# Part-I-Writeup

*Jingyi Zhang, Jonathan Klus, Bin Han*

## 1. Introduction:

In this study, we are looking at the auction prices of paintings in 18th century Paris. Specifically, through the assistance of model built based on existing training data, we wish to understand the factors that drive the prices of the paintings, and then be able to predict auction prices based on characteristics of a certain painting. After fitting appropriate model, we also intend to detect specific paintings that are either underpriced or overpriced based on the selected model.

One of the main task and challenge is to narrow down the number of potential predictors from 59 to less than 20 while maintaining a high performance of the model. But being able to explain the results and provide some recommendations to indivisuals without statistical background is equally important and challenging. Therefore, we aim at balancing the performance of model prediction, closeness to true model, simplicity, and interprebility.
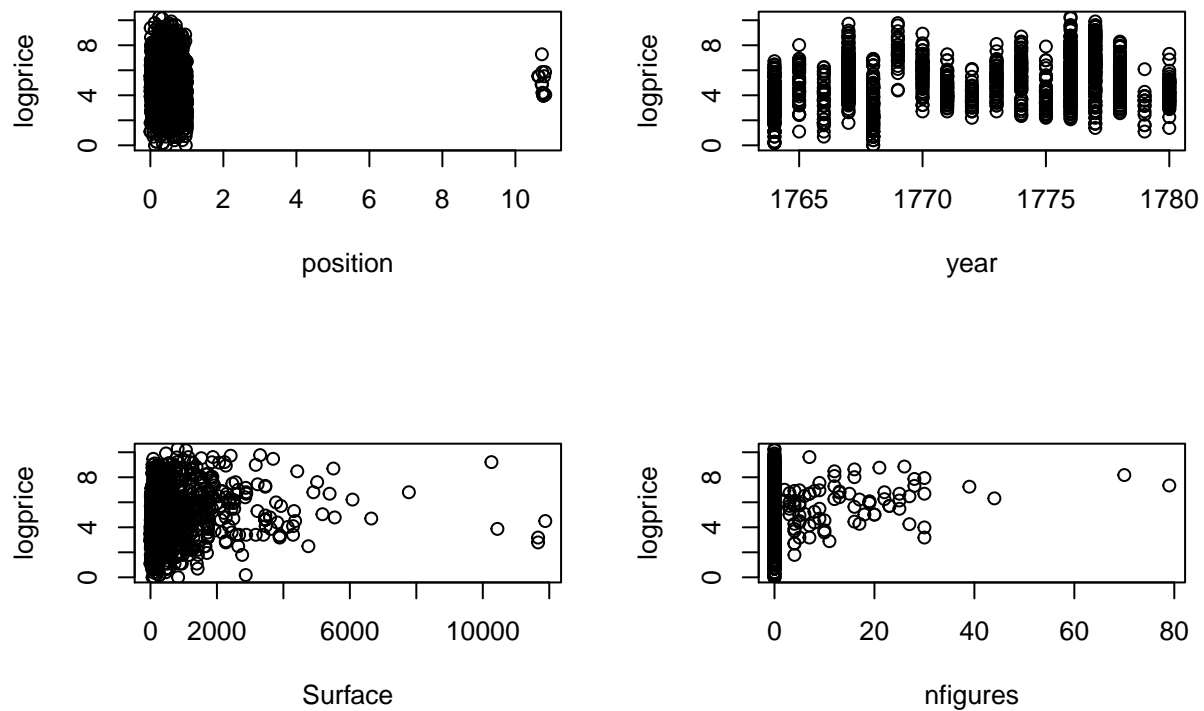
## 2. Exploratory data analysis:

## A) Data summary & cleaning

To start with, we looked at the summary of the original trainig data. There are few numeric variables and a lot of binary variables. Some variables, such as `Interm`, `Surface`, `Height_in` etc. have mising values, which need to be taken care of. The followings steps are how we cleaned the data:

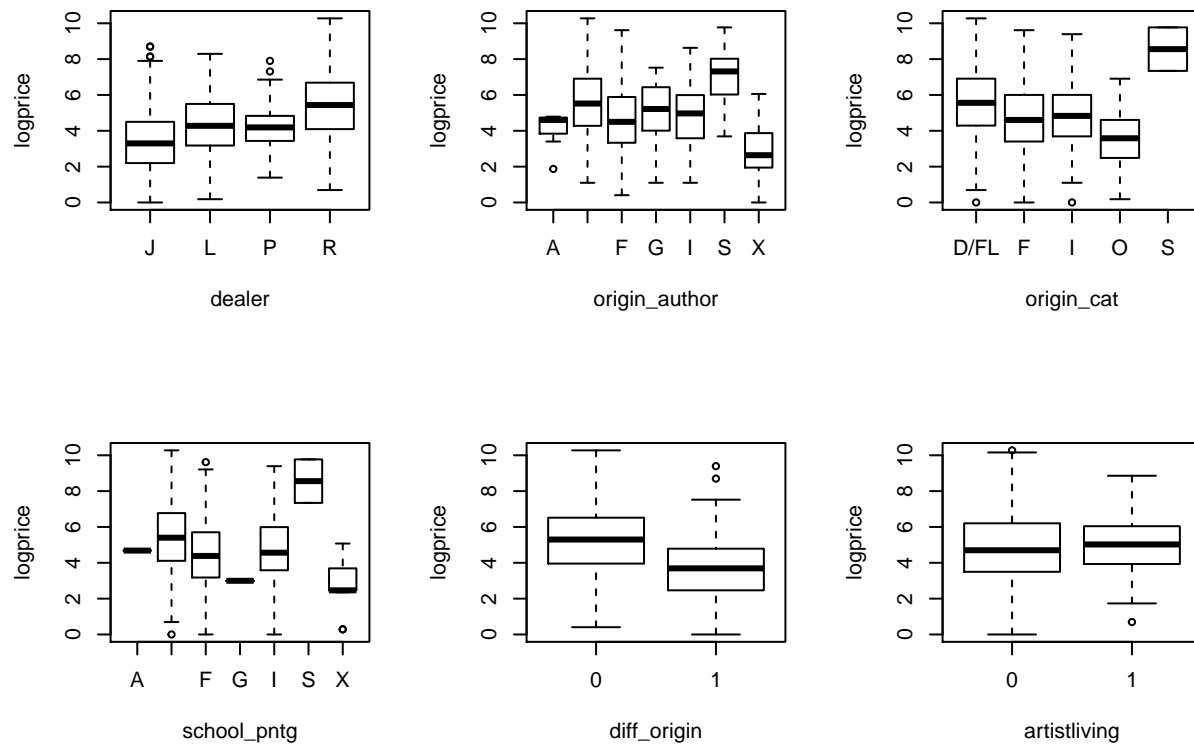a. The first step we did was to get rid of intuitivelly useless variables, including: `lot`, `sale`, `price`, `count`, `subject`, `authorstandard`, `author`, `winningbidder`, and 'other. The are not useful in predicting the response variable (such as names)

b. By further screening the variables, we found out that `Surface` and `Surface_Rnd`, `Surface_Rect` are corerlated, which are based on the value of `Height_in`, `Width_in`, and `Diam_in`. We decided to use `Surface` in our initial model. The same issue happened to `material`, `mat`, and `materialCat`. The latter one recodes the previous one. Therefore, we used `materialCat`. We applied the same strategy to keep `landsALL` and get rid of other variables related with landscape.

c. For those variables that have multiple levels, to be consistent with how the data was originally coded, we recoded the missing levels as "X", which stands for "no information". For `materialCat` and `Shape`, since there are so many levels, we grouped some levels with few observations together, coded as "other" group. The rest binary vairables are changed into factor.

d. Then we dealt with the missing values in `Surface` and `Interm`. We used the package "mice" to address this problem, which uses the observed values in the dataset to impute the missing values. It prevents directly throwing away the missing values, which results in lossing a large amont of information for prediction.

