

FORECASTING DAMAGES CAUSED BY PIPE BURSTS

PROBLEM

Oil pipeline leaks and spills can have severe and far-reaching consequences, both environmentally and economically.

Environmental Damage

- Soil and Groundwater Contamination
- Water Pollution
- Air Pollution
- Habitat Destruction

Economic Loss

- Cleanup Costs
- Private/Public Property Damage Costs
- Legal Liabilities
- Environmental Remediation

PURPOSE

We will be focused on economic loss and our aim is to develop a predictive model that can estimate the potential cost of damage based on various features and factors present in the dataset.



DATA OVERVIEW

TARGET VARIABLES

Property
Damage Costs

Lost
Commodity
Costs

Private/Public
Property
Damage Costs

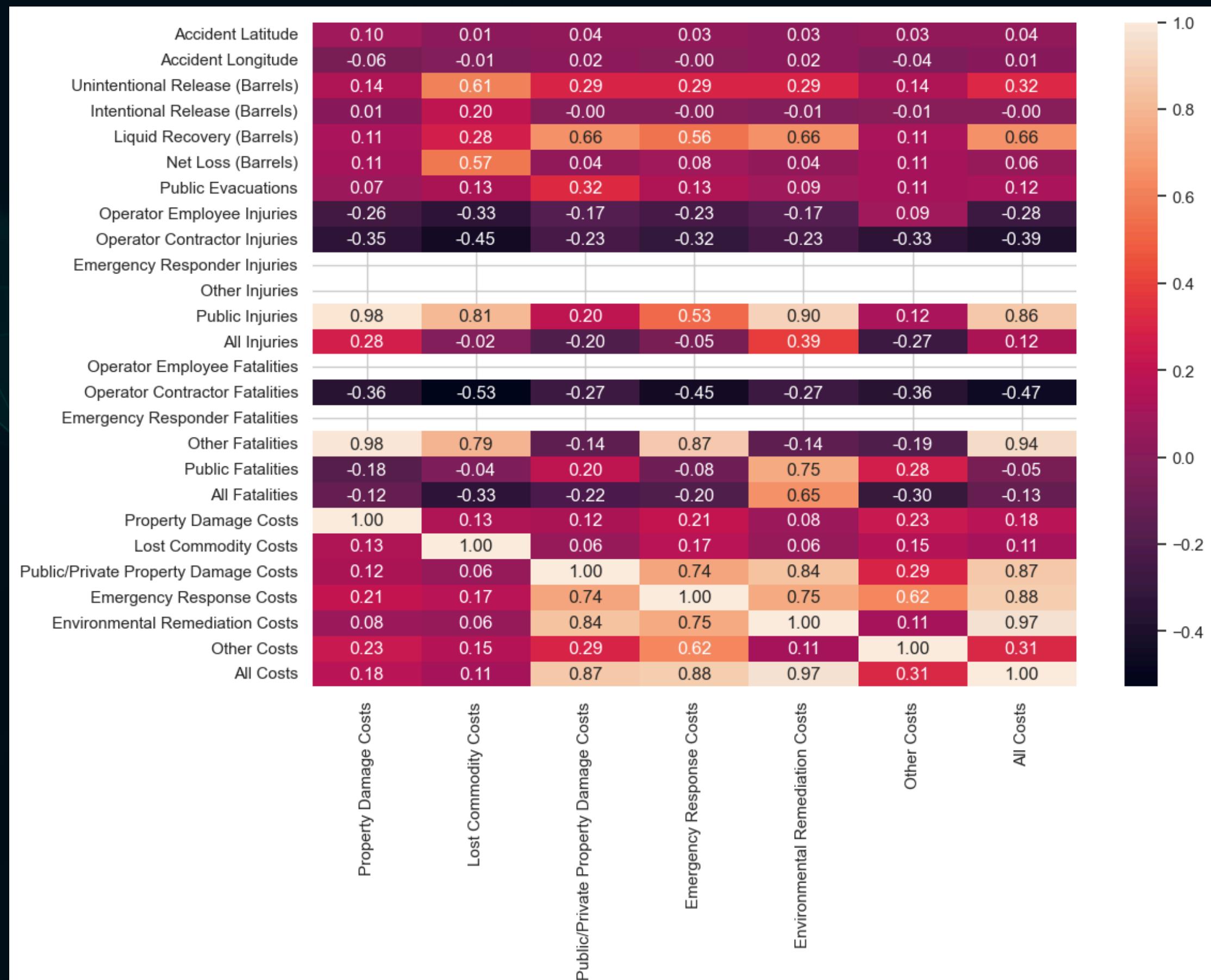
Emergency
Response
Costs

Other Costs

Environmental
Remediation
Costs

All Costs

CORRELATION MATRIX



EDA RESULTS

- 1 A lot of outliers
- 2 A lot of missing data
- 3 Class imbalance
- 4 Target variables distribution is skewed to right
- 5 Slight correlation between target variables and features

DATA CLEANING

- 1 Droping of columns (Identification, a lot of unique values)
- 2 Deletion of rows with missing values in target variables
- 3 Imputing of missing values in other variables

