Prediction of Health Evaluation Indices for Aquatic Ecosystem using Extreme Gradient Boosting Tree (XGBoost) and SWAT

Chung Gil Jung

(wjd0823@konkuk.ac.kr)

So Young Woo and Seong Joon Kim

School of Civil and Environmental Engineering, College of Engineering,

Konkuk University, South Korea



Introduction

- In South Korea, since 1990s, the people became interested in restoration of river environmental function for their life quality with water.
- Since 2007, the Ministry of Environment has monitored the aquatic ecosystem health (AEH) and evaluated the stream and river AEH.
- By the limitations of monitoring sites, we need some technique to develop AEH evaluation for the whole country streams.
- This study is to predict the AEH indices (FAI, TDI, and BMI) of ungauged streams for a standard watershed scale (about 500 km²) using SWAT results and Ensemble Machine Learning algorithm.

Assessment Procedure



• stream discharge (SD)

Study area



GIS data



4 Multipurpose dam data (area-level and storage-level relationship curve)



3 Multifunction weir data (area-level and storage-level relationship curve)



4 <u>Multipurpose dam</u> data (release and storage : 1984-2014)



Chungju dam (CJD)



Paldang dam (PDD)

2009

150

120

90

60

30

Storage (10⁶m³)

Fotal Release Storage

2012



3 Multifunction weir data (release and storage : 2012-2014)



Kangcheon wier (KCW)



Hydrology Results

Ahn et al. (2015)

Observed vs. simulated streamflow results of model calibration and validation

✓ Calibration : 5 years (2005-2009) / Validation : 5 years (2010-2014)



• Hydrology Results

Ahn et al. (2015)

Observed vs. simulated streamflow results of model calibration and validation

✓ Calibration : 2 years (2012-2013) / Validation : 1 year (2014)



Hydrology Results

Ahn et al. (2015)

Observed vs. simulated ET & SM results of model calibration and validation

✓ Calibration : 3 years (2009-2011) / Validation : 2 years (2012-2013)



Water quality

Ahn et al. (2015) Observed vs. simulated sediment results of SWAT model calibration and validation

✓ Calibration : 5 years (2005-2009) / Validation : 5 years (2010-2014) at 7 stations



Waterquality

Ahn et al. (2015)

Observed vs. simulated <u>T-N</u> results of SWAT model calibration and validation ✓ Calibration : 5 years (2005-2009) / Validation : 5 years (2010-2014) at 7 stations



Waterquality

Ahn et al. (2015)

Observed vs. simulated <u>T-P</u> results of SWAT model calibration and validation

✓ Calibration : 5 years (2005-2009) / Validation : 5 years (2010-2014) at 7 stations



Data for Aquatic Ecology Health Index

 South Korea has the National Aquatic Ecological Monitoring Program (NAEMP) operated by the Ministry of Environment and the National Institute of Environmental Research, Korea.



- since 2007 for entire country
- spring (April to May) and autumn (September to October) in twice a year
- Measurement components : water temperature, pH, DO, BOD, NH₄, NO₃, T-N, T-P, PO₄, Chlorophyll-a, and so on.
- from that components, TDI (Trophic Diatom Index), FAI (Fish Assessment Index), and BMI (Benthic Macroinvertebrate Index) have be estimated

Locations of the sampling sites in five major watersheds in South Korea

Int. J. Environ. Res. Public Health 2012, 9, 3629-3653

Fish Assessment Index, FAI (U. S. EPA, 1993) ٠

: Calculate the score for 4 metrics (M1, M2, M3 and M7) that depends on the stream order of Korea and the other 4 metrics (M4, M5, M6 and M8), using class division "0", "6.25", "12.5"

FAI = M1 + M2 + M3 + M4 + M5 + M6 + M7 + M8

- M1 : Total number of domestic species
- M2 : Number of riffle benthic species
- M3 : Number of sensitive species
- *M7*: Total population of sampled domestic species
- M4 : Population ratio of Tolerant species M5 : Population ratio of Omnivores M6 : Population ratio of domestic Insectivores
- M8 : Population ratio of Anormalities





Blue gil



Lon-nosed barbel

Trophic Diatom Index, TDI (Kelly and Withon, 1995)

