# Project Installation Guideline

*Author:* HoangBamberg@Gmail.Com

## 1. Project Description

### A. Purpose of the project

The Python GUI project was developed as the programming part for fulfilling the thesis (topic "*Evaluating the Impact of Sampling on Activity Monitoring in Dairy Cattle*") in MSc. Software Systems Science program at the University of Bamberg (https://www.uni-bamberg.de/en/ma-isosysc/). This is also a research part of the FutureIoT/Rindertracking project https://www.futureiot.de/portfolio/rindertracking/

The topic is about building and evaluating Activity recognition (https://en.wikipedia.org/wiki/Activity_recognition) (Machine Learning) models from triaxial sensors data (accelerometer and gyroscope carried by objects (cattle in this case) of the study). The project, however, can be used to build Machine Learning models in recognizing human activity as well.
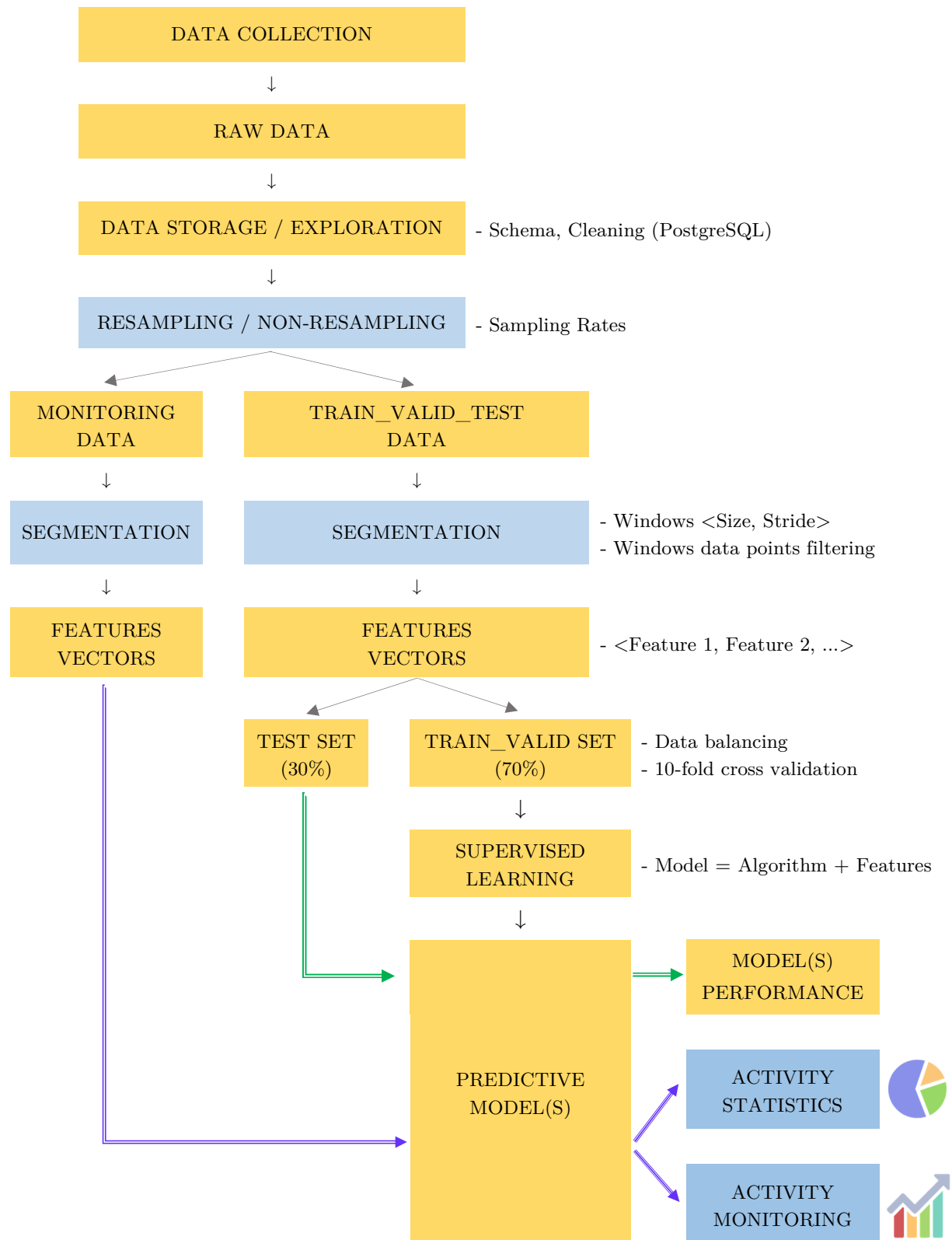
### B. Development Language and Environment