## B). Plots
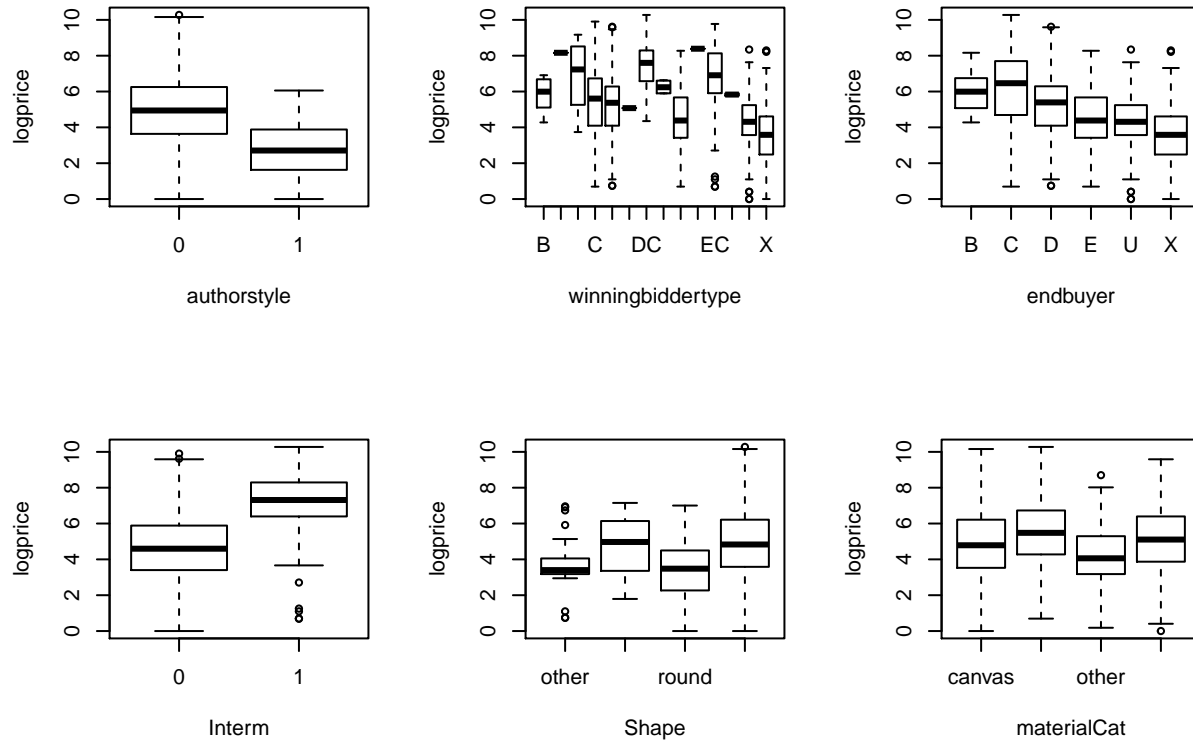
Then we analyed the relationship between those left features and the response variable. With the scatter plots, we can roughly determine which variables can be put into the initial model. For categorical variables, we want to check if the `logprice` spans different ranges in different levels. For numeric variables, we want to check if there is a clear relationship between them and `logprice`.

For numeric variables, we see that `Surface` and `nfigures` seem to show some weak but positive relationship with `logprice`. Since there are several extremely large values in `position` (potentially outliers), it is hard to see that real pattern between the majority of points and `logprice`. But we'll keep it in the model first.

Since there are 33 categorical variables, we don't show the boxplots for all of them. But applied the same method to check all the categorical variables. The following variables show some differences in `logprice` at different levels (not considering the magnitude of the difference at this time): `dealer`, `origin_author`, `origin_cat`, `school_pntg`, `diff_origin`, `authorstyle`, `endbuyer`, `Interm`, `Shape`, `materialCat`, `engraved`, `prevcoll`, `figures`, `finished`, `Irgfont`, `othgenre`, `discauth`, and `still_life`.

