







ESCAP/WMO Typhoon Committee 19th INTEGRATED WORKSHOP / AP-TCRC FORUM

The Deep-Learning Algorithm for Hydrological Data Quality Control and Flood Forecasting in TC member countries

20 Nov. 2024

Ph.D. Chung-Soo Kim WGH AOP 2 &3 Leader



Contents

Hydrological Data Quality Control using AI Deep-Learning Techniques

2

3

Application of AI Hydrological Data Quality Control in TC member countries

Flood Forecasting using AI Deep-Learning Techniques

Hydrological Data Quality Control using AI Deep-Learning Techniques

Hydrological Data Quality Control (HDQC)



<Systemic framework of hydrological data quality management >

- \checkmark Semi-automatic quality control
- ✓ Statistical data analysis, quality ratings & flag
- ✓ Historical data analysis and establish criteria for screening abnormal data

✓ Criteria:

- Rainfall: Max. Rainfall, RDS
- Water Level: Abrupt change,

Constant



Hydrological Data Quality Control System (HDQCS) Rainfall Data Import

2023-11-06 11:45:25					Hydrolog	jical Data	Quality	Cont	rol Sy	vsten	ı								
Rainfall	Rainfall / Verificati Data Max. Rai	on Rainfall	/ Data correction	Rainfall / I	Report														
Verification	Time Interval : 🤇	🕽 daily 🔵 hou	urly 🔿 10-minute	•	Data Import	al Import								Stat	ion 🕴	Station			-
– Max. Rainfall	Date	Station	RDS1	RDS2	RDS3	RDS4		Rainfall [S											
– RDS	2019-01-01	0	0	0	0	0				Ra	Rainfall [Station]								
[AI] Abnormal Dect.	2019-01-02	0	0	0	0	0		120								_	St	ation	
Supervised Learning	2019-01-03	0	0	0	0	0	120												
Unsupervised Learning	2019-01-04	0	0	0	0	0													
Data Correction	2019-01-05	0	0	0	0	0		100 -											
- Data	2019-01-06	0	0	0	0	0													
- Data Correction	2019-01-07	0	0	0	0	0	-	80 -											
Report	2019-01-08	0	0	0	0	0	[mu												
Report	2019-01-09	0	0	0	0	0	Ifall	60 -											
	2019-01-10	0	0	0	0	0	Rair	Rair											
Water Level	2019-01-11	0	0	0	0	0		40 -											
Verification	2019-01-12	0	0	0	0	0													
- Data	2019-01-13	0	0	0	0	0		20 -											
- Constant Value			RDS1	RDS2	RDS3	RDS4													
[AI] Abnormal Dect.	Distance (km	1)	2	2	3	4		₀⊥	ل ال	ш.Ц.					_ <u></u>		╷║║╢	ļ.	1
Supervised Learning Unsupervised Learning									01-01 04-11	07-20	10-28	02-05	05-15	08-23	12-01	03-11	06-19	09-27	
Data Correction									2019-	2019-	2019-	2020-	2020-	2020-	2020-	2021-	2021-	2021-	
Data Correction						Apply													

Hydrological Data Quality Control System (HDQCS)

Rainfall Criteria (Max. Rainfall)

Data Correction

2023-11-06 11:45:25				Hyd	rological Dat	a Quality	Control Sy	stem					
Rainfall	Rainfall / Verificat	ion Rainfall / Da	ata correction	Rainfall / Report									
Verification	Data Max. Ra Rainfall distribut	infall RDS A	1										
— Max. Rainfall		10mm~	20mm~	30mm~	40mm~	50mm~	60mm~	70mm~	80mm~	90mm~	100mm		
– RDS	No. of Count	94.0	51.0	30.0	19.0	12.0	8.0	6.0	3.0	2.0	2.0		
[AI] Abnormal Dect. Supervised Learning Unsupervised Learning	% of Count	100.0	54.3	31.9	20.2	12.8	8.5	6.4	3.2	2.1	2.1		
Data Correction Data Data Data Correction [AI] Data Correction Report	Statistics												
•	Mean	Std.	Mean+2*Std. N	Mean-2*Std. Max	(mm) 0.9max (mm)								
Water Level	13.5	18.2	49.8	-22.8 12	1.0 108.9			20.2	<u>8 8.5 6.4</u> 3.	2 2.1 2.1	· 0		
Verification – Data – Abrupt Change – Constant Value [Al] Abnormal Dect. Supervised Learning Unsupervised Learning	Outlier Verificat % of count Max. rainfall	tion Criteria 2.1 % 90.0 mm	Ca	alculation Apply		Count [#]	80 40 54.3 40 25 54.3 54.3 54.3 50 50 50 50 50 50 50 50 50 50						
Data Correction								Rain	fall [mm]				