The project is developed in Python 3.8 (64 bit) using PyCharm IDE. The required packages (dependencies) are defined in the file **requirements.txt** (located in the root folder of the Repository). In case the project is loaded by PyCharm IDE, then it is automatically detected, and a pop-up will ask for the installation (see https://www.jetbrains.com/help/pycharm/managing-dependencies.html).
Alternatively, this can be done by manual command python -m pip install -r /requirements.txt.

### C. The project workflow

After researching the previous studies, the author proposed an implemented workflow as shown in **Figure 1**.

```
┌─────────────────────────┐
│    DATA COLLECTION      │
└─────────────────────────┘
             ↓
┌─────────────────────────┐
│       RAW DATA          │
└─────────────────────────┘
             ↓
┌─────────────────────────┐
│ DATA STORAGE / EXPLORATION │   - Schema, Cleaning (PostgreSQL)
└─────────────────────────┘
             ↓
┌─────────────────────────┐
│ RESAMPLING / NON-RESAMPLING │   - Sampling Rates
└─────────────────────────┘
       ↙             ↘
┌──────────────┐  ┌──────────────────┐
│  MONITORING  │  │ TRAIN_VALID_TEST │
│     DATA     │  │       DATA       │
└──────────────┘  └──────────────────┘
       ↓                  ↓
┌──────────────┐  ┌──────────────────┐
│ SEGMENTATION │  │   SEGMENTATION   │   - Windows <Size, Stride>
└──────────────┘  └──────────────────┘   - Windows data points filtering
       ↓                  ↓
┌──────────────┐  ┌──────────────────┐
│   FEATURES   │  │     FEATURES     │   - <Feature 1, Feature 2, ...>
│   VECTORS    │  │     VECTORS      │
└──────────────┘  └──────────────────┘
                      ↙        ↘
              ┌───────────┐ ┌──────────────────┐
              │ TEST SET  │ │ TRAIN_VALID SET  │   - Data balancing
              │  (30%)    │ │      (70%)       │   - 10-fold cross validation
              └───────────┘ └──────────────────┘
                                   ↓
                            ┌──────────────────┐
                            │    SUPERVISED    │   - Model = Algorithm + Features
                            │     LEARNING     │
                            └──────────────────┘
                                   ↓
                            ┌──────────────────┐  →  ┌──────────────────┐
                            │                  │     │     MODEL(S)     │
                            │                  │     │   PERFORMANCE    │
                            │   PREDICTIVE     │     └──────────────────┘
                            │    MODEL(S)      │  →  ┌──────────────────┐
                            │                  │     │    ACTIVITY      │
                            │                  │     │   STATISTICS     │
                            │                  │     └──────────────────┘
                            │                  │  →  ┌──────────────────┐
                            │                  │     │    ACTIVITY      │
                            │                  │     │   MONITORING     │
                            └──────────────────┘     └──────────────────┘
```

Figure 1: Implemented workflow

2

## D. Database Management System

The data for the application come from the device Bosch BNO055 which contains two sensors (accelerometer and gyroscope) originally sampled at 10Hz.

The database management system used in the project is PostgreSQL version 9.5.17. The input data table has the structure described in **Table 1.**

The two data tables **Train_Valid_Test** and **Monitoring** must contain <u>at least</u> 9 columns (i.e., *cattle_id, label, gx, gy, gz, ax, ay, az, timestamp*) with data types as follow:

| Column name | Data type | Explanation |
| --- | --- | --- |
| cattle_id | text | ID number of the cow (e.g., Lilith/*Hanna*) |
| label | text | Labelled activity of cows (e.g., *Liegen*) |
| gx | double precision | Gyroscope x axis signal value |
| gy | double precision | Gyroscope y axis signal value |
| gz | double precision | Gyroscope z axis signal value |
| ax | double precision | Accelerometer x axis signal value |
| ay | double precision | Accelerometer y axis signal value |
| az | double precision | Accelerometer z axis signal value |
| timestamp | bigint | Timestamp in Unix Epoch time format |

