etichetei y_{N+1} pentru o observație nouă x_{N+1} .

Învățarea **supervizată** poartă acest nume deoarece necesită supraveghere umană pentru a putea funcționa. Un specialist în știința datelor trebuie să parcurgă fiecare observație din setul de date și să îi atribuie o **etichetă** corespunzătoare. Acest lucru este evident dificil, având în vedere că seturile de date cu minim sute de mii de entități sunt des întâlnite.

Această metodă este folosită în probleme de clasificare sau de prezicere a unor fenomene.

2.1.2 Nesupervizată

Formal, avem un set de date cu N elemente:

$$S = \{x_1, x_2, \dots, x_N\} \quad (2.2)$$

unde observațiile $x_i \in \mathbb{R}^d$, dar nu mai apar etichetele y_i . În lipsa etichetelor, problema

generalizată nu mai poate fi descrisă la fel de ușor.

Învățarea **nesupervizată** este opusul celei supervizate, deci implică faptul că interacțiunea umană nu este necesară în pregătirea setului de date. Cu toate acestea, rezultatele trebuie să fie interpretate de o persoană pentru a fi relevante, întrucât nu mai avem etichete pe care să le folosim în evaluarea automată a performanței.

Această metodă este folosită pentru a grupa datele în funcție de similaritate, a înțelege relația dintre punctele din setul de date și pentru a face o analiză inițială a datelor.

Toți algoritmi din această lucrare aparțin acestei metode de învățare automată.

2.1.3 Semi-supervizată

Formal, avem un set de date cu N elemente, partiționat în 2 mulțimi cu N_1 elemente sub formă de perechi, respectiv N_2 elemente:

$$S_1 = \{(x_1, y_1), (x_2, y_2), \dots, (x_{N_1}, y_{N_1})\} \quad (2.3)$$

$$S_2 = \{x_1, x_2, \dots, x_{N_2}\} \quad (2.4)$$

unde observațiile $x_i \in \mathbb{R}^d$ și etichetele $y_i \in \mathbb{R}$ sau $y_i \in \{0, 1\}$. La fel ca în cazul supervizat, problema este în general aproximarea cât mai bună a etichetei y_{N+1} pentru o observație nouă x_{N+1} .

Învățarea **semi-supervizată** îmbină ambele paradigme prezentate anterior, astfel că necesită un număr mic de date adnotate, lucru ce este evident mai ușor de obținut față de un întreg set, și un număr mare de date fara etichetă.

Această metodă este folosită spre exemplu în modelele ce se antrenează singure, folosind un algoritm supervizat antrenat pe datele adnotate ce este apoi folosit pe datele fara etichetă pentru a obține un nou set de date adnotat.

Deși această metodă nu este folosită în lucrarea noastră, am inclus-o deoarece ideea de bază în novelty detection cu metode nesupervizate este similară. Nu avem nevoie de date adnotate la antrenare și este suficient un set mic pentru testarea ulterioară a performanței.

2.2 One Class SVM

2.2.1 Ideea algoritmului

Această metodă este inspirată din clasificatorul cu vectori suport. Ideea este să găsim un **hiperplan cu margine maximă**, posibil într-un spațiu cu mai multe dimensiuni decât cel inițial, în funcție de kernel, care să separe originea (se presupune că punctele sunt centrate) spațiului de trăsături de restul punctelor din setul de date[scholkopf2000support].

Un alt mod de a privi algoritmul este găsirea celei mai **mici hipersfere** care să includă toate punctele din setul de date[**tax2004support**].

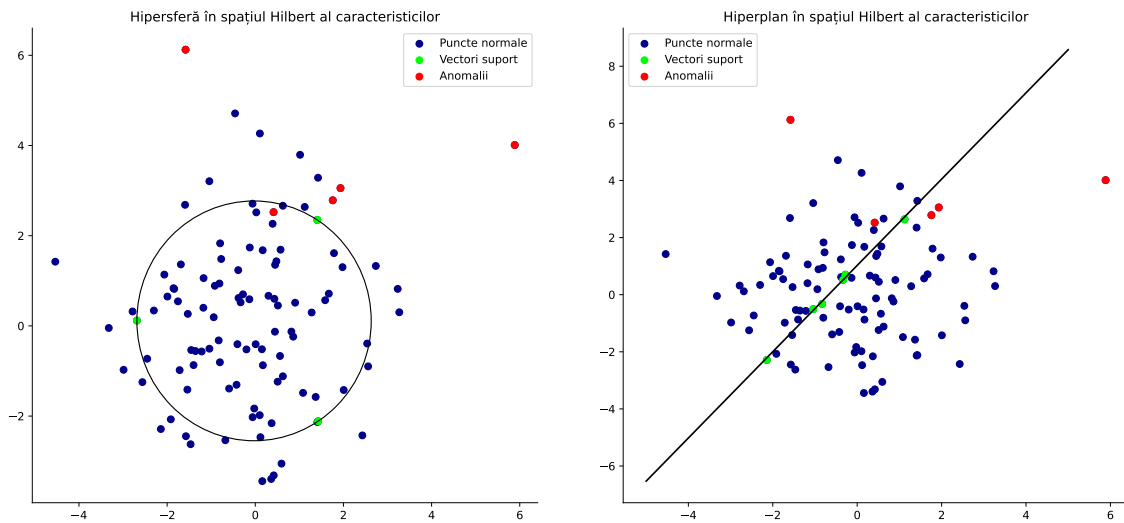


Figura 2.1: Varianta Schölkopf et al în dreapta și varianta Tax et al în stânga

2.3 Caracteristici

One Class SVM ne ajută să transformăm un model de clasificare cu mai multe clase, anume SVM, într-unul cu o singură clasă, păstrând abilitatea de a putea introduce ne-liniarități cu ajutorul funcțiilor de kernel ce au fost studiate extensiv. Prin urmare, are majoritatea avantajelor celui din urmă, precum garanția optimului global datorită funcției convexe pe care o minimizăm, faptul că este relativ robust în spații cu multe dimensiuni și că poate fi folosit chiar și dacă numărul de caracteristici al punctelor din set este mai mare decât cardinalul setului.

Totuși, acest model devine imposibil de folosit pentru seturi mari de date, complexitatea de timp lasă de dorit la antrenare, iar găsirea hiperparametrilor optimi este dificilă, atât din lipsa unor reguli stricte după care să ne ghidăm, dar și din faptul că o căutare exhaustivă ar dura prea mult din cauza complexității de timp ridicate la antrenare. De asemenea, rezultatele date de acest model nu pot fi ușor interpretate, deoarece ele reprezintă distanțe, nu probabilități.

2.3.1 Formularea matematică

Prima formulare ce implică un hiperplan de separare

$$\begin{aligned}
& \min_{w, \rho, \xi} \quad \frac{1}{2} \|w\|^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i - \rho \\
& \text{cu constrângerea} \quad \langle w, \phi(x_i) \rangle \geq \rho - \xi_i, \quad i = 1, 2, \dots, n \\
& \quad \xi_i \geq 0, \quad i = 1, 2, \dots, n
\end{aligned} \tag{2.5}$$

- w este vectorul de pondere al hiperplanului
- ρ este termenul de influență
- ξ_i sunt variabilele de relaxare pentru încălcarea marginii
- $\phi(x_i)$ este funcția de scufundare pentru x_i .
- n este numărul total de puncte
- ν este marginea superioară pentru ponderea de anomalii și marginea inferioară pentru ponderea de vectori suport

împreună cu forma sa duală ce implică folosirea unei funcții kernel pentru găsirea unui hiperplan de separare într-un spațiu cu mai multe dimensiuni decât cel inițial

$$\begin{aligned}
& \min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j) \\
& \text{cu constrângerea} \quad 0 \leq \alpha_i \leq \frac{1}{\nu n}, \quad i = 1, 2, \dots, n \\
& \quad \sum_{i=1}^n \alpha_i = 1
\end{aligned} \tag{2.6}$$

- α_i sunt variabilele duale asociate punctelor
- $K(x_i, x_j)$ este funcția kernel

A doua formulare ce implică găsirea hipersferei minime

$$\begin{aligned}
& \min_{R, c, \xi} \quad R^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i \\
& \text{cu constrângerea} \quad \|\phi(x_i) - c\|^2 \leq R^2 + \xi_i, \quad i = 1, 2, \dots, n \\
& \quad \xi_i \geq 0, \quad i = 1, 2, \dots, n
\end{aligned} \tag{2.7}$$

- R este raza hipersferei
- c este centrul hipersferei

2.3.2 Tehnici de optimizare

Toate funcțiile prezentate mai sus sunt convexe, iar constrângerile lor sunt mulțimi convexe. Chiar și prima constrângere din formularea cu hipersferă care conține un termen la pătrat în dreapta inegalității poate fi redusă la o constrângere liniară introducând o nouă variabilă $R_2 = R^2$ cu constrângerea $R_2 \geq 0$.

Deci știm că problemele de optimizare sunt convexe și astfel avem garanția unui unic minim global. Prin urmare, putem folosi tehnicile populare de optimizare convexă.

O metodă des întâlnită de ordinul întâi este coborârea pe gradient ce este avantajoasă atunci când setul de antrenare este mare, întrucât complexitatea algoritmului nu depinde de mărimea setului de date. Totuși, această tehnică întâmpină dificultăți la optimizarea funcțiilor cu constrângeri. Există modificări ale tehnicii de bază, precum coborârea pe gradient proiectată care poate lucra și cu constrângeri liniare simple sau algoritmul Pegasos[**Pegasos**] care poate rezolva varianta primitivă a primei forme prezentate mai sus, dar pentru formulările complexe nu se poate aplica această metodă.

