**Introduction**

In this project, we implemented linear regression with $L_2$ (quadratic) regularization which learns an weight vector w to optimize Sum of Squared Error objective(SSE) using Batch Gradient Algorithm. We aim to understand the difference between linear regression without regularization and with regularization in the value of SSE through this project.

      In part 0, we discuss preprocessing and simple analysis. Preprocessing helps us in utilizing data in training by making us think which feature might be useful, in which way we can transform the data to get a better result such as binarization and normalization.

      In part 1, we explore the influence of learning rate. In this section, we found the overfitting problem and found that when the learning rate is big, it is easy to get explode issue. However, when the learning rate is small, it is extremely to get the answer in life time. Surprisingly, we found that most of the converged answer without regularization is overfitted.

      Then we tried a lot in part 2 for the experiment with regularization. We trained the data with different learning rate and record different training and validation SSE and the number of iterations does the training need. Amazingly and surprisingly, we were excited about our observation just prove the equation. Any two number of iterations make same ratio with the ratio of their learning rate with a same regularization. Furthermore, we found that when  is increasing the training SSE is also increasing monodically, while the validation SSE would increase first then drop down with a relative large . However, sometimes, a really large  will also makes algorithm get exploded as our observations. In fact, with a good regularization, the converged model could behave much better on validation SSE.

      In part 3, we used non-normalization data to train. In this section, we found that it is really hard to find a enough small learning rate for training without normalization. Beside, without normalization, the gradient also becomes very huge and needs more time for running the algorithm. In addition, at the end of this project, we performed training using the binarized data, which is the idea mentioned in preprocessing part 0.b for exploring more details about the idea in Part 3. b (Extra). We binarized some features like date and zipcode. Binarizing date doesn't improve the result a lot, however, binarizing zipcode decrease our validation SSE over 30% off. Finally we used the model which are binarized in both date and zipcode feature for predicting the result of test file.

## Part 0. Preprocessing and simple analysis
### a. Remove the ID feature. Why do you think it is a bad idea to use it in learning?

: Because the ID feature almost doesn't have any relationships with the house price like using the feature of today's weather to predict the result of tossing a normal coin. Extremely speaking, if we just keep and only keep this feature for predicting a toss, we also could get a non-zero weight on this feature, which means this feature will cause some influence on our result. Since we have already known the background knowledge of there is no or very rare connections between the ID feature and the house price, just like the connection between weather and toss a coin, so the influence on our result must be a bad influence. So we have to remove this feature for a better result. On the other hand, we also could treat it as a overfitting problem. Since we have one more useless feature to the result, while we have to adjust our model to fit on this useless feature. As a result, our new model must be more overfit than the removed-model. From this perspective, the result also can be improved.

### b. Split the date feature into three separate numerical features

: Yes. Binarizing these separated features would be one way. After normalizing these date feature, there is still a inner linear relationship within the data which means that December have more influence than the January (12 > 1, and when we calculate w*x, the result in December could be increased inappropriately. ) While, instead of separating and normalizing it, we could increase the number of features and binarize these data. Like separate date to 100 (approximate # of years) + 12 (# of month) + 31 (# of date) = 143 new features. If the date is "Oct 7th, 2018", then the "Oct", "7th", and "2018" should be set to 1 and the rest field should be set to 0.

### c. Build a table that reports the statistics for each feature



```
dummy : mean:  1.0 deviation:  0.0 range:  1.0 - 1.0
bedrooms : mean:  3.38 deviation:  0.94 range:  1.0 - 33.0
bathrooms : mean:  2.12 deviation:  0.77 range:  0.5 - 7.75
sqft_living : mean:  2080.22 deviation:  911.29 range:  370.0 - 9890.0
sqft_lot : mean:  15089.2 deviation:  41201.83 range:  572.0 - 1651359.0
floors : mean:  1.5 deviation:  0.54 range:  1.0 - 3.5
waterfront : mean:  0.01 deviation:  0.08 range:  0.0 - 1.0
view : mean:  0.23 deviation:  0.76 range:  0.0 - 4.0
condition : mean:  3.41 deviation:  0.65 range:  1.0 - 5.0
grade : mean:  7.67 deviation:  1.18 range:  4.0 - 13.0
sqft_above : mean:  1793.1 deviation:  830.82 range:  370.0 - 8860.0
sqft_basement : mean:  287.12 deviation:  434.98 range:  0.0 - 2720.0
yr_built : mean:  1971.12 deviation:  29.48 range:  1900.0 - 2015.0
yr_renovated : mean:  81.23 deviation:  394.36 range:  0.0 - 2015.0
zipcode : mean:  98078.29 deviation:  53.52 range:  98001.0 - 98199.0
lat : mean:  47.56 deviation:  0.14 range:  47.1559 - 47.7776
long : mean:  -122.21 deviation:  0.14 range:  -122.514 - -121.319
sqft_living15 : mean:  1994.33 deviation:  691.87 range:  460.0 - 6110.0
sqft_lot15 : mean:  12746.32 deviation:  28239.83 range:  660.0 - 871200.0
month : mean:  6.59 deviation:  3.11 range:  1.0 - 12.0
day : mean:  15.8 deviation:  8.62 range:  1.0 - 31.0
year : mean:  2014.32 deviation:  0.47 range:  2014.0 - 2015.0
```

**Figure 1.**



```
waterfront :  defaultdict(<class 'int'>, {0.0: 0.993, 1.0: 0.007})
condition :  defaultdict(<class 'int'>, {3.0: 0.653, 5.0: 0.0812, 4.0: 0.2569, 1.0: 0.0013, 2.0: 0.0076})
grade :  defaultdict(<class 'int'>, {9.0: 0.1182, 8.0: 0.2838, 7.0: 0.413, 6.0: 0.0933, 10.0: 0.0547, 5.0: 0.0105, 11.0:
0.021, 4.0: 0.0011, 12.0: 0.0039, 13.0: 0.0005})
```

