



UNIVERZITET U NOVOM SADU
FAKULTET TEHNIČKIH NAUKA U
NOVOM SADU



Igor Manojlović

**KRATKOROČNA PROBABILISTIČKA PROGNOZA
OPTEREĆENJA NA NISKOM NAPONU U
ELEKTRODISTRIBUTIVNIM MREŽAMA**

DOKTORSKA DISERTACIJA

Mentori

Prof. dr Goran Švenda

Prof. dr Aleksandar Erdeljan

Novi Sad, 2022



УНИВЕРЗИТЕТ У НОВОМ САДУ • ФАКУЛТЕТ ТЕХНИЧКИХ НАУКА
21000 НОВИ САД, Трг Доситеја Обрадовића 6

КЉУЧНА ДОКУМЕНТАЦИЈСКА ИНФОРМАЦИЈА

Редни број, РБР:	
Идентификациони број, ИБР:	
Тип документације, ТД:	Монографска документација
Тип записа, ТЗ:	Текстуални штампани материјал
Врста рада, ВР:	Докторска дисертација
Аутор, АУ:	Игор Манојловић
Ментор, МН:	др Горан Швенда, редовни професор др Александар Ердељан, редовни професор
Наслов рада, НР:	Краткорочна пробабилистичка прогноза оптерећења на ниском напону у електродистрибутивним мрежама
Језик публикације, ЈП:	српски
Језик извода, ЈИ:	српски
Земља публикавања, ЗП:	Република Србија
Уже географско подручје, УГП:	Аутономна покрајина Војводина
Година, ГО:	2022
Издавач, ИЗ:	Факултет техничких наука
Место и адреса, МА:	Нови Сад, Трг Доситеја Обрадовића 6
Физички опис рада, ФО: <small>(поглавља/страница/цитата/табела/слика/графика/прилога)</small>	5/93/199/6/55/0/0
Научна област, НО:	Електротехничко и рачунарско инжењерство
Научна дисциплина, НД:	Енергетика, електроника и телекомуникације
Предметна одредница/Кључне речи, ПО:	пробабилистичка прогноза, временске серије, профили оптерећења, машинско учење, дубоко учење, екстракција атрибута, одабир атрибута, кластеризација
УДК	
Чува се у, ЧУ:	Библиотека Факултета техничких наука
Важна напомена, ВН:	
Извод, ИЗ:	Предмет истраживања ове докторске дисертације је краткорочна пробабилистичка прогноза оптерећења на ниском напону у електродистрибутивним мрежама. Циљ истраживања је да се развије ново решење које ће уважити варијабилност оптерећења на ниском напону и понудити конкурентну тачност прогнозе уз високу ефикасност са становишта заузећа рачунарских ресурса. Предложено решење се заснива на примени статистичких метода и метода машинског (дубоког) учења у репрезентацији података (екстракцији и одабиру атрибута), кластеровању и регресији. Ефикасност предложеног решења је верификована у студији случаја над скупом реалних података са паметних бројила. Резултат примене предложеног решења је висока тачност прогнозе и кратко време извршавања у поређењу са конкурентним решењима из актуелног стања у области.
Датум прихватања теме, ДП:	
Датум одбране, ДО:	
Чланови комисије, КО:	Председник: др Душко Бекут, редовни професор
	Члан: др Драган Тасић, редовни професор
	Члан: др Лука Стрезоски, доцент
	Члан: др Милан Гаврић, доцент
	Члан, ментор: др Горан Швенда, редовни професор
	Члан, ментор: др Александар Ердељан, редовни професор
	Потпис ментора



KEY WORDS DOCUMENTATION

Accession number, ANO:	
Identification number, INO:	
Document type, DT:	Monograph documentation
Type of record, TR:	Textual printed material
Contents code, CC:	Ph.D. thesis
Author, AU:	Igor Manojlović
Mentor, MN:	Goran Švenda, Ph.D., Full Professor Aleksandar Erdeljan, Ph.D., Full Professor
Title, TI:	Probabilistic short-term load forecasting at low voltage in distribution networks
Language of text, LT:	Serbian
Language of abstract, LA:	Serbian
Country of publication, CP:	Republic of Serbia
Locality of publication, LP:	Autonomous Province of Vojvodina
Publication year, PY:	2022
Publisher, PB:	Faculty of Technical Sciences
Publication place, PP:	Trg Dositeja Obradovića 6, 21000 Novi Sad
Physical description, PD: <small>(chapters/pages/ref./tables/pictures/graphs/appendixes)</small>	5/93/199/6/55/0/0
Scientific field, SF:	Electrical and Computer Engineering
Scientific discipline, SD:	Power, Electronic and Telecommunication Engineering
Subject/Key words, S/KW:	probabilistic forecasting, time series, load profiles, machine learning, deep learning, feature extraction, feature selection, clustering
UC	
Holding data, HD:	The Library of the Faculty of Technical Sciences
Note, N:	
Abstract, AB:	<p>This Ph.D. thesis deals with the problem of probabilistic short-term load forecasting at the low voltage level in power distribution networks. The research goal is to develop a new solution that considers load variability and offers high forecasting accuracy without excessive hardware requirements. The proposed solution is based on the application of statistical methods and machine (deep) learning methods for data representation (feature extraction and selection), clustering, and regression. The efficiency of the proposed solution was verified in a case study on real smart meter data. The case study results confirm that the application of the proposed solution leads to high forecast accuracy and short execution time compared to related solutions.</p>
Accepted by the Scientific Board on, ASB:	
Defended on, DE:	
Defended Board, DB:	President: Duško Bekut, Ph.D., Full Professor
	Member: Dragan Tasić, Ph.D., Full Professor
	Member: Luka Strezoski, Ph.D., Assistant Professor
	Member: Milan Gavrić, Ph.D., Assistant Professor
	Member, Mentor: Goran Švenda, Ph.D., Full Professor
	Member, Mentor: Aleksandar Erdeljan, Ph.D., Full Professor
	Mentor's sign

Zahvaljujem se

*Obnovljivom izvoru pozitivne energije, dobrote i razumevanja
tokom svih proteklih godina studiranja, mojoj dragoj Dajani.*

*Porodici i prijateljima, na dugogodišnjoj podršci
u naporima da istrajem na ovom putu.*

*Profesorima, na nesebičnom deljenju znanja koje mi je pomoglo da,
korak po korak, stignem do ove doktorske disertacije.*

„Nothing is too wonderful to be true, if it be consistent with the laws of nature.“

Michael Faraday

REZIME

Predmet istraživanja ove doktorske disertacije je kratkoročna probabilistička prognoza opterećenja na niskom naponu u elektrodistributivnim mrežama. Teoretski i praktično, takva prognoza je izazovnije od prognoze opterećenja na srednjem i visokom naponu zbog prirodno većih varijacija u opterećenju i veće količine podataka, koji su posledica većeg broja posmatranih čvorova mreže. Međutim, takav vid prognoze je takođe veoma važan elektrodistributivnim preduzećima koja pokušavaju da se prilagode rastućem trendu dekarbonizacije i decentralizacije u svakodnevnom radu. Kvantifikacija varijabilnosti opterećenja u vidu probabilističke prognoze je vitalni deo upravljanja rizicima i presudna je za smanjenje operativnih troškova, optimalno odlučivanje i upravljanje elektrodistributivnom mrežom. U skladu s tim, identifikovane su dve potrebe za istraživanjem: 1) potreba za proračunom kratkoročne probabilističke prognoze koji kvantifikuje varijabilnost opterećenja na niskom naponu, i 2) potreba za povećanjem tačnosti prognoze opterećenja na niskom naponu bez neopravdanog povećanja računarskih resursa.

Globalni cilj istraživanja koje je obuhvaćeno ovom doktorskom disertacijom je da se razvije novo rešenje za kratkoročnu probabilističku prognozu opterećenja na niskom naponu koje će uvažiti varijabilnost opterećenja i ponuditi konkurentnu tačnost prognoze uz visoku efikasnost sa stanovišta zauzeća računarskih resursa. Radi postizanja navedenog globalnog cilja, identifikovani su sledeći individualni ciljevi:

- 1) Razviti novo rešenje tako da redukuje model podataka o opterećenju na niskom naponu bez gubitka značajnih informacija o varijabilnosti opterećenja i na taj način smanji zauzeće računarskih resursa u prognozi.
- 2) Razviti novo rešenje tako da omogući primenu sofisticiranih regresionih metoda nad redukovanim modelom podataka i na taj način ponudi konkurentnu tačnost prognoze.
- 3) Verifikovati opravdanost primene predloženog rešenja nad skupom realnih podataka.
- 4) Uporediti predloženo rešenje sa konkurentnim rešenjima iz aktuelnog stanja u oblasti.

U skladu sa postavljenim ciljevima u doktorskoj disertaciji je razvijeno novo rešenje za kratkoročnu probabilističku prognozu opterećenja na niskom naponu – *Deep Centroid Learning*. Predloženo rešenje se zasniva na primeni statističkih metoda i metoda mašinskog (dubokog) učenja u reprezentaciji podataka (ekstrakciji i odabiru atributa), klasterizaciji i regresiji. Prognoza se vrši na osnovu vremenskih serija centralnih momenata opterećenja (proseka i standardne devijacije) na nivou grupe potrošača sa sličnim šablonima opterećenja. Ulazni podaci u prognozi su prethodno opterećenje, vremenske odrednice i meteorološki faktori. Rezultat primene predloženog rešenja je najpre probabilistička prognoza na nivou potrošačke grupe, a zatim i na nivou članova grupe (pojedinačnih potrošača).

Predloženo rešenje ukazuje na to da se varijabilnost opterećenja na nivou niskonaponskih potrošača može kvantifikovati probabilističkom prognozom na nivou potrošačkih grupa. Prema tome, model opterećenja se može redukovati tako da omogući primenu dubokog učenja za prognozu opterećenja na niskom naponu bez prekomernog zauzeća računarskih resursa. Zauzvrat, duboko učenje omogućava otkrivanje i najsloženijih nelinearnih veza između visoko varijabilnog opterećenja i promenljivih faktora koji na to opterećenje utiču. Na taj način, dodatno je povećana tačnost rezultata prognoze.

Efikasnost predloženog rešenja je verifikovana u studiji slučaja nad skupom realnih podataka koji su sakupljeni sa pametnih brojlara, u jednoj severnoameričkoj i jednoj australijskoj elektrodistributivnoj mreži. Rezultat primene predloženog rešenja je visoka tačnost prognoze i kratko vreme izvršavanja u poređenju sa konkurentnim rešenjima iz aktuelnog stanja u oblasti. Glavni doprinos ove doktorske disertacije je novo rešenje za kratkoročnu probabilističku prognozu opterećenja na niskom naponu koje je pogodno za primenu u savremenim sistemima za upravljanje distribucijom električne energije.

ABSTRACT

This Ph.D. thesis deals with the problem of probabilistic short-term load forecasting at the low voltage level in power distribution networks. It is theoretically and practically more challenging to obtain the forecast at the low voltage level than to obtain the forecast at medium and high voltage levels due to naturally high variations in load data and a larger amount of data, originating from a larger number of observed network nodes. However, this type of forecast is also very important for power distribution utilities that are trying to adapt to the growing need for decarbonization and decentralization. The quantification of load variability in the form of the probabilistic forecast is a vital part of risk management and crucial for reducing operating costs, optimal decision-making, and distribution management. Accordingly, two research needs have been identified: 1) the need to obtain a probabilistic short-term forecast that quantifies the load variability at the low voltage level, and 2) the need to increase the accuracy of load forecast at low voltage level without excessive hardware requirements.

The main research goal is to develop a new solution for probabilistic short-term load forecasting at the low voltage level, that considers load variability and offers high forecasting accuracy without excessive hardware requirements. To achieve the main research goal, the following subsidiary goals are introduced:

- 1) Develop a new solution to reduce the need for excessive hardware requirements by reducing the load data model for low voltage consumers without losing significant information on load variability.
- 2) Develop a new solution to achieve high forecast accuracy by enabling the application of sophisticated regression methods on the reduced data model.
- 3) Verify the proposed solution on real smart meter data.
- 4) Compare the proposed solution with existing solutions.

Accordingly, this Ph.D. thesis presents a new solution for probabilistic short-term load forecasting at the low voltage level – Deep Centroid Learning. The proposed solution is based on statistical methods and machine (deep) learning methods for data representation (feature extraction and selection), clustering, and regression. The forecast is based on time series of central moments (mean and standard deviation) of load values at the level of load groups. The forecast inputs are lagged historical loads, and weather and calendar features. The result of the application of the proposed solution is a probabilistic forecast both at the level of load groups and at the level of group members (individual consumers).

The proposed solution indicates that the load variability at the low voltage level can be quantified by a probabilistic forecast at the level of load groups. Therefore, the load data model can be reduced to enable the application of deep learning in low voltage load forecasting without excessive hardware requirements. In turn, deep learning enables the discovery of complex nonlinear relationships between the highly variable load and the variable factors affecting the load. This increases the forecast accuracy.

The efficiency of the proposed solution was verified in a case study on the real US and Australian smart meter data. The case study results confirm that the application of the proposed solution leads to high forecast accuracy and short execution time compared to related solutions. The main contribution of this Ph.D. thesis is a new solution for probabilistic short-term load forecasting at the low voltage level, convenient for application in modern distribution management systems.

SADRŽAJ

<i>Spisak skraćenica</i>	9
<i>Spisak simbola</i>	11
<i>Spisak slika</i>	14
<i>Spisak tabela</i>	16
1. <i>Uvod</i>	17
1.1. <i>Predmet istraživanja</i>	17
1.2. <i>Pregled stanja u oblasti</i>	20
1.3. <i>Potrebe za istraživanjem i ciljevi istraživanja</i>	21
1.4. <i>Pregled doktorske disertacije</i>	22
2. <i>Teoretske osnove</i>	23
2.1. <i>Reprezentacija</i>	23
2.1.1. <i>Ekstrakcija atributa</i>	24
2.1.2. <i>Odabir atributa</i>	25
2.2. <i>Klasterizacija</i>	27
2.3. <i>Regresija</i>	29
2.3.1. <i>Osnovne metode</i>	29
2.3.2. <i>Duboko učenje</i>	31
2.4. <i>Optimizacija</i>	35
2.5. <i>Verifikacija</i>	36
2.6. <i>Probabilistička prognoza</i>	39
3. <i>Predlog rešenja</i>	41
3.1. <i>HMTSR</i>	42
3.2. <i>HMLPSA</i>	43
3.3. <i>TSGA</i>	44
3.4. <i>LPR</i>	45
3.5. <i>DCLN</i>	48
3.5.1. <i>Ulazni i izlazni podaci</i>	48
3.5.2. <i>Struktura modela</i>	49
4. <i>Studija slučaja</i>	52
4.1. <i>Dizajn studije slučaja</i>	53
4.2. <i>Rezultati i diskusija</i>	54
4.2.1. <i>Međurezultati</i>	54
4.2.2. <i>Rezultati prognoze</i>	60
4.2.3. <i>Verifikacija prognoze</i>	65
4.2.4. <i>Zauzeće računarskih resursa</i>	71
4.2.5. <i>Rezime rezultata</i>	71
5. <i>Zaključak</i>	73
<i>Literatura</i>	75
<i>Biografija</i>	86
<i>Spisak radova</i>	87

SPISAK SKRAĆENICA

<i>Skraćenica</i>	<i>Pun naziv</i>
ACE	Average Coverage Error
ACF	Autocorrelation Function
Adam	Adaptive moment estimation
AHC	Agglomerative HC
AHDO	Aproksimirani HDO
ANN	Artificial Neural Network
CH	Calinski-Harabasz
CNN	Convolutional Neural Network
CRPS	Continuous Ranked Probability Score
CVRMSE	Coefficient of Variation of RMSE
DB	Davies-Bouldin
DCL	Deep Centroid Learning
DCLN	DCL Network
DCNN	Deep CNN
DeepAR	Deep Auto-Regression
DEL	Deep Ensemble Learning
DFT	Discrete Fourier Transform
DHC	Divisive HC
DM	Distributivna mreža
DMN	Deep Mixture Network
DMS	Distribution Management System
DNN	Deep Neural Network
DP	Distributivno preduzeće
DWT	Discrete Wavelet Transform
ECCF	Electricity Consumer Characterization Framework
EMD	Empirical Mode Decomposition
FNN	Feedforward Neural Network
FSC	Fast Spectral Clustering
GOA	Grasshopper Optimization Algorithm
GPU	Graphic Processing Unit
GRU	Gated Recurrent Unit
HBKM	Hybrid Bisect KM
HC	Hierarchical Clustering
HDO	Hronološki dijagram opterećenja
HMLPSA	Hierarchical Multiresolution Linear-function-based PSA
HMTSR	Hierarchical Multiresolution Time Series Representation
HTM	Hierarchical Temporal Memory
HWKM	Hartigan-Wong KM
JMIM	Joint Mutual Information Maximisation
KM	k -means

KnA	K-means and Agglomerative
KPNAHDO	Konkatenirani PNAHDO
LPR	Load Pattern Recognition
LSTM	Long Short-Term Memory
MAE	Mean Absolute Error
MAPE	Mean Absolute Percentage Error
MDN	Mixture Density Network
MIMO	Multi-Input Multi-Output
MTSMS	Multiresolution Time Series Management System
NAHDO	Normalizovani AHDO
PAA	Piecewise Aggregate Approximation
PACF	Partial ACF
PC	Partitional Clustering
PCA	Principal Component Analysis
PNAHDO	Prosečni NAHDO
PS	Pinball Score
PSA	Piecewise Statistical Approximation
RMSE	Root Mean Square Error
RNN	Recurrent Neural Network
S2S	Sequence-to-Sequence
SAE	Stacked Autoencoders
SCADA	Supervisory Control And Data Acquisition
SCL	SVM Centroid Learning
SDR	Sparse Distributed Representation
SGSC	Smart Grid Smart City
<i>sigm</i>	Sigmoidna funkcija
SOM	Self-Organizing Map
SP1	Skup podataka 1 (UMass)
SP2	Skup podataka 2 (SGSC)
SSA	Singular Spectrum Analysis
SVM	Support Vector Machines
<i>tanh</i>	Hiperbolični tangens
TSGA	Time Series Grouping Algorithm
UPGMA	Unweighted Pair Group Method with Average
VMD	Variational Mode Decomposition
WS	Winkler Score

SPISAK SIMBOLA

<i>Simbol</i>	<i>Opis</i>
\oplus	Sabiranje matrica
\otimes	Hadamardov proizvod matrica
\odot	Spajanje kompozitnih agregacionih funkcija
$\mathbb{1}$	Hevisajdova funkcija
$A_{(n,m)}$	AHDO za n -tog potrošača i m -tu fizičku veličinu
$\dot{A}_{(n,m)}$	NAHDO za n -tog potrošača i m -tu fizičku veličinu
$\ddot{A}_{(n,m)}$	PNAHDO za n -tog potrošača i m -tu fizičku veličinu
\ddot{A}_n	KPNAHDO za n -tog potrošača
\mathcal{A}	Skup AHDO
$\dot{\mathcal{A}}$	Skup NAHDO
$\ddot{\mathcal{A}}$	Skup PNAHDO
$\ddot{\mathcal{A}}'$	Skup KPNAHDO
a	Aktivaciona funkcija
b_i	Aditivni težinski faktor izlazne kapije
b_k	Aditivni težinski faktor kandidata za ulaznu kapiju
b_u	Aditivni težinski faktor ulazne kapije
b_z	Aditivni težinski faktor kapije zaborava
b_k	Aditivni težinski faktor za k -ti neuron
$C_{(k,m)}$	Centroid opterećenja za k -tu grupu i m -tu fizičku veličinu
\mathcal{C}	Skup centroida opterećenja
c_i	Stanje ćelije u i -tom vremenskom koraku
D	Broj vremenskih koraka na rezoluciji prognoze u jednom danu
D^*	Broj vremenskih koraka izmerenih vrednosti u jednom danu
\mathcal{D}	Skup karakterističnih tipova dana
d	Karakteristični tip dana
E	Entropija
e	Penalizacija greške tokom treniranja
$F_{(n,m)}$	Normalizacioni faktor za n -tog potrošača i m -tu fizičku veličinu
\hat{F}_i	Kumulativna prognozirana raspodela verovatnoće u i -tom vremenskom koraku
\mathcal{F}	Skup normalizacionih faktora
f	Reprezentaciona funkcija
G_k	k -ta potrošačka grupa
\mathcal{G}	Skup potrošačkih grupa
g	Kompozitna agregaciona funkcija
H	Dužina horizonta prognoze
\mathcal{H}	HMTSR model
h_i	Skriveno stanje ćelije u i -tom vremenskom koraku
I	Zajedničke informacije

K	Broj klastera (potrošačkih grupa)
K_{max}	Minimalan broj klastera (potrošačkih grupa)
K_{min}	Maksimalan broj klastera (potrošačkih grupa)
L_i	Donja granica intervala predviđanja u i -tom vremenskom koraku
\mathcal{L}	Skup karakterističnih tipova potrošača
M	Broj fizičkih veličina
N	Broj potrošača (broj objekata za klasterizaciju)
O	Kompleksnost algoritma
P	Raspodela verovatnoće
P_1	Broj naizmenično postavljenih konvolucionih slojeva i slojeva sažimanja
P_2	Broj konvolucionih filtera
P_3	Veličina konvolucionih filtera i koraka sažimanja
P_4	Dužina dvosmernog S2S konteksta
P_5	Broj SAE enkoder i dekoder slojeva
P_6	Broj SAE neurona u srednjem sloju
P_7	Procenat neurona za odbacivanje tokom treniranja
\mathcal{P}	Skup karakterističnih perioda
p	Karakteristični period
Q	Broj kvantila
R	Vremenska rezolucija prognoze
\mathcal{R}	Skup dužina vremenskih koraka
S	Tekuća vremenska serija
$S_{(n,m)}$	HDO za n -tog potrošača i m -tu fizičku veličinu
\mathcal{S}	Skup HDO
T	Broj vremenskih koraka
\mathcal{T}	Domen vremena
t	Vremenski trenutak
t_i	Vreme u i -toj tački vremenske serije
U_i	Gornja granica intervala predviđanja u i -tom vremenskom koraku
V	Broj atributa koji predstavljaju enkodirane vremenske odrednice
\mathcal{V}	Domen vrednosti
v_i	Vrednost u i -toj tački vremenske serije
W	Broj odabranih meteoroloških faktora
W^*	Broj dostupnih meteoroloških faktora
w_i	Multiplikativni težinski faktor izlazne kapije
w_k	Multiplikativni težinski faktor kandidata za ulaznu kapiju
w_u	Multiplikativni težinski faktor ulazne kapije
w_z	Multiplikativni težinski faktor kapije zaborava
$w_{(i,k)}$	Multiplikativni težinski faktor za i -ti ulazni atribut i k -ti neuron
\mathcal{X}	Skup odabranih ulaznih atributa
\mathcal{X}^*	Skup poznatih ulaznih atributa
\mathcal{X}_i	Skup vrednosti ulaznih atributa u i -tom vremenskom koraku

x	Ulazni atribut
x_i	Vrednost ulaznog atributa u i -tom vremenskom koraku
$x_{(i)}$	i -ti ulazni atribut
Y	Skup izlaznih atributa
Y_i	Skup vrednosti izlaznih atributa u i -tom vremenskom koraku
y	Izlazni atribut
y_i	Vrednost izlaznog atributa u i -tom vremenskom koraku
$y_{(i)}$	i -ti izlazni atribut
\hat{y}_i	Prognozirana vrednost izlaznog atributa y u i -tom vremenskom koraku.
$\hat{y}_{(i,q)}$	Prognozirana granična vrednost q -tog kvantila izlaznog atributa y u i -tom vremenskom koraku
Z	Odabrano zaostajanje
z	Zaostajanje
α	Statistički nivo značajnosti intervala predviđanja
γ	Širina SVM kernela
Δ	Dužina vremenskog koraka
δ	Najveći zajednički delilac skupa dužina vremenskih koraka
θ	Ugao enkodirane vremenske odrednice
λ	Vrednost vremenske odrednice
Λ	Period ponavljanja vremenske odrednice
μ	Prosečna vrednost
ν	Parametar kojim se kontrolišu donja granica broja <i>support</i> vektora i gornja granica broja trening primera koji se mogu smatrati pogrešnim
ρ	Tolerancija greške za rano prekidanje treniranja
σ	Standardna devijacija
τ	Redni broj kvantila izražen u opsegu (0, 1)
Φ_k	Nagib prave koja povezuje k -tu tačku vremenske serije sa prethodnom tačkom.
φ	Širina intervala predviđanja
ψ	Penalizacija odstupanja van intervala predviđanja
ω	Koeficijent mešanja normalnih raspodela
∇_k	Odsečak prave koja povezuje k -tu tačku vremenske serije sa prethodnom tačkom.

SPISAK SLIKA

Slika 1.1 – Podela prognoze opterećenja prema horizontu [16].....	18
Slika 1.2 – Podele prognoze opterećenja.....	20
Slika 2.1 – Zajedničke informacije i entropija [16].....	26
Slika 2.2 – Odnos između greške generalizacije i složenosti regresionog modela [16].....	29
Slika 2.3 – Matematički model neurona [141].....	30
Slika 2.4 – FNN [141].....	30
Slika 2.5 – RNN [57].....	30
Slika 2.6 – SAE [144].....	31
Slika 2.7 – LSTM ćelija [146].....	32
Slika 2.8 – GRU ćelija [146].....	32
Slika 2.9 – Prognoza za više vremenskih koraka unapred [143].....	33
Slika 2.10 – S2S arhitektura [147].....	34
Slika 2.11 – CNN [149].....	34
Slika 2.12 – MDN [151].....	35
Slika 2.13 – DMN [56].....	40
Slika 2.14 – DEL [60].....	40
Slika 3.1 – Glavni koraci predloženog DCL rešenja [65].....	41
Slika 3.2 – Ilustracija glavnih koraka predloženog DCL rešenja.....	41
Slika 3.3 – Algoritam za kreiranje predloženog HMTSR modela [70].....	43
Slika 3.4 – Primer aproksimacije jedne vremenske serije primenom predložene HMLPSA metode [70].	43
Slika 3.5 – Ilustracija koraka predloženog TSGA rešenja [28].....	45
Slika 3.6 – Primer redukcije količine podataka o opterećenju primenom predložene LPR metode [65]...	47
Slika 3.7 – Struktura predloženog DCLN modela.....	50
Slika 4.1 – Prosečan dnevni HDO za potrošače iz SP1 (a) i SP2 (b).....	53
Slika 4.2 – CH indeks za klastere potrošača iz SP1 (a) i SP2 (b).....	54
Slika 4.3 – Centroidi opterećenja za potrošače iz SP1.....	55
Slika 4.4 – Centroidi opterećenja za potrošače iz SP2.....	55
Slika 4.5 – Rezultat odabira zaostajanja za centroide opterećenja.....	56
Slika 4.6 – Uticaj meteoroloških faktora na agregirano opterećenje potrošača iz SP1.....	57
Slika 4.7 – Uticaj meteoroloških faktora na agregirano opterećenje potrošača iz SP2.....	58
Slika 4.8 – Uticaj promene odabranih metoda za pripremu podataka na tačnost prognoze.....	59
Slika 4.9 – Uticaj promene odabranih vrednosti hiperparametara DCLN modela na tačnost prognoze....	59
Slika 4.10 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 1 iz SP1.....	60
Slika 4.11 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 2 iz SP1.....	60
Slika 4.12 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 1 iz SP2.....	61
Slika 4.13 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 2 iz SP2.....	61
Slika 4.14 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 3 iz SP2.....	61
Slika 4.15 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 4 iz SP2.....	62
Slika 4.16 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 5 iz SP2.....	62
Slika 4.17 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 6 iz SP2.....	62
Slika 4.18 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 7 iz SP2.....	63
Slika 4.19 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 8 iz SP2.....	63
Slika 4.20 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 9 iz SP2.....	63
Slika 4.21 – PS za potrošače iz SP1.....	65
Slika 4.22 – PS za potrošače iz SP2.....	65
Slika 4.23 – PS na nivou potrošača iz SP1.....	66
Slika 4.24 – PS na nivou potrošača iz SP2.....	66

Slika 4.25 – PS na nivou percentila za potrošače iz SP1.....	67
Slika 4.26 – PS na nivou percentila za potrošače iz SP2.....	67
Slika 4.27 – PS i potrošena električna energija na dnevnom nivou za potrošače iz SP1	68
Slika 4.28 – PS i potrošena električna energija na dnevnom nivou za potrošače iz SP2	68
Slika 4.29 – WS, φ i ψ za potrošače iz SP1	69
Slika 4.30 – WS, φ i ψ za potrošače iz SP2.....	69
Slika 4.31 – Rezultati Diebold-Mariano testova	70
Slika 4.32 – Vreme izvršavanja naspram tačnosti prognoze	72

SPISAK TABELA

Tabela 4.1 – Skupovi podataka.....	52
Tabela 4.2 – Programske biblioteke	54
Tabela 4.3 – Specifikacija hardvera.....	54
Tabela 4.4 – Rezultat odabira meteoroloških faktora	56
Tabela 4.5 – PS na nivou meseci i tipova dana.	64
Tabela 4.6 – Vreme izvršavanja [min].....	71

1. UVOD

Povećanje energetske efikasnosti je esencijalno za smanjenje emisije štetnih gasova u više slojeve atmosfere i rešavanje problema globalnog zagrevanja. Da bi se ublažile katastrofalne posledice antropogenog globalnog zagrevanja, 196 zemalja je 2016. potpisalo dokument, poznat kao Pariski sporazum [1]. Zemlje potpisnice sporazuma su se obavezale da se ukupna emisija štetnih gasova smanji sa projektovanih 55 Gt na 40 Gt, sa ciljem da se globalno zagrevanje ograniči znatno ispod 2°C u odnosu na predindustrijsko doba, odnosno da se do 2030. godine ograniči na ispod 1,5°C. Poređenja radi, prosečna globalna temperatura je 2020. bila $1,2 \pm 0,1^\circ\text{C}$ veća u odnosu na period 1850-1900 [2]. Analiza potrošnje električne energije širom sveta pokazuje da je potrošnja električne energije u stambenim objektima jedan od najvećih uzročnika emisije štetnih gasova, kao i da je zagađenje rastuće zbog povećanja ljudske populacije i ekonomskog rasta [3–5]. Prema tome, povećanje efikasnosti elektroenergetskih sistema je od izuzetnog značaja za održivi razvoj civilizacije, a prognoza opterećenja distributivnih mreža (DM)¹ je jedno od sredstava koja su neophodna za postizanje tog cilja.

Problem prognoze opterećenja na niskom naponu u DM predstavlja osnovni predmet istraživanja ove doktorske disertacije i kao takav predstavljen je u nastavku ove glave. Nakon toga je dat pregled stanja u ovoj oblasti i objašnjenja koja su potrebna za razumevanje istraživanja u ovoj oblasti. Zatim su definisani konačni ciljevi i praktične potrebe za istraživanjima koja su realizovana u okviru ove doktorske disertacije. Na kraju ove glave predstavljena je organizacija doktorske disertacije.

1.1. Predmet istraživanja

DM predstavlja deo elektroenergetskog sistema koji vrši distribuciju električne energije od napojnih čvorova u transformatorskim stanicama visoki napon/srednji napon do krajnjih potrošača. Tradicionalne DM su se sastojale isključivo od pasivnih elemenata, kao što su: potrošači, nadzemni vodovi, kablovi, transformatori, kondenzatorske baterije i slično. Električna energija se tradicionalno nije proizvodila u DM. Zbog toga su tokovi snaga bili usmereni isključivo od napojnih čvorova prema potrošačima električne energije. Međutim, moderne DM doživljavaju značajne promene u samom dizajnu i principima upravljanja, sa ciljem da podrže koncept pametnih mreža (eng. *smart grid*). Glavne promene se ogledaju u povećanom broju distribuiranih i obnovljivih izvora, povećanom stepenu upravljive potrošnje i povećanom stepenu automatizacije i broja pametnih uređaja [6, 7]. Navedene promene omogućavaju da se zadovolje rastuće potrebe za potrošnjom električne energije, kao i brzo izvršavanje niza preventivnih i korektivnih akcija u upravljanju DM. Krajnji rezultat navedenih promena je veća efikasnost rada DM, kao i veća pouzdanost napajanja potrošača. Međutim, posledica navedenih promena jeste da je sam pogon DM znatno dinamičniji i podložan velikom broju pogonskih izazova, a samo neki od njih su: dvosmerni tok snaga, narušenost i nepouzdanost vrednosti sigurnosne margine, problemi sa stabilnošću u prelaznim režimima, višestruki izvori u režimima sa kvarom, itd.

Distributivna preduzeća (DP)² su kompanije odgovorne za tehničke i administrativne aktivnosti nad DM, koje obuhvataju: planiranje razvoja i proširenja DM, operativno planiranje pogona DM i upravljanje DM u realnom vremenu. Današnja DP se sreću sa novim izazovima, koji nisu postojali u tradicionalnim (pasivnim) DM, a koji su posledica sve aktivnijih DM u smislu distribuirane proizvodnje [7]. Osnovni izazov je kako efikasno obavljati sve potrebne aktivnosti, u različitim pogonskim i vremenskim uslovima, kada tradicionalne DM po pravilu imaju nizak stepen telemetrisanih podataka, što značajno otežava nadzor pogonskih prilika u realnom vremenu. Drugi izazov je kako odgovoriti na sve strožije zahteve vezane za kvalitet isporučene električne energije i pouzdanost napajanja, koji dolaze kako od individualnih potrošača (kupaca električne energije), tako i od regulatornih i vladinih agencija.

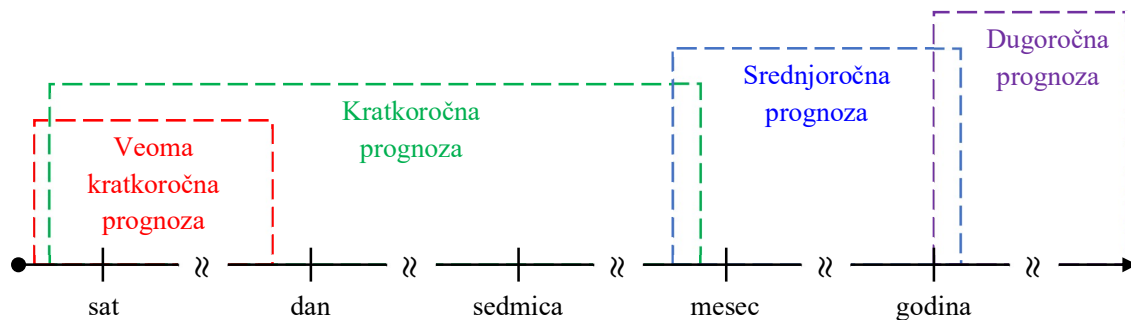
¹ Elektrodistributivna mreža – u nastavku skraćeno distributivna mreža, odnosno DM.

² Elektrodistributivno preduzeće – u nastavku skraćeno distributivno preduzeće, odnosno DP.

Postojeći sistemi za nadzor, upravljanje i prikupljanje podataka (eng. *Supervisory Control And Data Acquisition*, SCADA) nisu dovoljni za efikasno rešavanje navedenih izazova. Zbog toga DP nastoje prvo da investiraju u sisteme za prenos podataka u realnom vremenu, pametne uređaje i automatizaciju kako bi povećali pouzdanost nadzora DM i kako bi omogućili primenu efikasnih korektivnih, preventivnih i optimizacionih upravljačkih akcija. Zatim, DP nastoje da primene sofisticirana softverska rešenja za nadzor, upravljanje i optimizaciju pogona DM. Konačno, implementacija koncepta pametnih mreža uz dovoljno kvalitetnu telekomunikacionu infrastrukturu, sisteme merenja, daljinski kontrolisane uređaje i naravno povrhu svega integrisan sistem za upravljanje distribucijom električne energije (eng. *Distribution Management System*, DMS) omogućava postizanje ciljane efikasnosti [8–10]. DMS ima dvostruku ulogu u DP [7]. Prvo, da pomogne osoblju odgovornom za nadzor i upravljanje DM u donošenju odluka o primeni preventivnih i korektivnih upravljačkih akcija u realnom vremenu. Drugo, da automatski (eng. *closed loop*) optimizuje pogon DM bez intervencije osoblja.

Osnovni deo DMS-a su elektroenergetski proračuni za nadzor, analizu i optimizaciju pogona DM (DMS funkcije). Napredne DMS funkcije, kao što su optimalna kontrola napona i tokova reaktivnih snaga [11] i optimalna rekonfiguracija mreže [6], oslanjaju se na rezultate osnovnih DMS funkcija, kao što su proračun tokova snage [12], estimacija stanja [13, 14], kontrola napona [15, 16] i prognoza opterećenja [17].

Prognoza opterećenja predstavlja prognozu potrošnje aktivne i reaktivne snage na određenoj vremenskoj rezoluciji za određeni vremenski interval u budućnosti (horizont prognoze). Prognoza opterećenja se prema horizontu deli na [18]: veoma kratkoročnu (eng. *very short-term* ili *near-term*), kratkoročnu (eng. *short-term*), srednjoročnu (eng. *medium-term*) i dugoročnu (eng. *long-term*). Veoma kratkoročna prognoza se izvodi na najvišoj vremenskoj rezoluciji (na nivou minuta) i za najuži horizont prognoze (do nekoliko sati unapred). Kratkoročna prognoza se izvodi na nižoj vremenskoj rezoluciji (na nivou minuta ili sata) i za širi horizont prognoze (tipično do nekoliko dana unapred, ali ne više od par sedmica unapred). Srednjoročna prognoza se izvodi na još nižoj vremenskoj rezoluciji (na nivou dana ili sedmice) i za znatno širi horizont prognoze (od mesec dana do nekoliko godina unapred). Konačno, dugoročna prognoza se izvodi na najnižoj vremenskoj rezoluciji (na nivou meseca ili godine) i za najširi horizont prognoze (do više desetina godina unapred). Podela prognoze opterećenja prema horizontu je prikazana na slici 1.1.



Slika 1.1 – Podela prognoze opterećenja prema horizontu [18]

Srednjoročna i posebno dugoročna prognoza imaju strateški značaj za razvoj DM: proširenje kapaciteta DM, kupovinu i održavanje opreme, obnovu zaliha energetskih sirovina, zapošljavanje osoblja i slično. Međutim, kratkoročna i veoma kratkoročna prognoza imaju vitalni značaj za svakodnevne aktivnosti u DM: pravovremeno utvrđivanje optimalnog operativnog stanja, pripremu za buduća opterećenja, planiranje održavanja i slično. Dugoročna prognoza se koristi u svrhe planiranja već više od jednog veka, ali kako su DP počela da teže operativnoj izvrsnosti, kratkoročna prognoza je postepeno pronalazila sve veću primenu [17]. Zbog toga su se u literaturi vremenom pojavile finije podele kratkoročne prognoze³ [17]: prognoza za preostali deo tekućeg dana (eng. *intraday*), prognoza za naredni dan (eng. *day-ahead*), prognoza za narednu sedmicu (eng. *week-ahead*) i slično.

³ Rešenje koje je predloženo u ovoj disertaciji je generalno primenjivo u kratkoročnoj prognozi opterećenja na niskom naponu u DM, ali je verifikacija rešenja izvedena u odnosu na prognozu opterećenja za naredni dan.

Razvoj senzorskih mreža i računarskih sistema omogućio je prikupljanje velike količine podataka o opterećenju DM i prognozu opterećenja zasnovanu na istorijskim podacima [19]. Istovremeni razvoj u oblasti veštačke inteligencije omogućio je analizu podataka primenom metoda mašinskog učenja [20, 21]. Prognoza opterećenja zasnovana na statistici i mašinskom učenju nudi veću fleksibilnost u modelovanju opterećenja i manje zavisi od ekspertskog znanja, a više od dostupnosti i kvaliteta podataka, kao i načina obrade podataka. Osnovna pretpostavka prognoze zasnovane na istorijskim podacima je da će budućnost na neki način nalikovati prošlosti. Prema tome, otkrivanje šablona i skrivenih informacija iz istorijskih podataka je ključno za tačnost takve prognoze [17].

Podaci o opterećenju DM su inherentno hronološki dijagrami opterećenja (HDO) – vremenske serije aktivne i reaktivne snage u određenom vremenskom periodu [22]. Razvoj napredne metričke infrastrukture je omogućio da se takvi podaci prikupljaju kontinualno, u realnom vremenu i hronološkom redosledu, što ih dodatno karakteriše kao tekuće vremenske serije (eng. *streaming time series*) [23]. Dostupnost i kvalitet modela opterećenja zasnovanog na takvim podacima je od suštinskog značaja za kvalitetno upravljanje DM [24]. Međutim, HDO su tipično nelinearne i nestacionarne tekuće vremenske serije [25], tako da je analiza takvih podataka veoma složen zadatak [26–29]. Vrednosti HDO se tipično prikupljaju u jednakim vremenskim intervalima, ali pojava izuzetaka, nevažjećih i nedostajućih vrednosti nije neuobičajena. Pored toga, proširenje nadzora DM implicira povećanje broja posmatranih čvorova DM, a time i povećanje broja HDO, kao i količine podataka, što iziskuje značajno povećanje računarskih resursa koji su neophodni za obradu tih podataka. Na primer, ako se svake minute aktivna i reaktivna snaga milion potrošača beleže sa 4-bajtnom preciznošću, onda se za godinu dana prikupi ~3,8 TiB podataka [30], a za visoku tačnost prognoze često su potrebne dve do tri godine podataka [19].

Prognozu opterećenja koja se zasniva na opisanim podacima je praktično lakše izvoditi na srednjem i visokom naponu nego na niskom naponu jer je količina podataka i fluktuacija HDO prirodno manja [31]. Međutim, ako se nadzor DM sprovodi kod samih potrošača, na niskom naponu, to pruža nove mogućnosti za planiranje i upravljanje DM. Prognoza opterećenja na niskom naponu može doprineti upravljanju DM na više načina [32]: kroz pravovremene reakcije na potražnju električne energije, u integraciji distribuiranih energetske izvora, u kontroli skladištenja električne energije, itd. U mnogim DM širom sveta, niskonaponski deo mreže je organizovan tako da podstanice imaju 1-6 izvoda koji snabdevaju 1-150 potrošača [32]. Prema tome, niskonaponske podstanice mogu snabdevati raznovrsnu kombinaciju domaćinstava i komercijalnih potrošača. Pritom je očekivano da se raznovrsnost potrošača povećava sa povećanjem upotrebe zelenih tehnologija, kao što su: solarni paneli, električna vozila, toplotne pumpe, baterije za skladištenje električne energije, itd. Zbog toga prognoza na nivou agregiranih potrošača na višim naponskim nivoima (npr. na nivou podstanice ili izvoda) nije dovoljna da pruži potpun uvid u različite potrebe potrošača, ulično i saobraćajno osvetljenje, gubitke u mreži, itd. S druge strane, opterećenje na nivou pojedinačnih potrošača visoko varira, što čini prognozu izuzetno teškom.

Uobičajena deterministička prognoza opterećenja postaje praktično neprimenjiva na niskom naponu jer je zbog visoke varijabilnosti opterećenja veoma podložna greškama. Zbog toga se javlja rastuća potreba za probablističkom prognozom opterećenja [33–35]. Deterministička prognoza za jednog potrošača, jednu fizičku veličinu (aktivnu ili reaktivnu snagu) i jedan vremenski korak pruža jednu očekivanu vrednost. S druge strane, probablistička prognoza pruža uvid u neizvesnost očekivane vrednosti u vidu kvantila⁴, intervala predviđanja⁵ ili parametara određene raspodele verovatnoće⁶ [36]. Prema tome, probablistička prognoza omogućava kvantifikaciju varijabilnosti opterećenja na niskom naponu.

