

Programming Assignment 1

Write a program (you have to do that individually and submit individually! **NO** group submissions!) that implements a (batch) linear regression using the gradient descent method in *Python 3*. Use the following gradient calculation:

$$gradient = \sum_{i=1}^N \vec{x}_i (y_i - f(\vec{x}_i))$$

$$\vec{w} \leftarrow \vec{w} + \eta \cdot gradient$$

where \vec{x}_i is one data point (with N being the size of the data set), η the learning rate, y_i is the target output and $f(\vec{x}_i)$ is the linear function defined as $f(\vec{x}) = \vec{w}^T \vec{x}$ or equivalently $f(\vec{x}) = \sum_i w_i \cdot x_i$. Whereas \vec{w} and \vec{x} include the bias/intercept, i.e. w_0 (x_0 is always 1). All weights should be initialized as 0.

Given are the two random example data sets (uploaded in Moodle) named *random3* and *random4* as csv files. For random3 you have given the solution as well. Your task is to correctly implement the gradient descent method and return for each iteration the weights and sum of squared errors until a given threshold of **change** in the error is reached¹. The output of your algorithm should be printed onto the console/terminal and should look like this.

```
iteration_number,weight0,weight1,weight2,...,weightN,sum_of_squared_errors
```

Please do **NOT** print any extra information onto the output. The solution (It is rounded for readability to 9 decimals. You do NOT have to round!) for the data set is given with a learning rate of 0.00005 and a threshold of 0.0001. With that, you can check the correctness of your solution. Please be reminded, that small rounding errors are normal and will be treated as correct. If the program fails or the data format is incorrect you will get zero points.

Your program **must be** named `student.py` and **must** accept the following parameters:

1. **data** - The location of the data file (e.g. `/media/data/yacht.csv`).
2. **eta** - The learning rate of the gradient descent approach.
3. **threshold** - The threshold, that the change in error has to fall below, before the algorithm terminates.

The server will start your program in the following way:

```
python3 student.py --data random3.csv --eta 0.00005 --threshold 0.0001
```

The final program code must be uploaded to Moodle (in the respective VPL assignment) until Monday, the 13th of November 2023, 8:00am. The code will be automatically checked against randomly created data sets. You will get at most two points, if the output of your regressor is correct.

2 points

¹Meaning $e_t - e_{t+1} < \text{threshold}$