**Figure 2.**

|  | mean | standard deviation | range |
|---|---|---|---|
| **dummy** | 1.0 | 0 | [1.0,1.0] |
| **bedrooms** | 3.38 | 0.94 | [1.0,33.0] |
| **bathrooms** | 2.12 | 0.77 | [0.5,7.75] |
| **sqft_living** | 2080.22 | 911.29 | [370.0,9890.0] |
| **sqft_lot** | 15089.2 | 41201.83 | [572.0,1651359.0] |
| **floors** | 1.5 | 0.54 | [1.0,3.5] |
| **waterernt** | 0.01 | 0.08 | [0.0,1.0] |
| **view** | 0.23 | 0.76 | [0.0,4.0] |
| **condition** | 3.41 | 0.65 | [1.0,5.0] |
| **grade** | 7.67 | 1.18 | [4.0,13.0] |
| **sqft_above** | 1793.1 | 830.82 | [370.0,8860.0] |
| **sqft_basement** | 287.12 | 434.98 | [0.0,2720.0] |
| **yr_built** | 1971.12 | 29.48 | [1900.0,2015.0] |
| **yr_renovated** | 81.23 | 394.36 | [0.0,2015.0] |
| **zipcode** | 98078.29 | 53.52 | [98001.0,98199.0] |
| **lat** | 47.56 | 0.14 | [47.1559,47.7776] |
| **long** | -122.21 | 0.14 | [-122.514,-121.319] |
| **sqft_living15** | 1994.33 | 691.87 | [460.0,6110.0] |
| **sqft_lot15** | 12746.32 | 28239.83 | [660.0,871200.0] |
| **month** | 6.59 | 3.11 | [1.0,12.0] |
| **day** | 15.8 | 8.62 | [1.0,31.0] |
| **year** | 2014.32 | 0.47 | [2014.0,2015.0] |

**Table 1.**

**d. Which set of features are expected to be useful?**
: We think the feature of the feature of "sqft_above" should be the most useful feature in predicting the result. Based on our common knowledge, the price of house should be based on the equation: house price = total square * the quality of the house (unit price of sqft). Thus the newer could effect the quality of the house. However, the date at this domain doesn't have a rough deviation (all of house built in 2014-2015. ) Beside, some features like the condition, zipcode, waterfront, grade, floors, lat, long just have a small

deviation. Intuitively, I think those features would cause small influence on the result. Thus, the total square of the house is the most useful feature for predicting a house.

**e. Normalization**
: We looked up the internet about how to normalize the data, then we found four ways for normalization.

(1) $a_1$ / sqrt(sum($a_1^2$ +…+ $a_n^2$)) … an / sqrt(sum($a_1^2$ +…+$a_n^2$))

(2) $a_1$ / max($a_1$… $a_n$) … an / max($a_1$…$a_n$)

(3) $a = (a$ - min) / (max - min)

(4) z-score normalization: a = a - mean / standard deviation.

We will use the last one normalization method (4). If we have time, we will implement the second method and compare their different influence. Figure 2 shows a part of normalized data.



**Figure 3.**

# Part 1. Explore different learning rate for batch gradient descent
**a. Observation of different learning rates**
: $10^{-5}$ is the best learning rate in my observation. More specifically, $5 * 10^{-5}$ is even better and faster than $10^{-5}$.

$10^0, 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}$ and any learning rate greater than $5.8 * 10^{-5}$ will make the gradient decent explode.

x-rot is the number of iterations, y-rot is the value of SSE with each learning rate:
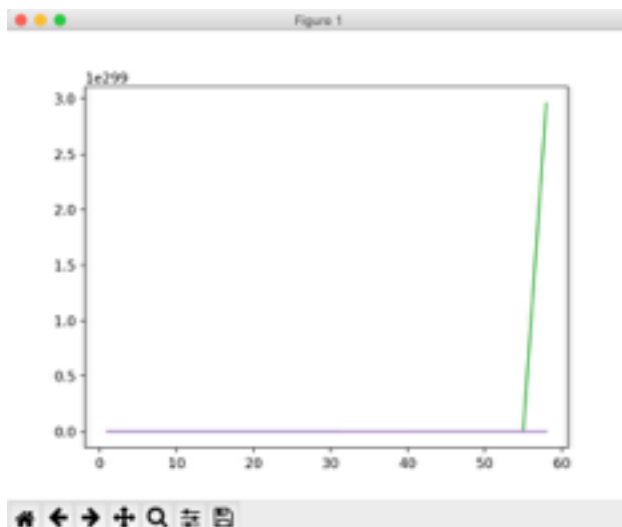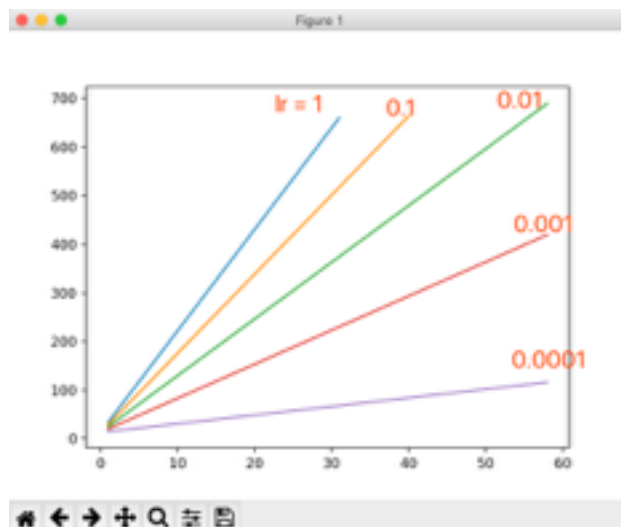
**Figure 4.**



**Figure 5.**

In Figure 4., due to the rapid growth of the result of SSE with the growth of iterations, most of data just become a line. So we take a log for shrinking the gap among the data.

In Figure 5., we can see that the bigger learning rate it is, the more rapid growth it is. Beside, if the growth is bigger than 700, the SSE value just become "inf" and stop growing anymore. The growth of learning rate at 0.0001 is the smallest learning rate is also grow slowest. The situation above is only taking about the explode situation.



**Figure 6.**

In Figure 6., we can see that the tendency of these three learning rate is convergence. The bigger learning rate always brings a faster convergence speed. However, there is another issue. When the number of iterations is growing up, the speed of convergence is decreasing. Thus, we should better pick up an as big as possible learning rate for decreasing the running time. Therefore, we can conclude that the biggest learning rate cannot make a convergence. But as long as the learning rate could make a convergence, we should better use a bigger learning rate for speeding up.

**b. SSE on the training data and validation data with iterations**

: For setting up epsilon at 0.5 with a learning rate 0.000057, the training time is 916s with 25734 iterations. When learning rate is 0.00001, it takes more than 5 hours and didn't get a result. Thus our solution is assigning w at the beginning instead of just assigning 0 to w.

