# HARVARD EXTENSION SCHOOL

EXT CSCI E-106 Model Data Class Group Project Template

Author Kaleo Pudim        Author Luciano Carvalho        Author Ibrahim Hashim
Author Bethun Bhowmik      Author Mohanish Kashiwar       Author Seymur Hasanov

22 November 2023

**Abstract**

This is the location for your abstract. It must consist of two paragraphs.

# Contents

# House Sales in King County, USA data to be used in the Final Project

| Variable | Description |
|---|---|
| id | **Unique ID for each home sold (it is not a predictor)** |
| date | *Date of the home sale* |
| price | *Price of each home sold* |
| bedrooms | *Number of bedrooms* |
| bathrooms | *Number of bathrooms, where ".5" accounts for a bathroom with a toilet but no shower* |
| sqft_living | *Square footage of the apartment interior living space* |
| sqft_lot | *Square footage of the land space* |
| floors | *Number of floors* |
| waterfront | *A dummy variable for whether the apartment was overlooking the waterfront or not* |
| view | *An index from 0 to 4 of how good the view of the property was* |
| condition | *An index from 1 to 5 on the condition of the apartment,* |
| grade | *An index from 1 to 13, where 1-3 falls short of building construction and design, 7 has an average level of construction and design, and 11-13 has a high-quality level of construction and design.* |
| sqft_above | *The square footage of the interior housing space that is above ground level* |
| sqft_basement | *The square footage of the interior housing space that is below ground level* |
| yr_built | *The year the house was initially built* |
| yr_renovated | *The year of the house's last renovation* |
| zipcode | *What zipcode area the house is in* |
| lat | *Latitude* |
| long | *Longitude* |
| sqft_living15 | *The square footage of interior housing living space for the nearest 15 neighbors* |
| sqft_lot15 | *The square footage of the land lots of the nearest 15 neighbors* |

# Instructions:

0. Join a team with your fellow students with appropriate size (Four Students total)
1. Load and Review the dataset named "KC_House_Sales'csv
2. Create the train data set which contains 70% of the data and use set.seed (1023). The remaining 30% will be your test data set.
3. Investigate the data and combine the level of categorical variables if needed and drop variables as needed. For example, you can drop id, Latitude, Longitude, etc.
4. Build a regression model to predict price.
5. Create scatter plots and a correlation matrix for the train data set. Interpret the possible relationship between the response.
6. Build the best multiple linear models by using the stepwise selection method. Compare the performance of the best two linear models.
7. Make sure that model assumption(s) are checked for the final model. Apply remedy measures (transformation, etc.) that helps satisfy the assumptions.
8. Investigate unequal variances and multicollinearity. If necessary, apply remedial methods (WLS, Ridge, Elastic Net, Lasso, etc.).
9. Build an alternative model based on one of the following approaches to predict price: regression tree, NN, or SVM. Check the applicable model assumptions. Explore using a logistic regression.
10. Use the test data set to assess the model performances from above.
11. Based on the performances on both train and test data sets, determine your primary (champion) model and the other model which would be your benchmark model.
12. Create a model development document that describes the model following this template, input the name of the authors, Harvard IDs, the name of the Group, all of your code and calculations, etc..:

## Due Date: December 18th, 2023 at 11:59 pm EST

**Notes No typographical errors, grammar mistakes, or misspelled words, use English language All tables need to be numbered and describe their content in the body of the document All figures/graphs need to be numbered and describe their content All results must be accurate and clearly explained for a casual reviewer to fully understand their purpose and impact Submit both the RMD markdown file and PDF with the sections with appropriate explanations. A more formal document in Word can be used in place of the pdf file but must include all appropriate explanations.**

Executive Summary

This section will describe the model usage, your conclusions and any regulatory and internal requirements. In a real world scneario, this section is for senior management who do not need to know the details. They need to know high level (the purpose of the model, limitations of the model and any issues).

# I. Introduction (5 points)

*This section needs to introduce the reader to the problem to be resolved, the purpose, and the scope of the statistical testing applied. What you are doing with your prediction? What is the purpose of the model? What methods were trained on the data, how large is the test sample, and how did you build the model?*

— in-progress —

This project aims to develop a predictive model for house prices in King County, USA. The model's purpose is to provide an analytical tool for understanding the key factors influencing property values in this area, which is significant for various stakeholders in the real estate sector.

