≣ 课程

✓ 返回到第5周

马上开始 ②

## In this assignment, you will use LASSO to select features, building on a pre-implemented solver for LASSO (using

Regression Week 5: LASSO Assignment 1

GraphLab Create, though you can use other solvers). You will:

Choose best L1 penalty using a validation set.

Run LASSO with different L1 penalties.

- Choose best L1 penalty using a validation set, with additional constraint on the size of subset.
- In the second assignment, you will implement your own LASSO solver, using coordinate descent.

## Pandas). You are free to experiment with other tools (e.g. R or Matlab), but they may not produce correct

numbers for the quiz questions. • If you are using GraphLab Create, download the IPython notebook and follow the instructions contained in the notebook.

kc\_house\_data.gl.zip

Download the companion IPython Notebook:

```
week-5-lasso-assignment-1-blank.ipynb....
```

file.

• Save both of these files in the same directory (where you are calling IPython notebook from) and unzip the data

If you are using scikit-learn with Pandas: Download the King County House Sales data csv file:

wk3\_kc\_house\_train\_data.csv.zip

wk3\_kc\_house\_valid\_data.csv.zip

Download the King County House Sales testing data csv file:

Useful resources You may need to install the software tools or use the free Amazon EC2 machine. Instructions for both options are

```
numpy-tutorial.ipynb.zip
```

**0**. Load the sales dataset using Pandas:

import pandas as pd

The instructions may apply to other tools, but the set of parameters are specific to scikit-learn.

'sqft\_living':float, 'floors':float, 'condition':int, 'lat':float, 'date':str,

If you are using GraphLab Create and the companion IPython Notebook

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,

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

```
1. Create new features by performing following transformation on inputs:
  1 from math import log, sqrt
   2 sales['sqft_living_sqrt'] = sales['sqft_living'].apply(sqrt)
    sales['sqft_lot_sqrt'] = sales['sqft_lot'].apply(sqrt)
    sales['bedrooms_square'] = sales['bedrooms']*sales['bedrooms']
   5 sales['floors square'] = sales['floors']*sales['floors']
```

"normalize=True" when creating the Lasso object. Refer to the following code snippet for the list of features.

**Note.** From here on, the list 'all\_features' refers to the list defined in this snippet.

model\_all = linear\_model.Lasso(alpha=5e2, normalize=True) # set parameters

testing['bedrooms\_square'] = testing['bedrooms']\*testing['bedrooms']

training['sqft\_living\_sqrt'] = training['sqft\_living'].apply(sqrt)

training['floors\_square'] = training['floors']\*training['floors']

Compute the RSS on VALIDATION for the current model (print or save the RSS)

Report which L1 penalty produced the lower RSS on VALIDATION.

where 'model' is an instance of linear\_model.Lasso.

**10.** Assign 7 to the variable 'max\_nonzeros'.

training['bedrooms\_square'] = training['bedrooms']\*training['bedrooms']

testing['floors\_square'] = testing['floors']\*testing['floors']

training['sqft\_lot\_sqrt'] = training['sqft\_lot'].apply(sqrt)

model all.fit(sales[all features], sales['price']) # learn weights

from sklearn import linear\_model # using scikit-learn

'yr\_built', 'yr\_renovated']

'sqft\_basement',

all\_features = ['bedrooms', 'bedrooms\_square', 'bathrooms', 'sqft\_living', 'sqft\_living\_sqrt',

```
3. Quiz Question: Which features have been chosen by LASSO, i.e. which features were assigned nonzero
weights?
4. To find a good L1 penalty, we will explore multiple values using a validation set. Let us do three way split into
train, validation, and test sets. Download the provided csv files containing training, validation and test sets.
     testing = pd.read_csv('wk3_kc_house_test_data.csv', dtype=dtype_dict)
   2 training = pd.read_csv('wk3_kc_house_train_data.csv', dtype=dtype_dict)
   3 validation = pd.read csv('wk3 kc house valid data.csv', dtype=dtype dict)
Make sure to create the 4 features as we did in #1:
       testing['sqft_living_sqrt'] = testing['sqft_living'].apply(sqrt)
    2 testing['sqft_lot_sqrt'] = testing['sqft_lot'].apply(sqrt)
```

```
validation['sqft_living_sqrt'] = validation['sqft_living'].apply(sqrt)
   validation['sqft_lot_sqrt'] = validation['sqft_lot'].apply(sqrt)
   validation['bedrooms_square'] = validation['bedrooms']*validation['bedrooms']
   validation['floors square'] = validation['floors']*validation['floors']
num=13).)
• Learn a model on TRAINING data using the specified I1_penalty. Make sure to specify normalize=True in the
  constructor:
   1 model = linear_model.Lasso(alpha=l1_penalty, normalize=True)
```

nonzero coefficients first, and add 1 if the intercept is also nonzero. A succinct way to do this is 1 np.count\_nonzero(model.coef\_) + np.count\_nonzero(model.intercept\_)

8. Quiz Question: Using the best L1 penalty, how many nonzero weights do you have? Count the number of

You are going to implement a simple, two phase procedure to achieve this goal: to have the desired number of non-zero weights.