Strict unequal equation: $|\Delta E(w)|$ (interpret as getting norm of a vector) $< \epsilon$

| Learning Rate | 0.000057 | 0.00001 | 0.000001 | 0.0000001 |
|---|---|---|---|---|
| Time | 916s | 5672s | - | - |
| Iterations | 25734 | 146684 | - | - |
| SSE on train | 38991.5180114 | 38991.5180414 | - | - |
| SSE on dev | 55026.301127 | 55026.2647565 | - | - |

**Table 2.**



**Figure 7.**

Figure 7. shows that the SSE of training set with learning rate 0.000057 which may be the fastest learning rate that I chose. However, we found that in this figure, after a few iterations, the algorithm just keep going, while the SSE didn't change a lot. Beside, with using this strict equation, it took really a long time for reaching the condition. (It took 1.5h with learning rate = 0.00001. We tried a long time with smaller learning rate, but it didn't come out answer in 2 hours, so we just skip trying that. ) However, there are too many iterations just improve a little bit for convergence. For learning rate = 0.000057, we just tried first 500 iterations for zoom-in the figure. Then;



**Figure 8.**

| Iterations | 500 | 25000 |
|---|---|---|
| SSE on train | 39651.03976671272 | 38991.5180114 |
| SSE on dev | 23080.3020957 | 55026.301127 |

**Table 3.**

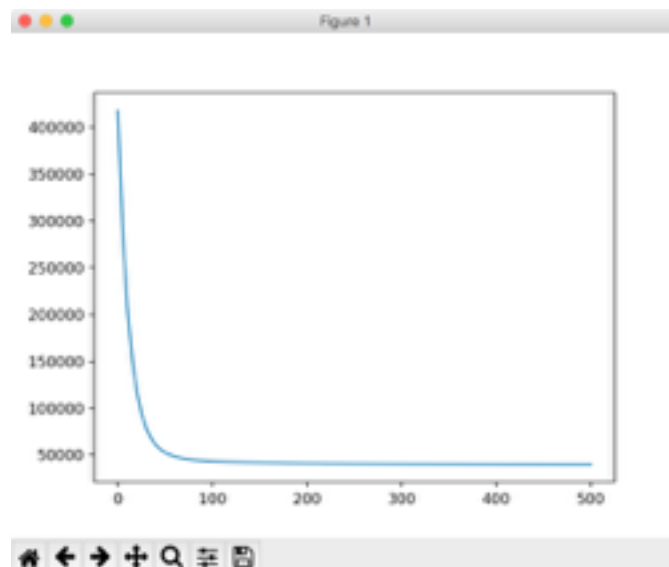In above data set, we found out there might be an overfitting problem. (less errors on training data with well-trained model but more errors on dev data. ) Hence, instead of setting an unknown loop stop point, we set up a fixed 25000 iterations, which costs 15 min per training.

learning rate = 0.000057, x-rot is # of iterations, y-rot is SSE, blue is on training data, yellow is on dev data. Figure 9. shows total SSE on each data set. Figure 10. gives average SSE on each data set.



**Figure 9.**                                                          **Figure 10.**

So at the first drop of SSE on dev data, SSE turns to increase a lot. When we take an average SSE, the gap becomes more obvious. We think that it could prove the existence of overfitting problem. Beside, more iterations and large learning rate doesn't guarantee a good model. For showing a more clear result, we take a log function for shrinking the huge gap and magnifying the small gap:
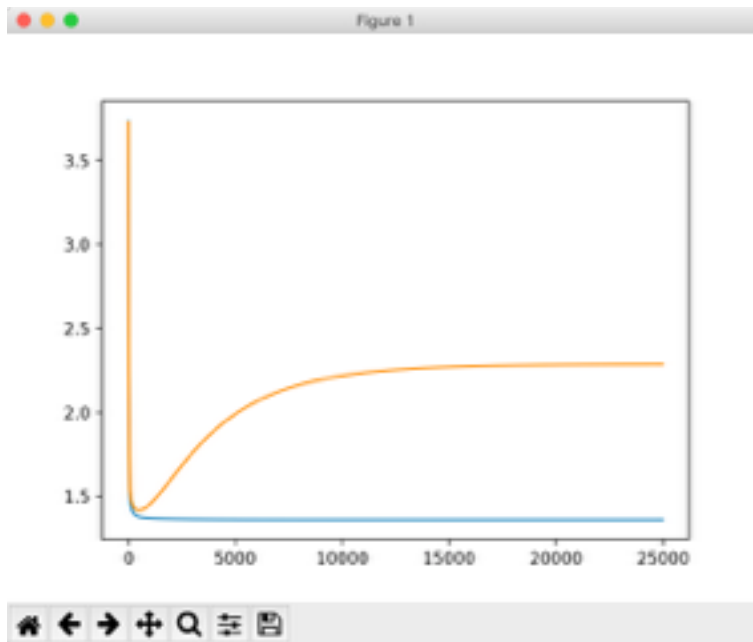
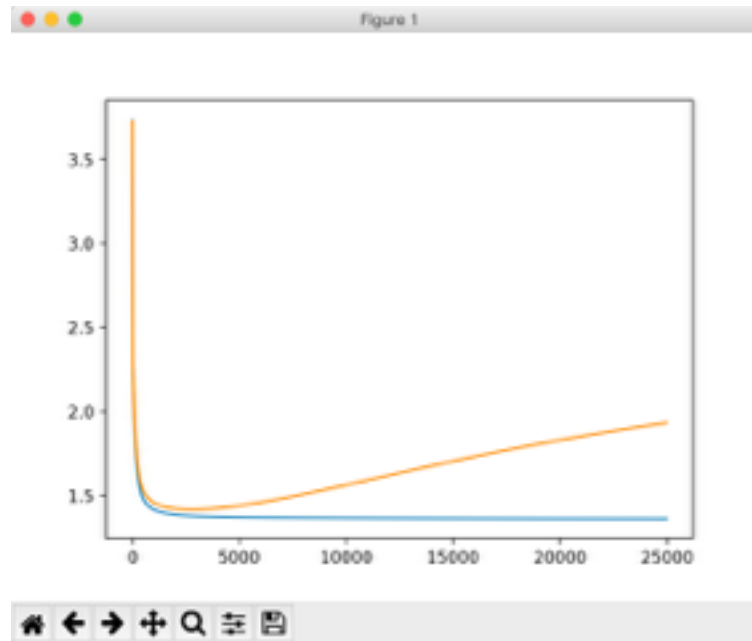**Figure 11. Learning rate = 0.000057, Y-rot = average log(SSE)**



