

2025년도 공공기관 용역과제  
AI개발 수행내역서

과제명	AI 기반 정수장 수질(pH) 예측모델 개발 및 음용 판단 여부 시스템 구축
담당자	류건우

2025년 12 월 28 일

## AI개발 수행내용

### 1. 사업과제 : AI 기반 정수장 수질(pH) 예측모델 개발 및 시각화 시스템 구축

#### 2. 개요 및 현황

##### 2.1 추진배경 및 목적

- 최근 기후 변화와 원수 수질 변동성 증가로 인해 정수장 운영의 불확실성이 확대되고 있음
- 기존 정수장 운영은 기준치 초과 여부에 대한 사후 대응 중심으로 이루어지고 있어, 약품 사용의 비효율이 발생
- 공공 수질 데이터(경도, 탁도, 잔류염소 등)를 활용한 AI 기반 예측 모델을 통해 pH 변화를 사전에 예측하고 약품 운전 부담 및 음용 안전성을 선제적으로 판단할 필요가 있음

##### 2.2 과제 범위

과제구분		내용
AI	AI기반 수질예측모델 구현	공공 정수장 수질 데이터 수집 및 DB 구축
		수질 항목 전처리 및 결측치 처리
		상관관계 분석 및 주요 변수 선정
		RandomForest 기반 PH 예측 모델 학습
		RMSE, MAE, R <sup>2</sup> 기반 모델 성능 평가
		예측모델 웹기반 시스템 구축
		테스트
시각화	수질 입력 기반 pH 시각화	수질 입력 기반 pH 예측 시연 화면 구성
		연·월·지역별 pH 예측 결과 시각화
		예측 결과에 대한 오차 및 잔차 분석
		테스트
		통합테스트 및 시운전

## 2.3 과제 추진 방법

### 1) 구축 대상 선정 기준

#### ○ 데이터 접근성 및 활용성

- 공공데이터 포털을 통해 제공되는 정수장 수질 데이터를 활용하여 데이터 수집 및 관리 용이
- 이미 구축된 api를 데이터베이스로 활용
- pH, 탁도, 잔류염소 등 정수 처리 공정에 직접적인 영향을 미치는 핵심 수질 변수가 포함되어 있어 예측 모델 학습에 적합함

#### ○ 예측모델 개발 효율성

- 수질 데이터는 연·월 단위로 정리되어 있어 시계열적 패턴 학습 및 예측이 가능
- 수질 항목 간 상관관계 분석을 통해 모델 입력 변수 축소 및 해석 가능성 확보