Our primary data source is the 'KC_House_Sales' dataset, which includes detailed information on various house attributes. **(Add here the hypothesize of initial exploration on house prices)**.

The methodology adopted for this project combines conventional statistical techniques with modern data analytics methods. Initial models will be based on linear regression, with further exploration into more complex approaches like neural networks and decision trees, depending on their performance in preliminary analyses.

To ensure robust model training and validation, the dataset has been divided into a training set (70%) and a test set (30%). This division is crucial for evaluating the model's effectiveness in making predictions on new, unseen data.

# I. Description of the data and quality (15 points)

*Here you need to review your data, the statistical test applied to understand the predictors and the response and how are they correlated. Extensive graph analysis is recommended. Is the data continuous, or categorical, do any transformation needed? Do you need dummies?*

**Data Overview**

```
df_house = read.csv("KC_House_Sales.csv")

head(df_house)
```

```
##           id            date         price bedrooms bathrooms sqft_living
## 1 7129300520 20141013T000000    $221,900.00        3      1.00        1180
## 2 6414100192 20141209T000000    $538,000.00        3      2.25        2570
## 3 5631500400 20150225T000000    $180,000.00        2      1.00         770
## 4 2487200875 20141209T000000    $604,000.00        4      3.00        1960
## 5 1954400510 20150218T000000    $510,000.00        3      2.00        1680
## 6 7237550310 20140512T000000  $1,225,000.00        4      4.50        5420
##   sqft_lot floors waterfront view condition grade sqft_above sqft_basement
## 1     5650      1          0    0         3     7       1180             0
## 2     7242      2          0    0         3     7       2170           400
## 3    10000      1          0    0         3     6        770             0
## 4     5000      1          0    0         5     7       1050           910
## 5     8080      1          0    0         3     8       1680             0
## 6   101930      1          0    0         3    11       3890          1530
##   yr_built yr_renovated zipcode     lat     long sqft_living15 sqft_lot15
## 1     1955            0   98178 47.5112 -122.257          1340       5650
## 2     1951         1991   98125 47.7210 -122.319          1690       7639
## 3     1933            0   98028 47.7379 -122.233          2720       8062
## 4     1965            0   98136 47.5208 -122.393          1360       5000
## 5     1987            0   98074 47.6168 -122.045          1800       7503
## 6     2001            0   98053 47.6561 -122.005          4760     101930
```