DATA PREPROCESSING



CATEGORICAL DATA

feature

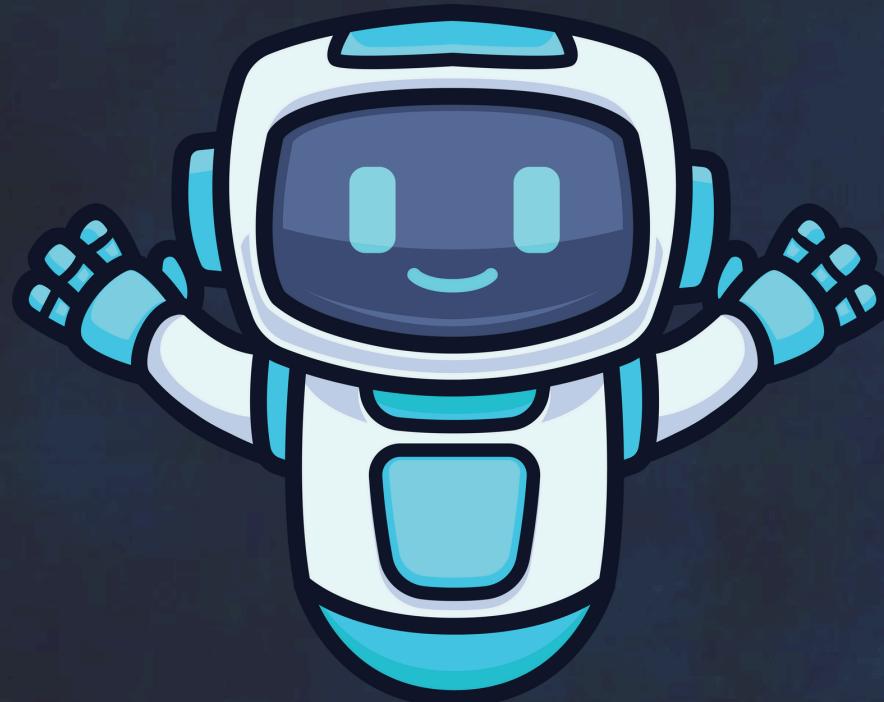
Pipeline Type
Accident State

Liquid Subtype
Liquid Name

preprocessing

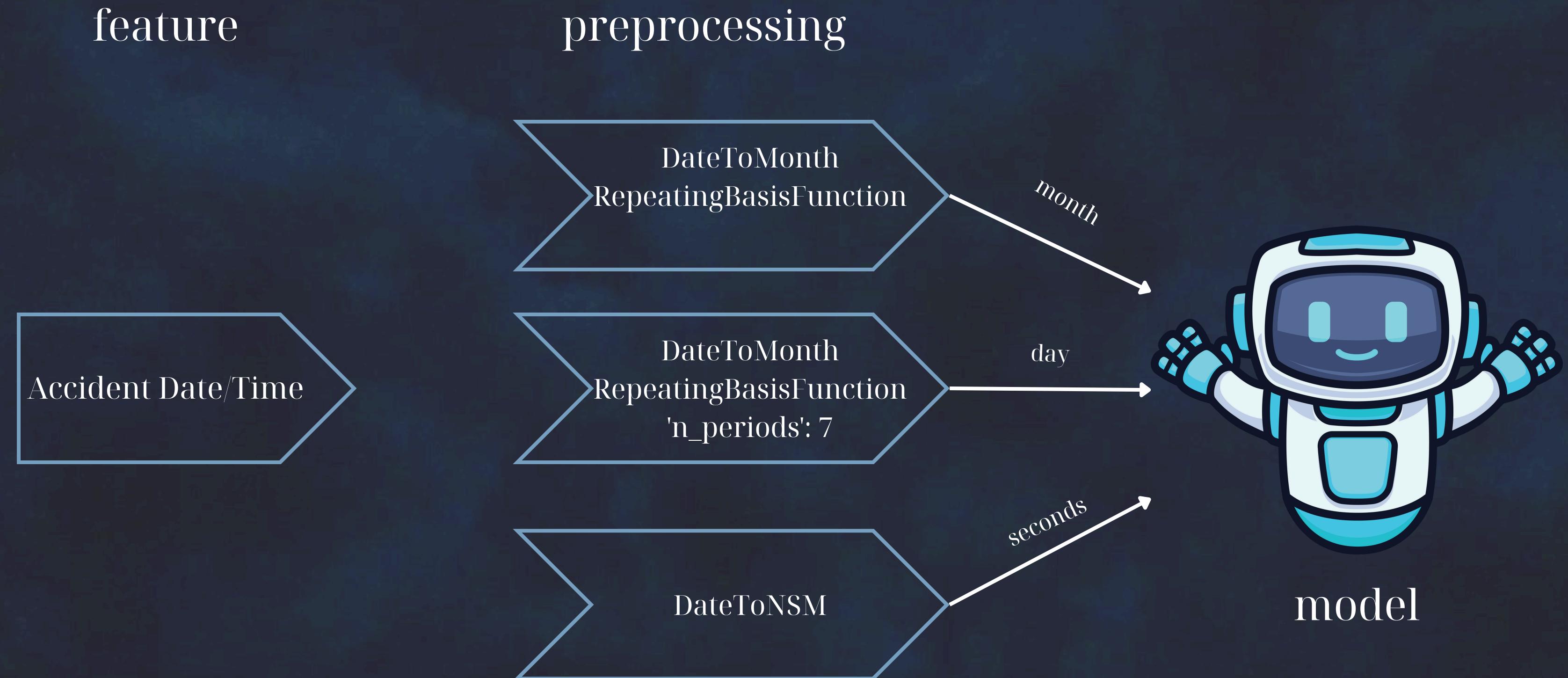
SimpleImputer
(most freq.)

BinaryEncoder



model

DATETIME



NUMERICAL DATA

feature

Pipeline Shutdown
Shutdown DT
Restart DT

Public Evacuations
Unintentional Release
Liquid Recovery
Net Loss

Injuries&Fatalities

Intentional Release

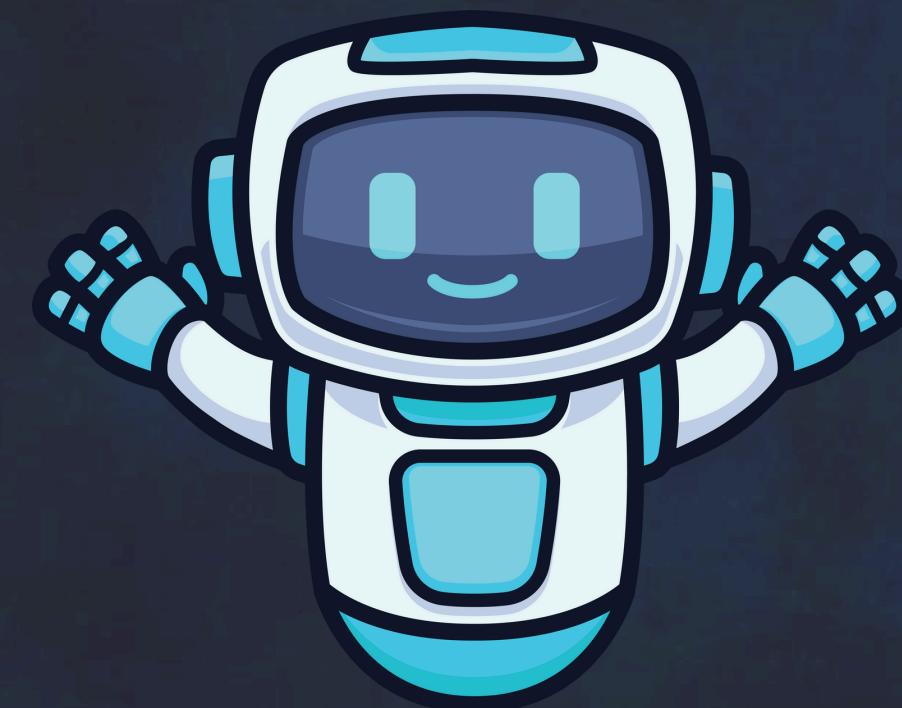
preprocessing

DateDiff

KNNImputer
FunctionTransformer
PolynomialFeatures

None to 0

BinaryEncoder



model

METRICS



For evaluating the performance of our regression model, we have selected the following metrics:

Metrics

R-square

MSE

MAE

By considering both R-squared and error metrics (MSE and MAE), we can make informed decisions about model selection and potential improvements.

MODELS



Hyperparameters tuning using Cost-Frugal Optimization

XGBoost

Hyperparameters:

- n_estimators,
- max_leaves,
- min_child_weight,
- learning_rate,
- subsample,
- colsample_bylevel,
- colsample_bytree,
- reg_alpha,
- reg_lambda.

RandomForest

Hyperparameters:

- n_estimators,
- max_features,
- max_leaves

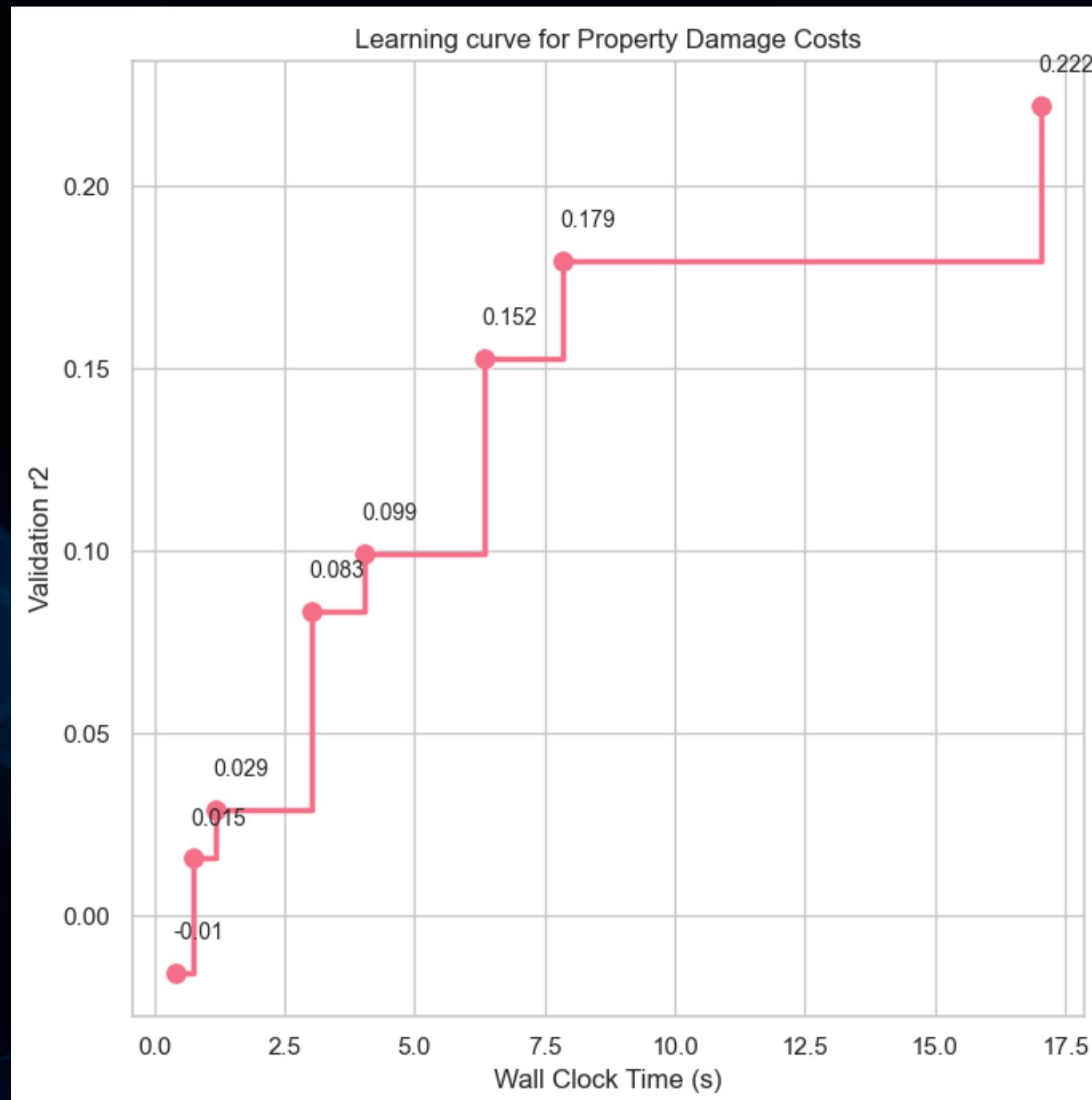
ExtraTree

Hyperparameters:

- n_estimators,
- max_features,
- max_leaves

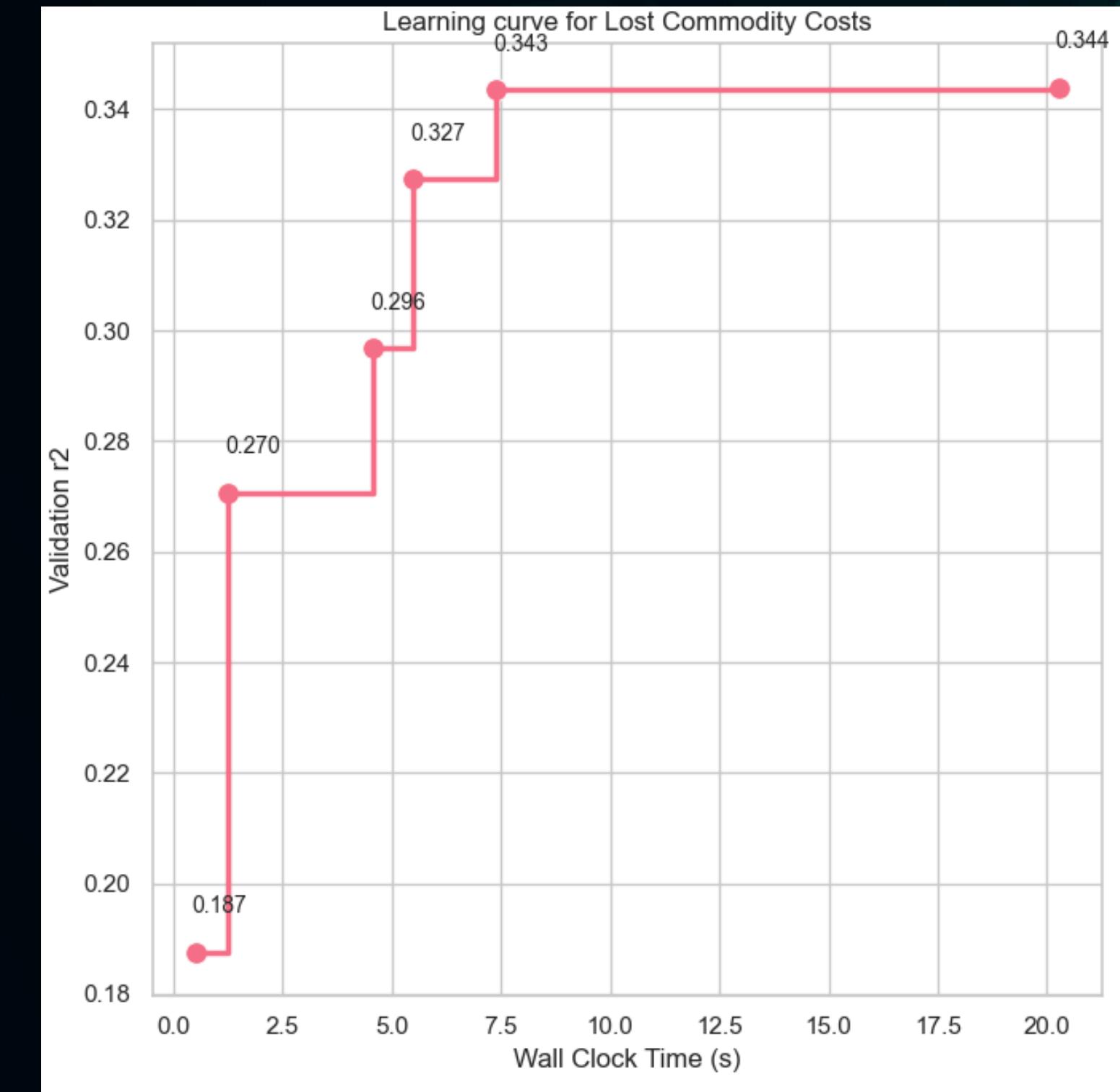
Property Damage Costs

XGBoost

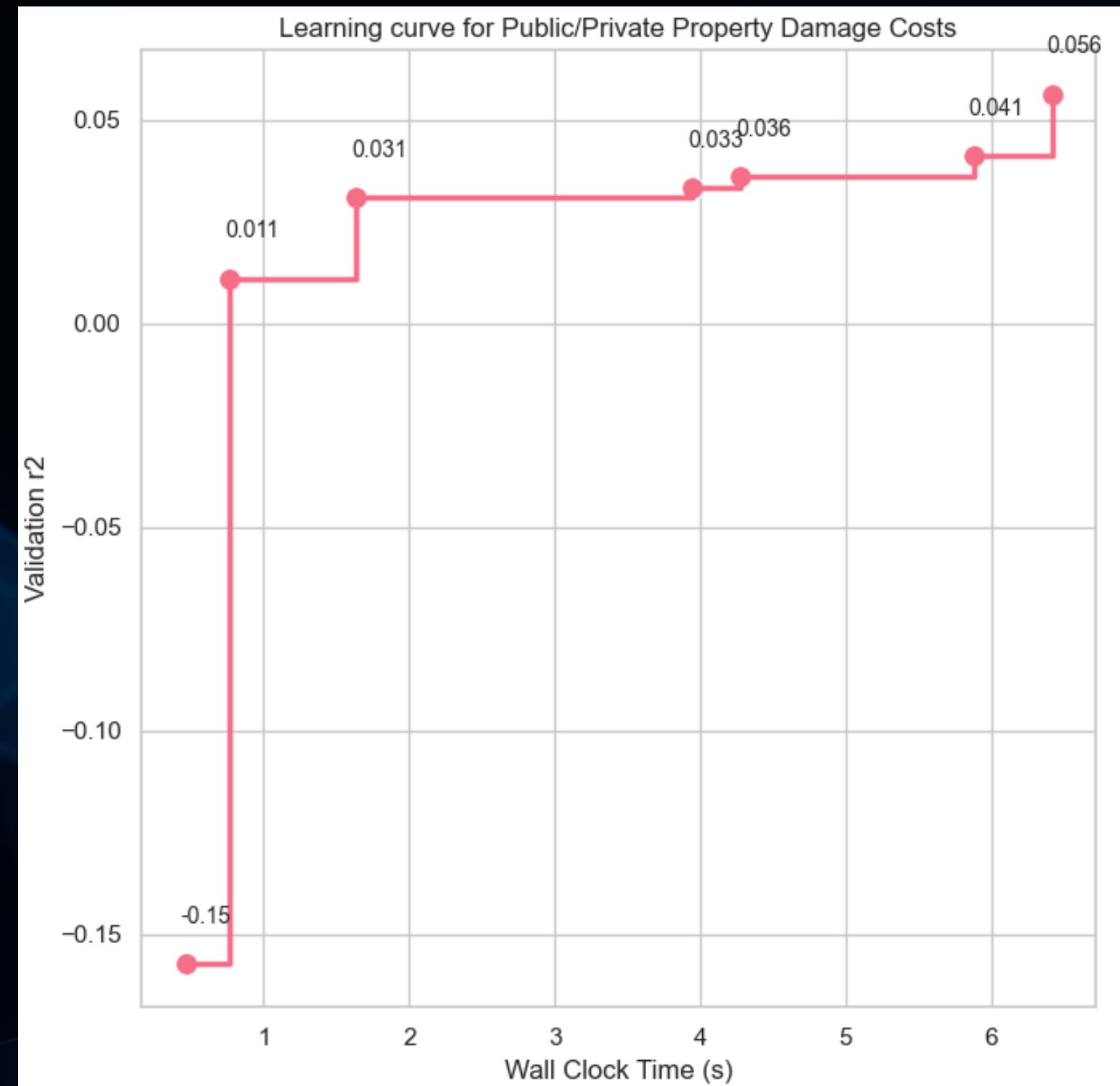


Lost Commodity Costs

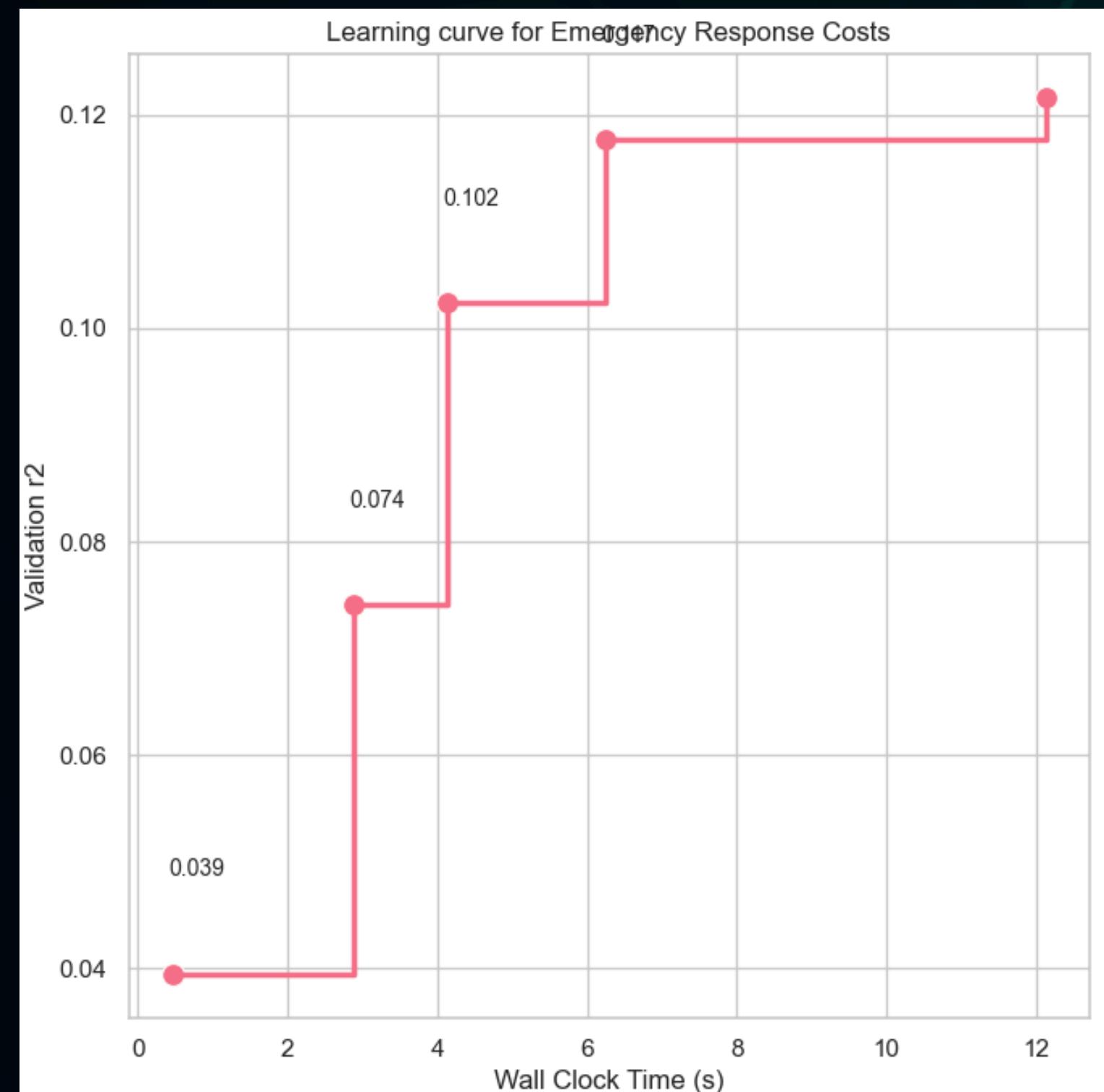
XGBoost



Public/Private Property Damage Costs XGBoost

