

## Regression Week 5: LASSO Assignment 2

In this notebook, you will implement your very own LASSO solver via coordinate descent. You will:

- Write a function to normalize features
- Implement coordinate descent for LASSO
- Explore effects of L1 penalty

## If you are doing the assignment with IPython Notebook

An IPython Notebook has been provided below to you for this quiz. This notebook contains the instructions, quiz questions and partially-completed code for you to use as well as some cells to test your code.

## What you need to download

If you are using GraphLab Create

- Download the King County House Sales data In SFrame format: [kc\\_house\\_data.gl.zip](#)
- Download the companion IPython Notebook: [week-5-lasso-assignment-2-blank.ipynb](#)
- Save both of these files in the same directory (where you are calling IPython notebook from) and unzip the data file.

If you are not using GraphLab Create

- Download the King County House Sales data csv file: [kc\\_house\\_data.csv](#)
  - Download the King County House Sales training data csv file: [kc\\_house\\_train\\_data.csv](#)
  - Download the King County House Sales testing data csv file: [kc\\_house\\_test\\_data.csv](#)
- **IMPORTANT: use the following types for columns when importing the csv files. Otherwise, they may not be imported correctly: [str, str, float, float, float, float, float, int, str, int, int, int, int, int, int, int, str, float, float, float, float]. If your tool of choice requires a dictionary of types for importing csv files (e.g. Pandas), use:**

```
dtype_dict = {'bathrooms':float, 'waterfront':int, 'sqft_above':int, 'sqft_living15':float, 'grade':int, 'yr_renovated':int, 'price':float, 'bedrooms':float, 'zipcode':str, 'long':float, 'sqft_lot15':float, 'sqft_living':float, 'floors':str, 'condition':int, 'lat':float, 'date':str, 'sqft_basement':int, 'yr_built':int, 'id':str, 'sqft_lot':int, 'view':int}
```

## Useful resources

You may need to install the software tools or use the free Amazon EC2 machine. Instructions for both options are provided in the reading for Module 1 (Simple Regression).

If you are following the IPython Notebook and/or are new to numpy then you might find the following tutorial helpful: [numpy-tutorial.ipynb](#)

## If you are using GraphLab Create and the companion IPython Notebook

Open the companion IPython notebook and follow the instructions in the notebook.

## If instead you are using other tools to do your homework

You are welcome to write your own code and use any other libraries, like Pandas or R, to help you in the process. If you would like to take this path, follow the instructions below.

1. If you're using SFrame, import GraphLab Create and load in the house data as follows:

```
sales = graphlab.SFrame('kc_house_data.gl/')
```

Otherwise, load all the three csv files.

2. If you're using Python: to do the matrix operations required to perform a coordinate descent we will be using the popular python library 'numpy' which is a computational library specialized for operations on arrays. For students unfamiliar with numpy we have created a numpy tutorial (see useful resources). It is common to import numpy under the name 'np' for short, to do this execute:

```
import numpy as np
```

3. Next, from Module 2 (Multiple Regression), copy and paste the 'get\_numpy\_data' function (or equivalent) that takes a data set, a list of features (e.g. ['sqft\_living', 'bedrooms']), to be used as inputs, and a name of the output (e.g. 'price'). This function returns a 'feature\_matrix' (2D array) consisting of first a column of ones followed by columns containing the values of the input features in the data set in the same order as the input list. It also returns an 'output\_array' which is an array of the values of the output in the data set (e.g. 'price').

e.g. in Python:

```
def get_numpy_data(data_sframe, features, output):
    ...
    return (feature_matrix, output_array)
```

4. Similarly, copy and paste the 'predict\_output' function (or equivalent) from Module 2. This function accepts a 2D array 'feature\_matrix' and a 1D array 'weights' and return a 1D array 'predictions'.

e.g. in Python:

```
def predict_output(feature_matrix, weights):
    ...
    return predictions
```

5. In the house dataset, features vary wildly in their relative magnitude: 'sqft\_living' is very large overall compared to 'bedrooms', for instance. As a result, weight for 'sqft\_living' would be much smaller than weight for 'bedrooms'. This is problematic because "small" weights are dropped first as l1\_penalty goes up.

To give equal considerations for all features, we need to normalize features as discussed in the lectures: we divide each feature by its 2-norm so that the transformed feature has norm 1.