**Table 1: Training / Monitoring data table structure**

## E. Repository Structure

| Folder/File | Purpose |
| --- | --- |
| app.ini | Stores setting for GUI building |
| car_model_building.py | Main Python module |
| csv_out | Folder containing csv output files for each run |
| db_credentials.ini | Stores credentials and tables' names for database connection |
| features.py | Contains functions for calculating features |
| requirements.txt | Stores versions' names of the required Python packages |
| sample_datasets | Folder contain sample data for some cattle |
| user_defined_functions.py | Contains other user-defined functions |

**Table 2: Repository Structure**

## 2. Step-by-step running



**Figure 2: Database connection screen**

– **Ini Filepath** section stores the path to the ini file that contains the database credentials (e.g., the template is the file **db_credentials.ini** shown in **Table 2**). In the very first run if it shows a message "*Could not find a proper credentials ini file for db connection*" then please update the file **db_credentials.ini** appropriately. The fields **traintable** and **monitortable** in this ini file should refer to the tables that exist in the PosgreSQL database system.

– **Host, Port, Database, User, Password** stores database connection and user credentials which are loaded (*updated*) from (*to*) ini file in **Ini Filepath** section.

– **Train_Valid_Test table** indicates the name of table which is used for the <u>train/valid and test</u> phrase. This table structure meets the requirements in **Table 1**.

– **Statistics_Monitoring table** indicates the name of table which is used for the <u>Statistics</u>/<u>Monitoring</u> phrase. This table structure meets the requirements in **Table 1**. If the checkbox is unchecked, then it only runs training/validation and testing phrase.

– **Table to store result** stores the result for each experiment, it is automatically created if not exists on the database system. Every derived model will be saved into this table as a new record. The structure of this table is described in **Structure of the result table**.
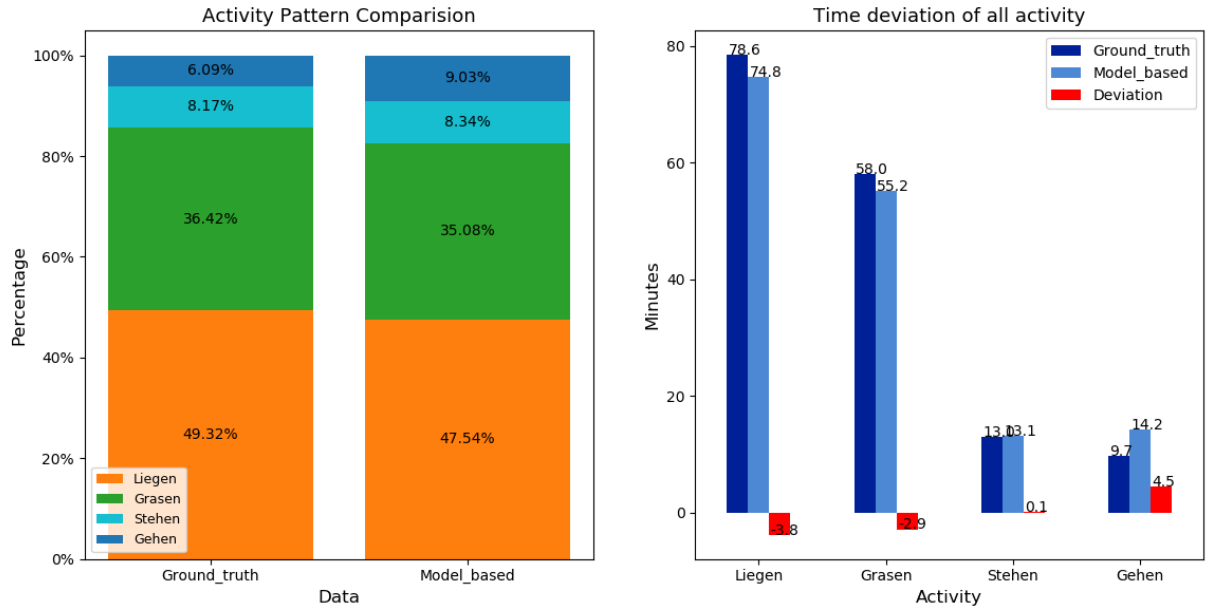
**Figure 3: Model building screen/tab**

– **Select labels for classifying**: Select **All labels** to train the model with all labels existing in train table, otherwise it will train the model with selected labels.