If we were to choose 10 best predictive variables for predicting, we would consider the magnitude of differences and the strength of relationships. The 10 variables we choose are: `Surface`, `dealer`, `school_pntg`, `diff_origin`, `authorstyle`, `endbuyer`, `Interm`, `prevcoll`, `engraved`, `Irgfont`.

## 3. Development and assessment of an initial model

## Initial Model

### JZS prior

```
##  [1] "Intercept"           "dealerL"              "dealerR"
##  [4] "year"                "origin_authorS"       "origin_authorX"
##  [7] "school_pntgD/FL"     "diff_origin1"         "artistliving1"
## [10] "authorstyle1"        "winningbiddertypeC"   "winningbiddertypeEC"
## [13] "winningbiddertypeU"  "winningbiddertypeX"   "endbuyerE"
## [16] "Interm1"             "Shaperound"           "Surface"
## [19] "materialCatother"    "materialCatwood"      "engraved1"
## [22] "prevcoll1"           "paired1"              "figures1"
## [25] "finished1"           "lrgfont1"             "portrait1"
## [28] "still_life1"         "discauth1"
```

3

**g-prior**

```
## [1] "Intercept"         "dealerL"           "dealerR"
## [4] "year"              "school_pntgD/FL"   "diff_origin1"
## [7] "artistliving1"     "authorstyle1"      "winningbiddertypeD"
## [10] "winningbiddertypeU" "endbuyerE"        "endbuyerX"
## [13] "Interm1"           "Shaperound"        "Surface"
## [16] "materialCatother"  "engraved1"         "prevcoll1"
## [19] "paired1"           "finished1"         "lrgfont1"
## [22] "portrait1"         "still_life1"       "discauth1"
```

The EDA process gives us an initial idea of which variables to drop out to reduce the dimension, and which variables might be significant in explaining the variation in logprice. But before we built the initial model, we applied BMA, Bayesian Model Averaging, to systemetically choose which base variables that have higher posterior probabilities to be in the initial model. We experimented two modelpriors, "JZS" and "g-prior", which gave us two sets of variables listed above. Then we picked up the common ones from Best Predictive Model(BPM).

