# Center for Building Infrastructure and Public Space (CBIPS) Water Sampling Data Science Team



### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

Mentor: Charles Shen Feniosky Peña-Mora Fredric M. Bell

## **Overview of Our Research Project**



### 01 Background

### 02 Machine Learning Based

- ✤ Water Sampling
- ✤ Water Monitoring

- How can ML help?
- Research & Case Study

# 03 Achievement 04 Summary

- Conversation with DEP
- Our Models

- Conclusions
- Future Work

# 01 Background:

Why Is Water Sampling & Monitoring So IMPORTANT?

# Background

countries.

The linkage between water and the economy is so compelling that decisions about water are rarely deferred. Decisions that are uninformed almost always have unintended consequences, with impacts on the environment, health, and society.

#### **Supply - Demand Issue**

Billion m3 Portion of gap Percent 8.000 Demand with no productivity improvements 7,000 Historical improvements 20% The water sector is severely under-funded, with in water productivity1 6,000 Remaining gap 60% 5,000 Increase in supply<sup>2</sup> under 20% business-as-usual Existing accessible. 3.000 reliable supply 2030 Today

#### Business-as-usual approaches will not meet demand for raw water

1 Based on historical agricultural yield growth rates from 1990-2004 from FAOSTAT, agricultural and industrial efficiency improvements from IFPR 2 Total increased capture of raw water through infrastructure buildout, excluding unsustainable extraction

3 Supply shown at 90% reliability and includes infrastructure investments scheduled and funded through 2010. Current 90%-reliable supply does not meet average demand SOURCE: 2030 Water Resources Group - Global Water Supply and Demand model; IFPRI: FAOSTAT

### **COLUMBIA** CBIPS Center for Buildings, Infrastructure and Public Space

Water challenge is tied to population growth,

especially serious financing gaps in developing

emerging developing countries.

(Source:https://www.mckinsey.com/~/media/mckinsey/dotcom/client service/sustainability 4 /pdfs/charting%20our%20water%20future/charting\_our\_water\_future\_full\_report\_.ashx)

# **Drinking Water Sampling**



- DEP continuously sample, and conduct analyses for numerous water quality parameters, including **microbiological, chemical, and physical measures**, throughout the watershed and as the water enters the distribution system.
- DEP also regularly tests water quality at nearly **1,000** water quality sampling **stations** throughout NYC.
- In 2019, DEP performed approximately 456,500 analyses on 36,300 samples from the distribution system, **meeting all state and federal monitoring requirements.** DEP also performed approximately 262,500 analyses on 15,000 samples from the reservoirs and their watersheds.

(Source:https://www1.nyc.gov/assets/dep/downloads/pdf/water/drinkingwater/drinking-water-supply-quality-report/2019-drinking-water-supplyquality-report.pdf)

# **Drinking Water Monitoring**



- DEP monitors the water in the **distribution system**, **upstate reservoirs and feeder streams**, **and wells** that are sources for New York City's drinking water supply.
- A growing network of **robotic monitoring** stations provided another 2 million water quality measurements last year, allowing DEP to optimize its operation of the reservoirs, support watershed protection efforts, and study water quality trends.

# Issues with Drinking Water @ NYC

# ASCE New York Infrastructure Grade (2015 the most recent)

Drinking Water:





### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

(Source:https://www.infrastructurereportcard.org/state-item/new-york/ Picture:https://www.smartcitiesdive.com/news/how-ai-and-data-turn-city-water-management-from-an-art-to-a-science/559424/)

# 02 Machine Learning Based

• How can Machine learning help?

Machine learning can automate, simplify and improve many aspects of water monitoring including:

Q Detect and correct equipment malfunctions

- Content of the second s
- Predict the effects of policy decisions
- Improve modeling and analysis
- Automate and control allocation and distribution

### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

(Source:Aquatic Informatics Workshop-Machine learning for water monitoring, hydrology and sustainability Kevin Swersky: http://www.cs.toronto.edu/~kswersky/wp-content/uploads/WorkshopPresentation.pdf)

9

# 02 Machine Learning Based

• Research & Case Study

# 1. Aquatic Informatics: software solutions provider





### One Water One Platform:

Water Data Management for the entire water cycle based on four software.