11. Exploring large range of l1\_penalty For I1\_penalty in np.logspace(1, 4, num=20):

1 model = linear\_model.Lasso(alpha=l1\_penalty, normalize=True)

rule of thumb" --- an interpretable model that has only a few features in them.

12. Out of this large range, we want to find the two ends of our desired narrow range of l1\_penalty. At one end, we will have I1\_penalty values that have too few non-zeros, and at the other end, we will have an I1\_penalty that has

later)

later) • The smallest I1\_penalty that has fewer non-zeros than 'max\_nonzeros' (if we pick a penalty larger than this value, we will definitely have too few non-zero weights)Store this value in the variable '11\_penalty\_max' (we will use it

appropriate boundaries. 13. Quiz Question: What values did you find for l1\_penalty\_min and l1\_penalty\_max?

VALIDATION set. For I1\_penalty in np.linspace(I1\_penalty\_min,I1\_penalty\_max,20):

Measure the RSS of the learned model on the VALIDATION set

and has sparsity equal to 'max\_nonzeros'? 16. Quiz Question: What features in this model have non-zero coefficients?

IMPORTANT: Choice of tools For the purpose of this assessment, you may choose between GraphLab Create and scikit-learn (with

• If you are using Pandas+scikit-learn combination, follow through the instructions in this reading.

What you need to download If you are using GraphLab Create:

Download the King County House Sales data In SFrame format:

kc\_house\_data.csv.zip

Download the King County House Sales training data csv file:

Download the King County House Sales validation data csv file:

```
wk3_kc_house_test_data.csv.zip
```

If you are following the IPython Notebook and/or are new to numpy then you might find the following tutorial helpful:

provided in the reading for Module 1 (Simple Regression).

'sqft\_basement':int, 'yr\_built':int, 'id':str, 'sqft\_lot':int, 'view':int} 5 sales = pd.read csv('kc house data.csv', dtype=dtype dict)

If you are using scikit-learn with Pandas:

```
• Squaring bedrooms will increase the separation between not many bedrooms (e.g. 1) and lots of bedrooms (e.g.
  4) since 1^2 = 1 but 4^2 = 16. Consequently this variable will mostly affect houses with many bedrooms.

    On the other hand, taking square root of sqft_living will decrease the separation between big house and small

  house. The owner may not be exactly twice as happy for getting a house that is twice as big.
2. Using the entire house dataset, learn regression weights using an L1 penalty of 5e2. Make sure to add
```

'sqft\_lot', 'sqft\_lot\_sqrt', 'floors', 'floors\_square', 'waterfront', 'view', 'condition', 'grade', 'sqft\_above',

```
5. Now for each I1_penalty in [10^1, 10^1.5, 10^2, 10^2.5, ..., 10^7] (to get this in Python, type np.logspace(1, 7,
```

6. Quiz Question: Which was the best value for the l1\_penalty, i.e. which value of l1\_penalty produced the lowest RSS on VALIDATION data? 7. Now that you have selected an L1 penalty, compute the RSS on TEST data for the model with the best L1 penalty.

• Explore a large range of 'l1\_penalty' values to find a narrow region of 'l1\_penalty' values where models are likely

• Further explore the narrow region you found to find a good value for 'l1\_penalty' that achieves the desired

9. What if we absolutely wanted to limit ourselves to, say, 7 features? This may be important if we want to derive "a

• Fit a regression model with a given I1\_penalty on TRAIN data. Add "alpha=I1\_penalty" and "normalize=True" to the parameter list.

#8, adding 1 whenever the intercept is nonzero. Save the number of nonzeros to a list.

sparsity. Here, we will again use a validation set to choose the best value for 'l1\_penalty'.

too many non-zeros. More formally, find:

Hint: there are many ways to do this, e.g.: Programmatically within the loop above

• Creating a list with the number of non-zeros for each value of l1\_penalty and inspecting it to find the

• The largest I1\_penalty that has more non-zeros than 'max\_nonzeros' (if we pick a penalty smaller than this value,

we will definitely have too many non-zero weights)Store this value in the variable 'l1\_penalty\_min' (we will use it

• Extract the weights of the model and count the number of nonzeros. Take account of the intercept as we did in

**14.** Exploring narrower range of l1\_penalty

We now explore the region of I1\_penalty we found: between 'I1\_penalty\_min' and 'I1\_penalty\_max'. We look for the

L1 penalty in this range that produces exactly the right number of nonzeros and also minimizes RSS on the

 Fit a regression model with a given I1\_penalty on TRAIN data. As before, use "alpha=I1\_penalty" and "normalize=True".

Find the model that the lowest RSS on the VALIDATION set and has sparsity equal to 'max\_nonzeros'. (Again, take account of the intercept when counting the number of nonzeros.) 15. Quiz Question: What value of I1\_penalty in our narrow range has the lowest RSS on the VALIDATION set

标记为完成