Download

Hydrological Data Quality Control System (HDQCS) Rainfall Data Criteria (RDS)

2023-11-06 11:45:25						H	ydro	ologi	cal C	Data 🕻	Quality	Control S	Syste	m						
Rainfall	Rainfall / Verification	Rainfall /	Data correc	tion	Rainfall	/ Report														
Verification Data	Data Max. Rainfa	RDS Criteria			- 50%		~70%	- 90%			_									
 Max. Rainfall RDS [Al] Abnormal Dect. 	No. of Count % of Count	0.0 0. 0.0 0.	0 0.0 0 0.0	~40% 0.0 0.0	2.0 2.4	~60% 5.0 6.1	~70% 5.0 6.1	~80% 8.0 9.8	~90% 9.0 11.0	~100% 15.0 - 18.3	14		.4 6. 5 .1		2.4 4.9 3.7	.2 1.2 1.2 2.4		••••	- 0	
Supervised Learning Unsupervised Learning	No. of Count	~110% ~12 10.0 8.	0% ~130% 0 6.0	~140% 4.0	~150% 2.0	~160% 3.0	~170% 1.0	~180% 2.0	~190% 1.0	~200% ⁻	12 - 10 -	: -) -	7.3 9.8 11.0 12.1		7.3 8		● N	o. of Count 5 of Count	$\begin{bmatrix} 5 \\ -10 \end{bmatrix}$	
Data Correction – Data – Data Correction	% of Count	12.2 9. ~210% ~22	8 7.3 0% ~230%	4.9 ~240%	2.4 ~250%	3.7 ~260%	1.2 ~270%	2.4 ~280%	1.2 ~290%	0.0 ~300%	Dent [- Count [- 9	; - ; -		18/3					- 15 Jo	
[AI] Data Correction Report	No. of Count % of Count	1.0 0. 1.2 0.	0 0.0 0 0.0	0.0 0.0	0.0	0.0 0.0	0.0 0.0	0.0 0.0	0.0 0.0	0.0	4 2				.	_			- 20 × - 25	
	Statistics Mean	Mean+2	Mean+2*Std. Mean-2*Std.			Max (mm) (0.9Max (mm)		0	· ,			ĻЩ				₃₀		
Water Level	102.2	32.1	166	.3	38.	1	20	8.6	1	87.8		~10% ~40%	% ~70%	~100% ~ R	130% ~1609 DS Criteri	% ~190% ~2 a [%]	20% ~250	% ~280%		
Verification – Data – Abrupt Change	Outlier Verificatio	on Criteria	9	6		Calculatio	'n	Un	per diff	erence	139.6	9	6	J.	Apply			al Dov	wnload	
[AI] Abnormal Dect. Supervised Learning Unsupervised Learning	Lower limit 5		9	6		Calculatio	'n	Lov	ver diff	erence	57.0	9	6	✓ 	Apply					
Data Correction																				

