

MINISTERUL EDUCAȚIEI



UNIVERSITATEA TEHNICĂ
DIN CLUJ-NAPOCA
FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE

MODELAREA CONFORTULUI TERMIC PENTRU OPTIMIZAREA SISTEMELOR HVAC

PROIECT DE SEMESTRU

Student: **Domide Maria**

Disciplina: **Sisteme bazate pe Cunoaștere**

2025

Cuprins

1	INTRODUCERE.....	3
1.1	CONTEXT GENERAL	3
1.2	OBIECTIVE.....	4
1.3	SPECIFICAȚII.....	5
2	CUNOAȘTEREA ȘI ANALIZA SETULUI DE DATE	6
2.1	Definiție.....	6
2.2	Introducere.....	6
2.3	Mediu de implementar	7
2.4	Citirea și explorarea inițială a datelor.....	7
2.5	Analiza Univariată.....	9
2.5.1	Scopul analizei univariat.....	9
2.5.2	Distribuția variabilelor numerice	9
2.5.3	Distribuția variabilelor categorice	10
2.6	Analiza bivariate	11
2.7	Analiza multivariată.....	13
2.7.1	PCA (Principal Component Analysis):	13
2.7.2	t-SNE (t-distributed Stochastic Neighbor Embedding):	14
2.7.3	Matricea de corelație	15
3	PRE-PROCESAREA SETULUI DE DATE	16
3.1	Introducere.....	16
3.2	Denoising (Eliminarea zgomotului)	16
3.2.1	Definiție și scop.....	16
3.2.2	Metode utilizate:.....	16
3.3	Detrending (Eliminarea tendințelor).....	18
3.3.1	Definiție și scop +grafice	18
3.4	Eliminarea valorilor aberante (Outlier Removal)	21
3.4.1	Definiție și scop	21
3.4.2	Metoda Z-Score (Z-Score Method).....	21
3.4.3	Metoda IQR (Interquartile Range)	23
3.4.4	MAD (Median Absolute Deviation)	24
3.4.5	DBSCAN (Density-Based Spatial Clustering of Applications with Noise).....	pag 25
3.4.6	Isolation Forest Method.....	pag 27

3.5. Interpolare.....	28
3.5.1. Definiție și importanță	28
3.5.2 Interpolare liniară	28
3.5.3 Interpolare cubică	29
3.5.4 LOESS (Locally Estimated Scatterplot Smoothing)	29
4. MODELAREA SISTEMULUI	31
4.1 ÎNTRUDUCERE	31
4.1.1 IMPORTANȚA MODELĂRII CORECTE ÎN CONTEXTUL APLICAȚIEI HVAC	31
4.1.2 UTILITATEA DATELOR PREPROCESATE PENTRU ÎMBUNĂȚĂȚIREA PERFORMANȚEI MODELULUI	31
4.2 METODOLOGIE ȘI TEHNICI APLICATE	32
4.2.1 REGRESIE LINIARĂ	32
4.2.2 REGRESIE LOGISTICĂ	33
4.2.3 REGRESII RIDGE ȘI LASSO	33
4.2.4 RANDOM FOREST	34
4.2.5 XGBOOST	36
4.3 ÎMPLIMENTAREA	37
4.3.1 MEDIUL DE ÎMPLIMENTARE	37
4.3.2 CUM SE POATE REPLICA LUCRAREA	37
4.4 EVALUAREA MODELELOR	38
4.4.1 METODE ȘI METRICI DE EVALUARE	38
4.4.2 ANALIZA REZULTATELOR	39
4.4.3 INTERPRETAREA REZULTATELOR	39
4.4.4 CONCLUZII PRIVIND PERFORMANȚA MODELELOR	39
4.5 CONCLUZII	40
5 CONCLUZII (1-3 PAG)	41
5.1 REZULTATE OBTINUTE	41
5.2 DIRECȚII DE DEZVOLTARE	42
6 BIBLIOGRAFIE	43

1 Introducere

1.1 Context general

Sistemele de încălzire, ventilație și aer condiționat (HVAC) joacă un rol esențial în asigurarea confortului termic și a sănătății mediului interior, fiind critice atât pentru locuințe, cât și pentru spațiile comerciale sau industriale.

Progresul tehnologic din ultimele decenii a facilitat integrarea unor tehnologii avansate bazate pe date și algoritmi de învățare automată, care pot îmbunătăți semnificativ performanța acestor sisteme.

În acest context, proiectul nostru se încadrează într-un domeniu multidisciplinar, la intersecția dintre ingineria mediului interior, analiza datelor și inteligența artificială.

Introducerea acestor tehnologii în sistemele HVAC aduce beneficii considerabile, inclusiv optimizarea consumului de energie, creșterea eficienței operaționale și îmbunătățirea calității vieții utilizatorilor prin personalizarea condițiilor de mediu.

Într-o eră în care schimbările climatice și sustenabilitatea devin priorități globale, optimizarea sistemelor HVAC este esențială pentru reducerea emisiilor de carbon și atingerea obiectivelor de eficiență energetică.

Datele climatice și subiective, cum ar fi temperatura, umiditatea, viteza aerului sau percepțiile utilizatorilor, au un impact direct asupra confortului termic. Astfel, dezvoltarea unor modele predictive care să coreleze aceste variabile este vitală pentru proiectarea unor sisteme HVAC inteligente și adaptative.

Pentru acest proiect, s-a utilizat baza de date **ASHRAE Global Thermal Comfort Database II**, o resursă cu peste 100.000 de înregistrări care oferă o combinație valoroasă de măsurători obiective și percepții subiective.

Prin analiza acestor date, proiectul urmărește să dezvolte un model avansat care să optimizeze reglajele sistemelor HVAC în timp real, contribuind la o gestionare mai eficientă a energiei și la crearea unui mediu interior plăcut.

Structura lucrării este organizată astfel:

Capitolul 1 introduce analiza exploratorie a datelor, incluzând analiza univaria(n)tă, bivariată și multivariată. Se prezintă distribuția variabilelor, corelațiile dintre acestea și vizualizările relevante pentru înțelegerea setului de date.

Capitolul 2 detaliază metodele de preprocesare a datelor utilizate, incluzând curățarea datelor, interpolarea valorilor lipsă și eliminarea anomaliilor pentru a asigura o bază de date curată și coerentă.

Capitolul 3 se concentrează pe analiza datelor, incluzând analiza statistică și vizualizările realizate pentru a înțelege relațiile dintre variabile și structura datelor.

Capitolul 4 detaliază metodologia utilizată în modelare, explicând tehnicile aplicate, evaluarea performanței modelelor predictive și selecția celui mai performant model.

Capitolul 5 prezintă concluziile generale ale proiectului, oferă o sinteză a rezultatelor obținute și propune direcții viitoare de dezvoltare pentru îmbunătățirea soluției.

Această lucrare nu doar că prezintă o soluție inovatoare pentru optimizarea sistemelor HVAC, dar oferă și o bază solidă pentru cercetări ulterioare.

1.2 Obiective

Obiectivele proiectului sunt următoarele:

Dezvoltarea unui model predictiv avansat pentru estimarea confortului termic, folosind variabile climatice obiective (ex. temperatura, umiditatea) și date subiective din baza de date ASHRAE.

Analiza factorilor determinanți ai confortului termic, precum temperatura, umiditatea, viteza aerului și alte condiții de mediu, pentru a înțelege modul în care aceștia influențează percepția utilizatorilor.

Optimizarea performanței modelelor predictive prin aplicarea mai multor algoritmi de învățare automată (regresie liniară, regresii Ridge și Lasso, Random Forest, XGBoost).

Evaluarea performanței modelelor utilizând metrici precum MSE (Mean Squared Error), R^2 Score și, pentru clasificare, acuratețea, precizia și F1-score.

Reducerea consumului energetic al sistemelor HVAC, utilizând modele predictive pentru reglarea automată a parametrilor sistemului, astfel încât să fie adaptate nevoilor utilizatorilor și condițiilor de mediu.

Integrarea unei metodologii reproductibile, care să poată fi implementată în alte contexte sau sisteme similare.

Aceste obiective subliniază abordarea sistematică a proiectului, de la analiza datelor la implementarea practică a unui sistem HVAC inteligent.

1.3 Specificații

Proiectul urmărește să îndeplinească următoarele specificații tehnice și funcționale:

1.3.1 Funcționalități principale:

- Dezvoltarea unui model predictiv pentru estimarea PMV (Predicted Mean Vote).
- Analizarea datelor climatice și subiective pentru identificarea relațiilor semnificative.
- Optimizarea reglajelor sistemului HVAC pe baza predicțiilor modelului.

1.3.2 Mediul de implementare:

- Limbaj de programare: Python 3.9.
- Biblioteci utilizate: pandas, numpy, matplotlib, scikit-learn, xgboost.
- Mediu de dezvoltare: Jupyter Notebook.
- Date: Baza de date ASHRAE Global Thermal Comfort Database II, disponibilă în format CSV.

1.3.3 Performanță și calitate:

- Modelele predictive vor fi evaluate utilizând metrici statistice pentru a asigura o precizie ridicată.
- Sistemul trebuie să ofere predicții în timp util pentru a putea fi aplicat în medii practice.

1.3.4 Limitări:

- Performanța modelelor depinde de calitatea și volumul datelor.
- Datele utilizate sunt limitate la înregistrările disponibile în baza ASHRAE, ceea ce poate afecta generalizarea în alte contexte.

1.3.5 Fiabilitate și securitate:

- Codul este structurat și documentat pentru a fi reproductibil.
- Proiectul respectă bunele practici în manipularea și stocarea datelor pentru a preveni pierderea sau compromiterea acestora.

Prin respectarea acestor specificații, proiectul oferă o soluție robustă, reproductibilă și eficientă pentru optimizarea confortului termic, cu potențial de aplicare în diverse medii comerciale și rezidențiale.

2 Cunoașterea și analiza setului de date

2.1 Definiție

Confortul termic reprezintă starea de satisfacție a ocupanților unui spațiu în raport cu mediul termic. Este influențat de factori subiectivi, precum percepțiile și preferințele individuale, dar și de factori obiectivi, cum ar fi temperatura aerului, umiditatea relativă și viteza aerului.

2.2 Introducere

În cadrul acestui proiect, am utilizat datele disponibile în **arhiva 1** a bazei de date **ASHRAE Global Thermal Comfort Database II**, care conține informații despre confortul termic colectate în diverse locații și condiții între 1995 și 2015. Setul de date include informații esențiale pentru înțelegerea relației dintre mediul interior și percepțiile subiective ale ocupanților. Cele mai relevante caracteristici selectate pentru analiza și modelarea ulterioară includ:

Caracteristici obiective ale mediului interior:

- *Air temperature (C)*: Temperatura aerului.
- *Relative humidity (%)*: Umiditatea relativă.
- *Air velocity (m/s)*: Viteza aerului.
- *Radiant temperature (C)*: Temperatura radiantă.
- *Operative temperature (C)*: Temperatura operativă.

Indici termici calculați:

- *PMV (Predictive Mean Vote)*: Indicele de confort termic prezis.
- *PPD (Predicted Percentage of Dissatisfied)*: Procentajul prezis de persoane nemulțumite.
- *SET (Standard Effective Temperature)*: Temperatura standard echivalentă.
- Preferințe și percepții ale ocupanților:
- *Thermal sensation*: Senzația termică percepută.
- *Thermal preference*: Preferințele ocupanților în ceea ce privește confortul termic.
- *Air movement preference*: Preferința pentru mișcarea aerului.

Scopul analizei este de a înțelege structura datelor, de a verifica calitatea acestora și de a identifica factorii principali care influențează confortul termic. Această etapă este esențială pentru a selecta variabilele relevante ce vor fi utilizate în modelele predictive ulterioare.

2.3 Mediu de implementare

Pentru implementarea și analiza setului de date utilizat în acest proiect, am folosit un mediu tehnic bine configurat, ce include atât biblioteci Python esențiale pentru prelucrarea datelor și modelare, cât și un mediu de dezvoltare interactiv.

Am utilizat biblioteca **pandas** pentru manipularea și prelucrarea eficientă a datelor. Aceasta mi-a permis să citesc, să filtrez și să restructurez setul de date în funcție de cerințele analizei. **NumPy** a fost folosit pentru efectuarea de calcule numerice și operații matriciale, fiind indispensabil în etapele de pregătire a datelor și în implementarea algoritmilor de învățare automată.

Vizualizările grafice au fost realizate cu ajutorul bibliotecilor **matplotlib** și **seaborn**, care au permis crearea de grafice clare și intuitive pentru interpretarea relațiilor dintre variabile. Aceste vizualizări au fost fundamentale în etapa de analiză exploratorie a datelor, oferind o perspectivă vizuală asupra distribuțiilor și corelațiilor dintre caracteristici.

Pentru partea de modelare a datelor, am utilizat biblioteca **scikit-learn**, care a furnizat funcționalitățile necesare pentru implementarea unor modele predictive precum regresia liniară, regresia logistică, Random Forest și Gradient Boosting.

Codul a fost scris și executat folosind **Python 3.13**, într-un mediu de dezvoltare **Jupyter Notebook**, care oferă un spațiu interactiv pentru rularea secvențială a codului și vizualizarea imediată a rezultatelor. Sistemul de operare utilizat a fost **Windows 10**, configurat cu toate pachetele și dependențele necesare pentru buna desfășurare a procesului.

Acest mediu de lucru a permis o implementare eficientă și o analiză detaliată, asigurând condițiile necesare pentru obținerea rezultatelor dorite.

2.4 Citirea și explorarea inițială a datelor

Importul setului de date

Pentru încărcarea și explorarea inițială a setului de date, am utilizat biblioteca **pandas**, care permite manipularea eficientă a datelor. Codul utilizat pentru citirea fișierului CSV și afișarea primelor rânduri este:

```
import pandas as pd
# Citirea datelor din arhiva 1
ashrae_data = pd.read_csv("C:\\Users\\asus\\Desktop\\SBCproiect\\archive1\\ashrae_db2.01.csv")
# Afișarea primelor rânduri pentru a verifica structura
print(ashrae_data.head())
```


Dimensiunea și structura datelor

Pentru a înțelege mai bine structura setului de date, am analizat numărul total de rânduri și coloane, precum și tipurile de date din fiecare coloană.

```
print(f"Număr de rânduri și coloane ASHRAE: {ashrae_data.shape}")  
print(ashrae_data.info())
```

Rezultatele obținute au fost : număr total de rânduri și coloane: 107,583 rânduri și 70 coloane ; variabile numerice (float64, int64) cum ar fi PMV, PPD, SET, Air temperature (C), Relative humidity (%), Air velocity (m/s), și variabile categorice (object) Season, Climate, Building type, Cooling strategy, Sex,

Verificarea calității datelor

Înainte de analiza detaliată a datelor, este important să verificăm existența valorilor lipsă și a valorilor duplicate.

```
# Verificarea valorilor lipsă  
print("Valori lipsă în setul de date:")  
print(ashrae_data.isnull().sum())  
  
# Verificarea valorilor duplicate  
print(f"Numărul de duplicate în setul de date: {ashrae_data.duplicated().sum()}")
```

Rezultate:

Valori lipsă: Setul de date conține valori lipsă în mai multe coloane, printre care: Publication (Citation): 1,655 valori lipsa ,Air velocity (m/s): 9,143 valori lipsa ,PMV: 13,780 valori lipsa ,SET: 15,000 valori lipsa ,Outdoor monthly air temperature (C): 28,245 valori lipsa , Door, Heater: peste 90,000 valori lipsa;

Valori duplicate: 5,429 rânduri duplicate au fost identificate în setul de date.