: Calculate the score using relative density, contamination sensitivity and indicative value of emerging species for about 400 species of Trophic Diatom ai: Relative abundance of species

TDI = 100-(25(
$$\sum ai * si * vi / \sum ai * vi) - 25$$
)

- in specimens, %
- si: Contamination sensitivity of species (1 - 5)
- vi: Indicative value of species (1 3)



Bacillariales

Genus Cocconeis

Benthic Macroinvertebrate Index, BMI (Merritt and Cummins 1984) • : Calculate the score using benthic macroinvertebrate organisms with contaminate and indicative weighted value *i*: the number of individuals in the sample

$$BMI = (4-\sum si * hi * gi / \sum hi * gi) * 25$$

si: pollution unit index of species *i hi*: appearance rank of species *i* gi: weight index of species i

ephemera orientalis Chironomidae spp. mclachlan

Index	A (Very Good)	B (Good)	C (Fair)	D (Poor)	E(Very Poor)
FAI	$100 \geq \sim \geq 80$	$80 > \sim \ge 60$	$60 > \sim \ge 40$	$40 > \sim \ge 20$	$20 > \sim \ge 0$
TDI	$100\geq \sim \geq 90$	$90 > \sim \ge 70$	$70>\sim \geq 50$	$50 > \sim \ge 30$	$30 > \sim \ge 0$
BMI	$100 \geq \sim \geq 80$	$80 > \sim \ge 65$	$65 > \sim \ge 50$	$50 > \sim \ge 35$	$35 > \sim \ge 0$

Fish Assessment Index, FAI :

It refers to the organism at the • top level of the food chain in the water body that represents omnivorous, herbivorous, inflorescence, and carnivorous at various nutritional stages.





버들치 Rhynchocypris oxycephalus









Coreoleuciscus splendidus



Zacco temminckii



쏘가리 Siniperca scherzeri





Zacco platypus



끄리 Opsarichthys uncirostris



모래무지 Pseudogobio esocinus







메기 Silurus asotus



Carassius auratus



임어 Cyprinus carpio



- As the primary producer of the river ecosystem food chain, it refers to the diatom attached to the stone (substrate) such as gravel or cobble stone in the bottom.
- It is responsible for energy transfer in the ecosystem.
- Also, It is sensitive to changes in water quality (TN, TP) and environment.



Benthic Macroinvertebrate Index, BMI :

- It is a biologic indicator that reflects local characteristics as a primary or secondary consumer of river ecosystem.
- It refers to a group of aquatic insects.
- As a sub-consumer linking producers and upper consumers in aquatic ecosystem food chain, they are sensitive to environmental changes and have excellent indicators and are used as indicators of water quality evaluation.



깔따구류(적색)

Chironomidae spp.(red type)

나방파리

Psychoda alternata

왼돌이물달팽이

Physa acuta

Correlation Analysis

· The relationship between the score and the water quality is not clear



Machine Learning

General Machine Learning



Ensemble Machine Learning





Machine Learning

eXtreme Gradient Boosting tree (XGBoost)

- Algorithm to make decision by mixing several models (tree boosting model)
- The algorithm that **learn**s results sequentially, with previous result and previous result affect the next model result in the current stage (additive training)
- This is a system in which the predicted value (Y) gradually approaches the target value (Ŷ) as the stage progresses for next stage
- So, this learn **weakly learning**, and gradually get closer to actual value unlike random forest model
- The method that weakly fit current data has a **high bias but low variance**.
- The high bias can be improved sufficiently by sequentially learning data weakly.



















Training technique (Random forest)

K-fold Cross validation



Training technique (XGBoost)

Training & Verification

- The k-fold cross validation & Grid search ٠
- Feature importance (parameter tuning) ٠
 - Max-depth : Maximum tree depth for base learners
 - Gamma : Minimum loss reduction required to make a further partition on a leaf node of the tree
- For optimization of parameters, max depth . applied from 1 to 10 with 1 intervals and Gamma applied from 0 to 4.5 with 0.5 intervals



FAI

TDI

BMI

Average accuracy of verification : 0.72 Average accuracy of verification : 0.66 Average accuracy of verification : 0.80





Watershed Health(FAI_spring)

2008 : Obs. (first), Predicted (second)



2011 : Obs. (first), Predicted (second)





2014 : Obs. (first), Predicted (second)





2009 : Obs. (first), Predicted (second)



2012 : Obs. (first), Predicted (second)



2015 : Obs. (first), Predicted (second)





2010 : Obs. (first), Predicted (second)





- ✓ 36 watersheds that have failed to predict.
- ✓ TP and NO₃ in failed watersheds were relatively greater than the overall average TP and NO₃ values.
 - TP was 19.3% greater than whole.
- ✓ PO_4 was 14.4% greater than whole.