- AQUARIUS: Analytics Software.
- WaterTrax: Water Quality Compliance Software.
- Linko: Regulatory Compliance Software.
- Tokay: Backflow Prevention Software.



# 1. Aquatic Informatics: software solutions provider



Reasons why we choose:

- → Clear to see their Data Data Processing Pipeline.
- → Provide us knowledge that how ML algorithms applied in applications in an established software development company.

### COLUMBIA | CBIPS

Center for Buildings, Infrastructure and Public Space

# 2. ML In Predicting Water Quality





| 1  | A          |       | c                | D                | E                    | F                       | G             | н         | 1    | J                    | K           | L         |
|----|------------|-------|------------------|------------------|----------------------|-------------------------|---------------|-----------|------|----------------------|-------------|-----------|
| 1  | Location I | Years | Measurement Date | Measurement Time | <b>Casing Volume</b> | <b>Dissolved Oxygen</b> | Flow (in gpm) | Oxidation | рН   | Specific Conductance | Temperature | Turbidity |
| 2  | 03-8-10    | 2006  | Jun              | 11:00            |                      | 34.2                    |               | 175.1     | 5.78 | 0.286                | 13.64       | -4        |
| 3  | 03-8-10    | 2006  | Aug              | 10:15            |                      |                         |               | 270.6     | 6.29 | 366                  | 15.2        | 153.5     |
| 4  | 03-8-10    | 2006  | Dec              | 10:53            |                      | 2.7                     | f             | 455.7     | 5.98 | 280                  | 13.5        | 44.1      |
| 5  | 03-8-10    | 2007  | Mar              | 10:51            |                      | 0.63                    |               | 346       | 6.12 | 1493                 | 12.4        | 15.       |
| 6  | 03-8-10    | 2007  | Jul              | 11:05            |                      | 0.81                    |               | 302       | 6.16 | 577                  | 14.6        | 41./      |
| 7  | 03-8-10    | 2007  | Sep              | 8:15             |                      | 3.41                    |               | 253       | 5.72 | 424                  | 14.7        | 77.       |
| 8  | 03-8-10    | 2007  | Dec              | 15:10            |                      | 2.42                    | 0.184         | 390       | 5.75 | 2.14                 | 12.8        | 8.0       |
| 9  | 03-8-10    | 2008  | Mar              | 10:40            |                      | 1                       | 0.18          | 221       | 5.75 | 1623                 | 12.3        | 33.       |
| 10 | 03-8-10    | 2008  | Jun              | 10:35            |                      | 0.81                    | 0.16          | 215       | 6.11 | 1210                 | 14          | 70.8      |
| 11 | 03-8-10    | 2008  | Sep              | 11:55            |                      | 0.1                     | 0.125         | 14        | 6.2  | 726                  | 13.9        | 7         |
| 12 | 03-8-10    | 2008  | Dec              | 13:05            |                      | 5.15                    | 0.13          | 278.4     | 5.98 | 384                  | 12.85       | 65.4      |
| 13 | 03-8-10    | 2009  | Mar              | 12:30            |                      | 2.28                    | 0.12          | 275       | 5.99 | 623                  | 13.21       | 50.       |
| 14 | 03-8-10    | 2009  | Jun              | 16:00            |                      | 1.27                    | 0.13          | 196.8     | 5.96 | 492                  | 13.09       | 100       |
| 15 | 03-8-10    | 2009  | Sep              | 11:00            |                      | 3.68                    | 0.69          | 166.4     | 6.1  | 268                  | 13.78       |           |
| 16 | 03-8-13    | 2006  | Jun              | 9:50             |                      | 30.8                    |               | 151.1     | 5.77 | 260                  | 14.29       | 73.5      |
| 17 | 03-8-13    | 2006  | Aug              | 11:50            |                      | 1.2                     |               | 225.4     | 6.22 | 328                  | 15.8        | 12        |
| 18 | 03-8-13    | 2006  | Dec              | 12:37            | -                    | 1                       |               | 448.1     | 5.91 | 356                  | 13.3        | 27.       |
| 19 | 03-8-13    | 2007  | Mar              | 13:00            |                      | 1.09                    |               | 275       | 6.05 | 1428                 | 11.9        | 10.5      |
|    |            |       |                  |                  |                      |                         |               |           |      |                      |             |           |

# 3. Pattern Extraction In Water Quality Prediction

Scholar Article

→ Aims to implore data mining technique for pattern extraction and model prediction of water quality in water reservoir using different parameters and water quality index.



