

HW5: Linear Regression

[Submit Assignment](#)

Due Tuesday by 2:35pm **Points** 100 **Submitting** a file upload **File Types** zip
Available Mar 2 at 12am - Mar 9 at 2:35pm 8 days

Assignment Goals

- Implement a linear regression calculation
- Examine the trends in real (messy) data

Summary

Percentage of body fat, age, weight, height, and ten body circumference measurements (e.g., abdomen) are recorded for 252 men. Body fat, one measure of health, has been accurately estimated by an underwater weighing technique. Fitting body fat to the other measurements using multiple regression provides a convenient way of estimating body fat for men using only a scale and a measuring tape. In this assignment, you will be looking at the [bodyfat dataset](http://jse.amstat.org/v4n1/datasets.johnson.html) (<http://jse.amstat.org/v4n1/datasets.johnson.html>) and build several models on top of it.

Program Specification

You will be using the bodyfat dataset ([bodyfat.csv](#)) for this assignment. Complete the following Python functions in this template [regression.py](#):

1. **get_dataset(filename)** — takes a filename and **returns** the data as described below in an n-by-(m+1) array
2. **print_stats(dataset, col)** — takes the dataset as produced by the previous function and **prints** several statistics about a column of the dataset; does not return anything
3. **regression(dataset, cols, betas)** — calculates and **returns** the mean squared error on the dataset given fixed betas
4. **gradient_descent(dataset, cols, betas)** — performs a single step of gradient descent on the MSE and **returns** the derivative values as an 1D array
5. **iterate_gradient(dataset, cols, betas, T, eta)** — performs T iterations of gradient descent starting at the given betas and **prints** the results; does not return anything
6. **compute_betas(dataset, cols)** — using the closed-form solution, calculates and **returns** the values of betas and the corresponding MSE as a tuple
7. **predict(dataset, cols, features)** — using the closed-form solution betas, **return** the predicted body fat percentage of the give features.

8. **synthetic_datasets(betas, alphas, X, sigma)** — generates two synthetic datasets, one using a linear model and the other using a quadratic model.
9. **plot_mse()** — fits the synthetic datasets, and **plots** a figure depicting the MSEs under different situations.

Get Dataset

The `get_dataset()` function should return an n -by- $(m+1)$ array of data, where n is the number of data points, and m is the number of features, plus an additional column of labels. The first column should be bodyfat percentage, which is the **target** that our regression model aims for. We denote it as **y** in the rest of the write up. Starting from the second column, there goes a list of **features**, including density, age, weight, and more. We use index 1 to represent density, 2 to represent age, and so on... You should ignore the "IDNO" column as it merely represents the individual id of each participant.

```
>>> get_dataset('bodyfat.csv')
=> array([[12.6 , 1.0708, 23. , ..., 32. , 27.4 , 17.1 ],
        [ 6.9 , 1.0853, 22. , ..., 30.5 , 28.9 , 18.2 ],
        [24.6 , 1.0414, 22. , ..., 28.8 , 25.2 , 16.6 ],
        ...,
        [28.3 , 1.0328, 72. , ..., 31.3 , 27.2 , 18. ],
        [25.3 , 1.0399, 72. , ..., 30.5 , 29.4 , 19.8 ],
        [30.7 , 1.0271, 74. , ..., 33.7 , 30. , 20.9 ]])
>>> dataset = get_dataset('bodyfat.csv')
>>> dataset.shape
(252, 16)
```

Dataset Statistics

This is just a quick summary function on one **feature**, given in the parameter `col`, from the above dataset. When called, you should **print**:

1. the number of data points
2. the sample mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
3. the sample standard deviation $\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

on three lines. Please format your output to include only **TWO digits** after the decimal point. For example:

```
>>> data = get_dataset('bodyfat.csv')
>>> print_stats(dataset, 1) # summary of density
252
1.06
0.02
```

You might find [this guide](https://pyformat.info/) [\(https://pyformat.info/\)](https://pyformat.info/) to python print formatting useful.

Note: Please implement the formula yourself, instead of using `np.mean()` or `np.std()`.