**Data Types, Categories and Cleaning**

```
#removing `id` column
df_house = subset(df_house, select = -id)

#converting `price` to numeric
df_house$price <- as.numeric(gsub("[\\$,]", "", df_house$price))

##cleaning `date` and adding `year`,`month`,`day` columns
df_house$date <- as.POSIXct(df_house$date, format = "%Y%m%d")
df_house$year <- as.numeric(format(df_house$date, "%Y"))
df_house$month <- as.numeric(format(df_house$date, "%m"))
df_house$day <- as.numeric(format(df_house$date, "%d"))

head(df_house)
```

```
##         date  price bedrooms bathrooms sqft_living sqft_lot floors waterfront
## 1 2014-10-13 221900        3      1.00        1180     5650      1          0
## 2 2014-12-09 538000        3      2.25        2570     7242      2          0
## 3 2015-02-25 180000        2      1.00         770    10000      1          0
## 4 2014-12-09 604000        4      3.00        1960     5000      1          0
```

```
## 5 2015-02-18  510000        3     2.00        1680      8080        1          0
## 6 2014-05-12 1225000        4     4.50        5420    101930        1          0
##   view condition grade sqft_above sqft_basement yr_built yr_renovated zipcode
## 1    0         3     7       1180             0     1955            0   98178
## 2    0         3     7       2170           400     1951         1991   98125
## 3    0         3     6        770             0     1933            0   98028
## 4    0         5     7       1050           910     1965            0   98136
## 5    0         3     8       1680             0     1987            0   98074
## 6    0         3    11       3890          1530     2001            0   98053
##        lat     long sqft_living15 sqft_lot15 year month day
## 1 47.5112 -122.257          1340       5650 2014    10  13
## 2 47.7210 -122.319          1690       7639 2014    12   9
## 3 47.7379 -122.233          2720       8062 2015     2  25
## 4 47.5208 -122.393          1360       5000 2014    12   9
## 5 47.6168 -122.045          1800       7503 2015     2  18
## 6 47.6561 -122.005          4760     101930 2014     5  12
```

```r
#converting all to numeric...ignores date column
ndf_house = df_house[sapply(df_house, is.numeric)]

head(ndf_house)
```

```
##      price bedrooms bathrooms sqft_living sqft_lot floors waterfront view
## 1  221900        3      1.00        1180     5650      1          0    0
## 2  538000        3      2.25        2570     7242      2          0    0
## 3  180000        2      1.00         770    10000      1          0    0
## 4  604000        4      3.00        1960     5000      1          0    0
## 5  510000        3      2.00        1680     8080      1          0    0
## 6 1225000        4      4.50        5420   101930      1          0    0
##   condition grade sqft_above sqft_basement yr_built yr_renovated zipcode
## 1         3     7       1180             0     1955            0   98178
## 2         3     7       2170           400     1951         1991   98125
## 3         3     6        770             0     1933            0   98028
## 4         5     7       1050           910     1965            0   98136
## 5         3     8       1680             0     1987            0   98074
## 6         3    11       3890          1530     2001            0   98053
##        lat     long sqft_living15 sqft_lot15 year month day
## 1 47.5112 -122.257          1340       5650 2014    10  13
## 2 47.7210 -122.319          1690       7639 2014    12   9
## 3 47.7379 -122.233          2720       8062 2015     2  25
## 4 47.5208 -122.393          1360       5000 2014    12   9
## 5 47.6168 -122.045          1800       7503 2015     2  18
## 6 47.6561 -122.005          4760     101930 2014     5  12
```

**Stat Summary**

**Graphical Analysis**

**Correlation Analysis**

**Data Transformation**

**Dummy Variables**

**Summary of the intro section**

## III. Model Development Process (15 points)

*Build a regression model to predict price. And of course, create the train data set which contains 70% of the data and use set.seed (1023). The remaining 30% will be your test data set. Investigate the data and combine the level of categorical variables if needed and drop variables. For example, you can drop id, Latitude, Longitude, etc.*

## IV. Model Performance Testing (15 points)

*Use the test data set to assess the model performances. Here, build the best multiple linear models by using the stepwise both ways selection method. Compare the performance of the best two linear models. Make sure that model assumption(s) are checked for the final linear model. Apply remedy measures (transformation, etc.) that helps satisfy the assumptions. In particular you must deeply investigate unequal variances and multicollinearity. If necessary, apply remedial methods (WLS, Ridge, Elastic Net, Lasso, etc.).*

## V. Challenger Models (15 points)

*Build an alternative model based on one of the following approaches to predict price: regression tree, NN, or SVM. Explore using a logistic regression. Check the applicable model assumptions. Apply in-sample and out-of-sample testing, backtesting and review the comparative goodness of fit of the candidate models. Describe step by step your procedure to get to the best model and why you believe it is fit for purpose.*

## VI. Model Limitation and Assumptions (15 points)

*Based on the performances on both train and test data sets, determine your primary (champion) model and the other model which would be your benchmark model. Validate your models using the test sample. Do the residuals look normal? Does it matter given your technique? How is the prediction performance using Pseudo R^2, SSE, RMSE? Benchmark the model against alternatives. How good is the relative fit? Are there any serious violations of the model assumptions? Has the model had issues or limitations that the user must know? (Which assumptions are needed to support the Champion model?)*

## VII. Ongoing Model Monitoring Plan (5 points)

*How would you picture the model needing to be monitored, which quantitative thresholds and triggers would you set to decide when the model needs to be replaced? What are the assumptions that the model must comply with for its continuous use?*

## VIII. Conclusion (5 points)

*Summarize your results here. What is the best model for the data and why?*

## Bibliography (7 points)

*Please include all references, articles and papers in this section.*

## Appendix (3 points)

*Please add any additional supporting graphs, plots and data analysis.*