În majoritatea cazurilor, un alt kernel decât cel liniar va fi folosit, și prin urmare vor apărea următoarele dificultăți: calcularea valorilor funcției kernel este o operație costisitoare, calcularea matricei kernel este de cele mai multe ori inutilă, întrucât ne interesează doar calculele ce implică un vector suport, și matricea kernel de cele mai multe ori nu va încăpea în memorie din cauza numărului excesiv de exemple de antrenare, având în vedere că dimensiunea ei este de $N \times N$ [**SVM-solvers**].

Astfel, apare nevoia unor algoritmi specializați. Cel mai popular dintre ei este SMO[**SMO**] care folosește tehnici de programare pătratică și care utilizează o descompunere în subprobleme mai mici ce pot fi rezolvate eficient și care încap în memorie. O implementare a acestei metode se poate găsi în biblioteca LIBSVM[**LIBSVM**].

2.4 Gaussian Mixture Model

2.4.1 Ideea algoritmului

Este prea restrictiv să presupunem ca fiecare punct dintr-un set de date provine din aceeași distribuție unimodală. Prin urmare, a apărut tehnica de bază "**Mixture model**" ce își propune să elimine presupunerea anterioară pentru a putea modela și distribuții multimodale care poate nu provin nici măcar din aceeași familie, utilizând o sumă ponderată de mai multe componente cu proprietăți cunoscute. Totuși, în practică, se folosesc componente din aceeași familie pentru modelarea distribuției doar că fiecare componentă are parametri diferiți.

Algoritmul încearcă să estimeze funcția densitate de probabilitate din care au fost generate datele folosind **mai multe distribuții Gaussiene**. Astfel, putem modela distribuții multimodale utilizând o distribuție bine cunoscută. Parametrii necesari sunt mediile,

matricele de covarianță și ponderile fiecărei componente.

Metoda este similară cu **k-means**, întrucât ambele folosesc măsuri de similaritate pentru a crea grupuri ce modelează datele. Totuși, față de k-means, Gaussian Mixture Model permite unui punct să aparțină mai multor grupuri, oferind probabilități pentru fiecare, și poate modela chiar și seturi de date complexe cu margini de decizie neliniare.

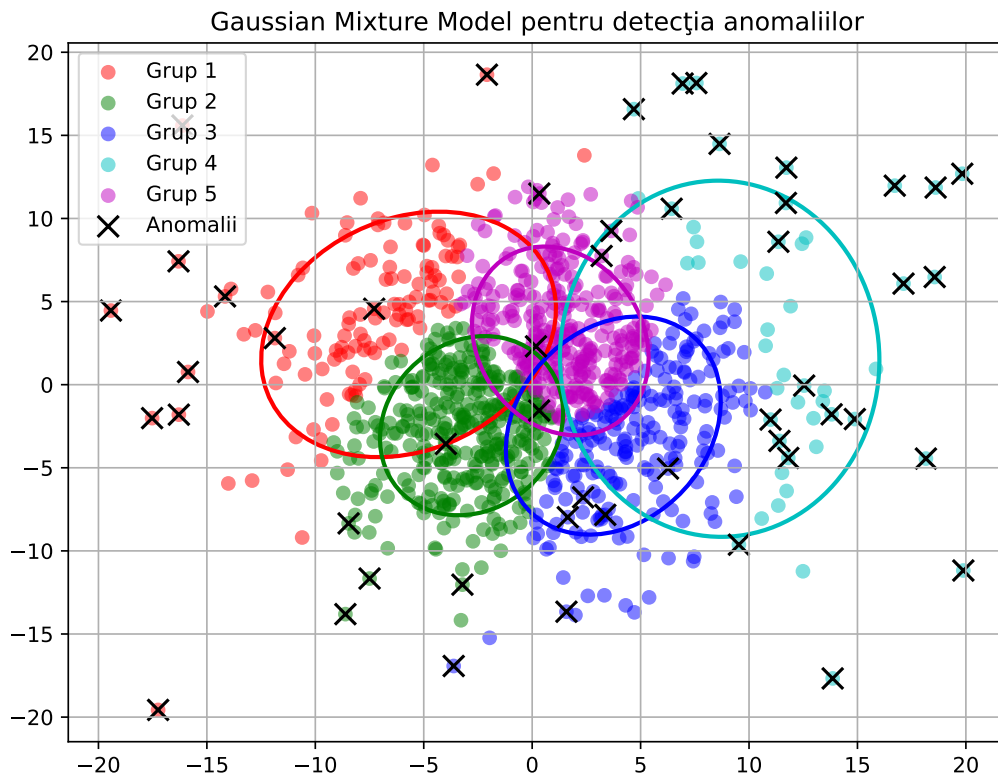


Figura 2.2: Detecția anomaliilor cu 5 componente Gaussiene

2.5 Caracteristici

Gaussian Mixture Model este util pentru modelarea distribuțiilor complexe și se descurcă bine chiar și când setul de date este compus din mulțimi suprapuse de puncte, întrucât află pentru fiecare componentă probabilitatea ca o observație să fie membra acesteia, deci chiar dacă un punct se află în mai multe mulțimi, nu va fi repartizat doar în una singură. De asemenea, timpul de antrenare este redus chiar și pentru seturi mari de date datorită tehnicii iterative de optimizare utilizată.

Din păcate, numărul de componente nu este ușor de ales pentru a obține un model optim, iar performanța algoritmului depinde foarte mult de inițializarea parametrilor inițiali, ceea ce poate duce la imposibilitatea convergenței către optim. De asemenea,

se presupune că datele pot fi modelate de o combinație de distribuții Gaussiene, dar dacă acest lucru nu este adevărat, algoritmul nu va oferi rezultate bune.

Pentru reducerea numărului de parametri ce trebuie învățați, se poate face presupunerea că caracteristicile datelor sunt independente unele de altele. Acest lucru este echivalent cu presupunerea că valorile ce nu se află pe diagonala principală a matricelor de covarianță pot fi considerate 0, dar în practică, datele pot fi corelate de-a lungul diferitelor dimensiuni, ceea ce ar conduce la învățarea unui număr de d^2 parametri pentru fiecare componentă, unde d este dimensiunea unui punct. Acest lucru ar putea duce la overfitting când setul de date este mic[[aggarwal2017outlier](#)]. Un alt dezavantaj este vulnerabilitatea la blestemul dimensionalității.

2.5.1 Formularea matematică

Funcția densitate de probabilitate estimată este dată de

$$p(x) = \sum_{i=1}^K \phi_i \mathcal{N}(x|\mu_i, \Sigma_i) \quad (2.8)$$

- K este numărul de componente Gaussiene
- $\mathcal{N}(x|\mu_i, \Sigma_i)$ este distribuția Gaussiană cu medie μ_i și matrice de covarianță Σ_i
- ϕ_i este ponderea componentei i

Totuși, pentru a detecta anomalii avem nevoie de

$$p'(x) = 1 - \prod_{i=1}^K (1 - p_i(x)) \quad (2.9)$$

- $p_k(x)$ este probabilitatea ca punctul x să aparțină componentei k

Această formulă ne indică probabilitatea ca un punct să fi fost generat de oricare dintre componentele Gaussiene implicate. Deci, o valoare cât mai mică sugerează o șansă mare ca punctul să fie anomalie.

2.5.2 Tehnici de optimizare

Pentru Gaussian Mixture Model nu se poate găsi nicio soluție în formă închisă din cauza faptului că valorile optime ale parametrilor depind de probabilitatea datelor de a aparține fiecărei componente, dar în același timp, calcularea acestor probabilități depinde de parametrii optimi.

Prin urmare, este nevoie de o soluție iterativă care aproximează alternativ valorile căutate mai sus. Cea mai populară metodă este algoritmul **Expectation-Maximization**

(EM). Acesta este compus din 2 pași care sunt repetați până când eroarea față de soluția optimă scade sub un anumit prag. La pasul E , folosim valorile actuale ale parametrilor pentru a calcula probabilitatea că un punct a fost generat de o anumită componentă, făcând acest lucru pentru toate punctele și toate componentele. La pasul M , utilizăm probabilitățile calculate la pasul E pentru a estima valorile parametrilor optimi cu metoda **Maximum Likelihood Estimation** care caută să maximizeze probabilitatea de observare a datelor sub un anumit model statistic.

Numărul de componente poate fi ales ori folosind anumite criterii, precum **Bayesian information criterion** și **Akaike information criterion**, ori folosind cunoștințe din domeniul problemei.

Pentru inițializarea parametrilor se poate folosi algoritmul k-means. Acesta ne va da direct valori inițiale pentru mediile grupurilor. Pentru matricele de covarianță putem folosi covarianțele calculate pentru punctele din fiecare grup, iar pentru ponderi vom folosi proporția de puncte din setul de date aparținând fiecărui grup[**EM-GMM-INIT**].

2.6 Kernel Density Estimation

2.6.1 Ideea algoritmului

Precum Gaussian Mixture Model, algoritmul încearcă să estimeze funcția densitate de probabilitate din care au fost generate datele.

Tehnica de bază utilizează un kernel funcție de distribuție probabilitate pe care ”îl plasăm” la fiecare punct din setul de date, căruia îi oferim o pondere de $\frac{1}{n}$, unde n este numărul de observații. Apoi, distribuția adevărată este aproximată prin adunarea tuturor rezultatelor precedente. Procesul este similar aflării ariei de sub graficul unei funcții cu ajutorul unei integrale ce apare dintr-o sumă de valori ale funcției într-o infinitate de puncte.

Altă tehnică similară de estimare a funcției densitate de probabilitate este cea a **histogramelor** care se bazează pe discretizarea datelor în coșuri de mărime fixă. Totuși, acest lucru irosește informația despre locațiile individuale ale punctelor, înlocuindu-le cu intervale ce corespund mai multor puncte. Astfel, graficul funcției va deveni discontinuu și constant pe fiecare interval. Kernel Density Estimation produce un grafic neted de cele mai multe ori și oferă o reprezentare mai bună a distribuției unei variabile continue[**KDE_paper**].