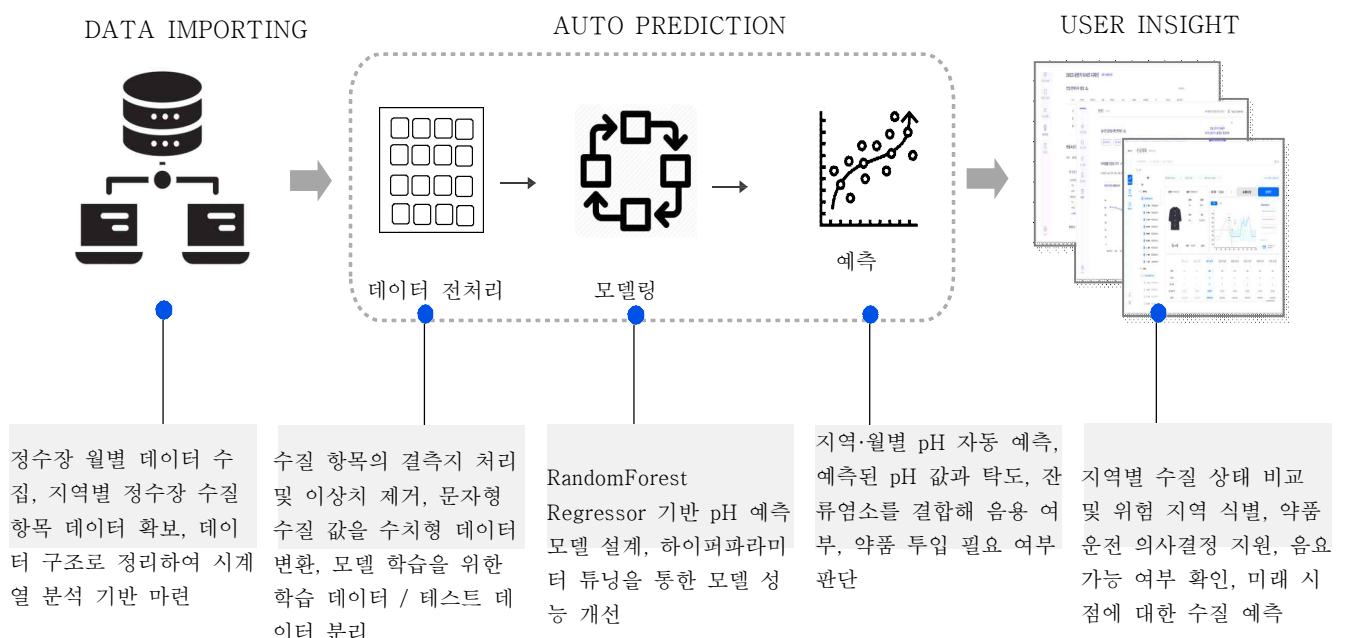
#### ○ 수질 관리 기여도 및 유지 비용 절감

- pH 예측 결과를 활용하여 약품 투입 여부 및 운전 조건 판단 가능
- 불필요한 약품 사용 감소를 통해 운영 비용 절감 효과 기대
- 수질 이상 징후를 사전에 파악함으로써 음용 안전성 확보 및 사회적 비용 감소에 기여

### 2) AI 예측 분석모델 적용 대상

	수집 데이터	예측모델인자(독립변수)	AI예측 분석 대상
정수장 수질 관리	<ul style="list-style-type: none"> <li>- 정수장 월별 수질 측정 데이터</li> <li>- 지역별 정수장 운영 데이터</li> </ul>	<ul style="list-style-type: none"> <li>- 경도(HR)</li> <li>- 증발잔류물(RE)</li> <li>- 질산성질소(MON)</li> <li>- 브롬산염(BRO)</li> <li>- 알루미늄(AL)</li> <li>- 클로로포름(CF)</li> <li>- 황산이온(SO), 탁도(TU), 잔류염소(RC)</li> </ul>	<ul style="list-style-type: none"> <li>- pH 예측</li> <li>- 약품 투입 필요 여부 판단</li> <li>- 음용 가능 여부 판단</li> <li>- 지역·월별 수질 상태 비교 분석</li> </ul>