Data

Data Correction

Hydrological Data Quality Control System (HDQCS) Water Level Data Import

2023-11-06 11:45:25			Hydrological Data Q	Quality Control System								
Rainfall	Water level /	Verification Criteria Water level / Data (Correction Water level / Report									
Verification	Time Inter	val : O daily () hourly () 10-minute	Data Import 🗾 Import									
— Data — Max. Rainfall		Date	Water level (m)									
– RDS	1	2018-01-01	0.85	Water Level								
[AI] Abnormal Dect.	2	2018-01-02	0.84	4.0 - Water level (m)								
Supervised Learning	3	2018-01-03	0.83									
Unsupervised Learning	4	2018-01-04	0.83	3.5 -								
Data Correction	5	2018-01-05	0.84									
- Data Data Correction	6	2018-01-06	0.83	3.0 -								
[AI] Data Correction	7	2018-01-07	0.83	E 25-								
Report	8	2018-01-08	0.84									
•	9	2018-01-09	0.83	<u>à</u> 2.0 -								
Water Level	10	2018-01-10	0.83	ater ater ater ater ater ater ater ater								
vvater Level	11	2018-01-11	0.82	₿ 1.5								
Verification	12	2018-01-12	0.82									
– Data	13	2018-01-13	0.82									
Onstant Value	14	2018-01-14	0.82	0.5 -								
[AI] Abnormal Dect.	15	2018-01-15	0.83									
Supervised Learning	16	2018-01-16	0.84	0.0								
Unsupervised Learning	17	2018-01-17	0.84	01								
Data Correction	18	2018-01-18	0.84	8-01- 8-01- 9-05- 9-12- 9-05- 9-12- 9-05- 9-12- 9-05- 9-12- 9-05- 9-12- 9-03- 12-09- 12-12- 12-12- 12-03-12-03- 12-03-12-03- 12-03-10-03-10-03-10-03-10-03-10-03-10-03-10-03-10-03-10-03-10-03-10-03-10								
Data Data Correction	19	2018-01-19	0.84	201 201 201 201 201 201 201 201 201 201								

Hydrological Data Quality Control System (HDQCS) Water Level Criteria (Abrupt change)

2023-11-06 11:45:25							ŀ	Hydro	ologi	cal D	ata Q	uality Co	ontrol S	ystem					
Rainfall	Water level / Veri	fication	Criteria	Water l	evel / Da	ata Correc	tion	Water le	evel / Rej	port									
Verification	Data Abrupt Water level ch	: Change lange rat	Cons io Ab	tant Value rupt chan	e Al ge criter	ria													
 Max. Rainfall RDS 	No. of Count	~(-X	5)	~(-X4) 2.0	~	(-X3)	~(-X 22.0	2) D	~(-X1) 94.0	3	~(X0) 339.0		0.8 0.3		2.2	0.4 0.	6 0.0 0.0	0.0 0.0 0.1	
[AI] Abnormal Dect. Supervised Learning	% of Count	0.8	2	0.27	2	2.05	3.0		12.82	4	46.25	350 -		3.0		2.2			
Unsupervised Learning Data Correction	No. of Count	(+X1)~ 203.0	(+X2)~ 16.0	(+X3)~ 16.0	(+X4)~ 3.0	(+X5)~ 4.0	(+X6)~ 0.0	(+X7)~ 0.0	(+X8)~ 0.0	(+X9)~ 0.0	(+X10)~ 1.0	300 -		12.8				No. of Count	- 10
 Data Data Correction 	% of Count	27.69	2.18	2.18	0.41	0.55	0.0	0.0	0.0	0.0	0.14	250 -						% of Count	
[AI] Data Correction Report	Mea 1 -0.2	n 6	3	6td. .84	Me	an+2*Std. 7.43		Mean-2*S - 7.95	itd.	Ma 47	ах 7	[# 200 - tur			27.7				- 20 - Connt [%]
Water Level															1				- 30 %
Verification – Data Abrupt Change	Outlier Verifi Upper limit Lower limit	ication C	riteria		%			Calculatio	on			100 -							- 40
Constant Value [Al] Abnormal Dect. Supervised Learning	Upper chang Lower chang	je 3.7 je -3.6				~	,	Apply				50 -		4	6.2				50
Unsupervised Learning Data Correction												0	~(-X5) ~	(-X3) ~(-X1)	(+x1)~ Change	(+x3)~ (+x5 e Criteria	5)~ (+X7))~ (+X9)~	30

Hydrological Data Quality Control System (HDQCS) Water Level Criteria (Constant value)