2.5 Analiza Univariată

2.5.1 Scopul analizei univariate

Analiza univariată reprezintă primul pas în explorarea setului de date și are ca scop înțelegerea distribuției fiecărei variabile luate individual. Aceasta ajută la detectarea unor posibile valori aberante, dezechilibre în distribuție și tipare care pot influența rezultatele analizei ulterioare.

Pentru această etapă, am analizat **variabilele numerice** utilizând histograme și grafice KDE, iar variabilele categorice au fost examinate prin countplot-uri pentru a observa frecvențele fiecărei categorii.

2.5.2 Distribuția variabilelor numerice

Scopul analizei distribuției variabilelor numerice este de a identifica tipare generale, cum ar fi media, variația sau posibilele valori aberante, care ar putea influența modelele predictive.

Pentru a vizualiza distribuția variabilelor numerice, s-au generat **histograme** și **grafice KDE** (Kernel Density Estimation). Aceste grafice permit observarea formei distribuției și identificarea eventualelor anomalii.

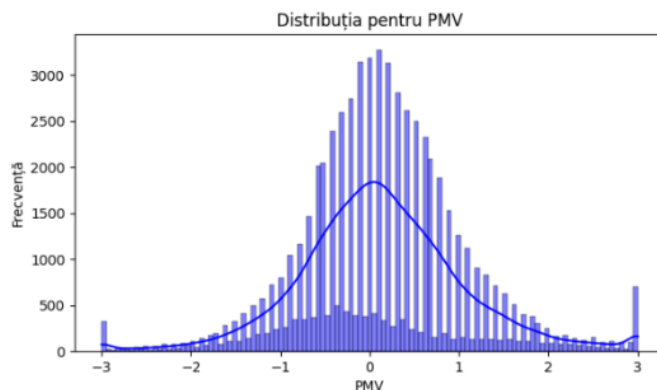


Figura 2.5.2 a) Distribuția pentru PMV

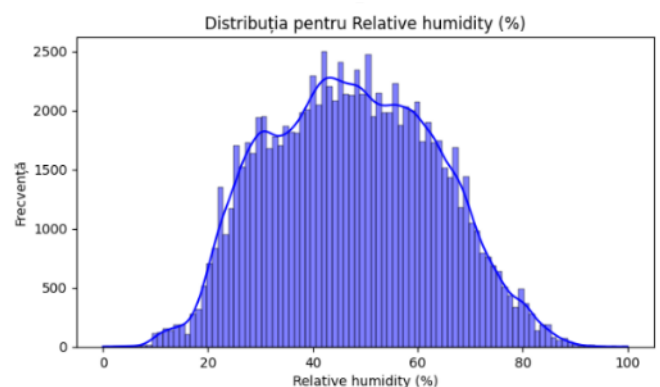


Figura 2.5.2 b) Distribuția pentru Umiditate

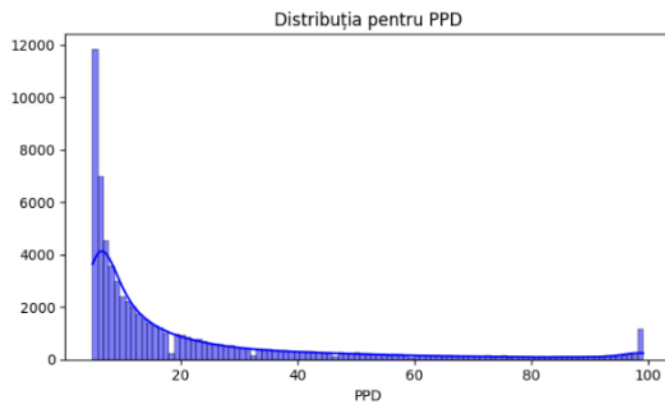


Figura 2.5.2 c). Distribuția pentru PPD (Procentul prezis de persoane nemulțumite)

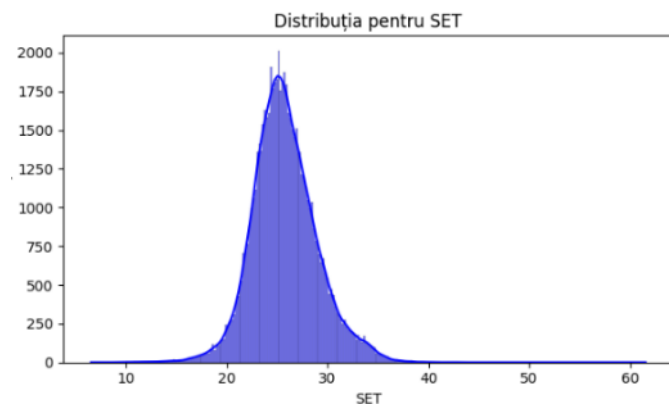


Figura 2.5.2 d). Distribuția pentru SET (Temperatura efectivă standard)

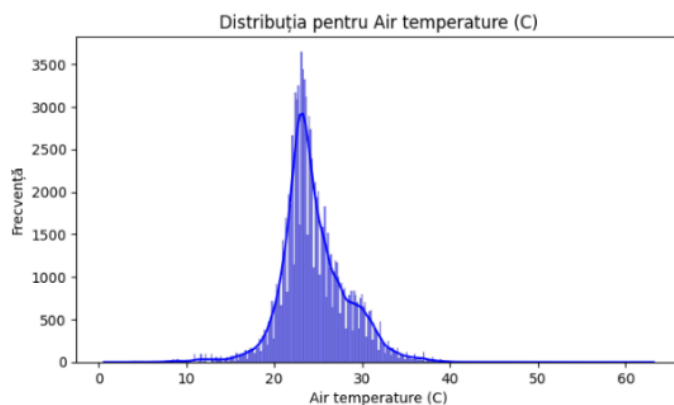


Figura 2.5.2 e). Distribuția pentru Temperatura

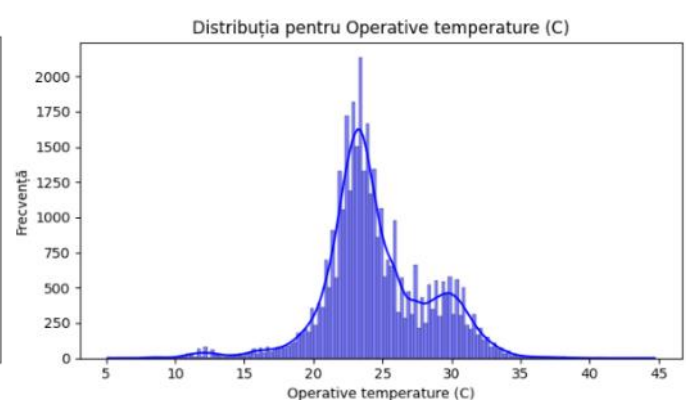


Figura 2.5.2 f). Distribuția pentru Temperatura

Se observă că distribuția variabilei PMV este aproape normală, cu cele mai multe valori în jurul intervalului de confort (-0.5 și 0.5).

Variabilele precum 'SET' și 'Air temperature (C)' prezintă o distribuție concentrată în jurul valorii medii, indicând o consistență în măsurătorile efectuate.

2.5.3 Distribuția variabilelor categorice

Analiza distribuției variabilelor categorice ajută la identificarea tiparelor și a frecvenței categoriilor distincte, oferind informații despre structura și echilibrul setului de date. Aceasta este esențială pentru a înțelege cum anumite caracteristici categorice, cum ar fi sezonul sau tipul clădirii, influențează variabilele numerice sau ținta modelului.

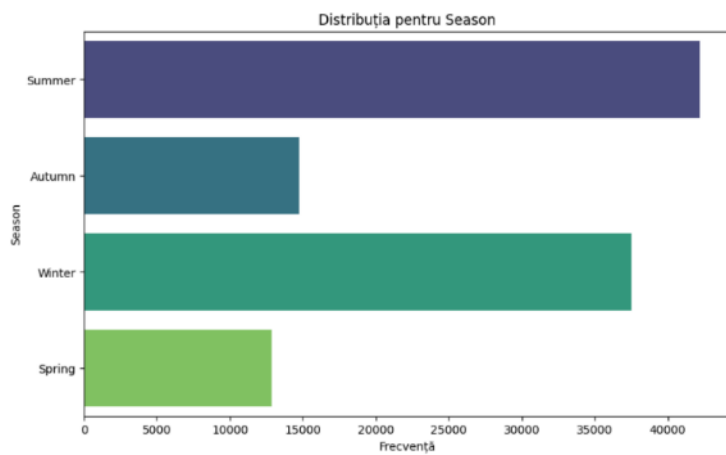


Figura 2.5.3 a). Distribuția pentru Sezon

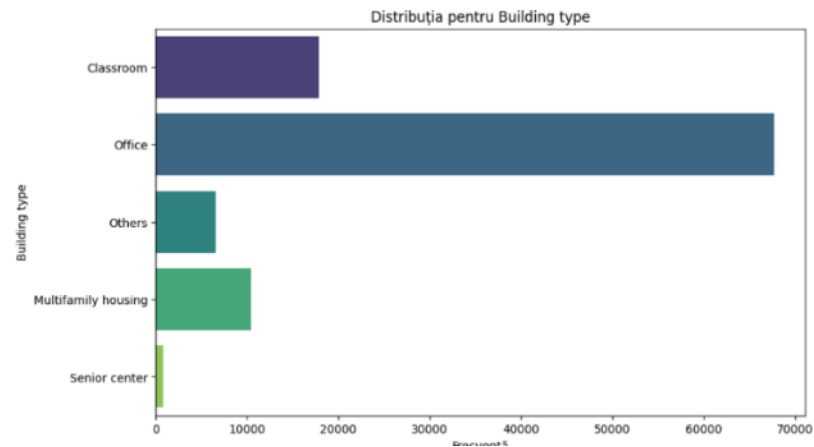


Figura 2.5.3 b). Distribuția pentru Tipul Clădirii

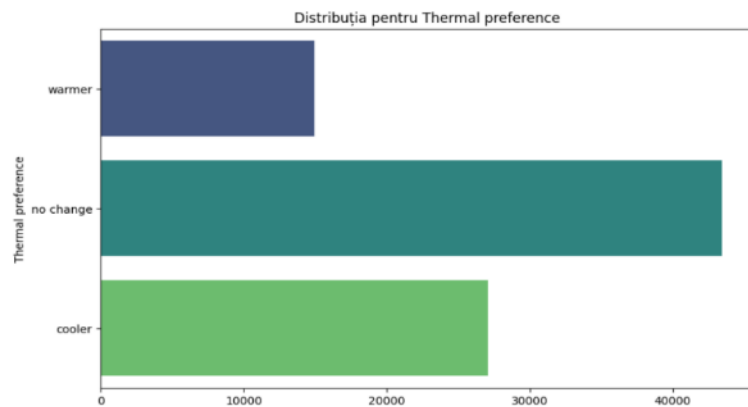


Figura 2.5.3 c). Distribuția pentru Preferința Termică

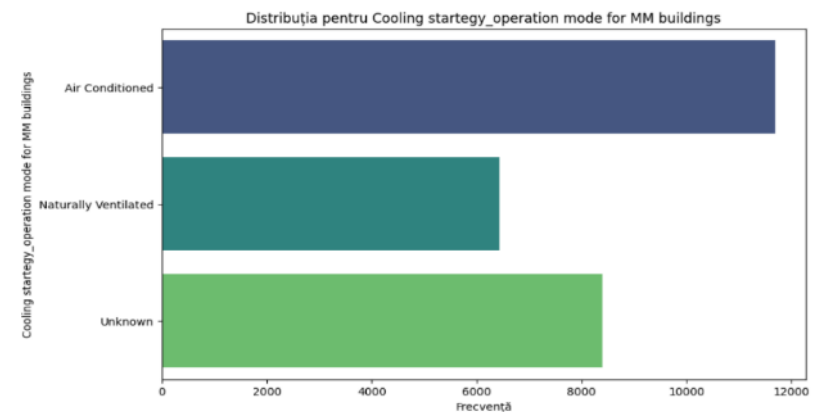


Figura 2.5.3 d). Distribuția pentru Modul de Operare al Strateției de Răcire

Distribuția variabilelor categorice evidențiază o predominanță a datelor din sezoanele extreme (vară și iarnă) și clădirile de birouri, majoritatea participanților percepend mediul ca fiind confortabil (fără preferințe de schimbare termică). De asemenea, sistemele de aer condiționat sunt cele mai frecvent utilizate, reflectând tendința generală de optimizare a confortului termic în spații controlate.

2.6 Analiza bivariata

Scopul analizei bivariata este de a explora relațiile dintre două variabile, în special pentru a identifica corelații sau tipare care ar putea influența rezultatele modelelor predictive. Pentru aceasta, s-au generat grafice de dispersie (scatter plots) pentru a vizualiza legătura dintre variabilele de interes.

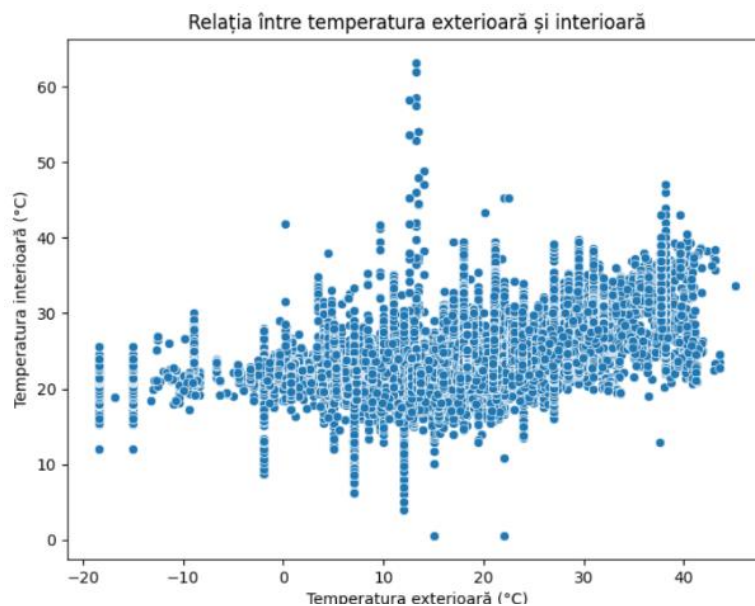


Figura 2.6 a). Relația între temperatura exterioară și temperatura interioară

Graficul arată o corelație pozitivă, ceea ce indică faptul că pe măsură ce temperatura exterioară crește, temperatura interioară tinde să crească, dar într-o manieră controlată.

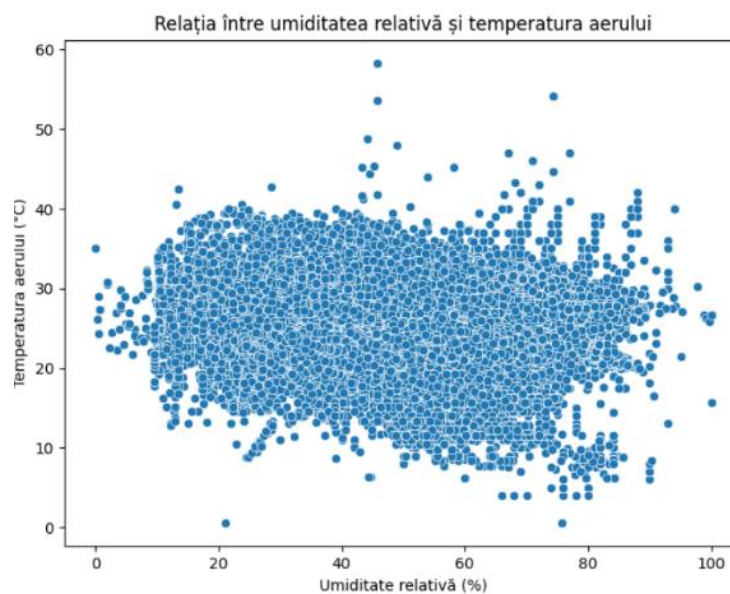


Figura 2.6 b). Relația între umiditatea relativă și temperatura aerului

Graficul nu evidențiază o corelație puternică între umiditatea relativă și temperatura aerului. Distribuția datelor sugerează o variabilitate mare a umidității la diferite temperaturi.