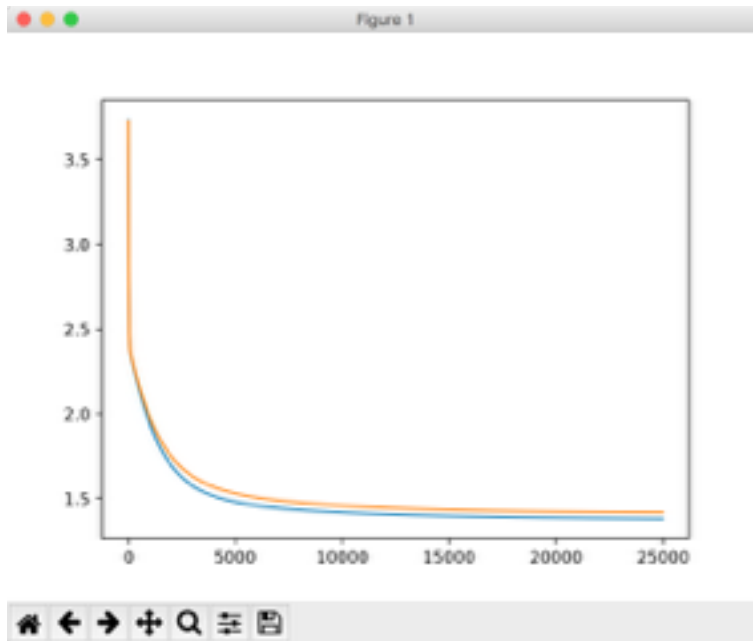**Figure 12. Learning rate = 0.00001, Y-rot = average log(SSE)**

**Figure 13. Learning rate = 0.000001, Y-rot = average log(SSE)**
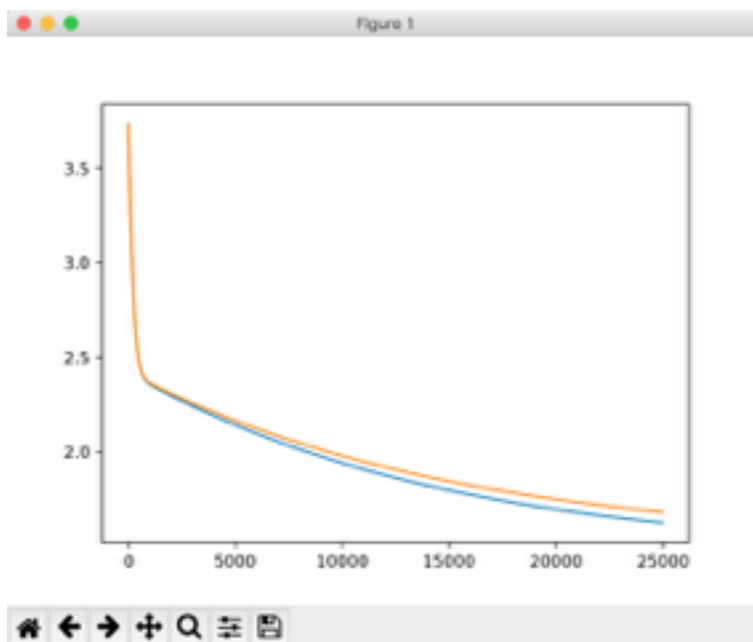


**Figure 14. Learning rate = 0.0000001, Y-rot = average log(SSE)**

From above observations, we found that some facts.

1.  The average SSE of dev data are always higher than the average SSE of training data. It's hard to surpass, otherwise it should be an overfitting problem on the dev data set.

2.  A lower learning rate needs more iterations and more time for convergence. From the two data we got, when learning rate a is 5.7 times to learning rate b. The iterations and the cost time does learning rate a needs only rough 1/5.7 to what the learning rate b needs. (a = 5.7* $10^{-5}$, b = 1*$10^{-5}$, while # of iterations ratio is $a/b = 146684/25734 = 5.7$). So it takes around 1.5h for getting the result with learning rate b = 1*$10^{-5}$, we assume that it may take 15h with learning rate 1*$10^{-6}$ and 150h with learning rate $10^{-7}$.

3.  The total SSE on training data is decreasing monotonically, while total SSE on dev data is decreasing first and then turn to increase a little bit with a fixed model learned from training set. We believe that is the case of overfitting problem, with the model more fitting on the training data, it starts to be not fitting on dev data.

4.  From 2, we know that a lower learning rate will decrease the speed of training, while it could extend the period of decreasing of SSE on dev data. When we want to train a model, we don't want a model which only could fit on current data. So our goal should be getting the lowest SSE of testing on dev data. However, the lowest SSE of dev data with a model which is trained with a higher learning rate cannot be found easily. In our examples, when learning rate is 5.7*$10^{-5}$, the lowest SSE value only lasts a short time. Thus, the higher learning rate isn't always mean a good thing. Sometimes, it cannot let us to find the lowest SSE value on dev data. (We almost missed it, because at the beginning, we only tried the model after convergence which may have already been the overfitted model. )

5.  When learning rate is equal to $10^{-6}$ and $10^{-7}$, there is not an increase of SSE on dev data. However, we assume that there isn't an increase of SSE on dev data so far, because there isn't enough iterations for them to get more fit on the training data.  So we assume may be after 25000 iterations ($10^{-5}$ change its tendency after around 2500 iterations), the SSE on dev data with learning rate $10^{-6}$ will turn to increase and SSE on dev data with learning rate $10^{-7}$ will change its tendency after 250000 iterations.

6.  The lower learning rate cannot guarantee a lower SSE on dev data set.

    LR = 5.7 * $10^{-5}$ -> minimum SSE 23078.986037
    LR = 1.0 * $10^{-5}$ -> minimum SSE 23080.3943119

**c. Learned weights for each feature**
: When learning rate is equal to $10^{-5}$, the number of iterations for convergence is over 140 thousand. In this feature weight, the most absolute biggest feature is the "number of bedrooms".

