

AUTONOMOUS OPEN DATA PREDICTION FRAMEWORK

Documentation

Riga Technical University

Faculty of Computer Science and Information Technology

Information Technology Institute

Department of Management Information Technology

Jānis Pekša

Research Assistant

<http://iti.rtu.lv/vitk/lv/katedra/darbinieki/janis-peksa>

INDEX

1. Introduction	5
2. Website Connection.....	6
Definitions:.....	6
Definition of get_table()	6
Definition of get_data_lists_from_table()	7
3. MySQL Connection.....	8
Definitions	8
Definition of MySQLClient __init__() also known as class constructor	8
Definition of MySQLClient connect_to_database().....	8
Definition of get_all_info_from_database()	9
Definition of get_info_by_station()	9
Definition of get_info_by_stations()	9
4. MongoDB Connection	10
I. Difference between database clients.....	10
Definitions	10
Definition of MongoClient __init__() also known as class constructor	10
Definition of get_connection().....	10
II. MongoDB Importing	10
Definitions	11
Definition of get_all_info_from_main_database()	11
Definition of get_data_from_collection().....	11
Definition of get_data_from_collections()	12
III. MongoDB Exporting.....	12
5. Data Handling.....	13
I. For exporting to MongoDB	13
Definitions	13
Definition of get_data_dicts().....	13
Definition of get_date()	13
II. Process data from MongoDB.....	14
Data from MongoDB database	14
Definitions	14
Definition of get_prepared_lists_for_estimation()	14
Definition of get_lists_of_chosen_values_with_station_names().....	14
Definition of get_lists_of_measurements_with_station_names()	15

Definition of is_possible_to_fill_missing_data()	15
Definition of get_indexes_for_filling()	16
Definition of fill_missing_data_in_list()	16
Definition of get_list_of_float_numbers_and_station()	16
Definition of get_station_names()	17
Definition of get_lists_of_measurements_without_station_names()	17
III. Process data from MySQL	17
Definitions	18
Definition of get_index_of_value()	18
Definition of get_exact_value_from_my_sql_records().....	19
Definition of get_exact_value_from_many_my_sql_records().....	19
Definition of replace_none_with_dash()	19
Definition of fill_missing_data_in_lists().....	19
Definition of get_station_codes()	20
Definition of zip_codes_and_measurements()	20
IV. Accuracy	20
Definitions	20
Definition of get_lists_of_ints().....	20
Definition of get_accuracies()	21
Definition of get_accuracy().....	21
6. Estimation.....	22
Definitions	22
Definition of KalmanFilter __init__() also known as class constructor	22
Definition of get_estimates()	22
Definition of make_basic_calculations()	22
Definition of __calculate_kalman_gain()	23
Definition of __calculate_estimate()	23
Definition of __calculate_error_in_estimate().....	23
Definition of get_lists_of_estimates()	23
Definition of __get_initial_estimate()	23
7. Plots	24
Definitions	24
Definition of GraphEditor __init__ also known as class constructor	24
Definition of create_est_and_meas_plot().....	24
Definition of create_est_plot().....	25

Definition of <code>create_meas_plot()</code>	25
Definition of <code>show_plot()</code>	25
Definition of <code>save_plot()</code>	25
8. API	26
I. <code>/get/estimates/all</code>	26
II. <code>/get/estimates</code>	26
For example	26
III. <code>/get/accuracies</code>	27
For example	27

1. INTRODUCTION

This is the detailed documentation of the program, which is part of the Autonomous Open Data Prediction Framework. Here are described methods to implement all essential functionality. Some of the methods aren't described because they are built-in in Python or taken from Python modules. It will be announced if some of Python files described below use other modules for Python. All these functions, methods, and modules written not by me have their references, description, and examples of usage on the Internet. All of them can be used freely, as provided and without almost any limitations.

2. WEBSITE CONNECTION

The first source of data is the website www.lvceli.lv/cms. The **Request.py** to implement the needed functionality. At the core of it are such modules as **requests** and **BeautifulSoup**. First, one is needed to make GET/POST requests to the website and the second one is needed to process HTML properly which is returned by requests.o

Examples of using and explanation:

```
table = Request.get_table()
```

Firstly, to get a table from the website. This method does everything it needs to return the correct table with observations that it contains. No need to pass the URL of the website because it's given as a default parameter, as shown in the definition below. Now to have a table with data, but before need to separate values from HTML because this method only returns part of the webpage.

```
data_lists = Request.get_data_lists_from_table(table)
```

Now using another method that will extract all data from HTML, save it to lists and return a list of records like `[[], [], []]`.

This data program can process correctly, export to the database and make estimations but before this operations, some more functionality will be described to understand what is going on inside fully.

Because working with databases to store vast amounts of data, this module is crucial for gathering data in real-time and has some simple methods to work with. Still, it is much more convenient to import/export data using clients of databases that are described next.

Definitions:

Definition of `get_table()`

```
def get_table(url='http://www.lvceli.lv/cms/'):
    with requests.session() as s:
        r = s.post(url, data=login_data, headers=headers)

    soup = BeautifulSoup(r.content, 'html5lib')
    return soup.find("table", attrs={"class": "norm", "id": "table-1"})
```

In this method, making a POST request to a website with data that contains login and password to get access to the table with observations. Because every successful application returns a whole webpage to separate our table from the whole piece of HTML to get by using unique attributes of the table that are needed. In details:

soup – an instance of **class BeautifulSoup**, where HTML will be stored.

Using method `find()`, and in arguments, passing the tag that looks for, and it's unique attributes to be sure that a correct table in the result is taken.

Definition of get_data_lists_from_table()

```
def get_data_lists_from_table(table) -> list:
    tmp_list = []
    data_list = []

    for row in table.find_all('tr'):
        for cell in row.find_all('td'):
            tmp_list.append(cell.text)

        while len(tmp_list) < 19:
            tmp_list.append('-')

        data_list.append(tmp_list)
        tmp_list = []

    data_list.pop(0)

    return data_list
```

In this method, iterating through every row in the table and through every cell of row (loop in loop, internal circuit) and taking text from HTML tags and adding it to the temporal list, which appends to the main index of data. At the end of the method, the first list of data is being deleted because it contains only titles of columns that is not needed in our main list.

3. MYSQL CONNECTION

To be able to get data from the MySQL database, **MySQLClient.py**, which has few methods to download data from a table in the database, and it uses **mysql.connector** module for Python to get a possibility to connect to the database.

A simple instruction to start:

1. Firstly, it is needed to create a client:

```
sql_client = MySQLClient('xx.xx.xx.xx', 'xxx', 'xxx', 'xxx')
```

In parameters, we pass the IP address of database, login, and password to get access and name of the table to get data from. In this case, it is a database with observations for one month.

2. Now it is connected to the database and have two ways: get all of the data stored in the database or data from exact station/s. It can be either use *get_all_info_from_database()* or *get_info_by_station()/get_info_by_stations()*. The difference is in the amount of returned data and which data is needed.

Example:

```
records = sql_client.get_info_by_station('LV01')
```

Now in variable *records* is stored all of observations from station with code “LV01”.

```
records = sql_client.get_all_info_from_database()
```

But using method to get all of stored observations, it will get all of them and store in variable *records*.

```
records = sql_client.get_info_by_stations(['LV01', 'LV02', 'LV03'])
```

By using this method it can pass list of station codes and return lists of data for every station that is asked.

Definitions

Definition of MySQLClient_init() also known as class constructor

```
def __init__(self, host, username, password, database):
    self.host = host
    self.username = username
    self.password = password
    self.database = database
    self.database_connection = self.connect_to_database()
    self.cursor = self.database_connection.cursor()
```

Here is just create an instance of the class with all the needed variables. *self.cursor* is a variable that is responsible for interaction with the database: making queries, executing commands, and returning results.

Definition of MySQLClient connect_to_database()

```
def connect_to_database(self):
    database_connection = mysql.connector.connect(
        host=self.host,
        user=self.username,
        passwd=self.password,
        database=self.database
    )
    return database_connection
```


In this method, an establishing the connection to the database using data provided to *init_()* and method of **mysql.connector** and returning it to variable in *init_()*.

