## Regression Week 5: LASSO Assignment 1

In this assignment, you will use LASSO to select features, building on a pre-implemented solver for LASSO (using GraphLab Create, though you can use other solvers). You will:

- Run LASSO with different L1 penalties.
- Choose best L1 penalty using a validation set.
- 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.

## 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
  (https://eventing.coursera.org/api/redirectStrict/mX9LT2DDyPv1WXdzeciZDNu_z-
  8_HQm2AwnTveOqAXuTqB3uVOEPPrKzPgzuoFOWFxRnBTZqzNsfhden8lSNgw.kIPn8unKIYjLnBDpHS_7iA.0Vio7Q
  hA1goWIf7G5mV9x1lvMBvrgVXY_N1GU0728CAd-
  ZE1rTAEtjfALIwS8qP0zV5RXVo5TLju8HBigxxbdb6f_eJknvnKu9sPGqyY9qj8RwUrQZFYuZraYDrdKeZV_WYlF06XhZa
  uni8lyWzSSieQVw_XbBMslA9n6W5cRX64zoao_o3GEFeGcLP-
  5BIDrKIIAormD9uzMd9_B5Hr2a7Zv1jesCOu9EmHdg4WNRpc4ymRi80TWwWDNddMfWuNt-
  kSXUF6VHDYUwu1FtaXvNZyA6kGsljqQTpY_f6Gn2vEiOpiStl08hZaRIGrzFwDApDIRZN6m3nMWKfJSZKPqQTdkiQIK
  NrrAuA3jMEEZ_55lnuHGCpPdNSJKtdR3IHvxgS63wrSCa_7iUIecYwpwxyu12x_cTE_2QPhmosVcgIB7eiIYFC9YzmTyY
  eHn0vB)

- Download the companion IPython Notebook: week-5-lasso-assignment-1-blank.ipynb
  (https://eventing.coursera.org/api/redirectStrict/aUIzV1F8vG1fr-6hsL5wEwRa4G73GMltEAWHQgN7JnZRlYU1mj7-
  FE3RUcdDe3qmWxct3zmJevXIHOQn88IBgg.R5AtsTLIYUN-mQPMLiQ2yA.I6EfVYRr3hw2f64hP5DooLO-
  P7wzlYih0HvKO-2hvhPS9Vj-j8Slm9q-
  5yLzyDcllRvJB3mngiX6XeoUayO0EWdt4AwZQfTLNa0gQkfcqz5T3lKplmuJZtOhNOymf1-
  FLlZoxJVEfpWqHFq1KrGL3oj6TKG7gXqg_Pf8JK8Qx-
  XApzizCpk6T3S49CsqyxWphRVgnme2GCJLMRYN0hU5z3PEwgrJIhZ6fTbeceWutPZ9ovfLQqn0R-
  7LDi0SZwob3pAnUXpzRNWvMcyvgQDgXDI9cbaEMAhCpGvC_H_lIuZjB68FE7Jn7-4R9bprE04ceu76hwc6jIJ4mmf3o
  ONsYf-kDyOqbDQBhIDDADM7bSpUyYk1DzRKzVIhZ5CcgZBHL89BU2yD_BedWK11SJk9Kck9o4cxg-B_OHj5j-
  hEOBO6VZDrJQo0r7VFuahxIowU2n277GOfi5UzAyMaqfpuw0HGT7bQS62ahhMuDRPwKLQ)

- 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 (https://eventing.coursera.org
  /api/redirectStrict/ohao-2qW-5N9AuXFL89SyTy-TUmdY3JeHlzd67UJrXZdk9Q86Q-z4retLhLU-62P5I-
  0kVP7PsPJawKo-U44sw.ZtbWUjGv1Jn1k2-wIMoRdw.-RnQT57yGGwTTqVeubHRuRriUW-
  5sl8NO7V8y2kXNJ5EYB9WyEq8sjSRn1QvonUmacsW-WDmyiKdjYvqOxzKw8xl5Qu7PD-
  8OUSKMa4MYqKgCNlBnMMO-UP3s8WVf4UXUMPUxJ4MiZ03S_OpN4X8Udu2eqvjtHILomepOrZrl7F8vAoaZBJ2B-
  aE7hfmcH_NuZFl7nFayZTVcY67z-VLwWPjITOaF-
  yxMDVbZkf4A976qqabQbRme7v_pw56SmpA4wy4ErBfA1nSUEpcAiVacZlHVF1Om00Lp-
  HlG5gwwUZZEhxFO7c4iZ00apMpLcpL4wZXVhBrzljudThpk8ct1n7H-
  A8pJ_ZloA1Z35ZeLBqOec7zKzwwGiSCVjijYaG6cwuQvookAo60OOkbmWPTz6iAUBpvZ2uSWKVwNthiSXveXGEQs7n
  61hTJWcmQAyOQ)

- NOTE: The following files are from Module 3 (Polynomial Regression).