Un parametru numit ”**lățime de bandă**” influențează netezimea distribuției rezultate. Cel mai des utilizat kernel este cel Gaussian și pe acesta îl vom folosi și noi.

Aici, pentru fiecare punct generăm o distribuție Gaussiană cu **media** egală cu punctul respectiv și **deviație** egală cu ”lățimea de bandă”. Apoi, adunăm toate distribuțiile obținute mai sus și le împărțim la numărul total de puncte.

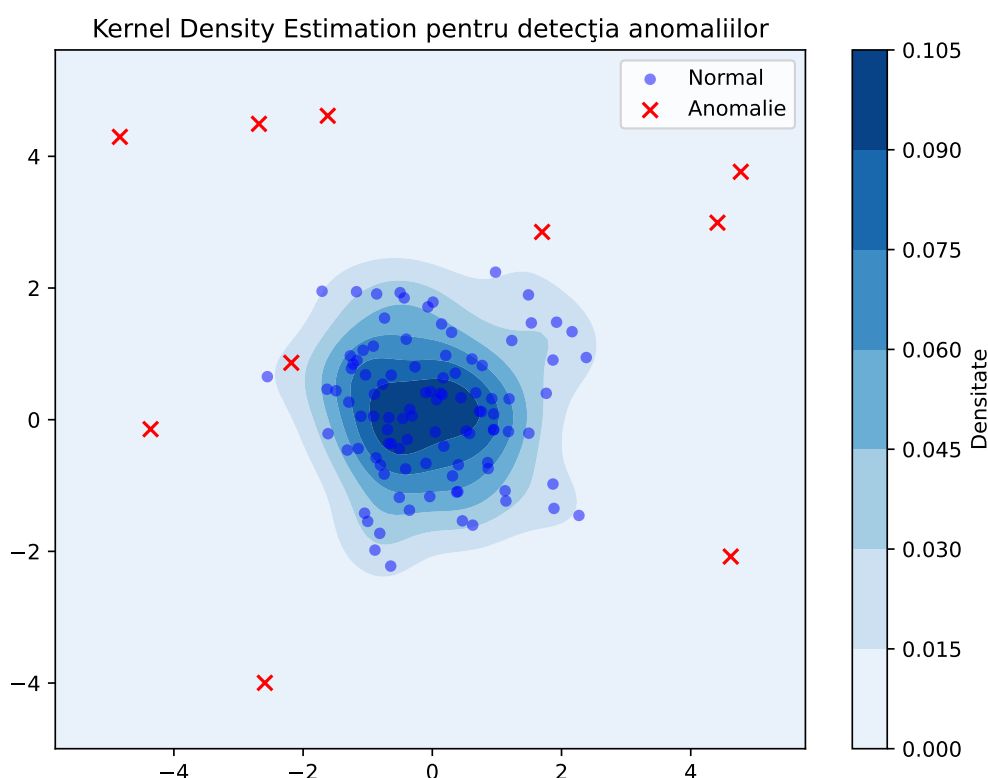


Figura 2.3: Anomaliile sunt reprezentate de punctele roșii

2.7 Caracteristici

Kernel Density Estimation are avantajul de a fi o metoda neparametrică și care nu face nicio presupunere asupra distribuției reale a datelor, ceea ce îl face flexibil pentru diverse tipuri de date. De asemenea, rezultatele sunt ușor de interpretat acestea fiind probabilități și beneficiem de informație locală despre densitatea datelor în jurul fiecărui punct din set, ceea ce ne ajută la găsirea regiunilor cu densitate ridicată sau scăzută.

Dezavantajele sunt găsirea dificilă a hiperparametrului optim pentru lățime de bandă, complexitatea de timp crescută, în special pentru seturi de date mari, vulnerabilitatea la blestemul dimensionalității, dar și tendința de a atribui o densitate ridicată punctelor din marginea setului de date, chiar dacă distribuția nu se mai întinde în direcția respectivă, iar acest lucru poate duce la netezirea excesivă a marginilor distribuției.

2.7.1 Formularea matematică

Densitatea estimată de kernel într-un punct x este dată de

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right) \quad (2.10)$$

- n este numărul total de puncte
- h este lățimea de bandă
- $K(u)$ este funcția kernel

2.7.2 Tehnici de optimizare

Kernel Density Estimation este un model din familia celor care se bazează pe utilizarea unui număr finit dintre cei mai apropiați vecini ai unui punct pentru a lua o decizie. Pentru a cuantifica cât de apropiate sunt 2 puncte, se pot folosi oricare dintre metricile populare, precum distanța euclidiană, sau funcțiile kernel în cazul acestui model, spre exemplu.

Apartenența la această familie de algoritmi implică nevoia de a memora setul de date pe care este antrenat, etapa de învățare, practic, nefiind existentă. Totuși, o abordare cu forță brută ar necesita calcularea funcției kernel între toate perechile de puncte memorate și cele pentru care dorim să aflăm densitatea de probabilitate, adică o complexitate de timp de $O(NM)$, unde N este numărul de puncte de antrenare și M este numărul de puncte evaluate.

O îmbunătățire ar putea fi utilizarea unei structuri de date ce oferă căutări rapide ale punctelor într-un anumit spațiu, precum **KD Trees** și **Ball trees**, aceasta din urmă fiind adesea utilizată pentru date cu dimensiuni $d < 20$. Pentru seturi de date mici sau pentru cazurile în care ne dorim să evaluăm un număr M mic de puncte, este mai util să folosim forța brută, întrucât crearea arborilor din KD Trees și Ball trees poate dura mai mult decât simpla calculare a funcției kernel între toate perechile. În restul cazurilor însă, timpul de preprocesare devine amortizat de evaluările ce se pot face în timp $O(M \log(N))$. Implementări pentru aceste structuri de date se găsesc în biblioteca scikit-learn[**scikit-learn**].

Pentru acest model nu este prea relevantă funcția kernel aleasă, hiperparametrul h fiind cel care influențează mai mult performanța.

Metodele prezentate mai sus pot estima funcția de densitate probabilitate doar pentru un singur h . Dacă dorim să testăm performanța pe mai multe valori h , trebuie să rulăm acei algoritmi de mai multe ori. Acest lucru va duce la o perioadă mult prea mare de căutare a hiperparametrului.

Astfel, a apărut și o variantă generalizată a modelelor ce foloseau un h fix, menită să scadă timpul de căutare a valorii optime pentru lățimea de bandă[**Rapid-KDE**].

2.8 Metrici de performanță

Setul de date prezentat în această lucrare este adnotat în întregime, așa că putem folosi aceleași tehnici de evaluare a performanței utilizate pentru clasificarea binară.

Prin urmare, are sens să folosim următoarele noțiuni ce vor fi utile pentru descrierea metricilor de performanță prezentate mai jos:

- **TP** (*true positive*) - reprezintă observațiile din clasa **pozitivă** ce au fost clasificate **corect**
- **TN** (*true negative*) - reprezintă observațiile din clasa **negativă** ce au fost clasificate **corect**
- **FP** (*false positive*) - reprezintă observațiile din clasa **pozitivă** ce au fost clasificate **greșit**
- **FN** (*false negative*) - reprezintă observațiile din clasa **negativă** ce au fost clasificate **greșit**

În cazul nostru, clasa **pozitivă** este reprezentată de **anomalii**, iar clasa **negativă** este reprezentată de datele **normale**.

Metricile descrise în această lucrare reprezintă o alegere personală pe care o facem, astfel încât să reflecte cât mai bine nevoile problemei expuse și nicidecum nu reprezintă singura sau cea mai bună cale de a evalua performanța algoritmilor. Putem pune în paralelă cu teorema **"No Free Lunch"** care ne spune că nu există un model care să fie cel mai bun în toate situațiile, ci că utilitatea acestuia depinde strict de context. La fel este și în cazul măsurilor de evaluare, iar acest fapt face găsirea unelei potrivite pentru problemă una nu tocmai simplă, ba chiar poate necesita un timp considerabil de gândire.

2.8.1 Accuracy

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (2.11)$$

Această metrică ne indică câte **clasificări făcute de model au fost corecte din totalul de puncte care trebuie clasificate**.

2.8.2 Precision

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2.12)$$

Precision ne indică **capacitatea modelului de a nu produce fals pozitive**, în cazul nostru, de a nu raporta o valoare normală ca fiind anomalie.

2.8.3 Recall

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2.13)$$

Recall ne indică **capacitatea modelului de a identifica toate observațiile pozitive**, în cazul nostru, de a detecta toate anomaliile.

2.8.4 F1 score

$$\text{F1 Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (2.14)$$

F1 score reprezintă **media armonică dintre precision și recall**. Prin urmare, f1 score va tinde către valoarea mai mică dintre aceste 2 metrice. Pentru a o maximiza, ar trebui să avem o valoare mare atât pentru precision, cât și pentru recall, fapt ce ar duce la un model ideal.

2.8.5 AUC și ROC Curve

ROC Curve ne ajută să evaluăm calitatea modelului prin reprezentarea grafică a ratei de fals pozitiv pe axa X și a ratei de adevărat pozitiv pe axa Y. **Punctul ideal al graficului se află în colțul din stânga sus** pentru că ne dorim o rată de fals pozitiv egală cu 0 și o rată de adevărat pozitiv egală cu 1. Prin urmare, ne dorim să maximizăm rata de adevărat pozitiv și de a minimiza rata de fals pozitiv.

Pentru a crea graficul, avem nevoie de **probabilitățile sau valorile de încredere** pentru fiecare observație din setul de testare, generate de funcția de decizie a modelului respectiv. Punctele de pe grafic pot fi văzute precum clasificatoare separate ce diferă prin **pragul** aplicat funcției de decizie. Prin urmare, dacă dorim să ilustrăm ROC Curve, avem nevoie de un algoritm care are ca valori de ieșire scoruri care pot fi comparate. One Class SVM, prin definiție, nu oferă astfel de rezultate, așa că va trebui să îl tratăm în mod diferit.

