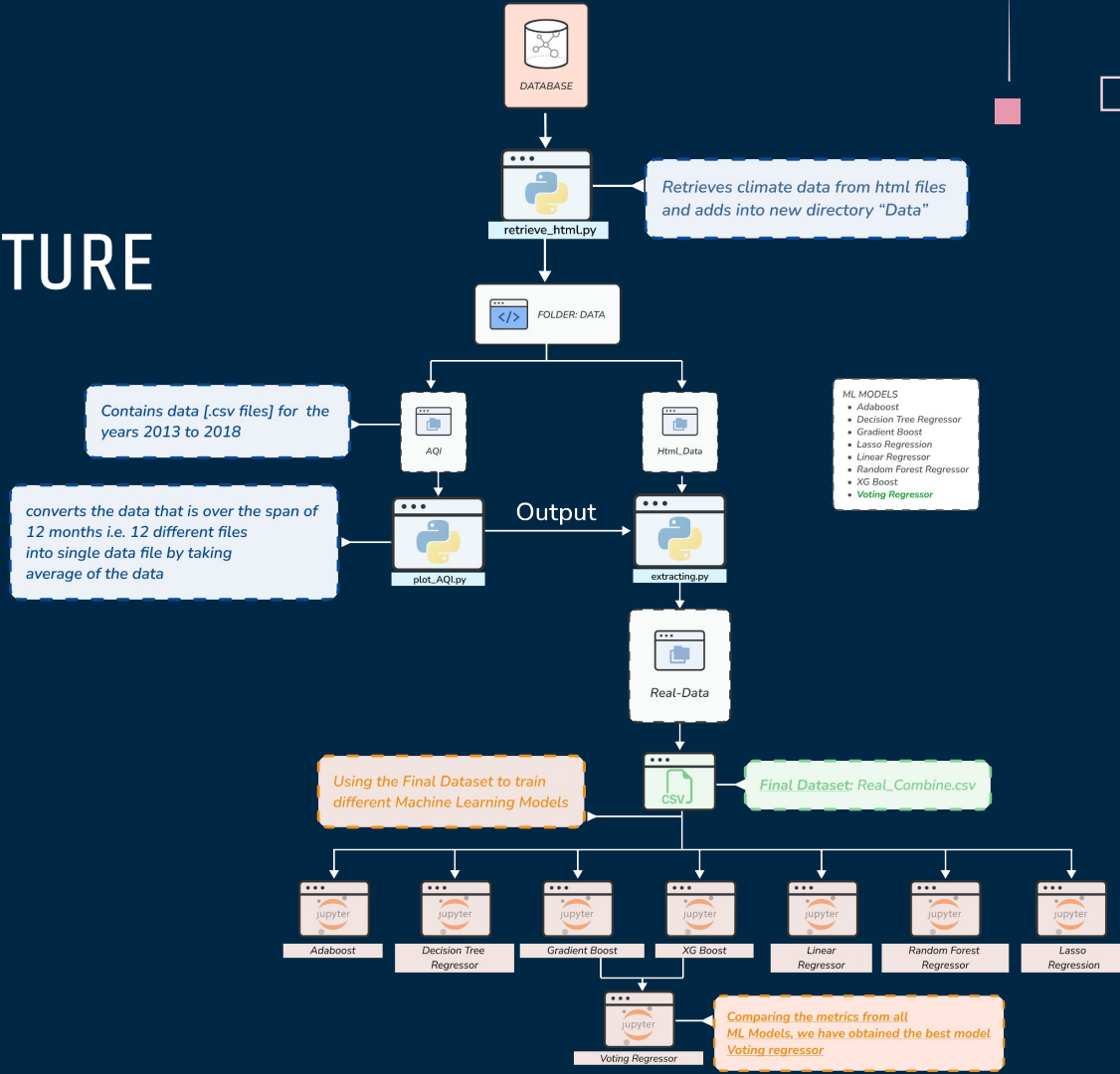


# SYSTEM ARCHITECTURE

02

# SYSTEM ARCHITECTURE



# PROPOSED MODEL

03

# PROPOSED MODEL

## DATA COLLECTION

Using Python libraries to parse and extract data from HTML files



## DATA PREPROCESSING

Cleaning, transforming, and organizing the data to make it suitable for further analysis or modeling.