Definition of get_all_info_from_database()

```
def get_all_info_from_database(self):
    command = 'SELECT * FROM data_records'
    self.cursor.execute(command)
    records = self.cursor.fetchall()
    return records
```

In this method, an executing simple query written and stored in variable *command* and returning all observations from the database.

Definition of get_info_by_station()

```
def get_info_by_station(self, station_code: str = 'LV01'):
    command = 'SELECT * FROM data_records WHERE stacijas_kods = %s'
    self.cursor.execute(command, (station_code,))
    records = self.cursor.fetchall()
    return records
```

In this method, an executing the query to the database and returning all rows in the database with data from the station passed in parameters or defined as the default parameter. In variable *command*, the SQL command is constructed to make a proper query.

Definition of get_info_by_stations()

```
def get_info_by_stations(self, station_codes: list) -> list:
    info = []
    for station_code in station_codes:
        command = 'SELECT * FROM data_records WHERE stacijas_kods = %s'
        self.cursor.execute(command, (station_code,))
        records = self.cursor.fetchall()
        info.append(records)
    else:
        return info
```

This method does almost everything the same as the previous one; however, for many stations and returns, a list of records from many stations.

4. MONGODB CONNECTION

I. *Difference between database clients*

Comparing with the MySQL module, the program has methods both for importing and exporting data if a user with a login/password that is passed to *MongoDBClient* has permission to do so.

This program is capable of working with MongoDB as well. Methods are located in **MongoDBClient.py** and **AutoDBFiller.py**. They use **pymongo** module to connect to the MongoDB cluster, where databases with collections of data are stored.

A simple instruction to start:

1. In the case with MySQL, first of all, is creating a client that will connect to the cluster:

```
mongo_db_client = MongoDBClient('login', 'password', 'MeteoInfoTable', 'LastInsertTable')
```

In arguments, passing login and password to connect to the cluster, name of the database where observations are stored, and name of the database where the time of the last insert is stored.

Definitions

Definition of MongoDBClient_init() also known as class constructor

```
def init(self, login: str, password: str, main_database_name: str, time_database_name: str):
    self.my_client = self.get_connection(login, password)
    self.main_database_name = main_database_name
    self._main_database = self.my_client[self.main_database_name]
    self.time_database_name = time_database_name
    self._time_database = self.my_client[self.time_database_name]
    self.is_it_first_filling = self.is_not_first_filling_made()
```

In this method, creating our client by establishing a connection and defining essential variables.

Definition of get_connection()

```
def get_connection(login, password):
    return pymongo.MongoClient(
        "mongodb+srv://xxxxx".format(login, password))
```

In this method an establishing connection using login and password provided to *init()* and returning it to class constructor.

II. *MongoDB Importing*

1. `records = mongo_db_client.get_all_info_from_main_database()`

Now it is stored all information from main database in list of lists with dictionaries inside. Consider this as two-dimension array of dictionaries like `[[{}], [{}], [{}], [{}], [{}], [{}]]`

2. `records = mongo_db_client.get_data_from_collection('A1 km 12 Ādaži')`

If needed to get info from one station, it is possible to use this method and pass the name of the station.

```
3. records = mongo_db_client.get_data_from_collections(['A1 km 12  
Ādaži', 'A1 km 21 Lilaste'])
```

In case needed observations from many stations it is possible to use this method and pass list of station names.

4. Please, be attentive that in the case of MySQL database, station codes were used but with MongoDB names of stations!
5. How to process this data correctly and use it in estimations is described later.

Definitions

Definition of `get_all_info_from_main_database()`

```
def get_all_info_from_main_database(self) -> List:  
    collection_names = self._main_database.list_collection_names()  
    list_of_info = []  
    temp_list = []  
    for name in collection_names:  
        for search_result in self._main_database[name].find():  
            search_result.pop('_id')  
            search_result['Station'] = name  
            temp_list.append(search_result)  
            list_of_info.append(temp_list)  
            temp_list = []  
  
    return list_of_info
```

In this method, calling for all collections present in the main database, and depending on names returned in list requests are being made in the loop. In *for* loop, we are deleting `_id` keyword with the value from every returned result because it's unnecessary and adding `Station` keyword with a value of the name of collection getting at this moment, so every result has the name of the station it relates to. As a result, getting a list of lists of dictionaries as it was described earlier.

Definition of `get_data_from_collection()`

```
def get_data_from_collection(self, station_name: str) -> List:  
    list_of_info = []  
    for search_result in self._main_database[station_name].find():  
        search_result.pop('_id')  
        search_result['Station'] = station_name  
        list_of_info.append(search_result)  
    else:  
        return list_of_info
```

In this method, passing the name of one station as a parameter and making a request to database, processing results of search the same as in previous method and returning list of info.

Definition of `get_data_from_collections()`

```
def get_data_from_collections(self, station_names: list) -> list:
    list_of_info = []

    for station_name in station_names:
        for search_result in self._main_database[station_name].find():
            search_result.pop('_id')
            search_result['Station'] = station_name
            list_of_info.append(search_result)
    else:
        return list_of_info
```

In this method, passing the list of station names, but the logic is the same as in two previous methods. The only difference is that it iterates through, the given list of names, that contain those names which the user-provided.

III. MongoDB Exporting

1. To use all of functionality it should recreate client and use `AutoDBFiller()`

```
mongo_db_client = AutoDBFiller('login', 'password',
                                'MeteoInfoTable', 'LastInsertTable')
```

Parameters are the same as in case with `MongoDBClient()` because *class **AutoDBFiller*** is child-class of *class **MongoDBClient*** and has all of it methods and some new.

2. As described earlier, it is possible to use requests to take data from the website, but this data has to be processed before it can be exported to the database, so be sure to check the Data Handling chapter and methods presented there.

5. DATA HANDLING

Created **DataHandler.py** with methods to process data in different ways to prepare it for different usage. This chapter will be separated in many small to systemize data and increase the readability of the document.

I. For exporting to MongoDB

After getting the list of lists of values from the website, need to turn it into a JSON-like document to be able to export it to MongoDB because this database contains all records in JSON. Python is a data type called the *dictionary*, so a set of methods that will help to prepare all the data to export.

```
data_dicts = DataHandler.get_data_dicts(data_lists)
```

Here using variable *data_lists* from the Website Connection chapter, which contains a list of lists of values from the website. After calling this method, turning a list of lists into a list of dictionaries like [{}, {}, {}]

Now data can be exported to the database. This process will be described in the MongoDB Exporting chapter.

Definitions

Definition of get_data_dicts()

```
def get_data_dicts(data_lists) -> list:
    data_dicts = []

    for data in data_lists:
        data_dict = {'Station': data[0],
                    'Time': data[1],
                    'Date': DataHandler.get_date(),
                    'Air Temperature': data[2],
                    'Air Temperature(-1 h)': data[3],
                    'Humidity': data[4],
                    'Dew Point': data[5],
                    'Precipitation': data[6],
                    'Intensity': data[7],
                    'Visibility': data[8],
                    'Road Temperature': data[9],
                    'Road Temperature(-1 h)': data[10],
                    'Road Condition': data[11],
                    'Road Warning': data[12],
                    'Freezing Point': data[13],
                    'Road Temperature 2': data[14],
                    'Road Temperature 2(-1 h)': data[15],
                    'Road Condition 2': data[16],
                    'Road Warning 2': data[17],
                    'Freezing Point 2': data[18]}
        data_dicts.append(data_dict)

    return data_dicts
```

This method iterates through a list of data, takes values from the list by indexes they are stored, and append created dictionary to list which will be returned in the result.

Definition of get_date()

```
def get_date() -> str:
    from datetime import datetime

    year = datetime.now().year
```

```

month = datetime.now().month
day = datetime.now().day

return 'Year: {} Month: {} Day: {}'.format(year, month, day)

```

This method takes current date-time, adds its parts to string (year, month, and day), and returns this string. It uses **datetime** module to obtain the current date and time.