6. Write a short function called 'normalize\_features(feature\_matrix)', which normalizes columns of a given feature matrix. The function should return a pair ('normalized\_features, norms'), where the second item contains the norms of original features. As discussed in the lectures, we will use these norms to normalize the test data in the same way as we normalized the training data.

e.g. in Python:

```
def normalize_features(features):
    ...
    return (normalized_features, norms)
```

Hint: If you are using Numpy, a handy shorthand for computing 2-norms of columns is

```
norms = np.linalg.norm(X, axis=0)
```

To normalize columns, perform element-wise division.

```
X_normalized = X / norms
```

## Review of Coordinate Descent

7. We seek to obtain a sparse set of weights by minimizing the LASSO cost function

```
SUM[ (prediction - output)^2 ] + lambda*( |w[1]| + ... + |w[k]| ).
```

(By convention, we do not include  $w[0]$  in the L1 penalty term. We never want to push the intercept to zero.)

The absolute value sign makes the cost function non-differentiable, so simple gradient descent is not viable (you would need to implement a method called subgradient descent). Instead, we will use coordinate descent: at each iteration, we will fix all weights but weight  $i$  and find the value of weight  $i$  that minimizes the objective. That is, we look for

```
argmin_{w[i]} [ SUM[ (prediction - output)^2 ] + lambda*( |w[1]| + ... + |w[k]| ) ]
```

where all weights other than  $w[i]$  are held to be constant. We will optimize one  $w[i]$  at a time, circling through the weights multiple times.

- Pick a coordinate  $i$
- Compute  $w[i]$  that minimizes the LASSO cost function
- Repeat the two steps for all coordinates, multiple times

8. For this assignment, we use cyclical coordinate descent with normalized features, where we cycle through coordinates 0 to  $(d-1)$  in order, and assume the features were normalized as discussed above. The formula for optimizing each coordinate is as follows:

```
w[i] = { (ro[i] + lambda/2)   if ro[i] < -lambda/2
        0                  if -lambda/2 <= ro[i] <= lambda/2
        (ro[i] - lambda/2)   if ro[i] > lambda/2
```

where

```
ro[i] = SUM[ [feature_i]*(output - prediction + w[i]*[feature_i]) ].
```

Note that we do not regularize the weight of the constant feature (intercept)  $w[0]$ , so, for this weight, the update is simply:

```
w[0] = ro[i]
```

## Effect of L1 penalty

9. Consider a simple model with 2 features: 'sqft\_living' and 'bedrooms'. The output is 'price'.

- First, run `get_numpy_data()` (or equivalent) to obtain a feature matrix with 3 columns (constant column added). Use the entire 'sales' dataset for now.
- Normalize columns of the feature matrix. Save the norms of original features as 'norms'.
- Set initial weights to [1,4,1].
- Make predictions with feature matrix and initial weights.
- Compute values of  $ro[i]$ , where

```
ro[i] = SUM[ [feature_i]*(output - prediction + w[i]*[feature_i]) ]
```

10. **Quiz Question: Recall that, whenever  $ro[i]$  falls between  $-l1\_penalty/2$  and  $l1\_penalty/2$ , the corresponding weight  $w[i]$  is sent to zero. Now suppose we were to take one step of coordinate descent on either feature 1 or feature 2. What range of values of  $l1\_penalty$  would not set  $w[1]$  zero, but would set  $w[2]$  to zero, if we were to take a step in that coordinate?**

11. **Quiz Question: What range of values of  $l1\_penalty$  would set both  $w[1]$  and  $w[2]$  to zero, if we were to take a step in that coordinate?**

So we can say that  $ro[i]$  quantifies the significance of the  $i$ -th feature: the larger  $ro[i]$  is, the more likely it is for the  $i$ -th feature to be retained.

## Single Coordinate Descent Step

12. Using the intercept (above), implement coordinate descent that minimizes the cost function over a single feature  $i$ . Note that the intercept (weight 0) is not regularized. The function should accept feature matrix, output, current weights, l1 penalty, and index of feature to optimize over. The function should return new weight for feature  $i$ .

e.g. in Python:

```
def lasso_coordinate_descent_step(i, feature_matrix, output, weights, l1_penalty):
    # compute prediction
    prediction = ...
    # compute ro[i] = SUM[ [feature_i]*(output - prediction + weight[i]*[feature_i]) ]
    ro_i = ...

    if i == 0: # intercept -- do not regularize
        new_weight_i = ro_i
    elif ro_i < -l1_penalty/2.:
        new_weight_i = ...
    elif ro_i > l1_penalty/2.:
        new_weight_i = ...
    else:
        new_weight_i = 0.

    return new_weight_i
```

If you are using Numpy, test your function with the following snippet:

```
# should print 0.425558846691
import math
print lasso_coordinate_descent_step(1, np.array([[3./math.sqrt(13),1./math.sqrt(10)],
                                                [2./math.sqrt(13),3./math.sqrt(10)]]), np.array([1., 1.]), np.array([1., 4.]), 0.1)
```

## Cyclical coordinate descent

13. Now that we have a function that optimizes the cost function over a single coordinate, let us implement cyclical coordinate descent where we optimize coordinates 0, 1, ...,  $(d-1)$  in order and repeat.