### 3) AI 분석모델 구축 프로세스



## 연구개발 주요 결과물

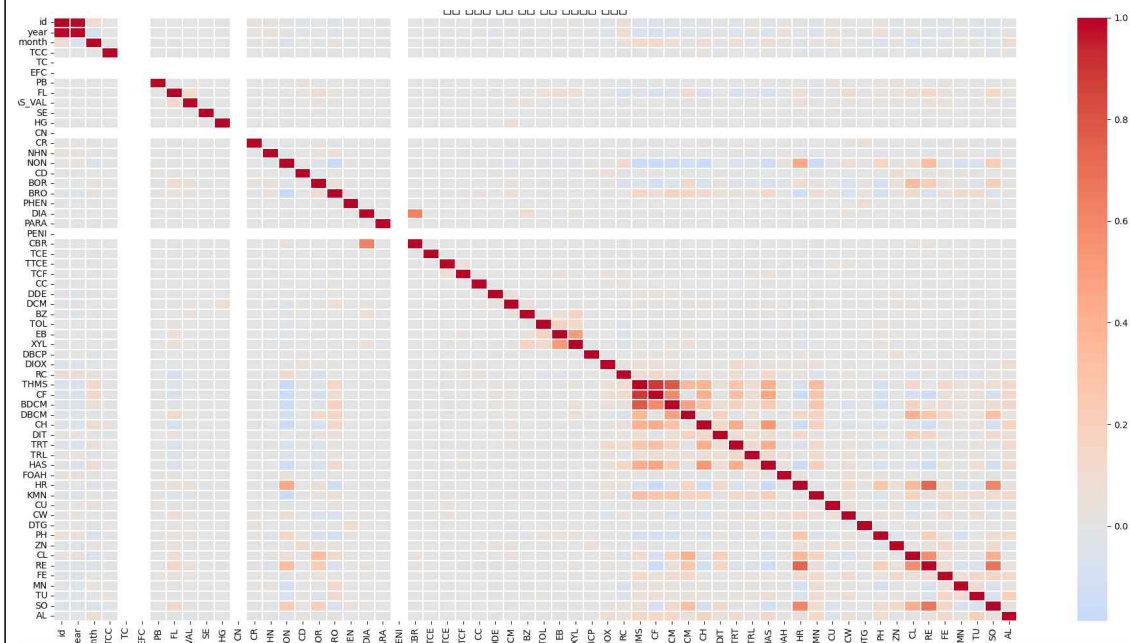
### 1. 데이터 수집

- 한국수자원공사 5년간 상수도법정수질정보(api) : 2020년 ~ 2025년

이름	타입	스키마
water_quality	INTEGER	CREATE TABLE water_quality (id INTEGER PRIMARY KEY AUTOINCREMENT, year REAL, month REAL, region TEXT, facility TEXT, TCC TEXT, TC TEXT, EFC TEXT, PB REAL, FL REAL, AS_VAL REAL, SE REAL, HG REAL, CN REAL, CR REAL, NHN REAL, NON REAL, CD REAL, BOR REAL, BRO REAL, PHEN REAL, DIA REAL, PARA REAL, PENI REAL, CBR REAL, TCE REAL, TTCE REAL, TCF REAL, CC REAL, DDE REAL, DCM REAL, BZ REAL, TOL REAL, EB REAL, XYL REAL, DBCP REAL, DIOX REAL, RC REAL, THMS REAL, CF REAL, BDCM REAL, CH REAL, DIT REAL, TRT REAL, TRL REAL, HAS REAL, FOAH REAL)
id	INTEGER	"id" INTEGER
year	REAL	"year" REAL
month	REAL	"month" REAL
region	TEXT	"region" TEXT
facility	TEXT	"facility" TEXT
TCC	TEXT	"TCC" TEXT
TC	TEXT	"TC" TEXT
EFC	TEXT	"EFC" TEXT
PB	REAL	"PB" REAL
FL	REAL	"FL" REAL
AS_VAL	REAL	"AS_VAL" REAL
SE	REAL	"SE" REAL
HG	REAL	"HG" REAL
CN	REAL	"CN" REAL
CR	REAL	"CR" REAL
NHN	REAL	"NHN" REAL
NON	REAL	"NON" REAL
CD	REAL	"CD" REAL
BOR	REAL	"BOR" REAL
BRO	REAL	"BRO" REAL
PHEN	REAL	"PHEN" REAL
DIA	REAL	"DIA" REAL
PARA	REAL	"PARA" REAL
PENI	REAL	"PENI" REAL
CBR	REAL	"CBR" REAL
TCE	REAL	"TCE" REAL
TTCE	REAL	"TTCE" REAL
TCF	REAL	"TCF" REAL
CC	REAL	"CC" REAL
DDE	REAL	"DDE" REAL
DCM	REAL	"DCM" REAL
BZ	REAL	"BZ" REAL
TOL	REAL	"TOL" REAL
EB	REAL	"EB" REAL
XYL	REAL	"XYL" REAL
DBCP	REAL	"DBCP" REAL
DIOX	REAL	"DIOX" REAL
RC	REAL	"RC" REAL
THMS	REAL	"THMS" REAL
CF	REAL	"CF" REAL
BDCM	REAL	"BDCM" REAL
CH	REAL	"CH" REAL
DIT	REAL	"DIT" REAL
TRT	REAL	"TRT" REAL
TRL	REAL	"TRL" REAL
HAS	REAL	"HAS" REAL
FOAH	REAL	"FOAH" REAL

