# Assignment 1&2 Take2 Report

## Name: Abdelrahman Ashraf Adel

## ID: 40-5927

This is a report on the first and the second assignments after applying techniques from lecture 5 and 6 to them, while changing the data set of the first one to the new house data set acquired.

The report is divided into 4 phases, starting with the data sets description, the techniques applied to the data, analysis of the results and finally the conclusion.

# Data set description and observations:

We have two different data sets each used for different models. The first data set is the house price data-set which was utilized in the linear regression model in assignment one. The second one is the student exam scores data set used in the second assignment which applied logistic regression models.

- House price data-set is a large one composed of 18 thousand samples and 21 features. OfCourse not all the data contribute have the same contribution to the prediction of the label as each other, so to make it simpler some features were dropped after a correlation matrix was conducted to measure the correlation between all the features and the label "Price". Then any feature with a correlation less then 0.4 was dropped.

- Student exam scores data set composed of 100 samples and 3 features. At which the last column was taken as the label.

There are two approaches for dividing the data sets, the first one is to divide it to 70% training set and 30% testing set and the other approach which was applied is to divide the data to 3 sets, data for training which takes 60%, data for cross validation 20% and the last 20% is for testing.

# Techniques applied to the data:

- **Data normalization**

The two data sets were normalized before applying any techniques to them to reduce their redundancy and improve the data integrity for better results.

- **Random sampling**

For the first data set using the linear regression, six models were created differing in their degrees. Which varied from one to six. The learning rate and the number of iterations parameters were fixed for all models as the learning rate was standardized at 0.1 and the iterations for 400 same as the old assignments. In the first-degree model all the features are first degree. The second all the features are first except the sqft living feature which is second degree. The third same as the second but with the sqft above feature third degree. The fourth goes with the same rhythm and have the sqft living 15 with the $4^{th}$ degree. The fifth have the grade feature with $5^{th}$ degree. And the last with the bathroom feature in $6^{th}$ degree.

The models trained and tested with two different approaches. The first-time using shuffling only and the other using random sampling and shuffling.


For the second data set using logistic regression, three models were created also differing in their degrees from one to three. The models were trained and tested twice with the same the data house price data set. The learning rate and the number of iterations parameters were fixed for all models as the learning rate was standardized at 0.1 and the iterations for 400 same as the old assignments. In the first-degree

model all the features were first degrees. In the second the exam 1 score feature was second degree and in the last the exam 2 score feature was 3$^{rd}$ degree where exam 1 score feature same as the second.

- **K-fold sampling**

The K-fold sampling was applied on the house price data set from the first exercise with K equal 5. The data was divided accordingly into 5 sets each containing 3600 sample and was used for the last model which have the 6$^{th}$ degree.

- **Regularization**

Regularization technique was applied on both data sets utilizing the linear and logistic regression. A set of lambdas were created containing 0 0.01 0.05 0.1 0.5 1 5 10. For the linear regression different degree models were trained using each lambda and tested. The logistic regression models with different degrees were trained and tested for each different lambda from the previous set.

# Analysing the results:

## 1. First Assignment results for different techniques

This is the results concluded from applying random sampling and shuffling techniques for each degree model.

| Degree | First | Second | Third | Fourth | Fifth | Sixth |
|--------|-------|--------|-------|--------|-------|-------|
| J (CV) | 343391 53416.0 14427 | 300146 36120.3 7325 | 300053 88815.9 3281 | 299850 26013.9 78413 | 272608 67640.7 91466 | 269272 11279.5 92083 |

| | | | | | | |
|---|---|---|---|---|---|---|
| J (test) | 240975 34927.2 9127 | 210499 39081.9 42837 | 210450 26694.7 35073 | 210361 56176.4 56947 | 179428 69579.0 33318 | 177777 34998.3 18657 |

And this is the results from testing the models without random sampling.

| Degree | First | Second | Third | Fourth | Fifth | Sixth |
|---|---|---|---|---|---|---|
| J (CV) | 252441 60206.9 76276 | 220341 53443.1 12682 | 220087 55640.0 55515 | 219808 88955.4 98722 | 191544 99854.2 37453 | 189377 56243.5 15877 |
| J (test) | 243695 72452.8 06427 | 212438 35443.1 1221 | 212463 99414.9 04533 | 212384 95286.1 5241 | 181956 76777.7 72175 | 180388 03725.8 8177 |

Although the J(CV) of the above table is slightly better than those of the first table but the J(test) in both are nearly the same.