II. Process data from MongoDB

Data from MongoDB database

```

lists_of_measurements, list_of_station_names =
DataHandler.get_prepared_lists_for_estimation(main_data, 'Dew Point')

```

Called method returns two lists: with measurements (observations) and with station names whose data returned. Because of the chance that sometimes stations haven't been gathering data program checks every list and, if it is possible to fill missing data, it fills, if not – deletes the list of data. Station names aren't needed for estimation but are essential for building plots. Plots are described later.

Definitions

Definition of `get_prepared_lists_for_estimation()`

```

def get_prepared_lists_for_estimation(main_data, value: str = 'Dew Point') -> tuple:
    lists_of_chosen_values = DataHandler.get_lists_of_chosen_values_with_station_names(main_data, value)

    lists_of_measurements_with_station_names = DataHandler.get_lists_of_measurements_with_station_names(
        lists_of_chosen_values)

    list_of_station_names = DataHandler.get_station_names(lists_of_measurements_with_station_names)

    lists_of_measurements_without_station_names = DataHandler.get_lists_of_measurements_without_station_names(
        lists_of_measurements_with_station_names)

    return lists of measurements without station names, list of station names

```

This method does a lot, but basically, it does what it has to. Here are four steps:

1. Getting a list of lists of chosen values or values defined as a default parameter, and at the end of every list, appending the name of the station to which it is related. For this goal is used `get_lists_of_chosen_values_with_station_names()`.
2. Because all values are stored as text(string type), to convert it to numbers to be able to work with data. For this goal is used `get_lists_of_measurements_with_station_names()` In this method are being checked lists of data is it possible to fill missing data and if it is possible, it is being filled. As a result, suitable lists returned. Those ones who weren't fillable were deleted.
3. Now the program starts dividing the list of data to two lists: list of station names and list_of_measurements. For first goal is used `get_station_names()`.
4. For last step is used `get_list_of_measurements_without_station_names()`.

Definition of `get_lists_of_chosen_values_with_station_names()`

```

def get_lists_of_chosen_values_with_station_names(lists_of_dicts: list, value_name: str) -> list:
    list_of_values = []
    temp_list = []

    for list_of_dicts in lists_of_dicts:
        for data_dict in list_of_dicts:
            temp_list.append(data_dict[value_name])
        else:
            temp_list.append(data_dict['Station'])
        list_of_values.append(temp_list)

```

```
temp_list = []
```

```
return list_of_values
```

Here is a simple method that fills the list with lists of values that are being passed as argument and station names values belong to.

Definition of `get_lists_of_measurements_with_station_names()`

```
def get_lists_of_measurements_with_station_names(lists_of_values: List) -> List:
    lists_of_measurements = []

    for list_of_values in lists_of_values:
        if '-' in list_of_values or '' in list_of_values:
            list_of_values.pop()
            if DataHandler.is_possible_to_fill_missing_data(list_of_values):
                indexes = DataHandler.get_indexes_for_filling(list_of_values)
                list_of_values = DataHandler.fill_missing_data(list_of_values, indexes)
            else:
                continue
        lists_of_measurements.append(DataHandler.get_list_of_float_numbers_and_station(list_of_values))

    return lists_of_measurements
```

This method iterates through lists of values, and if there are '-' or empty spots, these lists are sent to the special method to be checked if it is possible to fill missing data, and if it is, missing data will be filled. If it is not possible, these lists will not be added to the final list of data and then deleted from memory. Before appending list of values to main lists all values are turned to *float* type in `get_list_of_float_numbers_and_station()`.

Definition of `is_possible_to_fill_missing_data()`

```
def is_possible_to_fill_missing_data(list_of_values: List) -> bool:
    index = 0
    adder = 1

    while adder < 5 and index < len(list_of_values) - 1 and index + adder < len(list_of_values):
        if list_of_values[index] == '-' and list_of_values[index + adder] == '-':
            adder += 1
        else:
            index += adder
            adder = 1
    else:
        if '' in list_of_values:
            return False
        if adder > 4:
            return False
        elif 1 < adder < 5 and adder == len(list_of_values):
            return False
        else:
            return True
```

This method simply counts the quantity of missing data in a row, and if there is somewhere in the list, five and more missing data samples method returns *False*. It will return *False* also if there is just empty string because that means that all list is empty(a feature of the website).

Definition of `get_indexes_for_filling()`

```
def get_indexes_for_filling(list_of_values: list) -> tuple:

    index = 0
    adder = 1

    indexes = []

    while adder < 5 and index < len(list_of_values) - 1 and index + adder < len(list_of_values):
        if adder == 1 and list_of_values[index] == '-' and list_of_values[index + 1] != '-':
            indexes.append([index, 'one'])

        if list_of_values[index] == '-' and list_of_values[index + adder] == '-':
            adder += 1
        elif adder != 1:
            indexes.append([index - 1, index + adder, 'normal'])
            index += adder
            adder = 1
        else:
            index += adder
            adder = 1
    else:
        if adder != 1:
            indexes.append([index - 1, index + adder, 'in the end'])
```

This method searches for all spots in the list and creates a list of coordinates for another method to fill missing data later correctly. There are three types of cords: usual series of missing spots, when only one spot missing data (in real life there is very little chance that this situation will occur, but I left this variant just in case) and when series of missing spots are located at the end of the list. How these scenarios are processed is shown in the next method definition.

Definition of `fill_missing_data_in_list()`

```
def fill_missing_data_in_list(list_of_values: list, indexes: tuple) -> list:
    for index_list in indexes:
        if index_list[-1] == 'normal':
            average = (list_of_values[index_list[0]] + list_of_values[index_list[1]]) / 2

            for i in range(index_list[0] + 1, index_list[1]):
                list_of_values[i] = average
        elif index_list[-1] == 'one':
            average = (list_of_values[index_list[0] - 1] + list_of_values[index_list[0] + 1]) / 2
            list_of_values[index_list[0]] = average
        elif index_list[-1] == 'in the end':
            for i in range(index_list[0], index_list[1]):
                list_of_values[i] = list_of_values[index_list[0]]

    return list_of_values
```

This method returns the list with filled missing spots. It iterates through a list of coordinates for filling and depending on coordinates and type of variant of the positioning of missing spots filling occurs.

Definition of `get_list_of_float_numbers_and_station()`

```
def get_list_of_float_numbers_and_station(list_of_values: list) -> list:
    list_of_data = []

    for value in list_of_values:
        try:
            list_of_data.append(float(value))
        except ValueError:
            list_of_data.append(value)

    return list_of_data
```

The simple method to turn every value in the list from *str* to *float* except station name. The method returns a new list of data.

Definition of `get_station_names()`

```
def get_station_names(list_of_lists_of_data: list) -> list:
    list_of_station_names = []

    for list_of_data in list_of_lists_of_data:
        list_of_station_names.append(list_of_data[-1])

    return list_of_station_names
```

A simple method that just takes the last element from every list which is station name and appends it to list of station names.

Definition of `get_lists_of_measurements_without_station_names()`

```
def get_lists_of_measurements_without_station_names(
    list_of_lists_of_measurements_with_station_names: list) -> list:
    list_of_lists_of_measurements_without_station_names = []

    for list_of_measurements in list_of_lists_of_measurements_with_station_names:
        list_of_measurements.pop(len(list_of_measurements) - 1)
        list_of_lists_of_measurements_without_station_names.append(list_of_measurements)

    return list_of_lists_of_measurements_without_station_names
```

A simple method that iterates through every list deletes the last element of every list, which is station name and append these lists to the final list that will be returned.

III. Process data from MySQL

```
index = DataHandler.get_index_of_value('Dew Point')
```

If needed to estimate data from MySQL database first needs to be done is get index of value that are needed to make estimation with.

```
list_of_measurements = DataHandler.get_exact_value_from_my_sql_records(records, index)
```

Now after getting exact value from records returned from MySQL database.

IMPORTANT NOTE!

```
records = sql_client.get_info_by_stations(['LV01', 'LV02', 'LV03'])
index = DataHandler.get_index_of_value('Dew Point')
lists_of_measurements = DataHandler.get_exact_value_from_many_my_sql_records(records,
index)
```

If needed to get values from some exact stations than should use other method to extract value because previous method is only for processing one station.

```
lists_of_measurements = DataHandler.get_lists_of_floats(lists_of_measurements)
list_of_measurements = DataHandler.get_list_of_floats(list_of_measurements)
```

Depending on situation it should be using method for process info from one station or from many stations.