Jefferson Lerios - Technological Institute of the Philippines Mia Villarica - Laguna State Polytechnic University

International Journal of Mechanical Engineering and Robotics Research -Nov. 2019

### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

(Source: http://www.ijmerr.com/uploadfile/2019/0709/20190709114518923.pdf)

# 03 Achievement

- What we are aiming to achieve? By how?
  - Our Conversation with DEP
  - Model Explanation

# What we are aiming to do?



### **Data Source**

### GECCO 2019 Industrial Challenge: Monitoring of drinking-water quality

|    |  |                                   |    |    | -         | -    | -       | _         |          | -       |         |       |
|----|--|-----------------------------------|----|----|-----------|------|---------|-----------|----------|---------|---------|-------|
|    |  |                                   | 1  |    | Time      | Тр   | pН      | Cond      | Turb     | SAC     | PFM     | Event |
|    |  |                                   | 2  | 0  | 2017-07-0 | 6.94 | 8.60774 | 0.020953  | 0.125931 | 3.58683 | 43.7559 | FALSE |
|    |  |                                   | 3  | 1  | 2017-07-0 | 6.93 | 8.60589 | 0.0209654 | 0.127219 | 3.59025 | 43.4366 | FALSE |
|    |  |                                   | 4  | 2  | 2017-07-0 | 6.94 | 8.6022  | 0.020968  | 0.126482 | 3.58318 | 43.5994 | FALSE |
|    |  |                                   | 5  | 3  | 2017-07-0 | 6.94 | 8.6022  | 0.020971  | 0.126184 | 3.58769 | 43.3704 | FALSE |
|    |  | 32K testing data                  | 6  | 4  | 2017-07-0 | 6.94 | 8.60405 | 0.020973  | 0.127908 | 3.58287 | 43.1656 | FALSE |
|    |  | 100K training data<br>7 variables | 7  | 5  | 2017-07-0 | 6.94 | 8.60405 | 0.020980  | 0.126111 | 3.5878  | 43.4366 | FALSE |
|    | Thüringer                                      |                                   | 8  | 6  | 2017-07-0 | 6.95 | 8.60405 | 0.0209849 | 0.126327 | 3.58538 | 43.319  | FALSE |
|    | Fernwasserversorgung<br>Mehr als reines Wasser |                                   | 9  | 7  | 2017-07-0 | 6.95 | 8.6022  | 0.0209845 | 0.125598 | 3.59004 | 43.1656 | FALSE |
|    |  |                                   | 10 | 8  | 2017-07-0 | 6.95 | 8.6022  | 0.020998  | 0.126622 | 3.59189 | 42.5627 | FALSE |
| •• |  |                                   | 11 | 9  | 2017-07-0 | 6.94 | 8.6022  | 0.0209919 | 0.126646 | 3.58851 | 43.2824 | FALSE |
|    |  |                                   | 12 | 10 | 2017-07-0 | 6.95 | 8.60405 | 0.021001  | 0.12672  | 3.59556 | 43.5401 | FALSE |
|    |  |                                   | 13 | 11 | 2017-07-0 | 6.95 | 8.6022  | 0.021014  | 0.126704 | 3.59434 | 42.6052 | FALSE |
|    |  |                                   | 14 | 12 | 2017-07-0 | 6.95 | 8.60405 | 0.021016  | 0.1259   | 3.60614 | 43.3998 | FALSE |
|    |  |                                   | 15 | 13 | 2017-07-0 | 6.95 | 8.60589 | 0.021024  | 0.126167 | 3.59635 | 42.4146 | FALSE |
|    |  |                                   | 16 | 14 | 2017-07-0 | 6.96 | 8.60589 | 0.0210303 | 0.127636 | 3.585   | 42.2745 | FALSE |
|    |  |                                   | 17 | 15 | 2017-07-0 | 6.96 | 8.60405 | 0.0210264 | 0.124759 | 3.58992 | 43.8985 | FALSE |
|    |  |                                   | 18 | 16 | 2017-07-0 | 6.96 | 8.60589 | 0.021030  | 0.126349 | 3.58556 | 42.2745 | FALSE |
|    |  |                                   | 19 | 17 | 2017-07-0 | 6.96 | 8.60589 | 0.0210282 | 0.12642  | 3.57833 | 43.7185 | FALSE |