### 1. 데이터 분석

#### 2.1 수질데이터 상관관계(Heatmap)



Heatmap showing the correlation matrix for the variables HR, RE, NON, BRO, AL, CF, SQ, TU, PC, and PH. The color scale ranges from 0.0 (light blue) to 1.0 (dark blue). The diagonal elements are all 1.0. The off-diagonal elements represent the Pearson correlation coefficients between the variables.

	HR	RE	NON	BRO	AL	CF	SQ	TU	PC	PH
HR	1.00	0.73	0.45	-0.13	0.06	-0.15	0.61	-0.04	0.04	0.27
RE	0.73	1.00	0.34	-0.07	0.03	-0.15	0.68	-0.05	0.01	0.21
NON	0.45	0.34	1.00	-0.18	-0.02	-0.16	0.23	-0.08	0.12	0.15
BRO	-0.13	-0.07	-0.18	1.00	0.02	0.07	-0.05	0.10	-0.00	-0.12
AL	0.06	0.03	-0.02	0.02	1.00	0.14	0.03	0.20	0.03	0.11
CF	-0.15	-0.15	-0.16	0.07	0.14	1.00	-0.09	0.06	0.12	-0.10
SQ	0.61	0.68	0.23	-0.05	0.03	-0.09	1.00	-0.06	0.03	0.10
TU	-0.04	-0.05	-0.08	0.10	0.20	0.06	-0.06	1.00	-0.04	0.00
PC	0.04	0.01	0.12	-0.00	0.03	0.12	0.03	-0.04	1.00	0.01
PH	0.27	0.21	0.15	-0.12	0.11	-0.10	0.10	0.00	0.01	1.00

## 2. 데이터 학습 및 모델정의

### 3.1 모델정의 및 컴파일

- 모델 정의 : RandomForest 모델

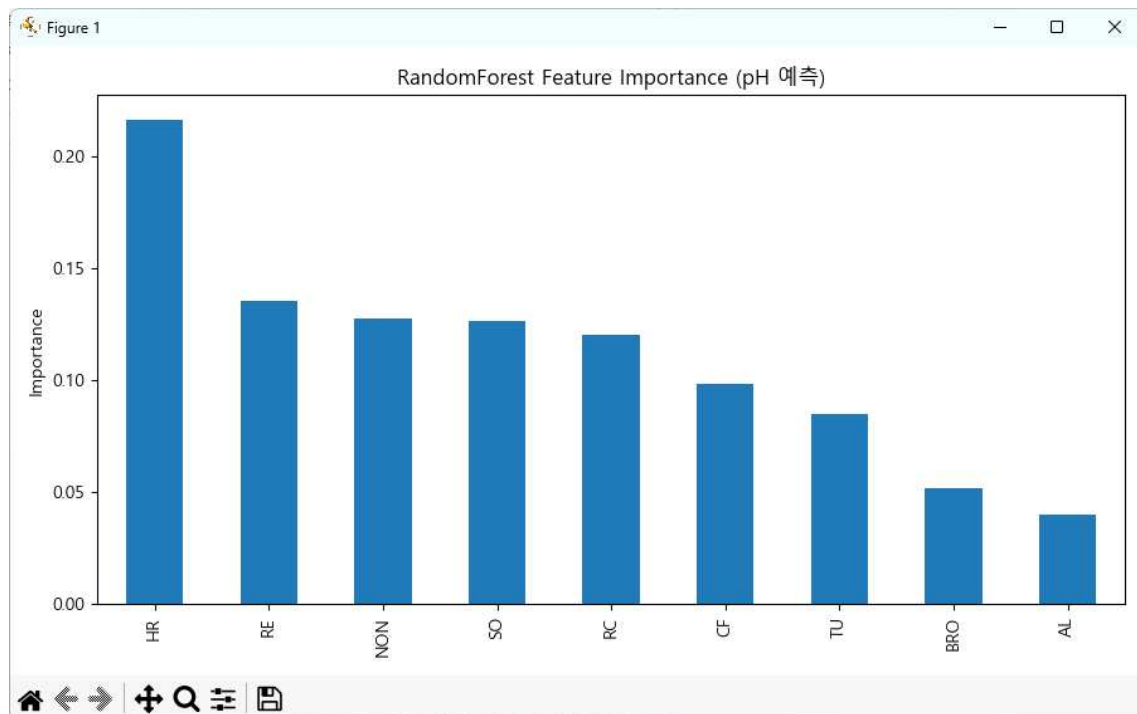
```
87 rf = RandomForestRegressor(  
88     n_estimators=300,  
89     min_samples_leaf=5,  
90     random_state=RANDOM_STATE,  
91     n_jobs=-1  
92 )  
93  
94 rf.fit(X_train, y_train)  
95
```