At this moment, it should have a list of values or lists of values. However, in this list/lists may be missing values of type *None*. A simple method to replace them with a dash of type *str*, so these lists could be suitable for usage with methods for filling missing data I described earlier.

```
lists_of_measurements =
DataHandler.replace_none_with_dash(lists_of_measurements)
```

This method accepts a list of lists and returns it after processing.

Now passing lists for filling missing data if it is possible.

```
lists_of_measurements = DataHandler.fill_missing_data_in_lists(lists_of_measurements)
```

This method accepts a list of lists as well and returns it. Some details of functionality will be described in the next Definition chapter.

However, in lists that returned are present elements that can't be used in estimation due to not every list, missing data were filled.

Now the only option is to check which lists contain *str* elements, dashes put instead of *None*, or empty at all and remove them from the main list of measurements. Besides, station codes should be updated because our data needs to be related to the correct station.

```
station_codes = DataHandler.get_station_codes()
```

The first step – just create a list of station codes using this simple method.

```
station_codes, lists_of_measurements =  
DataHandler.zip_codes_and_measurements(station_codes, lists_of_measurements)
```

Now passing our lists in the following method, and it will return corrected lists that can be used in estimation and plot building.

After these manipulations, data is ready for estimation. The process of evaluation is described later.

Definitions

Definition of `get_index_of_value()`

```
def get_index_of_value(value: str = 'Dew Point') -> int:  
    if value == 'id':  
        return 0  
    elif value == 'Station code':  
        return 1  
    elif value == 'Datetime':  
        return 2  
    elif value == 'Air Temperature':  
        return 3  
    elif value == 'Air Temperature(-1 h)':  
        return 4  
    elif value == 'Humidity':  
        return 5  
    elif value == 'Dew Point':  
        return 6  
    elif value == 'Precipitation':  
        return 7  
    elif value == 'Intensity':  
        return 8  
    elif value == 'Visibility':  
        return 9  
    elif value == 'Road Temperature':  
        return 10  
    elif value == 'Road Temperature(-1 h)':  
        return 11  
    elif value == 'Road Condition':  
        return 12  
    elif value == 'Road Warning':  
        return 13  
    elif value == 'Freezing Point':  
        return 14  
    elif value == 'Road Temperature 2':  
        return 15  
    elif value == 'Road Temperature 2(-1 h)':  
        return 16  
    elif value == 'Road Condition 2':  
        return 17  
    elif value == 'Road Warning 2':
```

```

        return 18
    elif value == 'Freezing Point 2':
        return 19

```

A simple method that returns *int* as an index of value passed in argument or which is defined as a default parameter.

Definition of get_exact_value_from_my_sql_records()

```

def get_exact_value_from_my_sql_records(records: tuple, index: int) -> list:
    values = []

    for record in records:
        values.append(record[index])

    return values

```

A simple method that just iterates through records and takes elements by index and adds them in a list that will be returned in the end.

Definition of get_exact_value_from_many_my_sql_records()

```

def get_exact_value_from_many_my_sql_records(list_of_records: list, index: int) -> list:
    values = []

    for records in list_of_records:
        values.append(DataHandler.get_exact_value_from_my_sql_records(records, index))

    return values

```

This method does the same as the previous one but for records from many stations. Returns list of lists.

Definition of replace_none_with_dash()

```

def replace_none_with_dash(lists_of_measurements: list) -> list:
    i = 0

    for list_of_measurements in lists_of_measurements:
        for _ in list_of_measurements:
            if list_of_measurements[i] is None:
                list_of_measurements[i] = '-'
                i += 1
            i = 0

    return lists_of_measurements

```

This simple method just iterates through every element of every row and, if it is *None* type, replaces it on '-' character.

Definition of fill_missing_data_in_lists()

```

def fill_missing_data_in_lists(lists_of_measurements: list) -> list:
    i = 0
    for list_of_measurement in lists_of_measurements:
        if DataHandler.is_possible_to_fill_missing_data(list_of_measurement):
            indexes = DataHandler.get_indexes_for_filling(list_of_measurement)
            list_of_measurement = DataHandler.fill_missing_data_in_list(list_of_measurement, indexes)
            lists_of_measurements[i] = list_of_measurement
            i += 1
        else:
            i += 1
    return lists_of_measurements

```

This method iterates through every list and checks if it is possible to fill missing data or not and if it is, it makes filling using already described methods.

IMPORTANT NOTE!

This method works with the whole list at a time, so if even there are only one series of missing data that couldn't be filled, the entire list will be ignored.

Definition of `get_station_codes()`

```
def get_station_codes() -> list:
    station_codes = []

    for i in range(1, 65):
        station_code = 'LV{:02d}'.format(i)
        station_codes.append(station_code)

    return station_codes
```

This method simply creates a list, fills it with station codes in *str* type, and returns it in the result.

Definition of `zip_codes_and_measurements()`

```
def zip_codes_and_measurements(station_codes: list, lists_of_measurements: list) -> tuple:
    stop = len(lists_of_measurements)
    i = 0

    while i < stop:
        if '-' in lists_of_measurements[i] or not bool(lists_of_measurements[i]):
            lists_of_measurements.pop(i)
            station_codes.pop(i)
            stop -= 1
            i = 0
        else:
            i += 1

    return station_codes, lists_of_measurements
```

This method checks every list is it containing dashes or is empty, and if it is, then it deletes this list from the main list and corresponding station code from the list of station codes. As a result, modified lists are returned.

IV. Accuracy

This chapter will be useful after getting familiar with the Estimation chapter because skill is counted using measurements and estimates. Anyway, the calculation of efficiency is handled by **DataHandler.py**, so this process will be described here.

Few steps should be made to calculate accuracy:

1. Cast all data values to *int* type

```
lists_of_estimates = DataHandler.get_lists_of_ints(lists_of_estimates)
lists_of_measurements = DataHandler.get_lists_of_ints(lists_of_measurements)
```

2. Use method, pass lists with data as argument and will get list of accuracies

```
accuracies = DataHandler.get_accuracies(lists_of_measurements, lists_of_estimates)
```

Now having list of accuracies of stations that were used in estimation process.

Definitions

Definition of `get_lists_of_ints()`

```
def get_lists_of_ints(lists_of_values) -> list:
    for values in lists_of_values:
        i = 0
        while i < len(values):
            values[i] = int(values[i])
            i += 1

    return lists_of_values
```

This method iterates through every element of every list and casts it to *int*.

Definition of get_accuracies()

```
def get_accuracies(lists_of_measurements: list, lists_of_estimates: list) -> list:
    accuracies = []
    for measurements, estimates in zip(lists_of_measurements, lists_of_estimates):
        accuracies.append(DataHandler.get_accuracy(measurements, estimates))

    return accuracies
```

This method iterates through every list of measurements, estimates, and appends the result of `get_accuracy()` method, which is described below.

Definition of get_accuracy()

```
def get_accuracy(measurements: list, estimates: list) -> float:
    from sklearn.metrics import accuracy_score

    return accuracy_score(measurements, estimates)
```

This method uses a method from **sklearn** module to calculate accuracy and returns the result.

6. ESTIMATION

For estimation purposes, **KalmanFilter.py** is created that has a basic implementation of the Kalman Filter for a one-dimension model.

```
kalman_filter = KalmanFilter(error_in_estimate, initial_estimate, error_in_estimate,
measurements)
```

The first step is to create Kalman filter object and pass needed arguments: error in the estimate, an initial estimate, error in the estimate, and list of measurements that will be used in the estimation process.

```
estimates = kalman_filter.get_estimates()
```

Now it can get a list of estimates by using only one method.

However, if working with a list of lists, then can use the static method from **KalmanFilter.py** that will handle essential processes and return us a list of lists of estimates.

```
lists_of_estimates = KalmanFilter.get_lists_of_estimates(lists_of_measurements,
error_in_estimate, error_in_measurement)
```

In arguments, passing lists of measurements, errors in the estimate, and error in measurement. As a result, getting our estimations that can be used for building plots. Plots are described later.