Linear Regression

This function will perform linear regression with the model

$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_m x_m$$

We first define the mean squared error (MSE) as the sum of squared errors divided by # data points:

$$MSE(\beta_0, \beta_1, \dots, \beta_m) = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_m x_{im} - y_i)^2$$

The first argument refers to the dataset, and the second argument is a list of features that we wish to learn on. For example, if we would like to study the relationship of body fat vs age and weight, cols should be [2, 3]. In this case, the model should be $f(x) = \beta_0 + \beta_2 x_2 + \beta_3 x_3$. The last argument (betas) of this function represents three betas, $\beta_0, \beta_2, \beta_3$. **Return** the corresponding MSE as calculated on your dataset.

```
>>> regression(dataset, cols=[2,3], betas=[0,0,0])
=> 418.50384920634923
>>> regression(dataset, cols=[2,3,4], betas=[0,-1.1,-.2,3])
=> 11859.17408611111
```

Gradient Descent

This function will perform gradient descent on the MSE. At the current parameter $(\beta_0, \beta_1, \dots, \beta_m)$, the gradient is defined by the vector of partial derivatives:

$$\begin{aligned} \frac{\partial MSE(\beta_0, \beta_1, \dots, \beta_m)}{\partial \beta_0} &= \frac{2}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_m x_{im} - y_i) \\ \frac{\partial MSE(\beta_0, \beta_1, \dots, \beta_m)}{\partial \beta_1} &= \frac{2}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_m x_{im} - y_i) x_{i1} \\ &\dots \end{aligned}$$

This function **returns** the corresponding gradient as a 1-D numpy array with the partial derivative with respect to β_0 as the first value.

```
>>> gradient_descent(dataset, cols=[2,3], betas=[0,0,0])
=> array([-37.87698413, -1756.37222222, -7055.35138889]) # order: [partial derivative of beta_0, beta_2, beta_3]
```

Iterate Gradient

Gradient descent starts from initial parameter $(\beta_0^{(0)}, \beta_1^{(0)}, \dots, \beta_m^{(0)})$ and iterates the following updates at time $t = 1, 2, \dots, T$:

$$\beta_0^{(t)} = \beta_0^{(t-1)} - \eta \frac{\partial \text{MSE}(\beta_0^{(t-1)}, \beta_1^{(t-1)}, \dots, \beta_m^{(t-1)})}{\partial \beta_0}$$

$$\beta_1^{(t)} = \beta_1^{(t-1)} - \eta \frac{\partial \text{MSE}(\beta_0^{(t-1)}, \beta_1^{(t-1)}, \dots, \beta_m^{(t-1)})}{\partial \beta_1}$$

and so on for the rest.

The parameters to this function are the dataset and a selection of features, the T number of iterations to perform, and η (eta), the parameter for the above calculations. Begin from the initial value as specified in the parameter betas.

Print the following for each iteration on one line, separated by spaces:

1. the current iteration number beginning at 1 and ending at T
2. the current MSE
3. the current value of beta_0
4. the current value of other betas

As before, all floating point values should be **rounded to two digits** for output.

```
>>> iterate_gradient(dataset, cols=[1,8], betas=[400,-400,300], T=10, eta=1e-4)
1 423085332.40 394.45 -405.84 -220.18 # order: T, mse, beta0, beta1, beta8
2 229744495.73 398.54 -401.54 163.14
3 124756241.68 395.53 -404.71 -119.33
4 67745350.04 397.75 -402.37 88.82
5 36787203.39 396.11 -404.09 -64.57
6 19976260.50 397.32 -402.82 48.47
7 10847555.07 396.43 -403.76 -34.83
8 5890470.68 397.09 -403.07 26.55
9 3198666.69 396.60 -403.58 -18.68
10 1736958.93 396.96 -403.20 14.65
```

Try different values for eta and a much larger T, and see how small you can make MSE (optional).

Compute Betas

Instead of using gradient descent, we can compute the closed-form solution for the parameters directly. For ordinary least-squares, this is

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}$$

This function returns the calculated betas (selected by cols) and their corresponding MSE in a **tuple**, as (MSE, beta_0, beta_1, and so on).

```
>>> compute_betas(dataset, cols=[1,2])
=> (1.4029395600144443, 441.3525943592249, -400.5954953685588, 0.009892204826346139)
```

Predict Body Fat

Using your closed-form betas, predict the body fat percentage for a given number of features. Return that value.