- Download the King County House Sales training data csv file: wk3_kc_house_train_data.csv
  (https://eventing.coursera.org/api/redirectStrict/tnOn6HwSMJ40UDeUEKgGXyDgdpJ5qDaSexuS4SaX0_kYYS-
  yJQSmtMKXsrdE_aMJN6JwvNnkOMNQ9-
  BoC4XDiw.RK6252buvYYqGUtu6XX5ZA.OoipGExfd7w44_hsR3CNr56DB4Z1ujMrFeOEBEhhYBqxaJdcepYChEFtRjV
  UB8WvSouXekE3O5O6SUObbnVRq8B-js5413FF7smMFyj5I-ttGXxI2-
  S6ClXy850B34Mzf4luAgHe5wiRmrtNRT4BcFdUqK-
  NYuuqj0I1ajgLtNyyRRsQSvX_VlMvYO5bpq4KwQmxGspzCaj5rLpxpq3RDrc_FEPwd4Ck50ahgJjanUPD6CFsZjhM4US
  jezbPZTKqFXDrwqrs7uKKYIrQch33Mml_Ce_SYW_h1bktlyPXUBMi4xtBWuSDQfGXVwwp2NJsMf_7cO_2dlNv7QBxT
  ES0LG2FM4EBTt6mKNYskToPhJ9tIo2RybLubPK-W_-MGxpFEZ8kfL5ghWJqzRH5L-
  82QlyvDoUjHboSMJoVQuenIsc0lEzBd5cd8rSaDzpocTFZS-FNQPv9sW4pB-x83eBEug)

- Download the King County House Sales validation data csv file: wk3_kc_house_valid_data.csv
  (https://eventing.coursera.org/api/redirectStrict/GupnKJBZC4a2kBGlfwNfh_tUW-FdjzJrA92Bxu8sc05gtBgB2ACCV-
  Nw639XSzFJ_XW81x1b-
  xVUzwhB5KwWWw.uNk3_iNdYRZAYPrq2XqP6w.Cd6S9Ib3Niej0yjUdA5nb5PrArEay0lvdbrNPhxy_ZTvxl9DCfsLmhw
  zipRKiNOLm2GYRKBQDE2Rl1nHmloex5mGK9ph54GRf1Gaq7GE7RYgOeMYOXQp-
  yIsKdkLJ4Z9XUX8xhmwyzClf8_S-
  QBAYVAKLGV50KdCsHwj2IiQj6T5OpVAPtxoqqcK6dg_qxklpDGAFi3PjxADib9ksSyFPXObpVWpC3t543KBPlvjQBbbrj
  IuGYc2-
  Uo_zW3m4CZTp48PA6jQWroJ8Ql9blVLzGF2lffjUm6XYs6qcgacLvnveYp3--wetFAwO1qzbfOpLXOLSuWn4DoTOXoP
  s0qQMexQg043bcVgDINbrigZhfUKBIgZ9AqN-
  J-534zAG45jC887U9Mw8tYXGNmLJbZu5wrN8heR48yIvh2UoP0clZ2WcfAER9Vc7p8Dr7FIJe60rn8ghqcKTMRefG_2g
  42qNg)

- Download the King County House Sales testing data csv file: wk3_kc_house_test_data.csv
  (https://eventing.coursera.org/api/redirectStrict/iqeabMmMxeaRXU74fGXHawq7YUoi-
  mb1XAljs3SB2QpZz1g40PmRLalgyh7MD44dJOW2nyL4GNVYwUmiBvxm-Q.gYnWGuTP71B9g-
  hnwVEVOQ.b3cRYTyxiM2Doe5bqxgTt2_D-
  C2vtiojNFbTtRezk24dgOZtmrbdTpFK7SUUTb6nvY5fMnaPphIFVfMpFRU1XDPnleWrxhpjKzgfz5PVSAa2Aau_SwR_v
  bsV7soWrpDX1mVinXWAK9kdwcMAxIBbMtHKkXDP1fAHtct8gpItnl2EYlProQLvlmw-
  RNoplJazbSTEoTUHGlt24hk87D7L-
  57hIfnHob6-02iDA7XM7Bh9EXKRUqWsUqn4Gwo4g2ALg7bvZBy4YYn0Is9ESBO3aoKrqaSlD7pL8rfGGCyG1HsNk5
  e85kKgUJCvObPCo-AzkwghZIASseHUyFKq0OQIfkC2YZ3hwpTHzrauFx88VRkXtwgMAvVzFOkXAjf35x_9idM-
  IvfNg60vsMeDEwjBExoAJpilR_6cPv9TOmjEHPpQNX7KXEKS-kNXzjXYPSvkHAP9Q3pTl4T3__nrRbgwZQ)

## 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. Create new features by performing following transformation on inputs:

- 'sqft_living_sqrt', by taking square root of 'sqft_living'

- 'sqft_lot_sqrt', by taking square root of 'sqft_lot'

- 'bedrooms_square', by squaring 'bedrooms'

- 'floors_square', by squaring 'floors'

2. Using the entire house dataset, learn regression weights using an L1 penalty of 1e10. Use a LASSO solver in your tool of choice.