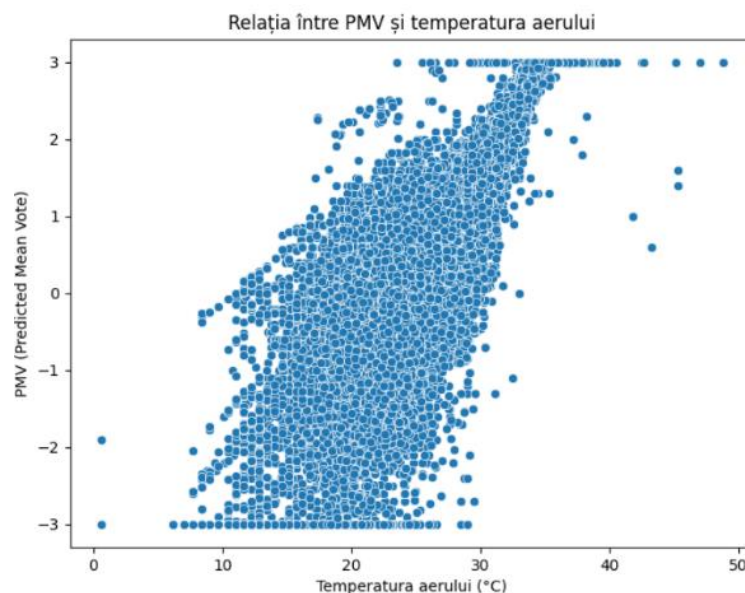


Figura 2.6 c). Relația între PMV și temperatura aerului

Graficul arată o relație clară, aproape liniară, între PMV și temperatura aerului. Pe măsură ce temperatura aerului crește, PMV devine mai mare, indicând o percepție de confort termic redus la temperaturi mai mari.

Aceste observații sunt utile pentru a înțelege cum variabilele cheie interacționează între ele, oferind o bază pentru selectarea variabilelor relevante pentru modelele predictive.

2.7 Analiza multivariată

Analiza multivariată este utilizată pentru a explora relațiile complexe dintre mai multe variabile simultan, identificând modele și structuri ascunse în date. În cadrul acestui proiect, am aplicat două tehnici principale de reducere a dimensionalității (PCA și t-SNE) și am analizat matricea de corelație extinsă pentru a înțelege mai bine interacțiunile dintre variabilele numerice.

Pentru a înțelege structura datelor și pentru a vizualiza în două dimensiuni distribuția sezonală a variabilelor, s-au aplicat următoarele metode:

2.7.1 PCA (Principal Component Analysis):

PCA este o tehnică liniară care reduce dimensionalitatea datelor, păstrând cât mai multă variație posibilă. Prima componentă principală explică cel mai mult din variația datelor, iar a doua componentă principală explică variația rămasă.

În acest proiect, PCA a fost aplicat pe un set de variabile numeric standardizate pentru a vizualiza distribuția în funcție de sezon.

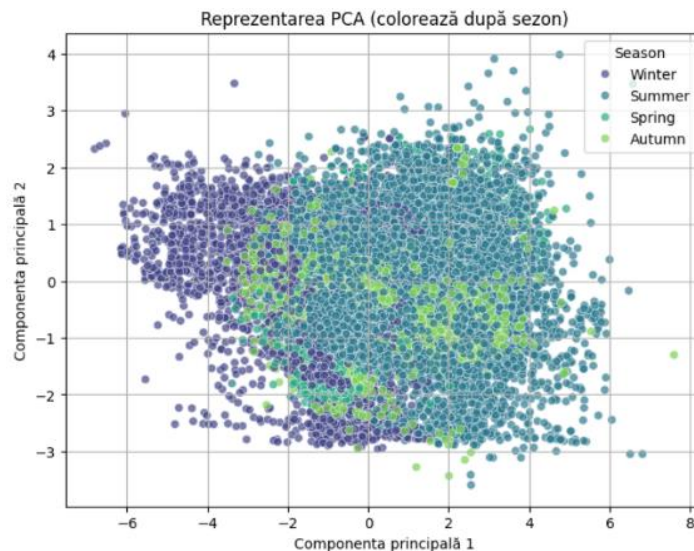


Figura 2.7.1 Reprezentarea PCA

Reprezentarea PCA arată o suprapunere a datelor pentru diferite sezoane, ceea ce sugerează că variația principală a datelor nu este strict influențată de anotimpuri. Totuși, această tehnică a permis o reducere semnificativă a dimensionalității datelor, păstrând structura general

2.7.2 t-SNE (t-distributed Stochastic Neighbor Embedding):

t-SNE este o tehnică de reducere a dimensionalității utilizată pentru vizualizarea relațiilor complexe într-un spațiu de dimensiuni reduse. Aceasta este mai potrivită pentru analiza clusterelor în date complexe.

Am utilizat un subset al datelor pentru a reduce timpul de calcul și am colorat rezultatele în funcție de anotările sezoniere.

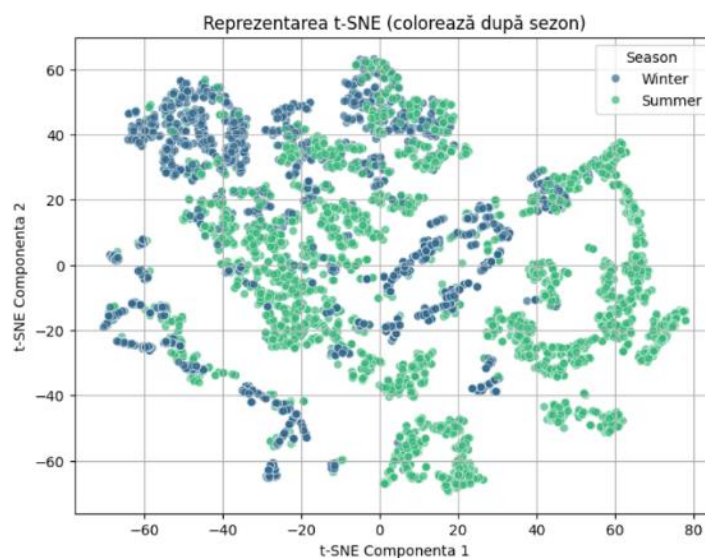


Figura 2.7.2 Reprezentarea t-SNE

Vizualizarea t-SNE evidențiază grupuri mai distincte între anumite sezoane, cum ar fi vara și iarna. Aceasta indică faptul că există variații sezoniere în datele colectate, dar și o anumită suprapunere între aceste grupuri.

2.7.3 Matricea de corelație

Pentru a evalua relațiile directe dintre variabilele numerice, am generat o matrice de corelație utilizând coeficientul Pearson. Această matrice oferă informații esențiale despre cât de puternică este legătura dintre două variabile.

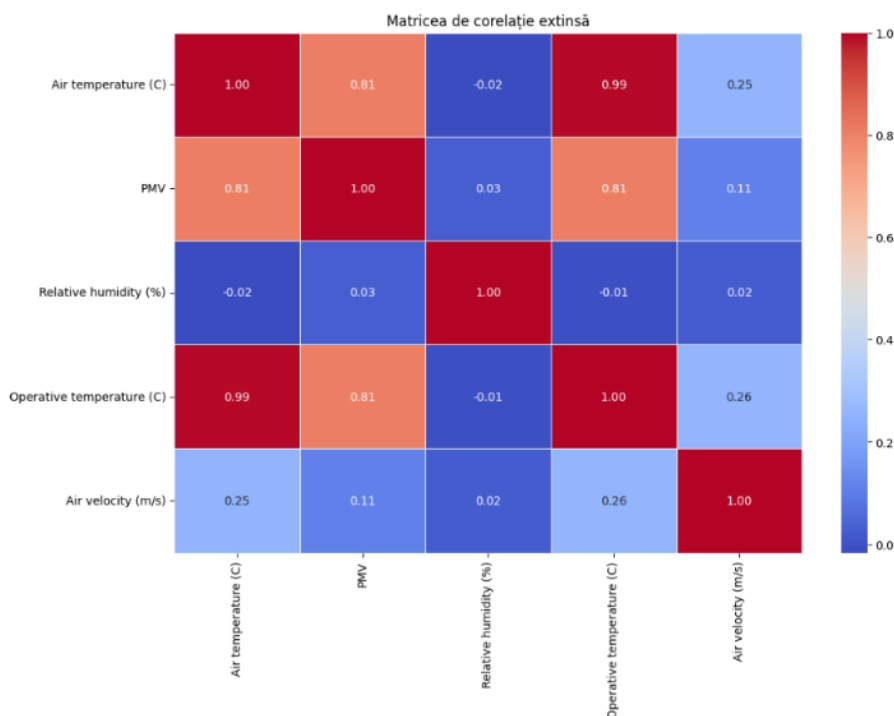


Figura 2.7.2 Matricea de corelație extinsă

Se observă o corelație foarte puternică între Air temperature (C) și Operative temperature (C) (coeficient de corelație = 0.99), ceea ce confirmă că temperaturile operative și aerului sunt aproape echivalente. În schimb, Relative humidity (%) prezintă corelații foarte slabe cu alte variabile, indicând o influență redusă asupra măsurătorilor legate de temperatură sau confort termic

Se observă o corelație foarte puternică între Air temperature (C) și Operative temperature (C) (coeficient de corelație = 0.99), ceea ce confirmă că temperaturile operative și aerului sunt aproape echivalente. În schimb, Relative humidity (%) prezintă corelații foarte slabe cu alte variabile, indicând o influență redusă asupra măsurătorilor legate de temperatură sau confort termic

3 Pre-procesarea setului de date

3.1 Introducere

Pre-procesarea datelor reprezintă o etapă esențială în cadrul oricărui proiect de analiză sau modelare a datelor, având scopul de a asigura calitatea și integritatea setului de date. Această etapă implică aplicarea unor metode și tehnici pentru eliminarea zgomotului, completarea valorilor lipsă, reducerea dimensionalității, eliminarea valorilor aberante și optimizarea caracteristicilor relevante.

Scopul pre-procesării este de a transforma datele brute într-un format adecvat pentru analize ulterioare și modelare predictivă. Aceasta nu doar că îmbunătățește performanța modelelor predictive, dar și asigură robustețea acestora prin reducerea efectelor erorilor, tendințelor sau valorilor extreme.

3.2 Denoising (Eliminarea zgomotului)

3.2.1 Definiție și scop

Denoising, sau eliminarea zgomotului, reprezintă procesul de reducere sau eliminare a variațiilor nedorite dintr-un set de date, care pot afecta acuratețea modelelor predictive sau analizele ulterioare. Zgomotul poate apărea din cauza erorilor de măsurare, fluctuațiilor neașteptate sau a altor factori externi. Scopul principal al denoising-ului este de a îmbunătăți calitatea datelor, păstrând semnalele relevante și eliminând variațiile inutile.

3.2.2 Metode utilizate

Eliminarea zgomotului se realizează prin aplicarea diverselor filtre, fiecare adaptat unui anumit tip de date sau zgomot. Aceste filtre sunt utile pentru a netezi valorile variabilelor numerice și pentru a reduce incertitudinea datelor, permițând o interpretare mai clară.

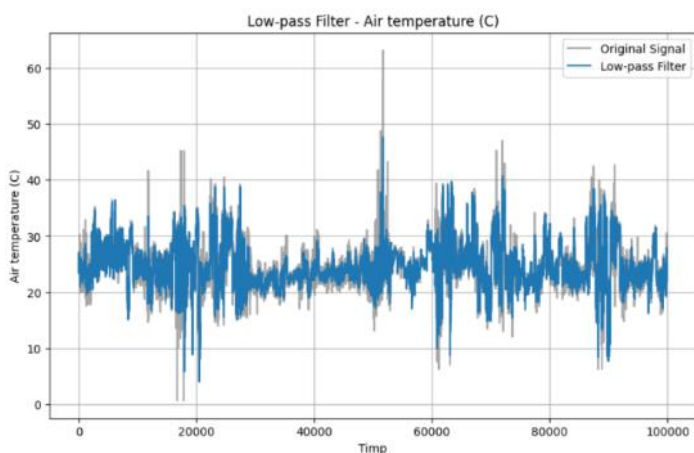


Figura 3.2.2 a). Low-pass Filter - Air temperature

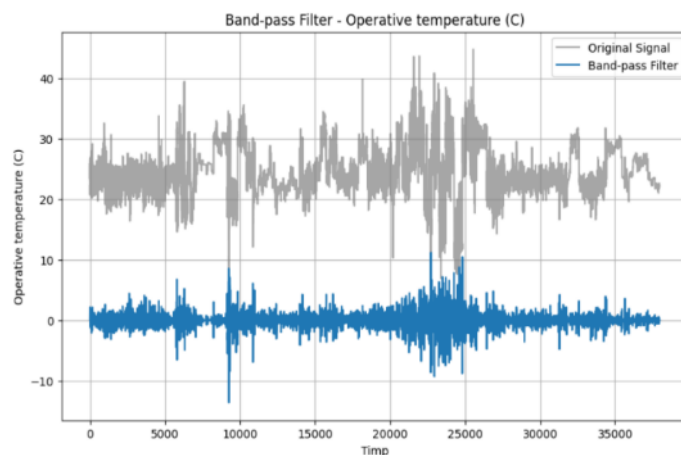


Figura 3.2.2 b). Band-pass Filter - Operative

Low-pass Filter - Air temperature (C)

Filtrul low-pass elimină variațiile de înaltă frecvență pentru a evidenția tendința generală a temperaturii aerului.

Semnalul filtrat arată o reducere clară a fluctuațiilor rapide, evidențiind variațiile lente ale temperaturii.

Band-pass Filter - Operative temperature (C)

Filtrul band-pass păstrează doar componentele de frecvență dintr-un interval specific, eliminând zgomotul de joasă și înaltă frecvență.

Semnalul rezultat scoate în evidență schimbările moderate ale temperaturii operative, eliminând zgomotul extrem.