20

21

22

18

19

20

2017-07-0

2017-07-0

2017-07-0

### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

(Source: https://gecco-2019.sigevo.org/index.html/Competitions)

6.96

6.96

6.96

8.6022 0.0210268 0.126454

8.6022 0.0210388 0.126788

8.60589 0.021042 0.126777 3.59976 43.3337

3.59874

3.58915

43.4661

42.591

FALSE

FALSE

FALSE

# Water Quality Indicators

# **Our Conversation with DEP**



<u>Conductivity:</u> is important is because it can tell how much dissolved substances, chemicals, and minerals are present in the water. Higher amounts of these impurities will lead to a higher conductivity.

#### Predict Conductivity

COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

(Source: New York City 2018 Drinking Water Supply and Quality Report, 18 GECCO 2019 Industrial Challenge)

# 03 Achievement

Models Explanation

 Ridge Regression

# Methodology

### **Ridge Regression:**

The objective function (also called the cost) to be minimized is the RSS (Residual Sum of Squares) plus the sum of square of the magnitude of weights. This can be depicted mathematically as:

 $Cost(W) = RSS(W) + \lambda * (sum of squares of weights)$ 

$$= \sum_{i=1}^{N} \left\{ y_i - \sum_{j=0}^{M} w_j x_{ij} \right\}^2 + \lambda \sum_{j=0}^{M} w_j^2$$

is the sum of the squares of residuals (deviations predicted from actual empirical values of data). It is a measure of the discrepancy between the data and an estimation model.

| The   | teri  | m   | wit | h l   | amb  | oda  | is | of  | ten | С   | alled |
|-------|-------|-----|-----|-------|------|------|----|-----|-----|-----|-------|
| 'Pena | alty' | sin | ce  | it ir | crea | ases | RS | SS. | We  | ite | erate |
| certa | in    | val | ues | 0     | nto  | th   | е  | lar | nbd | a   | and   |
| evalı | ıate  | the | mo  | odel  |      |      |    |     |     |     |       |

### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

(Source:https://towardsdatascience.com/ridge-regression-for-better-usage-2f19b3a202db)

# **Python Coding - Ridge Regression**



#### 28 Jf1 = pd.read\_excel("C:/Users/user/Desktop/ML/watersample-test.xlsx").drop(["E 29 df1 = df1.dropna() 30 df1 = df1[["Tp","pH","Turb","SAC","PFM","Cond"]] 31 test\_x = df1.loc[:,"Tp":"PFM"] 32 test\_x = scaler.fit\_transform(test\_x) 33 test\_y = df1.loc[:,"Cond"] 34

#### 36 def model(size):

- X = data.loc[:size,"Tp":"PFM"]
- Y = data.loc[:size,"Cond"]
- X\_train,X\_test,y\_train,y\_test = train\_test\_split(X,Y, test\_size = 0.2)
  #Use Ridge regression
- #Using cross validation to decide the best alpha for ridge
- alpha\_ridge = [0.1, 1, 10, 100, 1000,2000,4000,8000, 10000, 15000]
- model\_ridge = RidgeCV(alphas = alpha\_ridge).fit(X\_train,y\_train)

#### #Compute test error

- pred\_ridge = model\_ridge.predict(test\_x)
- test\_error = mean\_absolute\_error(test\_y,pred\_ridge)
- #Compute validation error
- val\_ridge = model\_ridge.predict(X\_test)
- val\_error = mean\_absolute\_error(y\_test,val\_ridge)
- #Compute train error
- train\_ridge = model\_ridge.predict(X\_train)
- train\_error = mean\_absolute\_error(y\_train,train\_ridge)
- return test\_error,val\_error,train\_error

58 #Computer the val\_error and train\_error with the interval of 500

59 gap = 200 60 index = 0 61 test\_error = [] 62 val\_error = [] 63 train error = []

### **Results**

The mean absolute error for the entire data set is(test\_error) 0.0003074708634444559 The mean absolute error for the entire data set is(train\_error) 0.00026323328485385234 The mean absolute error for the entire data set is (val\_error) 0.00026521310090718396 The prediction are: [0.02135556 0.02134536 0.02131694 ... 0.02086572 0.02083075 0.02085077]



