

Model Performance Metrics:

1- Multiple Linear Regression

Linear regression model applied to the dataset.

1-1-Evaluation Metrics:

MSE=10059525.716622159

MAE=1820.996476510666

R2=0.9108595703918305

1-2-Cross-Validation Results:

Mean cross validation test score: 0.17937827185733396

Mean cross validation train score: 0.907992832296129

Standard deviation in cv scores: 0.6116306629961016

2- Ridge regression:

The ridge regression model applied to the dataset.

2-1-Evaluation Metrics:

MSE=11305020.750448104

MAE=1906.9502517969547

R2=0.8998228708974764

2-2-Cross-Validation Results:

Mean cross validation test score: -0.8426050079251223

Mean cross validation train score: 0.9026701191934151

Standard deviation in cv scores: 2.443067784665445

3- Gradient Boosting Machines:

Gradient Boosting Regressor applied on the normalized data.

3-1-Evaluation Metrics:

MSE=4246017.906021098

MAE=1145.9909525630992

R2=0.9623747807870019

3-2-Cross-Validation Results:

Mean cross validation test score: 0.23718857442322427

Mean cross validation train score: 0.9442163359965724

Standard deviation in cv scores: 0.42081404036714865

4. Random Forest

4-1-Evaluation Metrics:

MSE=4448820.016800931

MAE=1145.9952714285714

R2=0.9605776913625486

4-2- Cross-Validation Results:

Mean cross validation test score: 0.23871470950619286

Mean cross validation train score: 0.9430095101005664

Standard deviation in cv scores: 0.5056310732557691

5. Support Vector Regression Model (SVR)

SVR model applied to the normalized data.

5-1-Evaluation Metrics:

MSE=125901292.44938798

MAE=4208.299780020598

R2=-0.11564855176201205

5-2- Cross-Validation Results:

Mean cross validation test score: -0.1972783418931082

Mean cross validation train score: -0.09710670835987818

Standard deviation in cv scores: 0.06234745815481293

5-4-Conclusion

Based on the above metrics results, Random Forest and Gradient Boosting are the best models among the ones compared. They have the highest R². However, the high values of MSE and MAE suggest that the actual predictions may still be quite far from the actual values. To make the process more accurate I want to use feature selection in the next step.

6-Features:

The Random Forest model and the Gradient Boosting model exhibit the same feature importance. Therefore features for both selected models include:

- GeoType_Citywide
- GeoType_Borough
- no2/Mean ppb
- pm/Mean mcg/m3
- Year

7- Hyperparameters:

7-1- Random Forest

max_depth: 10

max_features: sqrt

min_samples_leaf: 1

min_samples_split: 2

n_estimators: 200

-Best R2 Score: 0.9458677126090802

-Test Score: 0.9545611167355407

-Mean Squared Error on test set: 5127792.369216146

-R2 score on test set: 0.9545611167355407

7-2- Gradient Boosting

learning_rate: 0.1

max_depth: 5

min_samples_leaf: 1

min_samples_split: 2

n_estimators: 100

subsample: 0.8

-Best Score: 0.9527920253948207

-Test Score: 0.950376183226798

-Mean Squared Error on test set: 5600063.441260577

-R2 score on test set : 0.950376183226798