2023-11-06 11:45:25					Hydrologica	al Data 🕻	Quality	Contr	o <mark>l Syst</mark> e	m			
Rainfall	Water level / Verif	ication Criteria	Water level	/ Data Correction	Water level / Repo	rt							
Verification Data	Water level dist	tribution by cor 48hr	nstant time 72hr	Al 96hr	120hr	144hr	168	٦r	240hr	360hr	480hr	600hr	720hr
— Max. Rainfall — RDS	No. of Count	92.0	17.0	15.0	11.0	5.0	1.0		0.0	0.0	1.0	1.0	0.0
[AI] Abnormal Dect. Supervised Learning Unsupervised Learning	% of count	100.0	10.5	10.5	12.0	5.4			0.0	0.0			0.0
Data Correction	Statistics Mean		Std.	Mean+2*Std.	Mean-2*Std.	M	ах	100	T :	5	1.1 0.0 0).0 1.1 1.1 (0.0
 Data Correction [AI] Data Correction Report 	1 3.03		2.67	8.37	-2.31	2	5	80	18	12.0 .5 16.3		No. of Co	- 20 ount
Water Level								1t [#]					- 40 [%] tunc
Verification – Data	Outlier Verifica	tion Criteria						Ino ₂ 40	$\left \right $				- 60 Ŭ of %
 Abrupt Change Constant Value [Al] Abnormal Dect. Supervised Learning 	% of count 1 Max. period 7	.3	%	6 📕	Calculation Apply			20	100.0		_		- 80 - 100
Data Correction Data Data								0	48hr 72	hr 96hr 120hr 14	4hr 168hr 240hr 36 Period [hr]	50hr 480hr 600hr 72	20hr Download

Description of AI (Artificial Intelligent), Machine Learning, Deep Learning



Artificial Intelligence (AI):

The study of intelligence demonstrated by a machine manifested by its capability to perceive the environment and take actions to achieve its goals and tasks through flexible adaptation.

Machine Learning (ML)

A sub-set of AI, which are learning methods and algorithms that enable computers to automatically improve performance through experience.

Representation Learning (RL)

Techniques that automatically discover representation (or featur es) that are useful for subsequent learning tasks. Also known as feature learning.

Deep Learning (DL)

A class of machine learning algorithms based on artificial neural networks (ANNs) and using hierarchical architectures to extract higher level features from input data via representation learning

History of AI Application in Hydrology Sector

1980s: Artificial neural networks (ANNs) were first applied to hydrology, primarily for rainfall-runoff modeling.

1990s: Support vector machines (SVMs) were introduced and applied to hydrology, showing promising results for streamflow forecasting.

Early 2000s: Other machine learning techniques, such as decision trees and random forests, were applied to hydrology, mainly for flood risk mapping and water quality modeling.

Mid-2000s: Ensemble methods, such as bagging and boosting, were applied to hydrology modeling, showing improved accuracy compared to individual models.

Mid-2000s: Ensemble methods, such as bagging and boosting, were applied to hydrology modeling, showing improved accuracy compared to individual models.

Late 2000s: Deep learning algorithms, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), were introduced and applied to hydrology, showing improved accuracy for various hydrologic tasks, including streamflow prediction, rainfall-runoff modeling, and flood forecasting.

2010s: Transfer learning, a technique where knowledge learned in one task is transferred to another task, was applied to hydrology, showing improved accuracy for streamflow prediction.

2020s: With the increasing availability of large datasets, there has been a growing interest in applying ML to hydrology modeling, including applications of generative adversarial networks (GANs) for rainfall simulation and unsupervised learning for anomaly detection in hydrologic data.

Selection of AI Deep-Learning Technique for HDQC

1) Supervised Learning for Anomaly Detection Using XGBoost

- XGBoost, an extension of gradient boosting, can be used as a supervised anomaly detection algorithm by treating one class as anomalies and the other as normal data.
- Labelled data is required to enable the model to learn to distinguish between the abnormal and normal data. XGBoost cannot be carried out without labelled data.

2) Unsupervised Learning for Anomaly Detection Using Isolation Forest

- An ensemble method that isolates anomalies by constructing random forests and isolating data points that require fewer splits in the tree to be isolated.
- It is a simple yet effective approach for detecting anomalies. *No labelled data is required*, the anomaly can be recognized by just specifying the contamination level.

All the codes are written using open-source Python language

Application of AI HDQCS in TC member countries

Using Isolation Forest for Identifying Abnormal Rainfall Data in Republic of Korea (Hourly Data) & Lao PDR (Daily Data)

Max. Rainfall for hourly data in ROK (Yeoju Bridge) Zero data is excluded

Contamination Level=0.01, Threshold=23mm

Contamination Level=0.001, Threshold=50mm



- The concept of contamination level is similar with percentile where contamination level of 0.01 is equivalent to 99 percentile and contamination level of 0.001 is equivalent to 99.9 percentile.
- Low threshold value will be obtained for outlier identification when zero data is included due to increasing amount of irrelevant data sample.