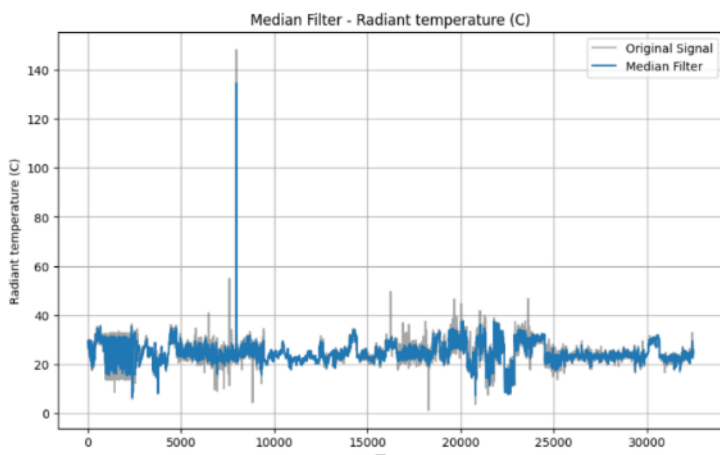


Figura 3.2.2 c). Median Filter - Radiant

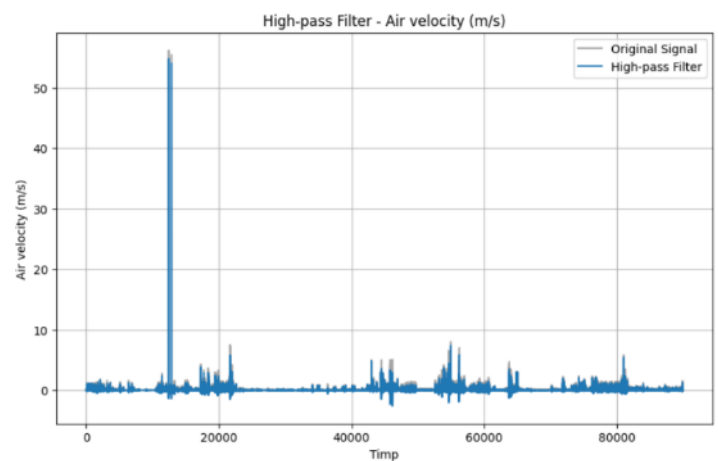


Figura 3.2.2 d). High-pass Filter - Air velocity

Median Filter - Radiant temperature (C)

Filtrul median este eficient pentru eliminarea valorilor aberante sau a zgomotului de tip "sare și piper."

Se observă o netezire semnificativă a semnalului, în special eliminarea vârfurilor extreme.

High-pass Filter - Air velocity (m/s)

Filtrul high-pass este utilizat pentru a păstra variațiile rapide și a elimina tendințele lente.

Zgomotul de frecvență joasă este eliminat, iar semnalul rezultat evidențiază variațiile rapide ale vitezei aerului.

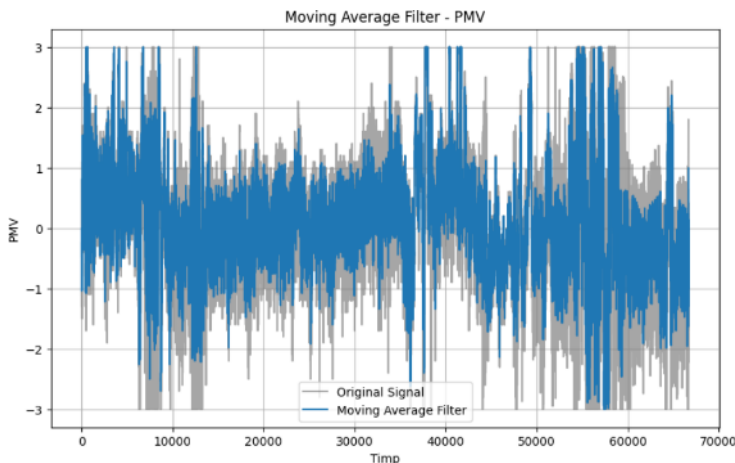


Figura 3.2.2 e). Moving Average Filter - PMV

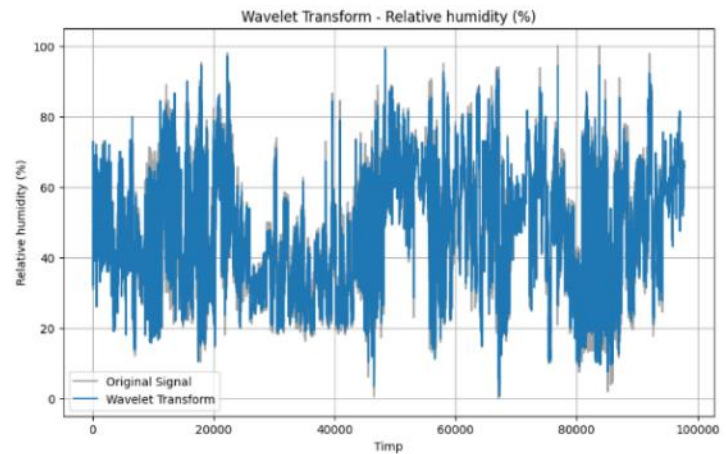


Figura 3.2.2 e). Wavelet Transform
Relative humidity (%)

Moving Average Filter - PMV

Filtrul de medie mobilă netezește semnalul pentru a elimina fluctuațiile mici, păstrând tendința generală.

Semnalul filtrat este mai neted, iar variațiile bruște ale valorilor PMV sunt reduse.

Wavelet Transform - Relative humidity (%)

Transformata wavelet este utilizată pentru a elimina zgomotul la mai multe scale de timp, păstrând structura semnalului.

Semnalul rezultat este netezit, iar zgomotul de înaltă frecvență este eliminat, evidențiind variabilitatea umidității relative.

3.3 Detrending (Eliminarea tendințelor)

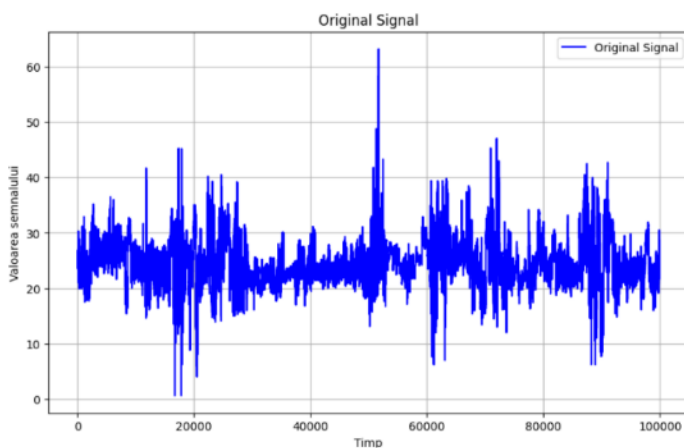
3.3.1 Definiție și scop

Detrending-ul reprezintă procesul de eliminare a tendințelor sau componentelor de ordin înalt dintr-un set de date, cu scopul de a analiza doar variațiile semnalului în jurul unei valori de bază. Tendințele pot apărea din cauza fenomenelor sezoniere, a erorilor de măsurare sau a variațiilor de lungă durată și pot masca informațiile relevante din date. Prin eliminarea acestora, putem obține un semnal mai curat, adecvat pentru analize mai detaliate.

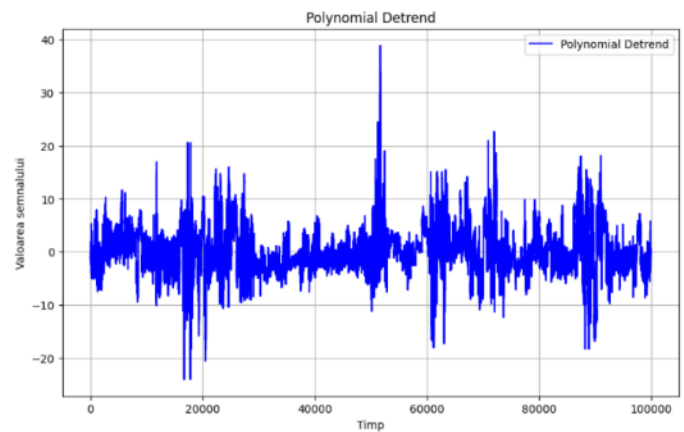
Metodele utilizate includ substrația mediei mobile, care elimină tendințele de ordin inferior prin calculul unei medii pe o fereastră glisantă. Detrending-ul polinomial este folosit pentru a îndepărta componentele de trend de ordin superior, ajustând datele cu ajutorul unor polinoame de diferite grade. Detrending-ul prin diferențiere permite extragerea variațiilor rapide prin calculul diferențelor dintre punctele succesive. De asemenea, detrending-ul folosind transformata Fourier ajută la identificarea și eliminarea componentelor de frecvență joasă, păstrând fluctuațiile mai rapide din date.

Air temperature (C):

Metodă utilizată: Detrending Polinomial pentru a elimina tendințele de ordin superior din datele privind temperatura aerului, a fost utilizată metoda detrending polinomial. Aceasta ajustează datele utilizând un polinom de gradul 2, eliminând fluctuațiile lente asociate cu trendurile.



*Figura 3.3 a) Semnal Original -
Air temperature (C)*

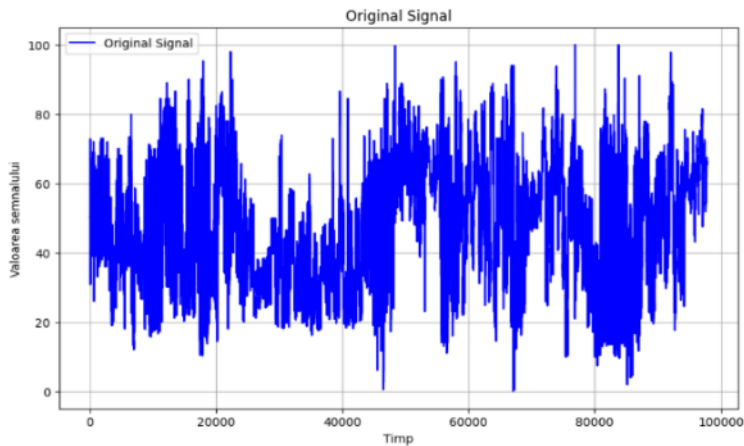


*Figura 3.3 b) Detrending Polinomial
- Air temperature (C)*

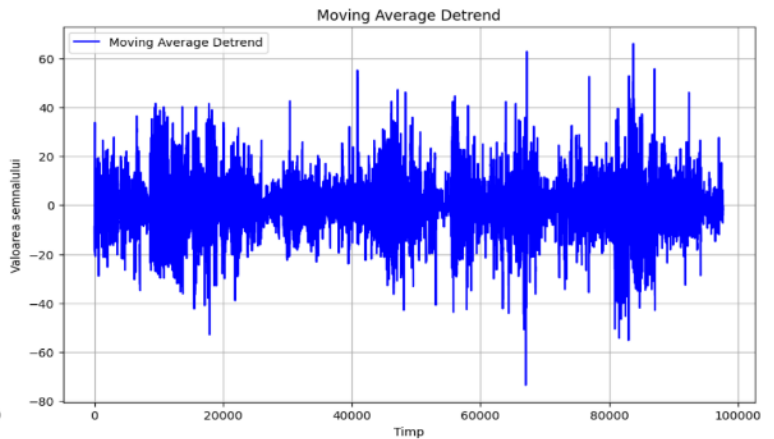
Graficul ilustrează compararea dintre semnalul original și semnalul rezultat după aplicarea detrending-ului polinomial. Se observă că variațiile de ordin inferior au fost eliminate, păstrând doar fluctuațiile rapide și relevante.

Relative humidity (%):

Metodă utilizată: Substrația Mediei Mobile pentru datele referitoare la umiditatea relativă, a fost aplicată metoda substrației mediei mobile. Aceasta elimină tendințele de ordin inferior prin calcularea unei medii pe o fereastră glisantă și scăderea acesteia din semnalul original.



*Figura 3.3 c) Semnal Original -
Relative humidity (%)*

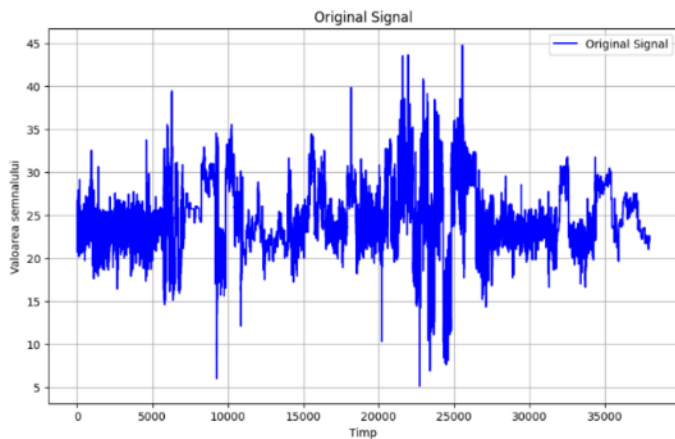


*Figura 3.3 d) Detrending prin
Substracția Mediei Mobile - Relative*

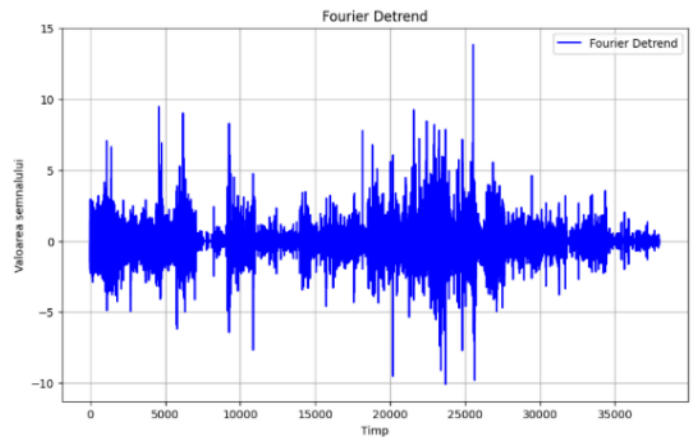
Graficul arată cum metoda substracției mediei mobile a redus componentele lente ale semnalului, păstrând variațiile rapide. Compararea cu semnalul original evidențiază eliminarea tendințelor pe termen lung.

Operative temperature (C):

Metodă utilizată: Transformata Fourier pentru temperatura operativă, a fost aplicată metoda detrending utilizând transformata Fourier. Aceasta a identificat și eliminat componentele de frecvență joasă, lăsând doar variațiile mai rapide și relevante pentru analiză.



*Figura 3.3 e) Semnal Original -
Operative temperature (C)*



*Figura 3.3 f) Detrending prin
Transformata Fourier - Operative*

Graficul afișează semnalul prelucrat după aplicarea metodei Fourier comparat cu semnalul original. Se observă că metoda a eliminat variațiile lente, păstrând detaliile rapide ale semnalului.

3.4 Eliminarea valorilor aberante (Outlier Removal)

3.4.1 Definiție și scop

Eliminarea valorilor aberante reprezintă un pas esențial în pre-procesarea datelor, având ca scop identificarea și eliminarea punctelor de date care diferă semnificativ de restul setului. Aceste valori pot apărea din cauza erorilor de măsurare, a introducerii incorecte a datelor sau a altor procese anormale, și pot afecta negativ calitatea modelelor predictive și analiza statistică.

Scopul eliminării valorilor aberante este de a îmbunătăți acuratețea și robustețea modelelor, reducând influența extremelor asupra rezultatelor analizei. Aceasta se realizează utilizând diverse metode statistice și algoritmice, fiecare fiind potrivită pentru tipuri specifice de seturi de date și structuri.

Urmează să analizăm și să aplicăm mai multe metode utilizate pentru identificarea și eliminarea valorilor aberante: Z-Score Method, IQR Method (Interquartile Range), MAD (Median Absolute Deviation), DBSCAN Method, Isolation Forest Method.

Pentru fiecare metodă, vom analiza implementarea și efectele asupra unui subset de date selectate, incluzând vizualizări grafice și interpretări ale rezultatelor.

3.4.2 Metoda Z-Score (Z-Score Method)

Metoda Z-Score este o tehnică utilizată pentru identificarea valorilor aberante (outliers) prin calcularea distanței fiecărui punct de date față de media dataset-ului, exprimată în termeni de deviații standard. Aceasta este potrivită pentru seturi de date care urmează o distribuție normală și permite detectarea valorilor extreme care se abat semnificativ de la majoritatea datelor..