It is similar to our pre-analysis because according to our approximate prediction equation: size of the house * quality of the house. Although we didn't think the most biggest feature is the "# of bedrooms", however, the most top 4 features are "# of bedroom", "sqft_living", "sqft_above" which are representing the size of the house. Then "grade" represents the quality of the house. Thus, this fit on our rough prediction.

| Features | weight |
|---|---|
| dummy | -2.02634726 |
| bedrooms | -9.1811135 |
| bathrooms | 3.33191851 |
| sqft_living | 7.26277625 |
| sqft_lot | 2.31366786 |
| floors | 0.11770157 |
| waterernt | 4.71107643 |
| view | 2.32284388 |
| condition | 1.25636608 |
| grade | 8.8744548 |
| sqft_above | 7.72517044 |
| sqft_basement | 1.3069606 |
| yr_built | -2.9363435 |
| yr_renovated | 0.28851062 |
| zipcode | -1.02065174 |
| lat | 3.76772534 |
| long | -2.67437416 |
| sqft_living15 | 1.33430012 |
| sqft_lot15 | -2.91738238 |
| month | 0.18228678 |
| day | -0.17278705 |
| year | 0.38104872 |

**Table 4.**

Beside, the feature weight of month, day, year are also very small and just like what we predicted. However, the lower feature weight somehow represent the relationship between the feature and prediction label. It's really hard to assert that the contribution of deviation, mean and range to the feature weight. (Like "year" feature has a low deviation while "day" feature has a relative higher deviation, but both of them don't have a larger weight. ) Thus, we guess that the lower feature weight only relate to the connection to the prediction label (in this domain, it should be the relationship between feature and predicted price. )
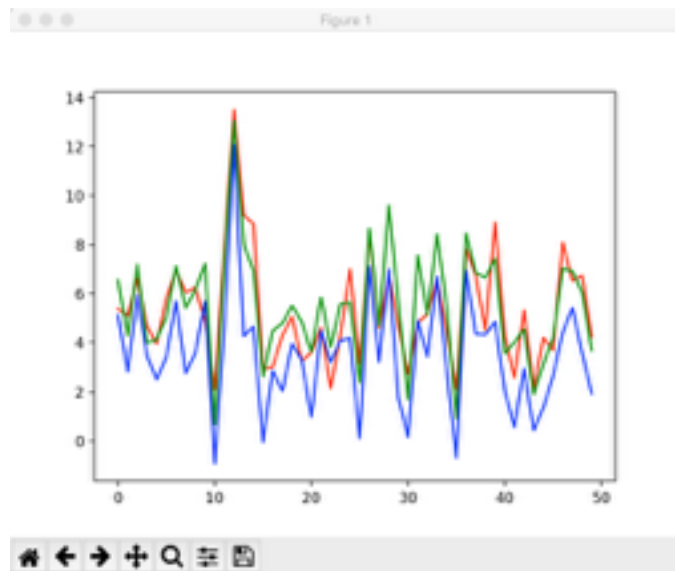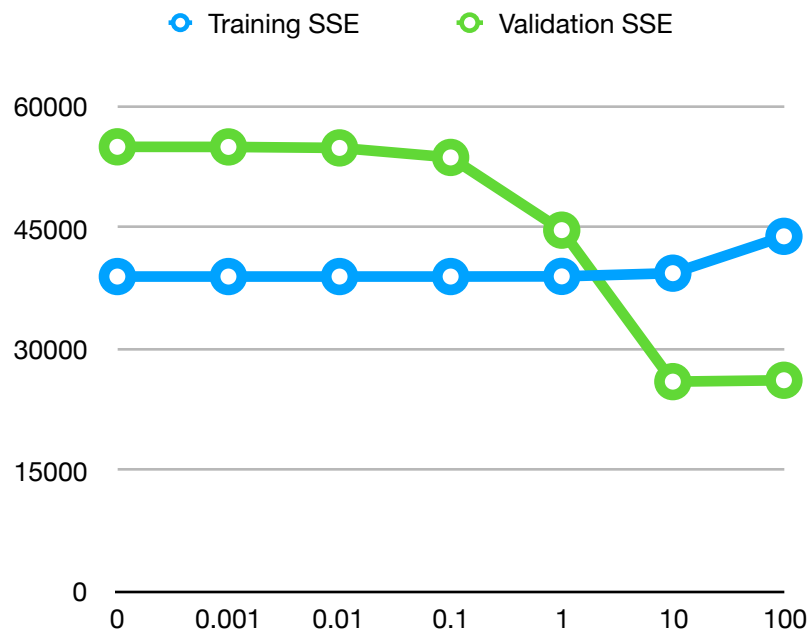


**Figure 15.**

Figure 15. shows that the first 50 predicted result of predicting in dev data. Red line is real data, green line is the model of 484 iterations model with learning rate 0.000057, blue line is the model of fully convergence 25734 iterations overfitting model with learning rate 0.000057.

## Part 2. Experiments with different λ values



**Graph 1.**

### SSE of Training Data

| Learning Rate \ λ | 0.0 | $10^{-3}$ | $10^{-2}$ | $10^{-1}$ | $10^{0}$ | $10^{1}$ | $10^{2}$ |
|---|---|---|---|---|---|---|---|
| $10^{0}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-1}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-2}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-3}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-4}$ | explode | explode | explode | explode | explode | explode | explode |
| $5.7 * 10^{-5}$ | 38991.5 | 38991.5 | 38991.5 | 38992.0 | 39012.3 | 39408.2 | 43962.9 |
| $10^{-5}$ | 38991.5 | 38991.5 | 38991.5 | 38992.0 | 39012.3 | 39408.2 | 43962.9 |
| $10^{-6}$ | too long | too long | too long | too long | too long | too long | 43962.9 |
| $10^{-7}$ | too long | too long | too long | too long | too long | too long | too long |

**Table 5.**

## SSE of Dev Data

| Learning Rate \ $\lambda$ | 0.0 | $10^{-3}$ | $10^{-2}$ | $10^{-1}$ | $10^{0}$ | $10^{1}$ | $10^{2}$ |
|---|---|---|---|---|---|---|---|
| $10^{0}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-1}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-2}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-3}$ | explode | explode | explode | explode | explode | explode | explode |
| $10^{-4}$ | explode | explode | explode | explode | explode | explode | explode |
| $5.7 * 10^{-5}$ | 55026.3 | 55012.8 | 54891.5 | 53714.1 | 44725.2 | 26003.0 | 26173.2 |
| $10^{-5}$ | 55026.3 | 55012.7 | 54891.5 | 53714.1 | 44725.1 | 26003.0 | 26173.2 |
| $10^{-6}$ | too long | too long | too long | too long | too long | too long | 26173.2 |
| $10^{-7}$ | too long | too long | too long | too long | too long | too long | too long |

**Table 6.**

## Iterations for Convergence

| | 0.0 | $10^{-3}$ | $10^{-2}$ | $10^{-1}$ | $10^{0}$ | $10^{1}$ | $10^{2}$ |
|---|---|---|---|---|---|---|---|
| $5.7 * 10^{-5}$ | 25734 | 25720 | 25592 | 24382 | 16763 | 5457 | 838 |
| $10^{-5}$ | 146684 | 146603 | 145874 | 138976 | 95551 | 31108 | 4784 |
| $10^{-6}$ | - | - | - | - | - | - | 47849 |

**Table 7.**

**a. Trend from the training SSE with change of $\lambda$ value**
: When $\lambda$ is going up, the training SSE is increasing monotonically.

