

# README : Searching for the Boson

## 1 Prerequisites

The code is written in Python3, and makes extensive use of Numpy, the fundamental package for scientific computing in Python. Another external module used is `csv` - it intervenes in importing the flow of data to and from the Python environment, in order to read the data from a `.csv` file and to write the predictions in a `.csv` file respectively. Therefore the user must have both packages to run this code. For installation information, please refer to the Python documentation.

## 2 Executing the code

Extract the `.zip` and open the terminal at the location of the extracted files. Then enter: `python.run.py`. The file `run.py` contains an executable Python script that takes train data to compute a model. To see more details about the objective of this script, one should read our report. The desired train and test sets are to be put in a `train.csv` and `test.csv`, which are provided as placeholders.

## 3 Description of the code

The script first loads the data and initializes the hyperparameters. The values provided initially coincide with the optimal values mentioned in the report.

First missing values are treated by deleting the columns where they occur. Then the script splits the training and test data into 3 categories according to \*\*\*\*\* (see justification in the report). Then, a polynomial extension up to the value of "degree" follows. Finally, pairwise products of some of the features are performed. After this treatment, data is normalized and put in an appropriate format for the learning process.

The initial weights  $w$  are random. The training is done using a regularized logistic regression, with parameter  $\lambda$ . The optimal  $w$  is approached via an iterative method: the gradient descent (see justification in report). The method outputs the optimal weights  $w$  and the loss that corresponds to it.

Then, the test data goes through the treatment described above: from replacing missing values, to standardization. The predictions are performed then converted to the appropriate format  $\{-1, 1\}$  according to the best threshold (found via Kaggle). This process is repeated for each category of \*\*\*\*\* values (0, 1, 2, 3). The predictions are then rearranged to match the order of the original sample and put out using the `csv` module. They can be found in the same folder as the `run.py` script.

## 4 Implementations.py

This file contains the methods that were to be implemented. All of them return the predicted weights  $w$  and the loss. More precisely :

- i) `least_squares_GD(y, tx, initial_w, max_iters, gamma)` : takes as input `y` which is the vector of flags ( $-1$  or  $1$ ); `tx` which is the matrix of features (in an order corresponding to `y` )
- ii) `least_squares_SGD(y, tx, initial_w, batch_size, max_iters, gamma)`
- iii) `least_squares(y, tx)`
- iv) `ridge_regression(y, tx, lambda_)`
- v) `logistic_regression(y, tx, initial_w, max_iters, gamma)`
- vi) `reg_logistic_regression(y, tx, initial_w, max_iters, gamma, lambda_)`