Cum funcționează:

Se calculează Z-Score-ul pentru fiecare punct de date utilizând formula:

$$Z = \frac{x - \mu}{\sigma} \quad (3.4.2)$$

Unde:

x :Valoarea punctului de date.

μ : Media dataset-ului.

σ : Deviația standard a dataset-ului.

Un punct de date este considerat aberant dacă valoarea absolută a scorului Z depășește un prag prestabilit (de exemplu, $Z > 3$ sau $Z < -3$).

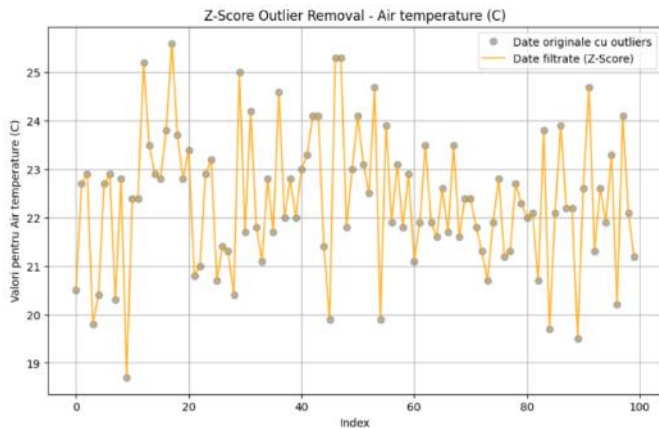


Figura 3.4.2 a) Eliminarea valorilor aberante folosind metoda Z-Score - Temperatura aerului (°C)

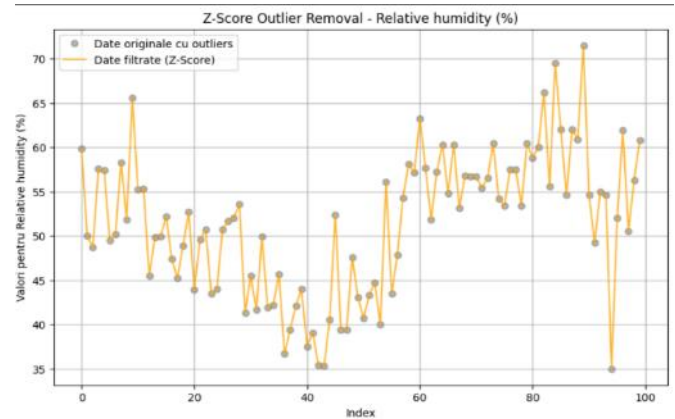


Figura 3.4.2 b) Eliminarea valorilor aberante folosind metoda Z-Score - Umiditatea relativă (%)

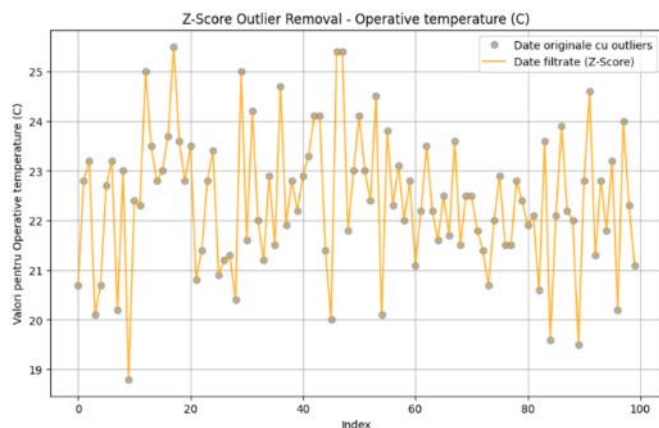


Figura 3.4.2 c) Eliminarea valorilor aberante folosind metoda Z-Score - Temperatura operativă (°C)

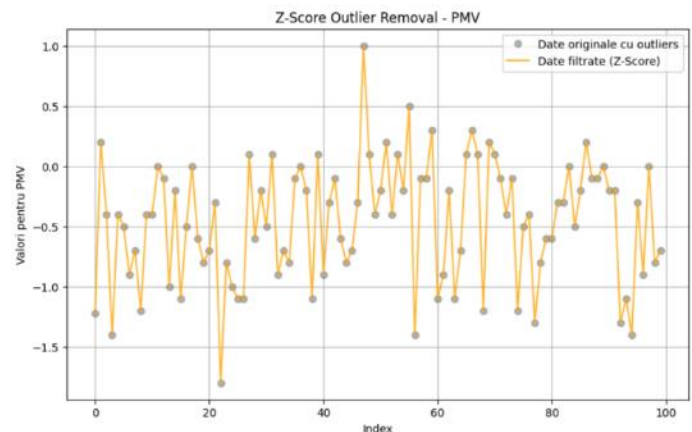


Figura 3.4.2 d) Eliminarea valorilor aberante folosind metoda Z-Score - PMV (Predicted Mean Vote)

Metoda Z-Score a fost aplicată pe cele patru coloane selectate, iar graficele evidențiază clar eliminarea valorilor aberante. Datele originale conțineau puncte care deviau semnificativ de la medie, în timp ce datele filtrate au eliminat aceste valori, rezultând un set de date mai curat și mai reprezentativ. Această abordare este ideală pentru date distribuite normal, iar rezultatele au demonstrat eficiența metodei în acest context.

3.4.3 Metoda IQR (Interquartile Range)

Metoda IQR (Interval Intercuartil) este utilizată pentru identificarea și eliminarea valorilor aberante dintr-un set de date. Aceasta se bazează pe calcularea intervalului dintre percentila 25 (Q1) și percentila 75 (Q3), iar valorile care depășesc acest interval extins cu un factor multiplicativ ($1.5 * IQR$, de obicei) sunt considerate valori aberante.

Cum funcționează:

Se calculează percentila 25 (Q1) și percentila 75 (Q3).

Intervalul intercuartil (IQR) se determină astfel: $IQR = Q3 - Q1$ (3.4.3.1)

Limitele pentru detectarea valorilor aberante sunt:

$Lower\ Bound = Q1 - 1.5 * IQR$ (3.4.3.2)

$Upper\ Bound = Q1 + 1.5 * IQR$ (3.4.3.3)

Orice punct de date care se află sub limita inferioară sau peste limita superioară este considerat o valoare aberantă.

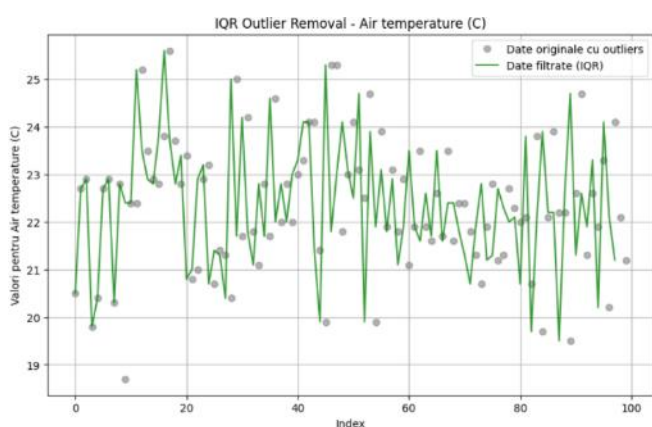


Figura 3.4.3 a) IQR Outlier Removal - Air Temperature (C)

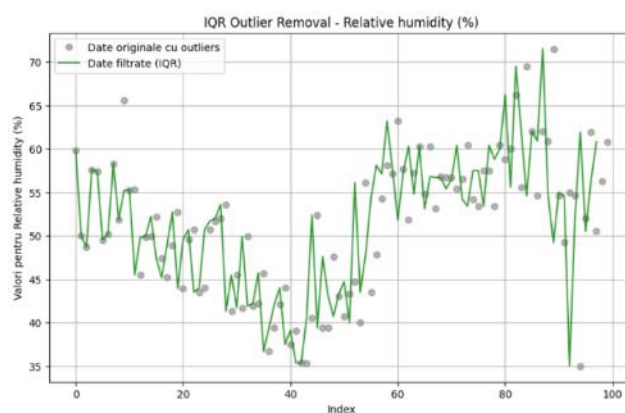


Figura 3.4.3 b) IQR Outlier Removal - Relative Humidity (%)

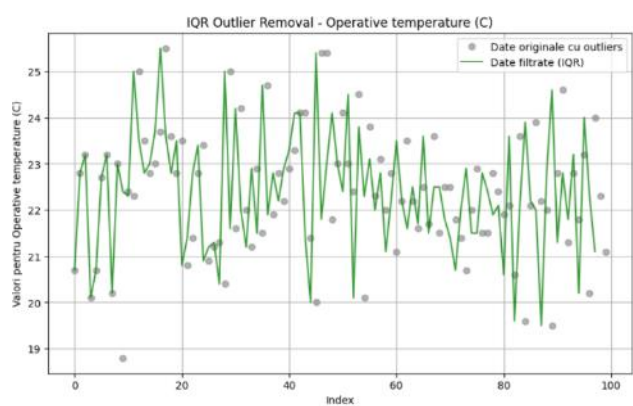


Figura 3.4.3 c) IQR Outlier Removal - Operative Temperature (C)

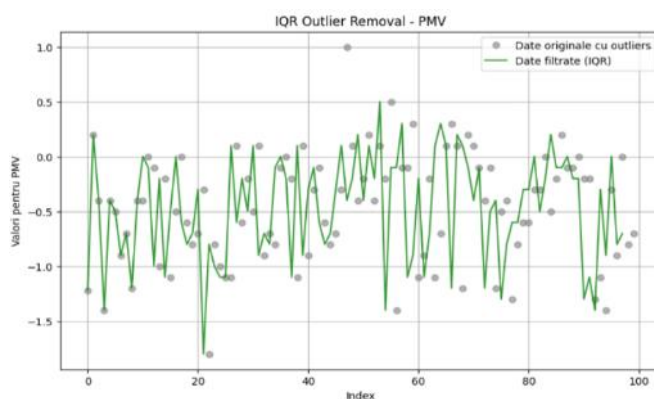


Figura 3.4.3 d) IQR Outlier Removal - PMV

Metoda IQR (Interquartile Range) a fost aplicată pentru cele patru coloane analizate, cu scopul de a elimina valorile aberante care se află în afara intervalului $[Q1 - 1.5 * IQR, Q3 + 1.5 * IQR]$. Graficele arată clar diferențele dintre datele originale și cele filtrate, evidențiind eliminarea punctelor aberante. Această metodă este deosebit de utilă pentru datele cu distribuții nesimetrice sau date cu extreme multiple, oferind un set de date mai curat și mai robust pentru analizele ulterioare.

3.4.4 MAD (Median Absolute Deviation)

Metoda MAD (Median Absolute Deviation) este utilizată pentru identificarea valorilor aberante într-un mod robust, fiind mai puțin sensibilă la extreme față de metodele bazate pe media aritmetică. Aceasta este ideală pentru seturi de date care prezintă distribuții asimetrice sau cu cozi lungi.

Cum Funcționează:

Se calculează mediana absolută a deviațiilor față de mediana datelor:

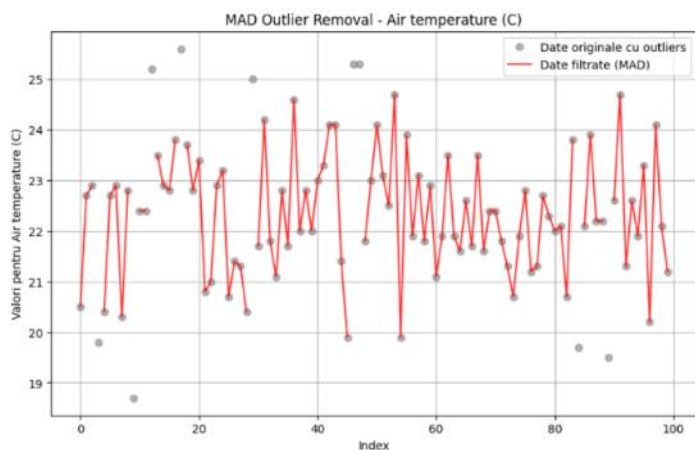
$$MAD = \text{median}(|x_i - M|) \quad (3.4.4.1)$$

unde M este mediana setului de date.

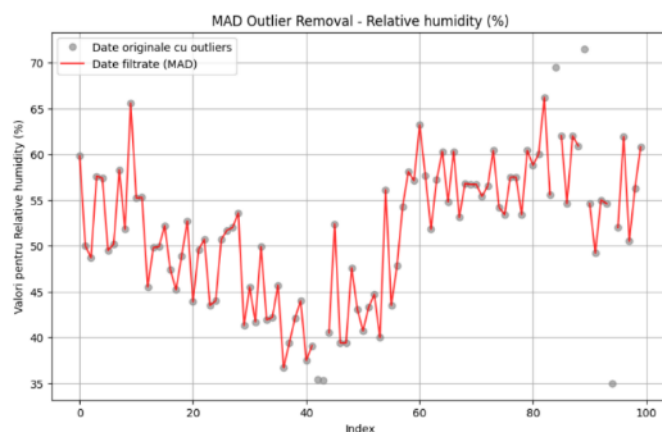
Se definește o valoare de prag kkk (de obicei 1.5 sau 3) pentru detectarea valorilor aberante. Un punct este considerat aberant dacă se află în afara intervalului:

$$[M - k * MAD, M + k * MAD] \quad (3.4.4.2)$$

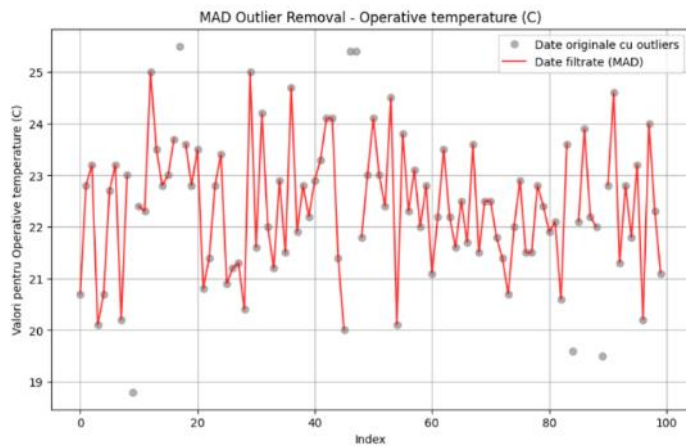
Această metodă este robustă deoarece medianele sunt mai puțin afectate de valorile extreme, comparativ cu metodele care folosesc media.



*Figura 3.4.4 a) MAD Outlier Removal
- Air temperature (C)*



*Figura 3.4.4 b) MAD Outlier Removal
- Relative humidity (%)*



*Figura 3.4.4 c) MAD Outlier Removal
- Operative temperature (C)*

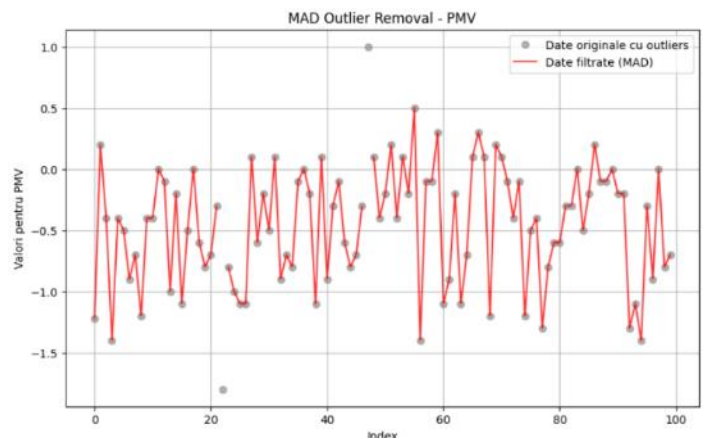


Figura 3.4.4 d) MAD Outlier Removal - PMV

Metoda MAD (Median Absolute Deviation) a fost aplicată pentru cele patru coloane selectate, demonstrând eficiența sa în identificarea și eliminarea valorilor aberante.