Definitions

Definition of KalmanFilter.__init__() also known as class constructor

```
def __init__(self, initial_error_in_estimate, initial_estimate, __error_in_measurement,
__measurements):
    self.__measurements = __measurements
    self.__estimate = __initial_estimate
    self.__error_in_estimate = initial_error_in_estimate
    self.__error_in_measurement = __error_in_measurement
    self.__measurement = None
    self.__kalman_gain = None
```

In this method, is just declare essential variables for the estimation process.

Definition of get_estimates()

```
def get_estimates(self) -> List:
    list_of_estimates = []

    for measurement in self.__measurements:
        self.__measurement = measurement
        self.make_basic_calculations()
        list_of_estimates.append(self.__estimate)

    return list_of_estimates
```

This method iterates through all measurements in the list of measurements that are passed when having been creating KalmanFilter() object. *make_basic_calculations()* does exactly what it stands for; its definition is provided below. After every calculation new estimate is appended in the list of estimates that will be returned in the result.

Definition of make_basic_calculations()

```
def make_basic_calculations(self):
    self.__kalman_gain = self.__calculate_kalman_gain()
    self.__estimate = self.__calculate_estimate()
    self.__error_in_estimate = self.__calculate_error_in_estimate()
```

This method uses other simple methods to produce Kalman Filter calculations for a one-dimension model and make an estimation. The description of every method is given below.

Definition of calculate_kalman_gain()

```
def calculate_kalman_gain(self) -> float:  
    return self._error_in_estimate / (self._error_in_estimate + self._error_in_measurement)
```

Definition of calculate_estimate()

```
def calculate_estimate(self) -> float:  
    return self._estimate + self._kalman_gain * (self._measurement - self._estimate)
```

Definition of calculate_error_in_estimate()

```
def calculate_error_in_estimate(self) -> float:  
    return (1 - self._kalman_gain) * self._error_in_estimate
```

Definition of get_lists_of_estimates()

```
def get_lists_of_estimates(lists_of_measurements: list, error_in_estimate: float,  
                           error_in_measurement: float) -> list:  
    lists_of_estimates = []  
  
    for measurements in lists_of_measurements:  
        initial_estimate = KalmanFilter.__get_initial_estimate(measurements)  
        kf = KalmanFilter(error_in_estimate, initial_estimate, error_in_measurement, measurements)  
        lists_of_estimates.append(kf.get_estimates())  
  
    return lists_of_estimates
```

This method iterates through lists of measurements, automatically creates initial estimate, KalmanFilter object, and append list of estimates to final list that will be returned in the result.

Definition of get_initial_estimate()

```
def get_initial_estimate(measurements: list) -> float:  
    return measurements[-1] + (  
        measurements[-1] - measurements[-2])
```

This method returns an average of last two elements.

7. PLOTS

This program has **GraphEditor.py** that uses **matplotlib** module to create plots.

To create one plot or plots for a list of estimates, it is needed to create an object of *class* **GraphEditor** and use correct proper methods:

```
graph = GraphEditor(estimates, measurements, value, station_code, period)
graph.create_est_and_meas_plot()
```

Also, there are variants of plots that can be created:

```
graph.create_est_plot()
graph.create_meas_plot()
```

Creating separate plots for estimates and measurements.

To show the created plot, the following method should be used:

```
graph.show_plot()
```

It will display a plot in the program that are used, but if need to save the created plot, this method should be used for saving:

```
graph.save_plot()
```

IMPORTANT NOTE!

Don't show plot and then try to save it, it will not work. Separately these methods work just fine.

Definitions

Definition of GraphEditor_init also known as class constructor

```
def __init__(self, estimates: list, measurements: list, value: str, station_name: str,
period: str, accuracy=None):
    self.estimates = estimates
    self.measurements = measurements
    self.value = value
    self.station_name = station_name
    self.period = period
    self.accuracy = accuracy
    self.plt = None
```

Initialization of class object and defining variables.

Definition of create_est_and_meas_plot()

```
def create_est_and_meas_plot(self):
    plt.plot(self.estimates, label='estimates', color='blue')
    plt.plot(self.measurements, label='measurements', color='orange')
    plt.title(self.station_name + ', Period: ' + self.period)
    plt.ylabel(self.value)
    plt.xlabel('Observations')
    plt.legend()
    self.plt = plt
```

This method setups settings for plot such as title, legend, labels for x-y axis, and saves plot for further usage like saving or showing in the program.

Definition of create_est_plot()

```
def create_est_plot(self):
    plt.plot(self.estimateds, color='blue')
    plt.title('Estimates: ' + self.station_name + ', Period: ' + self.period)
    plt.ylabel(self.value)
    plt.xlabel('Observations')
    self.plt = plt
```

This method does the same things but only for estimations.

Definition of create_meas_plot()

```
def create_meas_plot(self):
    plt.plot(self.measurements, color='orange')
    plt.title('Measurements: ' + self.station_name + ', Period: ' + self.period)
    plt.ylabel(self.value)
    plt.xlabel('Observations')
    self.plt = plt
```

This method creates a plot for measurement in the way, like two previous methods.

Definition of show_plot()

```
def show_plot(self):
    self.plt.show()
```

Simple method that shows plot in program.

Definition of save_plot()

```
def save_plot(self):
    self.plt.savefig('Plot{}.pdf'.format(GraphEditor.index), dpi=300)
    GraphEditor.index += 1
```

A method that saves plot to project folder with settings that were passed as arguments for *savefig()* method.

8. FORECASTING MODELS

I. Models

For creating a forecasting model is used `Forecasting_model` class instance, where can be specified model name, steps for making forecast, optimize value, max p and max q value for models.

Definition of `__init__()`

```
def __init__(self, model_name, data_series, steps: int, station_code: str,
              max_p_value=6, max_q_value=6):

    self.model_name = model_name
    self.data_series = data_series
    self.steps = steps

    if len(self.data_series) <= 100:
        raise Exception("Data must contains more than 100 measurements")

    self.station_code = station_code

    self.min_aic_value = None
    self.data_points_to_use = None
    self.best_model_order = None

    # range (0, 3) == [ 0, 1, 2 ]
    if model_name == "ARIMA":
        self.d_values = range(0, 3)
        self.p_values = range(0, max_p_value)
        self.q_values = range(0, max_q_value)
    elif model_name == "ARMA":
        self.d_values = range(0, 1)
        self.p_values = range(0, max_p_value)
        self.q_values = range(0, max_q_value)
    elif model_name == "AR":
        self.d_values = range(0, 1)
        self.q_values = range(0, 1)
        self.p_values = range(0, max_p_value)
    elif model_name == "MA":
        self.d_values = range(0, 1)
        self.p_values = range(0, 1)
        self.q_values = range(0, max_q_value)
    else:
        raise Exception("The model name {} is incorrect".format(model_name))
```

II. Forecast

For making forecast with model is used **Forecast_model.py** with modules such as **pandas** for working with data frames and data series, **numpy** and **ARIMA** from **statsmodels.tsa.arima_model**. For making forecast and receive all information that was used for forecasting, it is needed to create class instance and use correct method as shown here:

```
arima_model = ForecastingModel("ARIMA", series_of_measurements, steps, station_code)
forecast_arima = arima_model.get_forecast()
arima_data_points = arima_model.data_points_to_use
arima_model_order = arima_model.best_model_order
```

```

arma_model = ForecastingModel("ARMA", series_of_measurements, steps, station_code)
forecast_arma = arma_model.get_forecast()
arma_data_points = arma_model.data_points_to_use
arma_model_order = arma_model.best_model_order

ar_model = ForecastingModel("AR", series_of_measurements, steps, station_code)
forecast_ar = ar_model.get_forecast()
ar_data_points = ar_model.data_points_to_use
ar_model_order = ar_model.best_model_order

ma_model = ForecastingModel("MA", series_of_measurements, steps, station_code)
forecast_ma = ma_model.get_forecast()
ma_data_points = ma_model.data_points_to_use
ma_model_order = ma_model.best_model_order

```

IMPORTANT NOTE!

series_of_measurements must contain only one data column and one index column, because for making forecast are used *series_of_measurements* values.