Funcția de decizie este cea care atribuie un scor pentru un punct din set cu scopul de a indica nivelul de normalitate sau de anomalie al acestuia. Generarea etichetei se face apoi folosind un prag în cazul probabilităților, precum Gaussian Mixture Model și Kernel Density Estimation, sau efectiv reducând valorile pozitive la +1 și pe cele negative la -1, precum One Class SVM.

Totuși, ROC Curve are caracter vizual și nu ne oferă o măsură concretă a performanței. De aceea, avem nevoie de **AUC**, valoare ce reprezintă aria de sub grafic. Cu cât aria este mai mare, cu atât modelul este mai bun.

2.8.6 Micro average vs Macro average

Aceste 2 tehnici pot fi folosite și în cazul clasificării cu 2 clase, dar ele au apărut în principal pentru a servi ca generalizare a metricilor clasice de performanță care de unele singure tratează doar cazul **One-versus-All**, anume când considerăm că există doar 2 clase, una care conține clasa țintă și una care conține toate celelalte clase diferite de cea din urmă. Astfel, putem combina măsurile de evaluare deja existente pentru a obține rezultate relevante în cazul general.

Nu există o singură metodă de evaluare a modelelor care să fie potrivită în toate cazurile. În schimb, metodele sunt alese astfel încât să reflecte cât mai bine nevoile problemei.

Macro average pentru o măsură de evaluare are forma:

$$B_{macro} = \frac{1}{q} \sum_{\lambda=1}^q B(tp_{\lambda}, tn_{\lambda}, fp_{\lambda}, fn_{\lambda}) \quad (2.15)$$

Micro average pentru o măsură de evaluare are forma:

$$B_{micro} = B\left(\sum_{\lambda=1}^q tp_{\lambda}, \sum_{\lambda=1}^q tn_{\lambda}, \sum_{\lambda=1}^q fp_{\lambda}, \sum_{\lambda=1}^q fn_{\lambda}\right) \quad (2.16)$$

$$L = \{\lambda_j : j = 1, \dots, q\}$$

este setul tuturor etichetelor asociate claselor, iar

$$B(tp, tn, fp, fn)$$

este o măsură de evaluare binară bazată pe noțiunile introduse mai sus [Asch2013MacroandME]. În cazul nostru, $\lambda = 2$.

Diferența între cele 2 metode este că macro average acordă o importanță **egală fiecărei clase**, pe când micro average acordă o importanță **egală fiecărei observații**. Prin urmare, varianta micro favorizează clasa **majoritară** la calcularea scorului, în timp ce varianta macro favorizează clasa **minoritară**.

Ilustrăm aceste diferențe printr-un exemplu ce are ca măsură de evaluare binară **accuracy**, și care arată cum metodele acoperă nevoi diferite.

Ponderea claselor este extrem de neechilibrată în setul nostru de date, anomalii reprezentând doar 0.017% din total. Prin urmare, dacă dorim să maximizăm micro average accuracy, putem alege un model trivial ce mereu prezice clasa majoritară. Am obține un accuracy de peste 99.9% cu un minim de efort!

Este impresionant, dar modelul de mai sus este practic inutil pentru problema noastră. Toate anomaliiile ar trece nedetectate. În schimb, dacă am evalua același model folosind varianta macro average, am obține un accuracy de doar 50%. Modelul este acum inutil,

întrucât suntem interesați să detectăm anomaliile, nu doar să punem eticheta corectă pe cât mai multe observații indiferent de clasă.

Astfel, vom folosi varianta macro average pentru accuracy, iar pentru precision, recall și f1 score, vom calcula rezultatul folosind clasa minoritară ca referință. Decizia este influențată de faptul că vrem să urmărim performanța modelului pe detectarea anomaliilor în special, cu ajutorul precision și recall, dar în același timp vrem să obținem un echilibru între fals pozitive și adevărat pozitive, utilizând accuracy împreună cu f1 score, metode ce iau în calcul performanța per total pe ambele clase. Este important să detectăm cât mai multe fraude, dar în același timp nu ne dorim să semnalăm un număr prea mare de tranzacții ca fiind problematice deoarece modelul ar deveni un inconvenient.

Capitolul 3

Explorarea setului de date

3.1 Descriere

Acest set de date conține informații despre 284,807 de tranzacții dintre care 492 sunt **frauduloase**. Conține doar variabile numerice.

Din motive de confidențialitate, dimensiunea trăsăturilor originale a fost redusă folosind **Principal Component Analysis** (PCA) la 28 de trăsături denumite "V1-V28" și încă 2 trăsături ce nu au fost transformate, anume suma de bani a tranzacției și timpul relativ, începând de la 0, când aceasta a fost făcută. De asemenea, pentru fiecare tranzacție avem eticheta 0 sau 1, dacă este normală sau, respectiv, frauduloasă.

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...	-0.018307	0.277838	-0.110474	0.066928	0.128539	-0.189115	0.133558	-0.021053	149.62	0
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.638672	0.101288	-0.339846	0.167170	0.125895	-0.008983	0.014724	2.69	0
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771679	0.909412	-0.689281	-0.327642	-0.139097	-0.055353	-0.059752	378.66	0
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321	-1.175575	0.647376	-0.221929	0.062723	0.061458	123.50	0
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...	-0.009431	0.798278	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.215153	69.99	0
...
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215	7.305334	1.914428	...	0.213454	0.111864	1.014480	-0.509348	1.436807	0.250034	0.943651	0.823731	0.77	0
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330	0.294869	0.584800	...	0.214205	0.924384	0.012463	-1.016226	-0.606624	-0.395255	0.068472	-0.053527	24.79	0
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827	0.708417	0.432454	...	0.232045	0.578229	-0.037501	0.640134	0.265745	-0.087371	0.004455	-0.026561	67.88	0
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180	0.679145	0.392087	...	0.265245	0.800049	-0.163298	0.123205	-0.569159	0.546668	0.108821	0.104533	10.00	0
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006	-0.414650	0.486180	...	0.261057	0.643078	0.376777	0.008797	-0.473649	-0.818267	-0.002415	0.013649	217.00	0

284807 rows x 31 columns

Figura 3.1: Cele 284,807 de tranzacții

3.2 Matricea de corelație

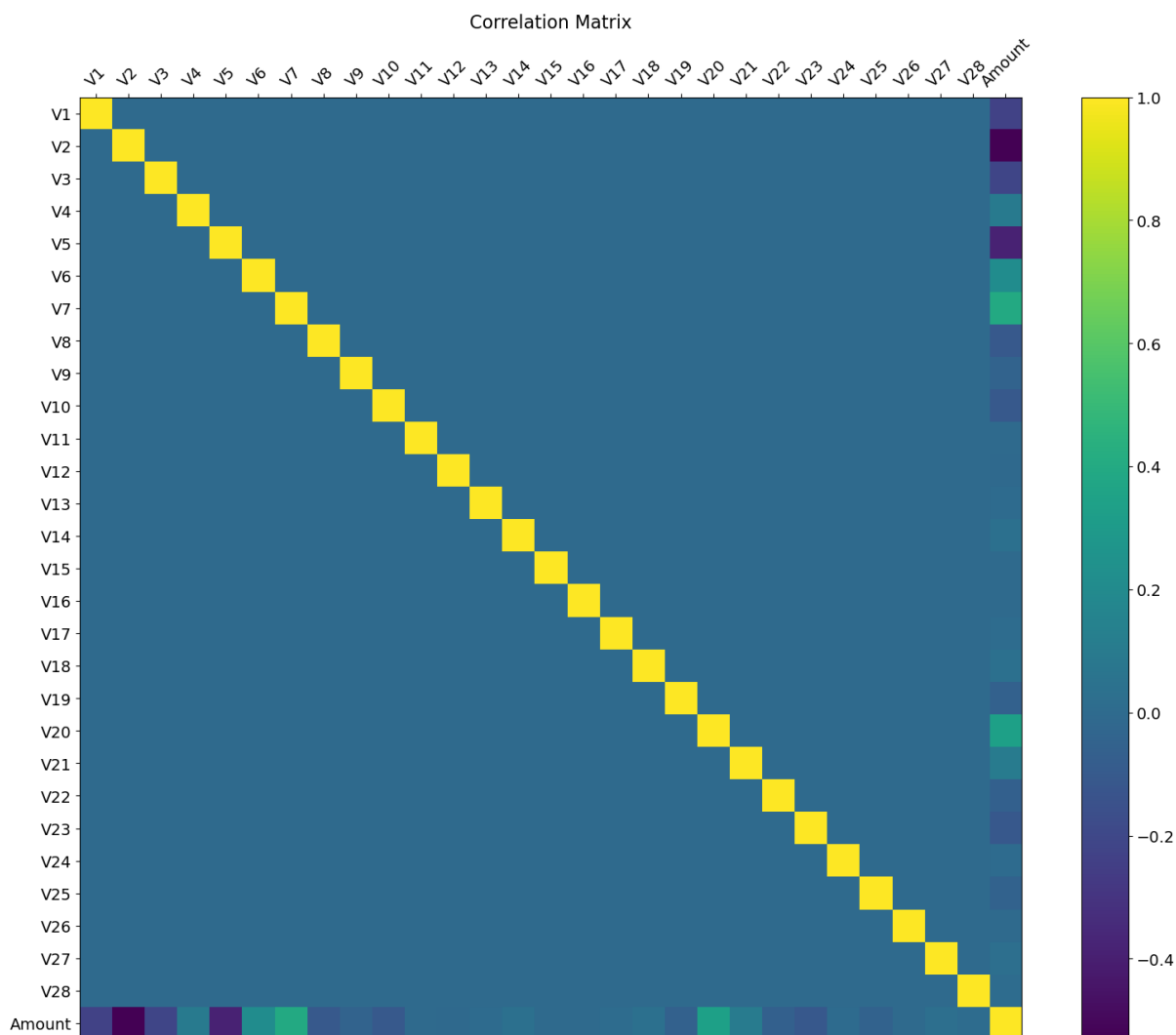


Figura 3.2: Matricea de corelație

Se observă lipsa corelației între variabilele anonime. Este de așteptat, totuși, având în vedere că variabilele au fost obținute prin PCA, care din definiție generează componente **necorelate liniar**.