Această metodă este robustă împotriva valorilor extreme și funcționează bine chiar și în cazul distribuțiilor asimetrice. În graficele prezentate, se observă eliminarea outlier-ilor, rezultând un set de date mai curat, potrivit pentru analize ulterioare.

3.4.5 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) este o metodă de clasificare bazată pe densitate care identifică grupurile de puncte dense în date și consideră punctele din regiunile de densitate scăzută ca fiind valori aberante (outliers). Această metodă este utilă pentru datele care prezintă structuri spațiale sau clustere naturale.

Cum funcționează:

Puncte de bază (core points): Punctele care au un număr minim de vecini într-un anumit radius definit de parametrul epsilon.

Vecini direcți: Punctele care sunt în vecinătatea imediată a unui punct de bază.

Outliers (noise points): Punctele care nu aparțin niciunui cluster din cauza densității scăzute.

Parametri principali:

epsilon (raza maximă): Definește distanța maximă dintre două puncte pentru a fi considerate vecine.

minPts (numărul minim de puncte): Numărul minim de puncte necesare pentru a forma un cluster.

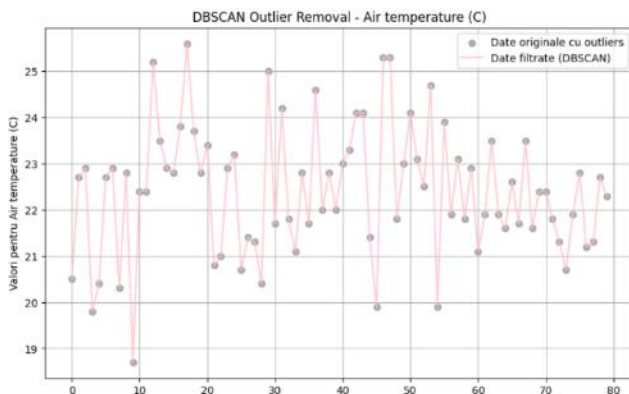


Figura 3.4.5 a) DBSCAN Outlier Removal - Air temperature (C)

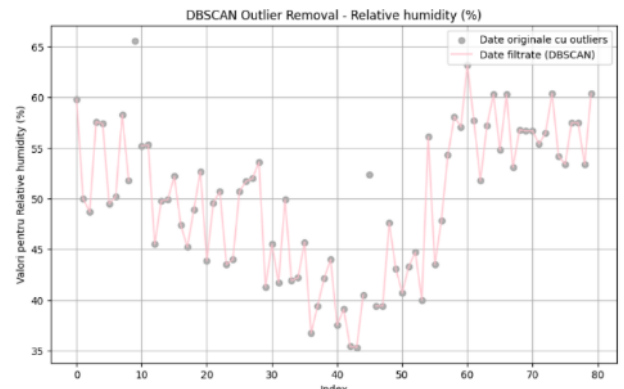


Figura 3.4.5 b) DBSCAN Outlier Removal - Relative humidity (%)

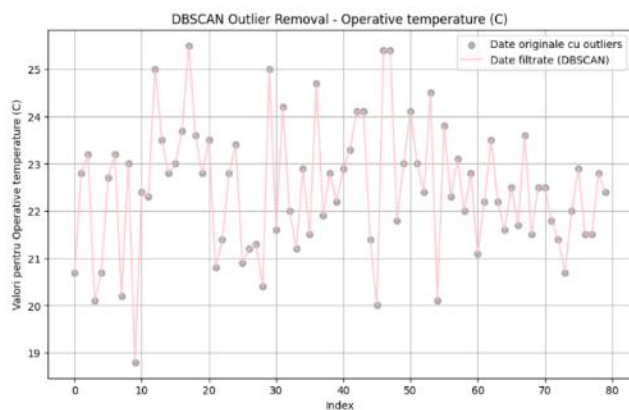


Figura 3.4.5 c) DBSCAN Outlier Removal - Operative temperature (C)

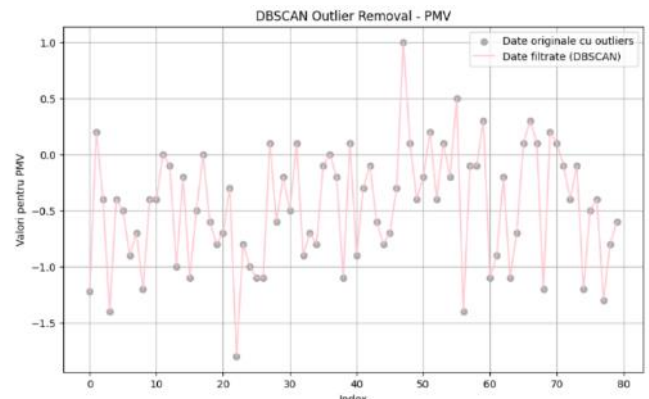


Figura 3.4.5 d) DBSCAN Outlier Removal - PMV

Metoda **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** a fost aplicată pentru identificarea și eliminarea valorilor aberante din cele patru coloane analizate. Graficele evidențiază cum această tehnică a identificat și filtrat valorile care se află în regiuni cu densitate scăzută, fiind considerate puncte aberante. Această metodă este deosebit de utilă pentru datele cu grupări naturale și distribuții neuniforme, oferind o abordare adaptivă pentru detecția valorilor aberante în funcție de parametrii densității (*eps*) și numărul minim de puncte (*min_samples*). Rezultatele arată un set de date mai curat, potrivit pentru analize ulterioare

3.4.6 Isolation Forest Method

Isolation Forest Method este o tehnică bazată pe învățarea ansamblului (ensemble learning) utilizată pentru detectarea valorilor aberante în seturi de date complexe. Această metodă izolează în mod iterativ punctele de date prin divizarea spațiului de date, utilizând arbori de decizie construiți aleatoriu.

Cum funcționează

Izolare prin partajare: Datele sunt partajate în mod repetat utilizând arbori de decizie, iar punctele care sunt izolate rapid (apar în ramuri mai îndepărtate) sunt considerate valori aberante.

Evaluarea scorului de anomalie: Fiecare punct primește un scor bazat pe cât de rapid este izolat în structura arborilor. Punctele cu un scor ridicat sunt marcate ca valori aberante.

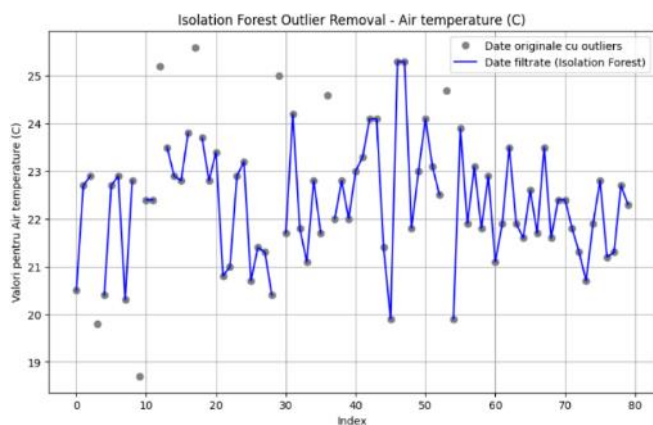


Figura 3.4.6 a) Isolation Forest Outlier Removal - Air temperature (C)

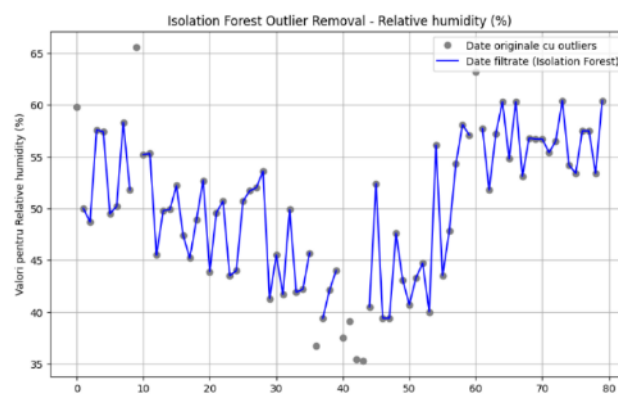


Figura 3.4.6 a) Isolation Forest Outlier Removal - Air temperature (C)

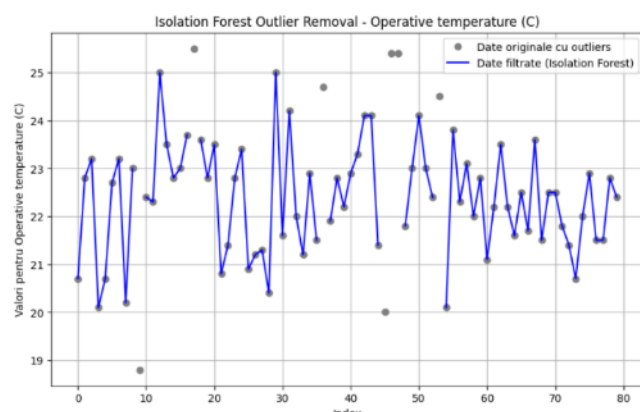


Figura 3.4.6 a) Isolation Forest Outlier Removal - Air temperature (C)

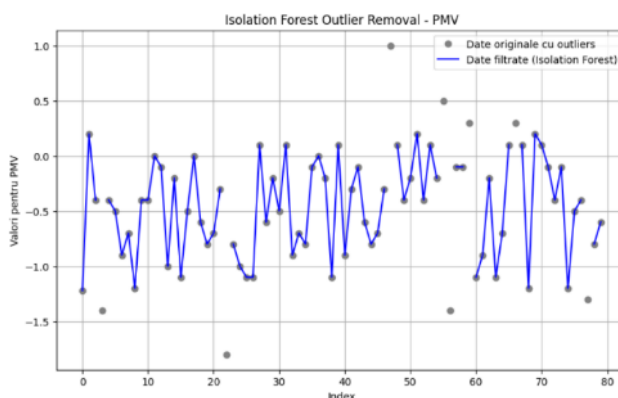


Figura 3.4.6 a) Isolation Forest Outlier Removal - Air temperature (C)

Graficele evidențiază diferențele între datele originale, care conțineau puncte izolate identificate drept outliers, și datele filtrate, care prezintă un set de valori mai uniform. Această metodă, fiind bazată pe un algoritm de tip ensemble ce utilizează păduri de decizie, este eficientă pentru seturi de date multidimensionale și poate identifica cu acuratețe punctele aberante. Rezultatele demonstrează capacitatea acestei tehnici de a îmbunătăți calitatea datelor pentru analize ulterioare.

3.5 Interpolare

3.5.1 Definiție și importanță

Interpolarea reprezintă procesul de estimare a valorilor lipsă într-un set de date, asigurând continuitatea și prevenind impactul negativ al datelor lipsă asupra analizei sau modelării.

Datele lipsă pot perturba analiza, mai ales în serii temporale sau seturi de date spațiale.

Interpolarea completează aceste lacune, facilitând o analiză mai fluidă și prevenind introducerea de erori sau biasuri.

3.5.2 Interpolare liniară

Această metodă conectează două puncte cunoscute cu o linie dreaptă, estimând valorile lipsă de-a lungul acestei linii.

Formulă utilizată:
$$y = y_1 + \frac{x - x_1}{x_2 - x_1} \times (y_2 - y_1) \quad (3.5.2)$$

Unde:

(x_1, y_1) și (x_2, y_2) sunt punctele cunoscute

x este poziția valorii lipsă.

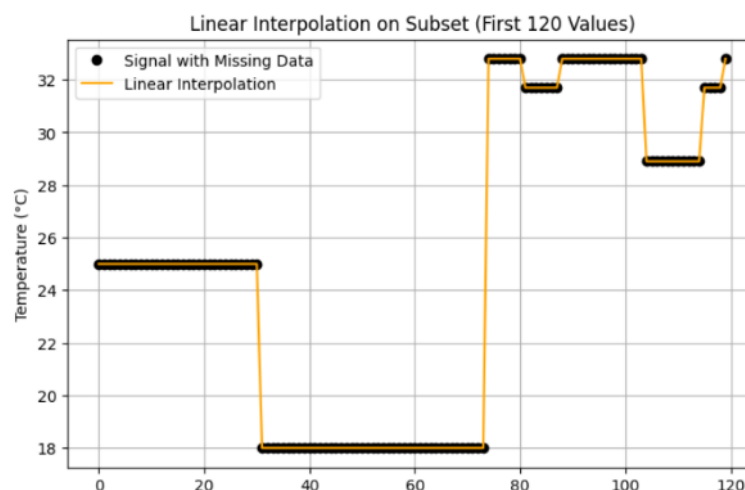


Figura 3.5.2 Linear Interpolation on Subset (First 120 Values)

Graficul arată tranziții directe și netede între punctele lipsă. Este simplu și eficient pentru serii temporale cu variații reduse

3.5.3 Interpolare cubică

Interpolarea cubică ajustează o funcție polinomială cubică între punctele cunoscute, creând o curbă continuă.

Formulă utilizată:

$$S_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3 \quad (3.5.3)$$

Unde coeficienții (a_i, b_i, c_i, d_i) asigură continuitatea derivatelor și curburii.

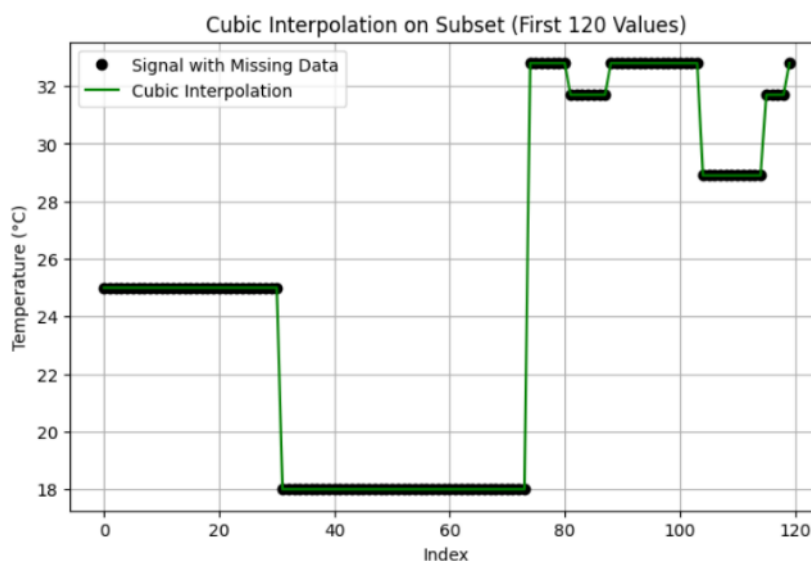


Figura 3.5.3 Cubic Interpolation on Subset (First 120 Values)

Rezultatele sunt similare interpolării liniare datorită variațiilor reduse din setul de date.

3.5.4 LOESS (Locally Estimated Scatterplot Smoothing)

LOESS aplică o regresie locală utilizând punctele dintr-un interval apropiat. Este flexibilă pentru date cu modele non-liniare.

$$\text{Formulă utilizată: } w_i = e^{-\frac{d(x, x_i)^2}{2h^2}} \quad (3.4.5)$$

Unde h este parametrul care controlează lățimea vecinătății.