– **Resampling data**: User can either choose **Keep original data** or resample the data with given sampling rates.

– **Window setting**: With the setting like in the <span style="color:blue">**Figure 3**</span>, it will run with the window sizes of 5, 7 and 9 seconds. The **Train (Test) stride** is set to 100% meaning all data in **Train_Valid_Test table** will be used for Train/Valid and Test phrase. The **Monitoring stride** option is set to 200% meaning only half of data in **Statistics_Monitoring table** will be used for Statistics/Monitoring metrics calculations.

– **Select functions** and **Select axes to be applied** help to choose features for the classifier(s) selected in the **Classifiers** section underneath. The number of features = number of functions * number of axes.

– After selecting **classifier(s)** and **Kfold** option, just click <mark>Model fitting</mark> button for running **Train_Valid_Test** phrase, the result will display on PyCharm console, log into *"Train_Test_Result.txt"* file in **csv_out** folder and update to Database.

– **Monitoring setting**: After fitting the model, user can view the Statistics and Monitoring metrics with the two buttons <mark>Statistics</mark> and <mark>Monitoring</mark>. These two buttons will be disabled if the checkbox in the <span style="color:blue">**Figure 2**</span> is unchecked.

– Because the statistics and monitoring metrics are generated under each classifier at a specific sampling rate and window size, in the **Monitoring setting** section user needs to choose these three options for viewing these metrics. <span style="color:red">It is important to select these options that fall into one combination of sampling rates and window sizes (types) selected previously right above.</span> To be more clearly, settings in **Figure 3** enables user to train the model with sampling rates of 1, 3, 5, 7, 9 (Hz) in combination with window sizes of 5, 7, 9 (seconds) under 4 classifiers (5 * 3 * 4 = 60 combinations). As a result, in the Monitoring setting section, if user selects Sampling rate of 10Hz

(which is not in the 60 combinations) then it does not show the statistics/monitoring metrics.

– The Statics function helps to see how well a model can predict unseen data. More clearly it will show a comparison between ground-truth data and the predicted data generated by the model. An example of this function is shown below



**Figure 4: Activity statistics comparison between ground truth and model-based prediction**

## 3. Other important notes

### A. How to train/valid/test on multiple (mixed) cows' data and testing on another unseen cow?

To train on cow 1,2,3 and test on cow 4, just create the **Train_Valid_Test table** containing data from cow 1,2,3 then create the **Statistics_Monitoring table** containing sensor data of cow 4. Data for each cow in the **Train_Valid_Test table** distinguished by **cattle_id** field**.** The structure of these tables must meet the requirements in **Table 1**.

The result will be shown in PyCharm console, updated into database and logged into *"Train_Test_Result.txt"* file located in **csv_out** folder as described in **(C)**.

**B.** *How to carry Binary classification (E.g.: Lying and Non-lying)*

The application is originally developed for multiclass classification, so in order to carry a binary classification as between a main label (**Liegen)** and some other **Non-main labels (Gehen, Grasen, Stehen),** the following requirements should be met

+ The data in both **Train__Valid__Test table** (and **Statistics__Monitoring table**) must contain all four labels **Gehen Grasen Liegen Stehen**

+ In the **app.ini** file/section [GENERAL SETTINGS], the following variables must be set like below:

> binarymode = 1
> mainlabel = Liegen
> sublabels = Gehen__Grasen__Stehen

Note: the sublabels should be separated by the underline (__) character. This setting can be used to build binary model to predict Rumination and Non-Rumination (Grasen, Stehen…) The number of sub labels could be 1,2,3…. The tool will automatically detect them.

**C.** *CSV log files*

With a specific sampling rate and window size (type), a text file called "Train__Test__Result.txt" will be created to store the experiment results for that configuration. This txt file is created in a sub folder of **csv__out** folder as shown **Figure 5**. This sub folder is created in every run (every click on Model fitting button).