### 3.2 모델학습 및 학습 시각화

- 모델 학습 및 성능 지표

```
(.venv) PS C:\Users\kwr51\Desktop\test> & C:/Users/kwr51/Desktop/test/.venv/Scripts/python.exe c:/Users/kwr51/Desktop/test/forestsave.py  
모델 학습 데이터 수: 29798  
  
RandomForest 회귀 성능  
RMSE : 0.3112  
MAE : 0.2348  
R2 : 0.3723
```

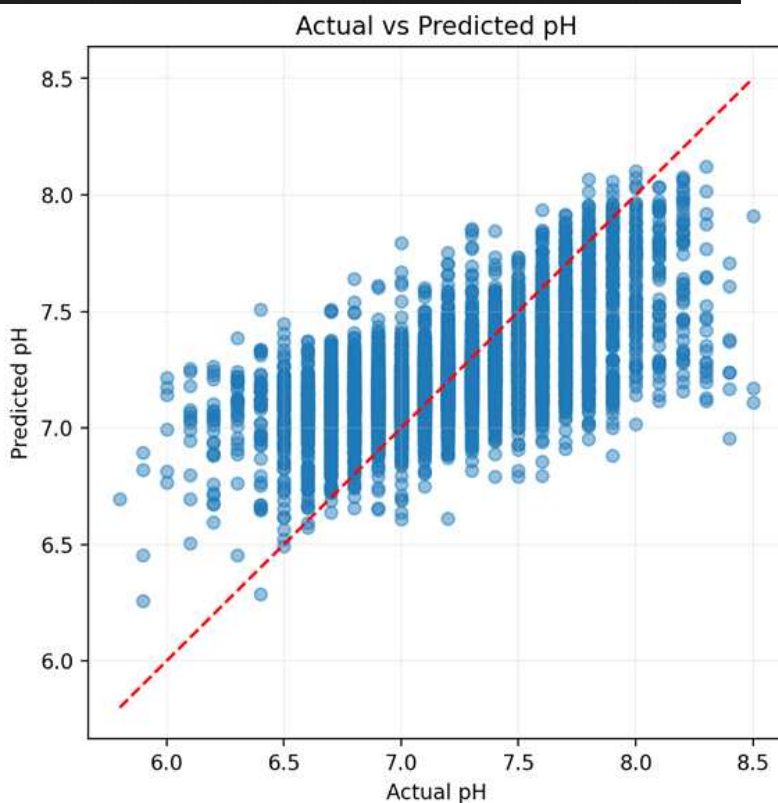
- 공공 수질 데이터에 있는 pH는 미세한 변화만으로도 의미가 있고, 실제 데이터 특성상 대부분의 정수장의 데이터 값은 정상 범주에 속해있기 때문에 R2 값이 낮게 나옴.