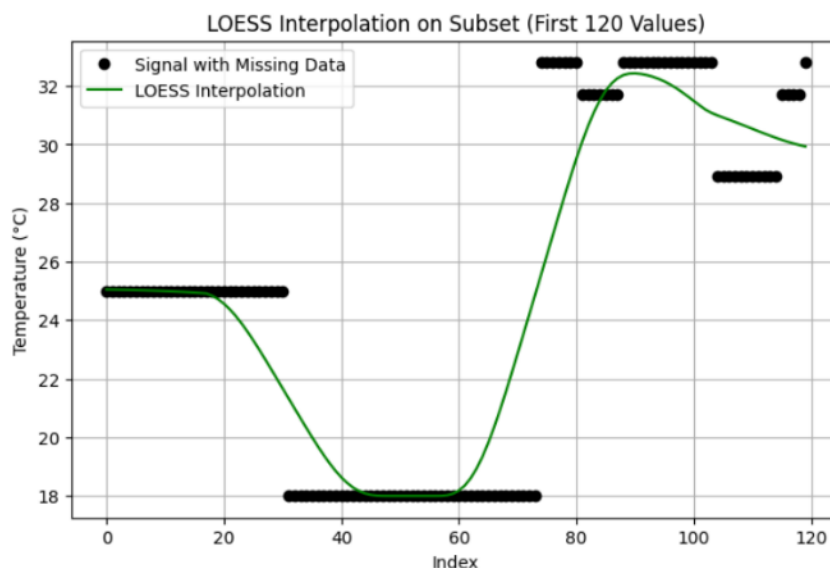


Figura 3.5.4 LOESS Interpolation on Subset (First 120 Values)

LOESS oferă un rezultat neted și flexibil, adaptându-se la variațiile locale ale datelor.

Analiza graficelor rezultate pentru metodele de interpolare evidențiază următoarele aspecte:

Interpolarea liniară și interpolarea cubică oferă rezultate similare în cazul datelor utilizate. Acest lucru se datorează caracterului constant sau variațiilor reduse ale datelor originale. În astfel de situații, curbele generate de aceste metode nu aduc diferențe semnificative, ceea ce face ca utilizarea ambelor să fie redundantă.

LOESS, în schimb, se remarcă prin flexibilitatea sa în adaptarea la variațiile locale și prin capacitatea de a genera o curbă mai netedă și mai realistă. Această metodă este mai potrivită pentru datele analizate, deoarece evidențiază diferențele subtile și oferă o interpolare mai naturală.

Astfel, pentru acest set de date cu variații reduse și valori constante în anumite intervale, **LOESS** reprezintă cea mai relevantă metodă de interpolare, fiind capabilă să captureze subtilitățile datelor, spre deosebire de metodele mai simple.

4 Modelarea sistemului

4.1 Introducere

Scopul acestei secțiuni este de a dezvolta modele predictive care să reliefeze relația dintre variabilele climatice și PMV (Predicted Mean Vote), o măsură esențială în evaluarea confortului termic. În cadrul proiectului, au fost utilizate diverse tehnici de regresie și modele bazate pe arbori de decizie, având ca obiectiv principal identificarea celui mai performant model din punct de vedere al preciziei și capacității de *generalizare*.

4.1.1 Importanța modelării corecte în contextul aplicației HVAC

Modelarea precisă a confortului termic este crucială pentru aplicațiile HVAC (Heating, Ventilation, and Air Conditioning), deoarece acestea trebuie să asigure condiții optime de confort termic în timp ce optimizează consumul de energie.

Un model precis permite sistemelor HVAC să anticipeze necesitățile utilizatorilor și să ia decizii informate pentru controlul variabilelor climatice, cum ar fi temperatura, umiditatea sau fluxul de aer. Acest lucru nu doar că îmbunătățește satisfacția utilizatorilor, dar contribuie și la sustenabilitatea energetică a clădirilor.

4.1.2 Utilitatea datelor preprocesate pentru îmbunătățirea performanței modelului

Datele preprocesate reprezintă o bază solidă pentru obținerea unor modele performante. Eliminarea outlierilor, interpolarea valorilor lipsă și normalizarea caracteristicilor au contribuit la crearea unui set de date consistent și uniform, reducând riscul de erori și asigurând condiții ideale pentru antrenarea modelelor.

În acest fel, modelele predictive dezvoltate au fost capabile să capteze mai bine relațiile dintre variabilele climatice și confortul termic, maximizând acuratețea predicțiilor.

Această secțiune prezintă metodologia, tehnicile aplicate, și rezultatele obținute pentru fiecare model, oferind o bază solidă pentru alegerea soluției optime.

4.2 Metodologie și tehnici aplicate

4.2.1 Regresie liniară

Regresia liniară este o tehnică de modelare statistică utilizată pentru a estima relația dintre o variabilă y și una sau mai multe variabile independente x .

Modelul presupune că relația este liniară și poate fi reprezentată prin ecuația:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n \quad (4.4.2)$$

unde:

- y : variabila dependentă (PMV în acest caz),
- β_0 : interceptul (valoarea y când toate variabilele independente sunt 0),
- $\beta_1, \beta_2 \dots \beta_n$: coeficienții care indică ponderea fiecărei variabile independente,
- $x_1, x_2 \dots x_n$: variabilele independente.

Scopul regresiei este să găsească valorile coeficienților care minimizează eroarea pătratică medie (MSE) între valorile observate și cele prezise.

4.2.1.1 Rezultate:

- **MSE (Mean Squared Error):** 0.0554
- **R² Score:** 0.9391

4.2.1.2 Interpretarea rezultatelor

Un **R² Score** de 0.939 sugerează că modelul explică aproximativ 93.9% din variația datelor, ceea ce indică o performanță bună.

Valoarea relativ mică a **MSE (0.0554)** arată că modelul are erori de predicție reduse.

Deși modelul oferă o predicție precisă, este important să fie comparat cu alte metode pentru a verifica dacă există soluții mai performante.

4.2.2 Regresie logistică

Regresia logistică este o metodă utilizată pentru clasificare binară, unde variabila dependentă y poate avea două clase (ex. confortabil/disconfortabil). Modelul prezice probabilitatea ca un punct de date să aparțină unei clase folosind funcția logistică:

$$P(y = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}} \quad (4.4.2)$$

Unde: $P(y = 1)$ reprezintă probabilitatea ca observația să aparțină clasei 1. Predicțiile sunt transformate în clase (0 sau 1) utilizând un prag (de obicei 0.5).

4.2.2.1 Rezultate:

- **Acuratețe:** 99.01%
- **Matricea de clasificare:**
- Clasa 0: Precision = 1.00, Recall = 0.98, F1-Score = 0.99
- Clasa 1: Precision = 0.98, Recall = 1.00, F1-Score = 0.99
- **F1-Score mediu ponderat:** 0.99

4.2.2.2 Interpretarea rezultatelor

Acuratețea de 99.01% arată că modelul clasifică foarte bine datele.

Precision și Recall sunt aproape 1.00, indicând un echilibru bun între clase și o clasificare precisă a confortului termic.

Regresia logistică este un model eficient pentru clasificarea binară a PMV în clase de confort termic. Performanța excelentă sugerează că este o soluție solidă pentru probleme de clasificare.

4.2.3 Regresii Ridge și Lasso

4.2.3.1 Ridge Regression

Utilizează regularizarea L2, adăugând o penalizare la suma pătratelor coeficienților în funcția de cost:

$$Cost = \sum (y_i - \hat{y}_i)^2 + \lambda \sum (\beta_j)^2 \quad (4.2.3.1)$$

Aceasta ajută la prevenirea suprapotrivirii prin reducerea magnitudinii coeficienților, păstrând toate caracteristicile.

4.2.3.2 Lasso Regression

Utilizează regularizarea L1, adăugând o penalizare la suma valorilor absolute ale coeficienților:

$$Cost = \sum (y_i - \hat{y}_i)^2 + \lambda \sum |\beta_j| \quad (4.2.3.2)$$

Această metodă poate elimina coeficienți, fiind utilă pentru selecția caracteristicilor.

4.2.3.3 Rezultate

Ridge Regression:

- MSE: 0.0554
- R² Score: 0.9391

Lasso Regression:

- MSE: 0.0563
- R² Score: 0.9381

4.2.3.4 Interpretarea rezultatelor

Ridge Regression: Performanțe similare cu regresia liniară, dar mai stabile datorită regularizării L2.

Lasso Regression: Ușor mai slabă decât Ridge (MSE puțin mai mare, R^2 mai mic), dar a fost utilă pentru identificarea caracteristicilor mai puțin relevante.

Ridge Regression este preferabilă pentru modele generale stabile, în timp ce Lasso este mai potrivită când selecția automată a caracteristicilor este necesară.

4.2.4 Random Forest

Random Forest este o tehnică de învățare automată bazată pe arbori de decizie, utilizând un ansamblu de arbori pentru a îmbunătăți precizia modelului și a reduce overfitting-ul.

Metoda principală din spatele Random Forest este Bagging (Bootstrap Aggregating), care presupune: generarea de subeșantioane aleatorii din datele de antrenament, construirea arborilor de decizie individuali pe aceste subeșantioane, agregarea predicțiilor arborilor prin: media predicțiilor pentru regresie, votul majoritar pentru clasificare.

4.2.4.2 Rezultate

- MSE (Mean Squared Error): 0.0041
- R^2 Score: 0.9955

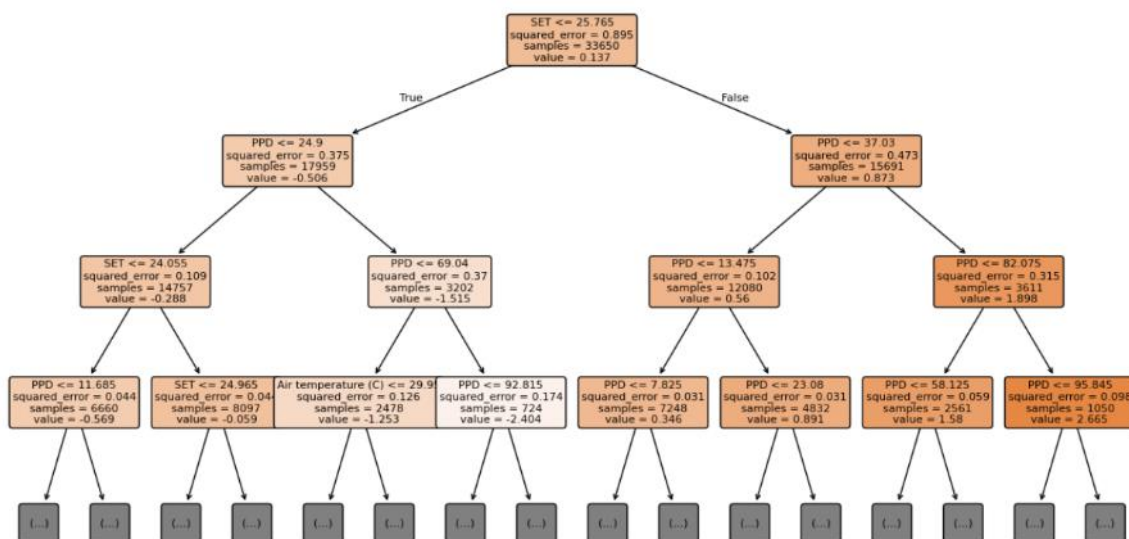


Figura 4.2.4 a) Vizualizare arobore nr.0 (max_depth=3)

Analiza importanței caracteristicilor:

- Air temperature (C): 0.008638
- Relative humidity (%): 0.000691
- Air velocity (m/s): 0.001731
- PPD: 0.365071
- SET: 0.623870

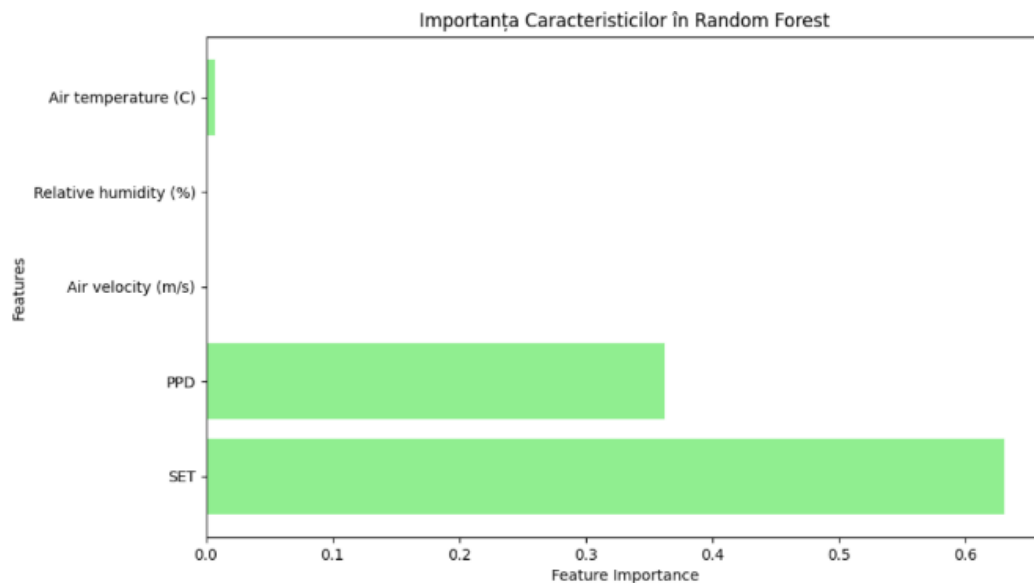


Figura 4.2.4 b) Importanța Caracteristicilor în Random Forest

4.2.4.2 Interpretarea rezultatelor

Performanță ridicată: Cu un **MSE de 0.0041** și un **R² Score de 0.9955**, Random Forest este unul dintre cele mai precise modele utilizate.

Feature Importance: Caracteristica SET este dominantă, contribuind cu peste 62% la predicțiile modelului, urmată de PPD cu 36.51%. Restul caracteristicilor (temperatura aerului, umiditatea relativă și viteza aerului) au importanțe semnificativ mai mici, ceea ce sugerează că modelul se bazează în principal pe factori integrați pentru a prezice PMV.

4.2.5 XGBoost

XGBoost (Extreme Gradient Boosting) este o variantă îmbunătățită a algoritmului Gradient Boosting, optimizată pentru performanță și eficiență. **Gradient Boosting** construiește arbori succesivi, fiecare încercând să corecteze erorile arborelui anterior, minimizând o funcție de cost prin gradient descent.

Avantaje ale XGBoost:

- Include regularizare L1 și L2 pentru prevenirea overfitting-ului.
- Construirea arborilor este optimizată pentru a utiliza eficient resursele hardware.
- Utilizează tehnici precum subsampling și shrinkage (înmulțirea predicțiilor arborilor cu un factor niu) pentru a îmbunătăți generalizarea.