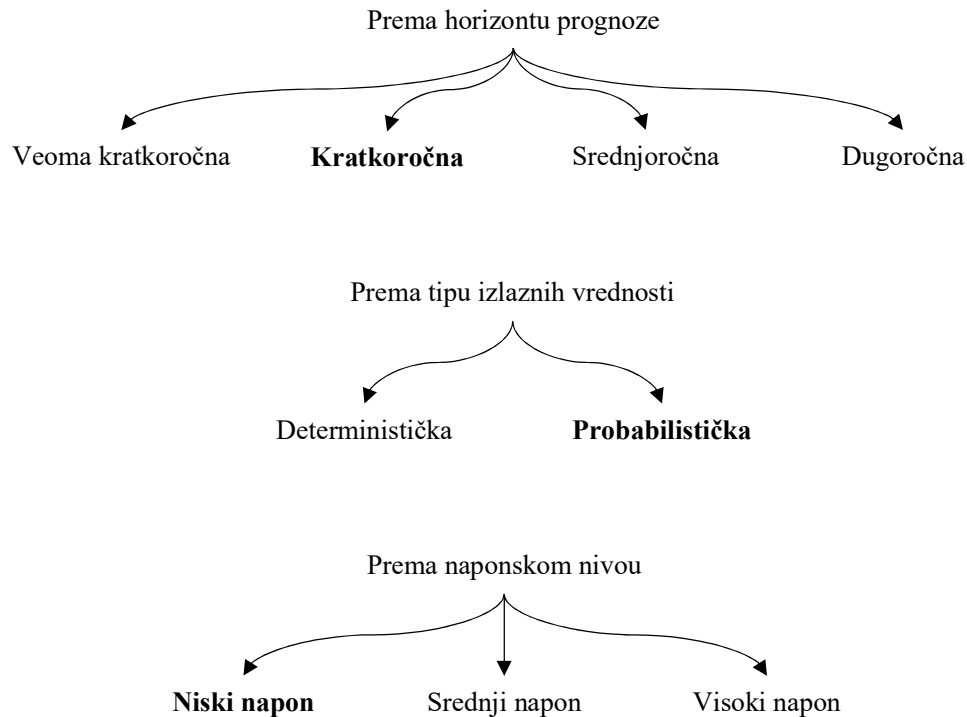
Kvantifikacija varijabilnosti opterećenja u vidu probablističke prognoze je vitalni deo upravljanja rizicima i presudna je za smanjenje operativnih troškova, optimalno odlučivanje i upravljanje DM [37]. Na primer, probablistička prognoza opterećenja omogućava primenu probablističkog proračuna tokova snage [38]. S druge strane, kratkoročna prognoza opterećenja je posebno važna za DP koja pokušavaju da se prilagode rastućem trendu dekarbonizacije i decentralizacije u svakodnevnom radu. Međutim, zbog velike

⁴ Kvantili su tačke podele uređenog skupa na zadat broj jednakih delova [199].

⁵ Intervali predviđanja su procena intervala koji će obuhvatiti buduće vrednosti sa određenom verovatnoćom [199].

⁶ Raspodela verovatnoće je matematička funkcija koja daje verovatnoću pojave određenog ishoda [199].

količine podataka i visoke varijabilnosti opterećenja, kratkoročna probabilistička prognoza opterećenja na niskom naponu u DM, koja je predmet istraživanja ove disertacije, i dalje predstavlja naučno-tehnološki izazov. Predmet istraživanja ove disertacije u kontekstu svih prethodno navedenih podela prognoze opterećenja istaknut je masnim slovima na slici 1.2.



Slika 1.2 – Podele prognoze opterećenja

1.2. Pregled stanja u oblasti

U literaturi se može naći veliki broj rešenja za prognozu opterećenja [39–41]. Pregled različitih rešenja za probabilističku prognozu opterećenja dat je u ovom delu. Od posebnog interesa za ovu disertaciju su rešenja koja su zasnovana na primeni različitih statističkih metoda i metoda mašinskog učenja.

Na temu determinističke prognoze opterećenja je napisano znatno više naučnih radova nego na temu probabilističke prognoze [17]. Međutim, deterministička prognoza inherentno ne pruža uvid u varijabilnost posmatranih procesa, zbog čega probabilistička prognoza pronalazi sve veću primenu [19]. Ako se posmatra iz perspektive teorije sistema, sistem probabilističke prognoze se može konstruisati izmenom delova sistema determinističke prognoze [17, 19]: ulaza, regresionog modela⁷ i izlaza.

Probabilistička prognoza se može ostvariti simulacijom u prostoru podataka i parametara regresionog modela na ulaznoj strani. Najčešće se simuliraju meteorološki uslovi [33]. Simulacija se izvodi na osnovu prethodno zabeleženih [42, 43] ili generisanjem novih scenarija [44–46]. Probabilistička prognoza se dobija statističkom analizom rezultata determinističkih prognoza u različitim scenarijima. Međutim, za pouzdanu statističku analizu je neophodno izvršiti prognozu mnogo puta, što zahteva značajne računarske resurse. Prognoza koja se zasniva na nekoliko desetina do nekoliko hiljada scenarija je uobičajena [33].

Probabilistička prognoza se može ostvariti analizom regresionog modela na izlaznoj strani. Neka od postojećih rešenja se zasnivaju na analizi grešaka u prognozi [47, 48]. Analiza grešaka omogućava da se određeni nedostaci determinističke prognoze umanje u njenom probabilističkom obliku, ali je ograničena

⁷ Regresioni model je model podataka koji opisuje povezanost datih ulaznih i izlaznih vrednosti i koji omogućava izračunavanje izlaznih vrednosti na osnovu novih ulaznih vrednosti [20].

u poboljšanju tačnosti probabilističke prognoze [33]. Druga rešenja, koja se sreću u većem broju radova, zasnivaju se na kombinovanju rezultata organizovanih skupova (ansambala) regresionih modela [49–55]. Osnovni problem ansambala je što linearne kombinacije rezultata često nisu zadovoljavajuće u praksi, dok nelinearne kombinacije zahtevaju složenu optimizaciju težinskih faktora koji se pridružuju rezultatima [39].

Konačno, probabilistička prognoza se može ostvariti izmenom regresionog modela tako da pruži uvid u raspodelu verovatnoće. Jedno rešenje je da se za svakog potrošača kreira ansambl regresionih modela u kojem je svaki model usmeren ka prognozi određenog kvantila [56]. Drugo rešenje je da se ansambl zameni jednim regresionim modelom koji omogućava prognozu svih kvantila odjednom [57]. Međutim, da bi uvid u raspodelu verovatnoće bio dovoljno detaljan, često se prognoziraju percentili⁸, čime se značajno povećava broj izlaza iz regresionog modela i vreme izvršavanja prognoze. Zbog toga je u [58] predloženo da se broj izlaza iz regresionog modela smanji kreiranjem modela tako da umesto kvantila pruža prognozu parametara mešavine zadanog broja normalnih raspodela⁹.

Novija istraživanja u oblasti prognoze vremenskih serija pokazuju da tačnost prognoze pojedinačnih vremenskih serija može biti poboljšana kreiranjem regresionih modela nad grupama sličnih vremenskih serija [59–61]. Slična istraživanja u oblasti prognoze opterećenja potvrđuju da grupisanje potrošača prema sličnosti šablona njihovih HDO vodi ka uklanjanju izuzetaka, što ima pozitivan uticaj na tačnost prognoze [40]. Na opisanom principu se zasniva rešenje za probabilističku prognozu opterećenja koje je predloženo u [62]. Međutim, prema [62], regresioni modeli se kreiraju na nivou potrošačkih grupa, ali svaki model i dalje služi za prognozu svih kvantila na nivou svih potrošača unutar grupe. Prema studiji slučaja koja je predstavljena u četvrtoj glavi ove disertacije, takvo rešenje vodi ka povećanju tačnosti prognoze, ali i ka visokom zauzeću računarskih resursa.

Teoretske osnove regresionih modela koji se često sreću literaturi i konkurentnih rešenja iz aktuelnog stanja u oblasti detaljno su predstavljene u drugoj glavi ove disertacije.

1.3. Potrebe za istraživanjem i ciljevi istraživanja

Prelazak sa determinističke na probabilističku prognozu je verovatno najvažniji iskorak u novijoj istoriji prognoze opterećenja [17]. Probabilistička prognoza opterećenja pruža uvid u moguću varijabilnost budućeg opterećenja i omogućava efikasnije upravljanje DM. Značaj kratkoročne probabilističke prognoze opterećenja je privukao veliku pažnju istraživača [33]. Međutim, prognoza opterećenja na niskom naponu je izuzetno malo istražena [63]. Teoretski i praktično, takva prognoza je izazovnija od prognoze opterećenja na srednjem i visokom naponu zbog prirodno većih varijacija u opterećenju i veće količine podataka, koji su posledica većeg broja posmatranih čvorova DM.

Postojeća rešenja za prognozu opterećenja na niskom naponu usmerena su pretežno na prognozu agregiranog opterećenja, npr. na prognozu na nivou izvoda [32], ili podstanice [50]. Međutim, iako je tačnost takve prognoze prirodno veća [31], ona ne pruža potpuni uvid u zahteve potrošača, gubitke na mreži i slično. Jedan od problema koji se često zanemaruje u prethodnim istraživanjima je prognoza opterećenja na nivou stambenih objekata [17]. Kako nadzor DM nastavlja da se širi, spušta sve više ka korisniku, tako i prognoza opterećenja na nivou stambenih objekata postaje sve veći interes elektroenergetskih i građevinskih inženjera [17]. Prema tome, **prva potreba** za istraživanjem je potreba za proračunom kratkoročne probabilističke prognoze koji kvantifikuje varijabilnost opterećenja na niskom naponu u DM.

Opterećenje DM nelinearno zavisi od velikog broja prostorno i vremenski promenljivih faktora, kao što su meteorološki faktori. Zbog toga tradicionalne statističke (linearne) metode, kao što su linearna regresija i autoregresija, često nisu dovoljne za dovoljno tačnu prognozu [20]. S druge strane, mašinsko učenje omogućava otkrivanje nelinearnih zavisnosti. Duboko učenje (podskup mašinskog učenja) omogućava otkrivanje i najsloženijih zavisnosti, što je posebno značajno za prognozu visoko varijabilnog opterećenja. Međutim, primena sofisticiranih regresionih metoda u prognozi opterećenja zavisi i od računarskih resursa

⁸ Percentili su kvantili koji dele uređen skup na 100 jednakih delova [199].

⁹ Normalna raspodela je raspodela verovatnoće čiji su parametri očekivana (prosečna) vrednost i standardna devijacija (prosečno odstupanje od očekivane vrednosti) [199].

koji su neophodni da se regresija završi u konačnom (prihvatljivom) vremenskom roku. Prema tome, **druga potreba** za istraživanjem je potreba za povećanjem tačnosti prognoze opterećenja na niskom naponu u DM bez neopravdanog povećanja računarskih resursa.

Globalni cilj istraživanja koje je obuhvaćeno ovom doktorskom disertacijom je da se razvije novo rešenje za kratkoročnu probabilističku prognozu opterećenja na niskom naponu koje će uvažiti varijabilnost opterećenja i ponuditi konkurentnu tačnost prognoze uz visoku efikasnost sa stanovišta zauzeća računarskih resursa. Radi postizanja navedenog globalnog cilja, identifikovani su sledeći individualni ciljevi:

- 1) Razviti novo rešenje tako da redukuje model podataka o opterećenju na niskom naponu bez gubitka značajnih informacija o varijabilnosti opterećenja i na taj način smanji zauzeće računarskih resursa u prognozi.
- 2) Razviti novo rešenje tako da omogući primenu sofisticiranih regresionih metoda nad redukovanim modelom podataka i na taj način ponudi konkurentnu tačnost prognoze.
- 3) Verifikovati opravdanost primene predloženog rešenja nad skupom realnih podataka.
- 4) Uporediti predloženo rešenje sa konkurentnim rešenjima iz aktuelnog stanja u oblasti.

1.4. Pregled doktorske disertacije

Doktorska disertacija je organizovana na sledeći način:

- U prvoj glavi je predstavljen predmet istraživanja, dat je pregled stanja u oblasti i predstavljene su potrebe za istraživanjem i ciljevi istraživanja.
- Druga glava sadrži teoretske osnove statističkih metoda i metoda mašinskog učenja na kojima se zasniva predlog novog rešenja za prognozu opterećenja, kao i teoretske osnove konkurentnih rešenja iz aktuelnog stanja u oblasti.
- Treća glava sadrži predlog novog rešenja za kratkoročnu probabilističku prognozu opterećenja na niskom naponu u DM – *Deep Centroid Learning* (DCL).
- U četvrtoj glavi su predstavljeni rezultati studije slučaja koja je izvedena sa ciljem verifikacije DCL rešenja nad skupom realnih podataka iz jedne severnoameričke i jedne australijske DM [64, 65]. Predloženo rešenje je pritom upoređeno sa konkurentnim rešenjima iz aktuelnog stanja u oblasti.
- Na kraju, u petoj glavi, predstavljen je zaključak disertacije i pravac daljeg istraživanja.

2. TEORETSKE OSNOVE

Rešenje za kratkoročnu probabilističku prognozu opterećenja, koje je predloženo u ovoj disertaciji, zasniva se na primeni nadgledanog i nenadgledanog mašinskog učenja i statistike. Nadgledano mašinsko učenje se zasniva na treniranju i testiranju. Treniranje je proces otkrivanja povezanosti između ulaznih i izlaznih promenljivih (atributa) na osnovu unapred zadatih vrednosti (primera). Testiranje je proces upotrebe otkrivenih veza za određivanje izlaznih vrednosti na osnovu novih ulaznih vrednosti. S druge strane, nenadgledano učenje se zasniva na otkrivanju povezanosti atributa bez unapred zadatih primera. Mašinsko učenje i statistika se u ovoj disertaciji koriste za rešavanje sledećih zadataka [20, 66]:

- 1) *Reprezentacija*. Reprezentacija ili učenje atributa je proces redukcije dimenzionalnosti podataka (broja atributa). Reprezentacija pripada zadacima i statističkih metoda i metoda mašinskog učenja.
- 2) *Klasterizacija*. Klasterizacija je proces organizacije objekata u grupe, koje se nazivaju klasterima, a koje su ujedno kompaktne i različite. Pritom, jedan objekat predstavlja jedan skup vrednosti svakog atributa u datom skupu podataka. Klasterizacija je zadatak nenadgledanog mašinskog učenja¹⁰.
- 3) *Regresija*. Regresija je proces treniranja i testiranja regresionog modela. Regresioni model je model podataka koji opisuje povezanost datih ulaznih i izlaznih vrednosti i koji omogućava izračunavanje izlaznih vrednosti na osnovu novih ulaznih vrednosti. Regresija je jedan od zadataka tradicionalnih statističkih metoda i metoda nadgledanog mašinskog učenja¹¹.
- 4) *Optimizacija*. Regresija tipično zahteva složenu parametrizaciju regresionog modela. Zbog toga je optimizacija parametara često sastavni deo regresije [5].

Reprezentacija i klasterizacija omogućavaju otkrivanje šablona i skrivenih informacija iz istorijskih podataka o opterećenju (eng. *Load Pattern Recognition*, LPR) [30, 67]. Zbog toga pronalaze veliku primenu u pripremi podataka za prognozu opterećenja [23, 68]. Regresija se ističe kao glavno sredstvo za prognozu opterećenja, a optimizacija kao neophodno sredstvo za treniranje regresionih modela. Konačno, verifikacija prognoze se izvodi merenjem tačnosti. Verifikacija prognoze se koristi i u svojstvu funkcije greške (gubitka) tokom optimizacije, što je čini izuzetno važnom za proračun prognoze.

U nastavku ove glave prvo su opisane metode reprezentacije, klasterizacije, regresije, optimizacije i verifikacije. Neke od njih su upotrebljene za formulisanje predloga rešenja, dok su ostale metode pobrojane radi kompletnosti pregleda teoretskih osnova. Nakon toga su detaljnije opisana rešenja za probabilističku prognozu, koja su upoređena sa predloženim rešenjem u okviru studije slučaja koja je predstavljena u četvrtoj glavi ove disertacije.

2.1. Reprezentacija

Dimenzionalnost podataka koji se koriste u analizi opterećenja DM ima različito značenje u kontekstu klasterizacije i regresije. Prilikom klasterizacije, potrošač se posmatra kao objekat, a vrednosti u njegovim HDO kao vrednosti njegovih atributa. U tom slučaju, dimenzionalnost podataka je određena brojem atributa koji pripadaju svim potrošačima. U slučaju regresije, broj vrednosti u HDO određuje broj trening primera, dok broj ulaznih atributa (faktora opterećenja) određuje dimenzionalnost podataka. Ako je dimenzionalnost suviše mala, onda posmatrani skup podataka ne nudi dovoljno informacija o posmatranom procesu. S druge strane, podaci velike dimenzionalnosti neizbežno postaju rasuti u višedimenzionalnom prostoru. To stvara rastuću potrebu za povećanjem količine podataka kako bi rezultati analize podataka bili statistički značajni. Opisani fenomen se naziva prokletstvom dimenzionalnosti [69].

¹⁰ Pored klasterizacije, nenadgledano mašinsko učenje se koristi za rešavanje drugih zadataka, kao što su pronalaženje asocijativnih pravila i detekcija anomalija [66], ali oni nisu sastavni deo predloženog rešenja u ovoj disertaciji.

¹¹ Pored regresije, nadgledano mašinsko učenje se koristi za rešavanje još jednog zadatka, a to je klasifikacija: proces označavanja ulaznih vrednosti unapred datim izlaznim vrednostima. Klasifikacione metode mogu biti prilagođene regresiji i obrnuto. Međutim, klasifikacija kao takva nije deo predloženog rešenja u ovoj disertaciji.

Redukcija podataka bez gubitka značajnih informacija je put ka smanjenju računarskih resursa koji su potrebni za analizu tih podataka [30]. Uopšteno, postoji mnogo načina za reprezentaciju vremenskih serija [70], kao i mnogi sistemi za upravljanje takvim podacima [71]. Međutim, dok su sofisticirane metode računarski zahtevne, suviše jednostavne metode mogu ukloniti značajne informacije o dinamici posmatranih procesa i ograničiti dalju analizu podataka [72]. Zbog toga analiza HDO tipično zahteva primenu više različitih metoda reprezentacije. Postojeće metode reprezentacije (učenja atributa) se prema nameni dele na metode ekstrakcije (transformacije i agregacije) atributa i metode odabira atributa.

2.1.1. Ekstrakcija atributa

HDO su tipično nelinearne i nestacionarne tekuće vremenske serije, što otežava identifikaciju fluktuacija od značaja za prognozu [25]. Vrednosti HDO se tipično prikupljaju u jednakim vremenskim intervalima, ali pojava izuzetaka, nevažjećih i nedostajućih vrednosti nije neuobičajena. Prema tome, očuvanje fundamentalnih osobina HDO u manjem modelu podataka zahteva ekstrakciju atributa koja uzima u obzir varijacije i prirodni kontinuitet posmatranih procesa [72].

Aproksimacija vremenskih serija je jedan način ekstrakcije atributa [70]. U zavisnosti od toga da li su aproksimirane vrednosti ravnomerno (regularno) ili neravnomerno (neregularno) raspoređene u vremenu, metode aproksimacije se dele na neadaptivne i adaptivne, respektivno. Adaptivne metode, kao što je *Piecewise Linear Approximation* [73] i mnoge njene adaptacije [74–76], su bolje u predstavljanju pojedinačnih segmenata vremenskih serija [70], ali otežavaju primenu klasterizacije i regresije. Na primer, izračunavanje rastojanja između dve neregularne vremenske serije zahteva interpolaciju vrednosti jedne vremenske serije za svaki vremenski trenutak u drugoj i obrnuto [77]. Druge neadaptivne metode obuhvataju aproksimaciju vremenskih serija u vidu teksta [78], slike [79] i heširanih vrednosti [80].

Piecewise Aggregate Approximation (PAA) [81] je tradicionalna neadaptivna metoda. Rezultat PAA nad segmentom vremenske serije je prosek vrednosti unutar segmenta. Prema tome, PAA opisuje očekivano ponašanje, ali ne i varijabilnost modelovanog procesa. Različite adaptacije PAA metode su predložene kako bi se adresirali njeni nedostaci. Jedna od takvih adaptacija je *Piecewise Cloud Approximation* [82, 83], koja opisuje segmente vremenskih serija prosekom, entropijom i devijacijom vrednosti unutar segmenata [84]. Druga adaptacija i ujedno generalizacija PAA metode je *Piecewise Statistical Approximation* (PSA) [85], koja opisuje segmente vremenskih serija centralnim momentima vrednosti unutar segmenata, kao što su prosek i standardna devijacija. Međutim, nijedna od navedenih metoda ne opisuje kontinuitet modelovanog procesa između vremenskih trenutaka zabeleženih vrednosti, jer se vrednosti posmatraju nezavisno jedna od druge. Zbog toga je u ovoj disertaciji predložena modifikovana verzija PSA metode, koja ujedno opisuje varijabilnost i kontinuitet modelovanog procesa (opterećenja DM) [72].

Neadaptivne metode aproksimacije zahtevaju da je vremenska rezolucija rezultujućih vrednosti unapred poznata. Međutim, u slučaju modelovanja HDO to ne mora biti slučaj. Distributivna preduzeća mogu prikupljati HDO i po nekoliko godina pre nego što se ti podaci upotrebe za prognozu. Zahtevi koji su vezani za vremensku rezoluciju prognoze, a koji potiču kako od distributivnih preduzeća tako i od regulatornih i vladinih agencija, se u međuvremenu mogu promeniti. Ako su HDO aproksimirani na jednoj vremenskoj rezoluciji, onda promena rezolucije zahteva retroaktivno procesiranje HDO. Međutim, takvo procesiranje je nepraktično za podatke o opterećenju na niskom naponu u velikim DM, sa stanovišta računarskih resursa koji bi bili potrebni za skladištenje i obradu podataka u prihvatljivom vremenskom roku. Jedan način za prevazilaženje problema nesigurnosti u odabiru vremenske rezolucije prilikom aproksimacije HDO je upotrebe više različitih vremenskih rezolucija (višerezoluciona aproksimacija) [86–89].

Generički model za višerezolucionu aproksimaciju vremenskih serija, *Multiresolution Time Series Management System* (MTSMS) model, je predložen u [90]. Prema MTSMS modelu, svaka vremenska serija se posmatra kroz predefinisane reprezentacione funkcije koja se aproksimira na zadatim vremenskim rezolucijama. Na taj način, MTSMS model omogućava da kontinuitet modelovanog procesa bude opisan odgovarajućom reprezentacionom funkcijom. Aproksimirane vrednosti se dobijaju iterativno, primenom predloženog algoritma koji održava dve generičke strukture podataka za svaku vremensku seriju i svaku

vremensku rezoluciju: bafer i disk. Bafer predstavlja privremeno skladište originalnih vrednosti vremenske serije u okviru poslednjeg vremenskog koraka respektivne vremenske rezolucije. S druge strane, disk predstavlja trajno skladište aproksimiranih vrednosti iz bafera u okviru svih prethodnih vremenskih koraka. Prema tome, predloženi model omogućava da se vremenske serije iterativno aproksimiraju, ali zahteva dvostruku obradu njihovih vrednosti: jednom kada se skladište u bafer i drugi put kada se aproksimiraju i skladište na disk. Slična rešenja su predložena u [23, 91–95]. Takođe postoje mnoga komercijalna rešenja koja se zasnivaju na sličnim principima [96–102]. Međutim, u slučaju aproksimacije HDO, takav pristup vodi ka neopravdanom povećanju vremena izvršavanja i kašnjenju u dostupnosti aproksimiranih vrednosti za prognozu opterećenja. Zbog toga je u ovoj disertaciji predložena modifikovana verzija MTSMS modela koja omogućava aproksimaciju HDO u jednom prolazu [72].

Dekompozicija vremenskih serija je još jedan način ekstrakcije atributa koji se, pored aproksimacije, često koristi u analizi HDO. Postojeće metode dekompozicije dele se na komponentne i frekvencijske [25]. Tradicionalne komponentne metode su *Principal Component Analysis* (PCA) i *Singular Spectrum Analysis* (SSA) [103]. PCA se uopšteno primenjuje nad matricom u kojoj redovi predstavljaju objekte, kao što su vremenske serije, dok kolone predstavljaju attribute objekata, kao što su vrednosti vremenskih serija. PCA se zasniva na dekompoziciji date matrice na sopstvene vektore i na otkrivanju linearno nekorelisanih kolona matrice (glavnih komponenti). S druge strane, SSA se zasniva na dekompoziciji pojedinačnih vremenskih serija primenom PCA metode. Tradicionalne frekvencijske metode su *Discrete Fourier Transform* (DFT) [104], *Discrete Wavelet Transform* (DWT) [105] i *Empirical Mode Decomposition* (EMD) [106]. EMD se pokazala bolje od DFT i DWT nad nestacionarnim vremenskim serijama, ali je osetljiva na mešanje bliskih komponenti [25]. *Variational Mode Decomposition* (VMD) [107] prevazilazi navedeni problem adaptivnim određivanjem opsega frekvencija, zbog čega se pokazala bolje od EMD u prognozi opterećenja [108, 109]. Predloženo rešenje u ovoj disertaciji obuhvata primenu PCA metode u pripremi podataka za klasterizaciju, jer PCA omogućava efikasnu klasterizaciju HDO na osnovu Euklidskog rastojanja [110].

2.1.2. Odabir atributa

Tačnost rezultata nadgledanog mašinskog učenja u velikoj meri zavisi od odabira atributa [111]. Ekspertska analiza je neophodna da bi se odredio početni skup ulaznih atributa, ali nedovoljna da bi se odredio konačan skup [112]. U početnom skupu se mogu naći ulazni atributi koji jesu u korelaciji sa izlaznim, ali i oni koji su nerelevantni, suvišni ili čak nisu u korelaciji sa izlaznim atributima. Na primer, ne moraju svi poznati meteorološki faktori biti od značaja za prognozu opterećenja.

Visoka tačnost prognoze se tipično postiže primenom metoda odabira atributa. Glavna odlika takvih metoda je njihova stabilnost [113]: konzistentnost odabira nad različitim podskupovima datog skupa podataka. Nestabilan odabir vodi ka pogrešnim zaključcima. Jedan od najvećih uzroka nestabilnosti je odbacivanje ulaznih atributa koji su u korelaciji sa izlaznim, ali i sa drugim ulaznim atributima. Metode odabira atributa se dele na filtrirajuće, ugrađene i obmotavajuće, u zavisnosti od toga da li se primenjuju pre, u toku ili nakon mašinskog učenja, respektivno [113]. Filtrirajuće metode se zasnivaju na određivanju statističkog značaja atributa. Ugrađene metode su sastavni deo regresionih modela. Obmotavajuće metode se zasnivaju na optimizaciji odabira atributa na osnovu tačnosti prognoze, kao u [114]. Ugrađene i obmotavajuće metode omogućavaju veće poboljšanje tačnosti prognoze, ali su filtrirajuće metode manje algoritamski složene i stabilnije u odabiru atributa [113].

Primena metoda odabira atributa takođe zavisi od toga da li se biraju endogeni ili egzogeni atributi. Endogeni atributi se odnose na vrednosti zavisne promenljive, dok se egzogeni atributi odnose na vrednosti nezavisne promenljive [37]. Na primer, ako se opterećenje sa određenim zaostajanjem (eng. *lag*) koristi kao ulazni atribut u prognozi opterećenja, onda je taj ulazni atribut endogen. S druge strane, meteorološki faktori su egzogeni ulazni atributi.

Odabir endogenih ulaznih atributa u prognozi vremenskih serija, kao i u prognozi opterećenja, tipično se vrši na osnovu analize autokorelacije i parcijalne autokorelacije [115]. Autokorelacija je korelacija između vrednosti vremenske serije u različitim vremenskim trenucima [116]. Parcijalna autokorelacija sa

zaostajanjem z je korelacija vrednosti vremenske serije za vremenske trenutke t i $t - z$ nakon eliminacije uticaja svih vrednosti između t i $t - z$ [116]. *Autocorrelation Function* (ACF) i *Partial ACF* (PACF) su funkcije koje kvantifikuju autokorelaciju i parcijalnu autokorelaciju, respektivno. Ako se ACF i PACF posmatraju kao filtrirajuće metode odabira endogenih ulaznih atributa, onda je PACF pogodniji od ACF za identifikaciju zaostajanja od najvećeg značaja za prognozu opterećenja [117].

Primena tradicionalnih (linearnih) koeficijenata korelacije tipično nije dovoljna za odabir egzogenih ulaznih atributa koji su nelinearno korelisani sa izlaznim atributima, kao što je to slučaj sa meteorološkim faktorima u prognozi opterećenja. Zbog toga se veliki broj filtrirajućih metoda za odabir egzogenih atributa zasniva na određivanju količine informacija koje ulazni i izlazni atributi dele – zajedničkih informacija.

Definicija zajedničkih informacija je izvedena iz teorije informacija, a određuje koliko poznavanje vrednosti jedne promenljive smanjuje neizvesnost neke druge promenljive. Formalno se definiše na sledeći način [18]: Neka su x i y slučajne kontinualne promenljive (ulazni i izlazni atributi, respektivno) i neka je P raspodela verovatnoće. Onda je neizvesnost atributa y definisana entropijom E na sledeći način:

$$E(y) = - \int P(y) \log(P(y)) dy \quad (2.1)$$

Dalje, neizvesnost atributa y u odnosu na x je definisana uslovnom entropijom:

$$E(y|x) = - \int P(x) \int P(y|x) \log(P(y|x)) dy dx \quad (2.2)$$

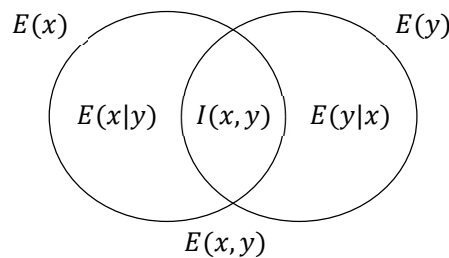
Zatim, zajednička neizvesnost atributa x i y definisana je zajedničkom entropijom:

$$E(x, y) = - \int \int P(x, y) \log(P(x, y)) dx dy \quad (2.3)$$

Konačno, zajedničke informacije atributa x i y su definisane na sledeći način:

$$I(x, y) = E(x) + E(y) - E(x, y) = E(x) - E(x|y) = E(y) - E(y|x) \quad (2.4)$$

Zajedničke informacije i entropija atributa x i y prikazani su Venn-ovim dijagramom na slici 2.1. Skup na levoj strani predstavlja entropiju atributa x , a skup na desnoj strani entropiju atributa y . Unija skupova predstavlja zajedničku entropiju. Razlika skupa s leve strane u odnosu na skup s desne strane predstavlja uslovnu entropiju atributa x . Razlika skupa s desne strane u odnosu na skup s leve strane predstavlja uslovnu entropiju atributa y . Presek skupova predstavlja zajedničke informacije.



Slika 2.1 – Zajedničke informacije i entropija [18]

Prednost filtrirajućih metoda zasnovanih na teoriji informacija je sposobnost otkrivanja nelinearnih povezanosti atributa bez prethodnog poznavanja strukture njihovih odnosa [18]. Jedna takva metoda je *minimal-Redundancy-Maximal-Relevance* [118], koja održava ravnotežu između redundanse i povećanja značaja odabranih atributa. Prednosti primene navedene metode u prognozi opterećenja su evidentne [119]. Međutim, u primeni navedene metode, kao i mnogih izvedenih metoda, nailazi se na dva problema [111]:

- 1) *Odbacivanje značajnog atributa*. Ulazni atribut, iako je *značajan* jer nudi nove informacije o izlaznom atributu, može biti odbačen kao *redundantan* zbog informacija koje deli sa prethodno odabranim atributima.

- 2) *Odabir redundantnog atributa*. Ulazni atribut, iako je *redundantan* jer sa nekima od ulaznih atributa deli veliku količinu informacija, može biti odabran kao *značajan* zbog velike količine informacija koje nudi o izlaznom atributu.

Joint Mutual Information Maximisation (JMIM) [111] adresira navedene probleme. JMIM se zasniva na određivanju količine informacija koju svaki kandidovani ulazni atribut i svaki prethodno odabrani ulazni atribut zajedno dele sa izlaznim atributom. Značaj svih ulaznih atributa je izražen realnim brojem u opsegu $[0, 1]$ (JMIM skorom). Opisana metoda vodi ka stabilnijem odabiru atributa i povećanju tačnosti prognoze [113]. Zbog toga je JMIM metoda sastavni deo rešenja koje je predloženo u ovoj disertaciji.

JMIM se formalno definiše na sledeći način. Neka je \mathcal{X}^* skup poznatih ulaznih atributa, \mathcal{X} skup odabranih ulaznih atributa, a y izlazni atribut. Onda se odabir članova skupa \mathcal{X} vrši primenom pohlepnog algoritma prema sledećem kriterijumu:

$$\arg \max_{x' \in \mathcal{X}^* \setminus \mathcal{X}} (\min_{x \in \mathcal{X}} I(x', x; y)) \quad (2.5)$$

gde $I(x', x; y) = I(x, y) + I(x', y|x)$ predstavlja udružene zajedničke informacije ulaznih atributa $\{x', x\}$ u odnosu na izlazni atribut y .

2.2. Klasterizacija

Klasterizacija je tipičan način za grupisanje potrošača na osnovu sličnosti njihovih HDO [66, 69]. Klasterizacija omogućava prepoznavanje šablona opterećenja uz uklanjanje izuzetaka, što pozitivno utiče na dalju analizu podataka [40].

Klasterizacija se zasniva na organizaciji objekata u klasterne na osnovu rastojanja između vrednosti njihovih atributa. Prema tome, odabir mere rastojanja je izuzetno važan. Od mnogih postojećih mera rastojanja između vremenskih serija, dve se ističu po širokoj primeni [70]: Euklidsko rastojanje i *Dynamic Time Warping*. Nedavno je predložena još jedna mera rastojanja, *Shape-Based Distance* [120], koja je primenjena u klasterizaciji vremenskih serija u [121]. Međutim, Euklidsko rastojanje se pokazalo izuzetno efikasnim u klasterizaciji regularnih vremenskih serija (sa ravnomerno raspoređenim vrednostima u vremenu) [70], zbog čega se koristi u okviru rešenja koje je predloženo u ovoj disertaciji [30].

Ako mera rastojanja omogućava prepoznavanje sličnih oblika HDO, onda normalizacija sprečava da HDO sličnih oblika budu svrstani u različite klasterne zbog razlike u njihovim apsolutnim vrednostima. Prema tome, odabir metode normalizacije je takođe izuzetno važan. HDO tipično sadrže mnogo izuzetaka. Zbog toga uobičajene metode, kao što su Z i *min-max* normalizacija, mogu rezultovati velikim rastojanjima između normalizovanih HDO iako su oni sličnih oblika. Deljenje vrednosti sadržanih u HDO sa prosečnim opterećenjem u toku godine je praktičnije u prisustvu izuzetaka [72] i zbog toga je sastavni deo predloženog rešenja u ovoj disertaciji [30].

Pored mere rastojanja i normalizacije, za klasterizaciju je veoma važna i mera (indeks) validnosti. Klasterizacija se može validirati poznavanjem tačnih rezultata (eksterno) ili bez tog znanja (interno) [122]. Klasterizacija HDO tipično zahteva primenu interne validacije, jer tačan rezultat nije unapred poznat. Opširno poređenje indeksa validnosti u [122] ukazuje na značajnu prednost tri indeksa u odnosu na ostale: *Silhouette* [123], Calinski-Harabasz (CH) [124] i Davies-Bouldin (DB) [125]. Sva tri indeksa pružaju uvid u kompaktnost i različitost klastera. Međutim, primena DB indeksa u klasterizaciji HDO vodi ka izuzetno iskrivljenoj raspodeli broja članova klastera [72] (relativno mali broj klastera sadrži izuzetno veliki broj članova, dok ostali klasteri sadrže izuzetno mali broj članova). S druge strane, *Silhouette* indeks se zasniva na upotrebi matrice rastojanja između objekata, što iziskuje značajne računarske resurse. Na primer, za 10^6 objekata je potrebna matrica rastojanja veličine $10^{6 \times 2} \times 4$ bajta $\approx 3,6$ TiB [30]. Nedavno je predložen i *Compact-Separate Proportion* indeks [126], ali takođe zahteva upotrebu matrice rastojanja. Zbog toga je CH praktičniji za klasterizaciju HDO u velikim DM, i kao takav primenjen je u okviru predloženog rešenja u ovoj disertaciji [30].

Metode klasterizacije tipično zahtevaju da je broj klastera unapred zadat [122]. Da bi se formirali ujedno najkompaktniji i najrazličitiji klasteri potreban je optimalan broj klastera [126]. Međutim, taj broj

tipično nije unapred poznat [127]. Prema tome, pronalaženje optimalnog broja klastera je izazovan zadatak [127], koji je neophodan za povećanje validnosti klastera [128, 129].

Postojeće metode klasterizacije se mogu podeliti u kategorije prema modelu podataka na kojem se zasnivaju. Dve kategorije se ističu po najširoj primeni [69, 70]: particionalna klasterizacija (eng. *Partitional Clustering*, PC) i hijerarhijska klasterizacija (eng. *Hierarchical Clustering*, HC). PC se zasniva na iterativnom razvrstavanju objekata u klastere [70]. Najrasprostranjeniji PC algoritam je k -means (KM) [70]. Klasični KM je veoma brz, ali i osetljiv na početan odabir predstavnika klastera (centroida) [130]. Mnoge adaptacije su predložene kako bi se prevazišao ovaj problem. Jedna od njih je KM++ [131], gde se opisani problem izbegava pažljivim odabirom početnih centroida. Još jedna poznata adaptacija je Hartigan-Wong KM (HWKM) [130], gde se opisani problem adresira naknadnom rekonfiguracijom klastera. HWKM je sastavni deo predloženog rešenja u ovoj disertaciji [30].

HC se zasniva na kreiranju hijerarhije klastera (binarnog stabla) [70]. Stablo se kreira ili od korena ka listovima (eng. *Divisive HC*, DHC) ili od listova ka korenu (eng. *Agglomerative HC*, AHC). DHC počinje sa jednim skupom svih objekata, koji se zatim rekurzivno deli do formiranja ciljanog broja klastera [132, 133]. AHC počinje kreiranjem matrice rastojanja između svih objekata, koji se zatim rekurzivno spajaju do korena stabla. Izgrađeno stablo se ne može dalje podešavati, a ciljani broj klastera se dobija sečom grana [134]. Zbog toga je izbor kriterijuma spajanja dva skupa objekata na određenom nivou hijerarhije veoma važan za validnost primene AHC. Autori u [127] su pokazali da je za klasterizaciju HDO najbolji izbor *Unweighted Pair Group Method with Average* (UPGMA). UPGMA se zasniva na spajanju onih skupova objekata čije je prosečno rastojanje između svih parova objekata najmanje na posmatranom nivou hijerarhije. AHC sa UPGMA kriterijumom spajanja je sastavni deo predloženog rešenja u ovoj disertaciji [30].

AHC postaje superioran u odnosu na DHC sa povećanjem broja klastera [135], ali nailazi na problem kreiranja matrice rastojanja. Zbog toga je u [136] predložena *K-means and Agglomerative* (KnA) metoda, po kojoj se AHC primenjuje nad centroidima klastera dobijenih KM algoritmom. Rezultati dobijeni na ovaj način su u visokoj korelaciji sa rezultatima standardne AHC metode, ali je upotrebljena matrica rastojanja značajno manja. Na primer, ako se 10^6 objekata zameni sa 10^3 centroida, onda će matrica rastojanja zauzimati $10^{3 \times 2} \times 4$ bajta $\approx 3,8$ MiB umesto $\sim 3,6$ TiB. Prema tome, u ovoj disertaciji je predložena nova metoda za pronalaženje optimalnog broja klastera vremenskih serija koja se zasniva na primeni HWKM i AHC metoda, kao i na iterativnoj seči grana dobijenih UPGMA kriterijumom spajanja [30].

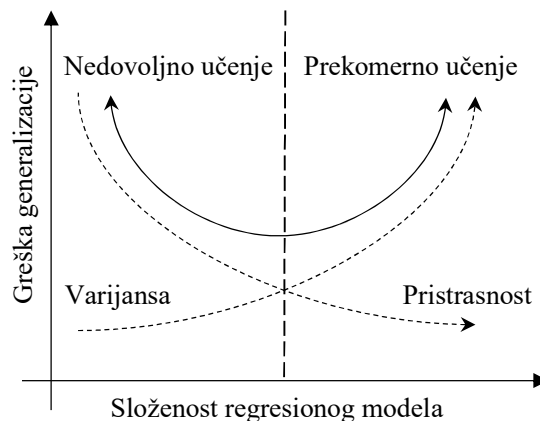
U prethodnim radovima su predložene mnoge druge kombinacije različitih metoda klasterizacije. U okviru *Hybrid Bisect KM* (HBKM) algoritma [132] se DHC i AHC koriste analogno KnA metodi. Na sličan način se u okviru *Fast Spectral Clustering* (FSC) algoritma [137] primenjuje prvo DHC, a zatim spektralna klasterizacija, koje se zasniva na teoriji grafova. Nasuprot tome, u okviru algoritma koji je predložen u [135] primenjuje se prvo AHC, pa zatim DHC, što inherentno podleže problemu veličine matrice rastojanja. U [129] je predložen *Electricity Consumer Characterization Framework* (ECCF), u okviru kojeg se prvo primenjuje *Self-Organizing Map* (SOM)¹² [127], kako bi se redukovao broj objekata za klasterizaciju, a zatim se redukovani broj objekata klasteruje koristeći KM.

U [138] je KM upotrebljen za kreiranje hijerarhije klastera od listova ka korenu, ali bez korišćenja kriterijuma spajanja kao kod AHC metode. Međutim, predloženi algoritam zahteva složenu parametrizaciju (početan broj klastera, broj hijerarhijskih nivoa, najmanje dozvoljeno rastojanje između centroida klastera na svakom hijerarhijskom nivou i slično). Nekoliko algoritama za pronalaženje optimalnog broja klastera je predloženo u [139–142], ali oni takođe zahtevaju složenu parametrizaciju. Sa velikim brojem HDO na niskom naponu, složenija parametrizacija potencijalno vodi ka značajnom povećanju vremena izvršavanja bez značajnog povećanja validnosti klastera. Zbog toga se metoda za pronalaženje optimalnog broja, koja je predložena u ovoj disertaciji, zasniva na jednostavnoj parametrizaciji [30].

¹² SOM je veštačka neuronska mreža koja služi za reprezentaciju visokodimenzionalnih podataka. Veštačke neuronske mreže su klasično rešenje iz oblasti veštačke inteligencije koje se može prilagoditi rešavanju različitih zadataka, kao što su reprezentacija, klasterizacija, klasifikacija i regresija. Rešenje koje je predloženo u ovoj disertaciji se zasniva na primeni veštačkih neuronskih mreža u regresiji i zbog toga su one detaljnije opisane u delu 2.2.

2.3. Regresija

Regresione metode nadgledanog mašinskog učenja su prirodno primenjive u prognozi opterećenja. Osnovni problem takvih metoda je generalizacija: prilagođavanje regresionog modela podacima u trening skupu, ali i u test skupu [18, 68]. Generalizacija se postiže balansiranjem između pristrasnosti i varijanse (eng. *bias-variance tradeoff*). Pristrasnost je pojednostavljenje modela zarad pojednostavljenja proračuna, dok je varijansa osetljivost modela na vrednosti ulaznih atributa. Prevelika pristrasnost se ogleda u nesposobnosti modela da se prilagodi trening skupu i uzrok je nedovoljnog učenja (eng. *underfitting*). Prevelika varijansa se ogleda u nesposobnosti modela da se prilagodi test skupu i uzrok je prekomernog učenja (eng. *overfitting*). Greška generalizacije se povećava u oba navedena slučaja. Na slici 2.2 je prikazan odnos između greške generalizacije i složenosti modela preko pristrasnosti i varijanse. U nastavku ovog dela su predstavljene najpre osnove poznatih regresionih metoda nadgledanog mašinskog učenja, a zatim i osnove metoda dubokog učenja.



Slika 2.2 – Odnos između greške generalizacije i složenosti regresionog modela [18]

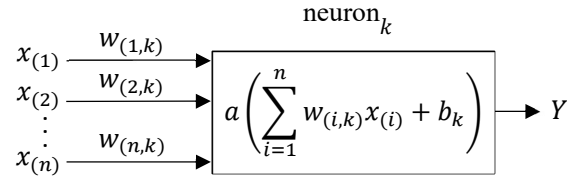
2.3.1. Osnovne metode

U prognozi opterećenja najbolje su se pokazala dva klasična rešenja iz oblasti veštačke inteligencije [3]: *Support Vector Machines* (SVM) i veštačke neuronske mreže (eng. *Artificial Neural Network*, ANN). Oba rešenja imaju visoku sposobnost predstavljanja složenih nelinearnih povezanosti atributa i generalno su primenjiva u rešavanju klasifikacionih i regresionih zadataka [143]. Zbog toga su oba rešenja primenjena u ovoj disertaciji. Međutim, ANN je sastavni deo predloženog rešenja, dok je SVM primenjen za potrebe verifikacije predloženog rešenja u okviru studije slučaja u četvrtoj glavi ove disertacije.