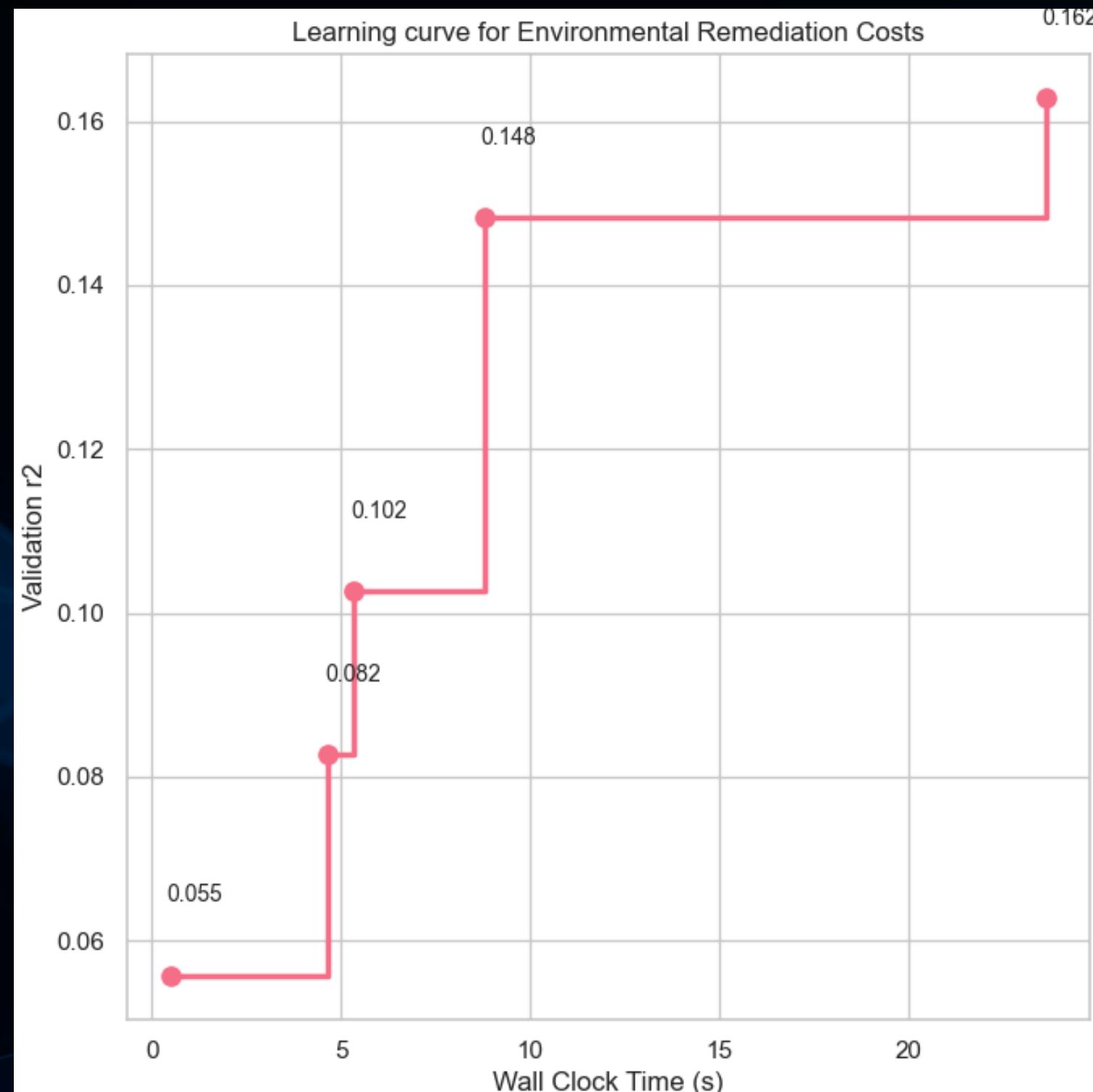


Emergency Response Costs XGBoost



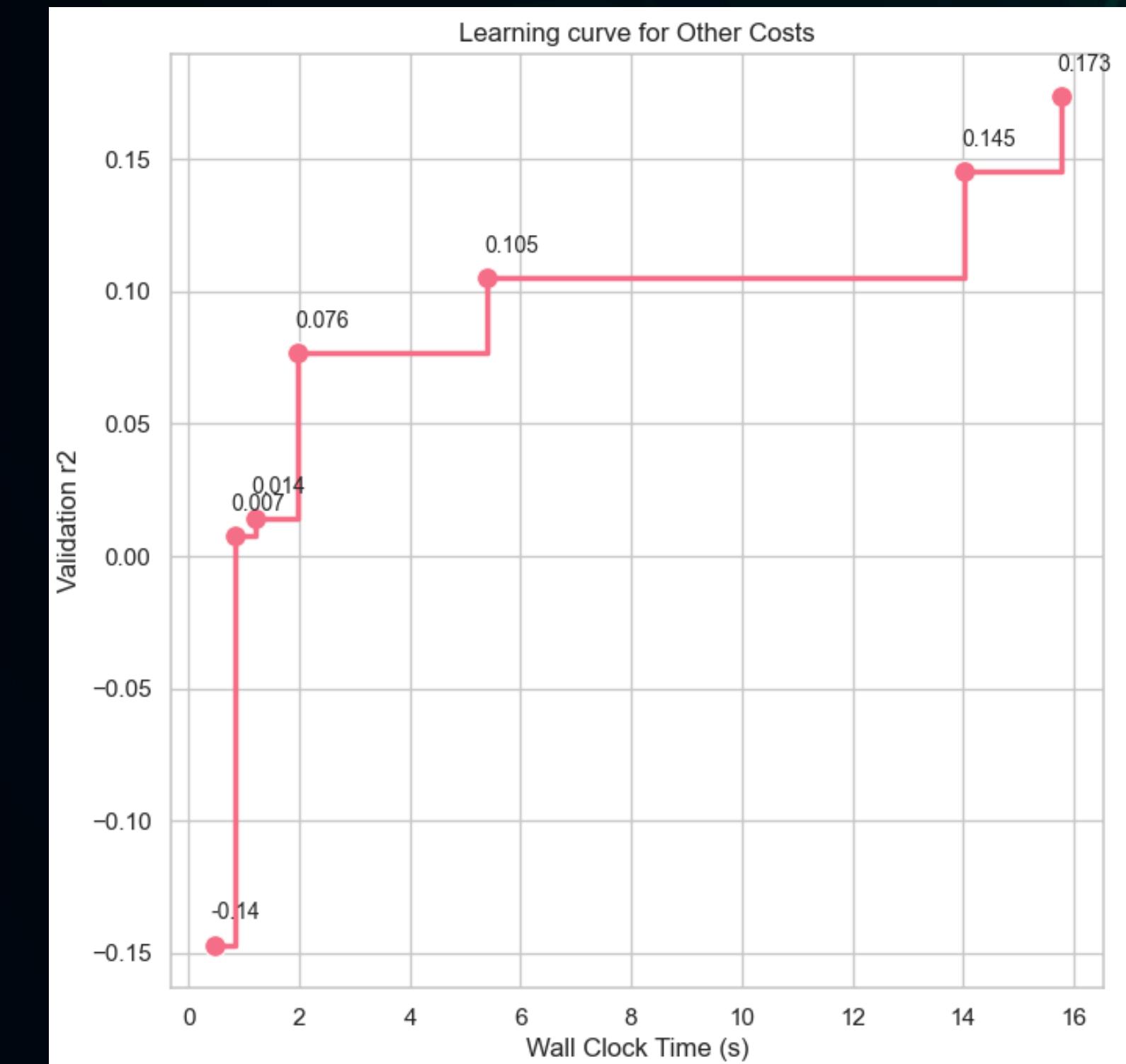
Environmental Remediation Costs

XGBoost

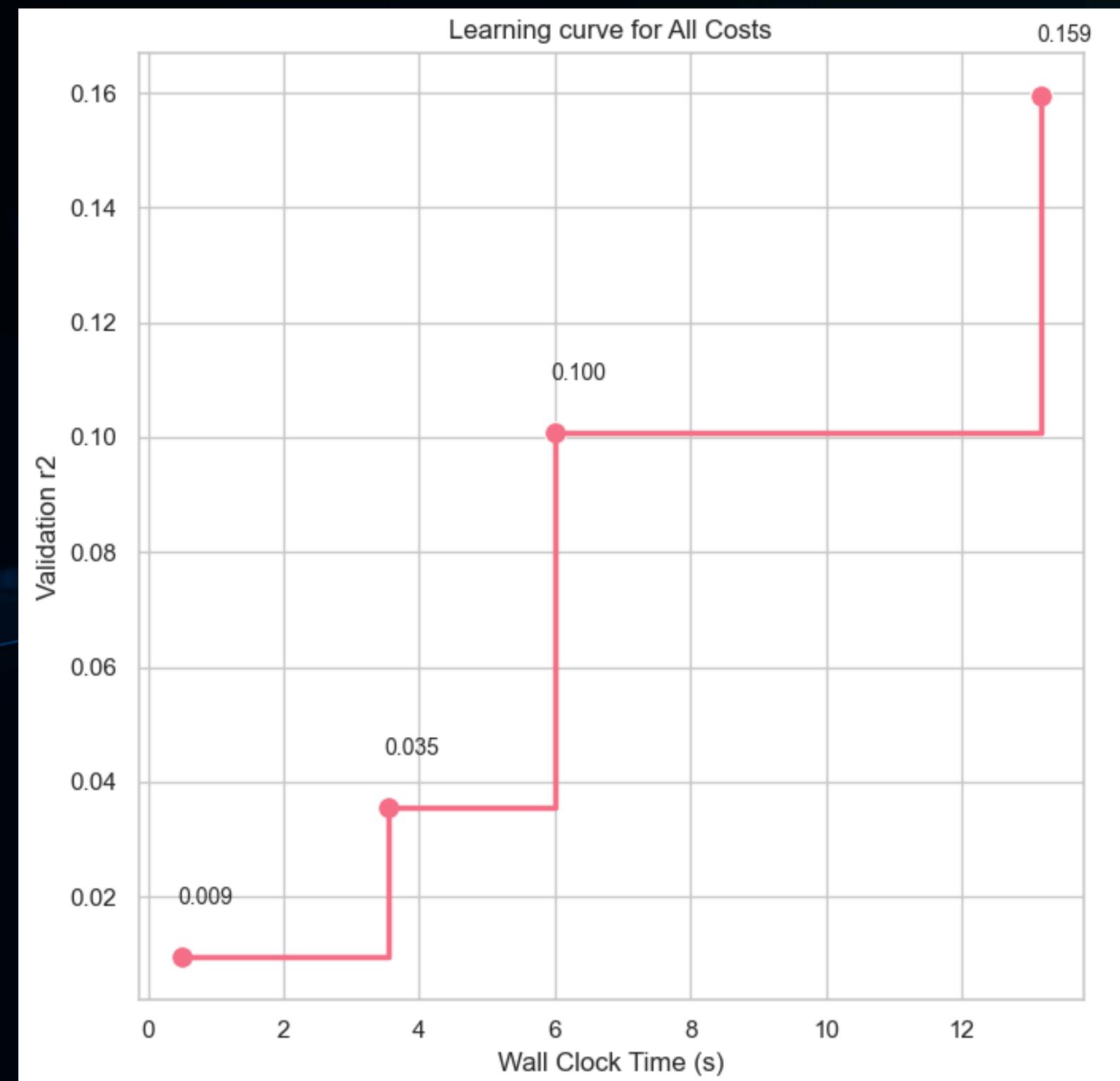


Other Costs

XGBoost

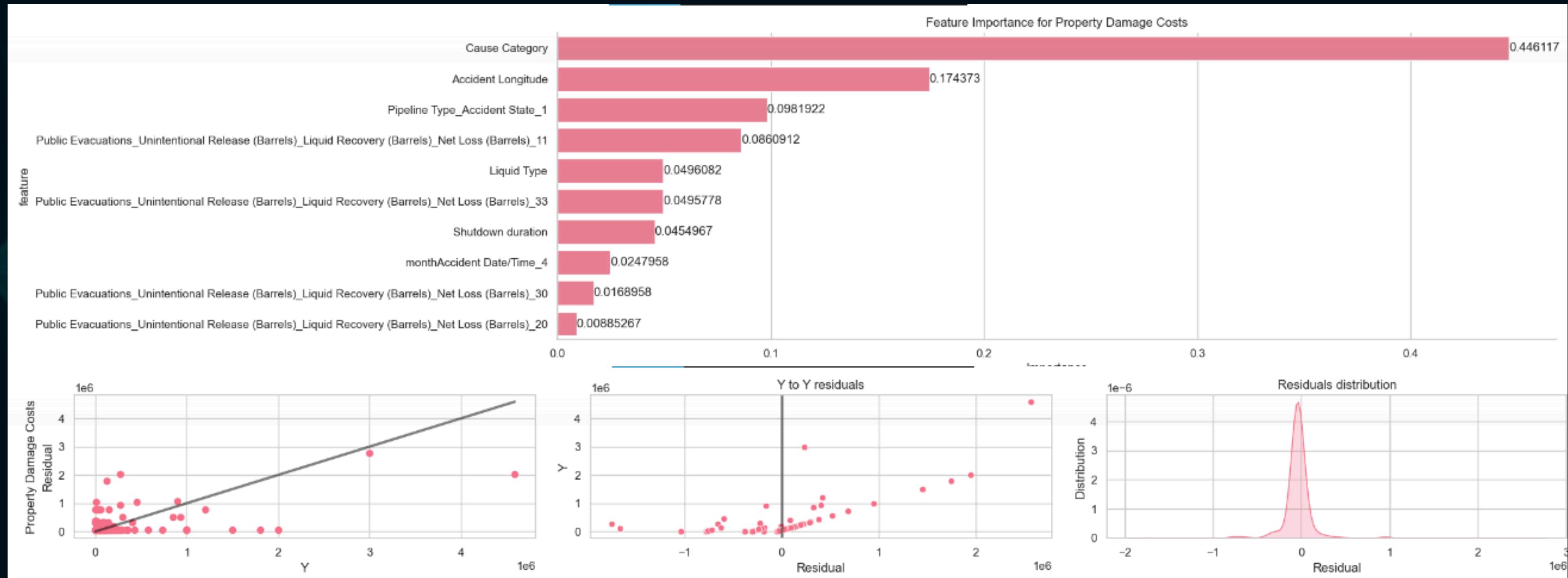


All Costs ExtraTree



RESULTS

PROPERTY DAMAGE COSTS