Singura corelație apare între variabila "Amount" și restul variabilelor. Corelațiile cu o magnitudine semnificativă apar între amount și V6, V7, V20, cu valoare pozitivă, și între amount și V1, V2, V5, cu valoare negativă.

3.3 Preprocesarea datelor

3.3.1 Eliminăm trăsăturile inutile

În viața reală, data și ora tranzacției ne pot ajuta să depistăm un comportament ciudat. Totuși, aceste informații sunt relevante numai atunci când am monitorizat pe o perioadă de măcar câteva zile activitatea din contul/cardul respectiv pentru a stabili ce reprezintă un comportament "normal".

În cazul nostru, nu avem la dispoziție nici măcar data și ora exactă. Prin urmare, vom elimina timpul din fiecare observație.

De asemenea, eliminăm etichetele. Acestea sunt utile numai la testare pentru a analiza performanța și la alegerea setului de antrenare. Metodele utilizate în această lucrare sunt exclusiv **nesupervizate**, deci nu avem nevoie de etichete la procesul de antrenare.

3.3.2 Validarea încrucișată

Împărțim setul de date, salvând 75% pentru antrenare. Din **setul de antrenare**, extragem toate datele cu eticheta de anomalie și le mutăm în porțiunea de 25% rămasă de la pasul anterior. Apoi, ultima porțiune este împărțită la rândul ei în jumătăți egale ce vor servi ca set de validare, respectiv testare.

Facem acest lucru deoarece algoritmi prezentați "învață" structura datelor normale și, ideal, nu am vrea niciun fel de anomalie în setul de antrenare. În practică, este o șansă destul de mare ca setul de antrenare să fie "**contaminat**" cu un procent mic de anomalii, fiind dificil să garantăm "**puritatea**" datelor. Totuși, am decis să presupunem că nu există nicio eroare în datele noastre pentru a păstra un număr mai mare de anomalii în setul de validare și cel de testare, pe care să evaluăm performanța modelelor.

De asemenea, o tehnică des utilizată în împărțirea setului de date pentru metode supervizate este **stratificarea** ce ne conferă o proporție a claselor, aproximativ egală cu cea din setul inițial, în porțiunile de antrenare, validare și testare. Acest lucru este important pentru a ne asigura că fiecare clasă este reprezentată corespunzător în eșantionul nostru. În cazul detecției anomaliilor, o vom ignora deoarece avem doar o singură clasă de referință, anume cea normală, la antrenare. Datele etichetate ca fiind anomalii sunt utile doar la testare.

Setul de validare va fi folosit pentru alegerea hiperparametrilor optimi, în timp ce **setul de testare** va fi folosit la final pentru a face o comparație imparțială a performanței modelelor.

3.3.3 Scalare

Având în vedere că nu cunoaștem dimensiunile, unitățile de măsură și nici numele variabilelor aleatoare asociate punctelor, este indicat să scalăm datele. Uneori poate fi util să lăsăm punctele așa cum sunt, dar aceasta ar necesita o decizie informată pe baza importanței variabilelor în calcularea distanței între puncte sau pe baza relevanței unităților variabilelor, de exemplu.

În lipsa acestor factori, vom standardiza datele, adică **vom scădea media și vom împărți la deviația standard**, operație aplicată fiecărui punct. Aceasta are rolul de a centra variabilele și de a obține varianță unitară pe fiecare variabilă în parte. Media și deviația standard sunt obținute din setul de antrenare. Le folosim pe acestea să transformăm setul de validare, cât și pe cel de testare. Astfel, fiecare trăsătură x_i va fi înlocuită de o noua trăsătură x'_i :

$$x'_i = \frac{x_i - \mu_i}{\sigma_i}$$

unde μ_i este media:

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)}$$

și σ_i este deviația standard:

$$\sigma_i = \sqrt{\frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2}$$

calculate luând trăsăturile x_i din toate cele m observații din setul de antrenare.

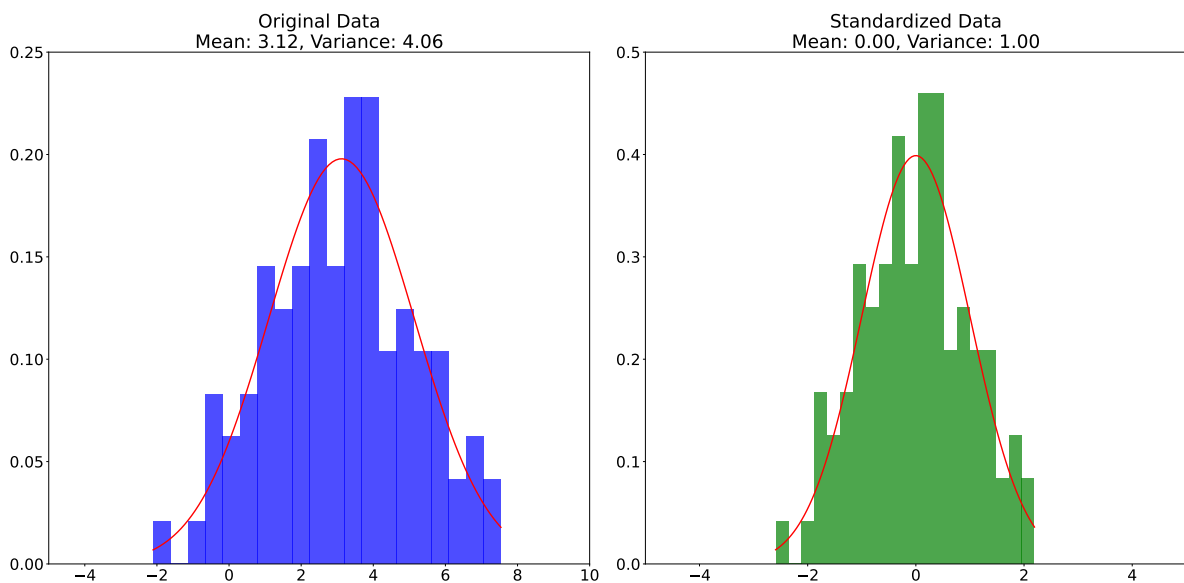


Figura 3.3: Media devine 0, varianța devine 1

Capitolul 4

Evaluarea modelelor

Căutarea **hiperparametrilor optimi** se va realiza pentru fiecare model în parte folosind setul de validare. La final, cele mai bune modele găsite sunt evaluate într-un mod imparțial pe setul de testare.

4.1 Modele nesupervizate

4.1.1 One Class SVM

Pentru a exploata capabilitățile de modelare a unei margini de decizie neliniare a SVM-urilor, avem nevoie de funcții kernel care să scufunde punctele din setul de date într-un spațiu cu mai multe dimensiuni unde să putem găsi mai ușor un hiperplan de separare.

- $K(x, y) = x^T y$ - **Liniar**
- $K(x, y) = (\gamma x^T y + c)^d$ - **Polinomial**
- $K(x, y) = \exp(-\gamma \|x - y\|^2)$ - **Gaussian**
- $K(x, y) = \tanh(\gamma x^T y + c)$ - **Sigmoid**

Un caz particular este kernel-ul **liniar** care păstrează punctele în spațiul inițial și încearcă să găsească un hiperplan de separare optim la fel ca în cazul clasic fără kernel al SVM-ului. Acest kernel este nepotrivit în cele mai multe cazuri, mai puțin când datele sunt aproape liniar separabile, deci nu ne așteptăm să performeze excepțional.

Totuși, putem exploata eficiența kernel-ului liniar folosind tehnica de optimizare **Stochastic Gradient Descent** (SGD) împreună cu metoda de aproximare **Nystroem** pentru a aplica o transformare neliniară asupra datelor de intrare și apoi să găsim o margine de decizie liniară în noul spațiu.

SGD este o metodă iterativă relativ simplă ce nu necesită un număr la fel de mare de calcule și nici la fel de multă memorie precum metodele clasice de rezolvare a unui sistem de ecuații folosind algebră liniară. De aceea, este potrivită atunci când mărimea setului de date trece de ordinul sutelor de mii, în ciuda faptului că sacrificăm acuratețea ponderilor estimate.

Metoda Nystroem aproximează matricea kernel pentru o anumită funcție kernel dată, folosind tehnica aproximării cu **matrice de rang scăzut** unde o fracțiune din punctele setului de antrenare este folosită ca bază vectorială pentru kernelul respectiv. Evident, calculul va fi mult mai rapid folosind o matrice de dimensiune redusă. Am comparat performanța acestei metode pentru diverse funcții kernel și pentru diverse fracțiuni de puncte din setul de antrenare, care erau folosite de algoritm. Numărul de puncte a fost selectat astfel încât să ocupe un procent, dat ca parametru, din memoria RAM. Prin urmare, rezultatele pot să difere în funcție de capacitatea memoriei pe care o avem la dispoziție.