Definition of get_forecast()

```

def get_forecast(self):
    models_for_all_data_points_dict =
self.get_possible_models_for_all_data_points_dict()

self.get_lowest_aic_value_and_data_points_from_dict(models_for_all_data_points_dict)

    if self.data_points_to_use is None:
        return None

    data_series_copy = self.data_series[-self.data_points_to_use:].copy()

    forecast = self.make_forecast(data_series_copy)

    return forecast

```

As a result we get a list of forecast values.

Definition of get_possible_models_for_all_data_points_dict()

```

def get_possible_models_for_all_data_points_dict(self) -> dict:
    if len(self.data_series) <= 1000:
        data_optimize_points = len(self.data_series)
        point_range = int(data_optimize_points / 10)
    else:
        data_optimize_points = 1010
        point_range = 100

    models_order_and_data_points_dict = dict()
    last_point = len(self.data_series)

    data_point_range = range(100, data_optimize_points, point_range)
    for data_points in data_point_range:
        first_point = last_point - data_points
        if first_point < 0:
            break

```

```

data_series_copy = self.data_series[first_point:last_point].copy()

order_and_min_aic_value = self.get_order_and_min_aic_value(data_series_copy)
if order_and_min_aic_value is None:
    continue
else:
    models_order_and_data_points_dict[data_points] = order_and_min_aic_value

return models_order_and_data_points_dict

```

This method returns a dictionary with data points and with the best ARIMA model order for this data set. It looks like this:

```
{ data_points1 : { (model_order_1) : aic_value_1 }, data_points2 : { (model_order_2) : aic_value_2 }, ... }
```

or

```
{ 100 : { (3, 1, 4) : 105 }, 110 : { (4, 1, 4) : 209 } ... }
```

Definition of make_forecast()

```

def make_forecast(self, data_series):

    model = ARIMA(data_series, order=self.best_model_order).fit(dispatch=False)
    forecast_list = model.forecast(steps=self.steps)[0].tolist()

    return forecast_list

```

Method uses the best ARIMA model order, that was find before, and returns list of forecast value. The number of received values in the list is equal with steps value.

III. Model order

To determine model order is used with modules **adfuller** test from **statsmodels.tsa.stattools** for finding p-value, that indicates if model is stationary or not, also can be used **plot_pacf**, **plot_acf** from **statsmodels.graphics.tsaplots** and **pyplot** for showing autocorrelation and partial autocorrelation graphs.

Definition of get_d_value_and_ADF_test()

```

def get_d_value_and_ADF_test(data_series):

    # ----- Augmented Dickey-Fuller test ----- #
    # ----- p-value > 0.05: the data has a unit root and is non-stationary ----- #
    # ----- p-value <= 0.05: the data does not have a unit root and is stationary --
    --- #
    # ----- The more negative is ADF Statistic, the more likely we have a stationary
    dataset ----- #

    d = 0
    try:
        adf_test_results = adfuller(data_series.values)
    except:
        return d

    data_series_diff = data_series
    while adf_test_results[1] > 0.05 or d == 0:
        if d > 2:
            return 0

```

```

# ----- make data stationary and drop NA values ----- #
data_series_diff = data_series_diff.diff().dropna()
d += 1
try:
    data_series_diff_values = data_series_diff.values
    adf_test_results = adfuller(data_series_diff_values)
except:
    d -= 1
    return d

return d

```

This method is used to define d value for ARIMA model (p, d, q) order. d represents the number of times that the data have to be “differenced” to produce a stationary signal (i.e., a signal that has a constant mean over time). This captures the “integrated” nature of ARIMA. If d=0, this means that our data does not tend to go up/down in the long term (i.e., the model is already “stationary”). In this case, then technically you are performing just ARMA, not AR-IMA. If p is 1, then it means that the data is going up/down linearly. If p is 2, then it means that the data is going up/down exponentially. To define d value is used Augmented Dickey-Fuller test with p-value:

- p-value > 0.05: the data has a unit root and is non-stationary.
- p-value <= 0.05: the data does not have a unit root and is stationary.
- The more negative is ADF Statistic, the more likely we have a stationary dataset.

Definition of get_order_and_min_aic_value()

```

def get_order_and_min_aic_value(self, data_series):
    aic_dict = dict()

    if len(self.d_values) != 1:
        d = self.get_d_value_and_ADF_test(data_series)
    else:
        d = 0

    for p in self.p_values:
        for q in self.q_values:
            order = (p, d, q)

            if order == (0, 0, 0):
                continue

            try:
                arima_model = ARIMA(data_series, order).fit(dispatch=False)
                aic = arima_model.aic
                aic_dict[order] = aic
                # print("Data Points: {}, Order: {}, AIC:
                {}.format(len(data_series), order, aic))
            except:
                continue

    # if aic_dict is empty
    # it is impossible to create arima model for this data set
    if len(aic_dict) == 0:
        return None
    else:
        min_aic_val = min(aic_dict.values())

        for key, value in aic_dict.items():

```

```
if value == min_aic_val:
    result_dict = {key: value}
    # order = (p, d, q)
    # {order: aic_value }
    return result_dict
```

This method returns best model order using Akaike information criterion for characterizing and d value from *d_value_and_ADF_test* method. AIC estimates the relative amount of information lost by a given model: the less information a model loses, the higher the quality of that model

9. TIME HANDLING

Is used to store time values and extract data using datetime indexes.

For **TimaHandler.py** is used such modules such as **pandas** for working with data frames and data series, **datetime** and **dateutil.relativedelta** for getting datetime values.

Definition of `__init__()`

```
def __init__(self, datetime_unit_of_measure: str, period_value: int,
             date_from: str = None, date_till: str = None, list_of_datetime: list =
             None):

    if date_from is None:
        if list_of_datetime is None:
            raise Exception("Missing list_of_datetime value")
        else:
            self.date_from = list_of_datetime[0]
    else:
        self.date_from = self.get_datetime_format(date_from)

    if date_till is None:
        if list_of_datetime is None:
            raise Exception("Missing list_of_datetime value")
        else:
            self.date_till = list_of_datetime[-1]
    else:
        self.date_till = self.get_datetime_format(date_till)

    # datetime delta number
    self.period_value = period_value

    # minute, year, second, hour ...
    self.datetime_unit_of_measure = datetime_unit_of_measure
    # 1 minute
    self.datetime_value_of_one_unit = self.get_datetime_value_of_one_unit()
    # 11 minutes, 1 hour, 15 days ...
    self.datetime_delta_value = self.get_datetime_delta_value()
```

Definition of `get_datetime_format()`

```
def get_datetime_format(date_time: str):
    if len(date_time) == 10:
        date_time += " 00:00:00"
    else:
        date_time += ":00"
    date_time = datetime.datetime.strptime(date_time, format("%Y-%m-%d %H:%M:%S"))
    return date_time
```

Returns datetime format.

Definition of `get_datetime_from_series()`

```
def get_datetime_from_series(self, data_series, date_from_request):

    max_number_of_iterations =
    TimeHandler.get_max_number_of_iterations(self.datetime_unit_of_measure)
    i = 0
    while data_series.loc[data_series.index == date_from_request].empty:
        date_from_request += self.datetime_value_of_one_unit
```

```

        i += 1
        if i >= max_number_of_iterations:
            return -1

    return date_from_request

```

By the increasing *data_from_request* by one *datetime_value_of_one_unit* program try to find this datetime index in the series until reach the max number of possible iterations.

Definition of get_max_number_of_iterations()

```

def get_max_number_of_iterations(datetime_unit_of_measure):
    if datetime_unit_of_measure == "second":
        return 60
    elif datetime_unit_of_measure == "minute":
        return 60
    elif datetime_unit_of_measure == "hour":
        return 24
    elif datetime_unit_of_measure == "day":
        return 31
    elif datetime_unit_of_measure == "month":
        return 12
    elif datetime_unit_of_measure == "year":
        return 10
    else:
        raise Exception("Parameter time_period_unit must be timedelta value")

```

Based on *datetime_unit_of_measure* values returns max number of iterations, that can be used for finding data series datetime index. Without it process of finding datetime index in the series can last forever.