For example:

```
>>> predict(dataset, cols=[1,2], features=[1.0708, 23])
=> 12.62245862957813
```

Synthetic Datasets

Now, let us create a few new datasets to study some characteristics of Linear Regression models. The function `synthetic_datasets(betas, alphas, X, sigma)` should return two synthetic datasets, a linear one and a quadratic one.

The linear dataset is defined as an $n \times 2$ array, where n is the size of the input array of data points X . You may assume that X is a 2D array with shape $(n,1)$. The second column of the linear dataset should be a copy of X . The first column is defined as follows:

$$y_i = \beta_0 + \beta_1 x_i + z_i, z_i \sim N(0, \sigma)$$

β_0 and β_1 are provided by the argument `betas`. z_i is an artificial error term to 'jitter' the labels so that they don't look too perfect. It is used to create a noise in the label. You should sample z 's from a normal distribution with mean 0 and standard deviation `sigma`, which is provided in the argument, too. You may use `numpy.random.normal()` detailed [here](https://numpy.org/doc/stable/reference/random/generated/numpy.random.normal.html)

(<https://numpy.org/doc/stable/reference/random/generated/numpy.random.normal.html>).

Similarly, the quadratic dataset is also defined as an $n \times 2$ array. The second column remains the same as above, while the first column is defined as follows:

$$y_i = \alpha_0 + \alpha_1 x_i^2 + z_i, z_i \sim N(0, \sigma)$$

The alphas and z 's are from the argument `alphas`, and a normal distribution with mean 0 and standard deviation `sigma`, respectively.

The function should return a tuple, with the first being the linear dataset array, followed by the quadratic dataset array.

```
>>> synthetic_datasets(np.array([0,2]), np.array([0,1]), np.array([[4]]), 1)
=> (array([[8.65003702, 4.          ]]), array([[15.5334939, 4.          ]]))
```

Compare and Plot

Using the two synthetic datasets, let us study the performance of linear regression under linear data and quadratic data respectively. In `plot_mse()`, do the following:

1. Create an input array X containing 1000 numbers within range [-100, 100].
2. Create couples of betas and alphas with non-zero values.
3. Set sigmas to be $10^{-4}, 10^{-3}, \dots, 1, 10, \dots, 10^5$.
4. Under each settings of sigmas, generate two synthetic datasets.
5. Fit both datasets using `compute_betas()`, obtain the corresponding MSEs.
6. Plot a figure showing how MSE changes while sigma is increasing:
 - You should use the plotting format `-o` (line with circle markers)
 - The x-axis should be different settings of sigma.
 - The y-axis should be the MSEs calculated from the linear and quadratic datasets.
 - Make both axes **log scale**.
 - Label your x- and y-axes.
 - Make a legend.
 - Save the figure as `mse.pdf`. Do NOT display it (i.e., no `plt.show()`).

Hint: You may expect your graph to contain two nearly straight lines (with a little curve).

Note that in order to run your code on CSL machines, you should invoke your program with an additional keyword `cs1`:

```
$ python3 regression.py cs1
```

Discovery

Experiment on different combinations of features, and learning methods (closed form solution, gradient descent). Try to answer these questions:

- What is the best feature that predicts well?
- What are the best set of features that predict well?
- Which learning method is the most effective in terms of training error (MSE)?
- Which learning method is the most efficient (takes the least time)?
- What will happen if our dataset has more than 1 million entries?

Observe the figure you plotted, and try to answer these questions:

- What do you see happening in the linear case as the noise variance goes up?
- What do you see happening in the quadratic case as the noise variance goes up?
- What happens if the noise is extremely small?
- What happens if the noise is extremely large?

We won't grade these questions, but feel free to share your thoughts on Piazza!

Submission

Please submit your code zipped in a file called `hw5_<netid>.zip`. Inside your zip file, there should be **only** one file named: `regression.py`. Do NOT submit a Jupyter notebook .ipynb file.

Be sure to **remove all debugging output** before submission. Failure to remove debugging output will be **penalized (10pts)**.

Cheating results in -100 pts and further punishment

This assignment due at 3/9/2021 2:30pm. Submitting right at 2:30pm will result in a late submission. It is preferable to first submit a version well before the deadline (at least one hour before) and check the content/format of the submission to make sure it's the right version. Then, later update the submission until the deadline if needed.