Pentru comparație am inclus și funcțiile **polinomială** și **sigmoidă**. Cu cea din urmă am obținut rezultatele cele mai bune în comparație cu celelalte funcții, dar nu mai bune decât kernelul Gaussian pe care îl vom folosi și care este și decizia des întâlnită în practică. De asemenea, atât kernel-ul liniar pentru valori mai mari ale lui ν , cât și cel polinomial începând cu gradul 7 nu terminau antrenarea nici măcar după 12 ore, așa că am abandonat căutarea hiperparametrilor pentru mai mult de atât, mai ales că acuratețea prezicerilor devenea din ce în ce mai îndoielnică.

d	γ	ν	Accuracy	Recall	Precision	F1 Score
7	0.001	0.007	0.474	0.947	0.006	0.0129
5	0.001	0.02	0.469	0.930	0.006	0.0128
7	1	0.0001	0.502	0.0243	0.008	0.0124
7	2	0.0001	0.502	0.024	0.008	0.0124

Tabela 4.1: Cele mai bune rezultate obținute pentru kernel-ul polinomial

γ	ν	Accuracy	Recall	Precision	F1 Score
0.03	0.005	0.804	0.613	0.476	0.536
0.02	0.005	0.792	0.589	0.448	0.509
0.03	0.001	0.668	0.337	0.674	0.449
0.02	0.007	0.756	0.520	0.347	0.416

Tabela 4.2: Cele mai bune rezultate obținute pentru kernel-ul sigmoid

ν	Accuracy	Recall	Precision	F1 Score
0.9	0.177	0.300	0.002	0.004
0.001	0.412	0.0813	0.002	0.004

Tabela 4.3: Cele mai bune rezultate obținute pentru kernel-ul liniar

γ	ν	Accuracy	Recall	Precision	F1 Score
0.2	0.5	0.7405	0.9796	0.0134	0.0264
0.03	0.5	0.7402	0.9837	0.0133	0.0263

Tabela 4.4: Cele mai bune rezultate obținute pentru kernel-ul liniar cu SGD și Nystroem

Utilizăm kernelul Gaussian, deci printre hiperparametrii optimi pe care îi căutăm se va regăsi și γ . Acest parametru influențează **aria zonei de influență a fiecărui vector suport**. O valoare prea mare ar cauza ca zona să includă numai vectorul suport și nimic altceva, ceea ce ar duce la o **varianță crescută** a modelului. La polul opus, o valoare prea mică ar cauza ca zona să includă toate punctele din setul de date, ceea ce ar duce la un **bias crescut**.

Parametrul ν este similar parametrului C din **Soft-Margin SVM**, cel din urmă fiind creat cu scopul de a rezolva problemele asociate parametrului C , anume că putea lua orice valoare pozitivă și nu avea o interpretare directă. ν se află în intervalul $(0, 1]$ și este interpretat ca marginea superioară a ponderii de anomalii și marginea inferioară a ponderii de vectori suport. Prin urmare, ν controlează mărimea **frontierei din jurul datelor normale** a modelului, unde o frontieră mai mică este asociată cu o varianță crescută, în timp ce o frontieră mai mare este asociată cu un bias crescut.

Folosim metoda **Grid Search** pentru a găsi parametrii favorabili. Afișăm doar valorile parametrilor pentru care obținem rezultate semnificative.

γ	ν	Accuracy	Recall	Precision	F1 Score
0.01	0.0001	0.822	0.646	0.676	0.661
0.01	0.001	0.822	0.646	0.676	0.661
0.01	0.005	0.875	0.756	0.531	0.624
0.02	0.0001	0.877	0.760	0.526	0.622
0.02	0.001	0.877	0.760	0.523	0.620

Tabela 4.5: Cele mai bune rezultate obținute pentru kernel-ul Gaussian

Se observă că γ este parametrul care aduce schimbările drastice în valorile metricilor, pe când ν doar crește sau scade relativ puțin aceste valori.

Urmărim f1 score, aceasta fiind metrica ce evaluează modelul oarecum echilibrat, în contrast cu precision și recall care favorizează minimizarea fals pozitivelor, respectiv a fals

negativelor. Astfel, alegem $\gamma = 0.01$ și $\nu = 0.0001$ pentru modelul final.

4.1.2 Gaussian Mixture Model

n	q	Accuracy	Recall	Precision	F1 Score
4	0.01	0.920	0.845	0.581	0.688
3	0.01	0.918	0.841	0.578	0.685
2	0.01	0.918	0.841	0.578	0.685
1	0.01	0.914	0.833	0.572	0.678
5	0.01	0.912	0.829	0.569	0.675

Tabela 4.6: Cele mai bune rezultate obținute pentru Gaussian Mixture Model

Pentru acest model, hiperparametrul optim căutat este numărul de componente Gaussiene n . Vom încerca pe rând fiecare valoare din $\{1, 2, 3, \dots, 16, 17\}$. De asemenea, pentru că modelul ne va oferi probabilitățile de apartenență a unui punct pentru fiecare componentă, vom avea nevoie și de un prag pentru a decide dacă punctul este sau nu anomalie. Ne vom folosi de cuantile calculate pe setul de validare pentru a găsi pragul optim.

Pentru $n = 1$ putem afla parametrii foarte eficient folosind **Maximum Likelihood Estimation**, întrucât problema se reduce la aflarea mediei și a matricei de covarianță pentru o distribuție Gaussiană.

Pentru $n > 1$ vom folosi **Expectation Maximization** pentru a găsi parametrii optimi, având în vedere că numărul de componente este precizat de la început.

Gaussian Mixture Model încearcă să estimeze o distribuție posibil multimodală folosind mai multe componente Gaussiene. Prin urmare, numărul optim de componente ne indică intuitiv numărul de **moduri** pe care le are distribuția ce a generat setul de date. Se observă că modelul se descurcă destul de bine chiar și cu o singură componentă. Aceasta ne indică faptul că distribuția ce a generat setul de date este similară cu una Gaussiană.

Pentru acest model, pragul este foarte important. Dacă avem același prag, dar număr de componente diferite, varianța rezultatelor nu este prea mare, dar dacă alegem greșit valoarea cuantilei, performanța modelului are de suferit. Totuși, dată simplitatea modelului, rezultatele sunt impresionante. Cu un timp de antrenare de sub câteva minute, obține rezultate mai bune decât One Class SVM pe setul de validare.

Alegem **numărul de componente** $n = 4$ pentru modelul final, chiar dacă am avut rezultate bune și cu $n = 1$ deoarece ne dorim un model puțin mai complex decât o banală distribuție Gaussiană. În practică, este o șansă mică să dăm peste un proces care să genereze date fix în acest mod. De asemenea, aici f1 score are valoarea cea mai mare și după aceasta ne vom ghida, având în vedere că ne interesează performanța modelului per total.