This are the results of the K-fold technique where the metric used is the cost function applied on the test set. The five sets were named S1,2,3,4 and 5. in each time one of the sets was taken as the test sample and the others were the train sets.

| Degree | First | Second | Third | Fourth | Fifth |
|---|---|---|---|---|---|
| Test Sample | S5 | S4 | S3 | S2 | S1 |
| J (CV) | 17044704401.954456 | 18624164667.758224 | 22694434694.140636 | 22350126942.5694 | 21184492467.779766 |

And those are the results from the regularization technique applied on different lambdas and different degrees using CV set.

| Degree /Lambdas | 0 | 0.01 | 0.05 | 0.1 | 0.5 | 1 | 5 | 10 |
|---|---|---|---|---|---|---|---|---|
| 1st | 252441602008.976276 | 5848471707166.148 | 29141161202905.555 | 58256526523380.75 | 2911595886874 38.1 | 5822387729 92865.5 | 290888 8078008638.5 | 5812246865343023.0 |
| 2nd | 220341534442.112682 | 6102080724821.959 | 30422030633811.633 | 60821436181703.68 | 3039954087 64866.1 | 6079097033 96084.8 | 303709 8948453307.0 | 6068281103914253.0 |
| 3rd | 220087556339.055515 | 6095689893653.899 | 30390176172410.902 | 60757747915152.914 | 3036768794 37048.7 | 6072721965 51472.8 | 303389 2647925022.0 | 6061821734392028.0 |
| 4th | 219808889536.498722 | 6101609367485.173 | 30419882668769.473 | 60817182926152.24 | 3039739321 06684.56 | 6078657455 46577.6 | 303683 7219570326.5 | 6067653109437821.0 |

| 5th | 19154<br>49985<br>9.2374<br>53 | 606964<br>273941<br>1.186 | 302713<br>663575<br>79.156 | 605230<br>048735<br>09.945 | 30251<br>54743<br>58475.<br>3 | 60495<br>44725<br>01224.<br>8 | 302240<br>453606<br>1575.5 | 6039070183<br>379429.0 |
|---|---|---|---|---|---|---|---|---|
| 6th | 18937<br>75624<br>7.5158<br>77 | 607706<br>351895<br>9.113 | 303093<br>372916<br>46.06 | 605991<br>636398<br>35.836 | 30289<br>71413<br>23682.<br>6 | 60571<br>80384<br>52755.<br>8 | 302622<br>381993<br>2054.0 | 6046710316<br>890518.0 |

## 2. <u>Second Assignment results for different techniques</u>

The table below shows the results of applying random sampling and shuffling on different degrees.

| Degree | First | Second | Third |
|---|---|---|---|
| J (CV) | 0.255875466422<br>8678 | 0.141146503130<br>7492 | 0.302652009371<br>9782 |
| J (test) | 0.274209620903<br>4244 | 0.092884700672<br>52748 | 0.351267954883<br>88396 |

And these are the results of applying the same technique without the random sampling on the same degrees.

| Degree | First | Second | Third |
|---|---|---|---|
| J (CV) | 0.167088457431<br>0742 | 0.157088457431<br>0742 | 0.350919421007<br>4731 |
| J (test) | 0.252445293885<br>6521 | 0.202445293885<br>6521 | 0.298893765009<br>10814 |

The last results in this exercise are the ones conducted from the regularization with different lambda values for different degrees using CV set.

| Degree /lambda | 0 | 0.01 | 0.05 | 0.1 | 0.5 | 1 | 5 | 10 |
|---|---|---|---|---|---|---|---|---|
| 1st | 0.28240653030885154 | 0.3095196628337426 | 0.41797218289175525 | 0.55353782420540411 | 1.638062429269027 | 2.9937169761457704 | 13.838904863962645 | 27.395267770469196 |
| 2nd | 0.31436472490637331 | 0.3143839703964817 | 0.31445565072230036 | 0.31454959879295929 | 0.311108740902115554 | 0.31761497658318516 | 0.325502998046254655 | 0.34461227507681336 |
| 3rd | 0.30973375730003858 | 0.3109055290958805 | 0.313412417788787377 | 0.316570237474968433 | 0.313298022571011497 | 0.32443726581422067 | 0.3283905505620721 | 0.3350141274674089 |

# Conclusion:

After applying the different techniques mentioned above to the two different sets of house prices and the student exam score sets using linear and logistic regression. The random sampling technique didn't improve the prediction accuracy, but in the logistic regression the random sampling did contribute and improved the results a little. Although the CV test results changed by changing between the sampling techniques using random and without the random sampling but the J test results didn't change.For the regularization technique the first lambda equal to zero resulted in the best accuracy for both regression models.