In [7]: print("The prediction are:",result[:20])
The prediction are: [0.02135556 0.02134536 0.02131694 0.02132844 0.02132125 0.02133675
0.02130564 0.02131144 0.02132444 0.0213268 0.0213325 0.02131131
0.02130093 0.0212625 0.02126737 0.02128587 0.02130741 0.02130025
0.02129438 0.02128749]

# 03 Achievement

Models Explanation

 Ridge Regression
 Boosting Algorithm

# Methodology

### **Boosting algorithm**

Boosting algorithm combines the outputs from weak learner and creates a strong learner which eventually improves the prediction power of the model.



### COLUMBIA | CBIPS

Center for Buildings, Infrastructure and Public Space

# **Python Coding -Boosting Algorithm**

# Coding Process —

### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space



25

### Results

| Collecti             | ng xgb        | oost                    |             |                   |               | 100           |                |      |
|----------------------|---------------|-------------------------|-------------|-------------------|---------------|---------------|----------------|------|
| Downlo               | ading         | xgboost-1.0.2-          | py3-none-wi | n_amdt            | 4. whi (24. 6 | MB)           | 24.6 MR 3.3    | MR/e |
| Requirem             | ent al        | ready satisfie          | d: numpy in | ıe:∖ar            | aconda\lib\   | site-packages | (from xgboost) | (1.) |
| Requirem             | ent al        | ready satisfie          | d: scipy in | e:\an             | aconda\lib\   | site-packages | (from xgboost) | (1.4 |
| Installi             | ig col        | lected package          | s: xgboost  |                   |               |               |                |      |
| Successi             | illy 1        | nstalled xgboo          | st-1.0.2    |                   |               |               |                |      |
| C:\Users             | liang         | yi\Desktop\Dat          | a ML Part 3 | pythe             | n "Haonan' s  | XGBOOST Model | . py"          |      |
| The vali             | lation        | mean absolute           | error is    |                   |               |               |                |      |
| 0.000130             | 504944        | 8457718                 | error in    |                   |               |               |                |      |
| 0.000130             | 280421        | 17721556                | 1101 13     |                   |               |               |                |      |
| The test             | mean          | absolute error          | is          |                   |               |               |                |      |
| 0.000475             | 084912        | 39910117                | _           |                   |               |               |                |      |
| Ine pred<br>[0_02091 | 169 0         | s are<br>02091169 0 020 | 91169 0 021 | 05999             | 0 02088082    | 0 02085814    |                |      |
| 0.02085              | 314 0.        | 02085814 0.020          | 85814 0.020 | 85814             | 0.02085814    | 0. 02091944   |                |      |
| 0.02085              | 314 0.        | 02105999 0.021          | 05999 0.021 | 05999             | 0.02105999    | 0.02105999    |                |      |
| 0.02105              | 999 O.        | 02105999 0.021          | 05999 0.021 | 05999             | 0.02115735    | 0.02105999    |                |      |
| 0.02115              | 735 0.        | 02115735 0.021          | 15735 0.021 | 15735             | 0.02115735    | 0.02115735    |                |      |
| 0.02115              | 735 0.        | 02115735 0.020          | 97929 0.020 | 97929             | 0.02097929    | 0.02097929    |                |      |
| 0.02097              | 929 0.        | 02097929 0.020          | 97929 0.020 | 97929             | 0.02097929    | 0.02102947    |                |      |
| 0.02097              | <u>929 O.</u> | 02097929 0.020          | 97929 0.020 | 97929             | 0.02097929    | 0.02105999    |                |      |
| The val              | idat          | ion mean so             | uared err   | or is             |               |               |                |      |
| 0 07100              | 1220          | 71005509                | aurea err   | 01 13             | '             |               |                |      |
| 2.0/122              | 1329          | 1199006-00              |             |                   |               |               |                |      |
| Ine tra              | .1n1n         | g mean squar            | red error   | 15                |               |               |                |      |
| 2.74619              | 8202          | 8902685e-08             |             |                   |               |               |                |      |
| The tes              | t me          | an squared (            | error is    |                   |               |               |                |      |
| 2.89797              | 9553          | 747056e-07              |             |                   |               |               |                |      |
| The pre              | dict          | ione are                |             |                   |               |               |                |      |
|                      | Q105          | 0 02007176              | 0 020991    | 05 0              | 0210770       | 0 0208624     | 0 00004000     |      |
| 0. 0200              | 4000          | 0.02001110              | 0.020001    | <del>3</del> 3 0. | 0210119       | 0.0200024     | 0.02004200     |      |
| 0.0208               | 4288          | 0.02084288              | 0.020842    | 88 0.             | 02084288      | 0.02084288    | 0. 02088597    |      |
| 0.0208               | 0864          | 0.0210779               | 0.021077    | 90.               | 0210779       | 0.0210779     | 0.0210779      |      |
| 0.0210               | 779           | 0.0210779               | 0.021077    | 9 0.              | 0210779       | 0.0210779     | 0.0210779      |      |