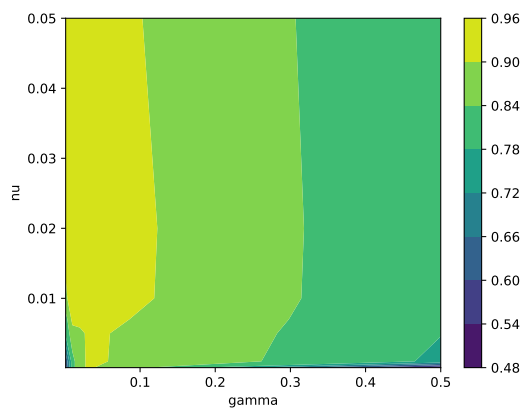


Figura 4.1: OCSVM Accuracy

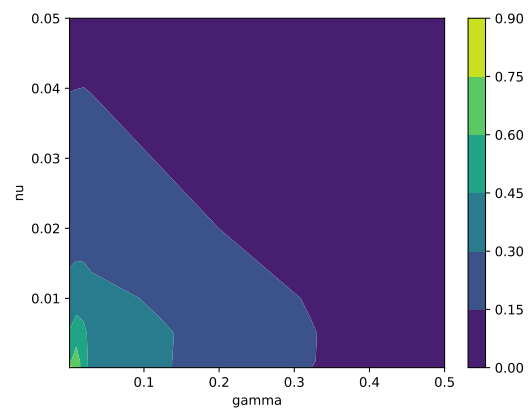


Figura 4.2: OCSVM Precision

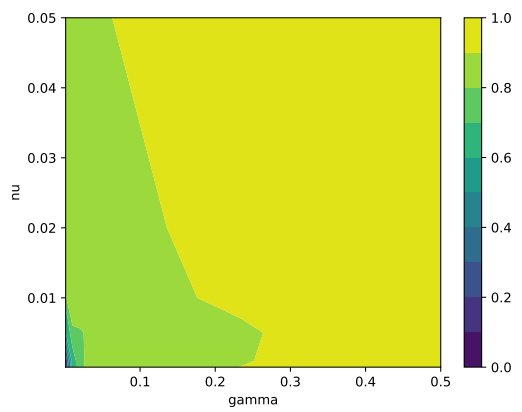


Figura 4.3: OCSVM Recall

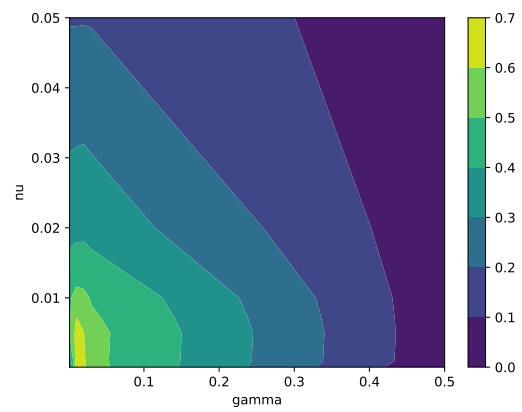


Figura 4.4: OCSVM F1 score

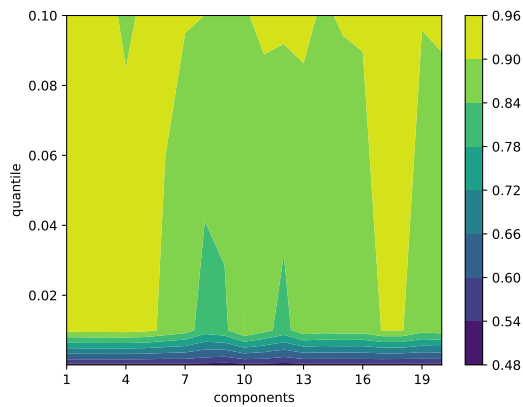


Figura 4.5: GMM Accuracy

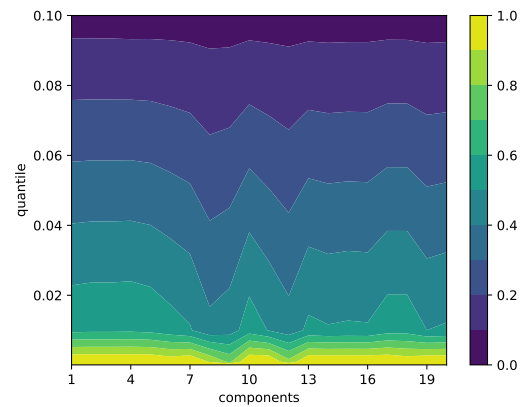


Figura 4.6: GMM Precision

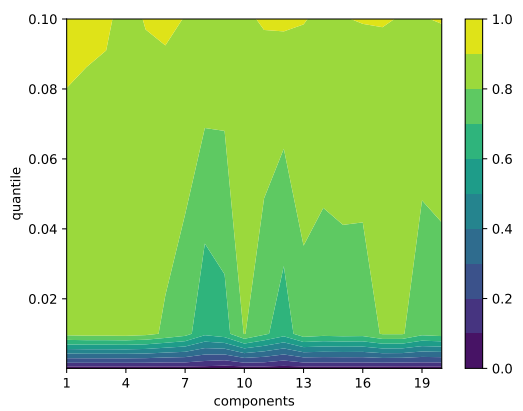


Figura 4.7: GMM Recall

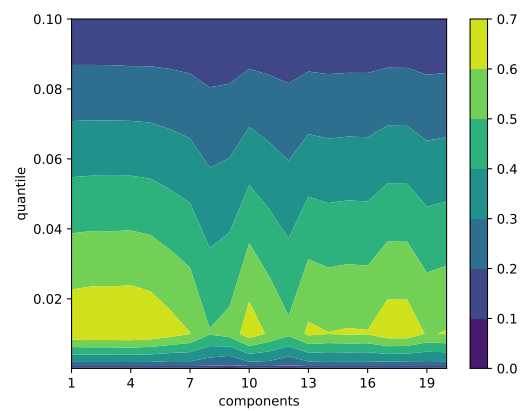


Figura 4.8: GMM F1 score

4.1.3 Kernel Density Estimation

<i>Kernel</i>	<i>h</i>	<i>q</i>	Accuracy	Recall	Precision	F1 Score
tophat	9	0.01	0.908	0.821	0.562	0.667
linear	9	0.01	0.881	0.768	0.527	0.625
parabolic	9	0.01	0.877	0.760	0.522	0.619
cosine	9	0.01	0.875	0.756	0.519	0.615
tophat	8	0.01	0.877	0.760	0.513	0.613

Tabela 4.7: Cele mai bune rezultate obținute pentru Kernel Density Estimation

Am testat mai multe funcții kernel, anume: **tophat**, **liniar**, **parabolic**, **cosinus**, **exponențial** și **gaussian**. Deși kernelul Gaussian este cel mai des întâlnit în practică datorită numeroaselor proprietăți utile pe care le deține, aici se observă că nu obține rezultate satisfăcătoare, kernelul **tophat** fiind cel care performează cel mai bine. De asemenea, trebuie să impunem un prag după care să decidem dacă un punct este sau

nu anomalie. Vom alege cuantile luate pe setul de validare, precum în cazul Gaussian Mixture Model.

Lăţimea de bandă este cea care stă la baza **bias-variance tradeoff** în acest model. Valorile prea mici implică **variance** mare, întrucât aria de sub grafic pentru fiecare punct este influenţată doar de punctele foarte apropiate de el, fapt ce duce la o distribuţie cu mulţi "ţepi". În schimb, valorile prea mari implică **bias** mare pentru că acum şi punctele aflate la distanţă mare joacă un rol important. În cel mai rău caz, o distribuţie multimodală ajunge să fie estimată ca una unimodală din cauza netezimii graficului.

Se observă că după 1.0, valoarea lăţimii de bandă nu mai aduce îmbunătăţiri semnificative. Acest lucru este ilustrat şi pe grafic sub forma unui maxim local pentru f1 score. La fel ca în cazul Gaussian Mixture Model, cuantila folosită are un puternic impact asupra performanţei modelului.

Din păcate, după un timp de antrenare ce depăşeşte 40 minute, timp ce se află între valorile pentru OCSVM şi GMM, nu obţinem rezultate mai bune decât niciunul dintre modelele precedente pe setul de validare.

Alegem **lăţimea de bandă** $bandwidth = 1.0$ pentru modelul final deoarece are cea mai bună valoare pentru f1 score. De asemenea, aceasta este şi cea mai mică valoare care aduce o performanţă decentă. O valoare prea mare a lăţimii de bandă duce la o netezire excesivă a particularităţilor distribuţiei.

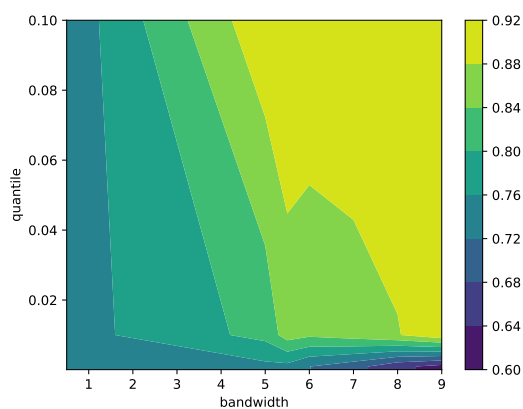


Figura 4.9: KDE Accuracy

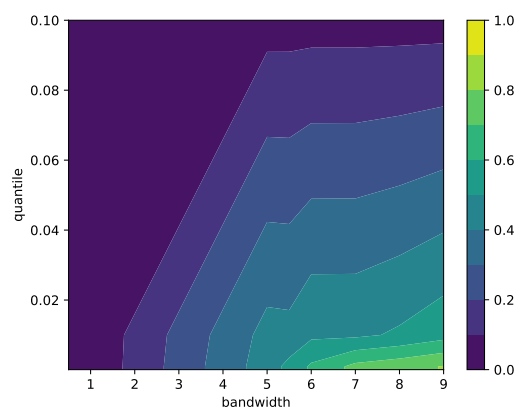


Figura 4.10: KDE Precision

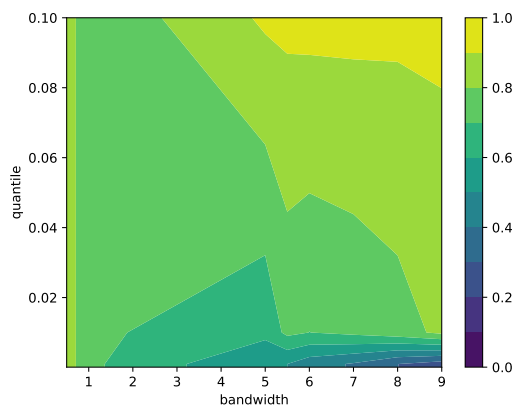


Figura 4.11: KDE Recall

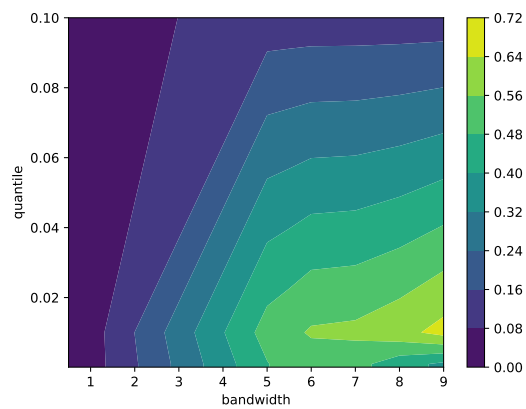


Figura 4.12: KDE F1 score

4.2 Modele supervizate

4.2.1 SVM Supervizat

Pentru a profita de faptul că avem la dispoziție etichetele datelor, vom trata problema și din punct de vedere supervizat. Gaussian Mixture Model și Kernel Density Estimation nu pot fi aplicate decât într-un mediu nesupervizat, așa că vom antrena doar SVM-ul, dar de data aceasta pentru clasificare binară. Astfel, putem observa performanța algoritmilor prezentați anterior și față de un algoritm supervizat care are de rezolvat o problemă relativ mai ușoară.

O modificare pe care o vom face la modul de împărțire al setului de date este că de această dată stratificarea este luată în calcul, așa că păstrăm ponderile claselor aproximativ la fel în toate cele 3 partiții.

Și aici, kernel-ul Gaussian este cel care oferă cele mai bune rezultate cu valori scăzute pentru γ , similar cu cazul OCSVM, dar parametrul ν este înlocuit de parametrul C care are rolul de regularizare în SVM-ul Soft-Margin.

Se observă o creștere substanțială a valorilor metricilor atunci când tratăm problema în mod supervizat, **F1 Score** depășind 0.87 pe setul de validare, în timp ce OCSVM pe setul de validare abia atinge 0.12.