**b. rend from the validation SSE with change of $\lambda$ value**
: When $\lambda$ is going up, the validation SSE is decreasing at the beginning and increasing at the end with a large $\lambda$.

**c. Explanation of observed behaviors**
: In part I, without any regularization, due to our batch gradient algorithm, the training SSE should decrease monodically for sure. So the convergence point is also the point where the lowest training SSE we can get. However, after adding a larger and larger $\lambda$, we are not only

speed up the training speed, but also decrease the accuracy on training data set. More specifically, the loss function is

$$\frac{1}{2}(\Sigma_{i=1}^{N}(x_i * w - y_i)^2 + \lambda ||w||^2)$$

then the gradient should be;

$$\Sigma_{i=1}^{N}(x_i * w - y_i)x_i + \lambda w$$

Then, we could use the result of learning rate by gradient for updating the value of w. However, when the learning rate is a fixed value and $\lambda$ is increasing, the model w is using larger and larger itself for updating itself. Thus, the term we used in class of "regularization" is for controlling how fast the w update itself. We list the # of iterations with 2 different learning rate. Thus, with a larger $\lambda$, the speed of training is faster for sure (because some iterations will be skipped when w update itself, then it only needs fewer iterations for convergence). At the beginning of this explanation, we also mentioned about that the $\lambda$=0 training could get the best result on training SSE for sure. When the w skipped some its iterations in $\lambda$=0, it must get a relative worse training SSE (It kind of like, when you are going down straight of the mountain with same speed, somehow, the lowest point of the mountain just increase a little bit, suddenly we just don't need too many time (iterations) for reaching the lowest point.) Thus, when $\lambda$ is going up, the training SSE is increasing monodically due to the regularization.
Beside, when $\lambda$ is going up, the validation SSE is decreasing at the beginning and increasing at the end with a large $\lambda$. Because the result we get in part 1.b is showing that there is some overfitting issue within the 25000 iterations training with learning rate both 0.000057 and 0.00001. The theory best model (get the lowest point on both training SSE and validation SSE) should only take around 480 iterations within a learning rate of 0.000057. So when the lambda is bigger, while the regularization is rough, the model just skipped a lot of useless updates and just try to keep and save the best model to the end of training. However, when the $\lambda$ is small, the influence of $\lambda$w is too small to the result, so the regularization is not strong and the result is trend to the result without any regularization. Our the third table shows the number of iterations does the model needs for convergence. Furthermore, when the $\lambda$ is too big (we tried 100000 as a $\lambda$), this large-$\lambda$ model just skipped almost updates, even include those best models (it only take 146 iterations when $\lambda$=100000, while the best model appears around 480th iterations). Therefore, we conclude this phenomenon in the issue of overfitting problem. When $\lambda$ is not big enough, the regularization is stronger and strong, then overfitting problem become less and less. As a result, the validation SSE become better and better. However, when $\lambda$ is too big, the regularization is too strong, then the model just too sparse and also cannot get a good validation SSE.
Sometimes, when $\lambda$ is really too big, the lower learning rate also got exploded.

**d. What features get turned off for $\lambda$= 10, $10^{-2}$ and 0 ?**
:

```
iterations:  146684
[-2.02634726 -9.1811135    3.33191851  7.26277625  2.31366786  0.11770157
  4.71107643  2.32284388  1.25636608  8.8744548   7.72517044  1.3069606
 -2.9363435   0.28851062 -1.02065174  3.76772534 -2.67437416  1.33430012
 -2.91738238  0.18228678 -0.17278705  0.38104872]
38991.5180414
55026.2647339
LR:  -1e-05   lambda:  0
```

```
iterations:  145874
[-2.02709724 -9.16190337  3.33098725  7.26006906  2.29894808  0.11813967
  4.7105344   2.32322219  1.2561597   8.87462106  7.72228459  1.30649313
 -2.93617532  0.28857949 -1.02049934  3.76766392 -2.67418807  1.33550008
 -2.90400529  0.18223169 -0.17280587  0.38100384]
38991.5418008
54891.4548422
LR:  -1e-05   lambda:  0.01
```

```
iterations:  31108
[-2.13879881 -2.66156923  3.07909106  6.01246993  0.36877585  0.40983885
  4.1366181   2.5177162   1.15140361  8.42127592  6.34451057  1.24037464
 -2.80987483  0.33655317 -0.93897075  3.70621651 -2.45438853  2.21914487
 -0.50729989  0.12500539 -0.191743    0.33943697]
39408.1769667
26002.9516732
LR:  -1e-05   lambda:  10
```

Figure 16.

Table 8.

**\*List absolute features less than 1 (small -> big)**

|       | turned off features |       |             |              |      |          |        |            |
|-------|---------------------|-------|-------------|--------------|------|----------|--------|------------|
| 0     | floor               | day   | month       | yr_renovated | year |          |        |            |
| 0.01  | floor               | day   | month       | yr_renovated | year |          |        |            |
| 10    | day                 | month | yr_renovated| year         | sqft_lot | floors | sqft_lot15 | zipcode zip |

```
[((0.17483381964327635, 'day'), (0.17527123320972934, 'month'), (0.35561144583333443, 'yr_renovated'), (
0.36754087405818531, 'year'), (0.37498949376340857, 'sqft_lot'), (0.39890142441959997, 'floors'), (0.50946801634478622,
'sqft_lot15'), (0.90307811246965486, 'zipcode zip'), (1.2291755179974513, 'condition'), (1.2347660380021661,
'sqft_basement'), (2.234251547351934, 'sqft_living15'), (2.3212545429545273, 'dummy'), (2.3956339202702299, 'long'),
(2.5133604292234928, 'view'), (2.5927247673156395, 'bedrooms'), (2.7660155631958574, 'yr_built'), (3.0734080662663033,
'bathrooms'), (3.7336570865356284, 'lat'), (4.1396290204717969, 'waterfront'), (5.9919285890350737, 'sqft_living'),
(6.323274033480331, 'sqft_above'), (8.4663957299532662, 'grade')]
```