When do we know to stop? Each time we scan all the coordinates (features) once, we measure the change in weight for each coordinate. If no coordinate changes by more than a specified threshold, we stop.

For each iteration:

- As you loop over features in order and perform coordinate descent, measure how much each coordinate changes.
- After the loop, if the maximum change across all coordinates is falls below the tolerance, stop. Otherwise, go back to the previous step.

Return weights

The function should accept the following parameters:

- Feature matrix
- Output array
- Initial weights
- L1 penalty
- Tolerance

e.g. in Python:

```
def lasso_cyclical_coordinate_descent(feature_matrix, output, initial_weights, l1_penalty, to_lerance):
    ...
    return weights
```

14. Let us now go back to the simple model with 2 features: 'sqft\_living' and 'bedrooms'. Using 'get\_numpy\_data' (or equivalent), extract the feature matrix and the output array from from the house dataframe. Then normalize the feature matrix using 'normalize\_features()' function.

Using the following parameters, learn the weights on the sales dataset.

- Initial weights = all zeros
- L1 penalty =  $1e7$
- Tolerance = 1.0

15. **Quiz Question: What is the RSS of the learned model on the normalized dataset?**

16. **Quiz Question: Which features had weight zero at convergence?**

## Evaluating LASSO fit with more features

17. Let us split the sales dataset into training and test sets. If you are using GraphLab Create, call 'random\_split' with .8 ratio and seed=0. Otherwise, please down the corresponding csv files from the downloads section.

18. Create a normalized feature matrix from the TRAINING data with the following set of features.

- bedrooms, bathrooms, sqft\_living, sqft\_lot, floors, waterfront, view, condition, grade, sqft\_above, sqft\_basement, yr\_built, yr\_renovated

Make sure you store the norms for the normalization, since we'll use them later.

19. First, learn the weights with  $l1\_penalty=1e7$ , on the training data. Initialize weights to all zeros, and set the tolerance=1. Call resulting weights 'weights1e7', you will need them later.

20. **Quiz Question: What features had non-zero weight in this case?**

21. Next, learn the weights with  $l1\_penalty=1e8$ , on the training data. Initialize weights to all zeros, and set the tolerance=1. Call resulting weights 'weights1e8', you will need them later.

22. **Quiz Question: What features had non-zero weight in this case?**

23. Finally, learn the weights with  $l1\_penalty=1e4$ , on the training data. Initialize weights to all zeros, and set the tolerance=5e5. Call resulting weights 'weights1e4', you will need them later. (This case will take quite a bit longer to converge than the others above.)

24. **Quiz Question: What features had non-zero weight in this case?**

## Rescaling learned weights

25. Recall that we normalized our feature matrix, before learning the weights. To use these weights on a test set, we must normalize the test data in the same way. Alternatively, we can rescale the learned weights to include the normalization, so we never have to worry about normalizing the test data:

In this case, we must scale the resulting weights so that we can make predictions with original features:

- Store the norms of the original features to a vector called 'norms':

```
features, norms = normalize_features(features)
```

- Run Lasso on the normalized features and obtain a 'weights' vector
- Compute the weights for the original features by performing element-wise division, i.e.

```
weights_normalized = weights / norms
```

Now, we can apply `weights_normalized` to the test data, without normalizing it!

26. Create a normalized version of each of the weights learned above. ('weights1e4', 'weights1e7', 'weights1e8'). To check your results, if you call 'normalized\_weights1e7' the normalized version of 'weights1e7', then

```
print normalized_weights1e7[3]
```

should print 161.31745624837794.

## Evaluating each of the learned models on the test data

27. Let's now evaluate the three models on the test data. Extract the feature matrix and output array from the TEST set. But this time, do NOT normalize the feature matrix. Instead, use the normalized version of weights to make predictions.

Compute the RSS of each of the three normalized weights on the (unnormalized) feature matrix.

28. **Quiz Question: Which model performed best on the test data?**