SVM se zasniva na transformaciji ulaznog skupa podataka u višedimenzionalni prostor pomoću odgovarajuće funkcije (kernela) i na pronalaženju hiperravni koja deli transformisane podatke u ciljanom prostoru [5]. SVM prvenstveno služi za klasifikaciju, ali se može koristiti i za regresiju [20]. S druge strane, ANN je inspirisan biološkim nervnim sistemom, gde se signali preko sinapsi prenose do neurona, koji se aktiviraju, obrađuju signal i prosleđuju ga dalje kroz nervni sistem. SVM ispoljava manju grešku generalizacije nego ANN, koji je složeniji i ima veću sklonost ka prekomernom učenju. Međutim, ANN omogućava formiranje složenih struktura podataka koje su sposobne za otkrivanje dubljih nelinearnih povezanosti atributa (duboko učenje) [144]. Zbog toga je u ovoj disertaciji predložen nov regresioni model koji se zasniva na primeni dubokog učenja.

U računarskom smislu, ANN je sačinjen od ulaznog, skrivenog i izlaznog sloja neurona koji redom predstavljaju ulazne, apstraktne i izlazne attribute. Neuroni su povezani težinskim faktorima koji se određuju u fazi treniranja. Multiplikativni težinski faktor (eng. *weight*) predstavlja jačinu sinapse, koja kontroliše uticaj ulazne vrednosti na izlaznu vrednost. Aditivni težinski faktor (eng. *bias*) predstavlja konstantu koja kontroliše da li će se i u kojoj meri neuron aktivirati. Neurone u skrivenom sloju čine aktivacione funkcije koje transformišu sumu ulaznih vrednosti sa pridruženim težinskim faktorima u izlazne vrednosti.

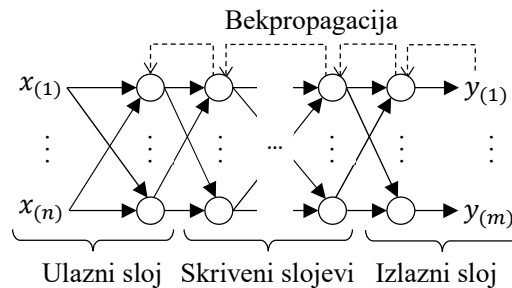
Na slici 2.3 je predstavljen matematički model jednog (k -tog) neurona u okviru jednog ANN modela. Neuron je predstavljen kvadratom, a sinapse strelicama. Ulazni i izlazni atributi su označeni skupovima $\mathcal{X} = \{x_{(1)}, \dots, x_{(n)}\}$ i $Y = \{y_{(1)}, \dots, y_{(m)}\}$, respektivno. Multiplikativni težinski faktor za i -ti ulazni atribut i k -ti neuron je označen sa $w_{(i,k)}$. Aditivni težinski faktor za k -ti neuron je označen sa b_k . Suma ulaznih vrednosti sa pridruženim težinskim faktorima prosleđuje se aktivacionoj funkciji a . Rezultat aktivacione funkcije se prosleđuje svim izlazima sa kojima je k -ti neuron povezan.



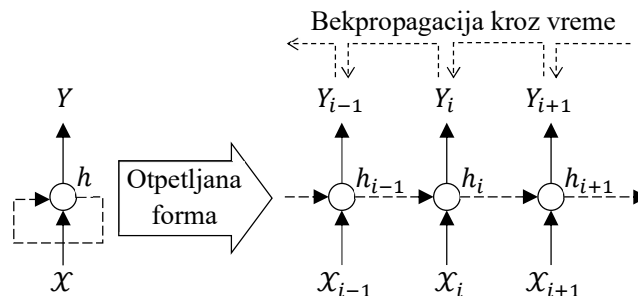
Slika 2.3 – Matematički model neurona [143]

U odnosu na vrstu veze između neurona, ANN može biti *Feedforward Neural Network* (FNN) ili *Recurrent Neural Network* (RNN). FNN prosleđuje informacije isključivo u jednom smeru kroz slojeve neurona (unapred), dok RNN sadrži i povratne (rekurentne) veze. FNN se tipično trenira bekpropagacijom [68] – algoritmom koji iterira kroz slojeve neurona od izlaza ka ulazu i u svakoj iteraciji određuje gradijent funkcije gubitka u odnosu na težine i zatim koriguje težine da bi se smanjio gubitak. S druge strane, RNN se trenira bekpropagacijom kroz vreme – algoritmom koji optimizuje težine rekurentne sekvence neurona u iteracijama od poslednjeg ka prvom elementu sekvence. RNN je pogodan za prognozu vremenskih serija, jer može da predstavi zavisnosti između vrednosti u datoj sekvenci [68, 145]. Predloženi regresioni model u ovoj disertaciji je kombinacija FNN i RNN modela.

FNN i RNN modeli su prikazani na slikama 2.4 i 2.5, respektivno. Neuroni su prikazani kao kružići, a sinapse kao strelice punih linija. Bekpropagacija je prikazana isprekidanim strelicama. Analogno slici 2.3, FNN povezuje ulazne attribute \mathcal{X} sa izlaznim atributima Y kroz više slojeva neurona. S druge strane, RNN transformiše ulazne vrednosti atributa \mathcal{X} u i -tom vremenskom koraku (X_i) u skriveno stanje h_i , koje se zatim transformiše u izlazne vrednosti atributa Y u i -tom vremenskom koraku (Y_i) i prenosi do vremenskog koraka $i + 1$. RNN je predstavljen u skupljenoj formi (onako kako izgleda njegova struktura) i u otpetljanoj formi (onako kako se posmatra tokom treniranja).



Slika 2.4 – FNN [143]

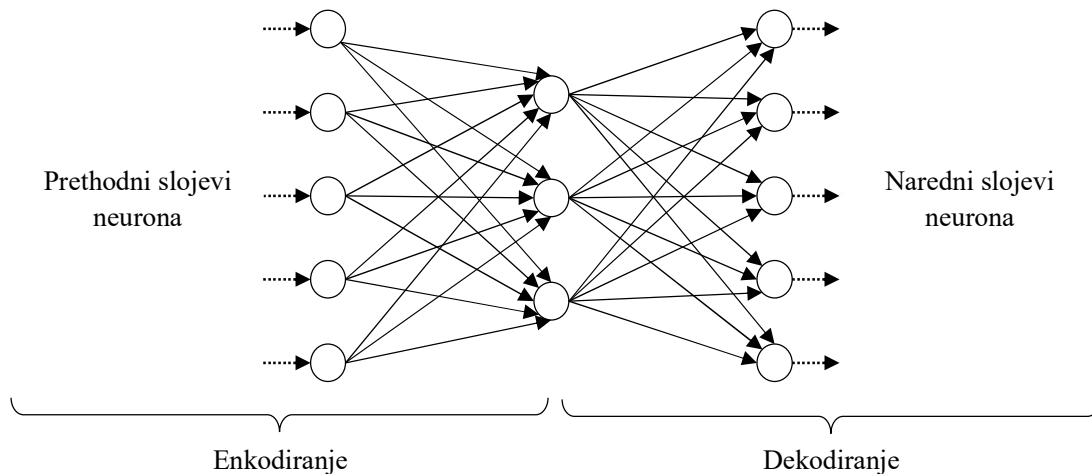


Slika 2.5 – RNN [59]

2.3.2. Duboko učenje

U odnosu na broj skrivenih slojeva, ANN može biti plitak (sa nekoliko skrivenih slojeva) ili dubok (sa znatno više skrivenih slojeva). Na primer, FNN na slici 2.4 je dubok. Duboko učenje se zasniva na dubokim ANN modelima (eng. *Deep Neural Network*, DNN). Dublji modeli imaju veći kapacitet za predstavljanje složenih nelinearnih povezanosti atributa uz uklanjanje izuzetaka, što je značajna prednost u prognozi opterećenja [115].

Stacked Autoencoder (SAE) [146] je višeslojni FNN koji ima jednak broj neurona na ulazu i izlazu, a manje neurona u sredini. Takva arhitektura omogućava rekonstrukciju ulaznih vrednosti na izlaznoj strani uz otklanjanje izuzetaka. Troslojni SAE je prikazan na slici 2.6. Neuroni s leve strane služe za redukciju dimenzionalnosti ulaznih podataka (enkodiranje), dok neuroni s desne strane služe za rekonstrukciju ulaznih podataka (dekodiranje). SAE je sastavni deo predloženog regresionog modela u ovoj disertaciji.



Slika 2.6 – SAE [146]

Uz veći kapacitet za predstavljanje složenih nelinearnih povezanosti atributa, DNN modeli su takođe podložniji prekomernom učenju. Prekomerno učenje se tipično sprečava regularizacijom. Regularizacija predstavlja mehanizme za redukciju greške generalizacije, kao što su [62]: rano prekidanje faze treniranja (eng. *early stopping*), penalizacija funkcije gubitka (L1 i L2 regularizacija), dodavanje šuma, odbacivanje neurona tokom treniranja (*dropout*), itd. Predloženo rešenje u ovoj disertaciji obuhvata *early stopping* i *dropout* regularizaciju.

Pored prekomernog učenja, postoje još dva problema koja su specifična za RNN, a koja se neizbežno javljaju kada se RNN koristi za predstavljanje dugih sekvenci vrednosti [145]: problemi eksplozivnog i nestajućeg gradijenta. S obzirom da se težine rekurentnih veza menjaju množenjem sa gradijentom funkcije gubitka, neizbežno je da se kod dugih sekvenci vrednosti ispolji jedan od navedenih problema. Problem eksplozivnog gradijenta nastaje ako je gradijent veći od 1, jer se onda i težine i gradijenti postepeno povećavaju, što vodi ka prenaplašavanju značaja starijih informacija. S druge strane, problem nestajućeg gradijenta nastaje ako je gradijent manji od 1, jer se onda i težine i gradijenti postepeno smanjuju, što vodi ka umanjivanju značaja starijih informacija.

Postoje dva tipa dubokih RNN modela koja su sposobna da prevaziđu opisane probleme [145]: *Long Short-Term Memory* (LSTM) i *Gated Recurrent Unit* (GRU). Oba se zasnivaju na mehanizmu kapija kojima se kontroliše zadržavanje informacija u modelu kroz vreme (stanje ćelije). Kapije predstavljaju skrivene slojeve neurona sa pažljivo odabranim aktivacionim funkcijama i operacijama nad matricama. Stanje ćelije predstavlja matricu težina koja na apstraktan način opisuje ulazne podatke. Stanje ćelije se menja prolaskom kroz kapije. LSTM sadrži tri kapije: kapiju zaborava, ulaznu i izlaznu kapiju. Ulazna kapija određuje koje informacije treba dodati u stanje ćelije. Kapija zaborava određuje koje informacije treba izbaciti iz stanja ćelije. Izlazna kapija određuje koje informacije treba zadržati u skrivenom stanju. S druge strane, GRU

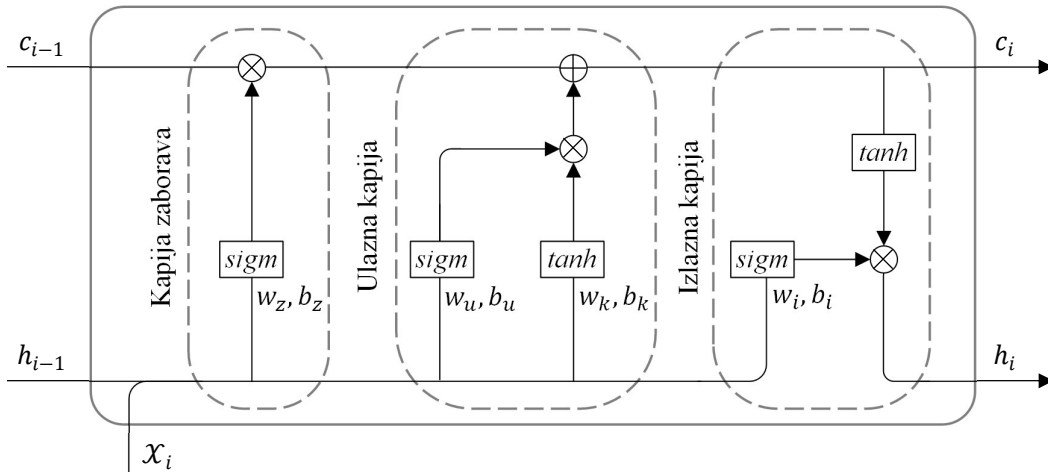
sadrži samo dve kapije, jer su kapija zaborava i ulazna kapija spojene u jednu kako bi se smanjio broj težina koje je potrebno podesiti u fazi treniranja. Međutim, na taj način dolazi do određenog pogoršanja tačnosti prognoze. Zbog toga su LSTM ćelije sastavni deo predloženog regresionog modela u ovoj disertaciji.

LSTM ćelija je predstavljena na slici 2.7. Neuroni su predstavljeni kvadratima, a sinapse strelicama. Spajanje i račvanje strelica ukazuju na uparivanje i kopiranje vrednosti, respektivno. Stanje ćelije, skriveno stanje i skup vrednosti ulaznih atributa u i -tom vremenskom koraku su označeni sa c_i , h_i i X_i , respektivno. Tokom treniranja se podešavaju težine kapije zaborava (w_z i b_z), ulazne kapije (w_u i b_u), kandidata za ulaznu kapiju (w_k i b_k) i izlazne kapije (w_i i b_i). Aktivacione funkcije u neuronima su sigmoidna funkcija ($sigm$) i hiperbolični tangens ($tanh$). Oznake \oplus i \otimes predstavljaju sabiranje i Hadamardov proizvod matrica, respektivno. Na slici 2.7 nije prikazan skup vrednosti izlaznih atributa u i -tom vremenskom koraku (Y_i), jer se te vrednosti uvek mogu dobiti transformacijom skrivenih stanja kao na slici 2.5. Radi kompletnosti, na slici 2.8 je prikazana i GRU ćelija (oznaka "1-" predstavlja oduzimanje od broja 1).

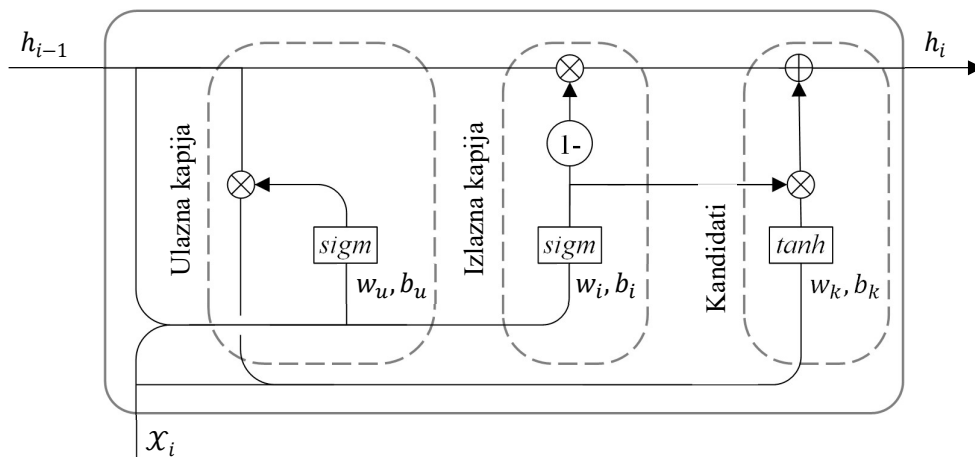
Formalno, aktivacione funkcije $sigm$ i $tanh$ su definisane na sledeći način [147]:

$$sigm(x) = \frac{1}{1 + e^{-x}} \quad (2.6)$$

$$tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (2.7)$$



Slika 2.7 – LSTM ćelija [148]



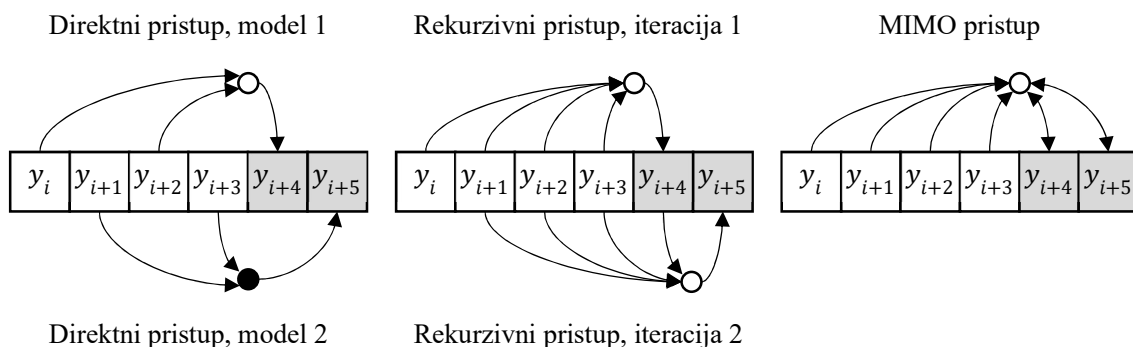
Slika 2.8 – GRU ćelija [148]

Stanje LSTM ćelije se menja na osnovu rezultata aktivacionih funkcija i operacija nad matricama. Sigmoidna funkcija otkriva značaj datih vrednosti tako što ih transformiše u opseg (0, 1): što je rezultat bliži jedinici, to je vrednost značajnija. Hiperbolični tangens kandiduje vrednosti za pamćenje i ujedno kontroliše skalu vrednosti tako što ih zadržava u opsegu (-1, 1). Hadamardov proizvod naglašava značaj vrednosti: vrednosti postaju značajnije ako se množe sa većim brojem. Hadamardov proizvod rezultata sigmoidne funkcije i hiperboličnog tangensa otkriva značaj kandidata i vraća pozitivan broj za kandidate koji se trebaju zapamtiti, a negativan broj za kandidate koji se trebaju zaboraviti. Pamćenje i zaboravljanje se kontrolišu sabiranjem matrica. Vrednost se pamti ako ostaje ista, urezuje u sećanje ako se povećava, bledi iz sećanja ako se smanjuje, a zaboravlja ako se anulira. Prema tome, u stanju ćelije će nakon kapije zaborava ostati naglašen značaj vrednosti. Značajne vrednosti će ostati zapamćene zajedno sa značajnim kandidatima nakon ulazne kapije, a sve ostalo će biti zaboravljeno. Novo stanje ćelije će biti kandidovano za novo skriveno stanje u izlaznoj kapiji. Pritom će značaj kandidata biti naglašen značajem prethodnog skrivenog stanja i tekućih ulaznih vrednosti u posmatranom vremenskom koraku. Na kraju će novo skriveno stanje biti pod većim utiskom novijih informacija, dok će novo stanje ćelije obezbediti dugo pamćenje značajnih informacija.

Jedna od najznačajnijih prednosti RNN modela u prognozi opterećenja je način na koji omogućava upotrebu prethodnog opterećenja za prognozu narednog. Praktičan problem koji se javlja u prognozi opterećenja za više vremenskih koraka unapred je to što podaci o prethodnom opterećenju postaju nedostupni već od drugog koraka. U literaturi se mogu pronaći tri uopštena pristupa za adresiranje ovog problema [145]:

- 1) *Direktni*. Direktan pristup se zasniva na kreiranju po jednog regresionog modela za svaki vremenski korak unapred. Međutim, rezultati različitih regresionih modela zajedno mogu biti nekoherentni i nepovezani.
- 2) *Rekurzivni*. Rekurzivni pristup se zasniva na iterativnom proračunu: svaka prognozirana vrednost postaje jedna od ulaznih vrednosti za prognozu sledeće. Rekurzivni pristup nudi uglačan prelaz između rezultujućih vrednosti, ali neizbežno vodi ka akumulaciji greške na kasnijim vremenskim koracima u horizontu prognoze.
- 3) *Multi-Input Multi-Output (MIMO)*. MIMO pristup se zasniva na primeni specijalizovanih regresionih modela koji izvedu proračun prognoze za sve vremenske korake odjednom. Takav pristup nudi uglačan prelaz između rezultujućih vrednosti bez akumulacije greške kao kod rekurzivnog pristupa. Zbog toga je MIMO pristup sastavni deo predloženog rešenja u ovoj disertaciji.

Slika 2.9 prikazuje sve navedene pristupe na primeru determinističke prognoze za dva vremenska koraka unapred. Vrednosti izlaznog atributa y u i -tom vremenskom koraku su označene sa y_i . Beli i sivi kvadrati predstavljaju istorijske i prognozirane vrednosti, respektivno. Regresioni modeli su predstavljeni kružićima. Strelice koje vode do i od kružića predstavljaju ulazne i izlazne vrednosti, respektivno.



Slika 2.9 – Prognoza za više vremenskih koraka unapred [145]

Sequence-to-Sequence (S2S) arhitektura je organizacija RNN modela koja omogućava otkrivanje vremenske zavisnosti između vrednosti ulazne i izlazne sekvence i na taj način pruža mogućnost za primenu MIMO pristupa u prognozi [145]. S2S arhitekturu čine dva RNN modela: enkoder i dekoder. Enkoder stvara kontekstni vektor od ulazne sekvence i prosleđuje ga dekoderu. Dekoder na osnovu kontekstnog vektora kreira izlaznu sekvencu. Opisana arhitektura je prikazana na slici 2.10. Primena MIMO pristupa u okviru rešenja koje je predloženo u ovoj disertaciji zasniva se na S2S LSTM arhitekturi.

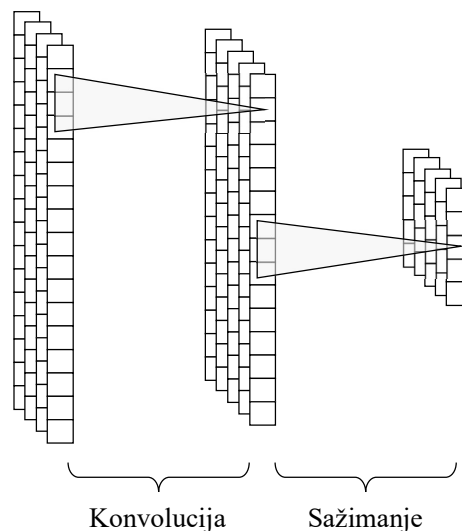


Slika 2.10 – S2S arhitektura [149]

S2S arhitektura sa dvosmernim rekurentnim vezama naglašava zavisnosti između prošlih i budućih stanja RNN modela i pomaže u stabilizaciji prognoze usled intenzivnih promena vrednosti u budućnosti [147]. S druge strane, S2S arhitektura gubi informacije o povezanosti redosleda vrednosti ulazne i izlazne sekvence ako se isti kontekstni vektor koristi za prognozu svih vrednosti. Jedan način da se opisani problem prevaziđe je da se primeni mehanizam pažnje [149]. Mehanizam pažnje se zasniva na upotrebi dodatnih težina koje se pridružuju kontekstnom vektoru za različite pozicije izlaznih vrednosti, što znatno poboljšava tačnost prognoze. Pored toga, RNN slojevi se uvek mogu kombinovati sa FNN slojevima u cilju [21, 145]: 1) redukcije dimenzionalnosti ulaznih podataka; 2) otkrivanja skrivenih informacija između RNN slojeva, ili 3) transformacije podataka u izlaznom sloju. Zbog toga se predloženi regresioni model u ovoj disertaciji zasniva na kombinovanju RNN i FNN slojeva uz primenu mehanizma pažnje.

Convolutional Neural Network (CNN) je FNN koji se u prognozi opterećenja koristi za redukciju dimenzionalnosti podataka (ekstrakciju ulaznih atributa) [150–152]. CNN se sastoji od konvolucionog sloja i sloja sažimanja (eng. *pooling*). Neuronu u konvolucionom sloju vrše konvoluciju ulaznih podataka na osnovu pridružene matrice težina (filtera). Neuronu u sloju sažimanja vrše disjunktnu podelu konvoluiranih podataka na manje celine, čije su dimenzije određene zadatim korakom sažimanja. U zavisnosti od toga da li se primenjuje sažimanje prosekom ili sažimanje maksimumom, iz svake celine se preuzima prosečna ili maksimalna vrednost, respektivno.

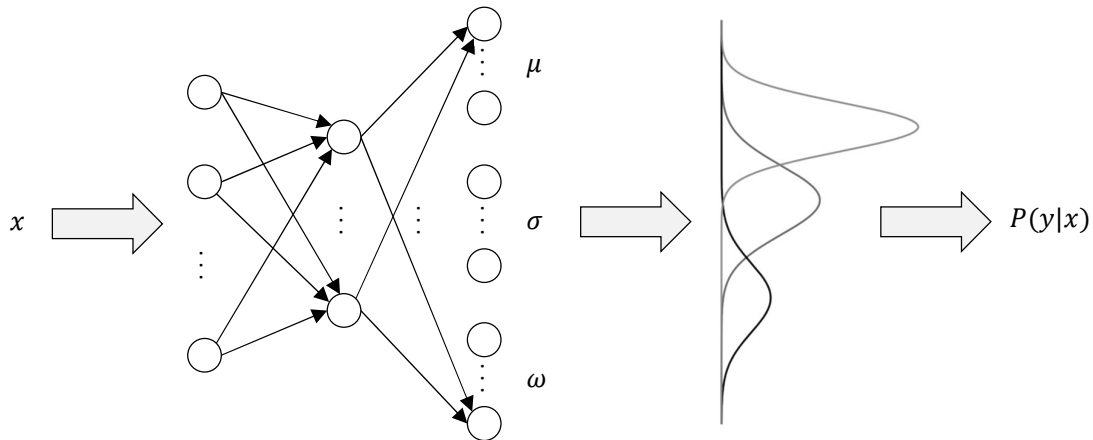
Na slici 2.11 je prikazan jednodimenzionalni CNN koji je primenjen na jednoj vremenskoj seriji. Segmenti vremenske serije su predstavljeni vertikalnim nizovima kvadrata, koji predstavljaju vrednosti vremenske serije. Na slici su prikazana 4 takva segmenta na ulazu, 5 segmenata u sredini i 5 na kraju. Broj segmenata koji nastaju konvolucijom i sažimanjem direktno zavise od broja konvolucionih filtera. S druge strane, dimenzionalnost segmenata zavisi od veličine filtera i koraka sažimanja. U ovom slučaju, CNN ima 5 filtera, a veličina filtera i koraka sažimanja je 3 (trouglovi predstavljaju konvoluciju i sažimanje).



Slika 2.11 – CNN [151]

Duboki CNN (DCNN) [150–152] se dobija naizmeničnim ulančavanjem konvolucionih slojeva i slojeva sažimanja. DCNN je koristan u prognozi opterećenja jer omogućava prepoznavanje lokalnih trendova u vremenskim serijama [150–152]. Kapacitet modela za predstavljanje ulaznih podataka u prostoru manje dimenzionalnosti se kontroliše brojem konvolucionih filtera. S druge strane, redukcija dimenzionalnosti podataka se kontroliše brojem konvolucionih slojeva i slojeva sažimanja, kao i veličinom filtera i koraka sažimanja. Manji broj filtera smanjuje reprezentativnost podataka, dok veći broj filtera povećava vreme treniranja. Ako se broj slojeva smanji, a veličina filtera i koraka sažimanja poveća, onda može doći do gubitka značajnih informacija. Ako se broj slojeva poveća, a veličina filtera i koraka sažimanja smanji, onda se povećava vreme treniranja. Prema tome, sa pažljivo odabranim parametrima, DCNN omogućava redukciju dimenzionalnosti ulaznih podataka bez gubitka značajnih informacija. Zbog navedenih prednosti, DCNN je sastavni deo predloženog regresionog modela u ovoj disertaciji.

Mixture Density Network (MDN) je duboki FNN koji se koristi u probabilističkoj prognozi za transformaciju podataka u izlaznom sloju [153]. MDN transformiše rezultat prethodnih slojeva pomoću specijalne funkcije gubitka tako da odgovara parametrima mešavine raspodela verovatnoće izlaznih vrednosti. Tip i broj raspodela verovatnoće se definiše unapred i određuje broj MDN izlaza. Na slici 2.12 je prikazan MDN koji je podešen tako da prognozira parametre mešavine normalnih raspodela: prosečne vrednosti (μ), standardne devijacije (σ) i koeficijente mešanja (preklapanja) normalnih raspodela (ω). Regresioni model koji je predložen u ovoj disertaciji je takođe dizajniran tako da prognozira parametre normalne raspodele, ali oni ne nastaju unutar predloženog modela već van njega, ekstrakcijom atributa i klasterizacijom. Zbog toga je MDN, kao sastavni kompletnog rešenja za probabilističku prognozu opterećenja koje je predloženo u [58], upoređen sa rešenjem koje je predloženo u ovoj disertaciji, u okviru studije slučaja u četvrtoj glavi.



Slika 2.12 – MDN [153]

Hierarchical Temporal Memory (HTM) je još jedan DNN koji se može primeniti na vremenske serije, a koji simulira rad biološkog neokorteksa [154]. HTM formira binarne reprezentacije podataka sa malim brojem aktivnih neurona: *Sparse Distributed Representation* (SDR). SDR slojevi su organizovani hijerarhijski i omogućavaju otkrivanje povezanosti ulaznih i izlaznih atributa u vremenu. Treniranje i testiranje se obavljaju simultano, dodavanjem i uklanjanjem veza između neurona na osnovu učestalosti njihove aktivacije. Međutim, povezivanje neurona između i unutar SDR slojeva, neizbežno zahteva značajne računarske resurse. Autori [148] ukazuju na to da zbog brzine izvršavanja i tačnosti prognoze vremenskih serija, LSTM i GRU i dalje imaju prednost u odnosu na HTM.

2.4. Optimizacija

Optimizacija je često sastavni deo regresije zbog složene parametrizacije regresionih modela [5]. Parametri regresionih modela se dele na interne (parametre modela) i eksterne (hiperparametre modela). Parametri modela se podešavaju na osnovu datog skupa primera u fazi treniranja i koriste se za kreiranje rezultata u fazi testiranja. Hiperparametri modela se koriste isključivo za potrebe treniranja. Parametrizacija tipično ima nelinearan uticaj na funkciju gubitka. Prema tome, optimalna deterministička parametrizacija

je praktično neizvodljiva. Umesto toga, koriste se optimizacione metode koje dovoljno brzo konvergiraju ka prihvatljivom rešenju [155]. U kontekstu optimizacionih metoda, parametri i hiperparametri modela se posmatraju kao optimizacione promenljive, a domen njihovih vrednosti kao optimizacioni prostor u okviru kojeg se vrši potraga za optimalnim rešenjem.

Adaptive moment estimation (Adam) [156] se može pronaći kao prvi izbor u optimizaciji parametara DNN modela u mnogim radovima [37]. Adam spada u grupu algoritama koji se zasnivaju na određivanju gradijenta funkcije gubitka. Adam adaptira veličinu koraka u potrazi za optimalnim rešenjem u odnosu na procenjenju vrednost proseka i varijanse gradijenta funkcije gubitka. Adam nije računarski zahtevan, zbog čega je pogodan za složene optimizacione probleme u pogledu količine podataka i broja parametara modela. Prema tome, rešenje koje je predloženo u ovoj disertaciji zasniva se na primeni Adam algoritma.

U prognozi opterećenja, hiperparametri modela se često optimizuju metaheurističkim algoritmima koji su inspirisani prirodnim optimizacionim procesima, kao što su evolucija i kolektivna inteligencija [5]. Najpoznatiji algoritmi koji su inspirisani evolucijom su: genetski algoritam [157] i diferencijalna evolucija [158]. Optimizacione promenljive se u navedenim algoritmima posmatraju kao geni, a tačnost prognoze kao sposobnost jedinke da opstane u smeni generacija. Međutim, osnovni problem ovih algoritama je to što prebrza smena generacija smanjuje raznolikost jedinki i vodi ka lokalnom optimumu, dok sporija smena generacija smanjuje skalabilnost algoritama.

Optimizacioni algoritmi koji su inspirisani kolektivnom inteligencijom zasnivaju se na globalnoj i lokalnoj pretrazi optimizacionog prostora u potrazi za optimalnim rešenjem. Globalna pretraga je potraga za obećavajućim regionima u optimizacionom prostoru. Lokalna pretraga je detaljnija inspekcija regiona pronađenih globalnom pretragom. Postizanje ravnoteže između globalne i lokalne pretrage je od suštinskog značaja za konvergenciju ovih algoritama i izbegavanje lokalnih optimuma. Najpoznatiji algoritmi koji su inspirisani kolektivnom inteligencijom su algoritmi koji simuliraju ponašanje životinja u kolektivu [41].

Particle Swarm Optimization [159] je jedan od najpoznatijih tradicionalnih algoritama zasnovanih na kolektivnoj inteligenciji [155]. Inspirisan je kretanjem ptica ili riba u jatu. Ravnoteža između globalne i lokalne pretrage se postiže razmenom informacija između jedinke i jata. Međutim, u prethodnoj deceniji su predloženi efikasniji algoritmi za rešavanje složenih optimizacionih problema kao što je parametrizacija regresionih modela za prognozu opterećenja [160].

Fruit Fly Optimization Algorithm je inspirisan ponašanjem roja voćnih mušica u potrazi za hranom [122, 161]. Voćne mušice odlaze u potragu za hranom u različitim smerovima, a roj se okuplja na onom mestu gde je miris hrane najjači. Glavni nedostatak opisanog algoritma je to što manji broj mušica sporije konvergira ka rešenju, dok su za veći broj mušica potrebni veći računarski resursi [161].

Grasshopper Optimization Algorithm (GOA) je inspirisan ponašanjem roja skakavaca u potrazi za hranom [162–164]. Kao i voćne mušice, skakavci odlaze u potragu za hranom u različitim smerovima i okupljaju se na onom mestu gde je hrana najbolja. Pored toga, kretanje skakavaca je uslovljeno njihovim društvenim interakcijama: silom privlačenja ili odbijanja, kao i težnjom ka postepenom zbližavanju jata oko hrane. Na taj način se postiže ravnoteža između globalne i lokalne pretrage, što čini opisani algoritam efikasnim u izbegavanju lokalnog optimuma i konvergenciji ka globalnom optimumu. Zbog prednosti koje nudi, GOA je primenjen u okviru studije slučaja koja je predstavljena u četvrtoj glavi ove disertacije.

U prethodnoj deceniji je predloženo još mnogo sličnih optimizacionih algoritama, koji su inspirisani ponašanjem životinja u potrazi za hranom ili tokom grupnog lova, kao što su *Grey Wolf Optimizer* [117, 165], *Whale Optimization Algorithm* [166–168] i mnogi drugi [155].

2.5. Verifikacija

Verifikacija prognoze zahteva primenu mera tačnosti. Mere tačnosti pružaju jednostavan i robustan uvid u grešku prognoze [20]. Deterministička prognoza se verifikuje uvidom u razlike između prognoziranih i ostvarenih vrednosti. Međutim, probabilistička prognoza zahteva posmatranje dodatnih kriterijuma.

Najvažniji kriterijumi probabilističke prognoze opterećenja su pouzdanost i oštrina [34]. Pouzdanost opisuje konzistentnost prognoziranе raspodele verovatnoće sa ostvarenim vrednostima. Oštrina opisuje bliskost prognoziranе raspodele verovatnoće i ostvarenih vrednosti. Na primer, 90% intervali predviđanja su pouzdani ako pokrivaju najmanje 90% vrednosti, ali to ne znači da su oštri jer mogu biti značajno udaljeni od vrednosti unutar intervala. Prema tome, cilj probabilističke prognoze se formuliše kao maksimizacija oštine sa ograničenjem pouzdanosti [34].

Mere tačnosti koje se često koriste u determinističkoj prognozi opterećenja su [20]: *Mean Absolute Error* (MAE), *Mean Absolute Percentage Error* (MAPE), *Root Mean Square Error* (RMSE) i *Coefficient of Variation of RMSE* (CVRMSE). Definisane su na sledeći način:

$$MAE = \frac{1}{T} \sum_{i=1}^T |y_i - \hat{y}_i| \quad (2.8)$$

$$MAPE = \frac{100}{T} \sum_{i=1}^T \left| \frac{\hat{y}_i - y_i}{y_i} \right| \quad (2.9)$$

$$RMSE = \sqrt{\frac{1}{T} \sum_{i=1}^T (\hat{y}_i - y_i)^2} \quad (2.10)$$

$$CVRMSE = 100 \frac{RMSE}{\frac{1}{T} \sum_{i=1}^T y_i} \quad (2.11)$$

gde je:

- T – broj vremenskih koraka,
- y_i – ostvarena vrednost izlaznog atributa y u i -tom vremenskom koraku,
- \hat{y}_i – prognozirana vrednost izlaznog atributa y u i -tom vremenskom koraku.

Kod svih navedenih mera tačnosti, veća vrednost ukazuje na veću grešku prognoze. MAE pruža najjednostavniji uvid u grešku, ali tretira sva odstupanja jednako. RMSE penalizira veća odstupanja. MAPE i CVRMSE imaju iste odlike kao MAE i RMSE, respektivno, ali ne zavise od skale podataka, već predstavljaju grešku u procentima.

Jedna od najjednostavnijih i najčešće korišćenih mera tačnosti probabilističke prognoze je *Average Coverage Error* (ACE) [34]:

$$ACE = 100 \left(\frac{1}{T} \sum_{i=1}^T \mathbb{1}\{L_i \leq y_i \leq U_i\} - (1 - \alpha) \right) \quad (2.12)$$

gde je:

- T – broj vremenskih koraka,
- y_i – ostvarena vrednost izlaznog atributa y u i -tom vremenskom koraku,
- L_i i U_i – donje i gornje granice intervala predviđanja u i -tom vremenskom koraku, respektivno,
- α – statistički nivo značajnosti intervala predviđanja
- $\mathbb{1}$ – Hevisajdova funkcija¹³.

ACE je jednak nuli ako intervali predviđanja pokrivaju tačno onoliko vrednosti koliko se očekuje. Ako je pokriveno manje vrednosti, onda je ACE veće od nule. Ako je pokriveno više vrednosti, onda je

¹³ Rezultat Hevisajdove funkcije je 0 za negativne vrednosti argumenta, a 1 za pozitivne vrednosti argumenta.

ACE manje od nule. Na taj način ACE ocenjuje pouzdanost probabilističke prognoze, ali ne i oštrinu. Zbog toga se ACE neće dalje razmatrati u ovoj disertaciji.

Mere tačnosti probabilističke prognoze koje ocenjuju i pouzdanost i oštrinu zasnivaju se na valjanim pravilima bodovanja (eng. *proper scoring rules*). Ta pravila predstavljaju funkcije koje pružaju uvid u razliku između prognozirane i ostvarene raspodele verovatnoće. Najpoznatije od takvih funkcija su: *pinball* i *winkler* funkcija i *Continuous Ranked Probability Score* (CRPS). One su definisane na sledeći način [34]:

$$\text{pinball}(y_i, \hat{y}_{(i,q)}) = \begin{cases} (1 - \tau)(\hat{y}_{(i,q)} - y_i), & y_i < \hat{y}_{(i,q)} \\ \tau(y_i - \hat{y}_{(i,q)}), & y_i \geq \hat{y}_{(i,q)} \end{cases}, \quad \tau = \frac{q}{Q + 1} \quad (2.13)$$

$$\text{winkler}(y_i, L_i, U_i) = \varphi + \psi, \quad \varphi = (U_i - L_i), \quad \psi = \begin{cases} 2(L_i - y_i)/\alpha, & y_i < L_i \\ 0, & L_i \leq y_i \leq U_i \\ 2(y_i - U_i)/\alpha, & y_i > U_i \end{cases} \quad (2.14)$$

$$\text{CRPS}(y_i, \hat{F}_i) = \int_{-\infty}^{\infty} (\hat{F}_i(y) - \mathbb{1}\{y_i \leq y\})^2 dy \quad (2.15)$$

gde je:

- Q – broj kvantila,
- y_i – ostvarena vrednost izlaznog atributa y u i -tom vremenskom koraku,
- $\hat{y}_{(i,q)}$ – prognozirana granična vrednost q -tog kvantila izlaznog atributa y u i -tom vremenskom koraku,
- τ – redni broj kvantila izražen u opsegu $(0, 1)$,
- L_i i U_i – donje i gornje granice intervala predviđanja u i -tom vremenskom koraku, respektivno,
- φ – širina intervala predviđanja
- ψ – penalizacija odstupanja van intervala predviđanja,
- α – statistički nivo značajnosti intervala predviđanja,
- \hat{F}_i – kumulativna prognozirana raspodela verovatnoće u i -tom vremenskom koraku,
- $\mathbb{1}$ – Hevisajdova funkcija.

Navedene funkcije ocenjuju tačnost prognoze za jedan (i -ti) vremenski korak u horizontu prognoze na osnovu prognoziranih kvantila, intervala predviđanja i raspodele verovatnoće, respektivno. Kod sve tri funkcije rezultati zavise od skale podataka i veći rezultat ukazuje na veću grešku u prognozi. U teoriji, CRPS pruža kompletan uvid u razliku između prognozirane i ostvarene raspodele verovatnoće. Međutim, u praksi, CRPS zahteva diskretizaciju integrala, čime postaje ekvivalentan *pinball* funkciji [34]. Zbog toga se u okviru studije slučaja koja je predstavljena u četvrtoj glavi ove disertacije razmatraju *pinball* i *winkler* funkcije, dok se CRPS neće dalje razmatrati.

Pinball funkcija penalizira svako odstupanje prognoziranih od ostvarenih kvantila u skladu sa rednim brojem kvantila. Pored toga, odstupanja ispod donjih kvantila se penaliziraju više nego odstupanja iznad njih, a odstupanja iznad gornjih kvantila više nego odstupanja ispod njih. *Winkler* funkcija se sastoji od dve komponente (sabarika): širine intervala predviđanja (φ) i penalizacije odstupanja van intervala (ψ). Međutim, *winkler* funkcija ne pruža dovoljno detaljan uvid u nedostatak oštine prognoze kao *pinball* funkcija, jer jednako tretira sva odstupanja unutar intervala predviđanja. S druge strane, komponente *winkler* funkcije se mogu analizirati zasebno, što može pružiti drugačiji uvid u tačnost prognoze.

Ako se rezultati funkcija *pinball* i *winkler* uproseče, omogućen je jednostavan uvid u tačnost probabilističke prognoze [34]. Tako nastaju dve mere tačnosti: *Pinball Score* (PS) i *Winkler Score* (WS). One su definisane na sledeći način:

$$PS = \frac{1}{TQ} \sum_{i=1}^T \sum_{q=1}^Q \text{pinball}(y_i, \hat{y}_{(i,q)}) \quad (2.16)$$

$$WS = \frac{1}{T} \sum_{i=1}^T \text{winkler}(y_i, L_i, U_i) \quad (2.17)$$

gde je:

- T – broj vremenskih koraka,
- Q – broj kvantila,
- y_i – ostvarena vrednost izlaznog atributa y u i -tom vremenskom koraku,
- $\hat{y}_{(i,q)}$ – prognozirana granična vrednost q -tog kvantila izlaznog atributa y u i -tom vremenskom koraku,
- L_i i U_i – donje i gornje granice intervala predviđanja u i -tom vremenskom koraku, respektivno.

Analogno *pinball* i *winkler* funkciji, PS i WS zavise od skale podataka i veći rezultat ukazuje na veću grešku u prognozi. Međutim, nijedna od navedenih mera tačnosti ne pruža uvid u statistički značaj razlike između tačnosti različitih prognoza. Na taj način se ostavlja prostor za slobodnu interpretaciju tačnosti. Statistički test koji formalizuje poređenje tačnosti različitih prognoza je Diebold-Mariano test [169]. To je statistički test koji proverava da li je prosečna razlika između tačnosti dveju prognoza jednaka nuli. S obzirom da između vremenskih koraka u horizontu prognoze tipično postoji visoka korelacija u odnosu na tačnost prognoze, izvođenje testa za sve prognozirane vrednosti odjednom dovodi do pogrešnih zaključaka [34]. Zbog toga je preporučljivo da se test izvodi posebno za svaki vremenski korak. Na primer, ako se dva rešenja poredi u odnosu na tačnost 24-časovnih prognoza za svaki dan u periodu od jedne godine, onda je potrebno izvršiti po jedan test za svaki od 24 sata, sa po 365 test slučajeva. Tako će se tačnost dve prognoze uporediti posebno za svaki sat u horizontu prognoze. Opisani pristup je primenjen u okviru studije slučaja u ovoj disertaciji.

