

Note: the phase variable has been set to global in main(), so that other functions can use it

Values of RMSE error for three cases:

closed solution (r1): 262987.51

gradient descent (r2): 262298.23

stochastic gradient descent (r3): 262302.06

1a. $\text{abs}(r1 - r2) = 689.28$

1b. If previous RMSE error on dev set was lesser than current epoch's error, we would stop gradient descent.

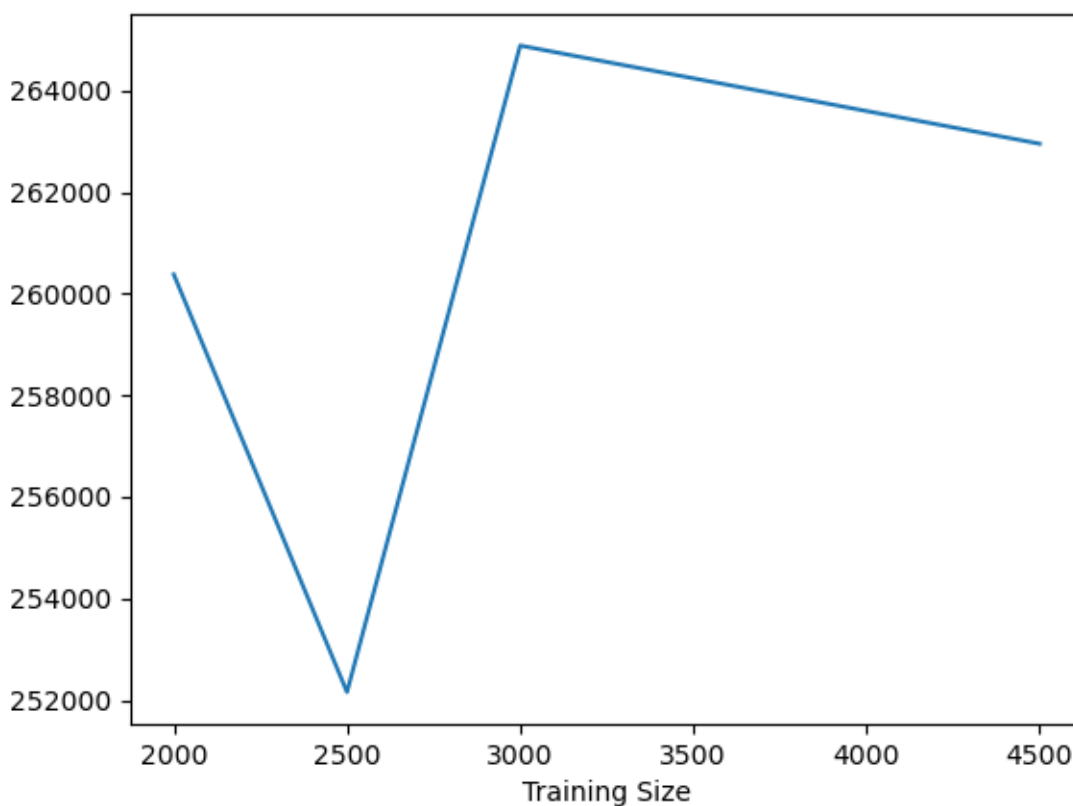
1c. $\text{abs}(r2 - r3) = 3.83$

2. pnorm2: 262955.05, pnorm4: 262299.22

3. basis: year was log scaled, since the years are all large numbers and 2010/2011 are very similar with respect to selling price. Km_driven was also log scaled because it is the scale and not the exact number that mattered. Furthermore, the km_driven filed had very large numbers hence their differences couldn't be accurately modelled linearly. Finally, the seats variable was exponentiated, in order to emphasise on the difference in number of seats.

RMSE: 259456.62

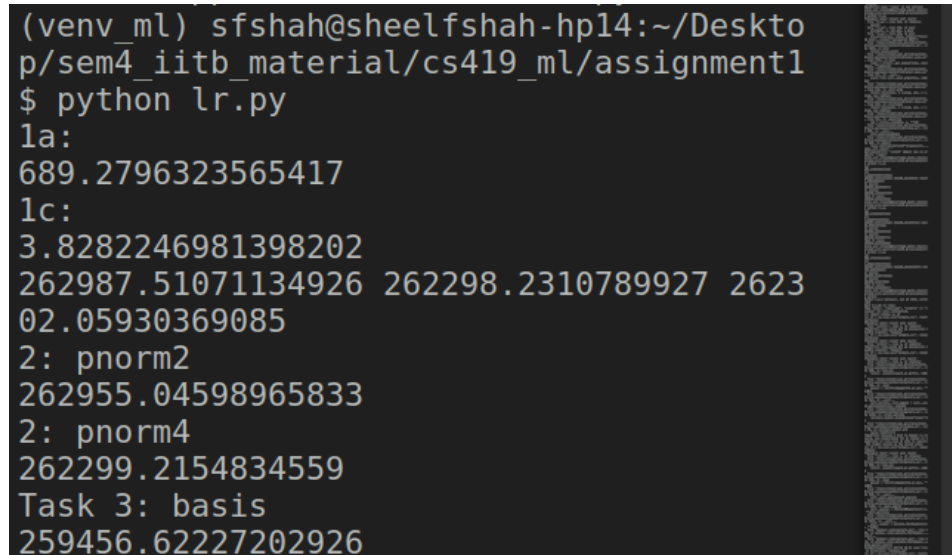
4.



5. The least useful feature is torque because it needed too much pre-processing due to no consistency in the format of the column. Hence torque had to be dropped from the dataset. After normalizing the dataset, the feature with the smallest magnitude of its corresponding weight was decided to be the second least useful. This feature was seller_type. (cumulative weight was seen for features that were one hot encoded)

6. pnorm with $p=2$ was used, along with the basis mentioned above. Also, the model tended to predict negative prices, and these were converted to their absolute value. (this isn't exactly ethical, but since the question allowed any enhancement whatsoever, I decided to go ahead with this)

Screenshot of results:



```
(venv_ml) sfshah@sheelfshah-hp14:~/Desktop/sem4_iitb_material/cs419_ml/assignment1
$ python lr.py
1a:
689.2796323565417
1c:
3.8282246981398202
262987.51071134926 262298.2310789927 2623
02.05930369085
2: pnorm2
262955.04598965833
2: pnorm4
262299.2154834559
Task 3: basis
259456.62227202926
```