# 03 Achievement

- Models Explanation
  - Ridge Regression
  - Decision Tree
  - Random Forest

# Methodology

### **Random Forest Algorithm**

**Decision Tree algorithm** belongs to the family of supervised learning algorithms. the decision tree algorithm can be used for solving **regression and classification problems**.

Random Forest consists of a large number of individual decision trees. Each individual tree in the random forest spits out a class prediction and the class with the most votes becomes its model's prediction.

### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space



(Source:https://towardsdatascience.com/understandingrandom-forest-58381eo6o2d2)

# Python Coding - Random Forest Regressor

# Coding Process —

### COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

import pandas as pd
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean\_squared\_error
from sklearn.model\_selection import train\_test\_split
#from sklearn.utils import shuffle
fp\_train = r'C:\Users\liangyi\Desktop\Data ML Part 3\watersample-train.csv'
fp test = r'C:\Users\liangyi\Desktop\Data ML Part 3\watersample-test.csv'

def Forest\_model(): df\_train = pd.read\_csv(fp\_train) df\_test = pd.read\_csv(fp\_test) #watersample\_data = shuffle(watersample\_data) y = df\_train.Cond[:10000] y\_test = df\_test.Cond[:10000] #print(y) watersample\_features = ['pH', 'Tp', 'Turb', 'SAC', 'PFM'] X = df\_train[watersample\_features][:10000] X\_test = df\_test[watersample\_features][:10000] #print(X) #watersample\_data = watersample\_data.dropna() train\_X, val\_X, train\_Y, val\_Y = train\_test\_split(X, y, test\_size = 0.2) model = RandomForestRepressor(max\_depth=5)

#watersample\_model.fit(X, y)

model.fit(train\_X, train\_y)
val\_predictions = model.predict(val\_X)
test\_predictions = model.predict(X\_test)
predicted\_watersample\_Cond = model.predict(X)

print("The validation mean squared error is")
print(mean\_squared\_error(val\_y, val\_predictions))

print("The training mean squared error is")
print(mean\_squared\_error(y, predicted\_watersample\_Cond))

print("The test mean squared error is")
print(mean\_squared\_error(y\_test, test\_predictions))
#print("Making\_predictions\_for\_the\_following\_5\_watersamples:")

Results

| The validation mean absolute error is  |
|--|
| 0. 00017848401165124973  |
| The training mean absolute error is  |
| 0. 0001804644757287498   |
| The test mean absolute error is  |
| 0. 00042701659617813703  |
| The predictions are  |
| [0. 02106792 0. 02106792 0. 02106792 0. 02106792 0. 02106792 0. 02106792 0. 02106792 |
| 0.02106792 0.02106792 0.02106792 0.02106792 0.02106792 0.02106792                    |
| 0.02106792 0.02106792 0.02106792 0.02106792 0.02106792 0.02106792                    |
| 0.02106792 0.02106792 0.02106792 0.02106792 0.02106792 0.02106792                    |

The validation mean squared error is 4.436288942087643e-08 The training mean squared error is 4.618563612324684e-08 The test mean squared error is 2.4465875705867854e-07 The predictions are [0.02106778 0.02106778

**Mean Absolute Error** is the average over the test sample of the absolute differences between prediction and actual observation.

|                  | <b>Ridge Regression</b> | Boosting    | Random Forest |
|------------------|-------------------------|-------------|---------------|
| MAE (Training)   | 0.000263233             | 0.000130280 | 0.000180485   |
| MAE (Validation) | 0.000265221             | 0.000130604 | 0.000182120   |
| MAE (Testing)    | 0.000307470             | 0.000475084 | 0.000425865   |

MAE (Training) : Ridge Regression > Random Forest > Boosting.MAE (Validation): Ridge Regression > Random Forest > Boosting.MAE (Testing): Boosting> Random Forest> Ridge Regression.