4.2.5.1 Rezultate

- **MSE (Mean Squared Error):** 0.0039
- **R² Score:** 0.995

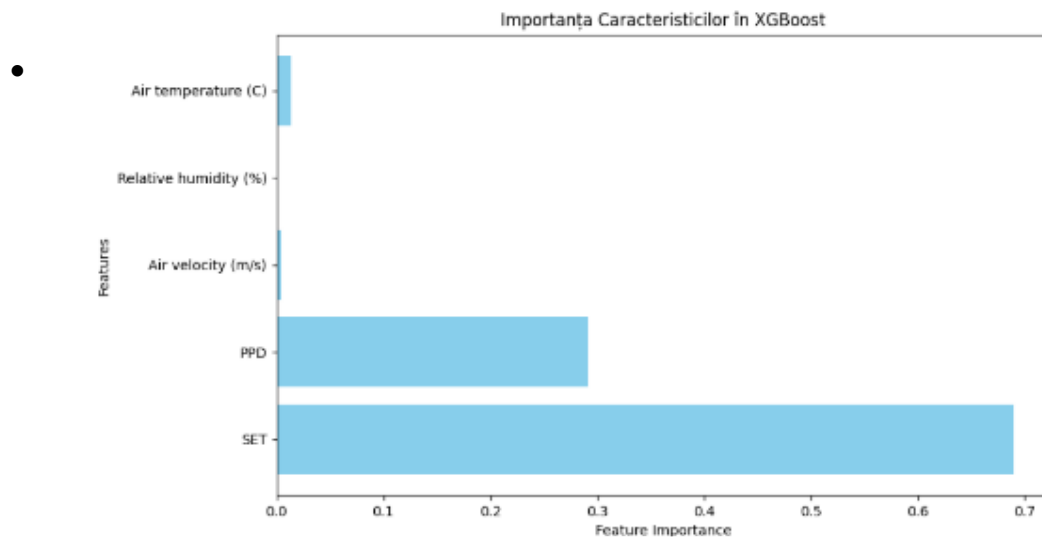


Figura 4.2.5 Importanța Caracteristicilor în XGBoost

Analiza importanței caracteristicilor:

- Air temperature (C): 0.013489
- Relative humidity (%): 0.000389
- Air velocity (m/s): 0.004639
- PPD: 0.291973
- SET: 0.689510

4.2.5.2 Interpretarea rezultatelor

Performanță excelentă: Cu un **MSE de 0.0039** și un **R² Score de 0.9956**, XGBoost depășește ușor Random Forest în precizie.

Feature Importance: Similar cu Random Forest, SET este cea mai importantă caracteristică, contribuind cu aproape 69% la predicții, urmată de PPD (29.20%). Celelalte caracteristici au importanțe neglijabile, confirmând relevanța dominantă a factorilor integrați.

4.3 Implementarea

4.3.1 Mediul de implementare

Pentru acest proiect, implementarea s-a realizat utilizând Python , iar codul complet al proiectului este salvat într-un fișier separat, care conține toate etapele necesare: de la preprocesarea datelor, la construirea și evaluarea modelelor. Acest fișier poate fi accesat și rulat pe orice sistem compatibil cu Python, asigurând reproducerea completă a lucrării.

4.3.2 Cum se poate replica lucrarea

Pentru a replica rezultatele acestui proiect, urmați pașii următori:

Pregătirea mediului:

Asigurați-vă că aveți instalate toate bibliotecile necesare:

- *scikit-learn*: pentru regresie și evaluarea modelelor.
- *xgboost*: pentru implementarea Gradient Boosting.
- *pandas*: pentru manipularea și analiza datelor.
- *numpy*: pentru operațiuni numerice.
- *matplotlib*: pentru vizualizarea rezultatelor.

Fișierul principal al proiectului:

Toate etapele proiectului, inclusiv codul pentru modele, sunt salvate într-un fișier Python

Acest fișier include: preprocesarea datelor, divizarea datelor în seturi de antrenament și testare, construirea și antrenarea modelelor, evaluarea performanței fiecărui model, vizualizarea rezultatelor și analiza importanței caracteristicilor.

4.4 Evaluarea modelelor

4.4.1 Metode și metrici de evaluare

Evaluarea performanței modelelor utilizate în acest proiect s-a bazat pe metrici specifice pentru regresie și clasificare, fiecare fiind alese pentru a măsura acuratețea și eficiența modelelor în diferite aspecte.

Pentru modelele de **regresie**, s-au utilizat următoarele metrici:

MSE (Mean Squared Error): Această metrică măsoară media pătratelor diferențelor dintre valorile reale și cele prezise de model. Un MSE mic indică o eroare de predicție redusă și, implicit, un model mai precis.

R² Score (Coeficientul de determinare): Această metrică arată proporția din variația totală a datelor care este explicată de model. Valori apropiate de 1 indică un model care explică bine datele, în timp ce valori apropiate de 0 sugerează că modelul nu capturează suficient de bine relațiile dintre variabile.

Pentru **clasificare binară**, s-au utilizat metrici precum:

Acuratețea (Accuracy): Proporția predicțiilor corecte în totalul datelor testate. Este o metrică globală, dar poate fi insuficientă în cazul unui dezechilibru între clase.

Precizia (Precision): Măsoară procentul de predicții corecte dintre toate predicțiile pozitive făcute de model.

Recall-ul: Măsoară capacitatea modelului de a identifica toate cazurile pozitive din setul de date.

F1-score: Combinația dintre precizie și recall, fiind o medie armonică a celor două. Este utilă în special când există un dezechilibru între clase, oferind o imagine mai completă a performanței modelului.

În acest proiect, **MSE** și **R² Score** au fost folosite pentru a evalua modelele de regresie, în timp ce **F1-score**, acuratețea, precizia și recall-ul au fost utilizate pentru evaluarea clasificării binare.

4.4.2 Analiza rezultatelor

Performanța fiecărui model utilizat a fost comparată și analizată în detaliu, utilizând metricele menționate mai sus. Tabelul de mai jos sintetizează rezultatele obținute pentru fiecare model:

Tabel 4.4.2 Performanța modelelor utilizate pentru analiza PMV

Model	MSE	R ² Score / Acuratețe
Linear Regression	0.0554	0.9391
Logistic Regression	N/A	0.9901 (Acuratețe)
Ridge Regression	0.0554	0.9391
Lasso Regression	0.0563	0.9381
Random Forest Regressor	0.0041	0.9955
XGBoost Regressor	0.0039	0.9956

4.4.3 Interpretarea rezultatelor

Modelele **XGBoost** și **Random Forest** au obținut cele mai bune rezultate pentru regresie, cu valori extrem de mici ale MSE (0.0039 și 0.0041, respectiv) și R² Score foarte apropiate de 1 (0.9956 și 0.9955). Aceste performanțe superioare se datorează capacității acestor algoritmi de a captura relațiile non-liniare dintre variabile și de a gestiona interacțiuni complexe între caracteristici.

Modelele liniare, cum ar fi **Linear Regression**, **Ridge** și **Lasso**, au obținut un R² Score de aproximativ 0.939, ceea ce arată că acestea pot explica ~93.9% din variația datelor. Totuși, aceste rezultate sunt mai slabe comparativ cu modelele bazate pe arbori de decizie, ceea ce sugerează că modelele liniare nu sunt la fel de eficiente în capturarea relațiilor complexe.

Pentru clasificare, **Logistic Regression** a demonstrat o performanță excelentă, cu o acuratețe de 99.01% și un F1-score aproape perfect (0.99). Acest lucru arată că modelul poate clasifica eficient confortul termic (PMV_binary) în funcție de variabilele de intrare.

4.4.4 Concluzii privind performanța modelelor

Modele bazate pe arbori de decizie: XGBoost și Random Forest s-au dovedit a fi cele mai performante modele, având capacitatea de a gestiona relații complexe și de a minimiza erorile de predicție. De asemenea, analiza importanței caracteristicilor a confirmat că variabilele compuse (SET și PPD) au fost dominante, contribuind semnificativ la predicții.

Modele liniare: Deși performanțele lor au fost bune, acestea nu au reușit să concureze cu modelele avansate din punct de vedere al preciziei. Totuși, aceste modele sunt mai simple și mai rapide, fiind utile în scenarii unde interpretabilitatea este esențială.

Clasificare logistică: Acest model a demonstrat o capacitate remarcabilă de a clasifica datele, fiind un instrument eficient pentru problemele de clasificare binară.

Această analiză arată că alegerea unui model depinde de natura problemei și de complexitatea datelor, iar modelele bazate pe arbori oferă cel mai bun echilibru între precizie și capacitatea de generalizare.

4.5 Concluzii

În concluzie, analiza performanței fiecărui model a evidențiat diferențe semnificative între metodele aplicate. Modelele liniare, precum regresia liniară, Ridge și Lasso, au oferit rezultate decente, cu un R^2 de aproximativ 93.9%, dar s-au dovedit mai puțin eficiente în capturarea relațiilor complexe dintre variabile. Pe de altă parte, Random Forest și XGBoost au excelat, cu scoruri R^2 de 0.9955 și 0.9956 și MSE-uri foarte mici, indicând o capacitate ridicată de predicție.

Modelul **XGBoost** a fost selectat ca model final datorită performanței superioare, reflectată de cel mai mic MSE și cel mai mare R^2 Score. Analiza importanței caracteristicilor a confirmat relevanța măsurilor integrate precum **SET** (62.39%) și **PPD** (36.51%), subliniind contribuția majoră a acestor variabile la predicția confortului termic.

Pentru îmbunătățirea viitoare a modelului, ar putea fi utile tehnici suplimentare de optimizare a hiperparametrilor, cum ar fi utilizarea căutării Bayesiene sau a rețelelor de căutare, pentru a ajusta performanța XGBoost și Random Forest. De asemenea, explorarea altor modele avansate, precum rețelele neuronale, ar putea aduce beneficii suplimentare în cazul unor seturi de date mai complexe.

5 Concluzii

5.1 Rezultate obținute

Proiectul a avut ca scop modelarea și analiza confortului termic utilizând date climatice și diverse tehnici de modelare, pentru a identifica cea mai potrivită metodologie.

Am parcurs etape clare, începând cu preprocesarea datelor, continuând cu aplicarea diferitelor modele și evaluarea rezultatelor, și încheind cu interpretarea performanțelor fiecărui model.

În etapa de preprocesare, am eliminat datele necorespunzătoare prin tehnici precum **Isolation Forest**, iar valorile lipsă au fost completate folosind metode diverse de interpolare, inclusiv LOESS, care s-a dovedit cea mai adaptată naturii datelor noastre. Acest proces a îmbunătățit calitatea setului de date, reducând impactul erorilor asupra performanței modelelor.

Pentru etapa de modelare, am aplicat mai multe metode, inclusiv **regresia liniară**, **Ridge și Lasso**, precum și modele avansate bazate pe arbori de decizie, cum ar fi **Random Forest și XGBoost**.

Rezultatele au demonstrat că XGBoost a avut cea mai bună performanță, cu un **MSE de 0.0039** și un **R² Score de 0.9956**, depășind semnificativ modelele liniare în capacitatea sa de a captura relațiile non-liniare din date.

Contribuțiile proiectului includ:

- Implementarea unui flux complet de preprocesare, modelare și evaluare, adaptat datelor climatice.
- Selectarea și validarea modelelor avansate pentru predicția PMV, subliniind relevanța caracteristicilor integrate precum SET și PPD.
- Oferirea unei metodologii detaliate și reproductibile, care poate fi aplicată în alte studii sau proiecte HVAC.

Comparativ cu alte lucrări din domeniu, proiectul nostru se remarcă prin integrarea mai multor modele, evaluarea detaliată a performanței și o atenție sporită acordată preprocesării datelor, aspect care a contribuit semnificativ la acuratețea rezultatelor.

5.2 Direcții de dezvoltare

Pentru a continua dezvoltarea proiectului, există mai multe direcții posibile:

Optimizarea hiperparametrilor: Tehnici precum căutarea Bayesienă sau optimizarea cu rețele neuronale pot îmbunătăți performanțele modelelor XGBoost și Random Forest.

Extinderea setului de date: Adăugarea mai multor caracteristici, precum date geografice sau istorice, ar putea îmbunătăți predicțiile și ar permite o analiză mai detaliată.

Explorarea modelelor avansate: Implementarea unor arhitecturi de rețele neuronale, cum ar fi rețelele neuronale convoluționale (CNN) sau rețelele recurente (RNN), ar putea deschide noi oportunități pentru analiza datelor climatice.

Implementarea în timp real: Integrarea unui sistem în timp real pe o platformă hardware, precum Raspberry Pi, pentru a monitoriza și controla confortul termic într-un mediu HVAC.

Validare pe seturi externe de date: Testarea modelelor pe alte seturi de date climatice, din diferite locații geografice, pentru a verifica robustețea soluțiilor propuse.

Proiectul oferă o bază solidă pentru analiza și predicția confortului termic, dar extinderea cercetării și aplicarea metodelor în medii practice ar putea aduce beneficii semnificative în domeniul sistemelor HVAC inteligente.

6 Bibliografie

- 1] **V. Moroșanu**, "Cercetări privind realizarea de sisteme integrate pentru clădiri cu consum redus de energie," Teză de doctorat, Universitatea Tehnică "Gheorghe Asachi" din Iași, 2024. [Online].
Disponibil:<https://doctorat.tuiasi.ro/wp-content/uploads/2024/09/Rezumat-Morosanu-Vlad.pdf>
- 2] **A. M. Măgurean**, "Analiza performanței energetice a clădirilor nerezidențiale prin tehnici de modelare numerică și Inteligență Artificială aplicată," *Teză de doctorat*, Universitatea Tehnică din Cluj-Napoca, 2017. [Online].
Disponibil:https://www.researchgate.net/publication/387504020_Analiza_performantei_energetice_a_cladirilor_nerezidentiale_prin_tehnici_de_modelare_numerica_si_Inteligenta_Artificiala_aplicata_Teza_de_doctorat
- 3] **IEEE Citation Reference**, 2009. [Online].
Disponibil:https://journals.ieeeauthorcenter.ieee.org/wp-content/uploads/sites/7/IEEE_Reference_Guide.pdf
- [4] **IEEE Editorial Style Manual**, 2016. [Online].
Disponibil:https://moodle.hs-augsburg.de/pluginfile.php/283344/mod_resource/content/2/IEEE-Editorial-Style_Manual.pdf

6.1 Specificații pentru bibliografie

Bibliografia proiectului a fost selectată pentru a include surse relevante din domeniul modelării sistemelor HVAC, al eficienței energetice și al metodologiilor utilizate în analiza și optimizarea performanței acestora. Referințele alese acoperă atât cercetări academice de actualitate, cât și standarde utilizate pentru redactarea documentației tehnice și științifice.

[1] V. Moroșanu, "Cercetări privind realizarea de sisteme integrate pentru clădiri cu consum redus de energie," Teză de doctorat, Universitatea Tehnică "Gheorghe Asachi" din Iași, 2024. [Online]. Disponibil:

Acest document oferă o perspectivă detaliată asupra strategiilor de integrare a sistemelor HVAC în clădiri cu consum redus de energie, subliniind metodele moderne de optimizare energetică. Informațiile incluse sunt utile pentru înțelegerea modului în care tehnologiile avansate pot contribui la eficiența termică a clădirilor și reducerea consumului de resurse.

[2] A. M. Măgurean, "Analiza performanței energetice a clădirilor nerezidențiale prin tehnici de modelare numerică și Inteligență Artificială aplicată," Teză de doctorat, Universitatea Tehnică din Cluj-Napoca, 2017. [Online]. Disponibil:

Această lucrare explorează aplicarea tehnicilor de modelare numerică și a algoritmilor de inteligență artificială în analiza performanței energetice a clădirilor. Studiul este relevant pentru proiectul actual, deoarece abordează utilizarea modelelor predictive pentru optimizarea consumului energetic în sistemele HVAC.

[3] IEEE Citation Reference, 2009. [Online]. Disponibil:

Acest ghid oferă regulile oficiale de citare conform standardelor IEEE. A fost utilizat pentru a asigura o structură corectă și uniformă a referințelor bibliografice din proiect.

[4] IEEE Editorial Style Manual, 2016. [Online]. Disponibil:

Manualul de stil IEEE a fost utilizat pentru redactarea documentului, asigurând respectarea standardelor de formatare, numerotare și organizare a conținutului tehnic.