Max. Rainfall for daily data in Lao PDR (Pakkayoung Station) Zero data is excluded



The lesser the contamination level specified, the higher the threshold value is, and the lesser data will be identified as abnormal data.

RDS for daily data in Lao PDR (Thalad Station)

Contamination Level=0.15, Upper Threshold=+55.5%

Contamination Level=0.15, Lower Threshold=-8.2%



- The RDS abnormal is identified from the differences between the rainfall value at the target station and the RDS rainfall value calculated from nearby rainfall stations.

Using Isolation Forest for Identifying Abnormal Water Level Data in ROK(Hourly Data) and Lao PDR (Daily Data)

Water Level Gradient (between current and previous one hour data) for hourly data in ROK (Yeoju Bridge) WL Gradient for Upper Limit, Contamination level = 0.001 WL Gradient for Lower Limit, Contamination level = 0.001



Constant water level for hourly data in **ROK** (Yeoju Bridge)



Contamination Level=0.01, Threshold=102 hours



Contamination Level=0.001, Threshold=407 hours

Water Level Gradient (between current and previous one day data) for daily data in Lao PDR (Pakkayoung) Dummy WL Gradient for Upper Limit, Contamination level = 0.001 Dummy WL Gradient for Lower Limit, Contamination level = 0.001



<u>Upper Threshold=+67.5%</u>

Constant water level for daily data in Lao PDR (Pakkayoung)



<u>Contamination Level=0.01, Threshold=22 days</u>

Contamination Level=0.001, Threshold=26 days

Using XGBoost for Identifying Abnormal Rainfall Data in ROK(Hourly Data) and Lao PDR (Daily Data)

Max. Rainfall for hourly data in ROK (Yeoju Bridge) Different Percentage of Training Data (40%) and Testing Data (60%) Dummy RF threshold >= 40mm (labelled data as abnormal)

Anomalies in Training (2013-2016)

Anomalies in Testing (2017-2022)



It is advised to specify longer period for training to allow the model to study and recognize the labelled data as much as possible for identifying the abnormal data in the testing period.

25



Anomalies in Training



RDS for daily data in Lao PDR (Thalad Station) (70% Training, 30% Testing) Dummy Upper Limit >= 50% Differences (labelled data as abnormal) Dummy Lower Limit <= -10% Differences (labelled data as abnormal)





Date

Using XGBoost for Identifying Abnormal Water Level Data in ROK(Hourly Data) and Lao PDR (Daily Data) Water Level Gradient (between current and previous one hour data) for hourly data in ROK (Yeoju Bridge) (70% Training, 30% Testing) Dummy WL Gradient Threshold, Upper Limit >= 100% Differences (labelled data as abnormal) Dummy WL Gradient Threshold, Lower Limit <= -30% Differences (labelled data as abnormal)



Constant water level for hourly data in ROK (Yeoju Bridge) (70% Training, 30% Testing) Dummy constant WL Threshold >= 102 hours (labelled data as abnormal)

Observed Water Level Observed WL **Input Data** 5 4 Water Level (m) з 2 1 0 2014 2016 Date 2020 2022 Longest Constant WL Periods in Training & Testing Training Histogram of 50000 Testing Input Data (Constant WL 40000 Periods) 00008 Count 20000 10000 0 20 60 80 0 40 100

Constant WL Periods

Anomaly Detection for Constant Water Level using XGBoost



Anomalies in Testing

Anomalies in Training

Anomaly Detection for Constant Water Level using XGBoost



Water Level Gradient (between current and previous one day data) for daily data in Lao PDR (Pakkayoung Station) (70% Training, 30% Testing) Dummy WL Gradient Threshold >= 99 percentile (labelled data as abnormal)





Constant water level for daily data in Lao PDR (Pakkayoung Station) (70% Training, 30% Testing) Dummy constant WL Threshold >= 22 days (labelled data as abnormal)

Anomalies in Training







Anomaly Detection for Constant Water Level using XGBoost

Hydrological Data Quality Control System (HDQCS) Rainfall Data Criteria (AI)