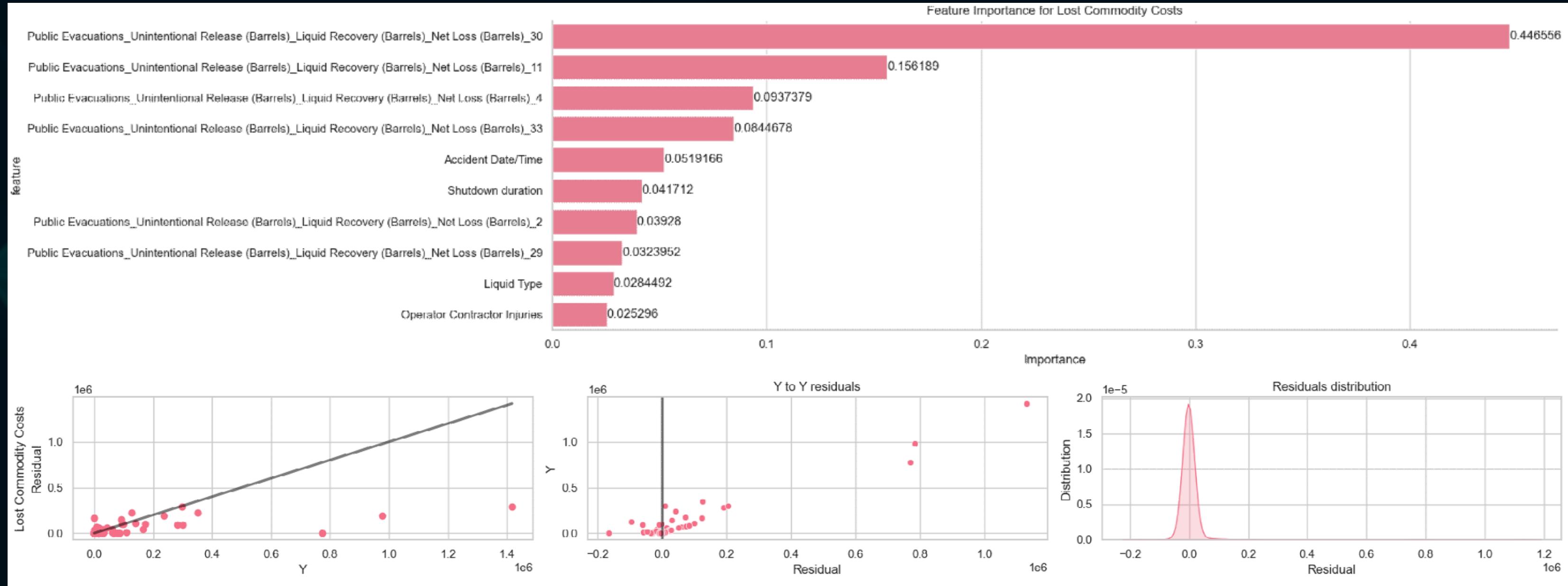


$R^2: 0.262$

$MSE: 6.4 \times 10^{10}$

$MAE: 1.02 \times 10^5$

LOST COMMODITY COSTS

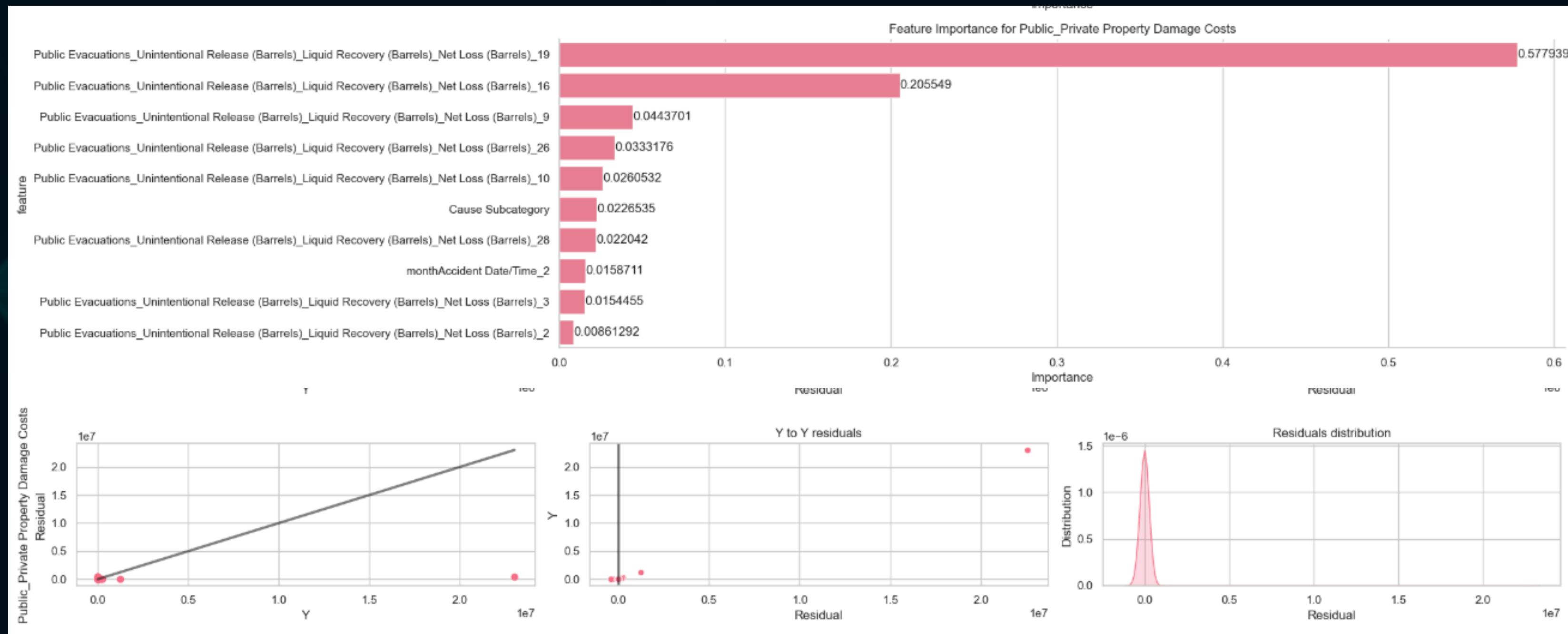


$R^2: 0.342$

$MSE: 4.84 \times 10^{10}$

$MAE: 1.16 \times 10^5$

PUBLIC/PRIVATE PROPERTY DAMAGE COSTS