### 3.3 모델 예측

#### ○ 예측값 vs 실제값 비교

```
# =====  
# 9. 예측 vs 실제 시각화  
# =====  
plt.figure(figsize=(6, 6))  
plt.scatter(y_test, y_pred, alpha=0.5)  
plt.plot([y_test.min(), y_test.max()],  
         [y_test.min(), y_test.max()],  
         'r--')  
  
plt.xlabel("Actual PH")  
plt.ylabel("Predicted PH")  
plt.title("Actual vs Predicted PH")  
plt.tight_layout()  
plt.show()
```



#### ○ 실제 pH가 증가할수록 예측 pH도 함께 증가

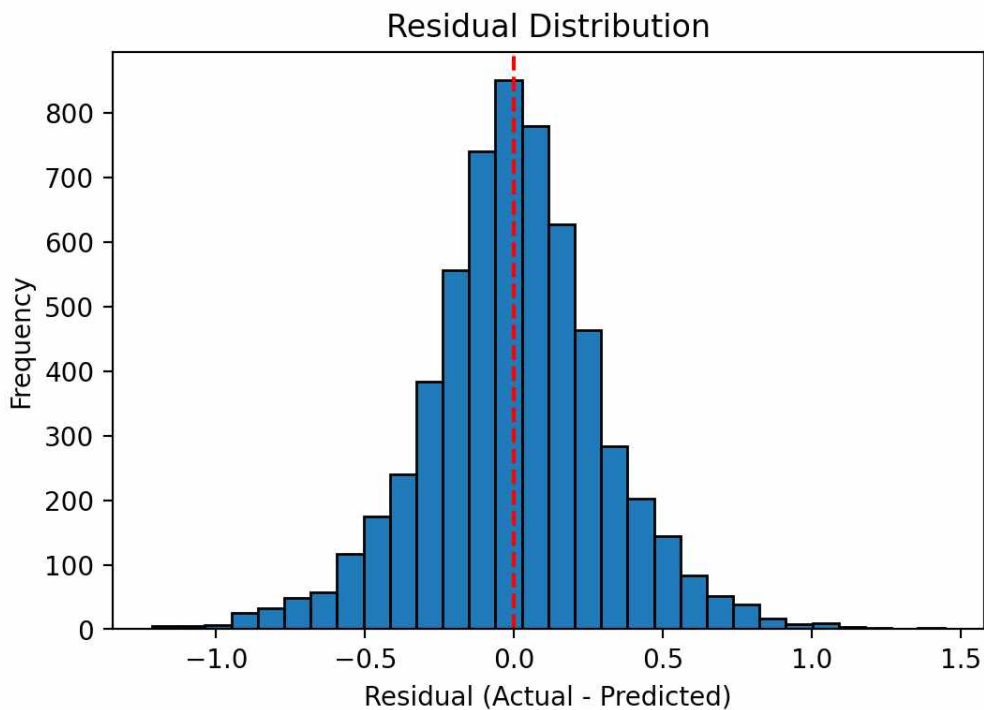
pH 값이 7.0 ~ 7.8 구간에 집중.

실제 수질 데이터 대부분이 정상 범위에 분포하기 때문.

그래서 극단값(6.0 이하, 8.5 이상)은 상대적으로 적음.

오차패턴, 일부 점들이 기준선 위 아래로 벗어나 있는 이유는,  
자연 환경 요인, 측정 오차, 약품 투입량과 수온, 유량 등의 환경적인 요인

○ Residual Distribution(잔차 분포)



○ 평균 잔차가 0에 가까움, 모델이 특정 방향으로 과대/과소 예측하지 않은 것.

○ 대부분의 잔차가 -0.3 ~ +0.3 구간에 집중  
극단값인 1.0 이상은 소수

○ 이는 일시적인 수질 이상, 측정 오차 등 외부요인에 의한 것.

### 3. 프로토타이핑(화면)

#### 4.1 수질 입력 후 pH 예측

- 각 요소에 영향을 받은 pH 예측
  - 경도, 증발잔류물 등 입력

pH 예측 시연 2026년 예측 분석 활용 안전판단

#### 수질 입력 → pH 예측



pH 예측

예측 pH : 7.16

정상 또는 미세 조정 수준

※ 본 예측은 과거 학습된 수질 범위 내 상대적 변화에 기반합니다.