2.6. Probabilistička prognoza

Probabilistička prognoza je generalizacija determinističke prognoze [17]. Deterministička prognoza za jednog potrošača, jednu fizičku veličinu (aktivnu ili reaktivnu snagu) i jedan vremenski korak pruža jednu očekivanu vrednost, dok probabilistička prognoza pruža uvid u neizvesnost očekivane vrednosti. Ako se posmatra iz perspektive teorije sistema, sistem probabilističke prognoze se može konstruisati izmenom delova sistema determinističke prognoze [17, 19]: generisanjem scenarija na ulazu [42–46], analizom grešaka na izlazu [47–55] i izmenom regresionog modela [56–58, 170–172]. Međutim, izmena determinističke prognoze na ulazu zahteva generisanje od nekoliko desetina do nekoliko hiljada scenarija, dok je analiza grešaka ograničena u poboljšanju tačnosti probabilističke prognoze [33]. Prema tome, izmena regresionog modela se prirodno ističe kao glavni pravac istraživanja u radovima iz aktuelnog stanja u oblasti.

Jedno rešenje je da se regresioni model izmeni zamenom funkcije gubitka. Na primer, LSTM model može biti treniran tako da pruži prognozu kvantila ako se u svojstvu funkcije gubitka koristi *pinball* funkcija [57]. Međutim, da bi uvid u raspodelu verovatnoće bio dovoljno detaljan, često se prognoziraju percentili, čime se značajno povećava broj izlaza iz modela, a time i vreme izvršavanja prognoze. Zbog toga je u [58] predloženo da se broj izlaza smanji primenom *Deep Mixture Network* (DMN) modela koji umesto kvantila pruža prognozu parametara mešavine zadatog broja normalnih raspodela. DMN je DNN koji je prikazan na slici 2.13, a koji se sastoji od CNN, GRU, FNN i MDN slojeva. Prethodno opterećenje i meteorološki faktori su upotrebljeni u [58] kao ulazi u DMN model, ali sam model nije ograničen na njihovu upotrebu (npr. vremenske odrednice takođe mogu biti ulaz u model). Izlaz iz DMN modela je buduće opterećenje u vidu mešavine zadatog broja normalnih raspodela. Slična rešenja su predložena u [173, 174].



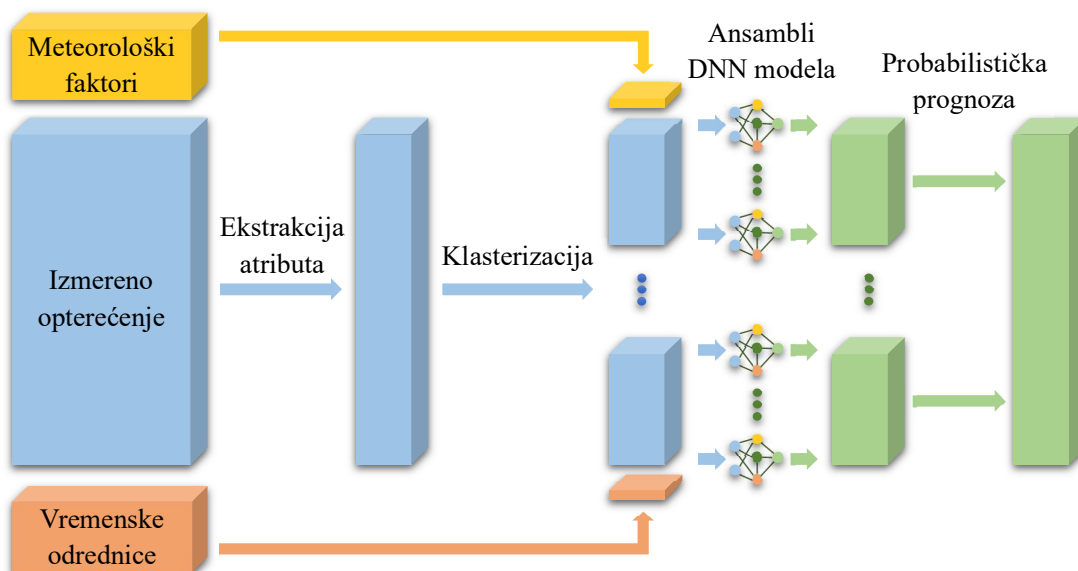
Slika 2.13 – DMN [58]

Novija istraživanja u oblasti prognoze vremenskih serija pokazuju da tačnost prognoze može biti poboljšana kreiranjem regresionih modela nad grupama sličnih vremenskih serija [59–61]. Jedno rešenje koje je generalno primenljivo nad vremenskim serijama i koje se zasniva na otkrivanju zajedničkih šablona vremenskih serija u datom skupu podataka je *Deep Auto-Regression* (DeepAR) [60]. DeepAR je autoregresioni RNN koji pruža probabilističku prognozu u vidu uzoraka normalne raspodele izlaznih vrednosti.

Istraživanja u oblasti prognoze opterećenja potvrđuju da grupisanje potrošača prema sličnosti šablona njihovih HDO i prognoza na nivou potrošačkih grupa vode ka uklanjanju izuzetaka, što ima pozitivan uticaj na tačnost prognoze [152, 175]. *Probabilistic Aggregated Load Forecasting* [176] pruža probabilističku prognozu opterećenja na nivou potrošačkih grupa na osnovu ansambla determinističkih prognoza dobijenih primenom različitih metoda klasterizacije i regresije. S druge strane, *Multitask Based Deep Learning* [177] pruža probabilističku prognozu opterećenja na nivou pojedinačnih potrošača na osnovu DNN modela koji se treniraju za prognozu kvantila na nivou potrošačkih grupa. Pritom se svaki DNN trenira na osnovu svih HDO unutar jedne grupe, a testira na nivou svakog potrošača pojedinačno. Opisano rešenje inherentno vodi ka znatnom povećanju vremena izvršavanja sa povećanjem broja potrošača. Vreme izvršavanja se može smanjiti *Graphic Processing Unit* (GPU) paralelizacijom [170–172]. Međutim, navedena rešenja i dalje mogu biti nepraktična za primenu u velikim DM.

Deep Ensemble Learning (DEL) [62] takođe pruža probabilističku prognozu opterećenja na nivou pojedinačnih potrošača, ali tako što se za svaku potrošačku grupu kreira ansambl DNN modela koji služe za prognozu svih kvantila za sve potrošače unutar grupe. Na primer, ako u grupi ima 10^3 potrošača i ako se prognoziraju percentili, onda svaki model ima 99×10^3 izlaza. Svaki DNN se sastoji od GRU i FNN slojeva koji su podešeni na različite načine unutar jednog ansambla. Izlazi modela se kombinuju na nivou ansambla uz optimizaciju težinskih faktora, čime se dobija konačan rezultat prognoze. Opisano rešenje je ilustrovano na slici 2.14.

Rešenje koje je predloženo u trećoj glavi ove disertacije je u okviru studije slučaja u četvrtoj glavi upoređeno sa DMN, DeepAR i DEL, kao reprezentativnim primerima standardnih rešenja.



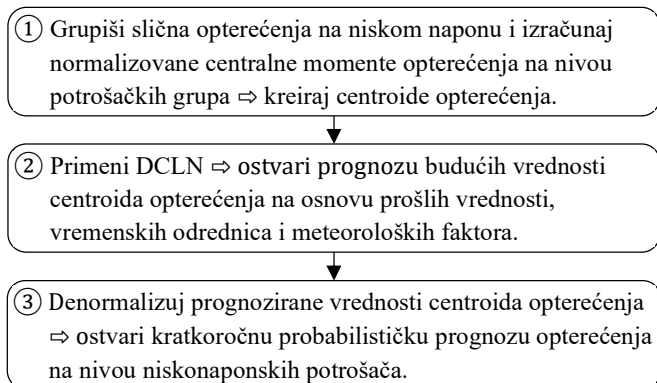
Slika 2.14 – DEL [62]

3. PREDLOG REŠENJA

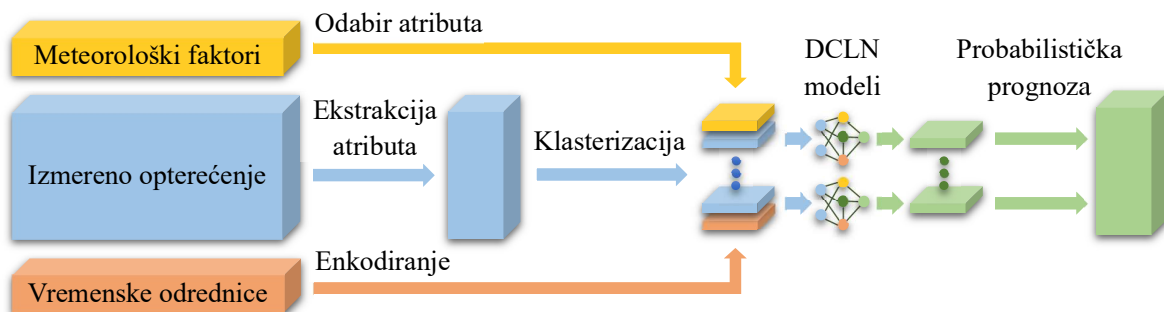
U ovoj disertaciji predloženo je novo rešenje za kratkoročnu probabilističku prognozu opterećenja na niskom naponu u DM – *Deep Centroid Learning* (DCL). Predloženo rešenje se zasniva na prognozi vremenskih serija centralnih momenta opterećenja na nivou potrošačkih grupa (centroida opterećenja) prema postupku datom u [67]. Centralni momenti (prosek i standardna devijacija) se prvo izvode na nivou potrošača, a zatim na nivou potrošačkih grupa. Rezultujući centriodi opterećenja se zatim koriste u regresiji, koja se zasniva na prethodnim vrednostima centroida opterećenja, vremenskim odrednicama i meteorološkim faktorima. Prognozirane vrednosti centroida opterećenja dobijene regresijom primenjuju se na sve potrošače unutar potrošačke grupe, čime se dobija probabilistička prognoza opterećenja na nivou potrošača.

Glavni koraci predloženog rešenja su prikazani na slici 3.1 i ilustrovani na slici 3.2, a u nastavku ove glave su predstavljeni:

- Novi generički model za višerezolucionu aproksimaciju tekućih vremenskih serija u jednom prolazu: *Hierarchical Multiresolution Time Series Representation* (HMTSR) model [72].
- Nova metoda aproksimacije na bazi HMTSR modela koja omogućava očuvanje informacija o fluktuaciji i kontinuitetu HDO: *Hierarchical Multiresolution Linear-function-based Piecewise Statistical Approximation* (HMLPSA) [72].
- Novi algoritam za grupisanje velikog broja aproksimiranih HDO u optimalan broj klastera u prihvatljivom vremenskom roku: *Time Series Grouping Algorithm* (TSGA) [30].
- Nova LPR metoda [67], koja se zasniva na primeni HMLPSA i TSGA sa ciljem redukcije HDO i kreiranja centroida opterećenja, uz očuvanje informacija koje su neophodne da bi se podržala kratkoročna probabilistička prognoza opterećenja.
- Novi regresioni model za prognozu budućih vrednosti centroida opterećenja, koji se zasniva na primeni dubokog učenja: *DCL Network* (DCLN).



Slika 3.1 – Glavni koraci predloženog DCL rešenja [67]



Slika 3.2 – Ilustracija glavnih koraka predloženog DCL rešenja

3.1. HMTSR

HMTSR model [72] je generički model za višerezolucionu aproksimaciju tekućih vremenskih serija u jednom prolazu. Analogno MTSMS modelu [90], koji je opisan u teoretskim osnovama ove disertacije, aproksimirane vrednosti se dobijaju iterativno, primenom novog algoritma koji održava dve generičke strukture podataka za svaku vremensku seriju i svaku vremensku rezoluciju: bafer i disk. Međutim, za razliku od MTSMS modela, HMTSR model ne zahteva da se podaci skladište u bafer u jednom prolazu, a zatim aproksimiraju u drugom prolazu, već se podaci aproksimiraju prilikom izmene bafera. Formalno, HMTSR model \mathcal{H} je definisan na sledeći način:

$$\mathcal{H}(S, \mathcal{R}, f, g) = \bigcup_{\forall \Delta \in \mathcal{R} \cup \{\delta\}} \{\text{disk}_\Delta, \text{bafer}_\Delta\} \bigcup \text{bafer}_S, \quad (3.1)$$

gde je:

S – tekuća vremenska serija u domenu vremena \mathcal{T} i domenu vrednosti \mathcal{V} :

$$S = \{(t_0, v_0), \dots, (t_n, v_n)\} \subset \mathcal{T} \times \mathcal{V} \quad (3.2)$$

$\mathcal{R} \subset \mathcal{T}$ – skup dužina vremenskih koraka za različite vremenske rezolucije,

$\delta \in \mathcal{T}$ – najveći zajednički delilac dužina vremenskih koraka iz skupa \mathcal{R} ,

$f: \mathcal{T} \rightarrow \mathcal{V}$ – reprezentaciona funkcija koja spaja tačke serije S ,

g – kompozitna agregaciona funkcija sa operatorom spajanja \odot , za koju važi:

$$(\forall t \in [t_i, t_j]) g(f|_{[t_i, t_i]}) \odot g(f|_{[t_i, t]}) = g(f|_{[t_i, t_j]}), \quad (3.3)$$

disk_Δ – disk koji predstavlja seriju S na vremenskim koracima dužine Δ :

$$\text{disk}_\Delta = \left\{ g \left(f|_{\left[\lfloor \frac{t_0}{\Delta} \rfloor \Delta, (\lfloor \frac{t_0}{\Delta} \rfloor + 1) \Delta \right]} \right), \dots, g \left(f|_{\left[(\lfloor \frac{t_n}{\Delta} \rfloor - 1) \Delta, \lfloor \frac{t_n}{\Delta} \rfloor \Delta \right]} \right) \right\}, \quad (3.4)$$

bafer_Δ – bafer koji predstavlja seriju S na vremenskim koracima dužine Δ :

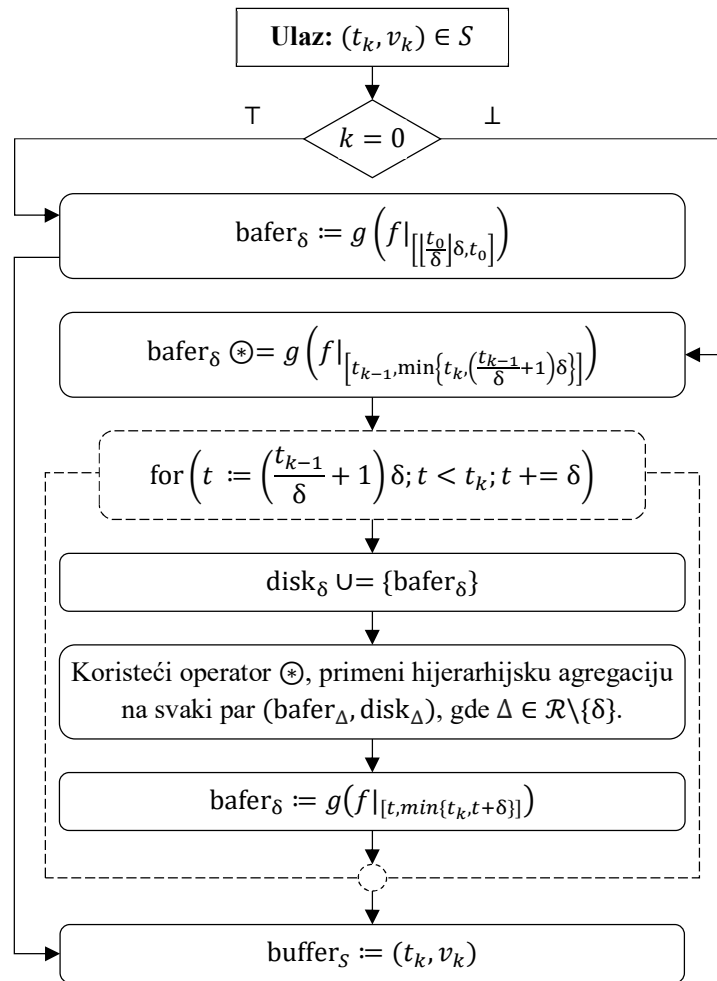
$$\text{bafer}_\Delta = g \left(f|_{\left[\lfloor \frac{t_n}{\Delta} \rfloor \Delta, (\lfloor \frac{t_n}{\Delta} \rfloor + 1) \Delta \right]} \right), \quad (3.5)$$

bafer_S – bafer koji sadrži poslednju tačku serije S :

$$\text{bafer}_S = (t_n, v_n). \quad (3.6)$$

Svaki par $(\text{disk}_\Delta, \text{bafer}_\Delta)$ čini aproksimaciju serije S na jednoj vremenskoj rezoluciji, sa vremenskim koracima dužine Δ , pri čemu važi $|\text{bafer}_\Delta| \ll |\text{disk}_\Delta|$. S druge strane, bafer_S sadrži samo poslednju tačku (vreme i vrednost) serije S za potrebe aproksimacije podataka u jednom prolazu. Funkcija f omogućava da kontinuitet procesa koji je modelovan serijom S bude opisan na odgovarajući način. S druge strane, funkcija g dozvoljavaju da predloženi model bude prilagođen odgovarajućim potrebama aproksimacije serije S . Međutim, da bi aproksimacija mogla da se izvrši u jednom prolazu, funkcija g mora podržati operaciju spajanja aproksimiranih vrednosti (npr. rezultat funkcije g može biti prosek, ali ne i medijana).

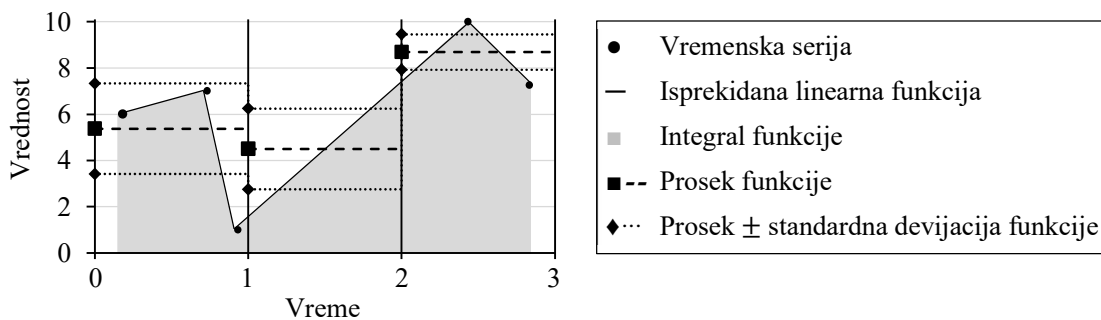
Algoritam koji je dizajniran za kreiranje HMTSR modela je prikazan na slici 3.3. Može se primetiti da se višerezoluciona aproksimacija serije S kreira u jednom prolazu, inkrementalnom izmenom bafera, kao i da aproksimirane vrednosti na kraju završavaju na diskovima. Aproksimacija serije S se kreira na najvišoj vremenskoj rezoluciji, sa vremenskim koracima dužine δ , kako podaci pristižu. Aproksimacija serije S na nižim rezolucijama, sa vremenskim koracima dužine Δ , dobija se hijerarhijskom agregacijom, koja je tipična za analitičku obradu podataka [93]. Zbog hijerarhijske agregacije, kompleksnost predstavljenog algoritma je u najgorem slučaju $O(|\mathcal{R}|)$. Međutim, amortizovana kompleksnost algoritma je $O(1)$, jer se aproksimacija najčešće vrši na najvišoj vremenskoj rezoluciji. Pored toga, predloženi model i algoritam omogućavaju paralelizaciju aproksimacije različitih vremenskih serija i na taj način postizanje visokih performansi obrade podataka sa većim računarskim resursima.



Slika 3.3 – Algoritam za kreiranje predloženog HMTSR modela [72]

3.2. HMLPSA

HMLPSA je modifikovana verzija PSA metode koja se zasniva na HMTSR modelu. Analogno PSA metodi, rezultat HMLPSA metode nad segmentom vremenske serije su prva dva centralna momenta, prosek (μ) i standardna devijacija (σ), vrednosti unutar segmenta. Međutim, za razliku od PSA metode, HMLPSA metoda zadržava informacije o kontinuitetu i varijabilnosti podataka na nižim vremenskim rezolucijama, posmatrajući vremensku seriju kroz isprekidanu linearnu funkciju koja povezuje tačke vremenske serije. Primer aproksimacije jedne vremenske serije na jednoj vremenskoj rezoluciji primenom opisanog rešenja je prikazan na slici 3.4.



Slika 3.4 – Primer aproksimacije jedne vremenske serije primenom predložene HMLPSA metode [72]

Formalno, HMLPSA je $\mathcal{H}(S, \mathcal{R}, f, g)$, gde su funkcije f i g definisane na sledeći način:

$$(\forall k \in [1..|S|]) f(t) = \Phi_k t + \nabla_k, \quad \Phi_k = \frac{v_k - v_{k-1}}{t_k - t_{k-1}}, \quad \nabla_k = v_k - \Phi_k t_k, \quad (3.7)$$

$$(\forall [t_i, t_j] \subset \mathcal{T}) g(f|_{[t_i, t_j]}) = (\mu, \sigma), \mu = \frac{1}{t_j - t_i} \int_{t_i}^{t_j} f(t) dt, \sigma = \sqrt{\frac{1}{t_j - t_i} \int_{t_i}^{t_j} (t - \mu)^2 f(t) dt} \quad (3.8)$$

gde je Φ_k nagib, a ∇_k odsečak prave koja povezuje k -tu tačku serije S sa prethodnom tačkom.

Funkcija f linearno povezuje tačke serije S i na taj način omogućava da trajanje vrednosti utiče na rezultat, čime se smanjuje uticaj kratkotrajnih izuzetaka i povećava uticaj očekivanih vrednosti. S druge strane, funkcija g izvodi centralne momente nad funkcijom f , čime se opisuju kontinuitet i varijabilnost modelovanog procesa. Numeričke metode koje se koriste za računanje integrala u okviru funkcije g ne menjaju složenost algoritma sa slike 3.3. Svaki integral se inkrementalno računa kao suma integrala između susednih tačaka serije S . S obzirom da HMTSR model čuva poslednju procesiranu tačku serije S , parametri i integral funkcije f se mogu odrediti svaki put kada se nova tačka prikupi. Rezultati funkcije g nad različitim segmentima serije S se zatim mogu spojiti primenom operatora \odot , koji je definisan na sledeći način:

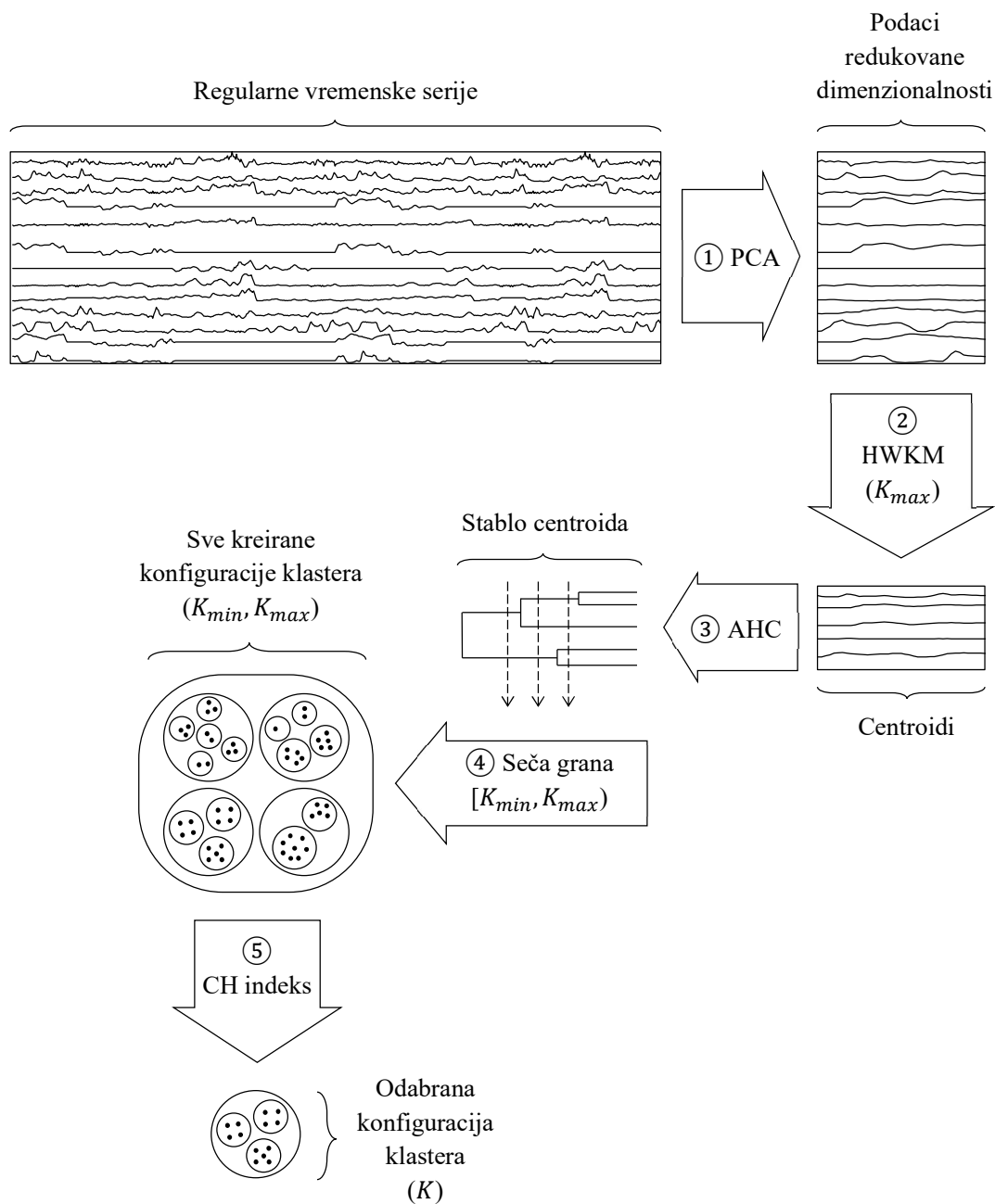
$$(\forall t \in [t_i, t_j]) g(f|_{[t_i, t]}) = (\mu_i, \sigma_i) \wedge g(f|_{[t, t_j]}) = (\mu_j, \sigma_j) \Rightarrow$$

$$g(f|_{[t_i, t]}) \odot g(f|_{[t, t_j]}) = \left(\frac{(t-t_i)\mu_i + (t_j-t)\mu_j}{t_j-t_i}, \sqrt{\frac{(t-t_i)(\sigma_i^2 + (\mu_i - \mu)^2) + (t_j-t)(\sigma_j^2 + (\mu_j - \mu)^2)}{t_j-t_i}} \right) \quad (3.9)$$

3.3. TSGA

TSGA [30] je algoritam za grupisanje velikog broja visokodimenzionalnih regularnih vremenskih serija u optimalan broj klastera u prihvatljivom vremenskom okviru na osnovu zadatog opsega broja klastera (K_{min}, K_{max}). Za klasterizaciju se koristi Euklidsko rastojanje, koje se pokazalo izuzetno efikasnim nad takvim podacima [70]. Na osnovu [126], potraga za optimalnim brojem klastera se izvodi u opsegu $(2, \lfloor \sqrt{N} \rfloor)$, gde je N broj objekata za klasterizaciju (broj vremenskih serija). TSGA je ilustrovan na slici 3.5, a sastoji se od sledećih koraka:

- 1) PCA se koristi za ekstrakciju glavnih komponenti vremenskih serija, koje kumulativno objašnjavaju 95% varijanse u podacima. PCA omogućava efikasnu klasterizaciju HDO na osnovu Euklidskog rastojanja [110]. Eksperimentalni rezultati primene PCA nad podacima koji su upotrebljeni u studiji slučaja koja je predstavljena u četvrtoj glavi ove disertacije, pokazuju da ovaj korak vodi ka značajnoj redukciji dimenzionalnosti (u proseku ~85%, od ~73% do ~97%).
- 2) HWKM [130] se koristi za formiranje maksimalnog broja klastera (K_{max}) i njihovih centroida. HWKM smanjuje osetljivost klasičnog KM algoritma na početan odabir centroida pomoću naknadne rekonfiguracije klastera.
- 3) Analogno KnA metodi [136], nad dobijenim centroidima se primenjuje AHC [70]. Za kreiranje stabla centroida se koristi UPGMA kriterijum spajanja [127] - spajanje onih skupova objekata (centroida) čije je prosečno rastojanje između svih parova objekata najmanje na posmatranom nivou hijerarhije. Autori u [127] su pokazali da je UPGMA najbolji izbor za klasterizaciju HDO.
- 4) Sečom grana [134] izgrađenog stabla centroida formiraju se sve ostale konfiguracije klastera u datom opsegu broja klastera. Algoritamska složenost seče grana izgrađenog stabla centroida je znatno manja od ponovne klasterizacije, što čini potragu za optimalnim brojem klastera računarski efikasnom.
- 5) Konačna konfiguracija, sa optimalnim brojem klastera (K), bira se na osnovu indeksa validnosti. Zbog prednosti koje su navedene u teoretskim osnovama ove disertacije, za validaciju se koristi CH indeks [124].



Slika 3.5 – Ilustracija koraka predloženog TSGA rešenja [30]

3.4. LPR

Predložena LPR metoda [67] (Algoritam 1) služi za redukciju podataka o opterećenju na niskom naponu uz očuvanje informacija koje su neophodne za kratkoročnu probabilističku prognozu opterećenja. Predložena metoda se zasniva na grupisanju potrošača sa sličnim šablonima opterećenja i na izvođenju vremenskih serija centralnih momenta opterećenja na nivou potrošačkih grupa. Centralni momenti (μ i σ) se prvo izvedu na nivou potrošača, a zatim na nivou potrošačkih grupa, koristeći HMLPSA [72] i TSGA [30]. Rezultat su centroidi opterećenja (očekivane vrednosti sa pridruženom standardnom devijacijom na nivou potrošačkih grupa). Na taj način se ujedno redukuje količina podataka i zadržavaju se informacije o fluktuacijama vrednosti sadržanih u HDO, što je neophodno za kreiranje regresionih modela (DCLN).

Ulazi:

- Skup HDO za N potrošača i M fizičkih veličina: $\mathcal{S} = \{S_{(1,1)}, \dots, S_{(N,M)}\}$, gde $S_{(n,m)} \in \mathcal{S}$ predstavlja HDO za n -tog potrošača i m -tu fizičku veličinu.
- Vremenska rezolucija prognoze opterećenja R (skup vremenskih koraka dužine Δ)
- Skupovi karakterističnih tipova dana \mathcal{D} , perioda \mathcal{P} i tipova potrošača \mathcal{L} .

Izlazi:

- Skup normalizacionih faktora \mathcal{F} .
- Skup potrošačkih grupa \mathcal{G} .
- Skup centroida opterećenja \mathcal{C} .

Koraci:

1. Kreiraj skup AHDO: $\mathcal{A} = \{A_{(1,1)}, \dots, A_{(N,M)}\}$, gde je $A_{(n,m)} \in \mathcal{A}$ rezultat primene HMLPSA metode na seriju $S_{(n,m)} \in \mathcal{S}$, odnosno serija μ i σ vrednosti za n -tog potrošača i m -tu fizičku veličinu na rezoluciji R (ukupno $N \times M \times 2 \times |R|$ vrednosti).
2. Kreiraj skup normalizacionih faktora: $\mathcal{F} = \{F_{(1,1)}, \dots, F_{(N,M)}\}$, gde je $F_{(n,m)} \in \mathcal{F}$ prosečno opterećenje izračunato na osnovu vrednosti sadržanih u $A_{(n,m)} \in \mathcal{A}$.
3. Kreiraj skup NAHDO: $\mathring{\mathcal{A}} = \{\mathring{A}_{(1,1)}, \dots, \mathring{A}_{(N,M)}\}$, gde je $\mathring{A}_{(n,m)} \in \mathring{\mathcal{A}}$ izračunato deljenjem vrednosti sadržanih u $A_{(n,m)} \in \mathcal{A}$ sa $F_{(n,m)} \in \mathcal{F}$.
4. Kreiraj skup PNAHDO: $\mathring{\mathring{\mathcal{A}}} = \{\mathring{\mathring{A}}_{(1,1)}, \dots, \mathring{\mathring{A}}_{(N,M)}\}$, gde je $\mathring{\mathring{A}}_{(n,m)} \in \mathring{\mathring{\mathcal{A}}}$ prosek vrednosti sadržanih u $\mathring{A}_{(n,m)} \in \mathring{\mathcal{A}}$ tokom karakterističnih perioda i tipova dana (d, p) , $d \in \mathcal{D}$, $p \in \mathcal{P}$ (ukupno $N \times M \times 2 \times D \times |\mathcal{D}| \times |\mathcal{P}|$ vrednosti, gde je D broj vremenskih koraka iz R u jednom danu).
5. Kreiraj potrošačke grupe $\mathcal{G} = \{G_1, \dots, G_K\}$ primenom TSGA na KPNAHDO $\mathring{\mathring{\mathring{\mathcal{A}}}} = \{\mathring{\mathring{\mathring{A}}}_1, \dots, \mathring{\mathring{\mathring{A}}}_N\}$ razvrstane u odnosu na \mathcal{L} , gde važi $(\forall \mathring{\mathring{A}}_n \in \mathring{\mathring{\mathcal{A}}}) \mathring{\mathring{A}}_n = \cup_m^M \mathring{\mathring{A}}_{(n,m)} \subset \mathring{\mathring{\mathcal{A}}} \wedge (\forall G_k \in \mathcal{G}) G_k \subset [1, N] \wedge K \in [2, \lfloor \sqrt{N} \rfloor]$.
6. Kreiraj skup centroida opterećenja: $\mathcal{C} = \{C_{(1,1)}, \dots, C_{(K,M)}\}$, gde je $C_{(k,m)} \in \mathcal{C}$ prosek vrednosti sadržanih u $\{\mathring{\mathring{A}}_{(n,m)}, \dots, \mathring{\mathring{A}}_{(n',m)}\} \subset \mathring{\mathring{\mathcal{A}}}$, $\{n, \dots, n'\} \in G_k$, $G_k \in \mathcal{G}$, odnosno serija normalizovanih μ i σ vrednosti za k -tu grupu i m -tu fizičku veličinu na rezoluciji R (ukupno $K \times M \times 2 \times |R|$ vrednosti).

Prvi korak je aproksimacija HDO na vremenskoj rezoluciji prognoze primenom HMLPSA metode [72]. Rezultat su aproksimirani HDO (AHDO) – vremenske serije centralnih momenata opterećenja (μ i σ) na nivou potrošača. Količina podataka o opterećenju se značajno redukuje u prvom koraku. Na primer, ako se svake minute aktivna i reaktivna snaga milion potrošača beleže sa 4-bajtnom preciznošću, onda se za godinu dana prikupi ~3,8 TiB podataka [30]. Međutim, ako se izmerene vrednosti aproksimiraju centralnim momentima na nivou sata, onda se količina podataka smanjuje sa ~3,8 TiB na ~131 GiB (za ~97%). Pritom, HMLPSA ne zahteva ~3,8 TiB skladišnog prostora već samo ~131 GiB, jer prema HMTSR modelu omogućava da se aproksimacija izvodi u jednom prolazu i može se primeniti u realnom vremenu.

Drugi i treći korak služe za normalizaciju AHDO. Na taj način se sprečava da AHDO sličnog oblika budu svrstani u različite klastere zbog razlike u njihovim apsolutnim vrednostima. Deljenjem vrednosti sadržanih u AHDO sa prosečnim opterećenjem na nivou potrošača dobijaju se normalizovani AHDO (NAHDO). Opisani postupak je praktičniji od drugih metoda normalizacije u prisustvu izuzetaka [72], što je u slučaju HDO očekivano. Normalizacioni faktori (proseci) koji su dobijeni u drugom koraku se takođe koriste za denormalizaciju prognoziranih vrednosti centroida opterećenja nakon primene DCLN modela.

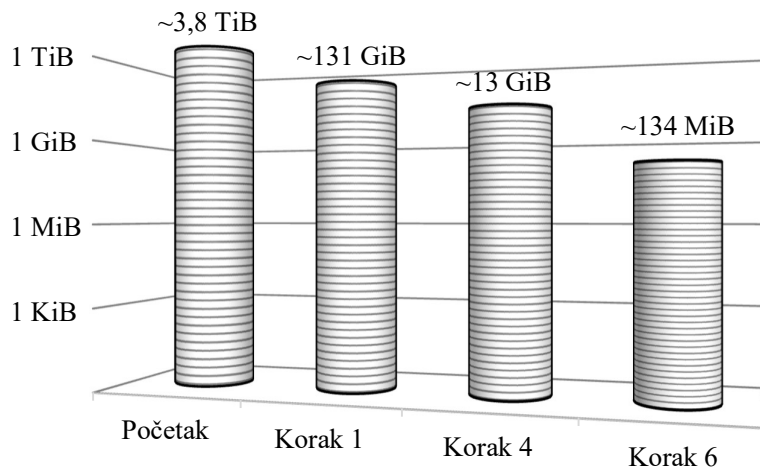
Četvrti korak služi da se podaci redukuju tako da se bez neopravdanog zauzeća računarskih resursa omogući praktična primena TSGA nad NAHDO na niskom naponu. Podaci se redukuju izvođenjem proseka NAHDO za unapred definisane karakteristične periode i tipove dana. Na taj način se dobijaju prosečni

NAHDO (PNAHDO). Model podataka o opterećenju koji se zasniva na karakterističnim periodima i tipovima dana vodi ka značajnoj redukciji količine podataka koji su potrebni za reprezentaciju opterećenja DM, ali i ka povećanju brzine obrade tih podataka bez značajne degradacije kvaliteta rezultata DMS funkcija [178]. Na primer, ako se prosek 365×10^6 NAHDO izvede za 12 perioda (meseći) i 3 tipa dana (radni dan, vikend i praznik), onda se količina podataka smanjuje sa ~ 131 GiB na ~ 13 GiB (za $\sim 90\%$).

U petom koraku se formira skup podataka za klasterizaciju, gde su potrošači predstavljani kao objekti, a redom spojene vrednosti PNAHDO kao njihovi atributi. Tako nastaju redom spojeni (konkatenirani) PNAHDO (KPNAHDO). Broj atributa KPNAHDO predstavlja dimenzionalnost podataka za klasterizaciju. Na primer, ako KPNAHDO za svakog potrošača sadrži 24 centralna momenta (μ i σ) aktivne i reaktivne snage za 12 karakterističnih perioda i 3 karakteristična tipa dana, onda rezultujući skup podataka ima $24 \times 2 \times 2 \times 12 \times 3 = 3.456$ dimenzija. S obzirom da potencijalno veliki broj objekata (npr. milion potrošača) i visoku dimenzionalnost podataka (npr. 3.456), za klasterizaciju se koristi TSGA [30], zbog prednosti koje su navedene u delu 3.3. Pritom se TSGA primenjuje posebno za svaki tip potrošača, kao što su domaćinstva i komercijalni potrošači. Razvrstavanje potrošača po tipu smanjuje algoritamsku složenost klasterizacije i sprečava da tačnost prognoze na nivou potrošačke grupe bude narušena svrstavanjem potrošača različitog tipa u istu grupu.

Šesti korak je izračunavanje proseka vrednosti NAHDO za potrošačke grupe. Rezultat ovog koraka su centriodi opterećenja – vremenske serije normalizovanih centralnih momenata opterećenja (μ i σ) na nivou potrošačkih grupa. Pronalaženje centroida opterećenja ima mnoge prednosti: 1) uklanjaju se izuzeci u opterećenju i otkriva se zajednički šablon opterećenja na nivou grupe; 2) zadržavaju se informacije o očekivanom opterećenju i varijabilnosti opterećenja; 3) zadržava se dovoljan primer za duboko učenje, i 4) značajno se redukuje količina podataka. Na primer, ako se 10^6 potrošača razvrsta u 10^3 potrošačkih grupa i ako se njihovi NAHDO sa satnim vrednostima za godinu dana predstave centroidima opterećenja, onda se količina podataka smanjuje sa ~ 131 GiB na ~ 134 MiB (za $\sim 99.9\%$). Primer koji je upotrebljen za teoretsku kvantifikaciju redukcije količine podataka u prvom, četvrtom i šestom koraku predložene LPR metode ilustrovan je na slici 3.6.

S obzirom da HMLPSA omogućava aproksimaciju HDO u jednom prolazu, Algoritam 1 se može ponavljati bez koraka 2, 4 i 5 u cilju inkrementalnog ažuriranja centroida opterećenja sa pojavom novih vrednosti HDO. To omogućava da se parametri (težine) respektivnih DCLN modela takođe inkrementalno ažuriraju. Prema tome, opisani postupak omogućava inicijalno treniranje i inkrementalno ažuriranje DCLN modela bez prekomerne upotrebe računarskih resursa (procesorske moći i skladišnog prostora). Zauzvrat, DCLN omogućava otkrivanje složenih nelinearnih povezanosti između visoko varijabilnog opterećenja i varijabilnih faktora koji na to opterećenje utiču. Posledično, predloženo rešenje vodi ka poboljšanju tačnosti prognoze i vremena izvršavanja u poređenju sa konkurentnim rešenjima iz aktuelnog stanja u oblasti, što pokazuje i studija slučaja koja je predstavljena u četvrtoj glavi ove disertacije.



Slika 3.6 – Primer redukcije količine podataka o opterećenju primenom predložene LPR metode [67]

3.5. DCLN

DCLN je DNN koji je dizajniran tako da omogući primenu MIMO pristupa u prognozi centroida opterećenja. Primena MIMO pristupa u prognozi opterećenja nudi uglacan prelaz između prognoziranih vrednosti (za razliku od direktnog pristupa) bez akumulacije greške na kasnijim vremenskim koracima u horizontu prognoze (kao kod rekurzivnog pristupa). Jedan DCLN se trenira za prognozu centroida opterećenja za jednu potrošačku grupu (trenira se onoliko modela koliko ima potrošačkih grupa). Svaki DCLN se trenira tako da prognozira buduće vrednosti centroida opterećenja na osnovu prethodnih vrednosti sa odabranim zaostajanjem i na osnovu enkodiranih vremenskih odrednica i odabranih meteoroloških faktora. U nastavku ovog dela su prvo opisani ulazni i izlazni podaci koji se koriste u regresiji, a zatim i struktura samog modela.

3.5.1. Ulazni i izlazni podaci

DCLN se koristi za prognozu centroida opterećenja na osnovu dva tipa ulaznih podataka koji se često koriste u prognozi opterećenja [115]: endogeni (vrednosti centroida opterećenja sa određenim zaostajanjem pre horizonta prognoze) i egzogeni (vremenske odrednice i meteorološki faktori u horizontu prognoze).

Vreme je jedan od najznačajnijih faktora opterećenja. Opterećenje se menja u toku dana, u korelaciji je sa opterećenjem u toku prethodnih dana i ponavlja se na sličan način u toku nedelje, meseca i godine [115]. Prema tome, zaostajanje za centroide opterećenja u danima koji prethode horizontu prognoze bira se primenom PACF [117]. S druge strane, ponavljanje šablona opterećenja tokom dana, nedelje i godine, kao i tokom praznika, je modelovano upotrebom vremenskih odrednica na ulazu u regresioni model. Upotreba vremenskih odrednica u regresiji zahteva da se one enkodiraju u numerički domen. Periodične vremenske odrednice (vreme u toku dana, dan u nedelji i mesec u godini) se enkodiraju u 2D koordinate Dekartovog koordinatnog sistema na jediničnom krugu da bi se zadržale informacije i njihovom redosledu i međusobnoj udaljenosti [115]. Formalno, vrednost λ vremenske odrednice sa periodom ponavljanja Λ se enkodira u par vrednosti $(\sin(\theta), \cos(\theta))$, gde je $\theta = 2\pi\lambda/\Lambda$. Na primer, ako se meseci u godini enkodiraju na opisan način (kao kazaljke na satu), onda januar uvek sledi nakon decembra, a pre februara, i jednako je udaljen od oba. S druge strane, za enkodiranje aperiodičnih vremenskih odrednica (praznika) se koristi binarni indikator (enkodirana vrednost je 1 ako je dan praznični, a 0 ako nije).