Additionally, for each configuration (e.g., sampling rate, window size/stride) the data set regarding **train__valid__test** and **monitoring** can be logged into .csv file for checking of correctness of features calculation. To enable .csv files saving, in the **app.ini** file/section [GENERAL SETTINGS], just set the variable **csvsaving** to 1 (Or *csvsaving = 1*)
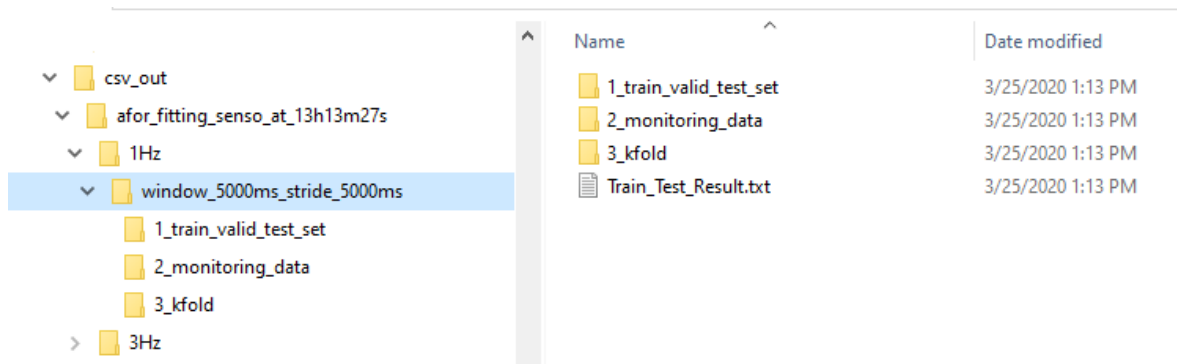


**Figure 5: csv__out folder structure**

## 4. Structure of the result table on the database

| Column name | Data Type | Description |
|---|---|---|
| model_title | text | This column is created for the purpose of showing the model to the end-user. The default value is an empty string, when it is given a string then it will be shown to the end-user. |
| model_init_name | text | The initialised name of the model, after a model is trained then its name is set at the default format: [date_of_creation]_[hhmmss]_username_[Binary/Multi]_[name_of_algorithm] An example is "20200520_005954_thaihp_Binary_RandomForest" |
| model_binary_content | bytea | The content of the model in binary format |
| model_comments | text | User comments of the model |
| train_table | text | The name of table for training |
| monitor_table | text | Table for the Staticstic_Monitoring metrics |
| no_of_predicted_classes | integer | Number of classes to be classified |
| list_of_predicted_classes | text | List of classes to be classified |
| original_sample_rate_in_hz | integer | The original sampling rate of the training data (train_table) |
| no_of_original_train_data_points | integer | Number of data points in the train_table |
| resampled_rate_in_hz | integer | resampling rate (if user chooses resample data in Figure 3 of readme.pdf file) |
| no_of_resampled_train_data_points | integer | The number of data points of training data after resampling |
| no_of_instances_for_each_class_in_resampled_train_table | integer | Number of instances for each class in the (resampled) training data |
| algorithm | text | The classifier selected |
| no_of_functions | integer | Number of functions |

| list_of_functions | text | List of functions |
|---|---|---|
| no_of_axes | integer | Number of axes selected for the training |
| list_of_axes | text | List of axes selected for the training |
| window_size | integer | The window size (in milliseconds) in of **training/monitoring** phrase |
| window_stride | text | The window stride for **training** phrase **only** |
| k_fold | integer | K fold |
| accuracy_train_valid | real | Accuracy of train_valid phrase |
| precision_train_valid | real | Precision of train_valid phrase |
| recall_train_valid | real | Recall of train_valid phrase |
| specificity_train_valid | real | Specificity of train_valid phrase |
| f1_train_valid | real | F1 score of train_valid phrase |
| accuracy_test | real | Accuracy on Test set (30% of the training data) |
| precision_test | real | Precision on Test set |
| recall_test | real | Recall on Test set |
| specificity_test | real | Specificity of Test set |
| f1_test | real | F1 score of Test set |
| monitoring_window_stride | text | The window stride for monitoring phrase |
| accuracy_monitor | real | Accuracy on Monitoring data |
| precision_monitor | real | Precision on Monitoring data |
| recall_monitor | real | Recall on Monitoring data |
| specificity_monitor | real | Specificity on Monitoring data |
| f1_monitor | real | F1 score on Monitoring data |
| start_time | timestamp | Starting time of the run |
| end_time | timestamp | Ending time of the run |
| running_time_in_minutes | text | The duration of the experiment |