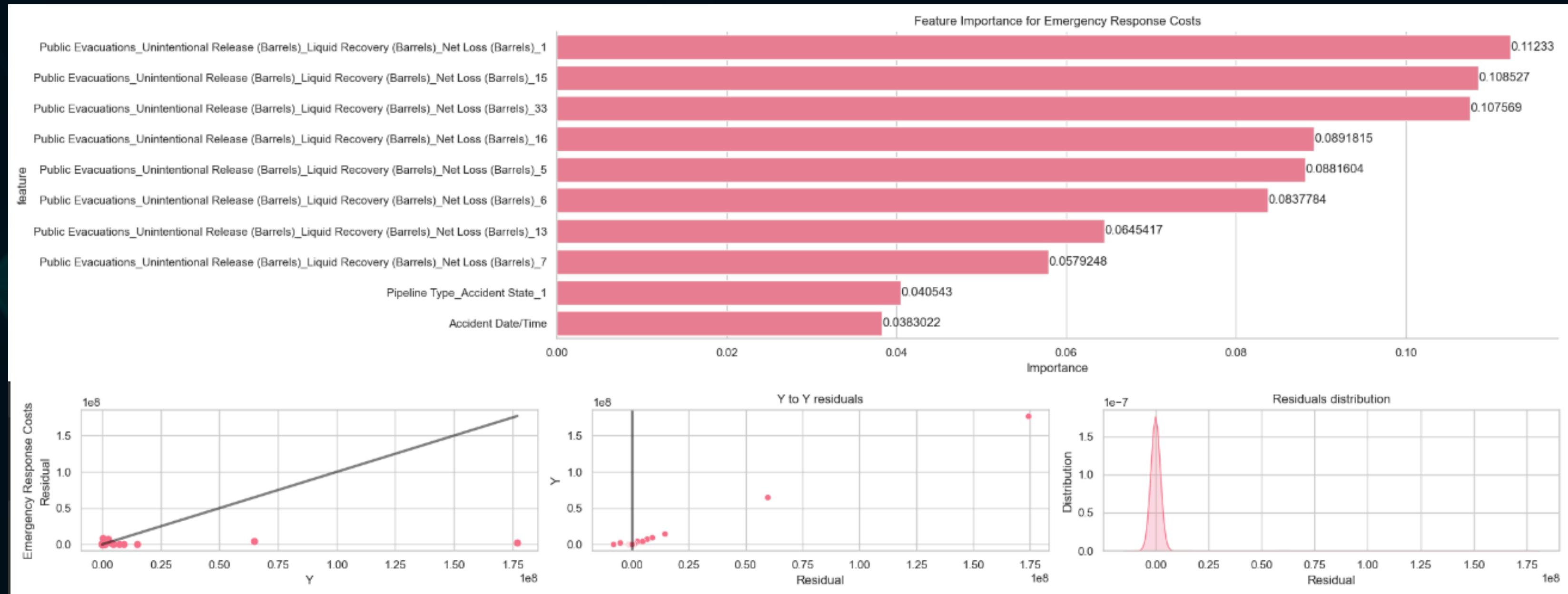


$R^2: 0.032$

$MSE: 9.21 \times 10^{11}$

$MAE: 6.20 \times 10^4$

EMERGENCY RESPONSE COSTS

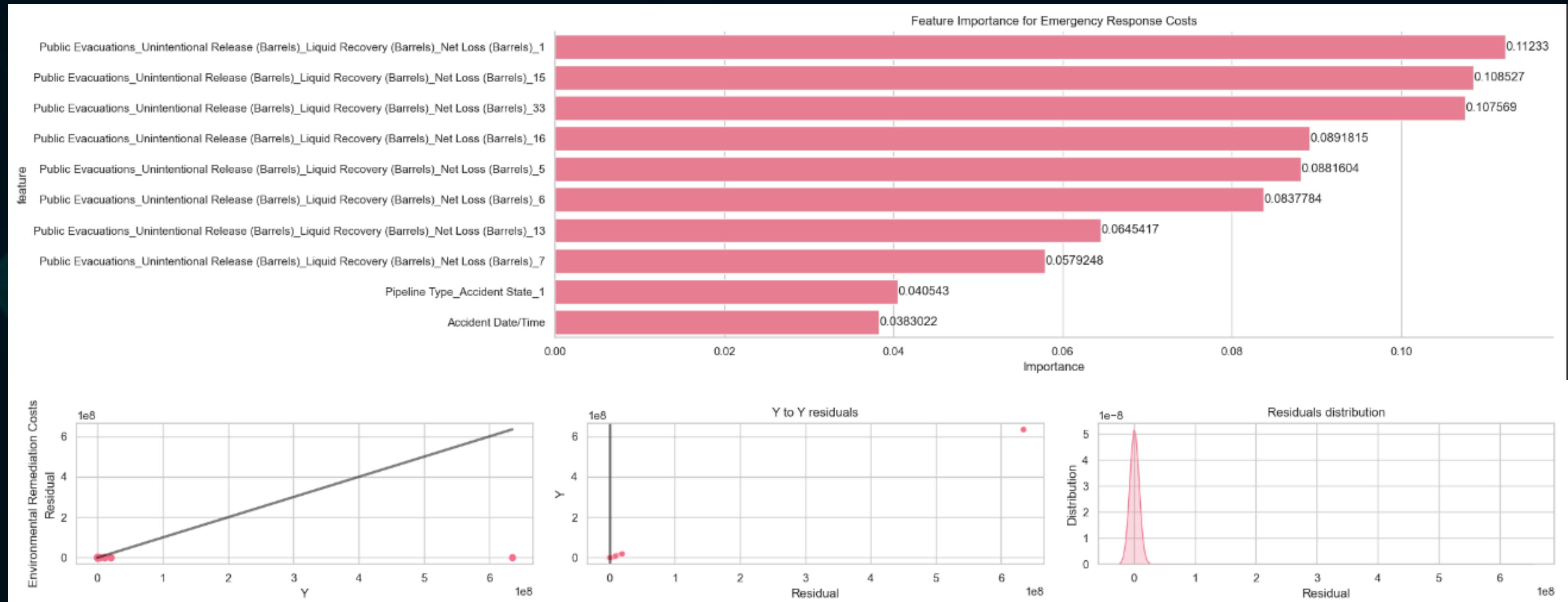


$R^2: 0.038$

$MSE: 6.18 \times 10^{13}$

$MAE: 6.67 \times 10^5$

ENVIRONMENTAL REMEDIATION COSTS

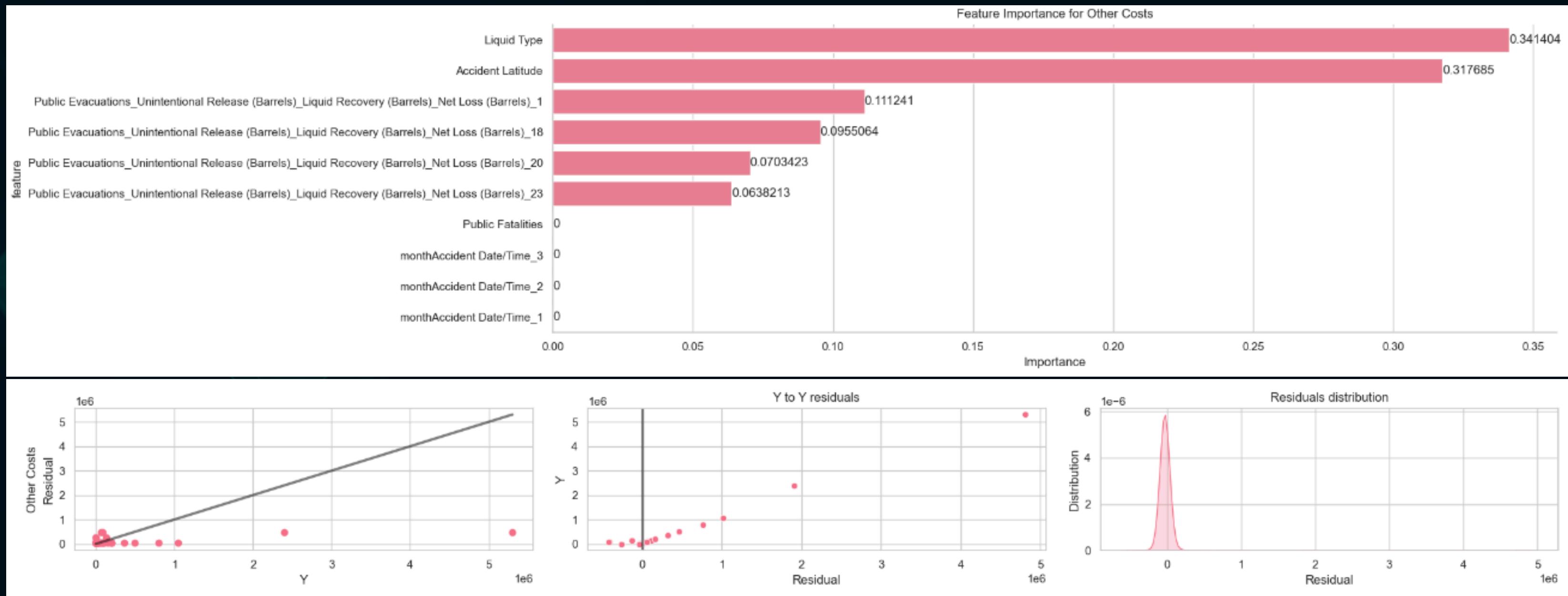