Watershed Health(TDI_spring)

A B C D E E Very good Good Normal Bad Very bad

2008 : Obs. (first), Predicted (second)



2011 : Obs. (first), Predicted (second)





2014 : Obs. (first), Predicted (second)





2009 : Obs. (first), Predicted (second)



2012 : Obs. (first), Predicted (second)



2015 : Obs. (first), Predicted (second)



2010 : Obs. (first), Predicted (second)





- ✓ 41 watersheds that have failed to predict.
- ✓ TP and PO₄ in failed watersheds were relatively smaller than the whole average TP and PO₄ values.
 - TP was -30.1% smaller than whole.
- ✓ PO₄ was -29.6% smaller than whole.

Watershed Health(BMI_spring)

2008 : Obs. (first), Predicted (second)



2011 : Obs. (first), Predicted (second)





2014 : Obs. (first), Predicted (second)





- 2009 : Obs. (first), Predicted (second)
 - 2012 : Obs. (first), Predicted (second)



2015 : Obs. (first), Predicted (second)



2010 : Obs. (first), Predicted (second)

Very good Good Normal Bad Very bad





E

2013 : Obs. (first), Predicted (second)



- ✓ 25 watersheds that have failed to predict.
- ✓ TP and PO₄ in failed watersheds were relatively greater than the overall flow and PO₄ values.
 - TP was 66.1% greater than whole.
- ✓ PO_4 was 81.8% greater than whole.

Watershed Health(FAI_autumn) Very good Good Normal Bad Very bad

2008 : Obs. (first), Predicted (second)



2011 : Obs. (first), Predicted (second)





2014 : Obs. (first), Predicted (second)







2012 : Obs. (first), Predicted (second)



2015 : Obs. (first), Predicted (second)



2010 : Obs. (first), Predicted (second)







- 19 watersheds that have failed to predict.
- NH₄ and NO₃ in failed watersheds were relatively greater than the overall values.
- NH_4 was 25.9% greater than whole. NO_3 was 28.0% greater than whole.
- Flow was -66.3% smaller than whole.

Watershed Health(TDI_autumn) Very good Good Normal Bad Very bad

2009 : Obs. (first), Predicted (second)

2008 : Obs. (first), Predicted (second)



2011 : Obs. (first), Predicted (second)





2014 : Obs. (first), Predicted (second)







2015 : Obs. (first), Predicted (second)









- 45 watersheds that have failed to predict.
- TP and NH_4 in failed watersheds were relatively smaller than the overall values.
- TP was -23.4% smaller than whole. NH_4 was -31.6% smaller than whole.
- Flow was -38.6% smaller than whole.

Watershed Health(BMI_autumn) Wery good Good Normal Bad Very bad

2008 : Obs. (first), Predicted (second)



2011 : Obs. (first), Predicted (second)





2014 : Obs. (first), Predicted (second)







2012 : Obs. (first), Predicted (second)



2015 : Obs. (first), Predicted (second)



2010 : Obs. (first), Predicted (second)









- 23 watersheds that have failed to predict.
- NO_3 and NH_4 in failed watersheds were relatively greater and smaller than the overall values.
- NO_3 was 50.6% greater than whole. NH_4 was -29.8% smaller than whole. Flow was -35.3% smaller than whole.

Findings and Future Researches

This study was to develop XGBoost which is one of ensemble machine learning algorithms (Random forest vs. XGBoost) for AEH indices prediction using SWAT results.



- We could predict AEH indices of ungauged streams via XGBoost with SWAT water quality results at a standard watershed scale.
- For further research, we will include environmental components such as river width, bankside cover treatment with porous material or vegetation cover etc..