γ	C	Accuracy	Recall	Precision	F1 Score
0.01	0.1	0.7972	0.5946	0.8302	0.6929
0.01	0.5	0.9053	0.8108	0.8571	0.8333
0.01	1.0	0.9053	0.8108	0.8824	0.8451
0.01	2.0	0.9053	0.8108	0.8955	0.8511
0.01	3.0	0.9054	0.8108	0.9375	0.8696
0.01	4.0	0.9054	0.8108	0.9524	0.8759
0.02	0.1	0.7094	0.4189	0.7750	0.5439
0.02	0.5	0.8850	0.7703	0.8769	0.8201
0.02	1.0	0.8918	0.7838	0.9355	0.8529
0.02	2.0	0.8919	0.7838	0.9508	0.8593
0.02	3.0	0.8919	0.7838	0.9508	0.8593
0.02	4.0	0.8919	0.7838	0.9667	0.8657
0.03	0.1	0.6148	0.2297	0.8947	0.3656
0.03	0.5	0.8243	0.6486	0.9412	0.7680
0.03	1.0	0.8783	0.7568	0.9492	0.8421
0.03	2.0	0.8851	0.7703	0.9661	0.8571
0.03	3.0	0.8851	0.7703	0.9661	0.8571
0.03	4.0	0.8851	0.7703	0.9661	0.8571

Tabela 4.8: Grid Search pentru SVM

Capitolul 5

Compararea modelelor finale

În final, antrenăm fiecare tip de model cu hiperparametrii optimi găsiți folosind setul de validare și folosim setul de testare pentru a evalua imparțial modelele finale.

Model	Accuracy	Recall	Precision	F1 Score
OCSVM	0.8913	0.0577	0.8823	0.1080
KDE	0.8911	0.0575	0.8823	0.1078
GMM	0.865	0.735	0.505	0.599

Tabela 5.1: Performanța fiecărui model pe setul de testare

Cele 2 modele care au o performanță similară sunt Gaussian Mixture Model și One Class SVM, ultimul fiind foarte puțin mai slab decât cel din urmă, dar cu un timp de antrenare mult mai mare. Kernel Density Estimation este de departe un model nepotrivit pentru acest set de date, necesitând un timp crescut de antrenare pentru niște rezultate mediocre.

Precision este scăzut și prin urmare și f1 score este scăzut. În schimb, recall are o valoare ridicată pentru toate modelele. Acest fapt ne arată că detectăm aproape toate anomaliiile, dar în același timp, semnalăm un număr mare de tranzacții obișnuite ca fiind frauduloase. Acest lucru poate deveni un inconvenient pentru client în cazul în care tranzacția este anulată ca urmare a raportării incorecte.

Gaussian Mixture Model este modelul cu valorile cele mai mari pentru toate metricile, chiar dacă nu este un model la fel de complex precum celelalte 2 tehnici nesupervizate. De asemenea, acesta a necesitat doar 3 minute pentru antrenare și prezicere a etichetelor, dar a produs rezultate aproape la fel de bune precum One class SVM și mult mai bune decât Kernel Density Estimation care au avut nevoie de mai mult de o oră și jumătate, respectiv 40 minute, cel din urmă neavând o performanță cu mult mai ridicată în raport cu complexitatea de timp.

Ilustrăm și ROC Curve alături de metrica AUC pentru modelele finale. Cu cât graficul tinde să se lipească de colțul stânga sus, cu atât aria de sub grafic crește rezultând într-

o performanță mai bună. Linia punctată pe diagonală din fiecare diagramă reprezintă graficul produs de clasificatorul aleatoriu și are rolul de a marca marginea inferioară a performanței față de această metrică. De asemenea, este de menționat că dacă graficul se află sub linia diagonală, atunci am trage concluzia că algoritmul nostru este inutil. Totuși, putem inversa clasa negativă cu cea pozitivă și astfel obținem simetricul graficului față de diagonală care acum, evident, se află deasupra liniei. Prin urmare, distanța între grafic și linie este o măsură mai relevantă decât poziția relativă. Se presupune că operația de aplicare a simetricului este folosită unde este cazul pentru diagramele de mai jos.

ROC Curve are nevoie ori de probabilități ori de valori ce exprimă încrederea pentru fiecare prezicere, fără aplicarea vreunui prag. Kernel Density Estimation și Gaussian Mixture Model produc deja probabilități, pe când One Class SVM nu produce decât distanța față de hiperplanul de separare. Totuși, vom folosi un truc pentru a transforma rezultatele OCSVM în valori de încredere [stackoverflow-auc]. Înlocuim fiecare valoare cu valoarea respectivă scăzută din valoarea maximă dată de funcția de decizie după cum urmează:

$$y_{score} = MAX - f(y)$$

unde

- y_{score} este valoarea folosită pentru ROC Curve
- MAX este $\max_{i=1}^n f(y_i)$ pentru n observații
- f este funcția de decizie

Astfel, putem compara modelele finale folosind și metrica AUC. Clasamentul este același și aici, cu One Class SVM fiind cel mai performant, urmat de Kernel Density Estimation, iar Gaussian Mixture Model la coadă.

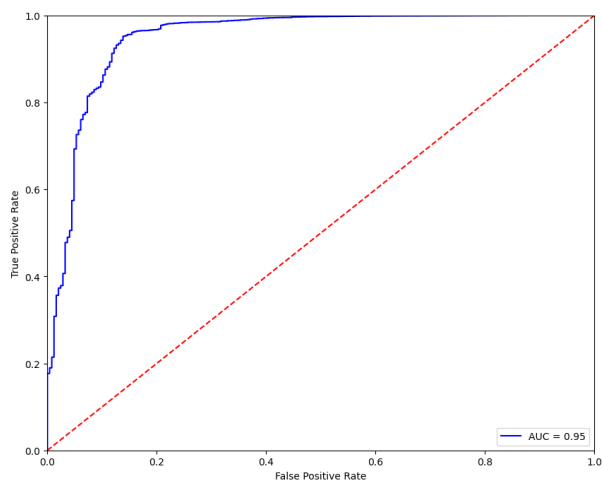


Figura 5.1: ROC Curve OCSVM

Figura 5.2: ROC Curve GMM

Figura 5.3: ROC Curve KDE

Capitolul 6

Concluzii

În această lucrare am ilustrat una din multitudinea de aplicații în industrie a detecției anomaliilor, anume identificarea tranzacțiilor frauduloase. Am explorat un set de date adnotat cu informații despre tranzacții bancare și am arătat cum putem evalua performanța unor metode de învățare automată nesupervizată folosind metrici din domeniul învățării supervizate, mai specific, clasificarea binară.

Deși tehnicile folosite aparțin învățării nesupervizate, am arătat cum putem exploata cunoștințele din alte două categorii, anume semi-supervizată și supervizată.

Prima a fost utilă pentru faptul că am antrenat algoritmi pe un set de date mare fără să utilizăm etichetele punctelor, dar apoi am folosit un set de date cu o mărime relativă mult mai mică pentru a analiza performanța algoritmilor.

A doua a fost de folos pentru că ne-a oferit o multitudine de metrici care într-un context strict nesupervizat nu ar fi existat, fiind nevoie de inspecția unui om pentru analiză.

Alegerea metricilor corespunzătoare este probabil una din cele mai dificile părți atunci când căutăm o soluție eficientă pentru problema dată, mai ales când proporția de clase este neechilibrată și acordăm o importanță mai mare unei clase față de cealaltă, întrucât o alegere greșită ne poate induce în eroare cu privință la utilitatea reală a algoritmului, după cum am și demonstrat.

De asemenea, am analizat trei algoritmi cu caracteristici diferite, dar cu scopuri similare, anume identificarea anomaliilor.

One Class SVM se folosește de un hiperplan de separare pentru a împărți spațiul Hilbert într-o porțiune care cuprinde toate punctele normale și una care cuprinde restul punctelor izolate. Astfel, putem spune că acesta face o clasificare dură a anomaliilor, întrucât ne indică doar dacă o observație este normală sau nu.

La polul opus, avem Gaussian Mixture Model și Kernel Density Estimation care încearcă să estimeze funcția densitate de probabilitate din care au fost generate datele și să atribuie fiecărei observații noi o valoare dată de această funcție. Peste aceste probabilități se aplică un prag și abia apoi obținem etichetele corespunzătoare. Astfel, putem spune că acestea fac o clasificare slabă a anomaliilor, întrucât ne indică cât de probabil este ca

punctul să fie normal sau nu, dar nu ne indică concret clasa de care aparține.

Deși am întâmpinat probleme în ilustrarea graficului ROC Curve pentru algoritmul cu clasificare dură, am demonstrat cum putem transforma distanțele produse de funcția de decizie în valori de încredere pentru a obține o clasificare slabă artificială.

La final, am observat că modelele au performat mai mult sau mai puțin la fel, cele mai complexe având un dezavantaj când vine vorba de complexitatea de timp.

Cu toate acestea, lucrarea nu a abordat problema seriilor de timp, lucru des întâlnit în practică, dar care creează o direcție de cercetare interesantă. Două cerințe relevante pentru acest subiect ar fi prognoza și clasificarea seriilor de timp.

Prima se referă la prezicerea evoluției în timp a seriei de date și este folosită ca suport pentru unele modele ce detectează anomalii folosindu-se de deviația prognozei față de valorile observate. Prognoza se bazează strict pe valorile punctelor din șir din contextul precedent și pe modelul antrenat anterior. Deseori, se utilizează o fereastră glisantă de mărime unitară pentru a crea contextul și a prezice câte un punct pe rând. Este similar cu abordarea clasică din domeniul Procesării Limbajului Natural, unde suntem interesați să prezicem cuvântul cel mai potrivit pentru continuarea unei propoziții, spre exemplu. Aici, putem privi cuvintele din text ca fiind puncte dintr-o serie de timp.

A doua se referă la atribuirea unei categorii pentru diverse subșiruri din seriile de timp. Astfel, metoda poate fi folosită ori pentru post-procesarea rezultatelor obținute după detecția anomaliilor pentru a categoriza diferitele valori deviate, ori chiar în cadrul modelului cu rolul de a clasifica componentele din seriile de date în grupuri relevante, implicit anomalii sau observații normale. Metoda ne aduce aminte de clasică problemă a clasificării unor puncte dintr-un set de date, întâlnită în contextul învățării automate supervizate [**time-series**].