U praksi, prognoza opterećenja zavisi od istorije i prognoze meteoroloških uslova: istorija se koristi za treniranje regresionog modela, a prognoza za testiranje modela. Međutim, verifikacija prognoze zahteva postojanje istorijskih vrednosti u oba slučaja. Zbog toga je u radovima koji se bave prognozom opterećenja uobičajeno da se verifikacija izvodi tako što se istorijske vrednosti dele na trening i test skup. Na taj način postupak verifikacije ne zavisi od prognoze meteoroloških uslova, ali se postojanje takve prognoze podrazumeva u praktičnoj primeni. Takođe se podrazumeva da su meteorološki podaci dostupni na vremenskoj rezoluciji prognoze opterećenja ili da se mogu aproksimirati na toj rezoluciji primenom metoda agregacije ili interpolacije. S obzirom da specifičnosti sistema za pružanje meteoroloških podataka prevazilaze okvire ove disertacije, dostupnost i vremenska rezolucija meteoroloških podataka se neće dalje razmatrati.

Aktivnosti potrošača takođe u velikoj meri zavise od meteoroloških faktora, kao što je temperatura, ali ne moraju svi dostupni meteorološki podaci biti od značaja za prognozu. Visoka tačnost prognoze se tipično postiže uz adekvatan odabir ulaznih atributa [111]. Zbog toga se za odabir meteoroloških faktora koristi JMIM metoda [111]. JMIM metoda za svaki ulazni atribut vraća realan broj u opsegu $[0, 1]$ (skor) koji kvantifikuje količinu zajedničkih informacija koju taj atribut deli sa izlaznim atributom. Rezultujućim skorovi se koriste da se odaberu meteorološki faktori koji kumulativno dele 95% zajedničkih informacija sa centroidima opterećenja (analogno primeni PCA metode u okviru TSGA).

Rezultat primene DCLN modela su prognozirane vrednosti centroida opterećenja (normalizovane μ i σ vrednosti na nivou potrošačkih grupa). Množenjem tih vrednosti sa normalizacionim faktorima koji su dobijeni primenom LPR metode iz dela 3.4 dobija se kratkoročna probabilistička prognoza opterećenja na niskom naponu (denormalizovane μ i σ vrednosti na nivou potrošača). Za potrebe prikaza i verifikacije, rezultujuća prognoza može biti predstavljena u vidu kvantila i intervala predviđanja, koji se mogu izračunati

na osnovu normalne raspodele koja je definisana prognoziranim centralnim momentima (što je primenjeno u okviru studije slučaja koja je predstavljena u četvrtoj glavi ove disertacije).

3.5.2. Struktura modela

Struktura DCLN modela je predstavljena na slici 3.7. Oblik višedimenzionalnog niza (eng. *tensor*) na izlazu iz svakog sloja je prikazan u uglastim zagradama. DCLN se sastoji od DCNN, S2S LSTM i SAE slojeva. S2S LSTM je osnovna podstruktura DCLN modela, dok su DCNN i SAE sporedne podstrukture. Navedene podstrukture se koriste na sledeći način:

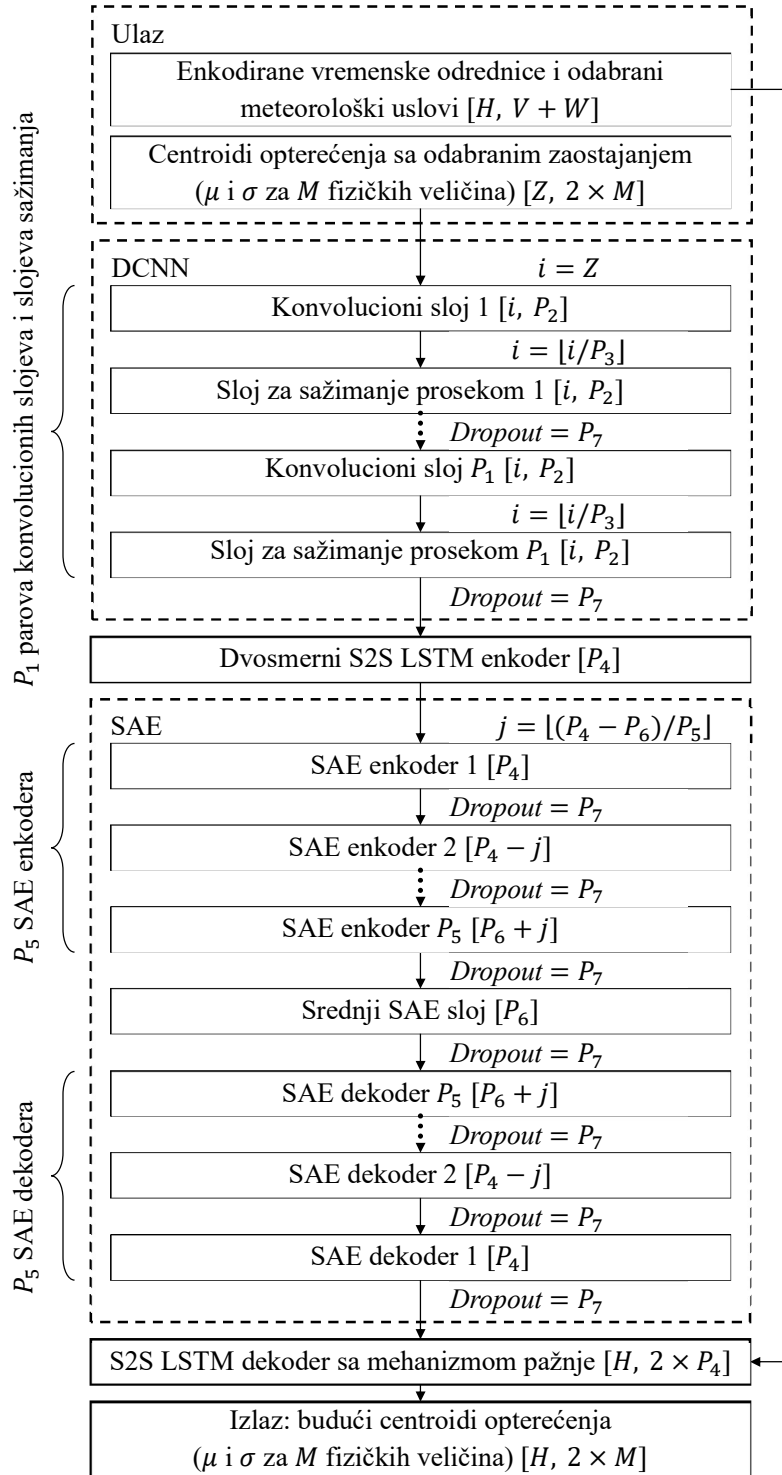
- 1) DCNN [150–152] se koristi za prepoznavanje lokalnih trendova u centroidima opterećenja na ulazu i redukciju dimenzionalnosti ulaznih podataka bez gubitka značajnih informacija sa ciljem kontrole osetljivosti predloženog modela na izuzetke. Dok je upotreba vremenskih odrednica i meteoroloških faktora u regresiji ograničena na horizont prognoze, upotreba prethodnog opterećenja je ograničena na odabrano zaostajanje. Na primer, u studiji slučaja koja je predstavljena u četvrtoj glavi ove disertacije za prognozu se koristiti opterećenje u prethodnih sedam dana. S obzirom na visoku varijabilnost opterećenja u potencijalno dugom vremenskom periodu, neophodno je redukovati dimenzionalnost ulaznih podataka i ukloniti izuzetke, što se postiže upotrebom naizmenično postavljenih konvolucionih slojeva i slojeva za sažimanje prosekom.
- 2) S2S LSTM [145] se koristi za otkrivanje vremenske zavisnosti između ulaznih i izlaznih vrednosti i za implementaciju MIMO pristupa u prognozi. Dvosmerni LSTM enkoder [147] služi za otkrivanje vremenske zavisnosti između vrednosti centroida opterećenja pre horizonta prognoze. S druge strane, (jednosmerni) LSTM dekoder sa mehanizmom pažnje [149] služi za prognozu centroida opterećenja u horizontu prognoze. Mehanizam pažnje služi da se centriodi opterećenja preko S2S arhitekture povežu sa enkodiranim vremenskim odrednicama i odabranim meteorološkim faktorima u horizontu prognoze.
- 3) SAE [146] se koristi za uklanjanje izuzetaka iz S2S konteksta i povećanje kapaciteta reprezentacije centroida opterećenja. SAE takođe omogućava da DCNN propusti više informacija do S2S LSTM slojeva i na taj način smanji rizik od gubitka značajnih informacija, ali i rizik od povećanja osetljivosti modela na izuzetke. Prema tome, SAE slojevi omogućavaju postizanje ravnoteže između uklanjanja izuzetaka i gubitka informacija u DCNN slojevima. Uloga SAE je da rekonstruiše S2S kontekst bez promene dimenzionalnosti podataka i zbog toga je simetričnog oblika: broj neurona na ulazu jednak je broju neurona na izlazu.

Parametri (težine) DCLN modela se podešavaju bekpropagacijom u slučaju FNN slojeva (DCNN i SAE) i bekpropagacijom kroz vreme u slučaju RNN slojeva (LSTM). Težine se tokom treniranja optimizuju primenom Adam algoritma, koji se može pronaći kao prvi izbor u optimizaciji parametara DNN modela u mnogim radovima [37]. S obzirom da se DCLN koristi za determinističku prognozu vrednosti centroida opterećenja (koje se kasnije transformišu u probablistički oblik na nivou potrošača), u optimizaciji se kao funkcija gubitka koristi RMSE. Za prevenciju prekomernog učenja se koriste *dropout* i *early stopping*. Pritom se *dropout* koristi za FNN, ali ne i za RNN slojeve. Eksperimentalno je ustanovljeno da *dropout* i drugi poznati tipovi regularizacije, kao što su penalizacija funkcije gubitka i dodavanje šuma, ne doprinose povećanju tačnosti prognoze kada se primene na RNN slojeve u okviru predloženog modela.

Struktura DCLN modela je potpuno određena sledećim hiperparametrima:

- 1) Dimenzije ulaznih podataka su određene brojem fizičkih veličina (M), zaostajanjem koje je odabrano za centroide opterećenja (Z), brojem atributa koji predstavljaju enkodirane vremenske odrednice (V), brojem odabranih meteoroloških faktora (W) i dužinom horizonta prognoze (H), izraženom u broju vremenskih koraka na vremenskoj rezoluciji prognoze.
- 2) DCNN je određen brojem naizmenično postavljenih konvolucionih slojeva i slojeva sažimanja (P_1), brojem konvolucionih filtera (P_2) i veličinom filtera i koraka sažimanja (P_3) [150, 152].

- 3) S2S LSTM je određen dužinom dvosmernog S2S konteksta – brojem skrivenih stanja i stanja ćelije u oba smera (P_4) [147].
- 4) SAE je određen brojem enkoder i dekoder slojeva (P_5) i brojem neurona u srednjem (eng. *code*) sloju (P_6) [146]. Zbog simetričnosti, broj neurona u ulaznim i izlaznim slojevima je P_4 .
- 5) *Dropout* regularizacija je određena procentom neurona za odbacivanje tokom treniranja (P_7).



Slika 3.7 – Struktura predloženog DCLN modela

Primitno je da bi navedeni hiperparametri takođe mogli da se optimizuju nekom od metoda koje su opisane u delu 2.5. Međutim, eksperimentalno je ustanovljeno da takva optimizacija vodi ka značajnom povećanju vremena izvršavanja bez značajnog povećanja tačnosti prognoze. Zbog toga se DCLN model oslanja na ekspertsko znanje u podešavanju hiperparametara i na Adam algoritam u podešavanju parametara modela za date hiperparametre.

Hiperparametri DCLN modela se mogu podešavati prema uobičajenim praksama dubokog učenja. Tokom podešavanja je uobičajeno da se samo 70-80% primera iz trening skupa koristi za treniranje modela, a da se preostalih 20-30% primera koristi za testiranje i verifikaciju modela. Analizom grešaka na svakom od ta dva skupa može se ustanoviti greška generalizacije, odnosno da li i u kojoj meri pristrasnost i varijansa vode ka nedovoljnom i prekomernom učenju, respektivno. Postepenim usložnjavanjem modela može se smanjiti greška generalizacije i povećati tačnost prognoze. Prema tome, DCLN model se može iterativno podesiti na sledeći način. Struktura modela se u prvoj iteraciji može podesiti tako da se sastoji samo od slojeva osnovne podstrukture (S2S LSTM). Dužina dvosmernog S2S konteksta (P_4) se može postaviti na najmanji stepen broja 2 koji je veći od broja izlaznih vrednosti centroida opterećenja ($H \times 2 \times M$). U drugoj iteraciji se mogu dodati slojevi sporednih podstrukture (DCNN i SAE). Hiperparametri koji određuju broj DCNN i SAE slojeva (P_1 i P_5) se mogu postaviti na 1. Broj konvolucionih filtera (P_2) se može postaviti na najmanji stepen broja 2 koji je veći od broja ulaznih vrednosti centroida opterećenja ($Z \times 2 \times M$). Veličina filtera i korak sažimanja (P_3) se mogu postaviti na 2 (ne mogu biti manji od 2). Broj neurona u srednjem SAE sloju (P_6) se može postaviti na najveći stepen broja 2 koji je manji od broja ulaznih i izlaznih vrednosti centroida opterećenja ($\min \{Z, H\} \times 2 \times M$). Procenat neurona za odbacivanje (P_7) se može postaviti na 50%. U narednim iteracijama hiperparametri modela se mogu podešavati na sledeći način: 1) P_1 , P_3 i P_5 se mogu povećavati za 1; 2) P_2 , P_4 i P_6 se mogu množiti ili deliti sa 2, i 3) P_7 se može smanjivati za 5-10% (nije uobičajeno da se odbaci više od 50% neurona, jer to ukazuje na nepotrebnu složenost modela). Nakon što se podešavanje hiperparametara završi (kada se postigne prihvatljiva greška generalizacije), uobičajeno je da se podešeni model istrenira na kompletnom trening skupu, kako bi zatim istrenirani model mogao da se primeni na stvarnom test skupu.

Otpornost DCLN modela na potencijalne greške u podešavanju hiperparametara je demonstrirana kroz analizu različitih podešavanja u okviru studije slučaja u četvrtoj glavi ove disertacije.

4. STUDIJA SLUČAJA

Studija slučaja je izvedena tako da se predloženo rešenje (DCL) verifikuje sa stanovišta tačnosti i zauzeća računarskih resursa nad skupom realnih podataka sa pametnih brojila. U studiji slučaja su korišćena dva skupa podataka, čije su detaljne karakteristike date u tabeli 4.1:

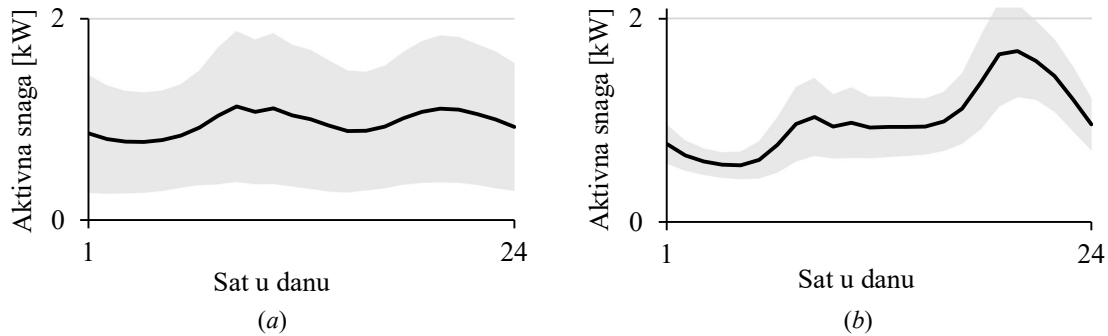
- 1) *UMass* podaci [64] (u nastavku skraćeno SP1). Podaci su prikupljeni u okviru Smart* projekta laboratorije za napredne softverske sisteme [179]. Prikupljeni podaci potiču sa pametnih brojila instaliranih kod niskonaponskih potrošača jedne severnoameričke DM.
- 2) *Smart Grid Smart City* (SGSC) podaci [65] (u nastavku SP2). Podaci su prikupljeni u okviru SGSC projekta australske vlade i industrijskog konzorcijuma predvođenim DP Ausgrid. Prikupljeni podaci potiču sa pametnih brojila instaliranih kod niskonaponskih potrošača jedne australske DM.

Tabela 4.1 – Skupovi podataka

Skup podataka	SP1 [64]	SP2* [65]
Lokacija	Nova Engleska, SAD	Njukasl, Australija
Period prikupljanja	~2 godine (~2015-2016)	~2 godine (~2012-2013)
Broj potrošača (N)	114 (stanovi u jednoj zgradi)	6.663 (potrošači sa najmanje godinu i po dana podataka u periodu prikupljanja)
Broj karakterističnih tipova potrošača	1 (domaćinstva)	
Broj potrošačkih grupa (K)	2 (opširnije u delu 4.2.1)	9 (opširnije u delu 4.2.1)
Broj fizičkih veličina (M)	1 (aktivna snaga)	
Broj vremenskih koraka izmerenih vrednosti u jednom danu (D^*)	96 – 1,440 (merenja na 1-15 minuta)	48 (merenja na 30 minuta)
Broj karakterističnih perioda	12 (meseći)	
Broj karakterističnih tipova dana	3 (radni dan, vikend i praznik [180, 181])	
Dužina horizonta prognoze (H), kao i broj vremenskih koraka na rezoluciji prognoze u jednom danu (D)	24 (vrednosti na satnom nivou)	
Odabrano zaostajanje (Z)	168 (satne vrednosti za sedam dana pre horizonta prognoze)	
Broj atributa koji predstavljaju enkodirane vremenske odrednice (V)	7 (jedan indikator praznika i tri 2D koordinate na jediničnom krugu za sat u danu, dan u nedelji i mesec u godini)	
Broj dostupnih meteoroloških faktora (W^*)	10 (numeričke vrednosti na nivou sata, opisane u delu 4.2.1)	
Broj odabranih meteoroloških faktora (W)	7 (opširnije u delu 4.2.1)	8 (opširnije u delu 4.2.1)
Veličina originalnih podataka ($N \times M \times D^* + D \times W^*$ vrednosti u jednom danu)	~488 MiB (~ $6,4 \times 10^7$ 4-bajtnih vrednosti i vremena)	~1,42 GiB (~ $1,9 \times 10^8$ 4-bajtnih vrednosti i vremena)
Veličina centroida opterećenja ($K \times M$ serija μ i σ vrednosti sa po sa po D vrednosti u svakom danu)	~0,27 MiB (~ 7×10^4 4-bajtnih vrednosti)	~1,2 MiB (~ 3×10^5 4-bajtnih vrednosti)
Veličina DCLN modela (K modela sa $2 \times M + V + W$ ulaznih i $2 \times M$ izlaznih atributa)	~20 MiB (~ $5,3 \times 10^6$ 4-bajtnih parametara modela)	~91 MiB (~ $2,4 \times 10^7$ 4-bajtnih parametara modela)

* Zbog nedostupnosti meteoroloških podataka za SP2 i anonimnosti potrošača, meteorološki podaci su preuzeti sa worldweatheronline.com.

Prosečni dnevni HDO iz SP1 i SP2 su prikazani na slici 4.1 (crnom linijom), sa svojom standardnom devijacijom (sivom površinom). SP1 i SP2 sadrže merenja potrošnje aktivne snage u periodu od približno dve godine. Prosečna potrošnja aktivne snage u jednom satu je ~ 1 kW kod potrošača iz oba skupa podataka. Prva godina podataka je upotrebljena za inicijalno treniranje regresionih modela za kratkoročnu probablističku prognozu opterećenja na nivou jednog sata za jedan dan unapred. Druga godina podataka je upotrebljena za inkrementalno testiranje i ažuriranje regresionih modela. Regresioni modeli su ažurirani na svakih 15 dana sa novim podacima, analogno [40]. Eksperimentalno je ustanovljeno da češće ažuriranje ne vodi ka značajnom poboljšanju tačnosti prognoze.



Slika 4.1 – Prosečan dnevni HDO za potrošače iz SP1 (a) i SP2 (b)

4.1. Dizajn studije slučaja

DCL rešenje je upoređeno sa tri standardna rešenja (DeepAR [60], DMN [58] i DEL [62]) i jednom varijantom DCL rešenja (eng. *SVM Centroid Learning*, SCL), koja se zasniva na primeni SVM umesto DCLN modela. SCL je primenjen u studiji slučaja kako bi se verifikovalo da ključan korak predloženog rešenja nije lako zamenjiv [17]. Svih pet rešenja su podešena tako da koriste iste potrošačke grupe, isto zaostajanje za opterećenje koje prethodi horizontu prognoze, iste vremenske odrednice i iste meteorološke faktore. Potrošačke grupe su dobijene primenom TSGA sa CH indeksom. Probablistička prognoza je verifikovana na nivou pojedinačnih potrošača primenom PS, WS i Diebold-Mariano testa.

DCLN je podešen na sledeći način: $P_1 = 5$, $P_2 = 512$, $P_3 = 2$, $P_4 = 256$, $P_5 = 5$, $P_6 = 16$ i $P_7 = 50\%$. Opravdanje za odabir navedenih vrednosti hiperparametara je dato u delu 4.2.1. Parametri DCLN modela su optimizovani primenom Adam algoritma u najviše 200 ponavljanja (epoha) uz rano prekidanje treniranja nakon 20 epoha ako se greška ne umanjuje. Razvrstavanje trening primera u pakete tokom optimizacije parametara modela je uobičajena praksa za uštedu RAM prostora kada je broj primera velik. Zbog toga su trening primeri razvrstani u pakete veličine 128 (eng. *batch size*). Prema tome, Adam je u svakoj epohi izvršen u onoliko iteracija koliko ima paketa.

SCL je podešen tako da koristi ν -SVM [182]. ν -SVM je jednostavniji za parametrizaciju od klasičnog ϵ -SVM i ima sledeće hiperparametre: parametar kojim se kontrolišu donja granica broja *support* vektora i gornja granica broja trening primera koji se mogu smatrati pogrešnim (ν), širina kernela (γ), penalizacija greške (e) i tolerancija greške za rano prekidanje treniranja (ρ). SCL je podešen da koristi Gausov kernel, $\nu \approx 0,05$, $\gamma \approx 0,01$, $e \approx 0,75$ i $\rho \approx 0,05$. Navedena podešavanja su dobijena primenom GOA nad nasumično odabranim uzorcima iz SP1 i SP2, kao u [163]. DeepAR, DMN i DEL su podešeni kao u [58, 60, 62].

Implementacija primenjenih rešenja je dostupna na GitHub platformi [183]. Za implementaciju su upotrebljena dva programska jezika, C# i R. C# je programski jezik opšte namene, a R je programski jezik za statističke proračune [184]. R je koristan alat za analizu podataka zbog jednostavnosti upotrebe širokog spektra programskih biblioteka za obradu podataka [185]. Programske biblioteke od najvećeg značaja za implementaciju navedene su u tabeli 4.2. Specifikacija hardvera na kojem su izvršeni svi proračuni data je u tabeli 4.3. Prateći trend razvoja aplikacija za analizu podataka [29, 186], duboko učenje je primenjeno na GPU. GPU sa hiljadama procesorskih jezgara omogućava visoku paralelizaciju obrade velikih količina podataka uz relativno malu potrošnju električne energije [187].

Tabela 4.2 – Programske biblioteke

Zadatak	Metod	Biblioteka
Ekstrakcija atributa	HMLPSA	MTSR [188]
	PCA	stats (ugrađena biblioteka)
Odabir atributa	PACF	stats (ugrađena biblioteka)
	JMIM	praznik [189]
Klasterizacija	TSGA, HBKM	stats, fastcluster [190] i clusterSim [191]
	ECCF	stats, kohonen [192] i clusterSim [191]
	FSC	stats, speccalt [193] i clusterSim [191]
Regresija i optimizacija parametara modela	SVM	e1071 [194]
	DeepAR	modeltime.gluonts [195]
	DMN, DEL, DCL	keras sa tensorflow pozadinom [196, 197]
Optimizacija hiperparametara	GOA	metaheuristicOpt [198]

Tabela 4.3 – Specifikacija hardvera

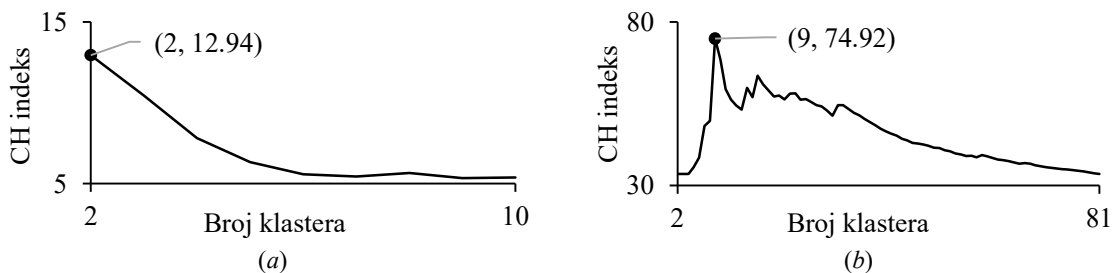
Jedinica	Naziv	Karakteristike
CPU	AMD FX-8350	8 jezgara, 4 GHz
RAM	Kingston HyperX	32 GB DDR3, 1866 MHz
GPU	NVIDIA GeForce GTX TITAN	2688 jezgara, 837 MHz, 6 GB GDDR5
SDD	Crucial MX200	250 GB, 555/500 MB/s

4.2. Rezultati i diskusija

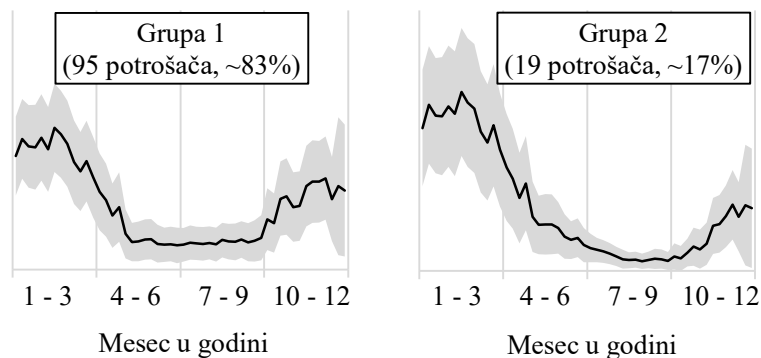
U ovom delu su predstavljeni i prodiskutovani: 1) rezultati nekih od međukoraka predloženog rešenja (međurezultati); 2) konačni rezultati prognoze; 3) rezultati primene prethodno navedenih mera tačnosti u cilju verifikacije prognoze, kao i 4) zauzeće računarskih resursa tokom izvršavanja proračuna.

4.2.1. Međurezultati

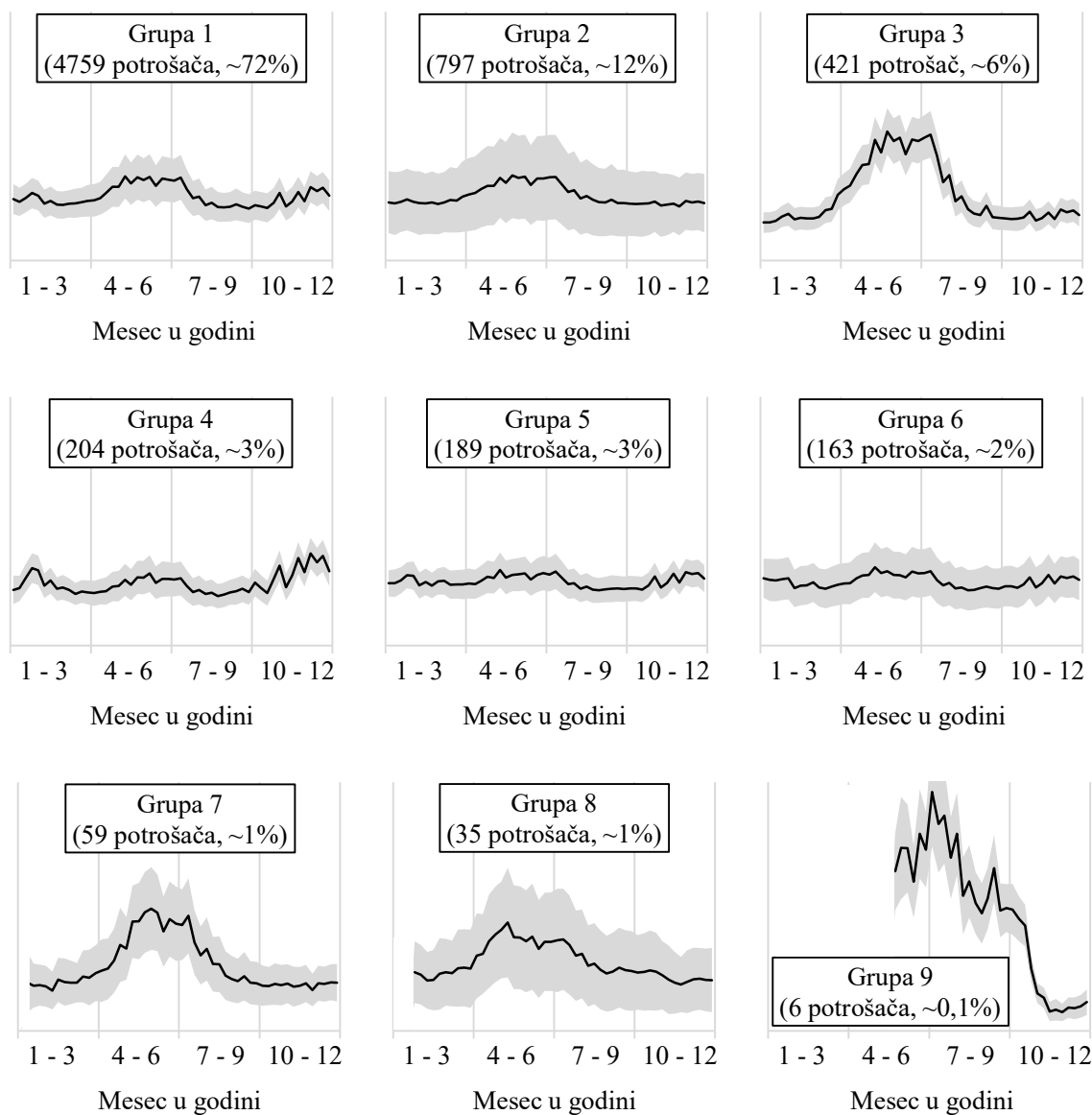
Prvih šest koraka predloženog rešenja služi za redukciju količine podataka i ekstrakciju centroida opterećenja. Centroidi opterećenja su dobijeni primenom TSGA sa CH indeksom. Slika 4.2 prikazuje CH indeks za različite konfiguracije klastera (veće vrednosti indeksa ukazuju na veću validnost klastera). U skladu s tim, dve potrošačke grupe su formirane za SP1, a devet potrošačkih grupa za SP2. Rezultujući centroidi opterećenja su prikazani na slikama 4.3 i 4.4, respektivno. Pritom su centroidi prikazani samo za prvu godinu podataka, koja je upotrebljena za inicijalno treniranje regresionih modela, dok je prognoza za narednu godinu prikazana u delu 4.2.2. Prosečne vrednosti su prikazane crnom linijom, a standardne devijacije sivom površinom. Radi jednostavnosti, vrednosti su prikazane na sedmodnevnom nivou. Zajedno sa centroidima je prikazan i broj potrošača u svakog grupi, kao i procenat broja potrošača u odnosu na ceo skup podataka. Za poslednje tri potrošačke grupe iz SP2 ne postoje podaci za celu godinu. U tabeli 4.1 se može primetiti da je veličina prikazanih centroida za ~99,9% manja od veličine originalnih podataka, kao i da je veličina istreniranih DCLN modela ~93-96% manja od veličine originalnih podataka. To ukazuje na visoku efikasnost predloženog rešenja sa stanovišta zauzeća memorijskog prostora.



Slika 4.2 – CH indeks za klasterne potrošača iz SP1 (a) i SP2 (b)

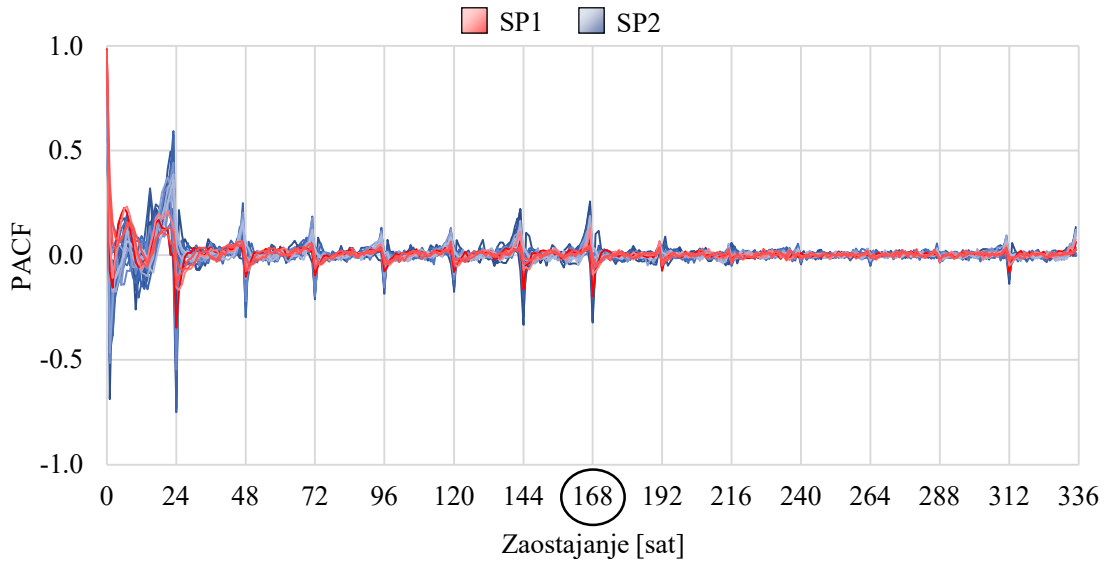


Slika 4.3 – Centroidi opterećenja za potrošače iz SP1



Slika 4.4 – Centroidi opterećenja za potrošače iz SP2

Rezultat primene PACF u odabiru zaostajanja za centroide opterećenja prikazan je na slici 4.5. Može se primetiti da kod svih centroida postoji značajna autokorelacija sa opterećenjem u prethodnih 7 dana, što opravdava odabrano zaostajanje od 168 sati, koje je zaokruženo na slici 4.5 i navedeno u tabeli 4.1.



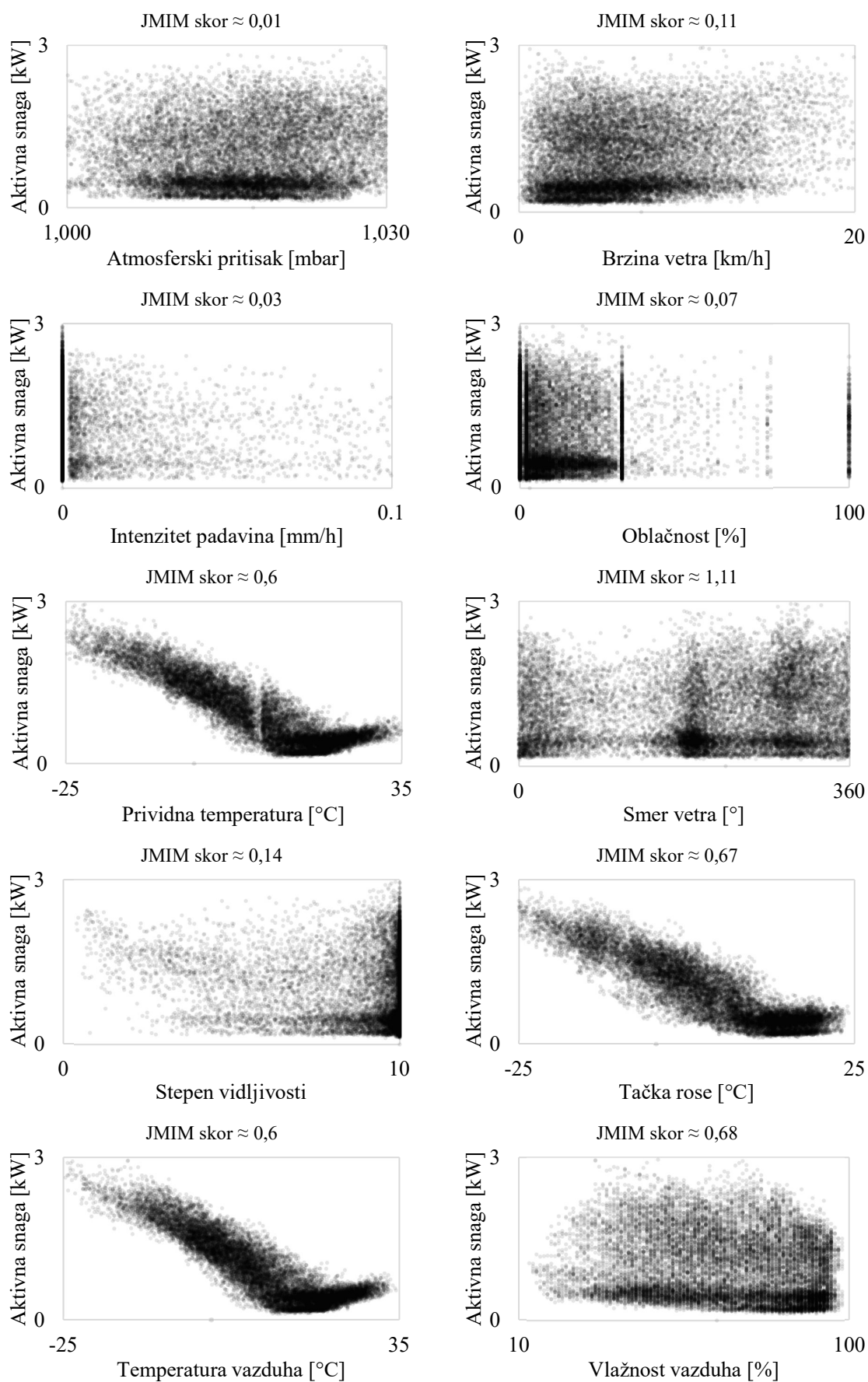
Slika 4.5 – Rezultat odabira zaostajanja za centroide opterećenja

Rezultat primene JMIM metode u odabiru meteoroloških faktora prikazan je u tabeli 4.4. JMIM skor je prikazan u vidu toplotne mape, a odabrani meteorološki faktori su obeleženi tačkom. Odabir je prikazan za svaki centralni momenat sadržan u centroidima opterećenja potrošačkih grupa iz SP1 i SP2. Detaljniji uvid u odnos između meteoroloških podataka i podataka o opterećenju iz SP1 i SP2 dat je na slikama 4.6 i 4.7, respektivno. Pritom je za svaki meteorološki faktor prikazan prosečan JMIM skor u odnosu na sve centroide opterećenja u okviru posmatranog skupa podataka. Kod oba skupa podataka se može primetiti da se opterećenje značajno menja s promenom vrednosti temperaturnih promenljivih (temperature vazduha, prividne temperature i tačke rose) i vlažnosti vazduha, što se takođe ogleda u visokom JMIM skor. Pored toga, JMIM skor ukazuje na to da smer vetra doprinosi značajnom količinom novih informacija u slučaju SP1, a oblačnost u slučaju SP2. U tabeli 4.4 se takođe može primetiti da intenzitet padavina nije odabran ni u jednom slučaju, dok se na slikama 4.6 i 4.7 može primetiti da promene u intenzitetu padavina nemaju značajan uticaj na promene opterećenja.

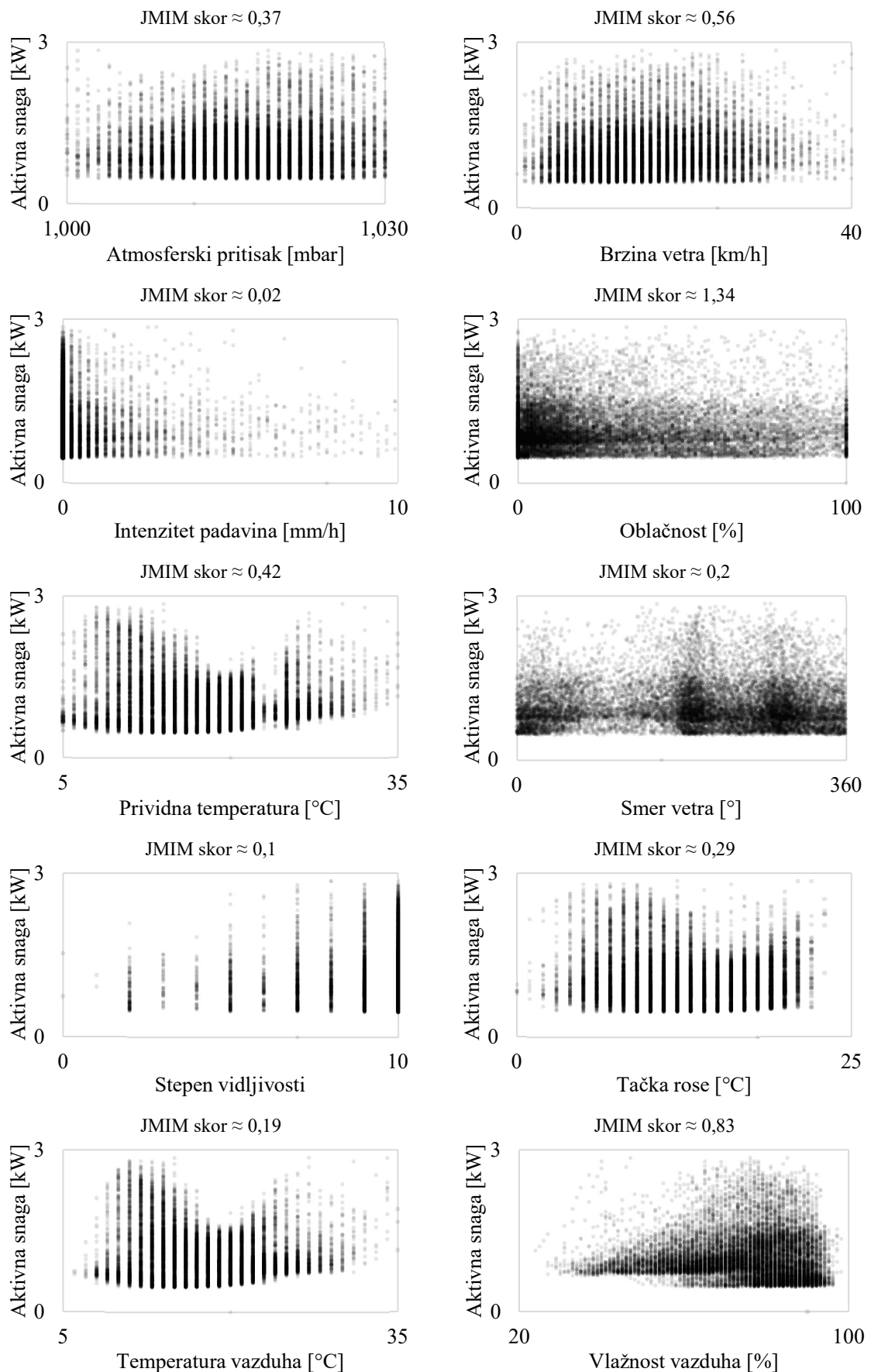
Tabela 4.4 – Rezultat odabira meteoroloških faktora

Skup podataka	SP1				SP2																		
	1		2		1		2		3		4		5		6		7		8		9		
Potrošačka grupa	μ	σ	μ	σ	μ	σ	μ	σ	μ	σ	μ	σ	μ	σ	μ	σ	μ	σ	μ	σ	μ	σ	
Atmosferski pritisak					•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Brzina vetra	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Intenzitet padavina																							
Oblačnost					•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Prividna temperatura	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Smer vetra	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Stepen vidljivosti	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Tačka rose	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Temperatura vazduha	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Vlažnost vazduha	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•

Nizak JMIM skor Visok

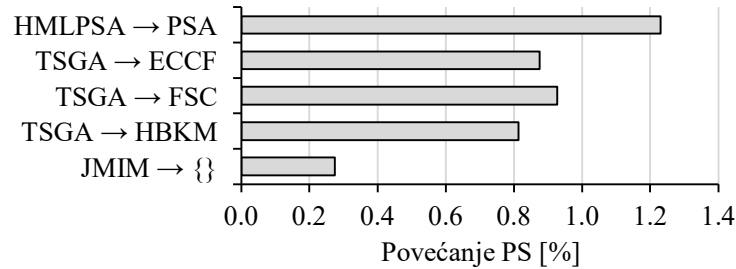


Slika 4.6 – Uticaj meteoroloških faktora na agregirano opterećenje potrošača iz SP1



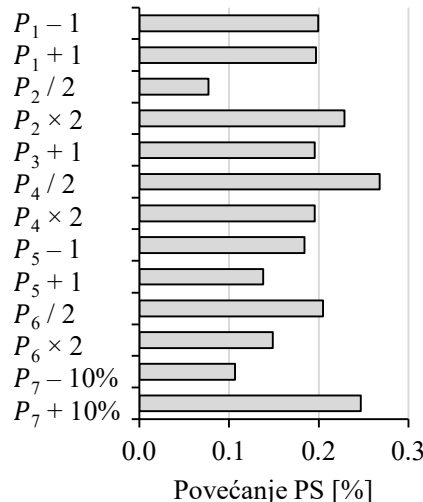
Slika 4.7 – Uticaj meteoroloških faktora na agregirano opterećenje potrošača iz SP2

Opravdanost primene odabranih metoda za pripremu podataka u okviru DCL rešenja kvantifikovana je degradacijom tačnosti prognoze koja se dobija nakon drugačijeg odabira metoda. Degradacija tačnosti prognoze je prikazana na slici 4.8 u vidu prosečnog povećanja PS za sve potrošače iz SP1 i SP2. Na slici je prvo prikazan uticaj zamene HMLPSA metode, koja se koristi za aproksimaciju HDO, sa konkurentnom PSA metodom. Zatim je prikazan uticaj zamene TSGA algoritma, koji se koristi za kreiranje potrošačkih grupa, sa konkurentnim ECCF, FSC i HBKM algoritmima. Na kraju je prikazan uticaj zamene odabira meteoroloških faktora, koji se dobijaju primenom JMIM metode, sa odabirom svih meteoroloških faktora, bez primene metoda odabira atributa. Može se primetiti da zamena bilo koje od navedenih metoda vodi ka degradaciji tačnosti prognoze, što opravdava njihovu primenu. Takođe se može primetiti da zamena metoda koje se primenjuju ranije u okviru predloženog rešenja, kao što je HMLPSA, ima veći uticaj nego zamena metoda koje se primenjuju kasnije, kao što je JMIM. Takav rezultat je očekivan s obzirom da sekvencijalna obrada podataka u okviru predloženog rešenja prirodno vodi ka akumulaciji grešaka.