## 4.2 미래 연도 예측 분석

- 예측 대상 연도, 예측 대상 월 입력

### 연도·월 선택 지역별 pH 예측 (미래 예측)

예측 대상 연도

2026

예측 대상 월

1

#### 2026년 1월 지역별 pH 예측

지역	예측 pH	탁도(TU)	잔류염소(RC)	약품 판단	음용 안전
15 제주특별자치도		7.88	0.05	0.48 ● 정상 또는 미세 조정 수준	☑ 음용 가능
3 인천광역시		7.41	0.06	0.77 ● 정상 또는 미세 조정 수준	☑ 음용 가능
7 경기도		7.34	0.07	0.75 ● 정상 또는 미세 조정 수준	☑ 음용 가능
10 충청남도		7.32	0.15	0.67 ● 정상 또는 미세 조정 수준	☑ 음용 가능
8 강원특별자치도		7.31	0.08	0.73 ● 정상 또는 미세 조정 수준	☑ 음용 가능
9 충청북도		7.31	0.07	0.77 ● 정상 또는 미세 조정 수준	☑ 음용 가능
6 울산광역시		7.29	0.07	0.64 ● 정상 또는 미세 조정 수준	☑ 음용 가능
13 경상북도		7.27	0.1	0.8 ● 정상 또는 미세 조정 수준	☑ 음용 가능
0 서울특별시		7.23	0.04	0.42 ● 정상 또는 미세 조정 수준	☑ 음용 가능
12 전라남도		7.2	0.15	0.69 ● 정상 또는 미세 조정 수준	☑ 음용 가능

#### 약품 사용 부담 증가 예상 TOP5

지역	예측 pH	탁도(TU)	잔류염소(RC)	약품 판단	음용 안전
15 제주특별자치도		7.8800	0.0500	0.4800 ● 정상 또는 미세 조정 수준	☑ 음용 가능
3 인천광역시		7.4100	0.0600	0.7700 ● 정상 또는 미세 조정 수준	☑ 음용 가능
7 경기도		7.3400	0.0700	0.7500 ● 정상 또는 미세 조정 수준	☑ 음용 가능
10 충청남도		7.3200	0.1500	0.6700 ● 정상 또는 미세 조정 수준	☑ 음용 가능
8 강원특별자치도		7.3100	0.0800	0.7300 ● 정상 또는 미세 조정 수준	☑ 음용 가능

※ 본 예측은 2025년 동일 월 수질 조건을 기반으로 2026년 1월을 가장한 시나리오 예측입니다.

## 4.3 음용 안전 판단 분석

pH 예측 시연 2026년 예측분석 음용 안전 판단

### 수돗물 음용 안전 판단

판단 기준 (요약)

- pH: 6.5 ~ 8.5
- 탁도(TU):  $\leq 1.0$  NTU
- 잔류염소(RC):  $\geq 0.4$  mg/L



## 4. 결론 및 제언

○ pH 예측 성능( $R^2$  약 0.31, RMSE 약 0.3 수준)을 확보하였으나, 일부 구간에서 오차가 발생하여 추가적인 보완이 필요함

→ 운영 변수(약품 투입량), 환경 변수(계절·강수량 등) 추가 필요

○ 지역·월별 예측 결과에서 pH 값의 변동 폭이 존재함

→ 원수 수질 특성 및 정수장 운영 조건 차이에 따른 영향으로 판단됨

○ Feature Importance 분석 결과, 탁도(TU), 잔류염소(RC) 등 공정 핵심 지표가 pH 예측에 주요한 영향을 미침

○ 본 분석을 통해 약품 사용 부담이 예상되는 지역 및 시점을 사전에 식별 가능

→ 지역·계절별 특성을 반영한 세분화된 데이터셋 구축 필요

→ 향후 약품 사용량 예측 및 정수장 운영 최적화 전략 수립에 활용 가능