Definition of get_datetime_value_of_one_unit()

```

def get_datetime_value_of_one_unit(self):
    delta_value = 1
    if self.datetime_unit_of_measure == "second":
        datetime_value_of_one_unit = datetime.timedelta(seconds=delta_value)
        return datetime_value_of_one_unit
    elif self.datetime_unit_of_measure == "minute":
        datetime_value_of_one_unit = datetime.timedelta(minutes=delta_value)
        return datetime_value_of_one_unit
    elif self.datetime_unit_of_measure == "hour":
        datetime_value_of_one_unit = datetime.timedelta(hours=delta_value)
        return datetime_value_of_one_unit
    elif self.datetime_unit_of_measure == "day":
        datetime_value_of_one_unit = datetime.timedelta(days=delta_value)
        return datetime_value_of_one_unit
    elif self.datetime_unit_of_measure == "month":
        datetime_value_of_one_unit =
dateutil.relativedelta.relativedelta(months=delta_value)
        return datetime_value_of_one_unit
    elif self.datetime_unit_of_measure == "year":
        datetime_value_of_one_unit =
dateutil.relativedelta.relativedelta(years=delta_value)
        return datetime_value_of_one_unit
    else:
        raise Exception("Parameter time_period_unit must be timedelta value")

```

Returns one unit of datetime value.

Definition of get_datetime_and_measurements_series()

```
def get_datetime_and_measurements_series(self, list_of_measurements,
list_of_datetime):

    data_series = pd.Series(list_of_measurements, index=list_of_datetime,
dtype=float)

    date_from = self.get_datetime_from_series(data_series, self.date_from)
    date_till = self.get_datetime_from_series(data_series, self.date_till)

    if date_from == -1 or date_till == -1:
        return pd.Series()

    data_series = data_series.loc[date_from:date_till]

    return data_series
```

Returns series with datetime index format. If in series can't be found date will be returned empty series.

10. DATA EMPTY VALUES FILLING

Is used to fill empty data values – None or NaN values.

For **DataFiller.py** is used such modules such as **math** for detecting Nan values, **pandas** for working with data frames and data series, **DataHandler**, **MySQLClient** for extracting additional data and **TimeHandler**.

DattaFiller use data from nearest stations, that was defined using clustering method, for filling None or NaN values in the series.

Definition of __init__()

```
def __init__(self, data_series, station_clusters_df, station_code,
             time_handler: TimeHandler, sql_client: MySQLClient):

    self.station_code = station_code

    self.sql_client = sql_client
    self.time_handler = time_handler

    self.data_series = data_series
    self.station_clusters_df = station_clusters_df
```

Definition of get_list_of_stations_codes_to_use()

```
def get_list_of_stations_codes_to_use(self):
    cluster_position = self.station_clusters_df.loc[self.station_code]["cluster"]
    stations_to_use = self.station_clusters_df[(self.station_clusters_df["cluster"]
== cluster_position)
                                              & (self.station_clusters_df.index !=
self.station_code)]

    return list(stations_to_use.index)
```

Returns list of nearest stations that are located in the same cluster with the station witch data we try to fill.

Definition of fill_the_data_none_values()

```
def fill_the_data_none_values(self, value):
    stations_codes_to_use_list = self.get_list_of_stations_codes_to_use()

    series_s = []

    for station_code in stations_codes_to_use_list:
        records = self.sql_client.get_info_by_station(station_code)
        if not records:
            continue

        list_of_measurements = DataHandler.get_filled_list_of_measurements(records,
value)
        list_of_datetime =
DataHandler.get_exact_value_from_many_my_sql_records([records], 2)[0]
        series =
self.time_handler.get_datetime_and_measurements_series(list_of_measurements,
list_of_datetime)

        if not series.empty:
```

```

        series_s.append(series)

    for datetime, measurement in
self.data_series[self.data_series.isnull()].items():
        temp_measurement = 0
        i_temp = 0
        for series in series_s:
            temp_datetime = self.time_handler.\
                get_datetime_from_dataframe_with_delta_limitation(series, datetime,
self.time_handler.period_value)

            if temp_datetime != -1:
                temp_value = series.loc[temp_datetime]
                if type(temp_value) == pandas.Series:
                    continue
                if temp_value is not None and not math.isnan(temp_value):
                    temp_measurement += temp_value
                    i_temp += 1

        if i_temp != 0:
            temp_measurement = temp_measurement / i_temp
            self.data_series.loc[datetime] = temp_measurement

    return self.data_series

```

Goes through fo all data series that are located in the same cluster and extracts data value with the same datetime index or with the datetime index that was taken shortly before (in this way in used limitation for finding datetime index from **TimeHandler**)

As a result returns series with filled data points.

11. CLUSTERING

Clustering method is used for **DataFiler** for finding the stations. That are located in the one cluster. In the **Cluster.py** is used models such as **KMeans** for clustering algorithm that aims to partition n observations into k clusters, **MinMaxScaler** for transforming features by scaling each feature to a given range (This estimator scales and translates each feature individually such that it is in the given range on the training set, e.g. between zero and one), **math** for finding square root, **pandas** and **numpy**.

Definition of `__init__()`

```
def __init__(self, records, max_k_value=10, default_k_value=None,
path_to_save_results=None):
    self.k_range = range(1, max_k_value)
    self.records_df = self.get_records_df(records)
    self.lat_and_lng_array = self.get_lat_and_lng_array()

    self.path_to_save_results = path_to_save_results
    if path_to_save_results is not None:
        self.save_results = True
    else:
        self.save_results = False

    self.default_k_value = default_k_value

    self.dist_points_from_cluster_center = None
    self.distance_of_points_from_line = None

    if default_k_value is None:
        self.optimal_k_value = self.get_optimal_k_value()
    else:
        self.optimal_k_value = default_k_value
```

For the records value is used data with station codes and with latitude and longitude for each station. As a default value of *max_k_value* used integer 10. If *default_k_value* is not specified algorithm automatically finds the *optimal_k_value*.

Definition of `get_records_df()`

```
def get_records_df(records):
    records_df = pd.DataFrame.from_records(records, columns=["index", "code",
"name", "lat", "lng"])
    records_df = records_df[["code", "lat", "lng"]]

    records_df["lat"] = records_df["lat"].astype(float)
    records_df["lng"] = records_df["lng"].astype(float)

    return records_df
```

Returns a data frame with all stations with latitude and longitude as float numbers.

Definition of `get_lat_and_lng_array()`

```
def get_lat_and_lng_array(self):
    lat_and_lng_array = np.column_stack((self.records_df["lng"],
self.records_df["lat"]))
    return lat_and_lng_array
```

Returns 2D array for latitude and longitude using records data frame.

Definition of get_optimal_k_value()

```
def get_optimal_k_value(self):
    self.dist_points_from_cluster_center =
self.get_dist_points_from_cluster_center_using_k_range()
    self.distance_of_points_from_line = self.cal_distance_of_points_from_line()

    opt_value = self.distance_of_points_from_line.index(max(
        self.distance_of_points_from_line)) + 1
    return opt_value
```

Return the *optimal_k_value*.

Definition of make_clustering()

```
def make_clustering(self):
    km = KMeans(n_clusters=self.optimal_k_value)
    predicted = km.fit_predict(self.records_df[["lng", "lat"]])
    self.records_df["cluster"] = predicted
```

Make clustering and returns a data frame with station codes and with cluster nuber for each station.

Definition of make_scaler()

```
def make_scaler(self):
    scaler = MinMaxScaler()

    scaler.fit(self.records_df[["lat"]])
    self.records_df["lat"] = scaler.transform(self.records_df[["lat"]])

    scaler.fit(self.records_df[["lng"]])
    self.records_df["lng"] = scaler.transform(self.records_df[["lng"]])
```

12. KALMAN FILTER 1D AND 2D

For data smoothing is used Kalman Filter from **KF.py**. When data is one-dimensional it can be used class **KF_1D** and when your data contains position and velocity (x, v) must be used class **KF_2D**. For using with Kalman Filter must be installed modules such as **pykalman**.