$R^2: 0.002$

$MSE: 7.23 \times 10^{14}$

$MAE: 1.34 \times 10^6$

OTHER COSTS

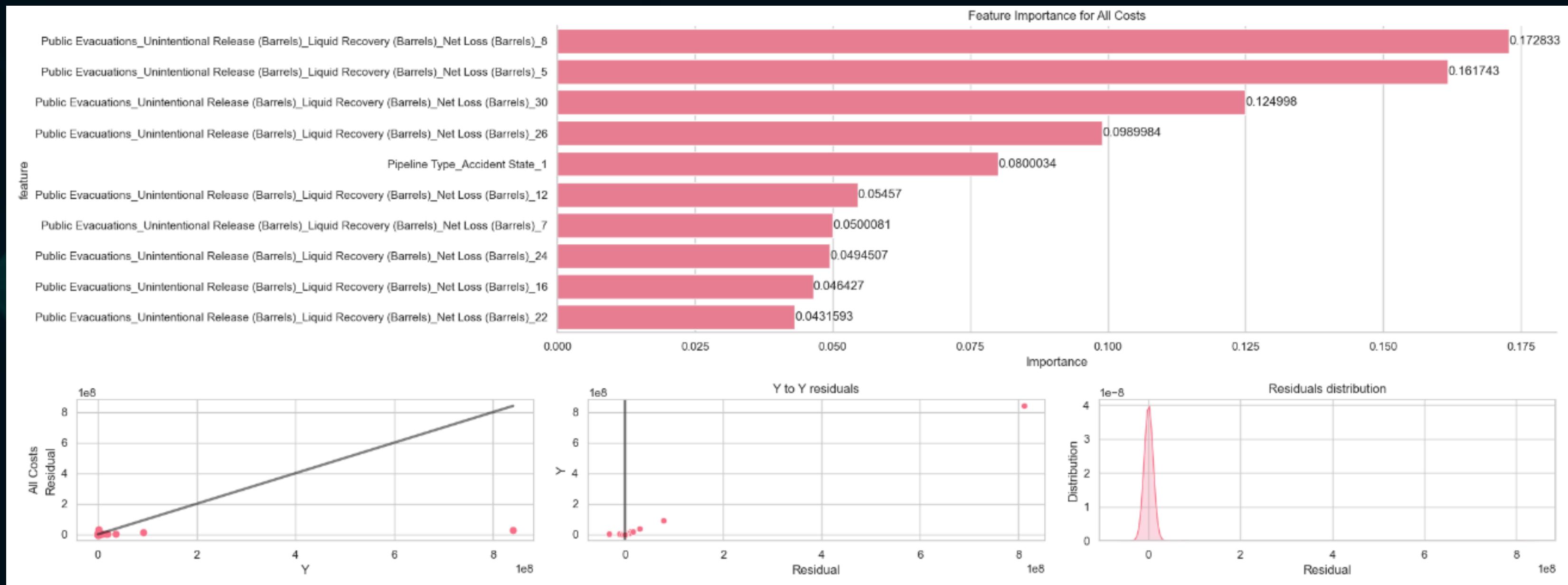


$R^2: 0.166$

$MSE: 5.38 \times 10^{10}$

$MAE: 5.46 \times 10^4$

ALL COSTS



$R^2: 0.066$

$MSE: 1.19 \times 10^{15}$

$MAE: 2.18 \times 10^6$

THANK YOU!

FOR YOUR ATTENTION