2023-11-06 11:45:25	Hydrological Data Q	uality Control System
Rainfall	Rainfall / Verification Rainfall / Data correction Rainfall / Report Data Max. Rainfall RDS Al	
Verification	Abnormal Detection using AI - Supervised Learning (S.L.)	Abnormal Detection using AI - Unsupervised Learning (U.L.)
— Data — Max. Rainfall — RDS	1. Select Training, Validation and Test period ratio (ex: 7 : 0 : 3) Training : 7 Validation : x Test : 3 Apply	1. Al model configuration (Contamination (outlier rate) ex.: 0.1%) Max. rainfall threshold: 52.7 Contamination : 1 % Image: Build a Model
EAI Abnormal Dect. Supervised Learning Unsupervised Learning	2. Al model configuration (Training -> Validation -> Test)	Outlier detection of rainfall verification data using Isolation Forest
Data Correction		120 - Observed Ar
Data Data Correction [AI] Data Correction Report	120 - Observed RF	100 -
Water Level	80 -	80 -
Verification Data Abrupt Change Constant Value [AI] Abnormal Dect. Supervised Learning Unsupervised Learning Data Correction	40 - 20 -	40 - 20 -
Data Data Correction [AI] Data Correction Report		
	3. Application the AI model to rainfall verification and correction data	d Graph Download
	Verification data Appy Correction data Apply Image: State of the sta	2. Application the AI model to rainfall correction data

Hydrological Data Quality Control System (HDQCS) Water Level Data Criteria(AI)



Flood Forecasting using AI Deep-Learning Techniques

AI Deep-Learning Techniques for Water level prediction



- Recurrent Neural Networks (RNN)
- Convolutional Neural Networks (CNN)
- Long Short-Term Memory (LSTM)
- Advantage of using Deep-Learning Techniques
 - Adaptability
 - Patten recognition capability
 - Automatic learning

Water level prediction Procedure using AI Deep-Learning Techniques



Hyperparameter tuning

1) Number of Nodes

(rule of thumb: 1 HL for simple problem and 2HL for complex problem; high accuracy with many nodes in a layer)

2) Number of hidden layers

(rule of thumb: 1-3 units or nodes per layer is a good base)

3) Dropout

(rule of thumb: The 20% (0.2) is widely accepted as the best compromise between overfitting and retaining model accuracy)

4) Input sequence length

5) Learning rate

(rule of thumb: Usually a decaying learning rate is preferred and this hyper-parameter is used in the training phase and has a small positive value, mostly between 0.0 and 0.1.)

6) Batch Size

(rule of thumb: Widely accepted, a good default value for batch size is 32. For experimentation, you can try multiples of 32, such as 64, 128 and 256.)

7) Input dimensionality

8) Number of Epochs

(rule of thumb: Widely accepted, a good default value for batch size is 32. For experimentation, you can try multiples of 32, such as 64, 128 and 256.)

Training evaluation

Application Case



No.	Station	Latitude	Longitude
1	Phiangluang	19°34'06" N	103°04′17″ E
2	Thalad	18°31'26″ N	102°30′54″ E
3	Pakkayoung	18°25′53″ N	102°32′16″ E
4	Veunkham	18°10'37" N	102°36′53″ E

	R	ainfall (mm	ı)	١	Water Level (m)				
Station	Annual Average	Max.	Average	Max.	Average	Min	– Rainfall (RF)	Water Level (WL)		
Phiangluang	1283	155.8	3.8	8.72	0.79	0.35	2019.01.01 ~2021.10.14	2019.01.01~2021.10.14		
Thalad	1515	118.5	4.5	10.95	6.09	0.64	2019.01.01~2021.10.11	2019.01.01~2021.10.11		
Pakkayoung	1629	111.5	4.8	8.72	3.95	2.28	2019.01.01~2021.10.11	2019.01.01~2021.10.11		
Veunkham	1651	142.8	4.9	9.10	2.71	0	2019.01.01~2021.10.10	2019.01.01~2021.10.10		

Application Case





Correlation analysis



Model Training



Model Testing





Validation Results











Prediction according to sequence length



Conclusion

\clubsuit The factors that should be considered

- Sufficient correlation of training data
- Selection of appropriate features
- Initial model
- Fine tuning
- Appropriate parameters

✤ Needs

- Non-time series data (Static variables) management
- Deterministic or Stochastic information?

Thanks for your attention.