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:

- First split sales into training_and_validation and testing with sales.random_split(0.9) use seed = 1.

- Next split training_and_validation into training and validation using .random_split(0.5) use seed = 1.

If you're not using SFrame, please download the provided csv files for training, validation and test data.

5. Now for each l1_penalty in [10^1, 10^1.5, 10^2, 10^2.5, ..., 10^7] (to get this in Python, type np.logspace(1, 7, num=13).)

- Learn a model on TRAINING data using the specified l1_penalty.

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

Report which L1 penalty produced the lower RSS on VALIDATION.

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.

8. Quiz Question: What is the RSS on TEST data for the model with the best l1_penalty?

9. Quiz Question: Using the best L1 penalty, how many nonzero weights do you have?

10. What if we absolutely wanted to limit ourselves to, say, 7 features? This may be important if we want to derive "a rule of thumb" --- an interpretable model that has only a few features in them.

You are going to implement a simple, two phase procedure to achieve this goal:

- Explore a large range of 'l1_penalty' values to find a narrow region of 'l1_penalty' values where models are likely to have the desired number of non-zero weights.

- Further explore the narrow region you found to find a good value for 'l1_penalty' that achieves the desired sparsity. Here, we will again use a validation set to choose the best value for 'l1_penalty'.

11. Assign 7 to the variable 'max_nonzeros'.

12. Exploring large range of l1_penalty

For l1_penalty in np.logspace(8, 10, num=20):

- Fit a regression model with a given l1_penalty on TRAIN data.

- Extract the weights of the model and count the number of nonzeros. Save the number of nonzeros to a list.

Hint: np.count_nonzero (https://eventing.coursera.org/api/redirectStrict
/v4r4xvycryK59WazRY3CxgMUgHRo9cSh8AvZPQ0x6j2v8JOh-
wu3pQKPmAV9RPcUWD3T6MWNBJ_nDjclDCxDlQ.TEX2uPde8FP6__trvLZtlA.wSAR74-
Cl2hwpU0oBKyE1d7lmrRV4pOrrUgOvmf20foH_vk3X-_Ln3gFNfTW-LcMBrBir7VcAo-s9C4CNgx9iMG-
z_doXOlcAkfoGE3ro2KXGzLDzxWFMYUg74Yk0WHl7Ac_J9mwD3zaieciVfhdOwvrDapLIdkLB4xBdilWRtijfLU9ug5XA3X
qn_LY1y1AhB7Fj3SbOEAmCiMVbXWWX_nzI8GHf_-TGN-LwfgeWOUCl57p_qk6zy0e_rLZ0z7hzyraWxWw-
rqcasK2CAi4tS098NLesrVo3HTKcHeEc25Mudn4suzLYd8E8G9YEb7pwnfXBuAHhLD1tn46I52yutgxV547j18S6sKQQM
eNFJLZyZGVbXu2awNOSxilBjMHUUNcDuVwpSbmUbib0bkVPEni2SuRNQNSB5tGcPDB2iE) may be helpful.

13. 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 l1_penalty values that have too few non-zeros, and at the other end, we will have an l1_penalty that has too many non-zeros.

More formally, find:

- The largest l1_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 later)

- The smallest l1_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 'l1_penalty_max' (we will use it later)

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 appropriate boundaries.

14. Quiz Question: What values did you find for l1_penalty_min andl1_penalty_max?

15. Exploring narrower range of l1_penalty

We now explore the region of l1_penalty we found: between 'l1_penalty_min' and 'l1_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 VALIDATION set.

For l1_penalty in np.linspace(l1_penalty_min,l1_penalty_max,20):

- Fit a regression model with a given l1_penalty on TRAIN data.

- Measure the RSS of the learned model on the VALIDATION set

Find the model that the lowest RSS on the VALIDATION set and has sparsity equal to 'max_nonzeros'.

16. Quiz Question: What value of l1_penalty in our narrow range has the lowest RSS on the VALIDATION set and has sparsity equal to 'max_nonzeros'?

17. Quiz Question: What features in this model have non-zero coefficients?