## FEATURE SELECTION

Remove features that are noisy or irrelevant by dropping the columns and get final dataset



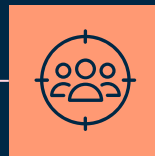
## MODEL BUILDING

Build individual ML models and use the final dataset for training the model

## MODEL

### EVALUATION

Using bagging and boosting techniques and evaluating performance metrics



## ANALYSIS & IDENTIFICATION

Comparing the performance metrics, identifying and finalising the best ML Model

# OUR PROCESS

Collect Climate data from the internet using **retrieve\_html.py**

DATABASE

01

02

DATA PREPROCESSING

Use **plot\_AQI.py** to combine 12 data files and use its output in **extracting.py** and create final dataset: **Real\_Combine.csv**

Use Bagging and Boosting Techniques in different ML Model and evaluate the performance metrics

BUILD MODEL

03

04

FINAL MODEL

Obtain final model by identifying and comparing best **rmse** values

# DATA COLLECTION

- Creates all data files, each year has 12 data files representing each month's data

## retrieve\_html.py

```
1 import os
2 import time
3 import requests
4 import sys
5
6
7 def retrieve_html():
8     for year in range(2013,2019):
9         for month in range(1,13):
10             if(month<10):
11                 url='http://en.tutiempo.net/climate/0{}-{}-ws-421820.html'.format(month
12                                     ,year)
13             else:
14                 url='http://en.tutiempo.net/climate/{}-{}-ws-421820.html'.format(month
15                                     ,year)
16             texts=requests.get(url)
17             text_utf=texts.text.encode('utf=8')
18
19             if not os.path.exists("Data/Html_Data/{}".format(year)):
20                 os.makedirs("Data/Html_Data/{}".format(year))
21             with open("Data/Html_Data/{}/{}.html".format(year,month),"wb") as output:
22                 output.write(text_utf)
23
24             sys.stdout.flush()
25
26 if __name__=="__main__":
27     start_time=time.time()
28     retrieve_html()
29     stop_time=time.time()
30     print("Time taken {}".format(stop_time-start_time))
```

# DATA COLLECTION

- plot\_AQI.py does the averaging and converts each year(12 data files) into one single file.
- Each year is averaged to one file and finally we get 4 files from year 2013 to 2016.

plot\_AQI.py

```
1
2 import ...
3
4
5
6 def avg_data_2013():
7     temp_i=0
8     average=[]
9     for rows in pd.read_csv('Data/AQI/aqi2013.csv', chunksize=24):
10         add_var=0
11         avg=0.0
12         data=[]
13         df=pd.DataFrame(data=rows)
14         for index,row in df.iterrows():
15             data.append(row['PM2.5'])
16             for i in data:
17                 if type(i) is float or type(i) is int:
18                     add_var=add_var+i
19                 elif type(i) is str:
20                     if i!='NoData' and i!='PwrFail' and i!='---' and i!='InVld':
21                         temp=float(i)
22                         add_var=add_var+temp
23             avg=add_var/24
24             temp_i=temp_i+1
25
26         average.append(avg)
27     return average
28
29 def avg_data_2014():
30     temp_i=0
31     average=[]
```

# DATA COLLECTION

extracting.py

- extracting.py combines all the 4 years data into one single file which is going to be used in various models.
- After combining 4 files, we get one single file, i.e, "Real\_combine.csv".
- Real\_Combine.csv file is then used to feed the data into various models like: Linear Regression, Decision Tree regressor, Lasso regression, XGBoost, Gradient Boost and AdaBoost

```
1
2 import ...
9
10 def met_data(month, year):
11     file_html = open('Data/html_Data/{}/{}.html'.format(year, month), 'rb')
12     plain_text = file_html.read()
13
14     tempD = []
15     finalD = []
16
17     soup = BeautifulSoup(plain_text, "lxml")
18     for table in soup.findAll('table', {'class': 'medias mensuales numspan'}):
19         for tbody in table:
20             for tr in tbody:
21                 a = tr.get_text()
22                 tempD.append(a)
23     rows = len(tempD)/15
24
25     for times in range(round(rows)):
26         newtempD = []
27         for i in range(15):
28             newtempD.append(tempD[0])
29             tempD.pop(0)
30         finalD.append(newtempD)
31
32     length = len(finalD)
33
34     finalD.pop(length-1)
35     finalD.pop(0)
36
```



THANK YOU!!