Slika 4.8 – Uticaj promene odabranih metoda za pripremu podataka na tačnost prognoze

Opravdanost odabranih vrednosti hiperparametara DCLN modela, koje su navedene u delu 4.1, je takođe kvantifikovana degradacijom tačnosti prognoze koja se dobija nakon drugačijeg odabira vrednosti. Degradacija tačnosti prognoze je prikazana na slici 4.9 u vidu prosečnog povećanja PS za sve potrošače iz SP1 i SP2. Hiperparametri koji određuju broj DCNN i SAE slojeva (P_1 i P_5) su promenjeni za 1 u odnosu na odabranu vrednost ($P_1 = P_5 = 5$). Broj konvolucionih filtera (P_2), dužina dvosmernog S2S konteksta (P_4) i broj neurona u srednjem SAE sloju (P_6) su prema uobičajenoj praksi promenjeni za stepen broja 2 u odnosu na odabrane vrednosti ($P_2 = 512$, $P_4 = 256$ i $P_6 = 16$). Veličina konvolucionih filtera i koraka sažimanja (P_3) je povećana za 1 u odnosu na odabranu vrednost ($P_3 = 2$), od koje ne može biti manja. Procenat neurona za odbacivanje (P_7) je promenjen za 10% u odnosu na odabranu vrednost ($P_7 = 50\%$). Može se primetiti da svaka pojedinačna promena odabranih vrednosti negativno utiče na tačnost prognoze. Takođe se može primetiti da je DCLN model veoma otporan na odabir vrednosti hiperparametara, s obzirom da je u svim analiziranim slučajevima degradacija tačnosti relativno mala (ispod 0,3%). Prema tome, predloženi model se može dovoljno dobro podesiti na osnovu preporuka koje su navedene u delu 3.5.2, bez primene algoritama za optimizaciju hiperparametara, koji bi degradirali vreme izvršavanja.

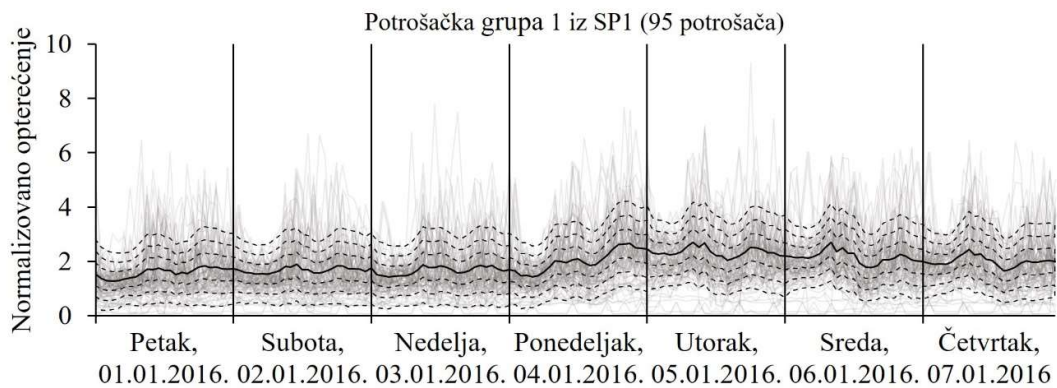


Slika 4.9 – Uticaj promene odabranih vrednosti hiperparametara DCLN modela na tačnost prognoze

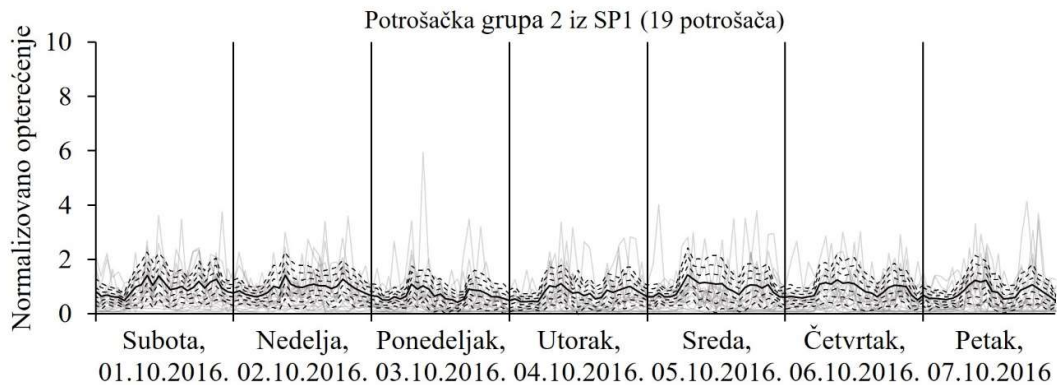
4.2.2. Rezultati prognoze

Neki od rezultata primene predloženog rešenja na SP1 i SP2 prikazani su na slikama 4.10-4.20. Na svakoj slici su prikazane ostvarene vrednosti (NAHDO) i prognozirane vrednosti centroida opterećenja za jednu potrošačku grupu tokom prvih sedam dana u jednom, nasumično odabranom mesecu. Ostvarene vrednosti su predstavljene sivim linijama, a prognozirane vrednosti crnim linijama. Pune crne linije predstavljaju prognozu prosečnih vrednosti na nivou potrošačke grupe, dok isprekidane crne linije predstavljaju intervale predviđanja širine jedne, dve i tri prognozirane standardne devijacije. Prema empirijskom pravilu „68-95-99,7“, očekivano je da intervale predviđanja širine jedne, dve i tri standardne devijacije obuhvate 68%, 95% i 99,7% ostvarenih vrednosti, respektivno. Može se primetiti da prognoza jeste prilagođena promenama opterećenja i da intervale predviđanja obuhvataju većinu ostvarenih vrednosti, kao i da postoje određena odstupanja (greške u prognozi).

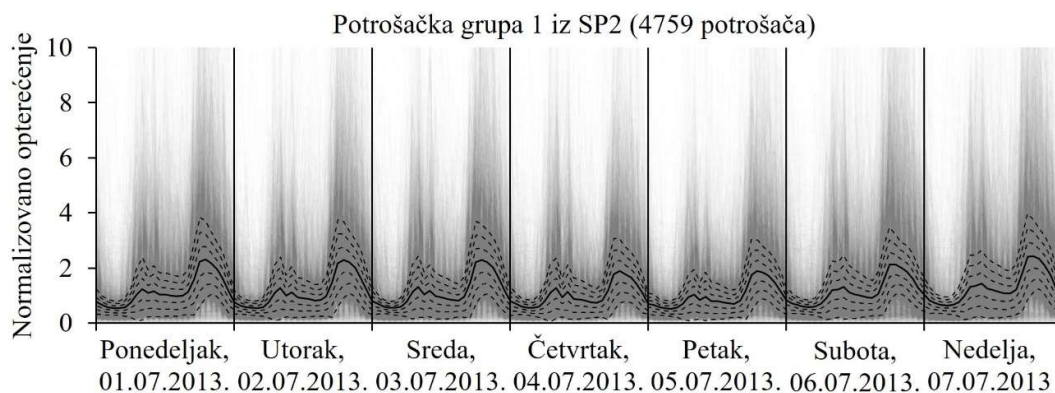
Radi kompletnosti prikaza, u tabeli 4.5 je prikazan PS za svaku potrošačku grupu, mesec i tip dana (manji PS ukazuje na tačniju prognozu). Na osnovu slika 4.10-4.20, kao i na osnovu tabele 4.5, može se zaključiti da prognoza najviše odstupa od ostvarenih vrednosti kada unutar potrošačke grupe dođe do nagle promene u opterećenju. Takva odstupanja se pretežno javljaju u periodima godine kada je opterećenje veće, što su proleće i zima za potrošače iz SP1, a leto i jesen za potrošače iz SP2. Takođe, takva odstupanja se manje javljaju tokom radnih dana, a više tokom vikenda i praznika. PS je detaljnije analiziran u delu 4.2.3.



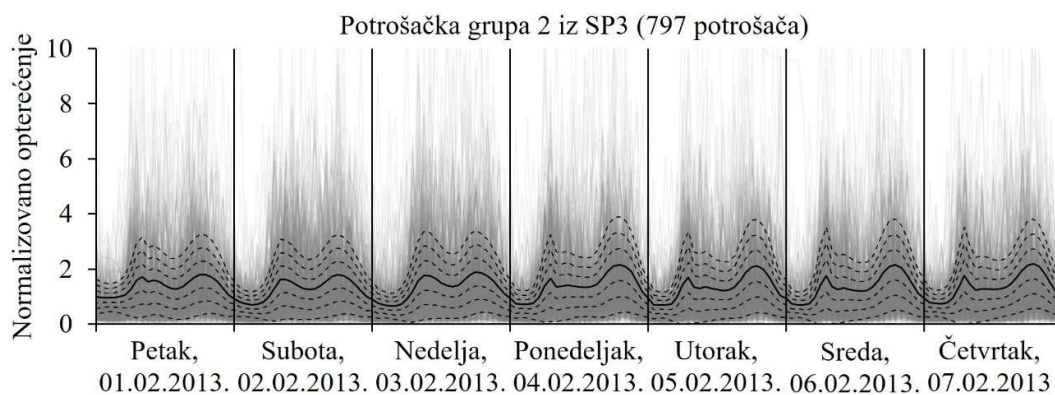
Slika 4.10 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 1 iz SP1.



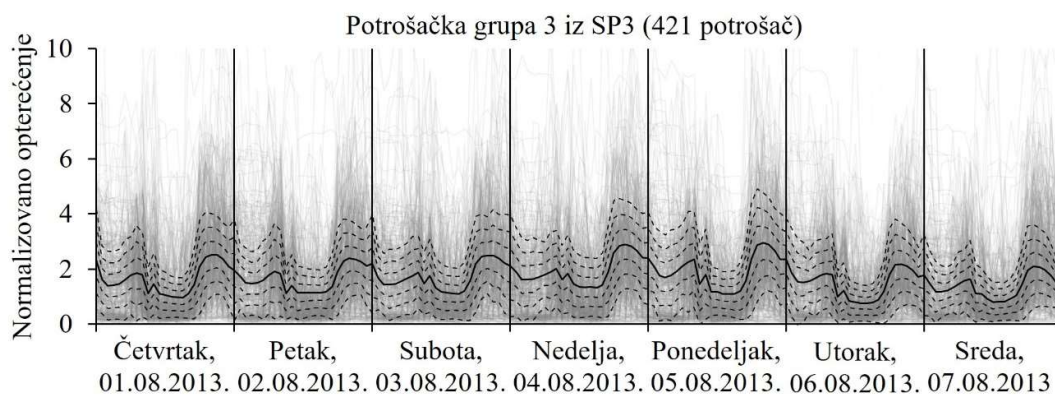
Slika 4.11 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 2 iz SP1.



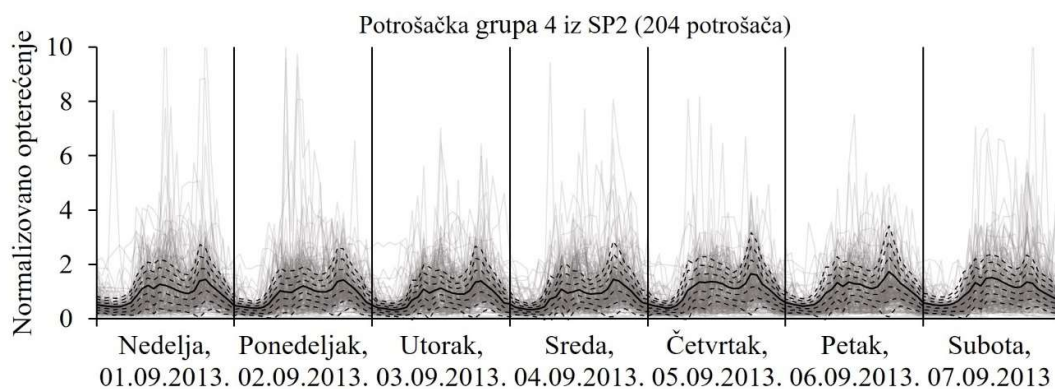
Slika 4.12 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 1 iz SP2.



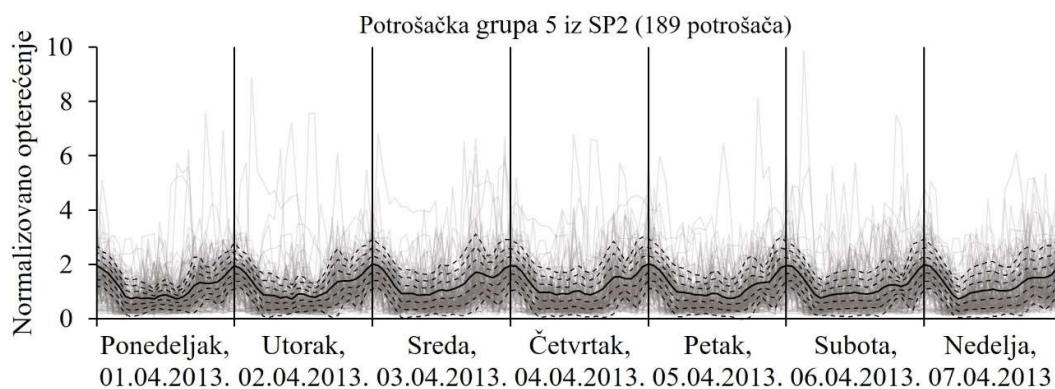
Slika 4.13 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 2 iz SP2.



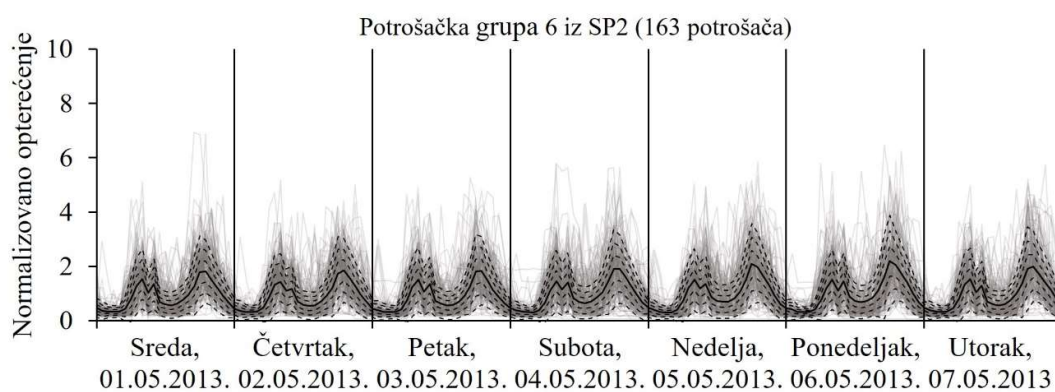
Slika 4.14 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 3 iz SP2.



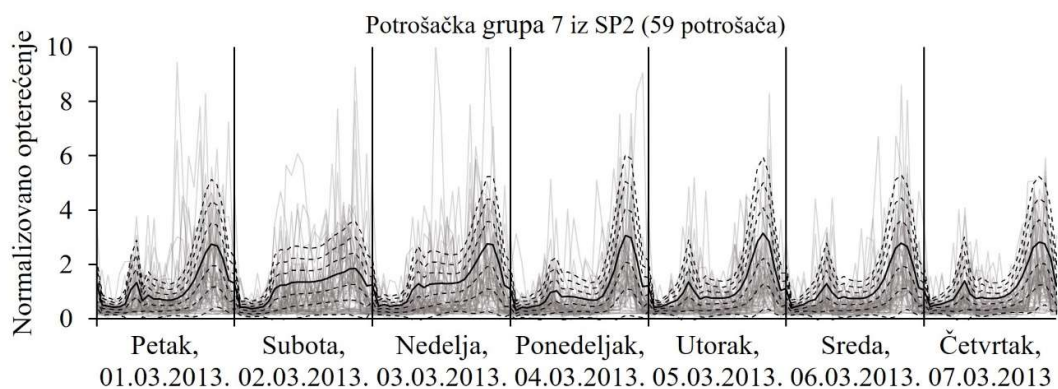
Slika 4.15 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 4 iz SP2.



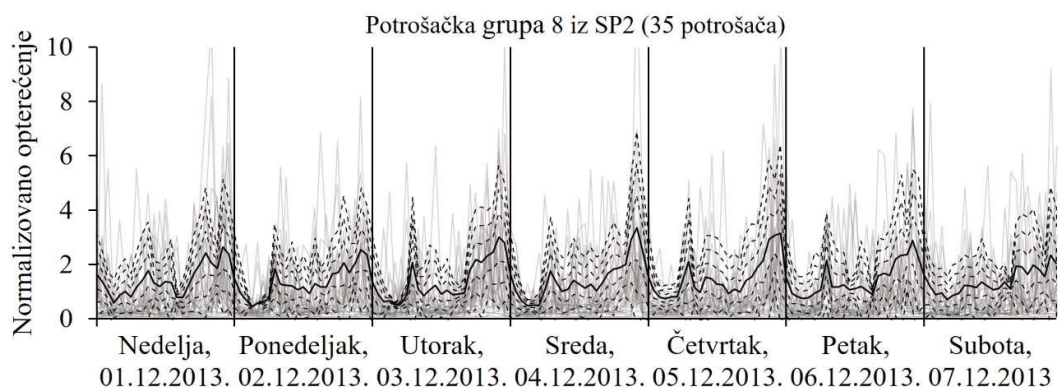
Slika 4.16 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 5 iz SP2.



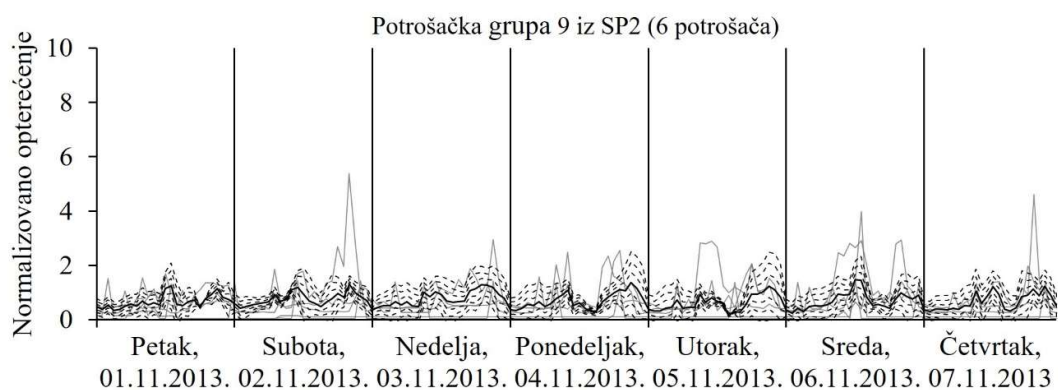
Slika 4.17 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 6 iz SP2.



Slika 4.18 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 7 iz SP2.



Slika 4.19 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 8 iz SP2.



Slika 4.20 – Poređenje ostvarenih i prognoziranih vrednosti za potrošačku grupu 9 iz SP2.

Tabela 4.5 – PS na nivou meseci i tipova dana.

Skup podataka	Grupa	Tip dana	Januar	Februar	Mart	April	Maj	Jun	Jul	Av gust	Septembar	Oktober	Novembar	Decembar
SP1	1	r. d.	0,312	0,311	0,282	0,256	0,212	0,157	0,156	0,146	0,155	0,223	0,272	0,303
		v.	0,294	0,327	0,300	0,260	0,205	0,159	0,154	0,156	0,156	0,234	0,276	0,320
		p.	0,305	0,394	×	0,268	0,202	0,171	0,158	×	0,152	0,236	0,299	×
	2	r. d.	0,275	0,270	0,265	0,231	0,207	0,153	0,149	0,139	0,138	0,180	0,241	0,254
		v.	0,271	0,296	0,275	0,251	0,201	0,162	0,149	0,148	0,145	0,191	0,246	0,271
		p.	0,261	0,321	×	0,258	0,224	0,164	0,148	×	0,146	0,194	0,266	×
SP2	1	r. d.	0,071	0,076	0,061	0,059	0,057	0,074	0,097	0,098	0,083	0,061	0,059	0,058
		v.	0,072	0,085	0,071	0,064	0,063	0,080	0,099	0,105	0,079	0,061	0,061	0,063
		p.	0,067	0,069	0,059	0,056	0,061	0,075	0,090	×	0,090	0,068	0,055	0,053
	2	r. d.	0,065	0,070	0,058	0,056	0,053	0,070	0,094	0,095	0,080	0,058	0,054	0,053
		v.	0,067	0,078	0,068	0,061	0,059	0,076	0,097	0,103	0,076	0,058	0,057	0,058
		p.	0,060	0,063	0,055	0,052	0,059	0,071	0,087	×	0,088	0,064	0,053	0,048
	3	r. d.	0,075	0,081	0,065	0,062	0,059	0,078	0,100	0,101	0,088	0,066	0,063	0,062
		v.	0,077	0,090	0,074	0,068	0,065	0,084	0,102	0,109	0,084	0,066	0,064	0,067
		p.	0,070	0,073	0,061	0,060	0,064	0,079	0,096	×	0,095	0,071	0,059	0,059
	4	r. d.	0,072	0,076	0,061	0,058	0,056	0,078	0,103	0,104	0,089	0,063	0,060	0,058
		v.	0,075	0,086	0,072	0,064	0,063	0,085	0,107	0,112	0,086	0,064	0,062	0,063
		p.	0,071	0,070	0,057	0,055	0,061	0,080	0,096	×	0,097	0,071	0,056	0,056
	5	r. d.	0,074	0,077	0,059	0,056	0,056	0,074	0,098	0,098	0,083	0,060	0,059	0,058
		v.	0,073	0,084	0,070	0,063	0,061	0,079	0,101	0,105	0,078	0,060	0,062	0,062
		p.	0,069	0,070	0,055	0,054	0,061	0,073	0,089	×	0,094	0,065	0,055	0,052
	6	r. d.	0,073	0,081	0,066	0,064	0,062	0,085	0,111	0,112	0,094	0,069	0,061	0,064
		v.	0,077	0,089	0,077	0,070	0,069	0,090	0,115	0,122	0,087	0,069	0,064	0,068
		p.	0,072	0,073	0,064	0,060	0,071	0,084	0,106	×	0,105	0,080	0,058	0,056
	7	r. d.	0,078	0,080	0,061	0,060	0,056	0,071	0,091	0,091	0,080	0,058	0,059	0,061
		v.	0,079	0,087	0,072	0,064	0,060	0,074	0,089	0,100	0,078	0,059	0,063	0,067
		p.	0,072	0,076	0,063	0,055	0,061	0,073	0,081	×	0,084	0,061	0,056	0,054
	8	r. d.	0,070	0,082	0,068	0,064	0,058	0,079	0,105	0,116	0,101	0,066	0,063	0,062
		v.	0,075	0,093	0,078	0,068	0,066	0,086	0,113	0,120	0,100	0,067	0,066	0,064
		p.	0,069	0,072	0,066	0,058	0,065	0,075	0,099	×	0,109	0,077	0,060	0,056
	9	r. d.	0,043	0,061	0,036	0,035	0,047	0,064	0,105	0,108	0,086	0,061	0,041	0,037
		v.	0,047	0,071	0,043	0,040	0,054	0,064	0,102	0,117	0,086	0,056	0,040	0,041
		p.	0,048	0,048	0,030	0,031	0,047	0,064	0,091	×	0,092	0,077	0,046	0,033

r. d. – radni dan
v. – vikend
p. – praznik

Manji



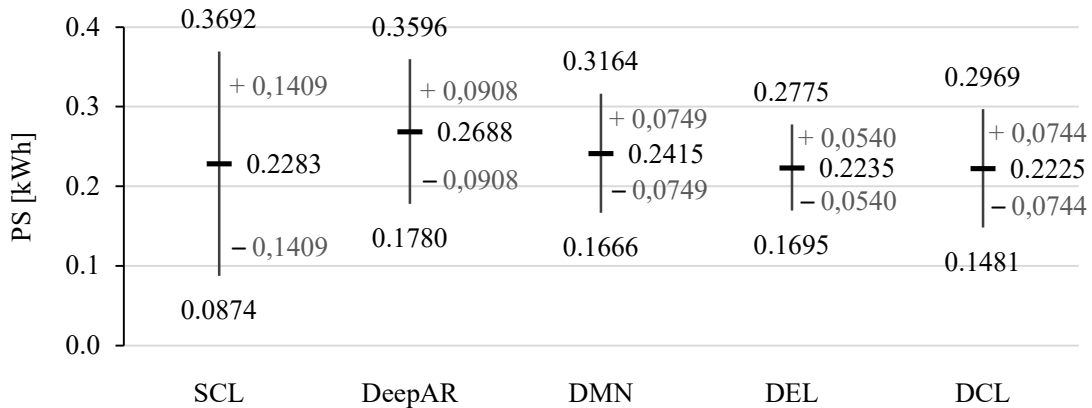
Veći

4.2.3. Verifikacija prognoze

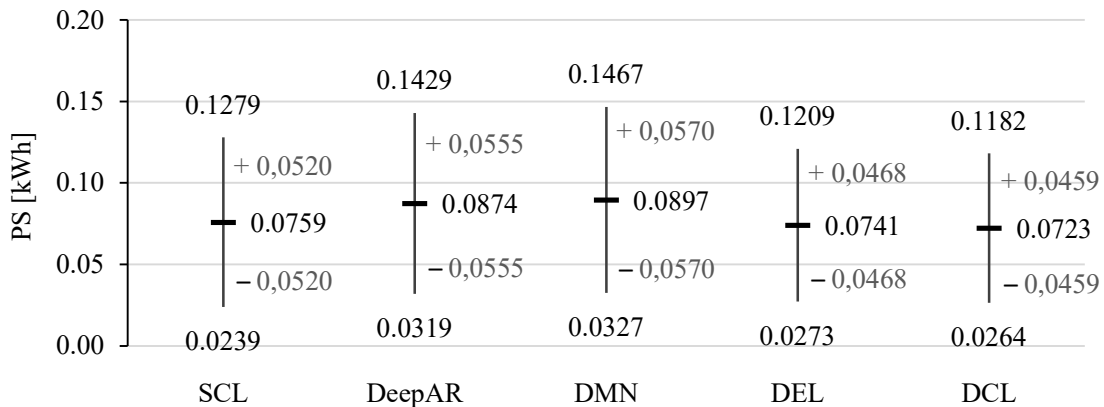
Predloženo rešenje (DCL) upoređeno je sa tri standardna rešenja (DeepAR, DMN i DEL) i jednom varijantom predloženog rešenja (SCL). Probabilistička prognoza je verifikovana na nivou pojedinačnih potrošača primenom PS, WS i Diebold-Mariano testa (manji PS i WS ukazuju na tačniju prognozu). Pritom je PS analiziran detaljnije od WS jer pruža detaljniji uvid u oštrinu probabilističke prognoze. Rezultati poređenja su prikazani na slikama 4.21-4.31 i prodiskutovani u nastavku.

PS za potrošače iz SP1 i SP2 je prikazan na slikama 4.21 i 4.22, respektivno. Prosek i standardna devijacija PS za različite potrošače su predstavljeni horizontalnim i vertikalnim linijama, respektivno. Može se primetiti da se primenom DCL rešenja dobija u proseku najmanji PS, kao i da je njegova standardna devijacija jedna od manjih. Takođe se može primetiti da se primenom SCL rešenja dobija u proseku mali PS u poređenju sa drugim rešenjima, ali da je njegova standardna devijacija jedna od većih. Prema tome, prikazani PS pokazuje da je primena dubokog učenja u okviru predloženog rešenja opravdana. S druge strane, primenom DeepAR i DMN rešenja se dobija u proseku najveći PS, a primenom DEL rešenja se dobija PS koji je u proseku manji nego za SCL, ali i dalje veći nego za DCL.

S obzirom da ne postoji standardizovana mera tačnosti koja izražava PS nezavisno od skale vrednosti, PS se može posmatrati u odnosu na prosečnu ostvarenu potrošnju u jednom satu, koja iznosi ~1 kW za potrošače iz SP1 i SP2. Iz toga sledi da je prosečna greška u prognozi koja je dobijena primenom predloženog rešenja ~22% za SP1 i ~7% za SP2. Dobijena razlika u tačnosti je u saglasnosti sa razlikom između SP1 i SP2 u odnosu na validnost klastera (slika 4.2). S obzirom da SP2 sadrži podatke za ~60 puta više potrošača nego SP1, može se zaključiti da se prednost primenjenih rešenja povećava sa povećanjem broja potrošača zbog prirodno veće mogućnosti za pronalaženjem kompaktnijih i različitijih klastera. Detaljniji uvid u PS je dat na slikama 4.23-4.28.

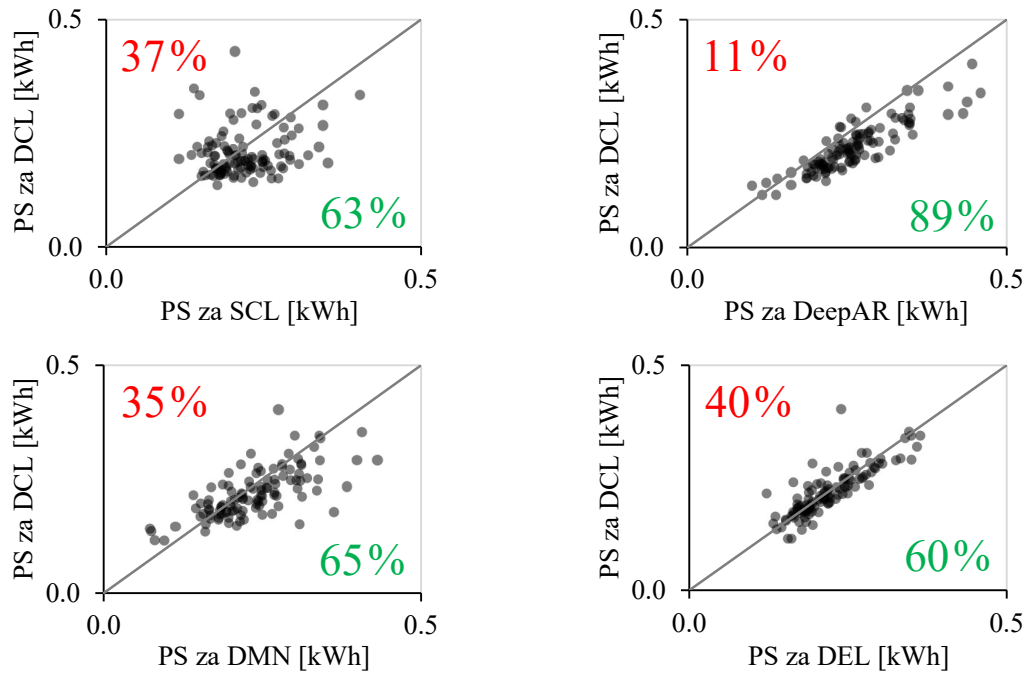


Slika 4.21 – PS za potrošače iz SP1

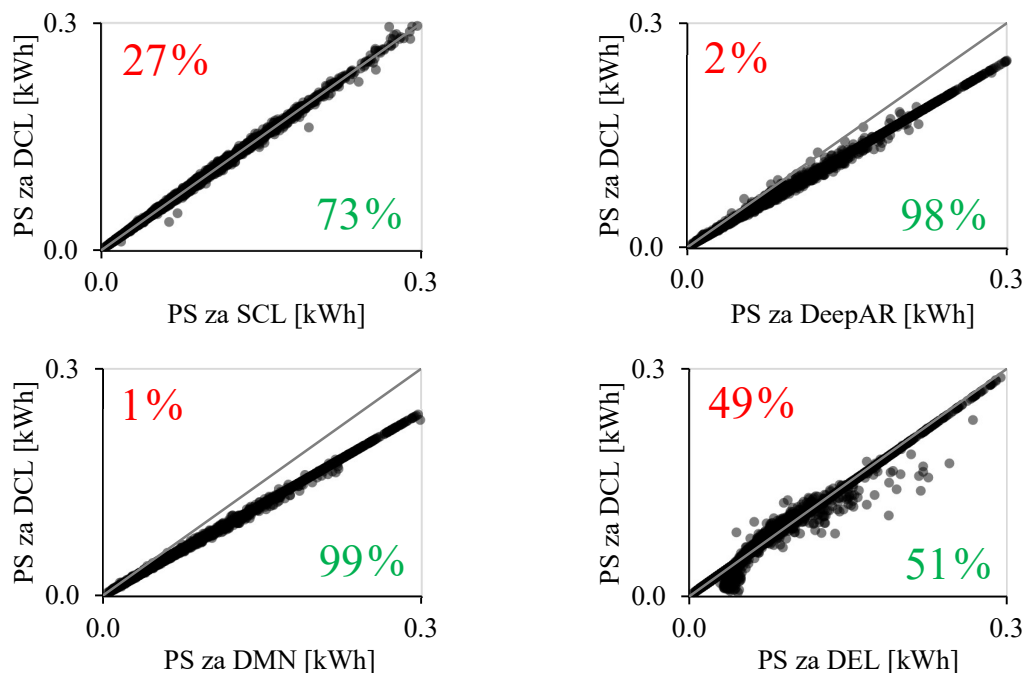


Slika 4.22 – PS za potrošače iz SP2

Slike 4.23 i 4.24 prikazuju poređenje između DCL rešenja i konkurentnih rešenja u odnosu na PS na nivou potrošača iz SP1 i SP2, respektivno. Vrednosti ispod sive dijagonalne linije prikazuju PS na nivou potrošača kod kojih DCL u proseku nudi tačniju prognozu. Vrednosti iznad linije prikazuju PS na nivou potrošača kod kojih posmatrano konkurentno rešenje u proseku nudi tačniju prognozu. Broj potrošača na koje se odnosi PS sa svake strane dijagonale je prikazan u procentima. Može se uočiti da je veći procenat potrošača kod kojih primena DCL rešenja rezultuje u manjim PS nego primena drugih rešenja, što je u saglasnosti sa slikama 4.21 i 4.22. Takođe se može uočiti da je u poređenju sa DEL rešenjem, razlika u procentima najmanja, kao i da je većina prikazanih vrednosti pozicionirana blizu dijagonale. Prema tome, razlike između DCL i DEL na nivou istih potrošača su manje nego razlike između DCL i drugih rešenja.

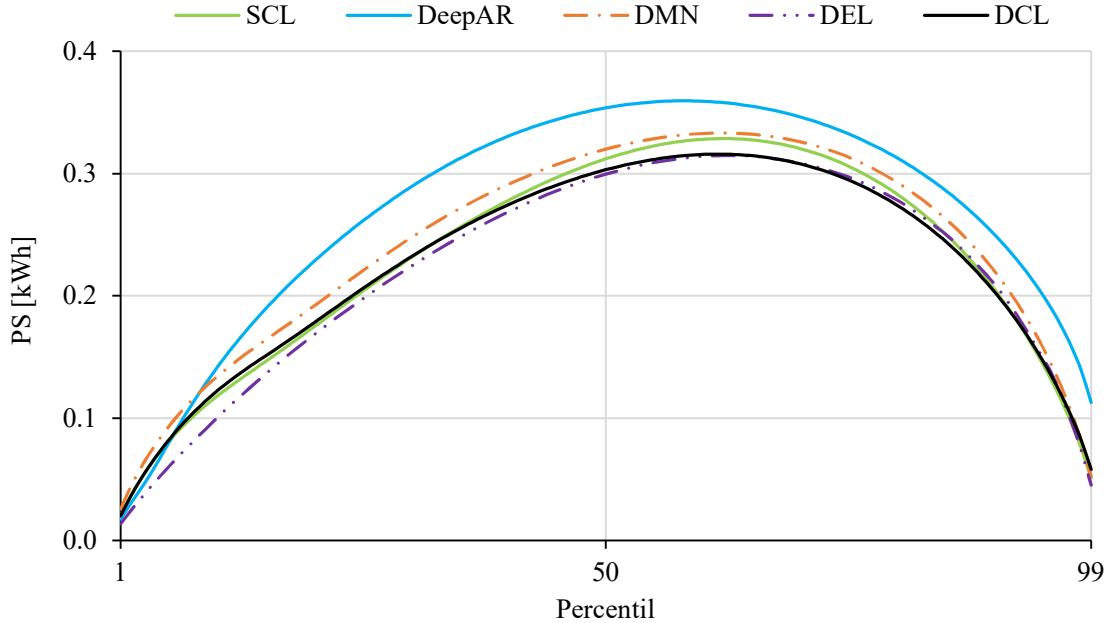


Slika 4.23 – PS na nivou potrošača iz SP1

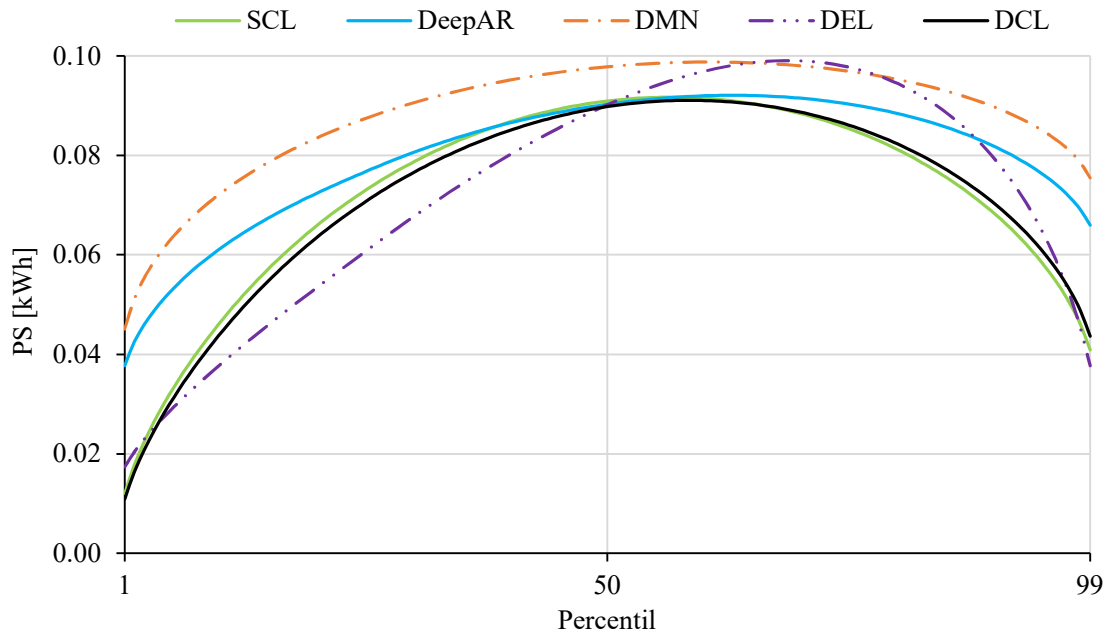


Slika 4.24 – PS na nivou potrošača iz SP2

PS na nivou percentila za potrošače iz SP1 i SP2 je prikazan na slikama 4.25 i 4.26, respektivno. Može se uočiti da je PS veći u sredini nego na krajevima u skladu sa statističkim značajem intervala između percentila. Takođe se može uočiti da je PS kod svih primenjenih rešenja nakrivljen ka gornjim percentilima. PS je nakrivljen zbog teško predvidivih iznenadnih skokova u opterećenju koji se javljaju kod pojedinačnih potrošača, što se slaže sa rezultatima iz literature [34]. Takvi skokovi se mogu primetiti na skoro svim slikama u poglavlju 4.3.2. Ako se PS posmatra na nivou percentila, onda DCL ima blagu prednost u odnosu na SCL, a znatnu prednost u odnosu na DeepAR i DMN. S druge strane, DCL ima prednost u odnosu na DEL na gornjim percentilima, dok DEL ima prednost u odnosu na DCL na donjim percentilima. Međutim, PS za DCL je manje nakrivljen ka gornjim percentilima nego PS za DEL, što implicira da se primenom DCL rešenja može predvideti više naglih skokova u opterećenju nego primenom DEL rešenja.