**Mean Squared Error** is the average squared difference between the estimated values and the actual value.

|                  | <b>Ridge Regression</b> | Boosting    | Random Forest |
|------------------|-------------------------|-------------|---------------|
| MSE (Training)   | 1.0580*e^-7             | 2.7461*e^-8 | 4.6185*e^-8   |
| MSE (Validation) | 9.9343*e^-8             | 2.8712*e^-8 | 4.4362*e^-8   |
| MSE (Testing)    | 1.5453*e^-7             | 2.8979*e^-7 | 2.4465*e^-7   |

MSE (Training) : Ridge Regression > Random Forest > Boosting.MSE (Validation): Ridge Regression > Random Forest > Boosting.MSE (Testing): Boosting > Random Forest> Ridge Regression.

COLUMBIA | CBIPS Center for Buildings, Infrastructure and Public Space

# 04 Summary

- Conclusion
- Improvement
- Current Trend
- Future Opportunity

# **Conclusion & Improvement**



- After comparing with other data sets eg: NYC open data, we choose the industrial challenge data because the overall number of the variables and the completeness of the dataset was appropriate.
- Given three models, in term of the testing error, ridge regression performs better than other two models.



- → More feature analysis and cross validation will be needed to further strengthen our conclusions and complete comparison among the different algorithms.
- → On the one hand, we shall see whether the model can generalize to other data. On the other hand, we could apply other algorithms.
- → Classification problems, which differs from regressions, could be our next focus.



- Monitor water quality data.
- Anticipate water supply flow.
- Optimize management of sewage, treatment plants, and desalination plants.



- → The environmental time series data can contain outliers, anomalies and gaps. Our models need to be flexible, fast and general.
- → Be able to develop good models for anomaly and fault detection for advanced study.
- → Apply ML regression technique to predict even harder-to-measure metrics which would require appropriate datasets being available.
- → Apply Deep Learning and more advanced AI technologies.

### COLUMBIA | CBIPS

Center for Buildings, Infrastructure and Public Space

### Resources

- Ridge Regression for Better Usage <u>https://towardsdatascience.com/ridge-regression-for-better-usage-2f19b3a202db</u> <u>https://www.analyticsvidhya.com/blog/2016/01/ridge-lasso-regression-python-complete-tutorial/</u>
- Python Model Tuning Methods Using Cross Validation and Grid Search <u>https://towardsdatascience.com/ridge-regression-for-better-usage-2f19b3a202db</u> <u>http://karlrosaen.com/ml/learning-log/2016-06-20/</u> <u>https://towardsdatascience.com/machine-learning-algorithms-part-14-cross-validation-and-ridge-regression-example-in-python-2d9e0c4de875</u>
- Statistical Learning Theory https://en.wikipedia.org/wiki/Statistical learning theory
- NY Codes for Water Quality <u>https://govt.westlaw.com/nycrr/Document/I4ed90412cd1711dda432a117e6e0f345?viewType=FullText&originationContext=documenttoc&transi tionType=CategoryPageItem&contextData=(sc.Default)</u>
- Water Quality Standards http://mrccc.org.au/wp-content/uploads/2013/10/Water-Quality-Salinity-Standards.pdf
- How machine learning can help?
  - http://www.cs.toronto.edu/~kswersky/wp-content/uploads/WorkshopPresentation.pdf
- Data Source:

https://www.th-koeln.de/mam/downloads/deutsch/hochschule/fakultaeten/informatik\_und\_ingenieurwissenschaften/rulesgeccoic2019.pdf

• Drinking Water:

https://www.smartcitiesdive.com/news/how-ai-and-data-turn-city-water-management-from-an-art-to-a-science/559424/ https://www.infrastructurereportcard.org/state-item/new-vork/

# COLUMBIA | CBIPS

#### Center for Buildings, Infrastructure and Public Space

# **Thank You!**

### Grateful if any comment comes up!

You can find us at: hy2636@columbia.edu yl4322@columbia.edu sf2970@columbia.edu