Figure 17. $\lambda = 10$

```
[((0.11812950999016128, 'floors'), (0.1727894724265113, 'day'), (0.18228166681719901, 'month'), (0.28859923453935549,
'yr_renovated'), (0.38103178579858059, 'year'), (1.0204618766806515, 'zipcode zip'), (1.2562370408108545, 'condition'),
(1.306483843300684, 'sqft_basement'), (1.3355185147133308, 'sqft_living15'), (2.0272870555420974, 'dummy'),
(2.2990133446937664, 'sqft_lot'), (2.3232183626981153, 'view'), (2.6741238187990182, 'long'), (2.9040404816066925,
'sqft_lot15'), (2.9361315804034489, 'yr_built'), (3.3309735388815724, 'bathrooms'), (3.767690333764337, 'lat'),
(4.7105398310390694, 'waterfront'), (7.2600299383293585, 'sqft_living'), (7.7222436936534908, 'sqft_above'),
(8.8746793648328364, 'grade'), (9.161710811419006, 'bedrooms')]
```

Figure 18. $\lambda = 0.01$

```
[((0.11770150187029091, 'floors'), (0.17278705798409899, 'day'), (0.18228678432588732, 'month'), (0.2885106401240049,
'yr_renovated'), (0.3810487241636854, 'year'), (1.020651735671799, 'zipcode zip'), (1.2563661538447024, 'condition'),
(1.3069607649321162, 'sqft_basement'), (1.3343003838865175, 'sqft_living15'), (2.0263471204595014, 'dummy'),
(2.3136918447512111, 'sqft_lot'), (2.32284373620547l, 'view'), (2.6743741022368051, 'long'), (2.9174000919391423,
'sqft_lot15'), (2.9363434573380336, 'yr_built'), (3.331918818732531, 'bathrooms'), (3.7677253765602519, 'lat'),
(4.7110764011203488, 'waterfront'), (7.2627766152924353, 'sqft_living'), (7.7251708005852802, 'sqft_above'),
(8.8744544648249732, 'grade'), (9.1811180884566408, 'bedrooms')]
```

Figure 19. $\lambda = 0$

When the $\lambda$ change from 0 to 0.01, there is not a big change. However, when the $\lambda$ from 0.1 to 10, there are much more features whose absolute weight is less than 1. Beside, when $\lambda = 10$, some features we discussed in 2.c are becoming smaller. Like the weight of # of bedrooms, from 9 to 3. Which is more fitting on our intuition: the square of above should have more connections with the price rather than the # of bedrooms. As a result, the validation SSE of $\lambda = 10$ is much better than the validation SSE without regularization. Therefore, after adding a moderate regularization as a punishment should improve the convergence result of the algorithm.

### e. Other observations
: We also have some other observations:
First, the learning rate which get explode in part 1 also cannot get converge with adding the regularization. However, some learning rate (like 0.000057) with a relative bigger $\lambda$ will also cause an explode.
Second, after adding the regularization, the ratio between # of iteration and learning rate is still same. Like it takes around 5.7 times iterations compare LR=0.000057 and 0.00001.
Third, after adding same regularization, both training SSE and validation SSE are almost same with different learning rate.

### Part 3. Training with non-normalized data and binarization
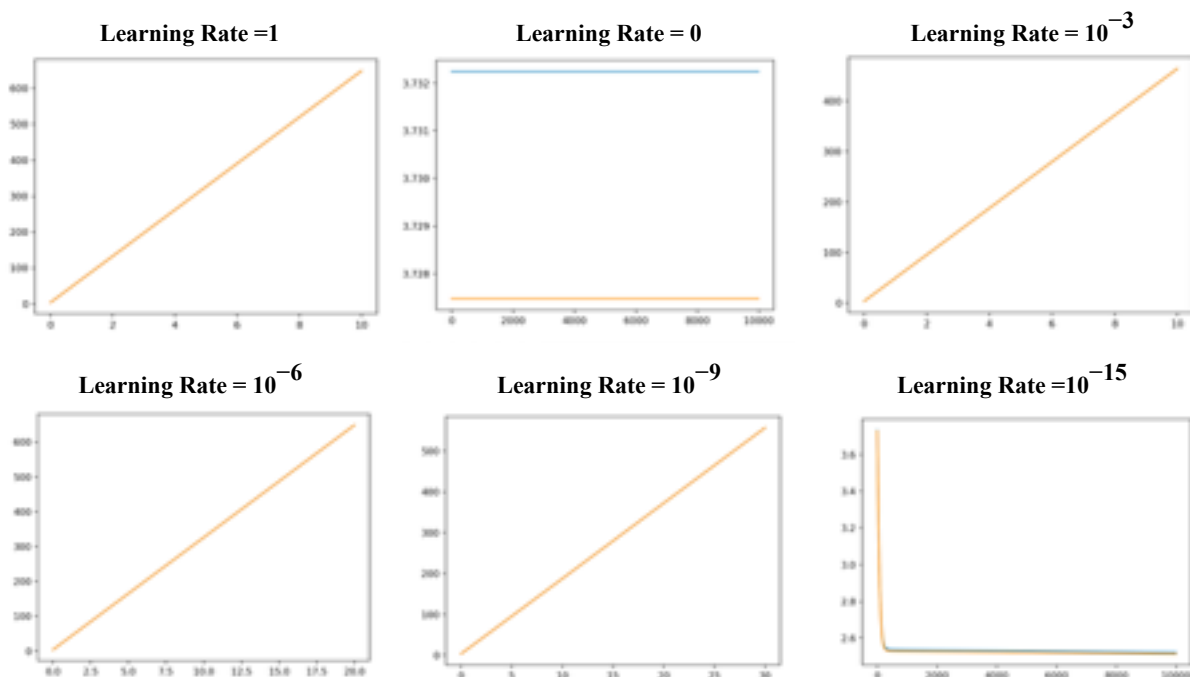### a. Observations from non-normalized data
:



Figure 20. ▬ Training Data ▬ Dev Data

Figure 20. shows the training SSE(Blue line) and validation SSE(Orange line) as a function of the number of iterations(x-rot). According to the graph, we can see that Learning Rate = 1, $10^{-3}$, $10^{-6}$, and $10^{-9}$ drives SSE both in training data and validation data explode. When learning rate is zero, SSE does not change since our program does not learn. However, when the learning rate is $10^{-15}$, which is very small, we can finally prevent the the gradient descent from exploding. As a proof, we can see the last graph in Figure 20., the only graph showing SSE decreasing with increase of iterations.

```
LR:  -1e-16
0 th iteration : |W| =  5427463725.481733
1 th iteration : |W| =  5372469891.144768
2 th iteration : |W| =  5318033765.29879
3 th iteration : |W| =  5264149695.083454
4 th iteration : |W| =  5210812084.954565
5 th iteration : |W| =  5158015396.103017
6 th iteration : |W| =  5105754145.879746
7 th iteration : |W| =  5054022907.226369
8 th iteration : |W| =  5002816308.111716
9 th iteration : |W| =  4952129030.973943
10 th iteration : |W| =  4901955812.168418
11 th iteration : |W| =  4852291441.421284
12 th iteration : |W| =  4803130761.28818
13 th iteration : |W| =  4754468666.618977
14 th iteration : |W| =  4706300104.027499
15 th iteration : |W| =  4658620071.366891
16 th iteration : |W| =  4611423617.210126
```
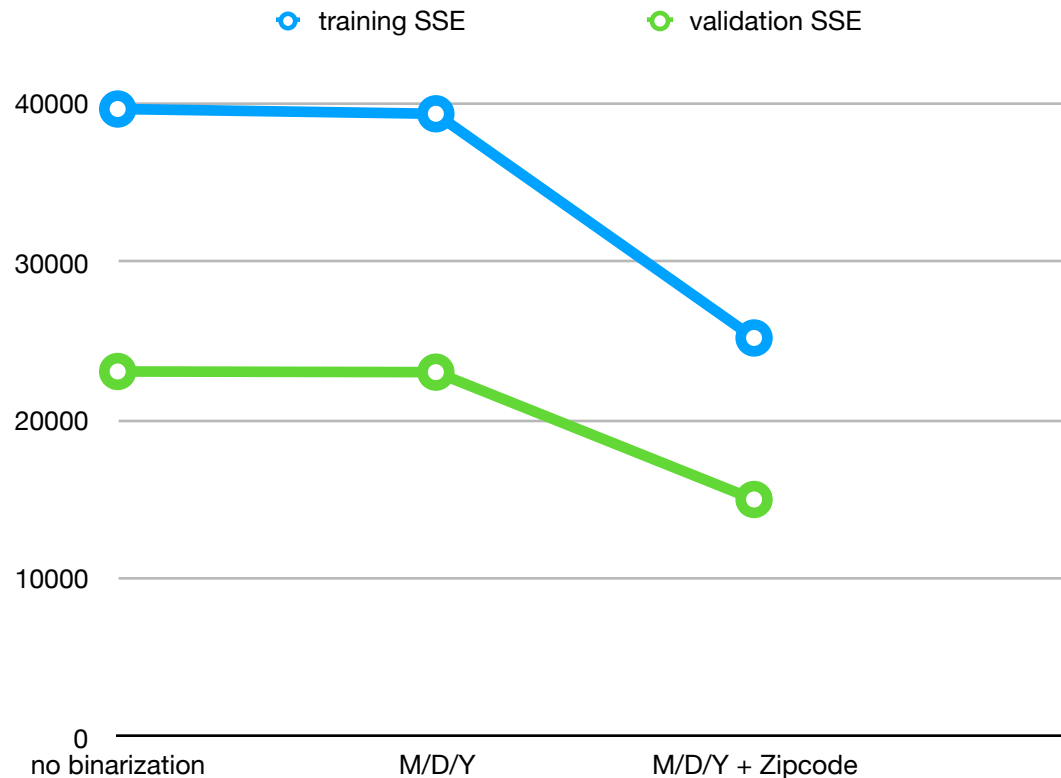
**Figure 21.**

```
[LR:  -1e-16
0 th iteration : |W| =  102640.40751457187
1 th iteration : |W| =  102640.40751422069
2 th iteration : |W| =  102640.40751386942
3 th iteration : |W| =  102640.40751351824
4 th iteration : |W| =  102640.40751316694
5 th iteration : |W| =  102640.40751281621
6 th iteration : |W| =  102640.40751246458
7 th iteration : |W| =  102640.4075121135
8 th iteration : |W| =  102640.4075117623
9 th iteration : |W| =  102640.40751141122
10 th iteration : |W| =  102640.40751106016
11 th iteration : |W| =  102640.40751070887
12 th iteration : |W| =  102640.40751035756
13 th iteration : |W| =  102640.40751000668
14 th iteration : |W| =  102640.4075096555
15 th iteration : |W| =  102640.40750930429
16 th iteration : |W| =  102640.40750895323
```

**Figure 22.**

Figure 21. shows the change of length of the gradient |∇E(w)| with the increase of integration on non-normalized data while Figure 22. presents it on normalized data. Those two figures informs that training the non-normalized data puts us in huge scale computing problem. At the 0-th iteration, the value of |∇E(w)| of non-normalized data is 5000 times bigger than the one of normalized data. This causes overflow in computing the norm of the vector, and requires smaller learning rate. For example, when we used normalized data, the best choice of our learning rate is 0.000057. However, in non-normalized data, when the learning rate reaches to $10^{-15}$, it finally started to work properly. $10^{-9}$ is not small enough for non-normalized data. We could find that normalization causes the gradient shrink more with larger learning rate than the non-normalization in the same fixed iterations. We realized that all those facts of normalization makes the training the data easier.

**b. Observations from binarized data (Extra)**
:



**Graph 2.**

       We mentioned about in part 0.b, we would like to try some binarization for the feature which doesn't have a linear positive correlation. So in this extra part, we binarize the date features into month*12 + day*31 +year*2 features. Then we found out that this binarization could increase in the result indeed. First, our criteria of the result is based on both the lowest validation SSE and the smaller gap between average training SSE and validation SSE. Then, after binarization, the training SSE from 39674.46 to 39373.04 while validation SSE from 23078.986 to 23036.46, which may not be a great improvement, but it decrease the total SSE of both training and validation data.

       Furthermore, we didn't use the category name given by professor. As our assumption, we believe that the zipcode feature should be more independent to each other zipcode. Thus, we binarize month, dat, year, zipcode in total. Then, the best training SSE drops dramatically, from 39373.04 to 25198.19 in training SSE and  23036.46 to 14985.41.

       Then we check the weight of each dimension of the w, we found out that the most important weight becomes a "zipcode" (we didn't mark which zipcode is the most important.) We assume that the house price based on zipcode between Corvallis and Silicon Valley must be different, then after binarizing the zipcode, this differences could be shown from the training process. Also, the date information also the relative small weight in w. Which may explain that

why just binarizing date information only improve a small SSE. Thus, we tested using this model to test the PA1_test.csv file. Figure 23. shows the part of the results.

| dummy | id | date | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | condition | grade | sqft_above | sqft_basement | yr_built | yr_renovated | zipcode | lat | long | sqft_living15 | sqft_lot15 | price |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Figure 23.**