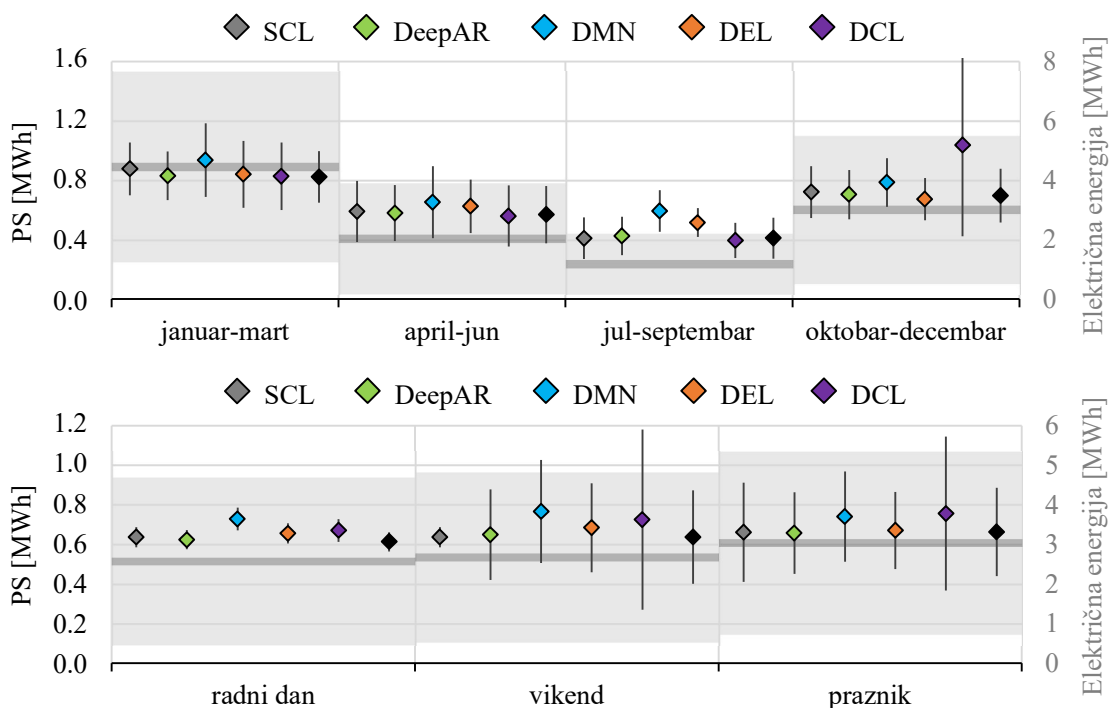


Slika 4.25 – PS na nivou percentila za potrošače iz SP1

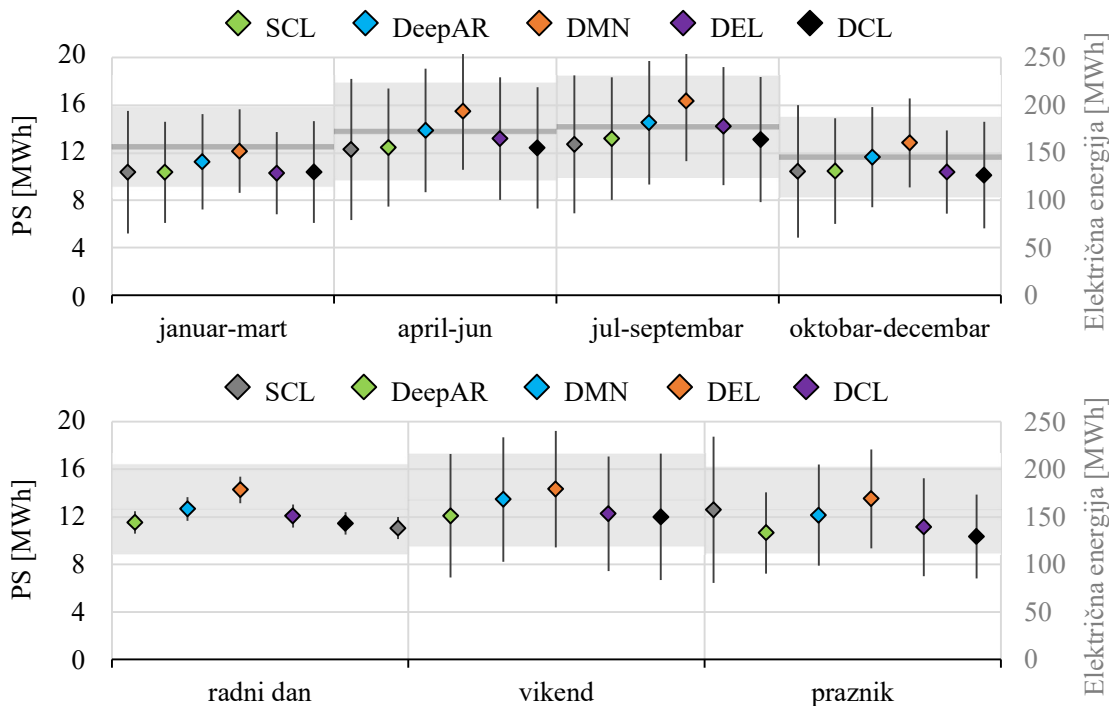


Slika 4.26 – PS na nivou percentila za potrošače iz SP2

PS i potrošena električna energija na dnevnom nivou za potrošače iz SP1 i SP2, tokom različitih meseci i tipova dana, prikazani su na slikama 4.27 i 4.28, respektivno. Prosek i standardna devijacija PS na dnevnom nivou u posmatranom periodu predstavljeni su romboidima i vertikalnim linijama, respektivno. Prosečna potrošena električna energija na dnevnom nivou u posmatranom periodu je predstavljena tamno sivom linijom, a standardna devijacija svetlo sivom površinom.



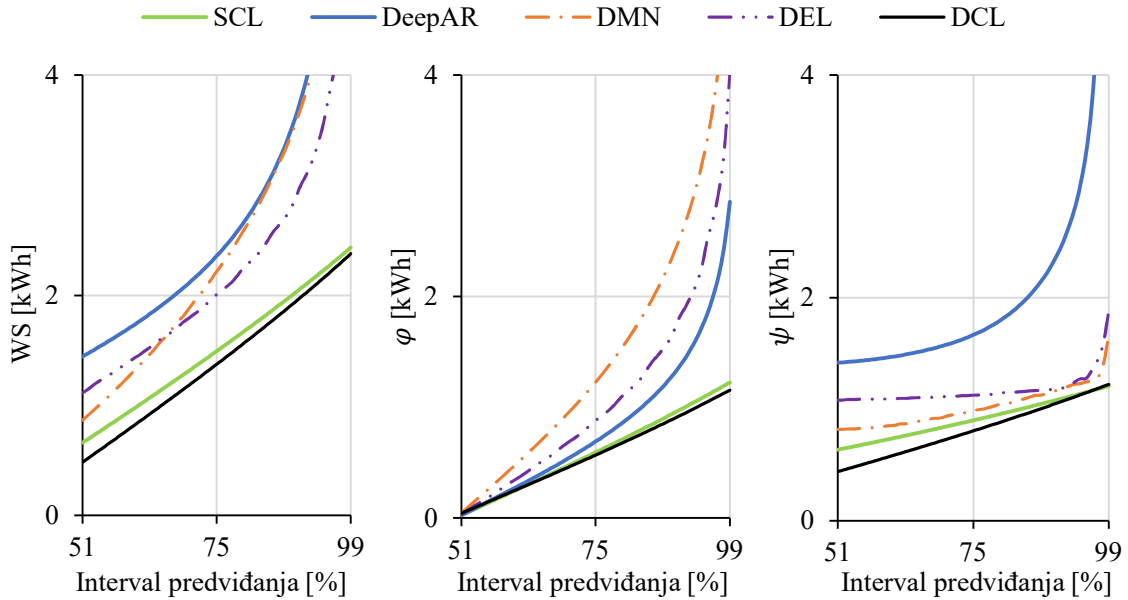
Slika 4.27 – PS i potrošena električna energija na dnevnom nivou za potrošače iz SP1



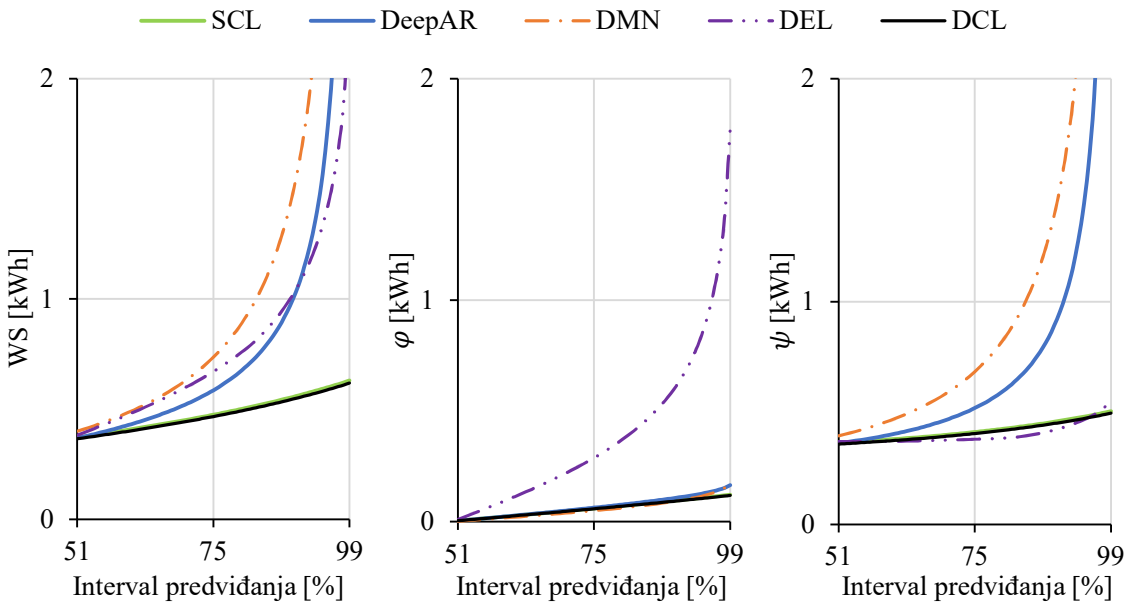
Slika 4.28 – PS i potrošena električna energija na dnevnom nivou za potrošače iz SP2

Na slikama 4.27 i 4.28 se može uočiti da prosek i standardna devijacija PS prate povećanje proseka i standardne devijacije potrošnje električne energije. Prognoza više odstupa od ostvarenih vrednosti tokom proleća i zime za potrošače iz SP1, a tokom leta i jeseni za potrošače iz SP2. Takođe, odstupanja su manja tokom radnih dana, a veća tokom vikenda i praznika, što je primetno na skoro svim slikama u poglavlju 4.3.2. Prema [4], jedan razlog za pojavu takvih rezultata je to što kod određenih potrošača s povećanjem potrošnje električne energije dolazi do teško predvidivih naglih promena u opterećenju. Pored toga, može se primetiti da je PS za DCL u proseku manji ili približno jednak kao PS za SCL, DeepAR, DMN i DEL u svim posmatranim periodima.

WS i njegove komponente, širina intervala predviđanja (φ) i penalizacija odstupanja van intervala predviđanja (ψ), za potrošače iz SP1 i SP2 prikazani su na slikama 4.29 i 4.30, respektivno. Intervali predviđanja su izraženi u procentima koji predstavljaju verovatnoću sa kojom je očekivano da prognozirani intervali obuhvate ostvarene vrednosti.



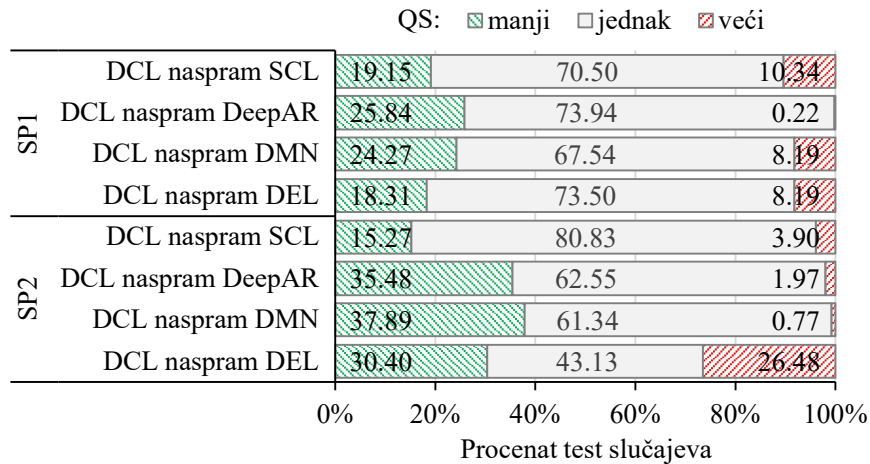
Slika 4.29 – WS, φ i ψ za potrošače iz SP1



Slika 4.30 – WS, φ i ψ za potrošače iz SP2

Na slikama 4.29 i 4.30 se može primetiti da WS, φ i ψ rastu sa porastom statističkog značaja intervala predviđanja kod svih primenjenih rešenja. Što je statistički značaj intervala predviđanja veći, to su intervali prirodno širi, da bi odstupanja van intervala bila manja. Međutim, čak i sa širim intervalima predviđanja dolazi do neočekivanih odstupanja, što rezultuje većom penalizacijom. Na taj način raste WS, što se slaže sa rezultatima iz literature [34]. Međutim, brzina kojom WS, φ i ψ rastu se razlikuje za različita rešenja. Može se primetiti da su WS, φ i ψ pretežno manji i da imaju sporiji (približno linearan) rast za DCL i SCL nego za DeepAR, DMN i DEL. Primenom DeepAR, DMN i DEL rešenja se dobijaju ili znatno širi intervali predviđanja kako bi se smanjila penalizacija ili uži intervali predviđanja koji ne obuhvataju visoke varijacije u opterećenju, pa je za njih penalizacija veća. Nasuprot njima, primenom DCL i SCL rešenja se dobijaju uži intervali predviđanja za koje je penalizacija manja jer obuhvataju više ostvarenih vrednosti. Pritom je DCL u blagoj prednosti u odnosu na SCL.

Statistički značaj prednosti koju ima DCL rešenje u odnosu na SCL, DeepAR, DMN i DEL u pogledu PS je potvrđen primenom Diebold-Mariano testa. Na slici 4.31 su prikazani rezultati testova. Po jedan test je izvršen za svaki prikazani par rešenja u oba smera, za svaki sat u horizontu prognoze i svakog potrošača. Prema tome, ukupno je izvršeno $4 \times 2 \times 24 \times (114 + 6.633) \approx 1,3 \times 10^6$ testova. Može se uočiti da postoji više test slučajeva u kojima se primenom DCL rešenja dobija značajno manji PS nego test slučajeva u kojima se primenom drugih rešenja dobija značajno veći PS. DCL ima najveću prednost u odnosu na DeepAR i DMN, a zatim u odnosu na SCL i DEL, što je u saglasnosti sa svim prethodno prikazanim rezultatima. Može se primetiti da je DEL najveći konkurent DCL rešenju kod potrošača iz SP2, gde DCL ima prednost kod ~4% više potrošača (30,40% – 26,48%). Takođe se može uočiti da postoji mnogo test slučajeva u kojima nema statistički značajne razlike između upoređenih rešenja. To je u velikoj meri posledica stohastičkih promena u opterećenju koja su teško predvidiva od strane svih rešenja, što se takođe može primetiti u prethodno prikazanim rezultatima.



Slika 4.31 – Rezultati Diebold-Mariano testova

Slike 4.21-4.31 potvrđuju da primena predloženog rešenja, koje se zasniva na prognozi centroida opterećenja primenom dubokog učenja, vodi ka povećanju tačnosti probabilističke prognoze opterećenja na nivou niskonaponskih potrošača. Primena SCL rešenja pokazuje da zamena predloženog DCLN modela sa SVM modelom u okviru predloženog rešenja smanjuje i destabilizuje tačnost prognoze, što opravdava primenu dubokog učenja. Primena DeepAR i DMN rešenja pokazuje da usmeravanje funkcije gubitka ka prognozi parametara normalnih raspodela ulaznih vrednosti tokom treniranja rezultuje većim greškama u prognozi u odnosu na direktnu prognozu centralnih momenata opterećenja predloženim pristupom. Primena DEL rešenja pokazuje da je prognoza percentila na nivou pojedinačnih potrošača manje stabilna u pogledu tačnosti nego prognoza centralnih momenata opterećenja na nivou potrošačkih grupa.

4.2.4. Zauzeće računarskih resursa

U pogledu zauzeća računarskih resursa, analizirani su zauzeće memorije (RAM) i vreme izvršavanja primenjenih rešenja. Posmatranjem zauzeća memorije tokom treniranja regresionih modela je ustanovljeno da je za primenu SCL rešenja zauzeto do 100 MB RAM, dok je za primenu DCL, DEL, DMN i DeepAR rešenja zauzet sav GPU RAM (6 GB). Prema tome, treniranje DNN modela iziskuje znatno više memorije nego treniranje SVM modela, a preostale razlike između DNN modela se ogledaju u vremenu izvršavanja.

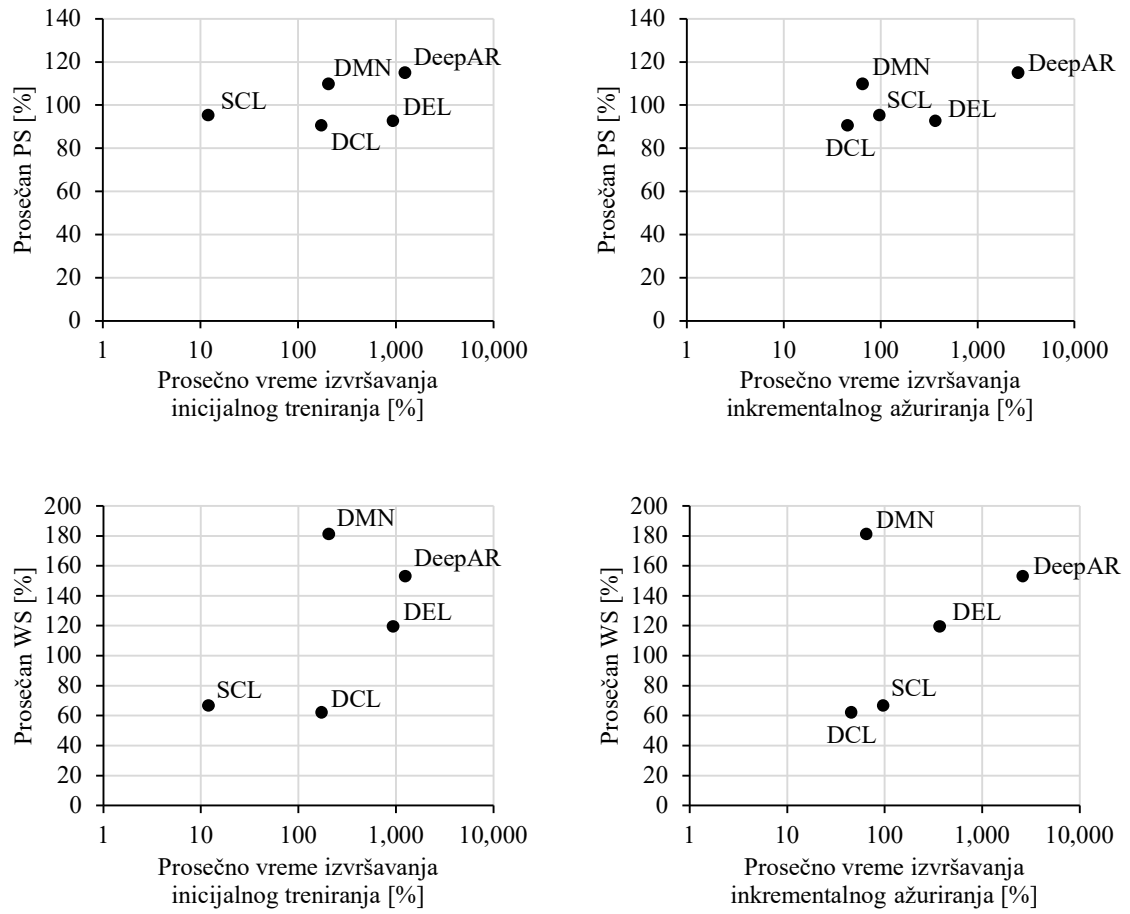
Vreme izvršavanja je prikazano u tabeli 4.6. Prikazano je ukupno vreme koje je potrebno za inicijalno treniranje regresionih modela i prosečno vreme koje je potrebno za inkrementalno ažuriranje regresionih modela. Vreme koje je potrebno za primenu (testiranje) istreniranih modela nije prikazano u tabeli 4.6 jer je zanemarljivo u odnosu na prikazana vremena izvršavanja. Može se uočiti da DCL rešenje iziskuje više vremena za inicijalno treniranje nego SCL rešenje, jer je DCLN kompleksniji regresioni model nego SVM. Međutim, inkrementalno ažuriranje težina DCLN modela iziskuje manje vremena nego ažuriranje *support* vektora SVM modela. S druge strane, za primenu DCL rešenja je potrebno manje vreme izvršavanja nego za primenu DeepAR, DMN i DEL rešenja u svim slučajevima, što je posledica visoke redukcije količine podataka (~99,9% prema tabeli 4.1). Na osnovu razlike između vrednosti koje su prikazane za SP1 i SP2 u tabeli 4.6, može se zaključiti da razlika u vremenu izvršavanja raste sa porastom broja potrošača. Primena DMN rešenja iziskuje više vremena nego primena DCL rešenja jer se funkcija gubitka tokom treniranja usmerava ka prognozi parametara mešavine normalnih raspodela ulaznih vrednosti umesto ka prognozi centroida opterećenja. Za primenu DeepAR i DEL rešenja je potrebno najviše vremena. DeepAR rešenje iziskuje približno isto vreme za inicijalno treniranje i inkrementalno ažuriranje regresionih modela zbog složenosti prilagođavanja regresionog modela novim šablonima HDO. DEL rešenje prirodno iziskuje više vremena jer se zasniva na prognozi percentila na nivou pojedinačnih potrošača.

Tabela 4.6 – Vreme izvršavanja [min]

Rešenje	SP1		SP2	
	Inicijalno treniranje	Inkrementalno ažuriranje	Inicijalno treniranje	Inkrementalno ažuriranje
SCL	4	5	20	27
DeepAR	145	118	640	527
DMN	27	4	112	18
DEL	23	4	953	162
DCL	22	3	103	15

4.2.5. Rezime rezultata

Rezultati studije slučaja pokazuju da primena predloženog rešenja vodi ka visokoj tačnosti prognoze, uz kratko vreme izvršavanja, u poređenju sa konkurentnim rešenjima iz aktuelnog stanja u oblasti. Na slici 4.32 su uporedo prikazane mere tačnosti prognoze (PS sa slika 4.21 i 4.22 i WS sa slika 4.29 i 4.30) i vreme izvršavanja (iz tabele 4.6). Radi zajedničkog prikaza tačnosti prognoze i vremena izvršavanja, svi rezultati su izraženi relativno. Relativne vrednosti su dobijene tako što je svako rešenje upoređeno sa svakim drugim rešenjem nad svakim skupom podataka, a zatim su rezultati poređenja izraženi u procentima i prosečeni. Može se primetiti da predloženo rešenje (DCL) i varijanta predloženog rešenja (SCL) nude bolju ravnotežu između tačnosti prognoze i vremena izvršavanja nego konkurentna rešenja (DEL, DMN i DeepAR). SCL nudi bolji odnos između tačnosti prognoze i vremena koje je potrebno za inicijalno treniranje regresionih modela. DCL nudi bolji odnos između tačnosti prognoze i vremena koje je potrebno za inkrementalno ažuriranje regresionih modela. DEL nudi visoku tačnost prognoze, ali uz duže vreme izvršavanja, dok za DMN važi obrnuto. DeepAR nudi nižu tačnost prognoze i duže vreme izvršavanja od svih ostalih rešenja. S obzirom da je očekivano da se inkrementalno ažuriranje u praktičnoj primeni izvodi znatno češće nego inicijalno treniranje regresionih modela, može se zaključiti da je predloženo rešenje znatno pogodnije od ostalih za primenu u DM sa većim brojem niskonaponskih potrošača.



Slika 4.32 – Vreme izvršavanja naspram tačnosti prognoze

5. ZAKLJUČAK

U ovoj doktorskoj disertaciji je adresiran problem kratkoročne probabilističke prognoze opterećenja na niskom naponu u DM. Takav vid prognoze je izazovan zbog visoke varijabilnosti opterećenja na niskom naponu i velike količine podataka o opterećenju, koji su posledica većeg broja posmatranih čvorova mreže. Problem istraživanja je adresiran razvojem novog rešenja – *Deep Centroid Learning* (DCL). Novo rešenje je razvijeno sa globalnim ciljem da uvaži varijabilnost opterećenja i ponudi konkurentnu tačnost prognoze uz visoku efikasnost sa stanovišta zauzeća računarskih resursa. Globalni cilj istraživanja je postignut kroz četiri individualna cilja na sledeći način.

Predloženo rešenje se zasniva na prognozi centroida opterećenja – vremenskih serija centralnih momenata opterećenja (proseka i standardne devijacije) na nivou potrošačkih grupa. Centralni momenti se prvo izvode na nivou potrošača, a zatim na nivou potrošačkih grupa. Za svaku potrošačku grupu se gradi po jedan regresioni model koji se zasniva na primeni dubokog učenja. Ulazni podaci u regresiji su prethodne vrednosti centroida opterećenja, vremenske odrednice i meteorološki faktori, a rezultat su prognozirane vrednosti centroida opterećenja na nivou svake grupe. Množenjem tih vrednosti sa kvantitativnim pokazateljima opterećenja, prognozirane, normalizovane vrednosti na nivou grupe se transformišu u vrednosti na nivou članova grupe, u apsolutnim jedinicama. Prema tome, konačan rezultat su centralni momenti opterećenja u apsolutnim vrednostima, na nivou potrošača, koji predstavljaju ciljanu probabilističku prognozu. Tako je postignut **prvi cilj** ove doktorske disertacije – razviti novo rešenje tako da redukuje model podataka o opterećenju na niskom naponu bez gubitka značajnih informacija o varijabilnosti opterećenja i na taj način smanji zauzeće računarskih resursa u prognozi.

Predloženo rešenje ukazuje na to da se varijabilnost opterećenja na nivou niskonaponskih potrošača može kvantifikovati probabilističkom prognozom na nivou potrošačkih grupa. Prema tome, model opterećenja se može redukovati tako da omogući primenu dubokog učenja za prognozu opterećenja na niskom naponu bez prekomernog zauzeća računarskih resursa. Zauzvrat, duboko učenje omogućava otkrivanje i najsloženijih nelinearnih veza između visoko varijabilnog opterećenja i promenljivih faktora koji na to opterećenje utiču. Primenom dubokog učenja je postignut **drugi cilj** ove doktorske disertacije – razviti novo rešenje tako da omogući primenu sofisticiranih regresionih metoda nad redukovanim modelom podataka i na taj način ponudi konkurentnu tačnost prognoze.

Efikasnost predloženog rešenja je verifikovana u studiji slučaja nad skupom realnih podataka koji su sakupljeni sa pametnih brojila, u jednoj severnoameričkoj i jednoj australijskoj DM. Tako je postignut **treći cilj** ove doktorske disertacije – verifikovati opravdanost primene predloženog rešenja nad skupom realnih podataka. Rezultat primene predloženog rešenja je visoka tačnost prognoze i kratko vreme izvršavanja u poređenju sa konkurentnim rešenjima iz aktuelnog stanja u oblasti. Na taj način je postignut i poslednji, **četvrti cilj** ove doktorske disertacije – uporediti predloženo rešenje sa konkurentnim rešenjima iz aktuelnog stanja u oblasti.

Studija slučaja pokazuje da predložen odabir i način primene statističkih metoda i metoda mašinskog (dubokog) učenja u reprezentaciji podataka (ekstrakciji i odabiru atributa), klasterizaciji i regresiji, vode ka poboljšanju tačnosti prognoze opterećenja. Primena predloženog rešenja vodi ka tačnijoj prognozi nego primena rešenja koja se ne zasnivaju na dubokom učenju ili na prognozi centroida opterećenja. Pored toga, primena dubokog učenja za direktnu prognozu opterećenja na nivou pojedinačnih potrošača iziskuje znatno veće računarske resurse. Zbog toga je predloženo rešenje pogodnije za prognozu opterećenja na niskom naponu u DM.

Prema tome, postoji mogućnost da se DCL primeni kao dodatan softverski alat u savremenom DMS. Na osnovu toga bi se podigla tačnost rešenja gotovo svih naprednih DMS funkcija za analizu, upravljanje i planiranje DM. Takav skup softverskih alata bi doprineo efikasnijem upravljanju DM na više načina: kroz pravovremene reakcije na potražnju električne energije; kroz integraciju distribuiranih energetskih izvora i kontrolu skladištenja električne energije; u dekarbonizaciji i decentralizaciji DM, itd. Kvantifikacija

varijabilnosti opterećenja na niskom naponu u vidu probabilističke prognoze koju pruža DCL omogućila bi efikasnije upravljanje rizicima, smanjenje operativnih troškova, optimalno odlučivanje i upravljanje DM.

Ova doktorska disertacija takođe otvara brojne pravce za buduća istraživanja. Neka od njih su:

1) *Detekcija anomalija*

Predloženo rešenje bi se moglo upotrebiti za detekciju anomalija u opterećenju. Anomalije u ostvarenom opterećenju na nivou potrošača bi se mogle detektovati u odnosu na prognozirano opterećenje na nivou potrošačke grupe. Primena predloženog rešenja bi na taj način mogla da pomogne DP u otkrivanju kvarova i sprečavanju krađe električne energije.

2) *Prenosno učenje*

Regresioni modeli koji su pripremljeni za prognozu centroida opterećenja u jednoj DM potencijalno bi se mogli upotrebiti u nekoj drugoj DM primenom prenosnog učenja. Predloženo rešenje bi na taj način omogućilo DP da efikasnije upravljaju DM, ili delovima DM za koje postoji manje podataka o opterećenju.

3) *Drugi sistemi*

Predloženo rešenje bi se moglo primeniti i u drugim, neelektroenergetskim sistemima koji su takođe suočeni sa velikim brojem tekućih vremenskih serija. Na primer, predloženo rešenje bi se moglo primeniti za prognozu opterećenja serverskih sistema, telekomunikacionog saobraćaja i slično.

LITERATURA

1. UNFCCC, *Paris Agreement, Decision 1/CP.21, Article 17*, Paris, 2016
2. World Meteorological Organization (WMO), *State of the global climate 2020*, Geneva, 2021
3. T. Ahmad, H. Chen, Y. Guo, J. Wang, *A comprehensive overview on the data driven and large scale based approaches for forecasting of building energy demand: A review*, Energy and Buildings, vol. 165, pp. 301–320, 2018, DOI: 10.1016/j.enbuild.2018.01.017
4. K. P. Amber, R. Ahmad, M. W. Aslam, A. Kousar, M. Usman, M. S. Khan, *Intelligent techniques for forecasting electricity consumption of buildings*, Energy, vol. 157, pp. 886–893, 2018, DOI: 10.1016/j.energy.2018.05.155
5. C. Deb, F. Zhang, J. Yang, S. E. Lee, K. W. Shah, *A review on time series forecasting techniques for building energy consumption*, Renewable and Sustainable Energy Reviews, vol. 74, pp. 902–924, 2017, DOI: 10.1016/j.rser.2017.02.085
6. N. Kovački, *Operativno planiranje rekonfiguracije distributivnih mreža primenom višekriterijumske optimizacije*, doktorska disertacija, Fakultet tehničkih nauka, Univerzitet u Novom Sadu, 2017
7. V. Krsman, *Specijalizovani algoritmi za detekciju, identifikaciju i estimaciju loših podataka u elektro-distributivnim mrežama*, doktorska disertacija, Fakultet tehničkih nauka, Univerzitet u Novom Sadu, 2017
8. V. Strezoski, D. Popović, D. Bekut, G. Švenda, *DMS - Basis for increasing of green distributed generation penetration in distribution networks*, Thermal Science, vol. 16, no. 1, pp. 189–203, 2012, DOI: 10.2298/TSCI120119071S
9. M. Radenković, J. Lukić, M. Despotović-Zrakić, A. Labus, Z. Bogdanović, *Harnessing business intelligence in smart grids: A case of the electricity market*, Computers in Industry, vol. 96, pp. 40–53, 2018, DOI: 10.1016/j.compind.2018.01.006
10. T. Wilcox, N. Jin, P. Flach, J. Thumim, *A big data platform for smart meter data analytics*, Computers in Industry, vol. 105, pp. 250–259, 2019, DOI: 10.1016/j.compind.2018.12.010
11. V. Sarfi, H. Livani, *Optimal Volt/VAR control in distribution systems with prosumer DERs*, Electric Power Systems Research, vol. 188, p. 106520, 2020, DOI: 10.1016/j.epsr.2020.106520
12. V. C. Strezoski, N. R. Vojnović, P. M. Vidović, *New bus classification and unbalanced power flow of large-scale networks with electronically interfaced energy resources*, International Transactions on Electrical Energy Systems, vol. 28, no. 3, p. e2502, 2018, DOI: 10.1002/etep.2502
13. G. Švenda, V. Strezoski, S. Kanjuh, *Real-life distribution state estimation integrated in the distribution management system*, International Transactions on Electrical Energy Systems, vol. 27, no. 5, pp. 2296–2296, 2017, DOI: 10.1002/etep.2296
14. G. Švenda, S. Kanjuh, *Automatically generated three-phase state estimation for unbalanced distribution power grids*, 2021 IEEE Power & Energy Society General Meeting – IEEE PESGM 2021, Session: AMPS Distribution System Analysis Poster Session, No. 21PESGM2111, 2021, DOI: 10.1109/PESGM46819.2021.9638188
15. G. Švenda, Z. Simendić, V. Strezoski, *Advanced voltage control integrated in DMS*, International Journal of Electrical Power & Energy Systems, vol. 43, no. 1, pp. 333–343, 2012, DOI: 10.1016/j.ijepes.2012.05.014
16. G. Švenda, Z. Simendić, *Adaptive on-load tap-changing voltage control for active distribution networks*, Electrical Engineering, vol. 104, no. 2, pp. 1041–1056, 2022, DOI: 10.1007/S00202-021-01357-8/FIGURES/15
17. T. Hong, P. Pinson, Y. Wang, R. Weron, D. Yang, H. Zareipour, *Energy forecasting: A review and outlook*, IEEE Open Access Journal of Power and Energy, vol. 7, pp. 376–388, 2020, DOI: 10.1109/oajpe.2020.3029979
18. M. M. Božić, *Kratkoročna prognoza potrošnje električne energije zasnovana na metodama veštačke inteligencije*, doktorska disertacija, Elektronski fakultet, Univerzitet u Nišu, 2014

19. Y. Wang, Q. Chen, T. Hong, C. Kang, *Review of smart meter data analytics: Applications, methodologies, and challenges*, IEEE Transactions on Smart Grid, vol. 10, no. 3, 2019, DOI: 10.1109/TSG.2018.2818167
20. M. Bourdeau, X. qiang Zhai, E. Nefzaoui, X. Guo, P. Chatellier, *Modeling and forecasting building energy consumption: A review of data-driven techniques*, Sustainable Cities and Society, vol. 48, p. 101533, 2019, DOI: 10.1016/j.scs.2019.101533
21. H. Wang, Z. Lei, X. Zhang, B. Zhou, J. Peng, *A review of deep learning for renewable energy forecasting*, Energy Conversion and Management, vol. 198, p. 111799, 2019, DOI: 10.1016/j.enconman.2019.111799
22. M. Hayn, V. Bertsch, W. Fichtner, *Electricity load profiles in Europe: The importance of household segmentation*, Energy Research & Social Science, vol. 3, pp. 30–45, 2014, DOI: 10.1016/j.erss.2014.07.002
23. Y. Hu, P. Ren, W. Luo, P. Zhan, X. Li, *Multi-resolution representation with recurrent neural networks application for streaming time series in IoT*, Computer Networks, vol. 152, pp. 114–132, 2019, DOI: 10.1016/j.comnet.2019.01.035
24. S. C. Huang, C. N. Lu, Y. L. Lo, *Evaluation of AMI and SCADA data synergy for distribution feeder modeling*, IEEE Transactions on Smart Grid, vol. 6, no. 4, pp. 1639–1647, 2015, DOI: 10.1109/TSG.2015.2408111
25. Z. Shao, F. Chao, S. L. Yang, K. le Zhou, *A review of the decomposition methodology for extracting and identifying the fluctuation characteristics in electricity demand forecasting*, Renewable and Sustainable Energy Reviews, vol. 75, pp. 123–136, 2017, DOI: 10.1016/j.rser.2016.10.056
26. Z. Lv, H. Song, P. Basanta-Val, A. Steed, M. Jo, *Next-generation big data analytics: State of the art, challenges, and future research topics*, IEEE Transactions on Industrial Informatics, vol. 13, no. 4, pp. 1891–1899, 2017, DOI: 10.1109/TII.2017.2650204
27. S. Sukhanov, R. Wu, C. Debes, A. M. Zoubir, *Dynamic pattern matching with multiple queries on large scale data streams*, Signal Processing, vol. 171, p. 107402, 2019, DOI: 10.1016/j.sigpro.2019.107402
28. P. D. Diamantoulakis, V. M. Kapinas, G. K. Karagiannidis, *Big data analytics for dynamic energy management in smart grids*, Big Data Research, vol. 2, no. 3, pp. 94–101, 2015, DOI: 10.1016/j.bdr.2015.03.003
29. R. Sahal, J. G. Breslin, M. I. Ali, *Big data and stream processing platforms for Industry 4.0 requirements mapping for a predictive maintenance use case*, Journal of Manufacturing Systems, vol. 54, pp. 138–151, 2020, DOI: 10.1016/j.jmsy.2019.11.004
30. I. Manojlović, G. Švenda, A. Erdeljan, M. Gavrić, *Time series grouping algorithm for load pattern recognition*, Computers in Industry, vol. 111, pp. 140–147, 2019, DOI: 10.1016/j.com-pind.2019.07.009
31. R. Sevlian, R. Rajagopal, *A scaling law for short term load forecasting on varying levels of aggregation*, International Journal of Electrical Power & Energy Systems, vol. 98, pp. 350–361, 2018, DOI: 10.1016/j.ijepes.2017.10.032
32. S. Haben, G. Giasemidis, F. Ziel, S. Arora, *Short term load forecasting and the effect of temperature at the low voltage level*, International Journal of Forecasting, vol. 35, no. 4, pp. 1469–1484, 2019, DOI: 10.1016/j.ijforecast.2018.10.007
33. T. Hong, S. Fan, *Probabilistic electric load forecasting: A tutorial review*, International Journal of Forecasting, vol. 32, no. 3, pp. 914–938, 2016, DOI: 10.1016/j.ijforecast.2015.11.011
34. J. Nowotarski, R. Weron, *Recent advances in electricity price forecasting: A review of probabilistic forecasting*, Renewable and Sustainable Energy Reviews, vol. 81, pp. 1548–1568, 2018, DOI: 10.1016/j.rser.2017.05.234
35. D. W. van der Meer, J. Widén, J. Munkhammar, *Review on probabilistic forecasting of photovoltaic power production and electricity consumption*, Renewable and Sustainable Energy Reviews, vol. 81, pp. 1484–1512, 2018, DOI: 10.1016/j.rser.2017.05.212

36. S. Haben, S. Arora, G. Giasemidis, M. Voss, D. Vukadinović Greetham, *Review of low voltage load forecasting: Methods, applications, and recommendations*, Applied Energy, vol. 304, p. 117798, 2021, DOI: 10.1016/j.apenergy.2021.117798
37. A. Mashlakov, T. Kuronen, L. Lensu, A. Kaarna, S. Honkapuro, *Assessing the performance of deep learning models for multivariate probabilistic energy forecasting*, Applied Energy, vol. 285, p. 116405, 2021, DOI: 10.1016/j.apenergy.2020.116405
38. J. Á. González-Ordiano, T. Mühlpfordt, E. Braun, J. Liu, H. Çakmak, U. Kühnapfel, C. Döpmeier, S. Waczowicz, T. Faulwasser, R. Mikut, V. Hagenmeyer, R. R. Appino, *Probabilistic forecasts of the distribution grid state using data-driven forecasts and probabilistic power flow*, Applied Energy, vol. 302, p. 117498, 2021, DOI: 10.1016/j.apenergy.2021.117498
39. Z. Hajirahimi, M. Khashei, *Hybrid structures in time series modeling and forecasting: A review*, Engineering Applications of Artificial Intelligence, vol. 86, pp. 83–106, 2019, DOI: 10.1016/j.engappai.2019.08.018
40. S. V. Oprea, A. Bara, *Machine learning algorithms for short-term load forecast in residential buildings using smart meters, sensors and big data solutions*, IEEE Access, vol. 7, pp. 177874–177889, 2019, DOI: 10.1109/ACCESS.2019.2958383
41. M. A. Mat Daut, M. Y. Hassan, H. Abdullah, H. A. Rahman, M. P. Abdullah, F. Hussin, *Building electrical energy consumption forecasting analysis using conventional and artificial intelligence methods: A review*, Renewable and Sustainable Energy Reviews, vol. 70, pp. 1108–1118, 2017, DOI: 10.1016/j.rser.2016.12.015
42. J. Xie, T. Hong, *GEFCom2014 probabilistic electric load forecasting: An integrated solution with forecast combination and residual simulation*, International Journal of Forecasting, vol. 32, no. 3, pp. 1012–1016, 2016, DOI: 10.1016/j.ijforecast.2015.11.005
43. I. Dimoulkas, P. Mazidi, L. Herre, *Neural networks for GEFCom2017 probabilistic load forecasting*, International Journal of Forecasting, vol. 35, no. 4, pp. 1409–1423, Oct. 2019, DOI: 10.1016/j.ijforecast.2018.09.007
44. P. Gaillard, Y. Goude, R. Nedellec, *Additive models and robust aggregation for GEFCom2014 probabilistic electric load and electricity price forecasting*, International Journal of Forecasting, vol. 32, no. 3, pp. 1038–1050, 2016, DOI: 10.1016/j.ijforecast.2015.12.001
45. J. Xie, T. Hong, *Temperature scenario generation for probabilistic load forecasting*, IEEE Transactions on Smart Grid, vol. 9, no. 3, pp. 1680–1687, 2018, DOI: 10.1109/TSG.2016.2597178
46. A. Khoshrou, E. J. Pauwels, *Short-term scenario-based probabilistic load forecasting: A data-driven approach*, Applied Energy, vol. 238, pp. 1258–1268, 2019, DOI: 10.1016/j.apenergy.2019.01.155
47. Y. Wang, Q. Chen, N. Zhang, Y. Wang, *Conditional residual modeling for probabilistic load forecasting*, IEEE Transactions on Power Systems, vol. 33, no. 6, pp. 7327–7330, 2018, DOI: 10.1109/TPWRS.2018.2868167
48. Y. Wang, G. Hug, Z. Liu, N. Zhang, *Modeling load forecast uncertainty using generative adversarial networks*, Electric Power Systems Research, vol. 189, p. 106732, 2020, DOI: 10.1016/j.epsr.2020.106732
49. B. Liu, J. Nowotarski, T. Hong, R. Weron, *Probabilistic load forecasting via quantile regression averaging on sister forecasts*, IEEE Transactions on Smart Grid, vol. 8, no. 2, pp. 730–737, 2017, DOI: 10.1109/TSG.2015.2437877
50. Z. Cao, C. Wan, Z. Zhang, F. Li, Y. Song, *Hybrid ensemble deep learning for deterministic and probabilistic low-voltage load forecasting*, IEEE Transactions on Power Systems, vol. 35, no. 3, pp. 1881–1897, 2020, DOI: 10.1109/TPWRS.2019.2946701
51. T. Li, Y. Wang, N. Zhang, *Combining probability density forecasts for power electrical loads*, IEEE Transactions on Smart Grid, vol. 11, no. 2, pp. 1679–1690, 2020, DOI: 10.1109/TSG.2019.2942024

52. Y. Wang, N. Zhang, Y. Tan, T. Hong, D. S. Kirschen, C. Kang, *Combining probabilistic load forecasts*, IEEE Transactions on Smart Grid, vol. 10, no. 4, pp. 3664–3674, 2019, DOI: 10.1109/TSG.2018.2833869
53. C. Feng, M. Sun, J. Zhang, *Reinforced deterministic and probabilistic load forecasting via Q-learning dynamic model selection*, IEEE Transactions on Smart Grid, vol. 11, no. 2, pp. 1377–1386, 2020, DOI: 10.1109/TSG.2019.2937338
54. S. Zhang, Y. Wang, Y. Zhang, D. Wang, N. Zhang, *Load probability density forecasting by transforming and combining quantile forecasts*, Applied Energy, vol. 277, p. 115600, 2020, DOI: 10.1016/j.apenergy.2020.115600
55. G. Marcjasz, B. Uniejewski, R. Weron, *Probabilistic electricity price forecasting with NARX networks: Combine point or probabilistic forecasts?*, International Journal of Forecasting, vol. 36, no. 2, pp. 466–479, 2020, DOI: 10.1016/j.ijforecast.2019.07.002
56. A. J. Cannon, *Quantile regression neural networks: Implementation in R and application to precipitation downscaling*, Computers and Geosciences, vol. 37, no. 9, pp. 1277–1284, 2011, DOI: 10.1016/j.cageo.2010.07.005
57. Y. Wang, D. Gan, M. Sun, N. Zhang, Z. Lu, C. Kang, *Probabilistic individual load forecasting using pinball loss guided LSTM*, Applied Energy, vol. 235, pp. 10–20, 2019, DOI: 10.1016/j.apenergy.2018.10.078
58. M. Afrasiabi, M. Mohammadi, M. Rastegar, L. Stankovic, S. Afrasiabi, M. Khazaei, *Deep-based conditional probability density function forecasting of residential loads*, IEEE Transactions on Smart Grid, vol. 11, no. 4, pp. 3646–3657, 2020, DOI: 10.1109/TSG.2020.2972513
59. K. Bandara, C. Bergmeir, S. Smyl, *Forecasting across time series databases using recurrent neural networks on groups of similar series: A clustering approach*, Expert Systems with Applications, vol. 140, p. 112896, 2020, DOI: 10.1016/j.eswa.2019.112896
60. D. Salinas, V. Flunkert, J. Gasthaus, T. Januschowski, *DeepAR: Probabilistic forecasting with autoregressive recurrent networks*, International Journal of Forecasting, vol. 36, no. 3, pp. 1181–1191, 2020, DOI: 10.1016/j.ijforecast.2019.07.001
61. Y. Chen, Y. Kang, Y. Chen, Z. Wang, *Probabilistic forecasting with temporal convolutional neural network*, Neurocomputing, vol. 399, pp. 491–501, 2020, DOI: 10.1016/j.neucom.2020.03.011
62. Y. Yang, W. Hong, S. Li, *Deep ensemble learning based probabilistic load forecasting in smart grids*, Energy, vol. 189, p. 116324, 2019, DOI: 10.1016/j.energy.2019.116324
63. S. Haben, S. Arora, G. Giasemidis, M. Voss, D. Vukadinović Greetham, *Review of low voltage load forecasting: Methods, applications, and recommendations*, Applied Energy, vol. 304, p. 117798, Dec. 2021, DOI: 10.1016/J.APENERGY.2021.117798
64. UMass trace repository, *UMass Smart* dataset - 2017 release*, 2017, <http://traces.cs.umass.edu/index.php/Smart/Smart> (22.06.2022.)
65. Australian Government, *Smart-Grid Smart-City customer trial data*, 2018, <https://data.gov.au/data/dataset/smart-grid-smart-city-customer-trial-data> (22.06.2022.)
66. C. Miller, Z. Nagy, A. Schlueter, *A review of unsupervised statistical learning and visual analytics techniques applied to performance analysis of non-residential buildings*, Renewable and Sustainable Energy Reviews, vol. 81, no. 1, pp. 1365–1377, 2018, DOI: 10.1016/j.rser.2017.05.124
67. I. Manojlović, G. Švenda, A. Erdeljan, *Load pattern recognition method for probabilistic short-term load forecasting at low voltage level*, 2022 IEEE PES Innovative Smart Grid Technologies Conference Europe (ISGT-Europe), Session 20: Industry Application – Part I, Paper No. PAS-20-4
68. M. Q. Raza, A. Khosravi, *A review on artificial intelligence based load demand forecasting techniques for smart grid and buildings*, Renewable and Sustainable Energy Reviews, vol. 50, pp. 1352–1372, 2015, DOI: 10.1016/j.rser.2015.04.065
69. A. M. Tureczek, P. S. Nielsen, *Structured literature review of electricity consumption classification using smart meter data*, Energies, vol. 10, no. 5, pp. 1–19, 2017, DOI: 10.3390/en10050584