I. One-dimensional Kalman Filter

For getting smoothed values must be used class **KF_1D** and as a input paramentrs must be used list of measurements:

```
filtered_values = KF_1D(list_of_measurements).get_filtered_values()
```

Definition of get filtered values()

```
def get_filtered_values(self):  
  
    kf = KalmanFilter(transition_matrices=[1],  
                      observation_matrices=[1],  
                      initial_state_mean=0,  
                      initial_state_covariance=1,  
                      observation_covariance=1,  
                      transition_covariance=0.05)  
  
    state_means, state_covariances = kf.smooth(self.data_values_list)  
  
    state_means_list = [state_mean[0] for state_mean in state_means]  
  
    return state_means_list
```

Returns a list of smoothed values.

II. Two-dimensional Kalman Filter

For getting smoothed values with position and velocity (x, v) must be used class **KF_2D** and as a input paramentrs must be used list of measurements, starting position, starting velocity and period of time when data was measured (time delta):

```
filtered_values = KF_2D(list_of_measurements, position=position, velocity=velocity,  
                        time_delta=time_delta).get_filtered_values()
```

Definition of get filtered values ()

```
def get_filtered_values(self):  
  
    kf = KalmanFilter(n_dim_obs=1, n_dim_state=2,  
                      initial_state_mean=[self.position, self.velocity],  
                      transition_matrices=[[1, self.time_delta], [0, 1]],  
                      observation_matrices=[[1, 0]],  
                      transition_offsets=[0.5 * self.time_delta ** 2,  
self.time_delta])  
  
    state_means, state_covariances = kf.smooth(self.data_values_list)  
  
    state_means_list = [state_mean[0] for state_mean in state_means]  
  
    return state_means_list
```

For the two-dimensional Kalman Filter is used:

- `initial_state_mean` or \mathbf{x}_k as a (x, v)
- `transition_matrices` or \mathbf{F} as a matrix form of $\begin{bmatrix} 1 & \text{self.time_delta} \\ 0 & 1 \end{bmatrix}$
- `observation_matrices` or \mathbf{H} as a vector form of $\begin{bmatrix} 1 & 0 \end{bmatrix}$
- `transition_offsets` or \mathbf{G} as a matrix form of $\begin{bmatrix} 0.5 * \text{self.time_delta} ** 2 & \text{self.time_delta} \end{bmatrix}$

Returns a list of smoothed values.

13. API

The fundamental representation of API that can be used as Web Service of this framework.

Twelve routes are used:

```
/use/1d_kalman_filter
```

```
/use/2d_kalman_filter
```

```
/get/forecast
```

```
/get/forecast/time_period
```

```
/get/forecast/all_models
```

All these requests have essential arguments and secondary ones. Every request has an essential argument *value*.

I. */use/1d_kalman_filter*

Is used for getting smoothed values using one-dimensional Kalman Filter .

Essential arguments: *value*, *station*, *datetime_unit_of_measure*, *period_value*.

For *value* argument can be used “Dew Point”, “Air Temperature” and etc.

For *datetime_unit_of_measure* must be used str value, that defines unit of measure of time delta value (“second”, “minute”, “hour”, “day”, “month”, “year”)

For *period_value* must be used int value, that defines delta value.

If *period_value* is 11 and *datetime_unit_of_measure* is “minute”, it means delta value will be 11 minutes.

For example:

[/use/1d_kalman_filter?station=LV01&value=Dew%20Point&datetime_unit_of_measure=minute&period_value=11](#)

II. */use/2d_kalman_filter*

Is used for getting smoothed values using two-dimensional Kalman Filter .

Essential arguments: *value*, *station*, *datetime_unit_of_measure*, *period_value*.

For *value* argument can be used “Dew Point”, “Air Temperature” and etc.

Secondary argument: *position*, *velocity* and *time_delta*

As a default values for secondary arguments:

position = 0

velocity = 0

time_delta = 0.1

For *datetime_unit_of_measure* must be used str value, that defines unit of measure of time delta

value (“second”, “minute”, “hour”, “day”, “month”, “year”)

For *period_value* must be used int value, that defines delta value.

If *period_value* is 11 and *datetime_unit_of_measure* is “minute”, it means delta value will be 11 minutes.

For example:

[/use/2d_kalman_filter?station=LV01&value=Dew%20Point&datetime_unit_of_measure=minute&period_value=11](#)

In this case will be used default values for secondary arguments. This request will return JSON with smoothed values for stations and value that are set up in request arguments.

[/use/2d_kalman_filter?station=LV01&value=Dew%20Point&position=1&velocity=9&time_delta=0.2&datetime_unit_of_measure=minute&period_value=11](#)

III. /get/forecast

Is used for getting forecast values using station data.

Essential arguments: *model_name*, *value*, *station*, *steps*, *datetime_unit_of_measure*, *period_value*.

Secondary arguments: *using_Kalman_Filter*.

For the *model_name* value can be used ARIMA, ARMA, AR and MA models.

For *value* argument can be used “Dew Point”, “Air Temperature” and etc.

Steps value means how many points ahead will be made forecast.

For *datetime_unit_of_measure* must be used str value, that defines unit of measure of time delta value (“second”, “minute”, “hour”, “day”, “month”, “year”)

For *period_value* must be used int value, that defines delta value.

If *period_value* is 11 and *datetime_unit_of_measure* is “minute”, it means delta value will be 11 minutes.

For the value *using_Kalman_Filter* must be used only True or False values.

For example:

[/get/forecast?model_name=ARIMA&value=Dew%20Point&station=LV01&steps=5&datetime_unit_of_measure=minute&period_value=11](#)

[/get/forecast?model_name=AR&value=Dew%20Point&station=LV06&steps=10&datetime_unit_of_measure=hour&period_value=1&using_Kalman_Filter=True](#)

IV. /get/forecast/time_period

Is used for getting forecast values using station data and using time period with date from and date till.

Essential argument: *model_name* , *value*, *station*, *steps*, *datetime_unit_of_measure*, *period_value*.

Secondary arguments: *using_Kalman_Filter*.

For the *model_name* value can be used ARIMA, ARMA, AR and MA models.

For *datetime_unit_of_measure* must be used str value, that defines unit of measure of time delta value (“second”, “minute”, “hour”, “day”, “month”, “year”)

For *period_value* must be used int value, that defines delta value.

If *period_value* is 11 and *datetime_unit_of_measure* is “minute”, it means delta value will be 11 minutes.

For the value *using_Kalman_Filter* must be used only True or False values.

Date value can be inputted in two formats:

- YYYY-MM-DD and time value will be automatically set to 00:00
- YYYY-MM-DD_HH:MM

For example:

/get/forecast/time_period?model_name=ARMA&value=Dew%20Point&station=LV01&steps=5&date_from=2020-02-10_12:00&date_till=2020-04-15_14:00&datetime_unit_of_measure=minute&period_value=11

/get/forecast/time_period?model_name=ARIMA&value=Dew%20Point&station=LV05&steps=5&date_from=2020-02-10&date_till=2020-04-15&datetime_unit_of_measure=minute&period_value=11&using_Kalman_Filter=True

V. */get/forecast/all_models*

Is used for getting forecast values using station data for all models (ARIMA, ARMA, AR, MA) using one-dimensional .

Essential argument: value, station, steps datetime_unit_of_measure, period_value.

Secondary arguments: using_Kalman_Filter.

For *value* argument can be used “Dew Point”, “Air Temperature” and etc.

Steps value means how many points ahead will be made forecast.

For *datetime_unit_of_measure* must be used str value, that defines unit of measure of time delta value (“second”, “minute”, “hour”, “day”, “month”, “year”)

For *period_value* must be used int value, that defines delta value.

If *period_value* is 11 and *datetime_unit_of_measure* is “minute”, it means delta value will be 11 minutes.

For the value *using_Kalman_Filter* must be used only True or False values.

For example:

/get/forecast/all_models?value=Dew%20Point&station=LV01&steps=5&datetime_unit_of_measure=minute&period_value=11

/get/forecast/all_models?value=Dew%20Point&station=LV01&steps=5&datetime_unit_of_measure=minute&period_value=11&using_Kalman_Filter=True