70. S. Aghabozorgi, A. Seyed Shirkorshidi, T. Ying Wah, *Time-series clustering - A decade review*, Information Systems, vol. 53, pp. 16–38, 2015, DOI: 10.1016/j.is.2015.04.007
71. S. K. Jensen, T. B. Pedersen, C. Thomsen, *Time series management systems: A survey*, IEEE Transactions on Knowledge and Data Engineering, vol. 29, no. 11, pp. 2581–2600, 2017, DOI: 10.1109/TKDE.2017.2740932
72. I. Manojlović, G. Švenda, A. Erdeljan, M. Gavrić, D. Čapko, *Hierarchical multiresolution representation of streaming time series*, Big Data Research, vol. 26, p. 100256, 2021, DOI: 10.1016/j.bdr.2021.100256
73. E. J. Keogh, M. J. Pazzani, *An enhanced representation of time series which allows fast and accurate classification, clustering and relevance feedback*, KDD'98 - Proceedings of the 4th International Conference of Knowledge Discovery and Data Mining, 1998, pp. 239–241
74. D. Ruta, L. Cen, E. Damiani, *Fast summarization and anonymization of multivariate big time series*, Proceedings - 2015 IEEE International Conference on Big Data, IEEE Big Data 2015, 2015, pp. 1901–1904, DOI: 10.1109/BigData.2015.7363965
75. Y. Hu, P. Guan, P. Zhan, Y. Ding, X. Li, *A novel segmentation and representation approach for streaming time series*, IEEE Access, vol. 7, pp. 184423–184437, 2018, DOI: 10.1109/ACCESS.2018.2828320
76. R. P. Silveira, J. O. Trierweiler, M. Farenzena, H. C. Teixeira, *Systematic approaches for PI systemTM data compression tuning*, IFAC Proceedings Volumes, 2012, vol. 45, no. 15, pp. 309–313, DOI: 10.3182/20120710-4-SG-2026.00137
77. L. Duan, F. Yu, W. Pedrycz, X. Wang, X. Yang, *Time-series clustering based on linear fuzzy information granules*, Applied Soft Computing Journal, vol. 73, pp. 1053–1067, 2018, DOI: 10.1016/j.a-soc.2018.09.032
78. J. Lin, E. Keogh, L. Wei, S. Lonardi, *Experiencing SAX: A novel symbolic representation of time series*, Data Mining and Knowledge Discovery, vol. 15, no. 2, pp. 107–144, 2007, DOI: 10.1007/s10618-007-0064-z
79. S. Karimi-Bidhendi, F. Munshi, A. Munshi, *Scalable classification of univariate and multivariate time series*, Proceedings - 2018 IEEE International Conference on Big Data, Big Data 2018, 2019, pp. 1598–1605, DOI: 10.1109/BigData.2018.8621889
80. C. Zhang, Y. Chen, A. Yin, *Anomaly subsequence detection with dynamic local density for time series*, Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), 2019, vol. 11707, pp. 291–305
81. E. Keogh, K. Chakrabarti, M. Pazzani, S. Mehrotra, *Dimensionality reduction for fast similarity search in large time series databases*, Knowledge and Information Systems, vol. 3, no. 3, pp. 263–286, 2001, DOI: 10.1007/PL00011669
82. H. Li, C. Guo, *Piecewise cloud approximation for time series mining*, Knowledge-Based Systems, vol. 24, no. 4, pp. 492–500, 2011, DOI: 10.1016/j.knsys.2010.12.008
83. G. Si, K. Zheng, Z. Zhou, C. Pan, X. Xu, K. Qu, Y. Zhang, *Three-dimensional piecewise cloud representation for time series data mining*, Neurocomputing, vol. 316, pp. 78–94, 2018, DOI: 10.1016/j.neucom.2018.07.053
84. D. Li, *Knowledge representation in KDD based on linguistic atoms*, Journal of Computer Science and Technology, vol. 12, no. 6, pp. 481–496, 1997, DOI: 10.1007/BF02947201
85. Q. Cai, L. Chen, J. Sun, *Piecewise statistic approximation based similarity measure for time series*, Knowledge-Based Systems, vol. 85, pp. 181–195, 2015, DOI: 10.1016/j.knsys.2015.05.005
86. E. Spiliotis, F. Petropoulos, N. Kourentzes, V. Assimakopoulos, *Cross-temporal aggregation: Improving the forecast accuracy of hierarchical electricity consumption*, Applied Energy, vol. 261, p. 114339, 2020, DOI: 10.1016/j.apenergy.2019.114339
87. P. Nystrup, E. Lindström, P. Pinson, H. Madsen, *Temporal hierarchies with autocorrelation for load forecasting*, European Journal of Operational Research, vol. 280, no. 3, pp. 876–888, 2020, DOI: 10.1016/j.ejor.2019.07.061

88. M. Naimur Rahman, A. Esmailpour, J. Zhao, *Machine learning with big data an efficient electricity generation forecasting system*, Big Data Research, vol. 5, pp. 9–15, 2016, DOI: 10.1016/j.bdr.2016.02.002
89. R. Li, F. Li, N. D. Smith, *Multi-resolution load profile clustering for smart metering data*, IEEE Transactions on Power Systems, vol. 31, no. 6, pp. 4473–4482, 2016, DOI: 10.1109/TPWRS.2016.2536781
90. A. Llusà Serra, S. Vila-Marta, T. Escobet Canal, *Formalism for a multiresolution time series database model*, Information Systems, vol. 56, pp. 19–35, 2016, DOI: 10.1016/J.IS.2015.08.006
91. A. Cuzzocrea, D. Saccà, *Exploiting compression and approximation paradigms for effective and efficient online analytical processing over sensor network readings in data grid environments*, Concurrency and Computation: Practice and Experience, vol. 25, no. 14, pp. 2016–2035, 2013, DOI: 10.1002/cpe.2982
92. I. Manojlović, A. Erdeljan, *Efficient aggregation of time series data*, ICIST 2017 Proceedings, 2017, vol. 1, pp. 102–107
93. I. Manojlović, *Algoritmi za agregaciju vremenskih serija u realnom vremenu*, specijalistički rad, Fakultet tehničkih nauka, Univerzitet u Novom Sadu, 2018, DOI: 10.13140/RG.2.2.16164.94087
94. H. Nakanishi, M. Ohsuna, M. Kojima, S. Imazu, M. Nonomura, M. Emoto, M. Yoshida, C. Iwata, K. Ida, *Real-time data streaming and storing structure for the LHD's fusion plasma experiments*, IEEE Transactions on Nuclear Science, vol. 63, no. 1, pp. 222–227, 2016, DOI: 10.1109/TNS.2016.2515099
95. X. Gong, S. Fong, Y. W. Si, *Fast multi-subsequence monitoring on streaming time-series based on Forward-propagation*, Information Sciences, vol. 450, pp. 73–88, 2018, DOI: 10.1016/j.ins.2018.03.023
96. S. di Martino, L. Fiadone, A. Peron, A. Riccabone, V. N. Vitale, *Industrial internet of things: Persistence for time series with NoSQL databases*, 2019 IEEE 28th International Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE), 2019, pp. 340–345, DOI: 10.1109/WETICE.2019.00076
97. A. Struckov, S. Yufa, A. A. Visheratin, D. Nasonov, *Evaluation of modern tools and techniques for storing time-series data*, Procedia Computer Science, vol. 156, pp. 19–28, 2019, DOI: 10.1016/j.procs.2019.08.125
98. A. Stefanović, *Agregacija vremenskih serija u oblasti elektroenergetskih sistema na različitim tipovima OLAP servera*, master rad, Fakultet tehničkih nauka, Univerzitet u Novom Sadu, 2018
99. A. MacDonald, *PhilDB: the time series database with built-in change logging*, PeerJ Computer Science, vol. 2, p. e52, 2016, DOI: 10.7717/peerj-cs.52
100. B. Chardin, J.-M. Lacombe, J.-M. Petit, *Chronos: a NoSQL system on flash memory for industrial process data*, Distributed and Parallel Databases, vol. 34, no. 3, pp. 293–319, 2016, DOI: 10.1007/s10619-015-7175-0
101. C. Li, B. Li, M. Z. A. Bhuiyan, L. Wang, J. Si, G. Wei, J. Li, *FluteDB: An efficient and scalable in-memory time series database for sensor-cloud*, Journal of Parallel and Distributed Computing, vol. 122, pp. 95–108, 2018, DOI: 10.1016/j.jpdc.2018.07.021
102. J. R. Sutton, R. Mahajan, O. Akbilgic, R. Kamaleswaran, *PhysOnline: An open source machine learning pipeline for real-time analysis of streaming physiological waveform*, IEEE Journal of Biomedical and Health Informatics, vol. 23, no. 1, pp. 59–65, 2019, DOI: 10.1109/JBHI.2018.2832610
103. R. Vautard, P. Yiou, M. Ghil, *Singular-spectrum analysis: A toolkit for short, noisy chaotic signals*, Physica D: Nonlinear Phenomena, vol. 58, no. 1–4, pp. 95–126, 1992, DOI: 10.1016/0167-2789(92)90103-T
104. R. Agrawal, C. Faloutsos, A. Swami, *Efficient similarity search in sequence databases*, Foundations of Data Organization and Algorithms, 1993, pp. 69–84, DOI: 10.1007/3-540-57301-1_5

105. Kin-Pong Chan, Ada Wai-Chee Fu, *Efficient time series matching by wavelets*, Proceedings 15th International Conference on Data Engineering, 1999, pp. 126–133, DOI: 10.1109/ICDE.1999.754915
106. N. E. Huang, Z. Shen, S. R. Long, M. C. Wu, H. H. Shih, Q. Zheng, N.-C. Yen, C. C. Tung, H. H. Liu, *The empirical mode decomposition and the Hilbert spectrum for nonlinear and non-stationary time series analysis*, Proceedings of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences, vol. 454, no. 1971, pp. 903–995, 1998, DOI: 10.1098/rspa.1998.0193
107. K. Dragomiretskiy, D. Zosso, *Variational mode decomposition*, IEEE Transactions on Signal Processing, vol. 62, no. 3, pp. 531–544, 2014, DOI: 10.1109/TSP.2013.2288675
108. F. He, J. Zhou, Z. kai Feng, G. Liu, Y. Yang, *A hybrid short-term load forecasting model based on variational mode decomposition and long short-term memory networks considering relevant factors with Bayesian optimization algorithm*, Applied Energy, vol. 237, pp. 103–116, 2019, DOI: 10.1016/j.apenergy.2019.01.055
109. Z. Zhang, W. C. Hong, J. Li, *Electric load forecasting by hybrid self-recurrent support vector regression model with variational mode decomposition and improved cuckoo search algorithm*, IEEE Access, vol. 8, pp. 14642–14658, 2020, DOI: 10.1109/aACCESS.2020.2966712
110. N. Obrenović, G. Vidaković, I. Luković, *The choice of metric for clustering of electrical power distribution consumers*, Data Science – Analytics and Applications, pp. 71–76, 2017, DOI: 10.1007/978-3-658-19287-7_10
111. M. Bannasar, Y. Hicks, R. Setchi, *Feature selection using joint mutual information maximisation*, Expert Systems with Applications, vol. 42, no. 22, pp. 8520–8532, 2015, DOI: 10.1016/j.eswa.2015.07.007
112. L. Zhang, J. Wen, *A systematic feature selection procedure for short-term data-driven building energy forecasting model development*, Energy and Buildings, vol. 183, pp. 428–442, 2019, DOI: 10.1016/j.enbuild.2018.11.010
113. U. M. Khaire, R. Dhanalakshmi, *Stability of feature selection algorithm: A review*, Journal of King Saud University - Computer and Information Sciences, 2019, DOI: 10.1016/j.jksuci.2019.06.012
114. S. Ilić, *Kratkoročno predviđanje potrošnje električne energije u velikim elektroenergetskim sistemima*, doktorska disertacija, Fakultet tehničkih nauka, Univerzitet u Novom Sadu, 2013
115. X. Liu, Z. Zhang, Z. Song, *A comparative study of the data-driven day-ahead hourly provincial load forecasting methods: From classical data mining to deep learning*, Renewable and Sustainable Energy Reviews, vol. 119, p. 109632, 2020, DOI: 10.1016/j.rser.2019.109632
116. G. Box, G. Jenkins, G. Reinsel, G. Ljung, *Time series analysis, forecasting and control*, 1995
117. A. Yang, W. Li, X. Yang, *Short-term electricity load forecasting based on feature selection and least squares support vector machines*, Knowledge-Based Systems, vol. 163, pp. 159–173, 2019, DOI: 10.1016/j.knosys.2018.08.027
118. H. Peng, F. Long, C. Ding, *Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy*, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 27, no. 8, pp. 1226–1238, 2005, DOI: 10.1109/TPAMI.2005.159
119. Y. Liang, D. Niu, W. C. Hong, *Short term load forecasting based on feature extraction and improved general regression neural network model*, Energy, vol. 166, pp. 653–663, 2019, DOI: 10.1016/j.energy.2018.10.119
120. J. Paparrizos, L. Gravano, *Fast and accurate time-series clustering*, ACM Transactions on Database Systems, vol. 42, no. 2, pp. 1–49, 2017, DOI: 10.1145/3044711
121. J. Yang, C. Ning, C. Deb, F. Zhang, D. Cheong, S. E. Lee, C. Sekhar, K. W. Tham, *k-Shape clustering algorithm for building energy usage patterns analysis and forecasting model accuracy improvement*, Energy and Buildings, vol. 146, pp. 27–37, 2017, DOI: 10.1016/j.enbuild.2017.03.071
122. O. Arbelaitz, I. Gurrutxaga, J. Muguerza, M. J. Pérez, I. Perona, *An extensive comparative study of cluster validity indices*, Pattern Recognition, vol. 46, no. 1, pp. 243–256, 2013, DOI: 10.1016/j.patcog.2012.07.021

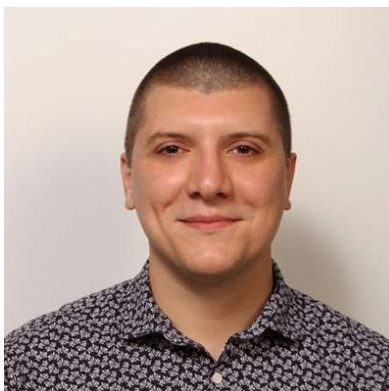
123. P. J. Rousseeuw, *Silhouettes: A graphical aid to the interpretation and validation of cluster analysis*, Journal of Computational and Applied Mathematics, vol. 20, pp. 53–65, 1987, DOI: 10.1016/0377-0427(87)90125-7
124. T. Calinski, J. Harabasz, *A dendrite method for cluster analysis*, Communications in Statistics - Theory and Methods, vol. 3, no. 1, pp. 1–27, 1974, DOI: 10.1080/03610927408827101
125. D. L. Davies, D. W. Bouldin, *A cluster separation measure*, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. PAMI-1, no. 2, pp. 224–227, 1979, DOI: 10.1109/TPAMI.1979.4766909
126. S. Zhou, Z. Xu, F. Liu, *Method for determining the optimal number of clusters based on agglomerative hierarchical clustering*, IEEE Transactions on Neural Networks and Learning Systems, vol. 28, no. 12, pp. 3007–3017, 2017, DOI: 10.1109/TNNLS.2016.2608001
127. G. Chicco, R. Napoli, F. Piglionone, *Comparisons among clustering techniques for electricity customer classification*, IEEE Transactions on Power Systems, vol. 21, no. 2, pp. 933–940, 2006, DOI: 10.1109/TPWRS.2006.873122
128. A. Mutanen, M. Ruska, S. Repo, P. Jarventausta, *Customer classification and load profiling method for distribution systems*, IEEE Transactions on Power Delivery, vol. 26, no. 3, pp. 1755–1763, 2011, DOI: 10.1109/TPWRD.2011.2142198
129. V. Figueiredo, F. Rodrigues, Z. Vale, J. B. Gouveia, *An electric energy consumer characterization framework based on data mining techniques*, IEEE Transactions on Power Systems, vol. 20, no. 2, pp. 596–602, 2005, DOI: 10.1109/TPWRS.2005.846234
130. N. Slonim, E. Aharoni, K. Crammer, *Hartigan's k-means versus Lloyd's k-means: Is it time for a change?*, IJCAI '13 Proceedings of the Twenty-Third international joint conference on Artificial Intelligence, 2013, pp. 1677–1684
131. D. Arthur, S. Vassilvitskii, *k-means++: The advantages of careful seeding*, SODA '07 - Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms, 2007, pp. 1027–1035
132. K. Murugesan, J. Zhang, *Hybrid bisect k-means clustering algorithm*, 2011 International Conference on Business Computing and Global Informatization, 2011, pp. 216–219, DOI: 10.1109/BCGIIn.2011.62
133. Y.-I. Kim, J.-M. Ko, J.-J. Song, H. Choi, *Repeated clustering to improve the discrimination of typical daily load profile*, Journal of Electrical Engineering and Technology, vol. 7, no. 3, pp. 281–287, 2012, DOI: 10.5370/JEET.2012.7.3.281
134. L. Jin, D. Lee, A. Sim, S. Borgeson, K. Wu, A. C. Spurlock, A. Todd, *Comparison of clustering techniques for residential energy behavior using smart meter data*, AAAI-17 Workshop on Artificial Intelligence for Smart Grids and Smart Buildings, 2017, pp. 260–266
135. H. Chipman, R. Tibshirani, *Hybrid hierarchical clustering with applications to microarray data*, Biostatistics, vol. 7, no. 2, pp. 286–301, 2006, DOI: 10.1093/biostatistics/kxj007
136. A. Bouguettaya, Q. Yu, X. Liu, X. Zhou, A. Song, *Efficient agglomerative hierarchical clustering*, Expert Systems with Applications, vol. 42, no. 5, pp. 2785–2797, 2015, DOI: 10.1016/j.eswa.2014.09.054
137. W. Zhu, F. Nie, X. Li, *Fast spectral clustering with efficient large graph construction*, 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2017, pp. 2492–2496, DOI: 10.1109/ICASSP.2017.7952605
138. T.-S. Xu, H.-D. Chiang, G.-Y. Liu, C.-W. Tan, *Hierarchical k-means method for clustering large-scale advanced metering infrastructure data*, IEEE Transactions on Power Delivery, vol. 32, no. 2, pp. 609–616, 2017, DOI: 10.1109/TPWRD.2015.2479941
139. G. Chicco, *Clustering methods for electrical load pattern classification*, The Scientific Bulletin of Electrical Engineering Faculty, pp. 5–13, 2010
140. G. Chicco, R. Napoli, F. Piglionone, P. Postolache, M. Scutariu, C. Toader, *Load pattern-based classification of electricity customers*, IEEE Transactions on Power Systems, vol. 19, no. 2, pp. 1232–1239, 2004, DOI: 10.1109/TPWRS.2004.826810

141. G. Chicco, I.-S. Ilie, *Support vector clustering of electrical load pattern data*, IEEE Transactions on Power Systems, vol. 24, no. 3, pp. 1619–1628, 2009, DOI: 10.1109/TPWRS.2009.2023009
142. G. Chicco, J. Sumaili Akilimali, *Renyi entropy-based classification of daily electrical load patterns*, IET Generation, Transmission & Distribution, vol. 4, no. 6, p. 736, 2010, DOI: 10.1049/iet-gtd.2009.0161
143. Y. Wei, X. Zhang, Y. Shi, L. Xia, S. Pan, J. Wu, M. Han, X. Zhao, *A review of data-driven approaches for prediction and classification of building energy consumption*, Renewable and Sustainable Energy Reviews, vol. 82, pp. 1027–1047, 2018, DOI: 10.1016/j.rser.2017.09.108
144. S. Aslam, H. Herodotou, S. M. Mohsin, N. Javaid, N. Ashraf, S. Aslam, *A survey on deep learning methods for power load and renewable energy forecasting in smart microgrids*, Renewable and Sustainable Energy Reviews, vol. 144, p. 110992, 2021, DOI: 10.1016/j.rser.2021.110992
145. C. Fan, J. Wang, W. Gang, S. Li, *Assessment of deep recurrent neural network-based strategies for short-term building energy predictions*, Applied Energy, vol. 236, pp. 700–710, 2019, DOI: 10.1016/j.apenergy.2018.12.004
146. I. F. Kao, J. Y. Liou, M. H. Lee, F. J. Chang, *Fusing stacked autoencoder and long short-term memory for regional multistep-ahead flood inundation forecasts*, Journal of Hydrology, vol. 598, p. 126371, 2021, DOI: 10.1016/j.jhydrol.2021.126371
147. X. Tang, Y. Dai, T. Wang, Y. Chen, *Short-term power load forecasting based on multi-layer bidirectional recurrent neural network*, IET Generation, Transmission and Distribution, vol. 13, no. 17, pp. 3847–3854, 2019, DOI: 10.1049/iet-gtd.2018.6687
148. J. Struye, S. Latré, *Hierarchical temporal memory and recurrent neural networks for time series prediction: An empirical validation and reduction to multilayer perceptrons*, Neurocomputing, vol. 396, pp. 291–301, 2020, DOI: 10.1016/j.neucom.2018.09.098
149. L. Sehovac, K. Grolinger, *Deep learning for load forecasting: Sequence to sequence recurrent neural networks with attention*, IEEE Access, vol. 8, pp. 36411–36426, 2020, DOI: 10.1109/ACCESS.2020.2975738
150. T. Y. Kim, S. B. Cho, *Predicting residential energy consumption using CNN-LSTM neural networks*, Energy, vol. 182, pp. 72–81, 2019, DOI: 10.1016/j.energy.2019.05.230
151. C. Tian, J. Ma, C. Zhang, P. Zhan, C. Tian, J. Ma, C. Zhang, P. Zhan, *A deep neural network model for short-term load forecast based on long short-term memory network and convolutional neural network*, Energies, vol. 11, no. 12, p. 3493, 2018, DOI: 10.3390/en1123493
152. M. Alhussein, K. Aurangzeb, S. I. Haider, *Hybrid CNN-LSTM model for short-term individual household load forecasting*, IEEE Access, vol. 8, pp. 180544–180557, 2020, DOI: 10.1109/ACCESS.2020.3028281
153. C. M. Bishop, *Mixture density networks*, Neural Computing Research Group Report: NCRG/94/004, Aston University, Birmingham, UK, 1994
154. E. N. Osegi, *Using the hierarchical temporal memory spatial pooler for short-term forecasting of electrical load time series*, Applied Computing and Informatics, 2018, DOI: 10.1016/j.aci.2018.09.002
155. L. Goel, *An extensive review of computational intelligence-based optimization algorithms: Trends and applications*, Soft Computing, pp. 1–31, 2020, DOI: 10.1007/s00500-020-04958-w
156. D. P. Kingma, J. L. Ba, *Adam: A method for stochastic optimization*, 3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings, 2015
157. J. H. Holland, *Adaptation in natural and artificial systems: An introductory analysis with applications to biology, control, and artificial intelligence*, MIT Press, 1992
158. R. Storn, K. Price, *Differential evolution - A simple and efficient heuristic for global optimization over continuous spaces*, Journal of Global Optimization, vol. 11, no. 4, pp. 341–359, 1997, DOI: 10.1023/A:1008202821328
159. J. Kennedy, R. Eberhart, *Particle swarm optimization*, Proceedings of ICNN'95 - International Conference on Neural Networks, 1995, vol. 4, no. 6, pp. 1942–1948, DOI: 10.1109/ICNN.1995.488968

160. R. Khalid, N. Javaid, *A survey on hyperparameters optimization algorithms of forecasting models in smart grid*, Sustainable Cities and Society, vol. 61, p. 102275, 2020, DOI: 10.1016/j.scs.2020.102275
161. W. T. Pan, *A new fruit fly optimization algorithm: Taking the financial distress model as an example*, Knowledge-Based Systems, vol. 26, pp. 69–74, 2012, DOI: 10.1016/j.knosys.2011.07.001
162. S. Saremi, S. Mirjalili, A. Lewis, *Grasshopper optimisation algorithm: Theory and application*, Advances in Engineering Software, vol. 105, pp. 30–47, 2017, DOI: 10.1016/j.advengsoft.2017.01.004
163. M. Barman, N. B. Dev Choudhury, S. Sutradhar, *A regional hybrid GOA-SVM model based on similar day approach for short-term load forecasting in Assam, India*, Energy, vol. 145, pp. 710–720, 2018, DOI: 10.1016/j.energy.2017.12.156
164. M. Talaat, M. A. Farahat, N. Mansour, A. Y. Hatata, *Load forecasting based on grasshopper optimization and a multilayer feed-forward neural network using regressive approach*, Energy, vol. 196, p. 117087, 2020, DOI: 10.1016/j.energy.2020.117087
165. S. Mirjalili, S. M. Mirjalili, A. Lewis, *Grey wolf optimizer*, Advances in Engineering Software, vol. 69, pp. 46–61, 2014, DOI: 10.1016/j.advengsoft.2013.12.007
166. S. Mirjalili, A. Lewis, *The whale optimization algorithm*, Advances in Engineering Software, vol. 95, pp. 51–67, 2016, DOI: 10.1016/j.advengsoft.2016.01.008
167. W. Sun, C. Zhang, *Analysis and forecasting of the carbon price using multi-resolution singular value decomposition and extreme learning machine optimized by adaptive whale optimization algorithm*, Applied Energy, vol. 231, pp. 1354–1371, 2018, DOI: 10.1016/j.apenergy.2018.09.118
168. F. S. Gharehchopogh, H. Gholizadeh, *A comprehensive survey: Whale optimization algorithm and its applications*, Swarm and Evolutionary Computation, vol. 48, pp. 1–24, 2019, DOI: 10.1016/j.swevo.2019.03.004
169. F. X. Diebold, R. S. Mariano, *Comparing predictive accuracy*, Journal of Business and Economic Statistics, vol. 13, pp. 253–265, 1995
170. W. Zhang, H. Quan, D. Srinivasan, *Parallel and reliable probabilistic load forecasting via quantile regression forest and quantile determination*, Energy, vol. 160, pp. 810–819, 2018, DOI: 10.1016/j.energy.2018.07.019
171. W. Zhang, H. Quan, D. Srinivasan, *An improved quantile regression neural network for probabilistic load forecasting*, IEEE Transactions on Smart Grid, vol. 10, no. 4, pp. 4425–4434, 2019, DOI: 10.1109/TSG.2018.2859749
172. W. Zhang, H. Quan, O. Gandhi, R. Rajagopal, C. W. Tan, D. Srinivasan, *Improving probabilistic load forecasting using quantile regression NN with skip connections*, IEEE Transactions on Smart Grid, vol. 11, no. 6, pp. 5442–5450, 2020, DOI: 10.1109/TSG.2020.2995777
173. Z. Xie, H. Hu, Q. Wang, R. Li, *ALGeNet: Adaptive Log-Euclidean Gaussian embedding network for time series forecasting*, Neurocomputing, vol. 423, pp. 353–361, 2021, DOI: 10.1016/j.neucom.2020.11.001
174. A. Brusaferrri, M. Matteucci, S. Spinelli, A. Vitali, *Probabilistic electric load forecasting through Bayesian Mixture Density Networks*, Applied Energy, vol. 309, p. 118341, 2022, DOI: 10.1016/j.apenergy.2021.118341
175. K. Aurangzeb, M. Alhussein, K. Javaid, S. I. Haider, *A pyramid-CNN based deep learning model for power load forecasting of similar-profile energy customers based on clustering*, IEEE Access, vol. 9, pp. 14992–15003, 2021, DOI: 10.1109/ACCESS.2021.3053069
176. Y. Wang, L. von Krannichfeldt, G. Hug, *Probabilistic aggregated load forecasting with fine-grained smart meter data*, 2021 IEEE Madrid PowerTech, PowerTech 2021 - Conference Proceedings, 2021, DOI: 10.1109/POWERTECH46648.2021.9494815
177. D. Kiruthiga, V. Manikandan, *Time series load forecasting using multitask deep neural network*, 2021 IEEE 2nd International Conference on Control, Measurement and Instrumentation, CMI 2021 - Proceedings, pp. 166–171, 2021, DOI: 10.1109/CMI50323.2021.9362936

178. S. Kuzmanovic, G. Švenda, Z. Ovcin, *Practical statistical methods in distribution load estimation*, CIREC 2009 - 20th International Conference and Exhibition on Electricity Distribution - Part 1, 2009, pp. 585–585, DOI: 10.1049/cp.2009.0857
179. Laboratory for Advanced Software Systems (LASS), *Smart*: Optimizing energy consumption in smart homes*, <http://lass.cs.umass.edu/projects/smart/> (22.06.2022.)
180. Wikipedia, *Federal holidays in the United States*, https://en.wikipedia.org/wiki/Federal_holidays_in_the_United_States (22.06.2022.)
181. Wikipedia, *Public holidays in Australia*, https://en.wikipedia.org/wiki/Public_holidays_in_Australia (22.06.2022.)
182. F. Zhang, C. Deb, S. E. Lee, J. Yang, K. W. Shah, *Time series forecasting for building energy consumption using weighted support vector regression with differential evolution optimization technique*, *Energy and Buildings*, vol. 126, pp. 94–103, 2016, DOI: 10.1016/j.enbuild.2016.05.028
183. I. Manojlović, *TimeSeriesR: Time series machine learning in R*, 2021, <https://github.com/igormanojlovic/TimeSeriesR> (22.06.2022.)
184. R. Ihaka, R. Gentleman, *R: A language for data analysis and graphics*, *Journal of Computational and Graphical Statistics*, vol. 5, no. 3, pp. 299–314, 1996, DOI: 10.1080/10618600.1996.10474713
185. F. Morandat, B. Hill, L. Osvald, J. Vitek, *Evaluating the design of the R language*, *ECOOP 2012 – Object-Oriented Programming*, 2012, pp. 104–131, DOI: 10.1007/978-3-642-31057-7_6
186. R. Elshawi, S. Sakr, D. Talia, P. Trunfio, *Big data systems meet machine learning challenges: Towards big data science as a service*, *Big Data Research*, vol. 14, pp. 1–11, 2018, DOI: 10.1016/j.bdr.2018.04.004
187. I. M. Coelho, V. N. Coelho, E. J. da S. Luz, L. S. Ochi, F. G. Guimarães, E. Rios, *A GPU deep learning metaheuristic based model for time series forecasting*, *Applied Energy*, vol. 201, pp. 412–418, 2017, DOI: 10.1016/j.apenergy.2017.01.003
188. I. Manojlović, *MTSR: Multiresolution time series representation*, 2020, <https://github.com/igormanojlovic/MTSR> (22.06.2022.)
189. M. B. Kurska, *praznik: Tools for information-based feature selection*, 2020, <https://cran.r-project.org/web/packages/praznik/praznik.pdf> (22.06.2022.)
190. D. Müllner, *fastcluster: Fast hierarchical, agglomerative clustering routines for R and Python*, *Journal of Statistical Software*, vol. 53, no. 9, pp. 1–18, 2013, DOI: 10.18637/jss.v053.i09
191. M. Walesiak, A. Dudek, *clusterSim: Searching for optimal clustering procedure for a data set*, 2018, <https://cran.r-project.org/web/packages/clusterSim/clusterSim.pdf> (22.06.2022.)
192. R. Wehrens, L. M. C. Buydens, *Self- and super-organizing maps in R: The kohonen package*, *Journal of Statistical Software*, vol. 21, no. 5, pp. 1–19, 2007, DOI: 10.18637/jss.v021.i05
193. P. Bruneau, *speccalt: Alternative spectral clustering, with automatic estimation of k*, 2015, <https://cran.r-project.org/web/packages/speccalt/speccalt.pdf> (22.06.2022.)
194. D. Meyer, *Support Vector Machines: The interface to libsvm in package e1071*. University of Applied Sciences Technikum Wien, Vienna, Austria, 2021.
195. M. Dancho, *modeltime.gluonts: GluonTS Deep Learning*, 2020, <https://ts.gluon.ai> (22.06.2022.)
196. D. Falbel, J. J. Allaire, F. Chollet, Y. Tang, W. van der Bijl, M. Studer, S. Keydana, *keras: R interface to Keras*, 2021, <https://github.com/rstudio/keras> (22.06.2022.)
197. D. Falbel, J. J. Allaire, Y. Tang, D. Eddelbuettel, N. Golding, T. Kalinowski, *tensorflow: R Interface to TensorFlow*, 2021, <https://github.com/rstudio/tensorflow> (22.06.2022.)
198. L. S. Riza, E. P. Nugroho, M. B. A. Prabowo, E. Junaeti, A. G. Abdullah, *metaheuristicOpt: Metaheuristic for optimization*, 2019, <https://cran.r-project.org/web/packages/metaheuristicOpt/metaheuristicOpt.pdf> (22.06.2022.)
199. F. M. Dekking, C. Kraaikamp, H. P. Lopuhaä, L. E. Meester, *A modern introduction to probability and statistics*, 2005, DOI: 10.1007/1-84628-168-7

BIOGRAFIJA



Igor Manojlović je rođen 25. maja 1990. godine u Novom Sadu. U Odžacima je završio osnovnu školu "Branko Radičević", 2005. godine, i opšti smer gimnazije "Jovan Jovanović Zmaj", 2009. godine. Na Departmanu za matematiku i informatiku na Prirodno-matematičkom fakultetu Univerziteta u Novom Sadu je završio osnovne akademske studije informacionih tehnologija, 2012. godine, i master akademske studije softverskog inženjerstva, 2014. godine. Na departmanu za energetiku, elektroniku i telekomunikacije na Fakultetu tehničkih nauka Univerziteta u Novom Sadu je završio specijalističke akademske studije, 2018. godine, i doktorske akademske studije, 2022. godine, položivši sve ispite predviđene planom i programom. Za vreme master studija je bio stipendista kompanije "Schneider Electric DMS NS" u Novom Sadu, gde je po završetku studija postao softverski inženjer, a kasnije i softverski arhitekta. U istoj kompaniji ostaje sve do 2022. godine, kada postaje *senior data scientist* u kompaniji "High Tech Engineering Center" u Novom Sadu. Autor je dva naučna rada objavljena u vrhunskim međunarodnim časopisima, kategorije M21.

SPISAK RADOVA

Radovi u vrhunskim međunarodnim časopisima (**M21**):

1. Igor Manojlović, Goran Švenda, Aleksandar Erdeljan, Milan Gavrić: *Time series grouping algorithm for load pattern recognition*, Computers in Industry 111 (2019) 140-147, DOI: [10.1016/j.com-pind.2019.07.009](https://doi.org/10.1016/j.com-pind.2019.07.009)
2. Igor Manojlović, Goran Švenda, Aleksandar Erdeljan, Milan Gavrić, Darko Čapko: *Hierarchical multiresolution representation of streaming time series*, Big Data Research 26 (2021) 100256, DOI: [10.1016/j.bdr.2021.100256](https://doi.org/10.1016/j.bdr.2021.100256)

Saopštenja sa međunarodnih skupova štampana u celini (**M33**):

1. Igor Manojlović, Aleksandar Erdeljan: *Efficient aggregation of time series data*, ICIST 2017 Proceedings 1 (2017) 102-107, <https://www.eventiotic.com/eventiotic/library/paper/258>
2. Igor Manojlović, Goran Švenda, Aleksandar Erdeljan: *Load pattern recognition method for probabilistic short-term load forecasting at low voltage level*, 2022 IEEE PES Innovative Smart Grid Technologies Conference Europe (ISGT-Europe), Session 20: Industry Application – Part I, Paper No. PAS-20-4 (ID-11137)

Овај Образац чини саставни део докторске дисертације, односно докторског уметничког пројекта који се брани на Универзитету у Новом Саду. Попуњен Образац укоричити иза текста докторске дисертације, односно докторског уметничког пројекта.

План третмана података

Назив пројекта/истраживања
Краткорочна пробабилистичка прогноза оптерећења на ниском напону у електродистрибутивним мрежама
Назив институције/институција у оквиру којих се спроводи истраживање
Универзитет у Новом Саду, Факултет техничких наука, Департман за енергетику, електронику и телекомуникације
Назив програма у оквиру ког се реализује истраживање
Докторске академске студије - Енергетика, електроника и телекомуникације
1. Опис података
<p>1.1 Врста студије</p> <p><i>Укратко описати тип студије у оквиру које се подаци прикупљају</i></p> <p>Ефикасност предложеног решења за краткорочну пробабилистичку прогнозу оптерећења на ниском напону у електродистрибутивним мрежама је верификована у студији случаја над скупом реалних и јавно доступних података који су сакупљени са паметних бројила, у једној северноамеричкој и једној аустралијској електродистрибутивној мрежи.</p> <p>1.2 Врсте података</p> <p>а) квантитативни</p> <p>б) квалитативни</p> <p>1.3. Начин прикупљања података</p> <p>а) анкете, упитници, тестови</p> <p>б) клиничке процене, медицински записи, електронски здравствени записи</p> <p>в) генотипови: навести врсту _____</p> <p>г) административни подаци: навести врсту _____</p> <p>д) узорци ткива: навести врсту _____</p> <p>ђ) снимци, фотографије: навести врсту _____</p>

е) текст, навести врсту _____

ж) мапа, навести врсту _____

з) остало: јавно доступни подаци који су сакупљени са паметних бројила у једној северноамеричкој и једној аустралијској електродистрибутивној мрежи

1.3 Формат података, употребљене скале, количина података

1.3.1 Употребљени софтвер и формат датотеке:

а) Excel фајл, датотека _____

б) SPSS фајл, датотека _____

с) PDF фајл, датотека _____

д) Текст фајл, датотека: .csv

е) JPG фајл, датотека _____

ф) Остало, датотека _____

1.3.2. Број записа (код квантитативних података)

а) број варијабли: 11 (1 варијабла за активну снагу и 10 варијабли за временске услове)

б) број мерења (испитаника, процена, снимака и сл.): 6777 (114 потрошача из једне и 6663 потрошача из друге електродистрибутивне мреже)

1.3.3. Поновљена мерења

а) да

б) не

Уколико је одговор да, одговорити на следећа питања:

а) временски размак између поновљених мера је 1-30 минута

б) варијабле које се више пута мере односе се на активну снагу и временске услове

в) нове верзије фајлова који садрже поновљена мерења су именоване као _____

Да ли формати и софтвер омогућавају дељење и дугорочну валидност података?

а) Да

б) Не

Ако је одговор не, образложити _____

2. Прикупљање података

2.1 Методологија за прикупљање/генерисање података

2.1.1. У оквиру ког истраживачког нацрта су подаци прикупљени?

а) експеримент, навести тип _____

б) корелационо истраживање, навести тип _____

ц) анализа текста, навести тип _____

д) остало, навести шта - У студији случаја су коришћена два јавно доступна скупа података: 1) *UMass* подаци су прикупљени у оквиру *Smart** пројекта лабораторије за напредне софтверске системе; 2) *Smart Grid Smart City (SGSC)* подаци су прикупљени у оквиру истог пројекта аустралијске владе и индустријског конзорцијума предвођеним дистрибутивним предузећем *Ausgrid*.

2.1.2 Навести врсте мерних инструмената или стандарде података специфичних за одређену научну дисциплину (ако постоје).

2.2 Квалитет података и стандарди

2.2.1. Третман недостајућих података

а) Да ли матрица садржи недостајуће податке? Да Не

Ако је одговор да, одговорити на следећа питања:

а) Колики је број недостајућих података? _____

б) Да ли се кориснику матрице препоручује замена недостајућих података? Да Не

в) Ако је одговор да, навести сугестије за третман замене недостајућих података

2.2.2. На који начин је контролисан квалитет података? Описати

2.2.3. На који начин је извршена контрола уноса података у матрицу?

3. Третман података и пратећа документација

3.1. Третман и чување података

3.1.1. Подаци су депоновани у UMass и SGSC репозиторијумима.

3.1.2. URL адресе:

<http://traces.cs.umass.edu/index.php/Smart/Smart>

<https://data.gov.au/data/dataset/smart-grid-smart-city-customer-trial-data>

3.1.3. DOI _____

3.1.4. Да ли ће подаци бити у отвореном приступу?

а) Да

б) Да, али после ембарга који ће трајати до _____

в) Не

Ако је одговор не, навести разлог _____

3.1.5. Подаци неће бити депоновани у репозиторијум, али ће бити чувани.

Образложење

3.2 Метаподаци и документација података

3.2.1. Који стандард за метаподатке ће бити примењен? _____

3.2.1. Навести метаподатке на основу којих су подаци депоновани у репозиторијум.

Ако је потребно, навести методе које се користе за преузимање података, аналитичке и процедуралне информације, њихово кодирање, детаљне описе варијабли, записа итд.

3.3 Стратегија и стандарди за чување података

3.3.1. До ког периода ће подаци бити чувани у репозиторијуму? _____

3.3.2. Да ли ће подаци бити депоновани под шифром? Да Не

3.3.3. Да ли ће шифра бити доступна одређеном кругу истраживача? Да Не

3.3.4. Да ли се подаци морају уклонити из отвореног приступа после извесног времена?

Да Не

Образложити

4. Безбедност података и заштита поверљивих информација

Овај одељак МОРА бити попуњен ако ваши подаци укључују личне податке који се односе на учеснике у истраживању. За друга истраживања треба такође размотрити заштиту и сигурност података.

4.1 Формални стандарди за сигурност информација/података

Истраживачи који спроводе испитивања с људима морају да се придржавају Закона о заштити података о личности (https://www.paragraf.rs/propisi/zakon_o_zastiti_podataka_o_licnosti.html) и одговарајућег институционалног кодекса о академском интегритету.

4.1.2. Да ли је истраживање одобрено од стране етичке комисије? Да Не

Ако је одговор Да, навести датум и назив етичке комисије која је одобрила истраживање

4.1.2. Да ли подаци укључују личне податке учесника у истраживању? Да Не

Ако је одговор да, наведите на који начин сте осигурали поверљивост и сигурност информација везаних за испитанике:

- а) Подаци нису у отвореном приступу
- б) Подаци су анонимизирани
- ц) Остало, навести шта

5. Доступност података

5.1. Подаци ће бити

а) јавно доступни

б) доступни само уском кругу истраживача у одређеној научној области

ц) затворени

Ако су подаци доступни само уском кругу истраживача, навести под којим условима могу да их користе:

Ако су подаци доступни само уском кругу истраживача, навести на који начин могу приступити подацима:

5.4. Навести лиценцу под којом ће прикупљени подаци бити архивирани.

6. Улоге и одговорност

6.1. Навести име и презиме и мејл адресу власника (аутора) података

6.2. Навести име и презиме и мејл адресу особе која одржава матрицу с подацима

6.3. Навести име и презиме и мејл адресу особе која омогућује приступ подацима другим истраживачима