Then we fit the linear regression model using the chosen features and all their possible interactions. From the summary table, the $R^2 = 0.5828$, which is fairly high. But we realized that lots of estimated coefficients for interactions are NAs, indicating that some levels in those variables have too few observations to be estimated. Therefore, we need to further reduce the dimention through variable selection.

```
##
## Call:
## lm(formula = logprice ~ dealer + school_pntg + diff_origin +
##     artistliving + endbuyer + authorstyle + Interm + Shape +
##     Surface + engraved + prevcoll + paired + finished + lrgfont +
##     portrait + discauth + still_life, data = paintings_train_2)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -4.3813 -0.7706  0.0252  0.7926  4.7942
##
## Coefficients:
##                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)     4.625e+00  1.336e+00   3.461 0.000553 ***
## dealerL         1.531e+00  1.346e-01  11.373  < 2e-16 ***
## dealerP         5.210e-01  1.650e-01   3.157 0.001626 **
## dealerR         1.176e+00  1.079e-01  10.894  < 2e-16 ***
## school_pntgD/FL -6.130e-01  1.258e+00  -0.487 0.626263
## school_pntgF    -1.364e+00  1.259e+00  -1.083 0.278843
## school_pntgG    -3.000e+00  1.773e+00  -1.692 0.090822 .
## school_pntgI    -1.273e+00  1.261e+00  -1.009 0.313224
## school_pntgS    -3.413e-01  1.548e+00  -0.221 0.825489
## school_pntgX    -1.953e+00  1.272e+00  -1.535 0.125033
## diff_origin1    -6.491e-01  9.138e-02  -7.103 1.89e-12 ***
## artistliving1    6.772e-01  1.049e-01   6.456 1.46e-10 ***
## endbuyerC       -1.379e-01  3.440e-01  -0.401 0.688626
## endbuyerD       -1.299e-01  3.410e-01  -0.381 0.703257
## endbuyerE       -7.968e-01  3.549e-01  -2.245 0.024910 *
## endbuyerU       -6.208e-01  3.507e-01  -1.770 0.076897 .
## endbuyerX       -1.305e+00  3.485e-01  -3.744 0.000188 ***
## authorstyle1    -1.075e+00  1.573e-01  -6.834 1.21e-11 ***
## Interm1          1.011e+00  1.400e-01   7.218 8.42e-13 ***
```

```
## Shapeoval          4.408e-01  3.858e-01   1.143 0.253346
## Shaperound        -3.799e-01  3.522e-01  -1.079 0.280889
## Shapesqu_rect      5.856e-01  2.576e-01   2.274 0.023138 *
## Surface            2.287e-04  3.353e-05   6.820 1.33e-11 ***
## engraved1          3.208e-01  1.514e-01   2.119 0.034264 *
## prevcoll1          1.212e+00  1.501e-01   8.071 1.43e-15 ***
## paired1           -3.377e-01  7.072e-02  -4.776 1.97e-06 ***
## finished1          6.293e-01  9.662e-02   6.514 1.01e-10 ***
## lrgfont1           1.088e+00  1.244e-01   8.744  < 2e-16 ***
## portrait1         -6.465e-01  1.762e-01  -3.670 0.000251 ***
## discauth1          3.747e-01  1.447e-01   2.588 0.009736 **
## still_life1       -7.202e-01  1.712e-01  -4.208 2.74e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.247 on 1469 degrees of freedom
## Multiple R-squared:  0.5858, Adjusted R-squared:  0.5773
## F-statistic: 69.25 on 30 and 1469 DF,  p-value: < 2.2e-16
```

## Model Selection

After completing the initial exploratory data analysis, methods including Stepwise Best Subset Selection using both AIC and BIC were used in order to assess more systematically which covariates were most important for predicting the logprice of paintings. While the number of relevant covariates was initially thinned by examining the data and determining which variables were best suited for modeling (e.g. via dimension reduction, elimination or recoding of categorical variables with too many levels or too few observations for a given level to be useful in estimating a coefficient), there still remained a large number of covariates from which to choose. The goal in using the above described methodology was to demonstrate among several methods, both frequentist and Bayesian, which covariates were routinely deemed to be the most important for modeling logprice.

The variable selection methods described above remain computationally intensive, particularly given the number of variables and potential two-way interactions that must be considered. In order to begin the analysis, The two-way interactions were considered using stepwise selection (AIC & BIC). The goal of this penalized selection process was to avoid overfitting and to deliver a model that was both interpretable and performed well in prediction. Then we compared the results from two methods and filtered out interactions that have NAs as coefficients, that are not significant, and that do not make sense to be interacted (such as $artistling * endbuyer$).

Ultimately, the following variables were selected using the above methods and were fit using OLS regression. The $R^2$ reduces to 0.6315, which is expected. All the estimated coefficients do not contain NAs.

```
##
## Call:
## lm(formula = logprice ~ Shape + school_pntg + dealer * Interm +
##     dealer * Surface + dealer * paired + dealer * finished +
##     dealer * discauth + diff_origin * Surface + diff_origin *
##     portrait + artistliving * endbuyer + artistliving * authorstyle +
##     Interm * Surface + Interm * lrgfont + Surface * lrgfont +
##     Surface * still_life + Surface * discauth + prevcoll * finished +
##     paired * lrgfont + paired * discauth + diff_origin * authorstyle +
##     diff_origin * still_life + finished * discauth + lrgfont *
##     discauth + artistliving * finished + Interm * portrait +
##     dealer * artistliving + authorstyle * prevcoll, data = paintings_train_2)
##
```

```
## Residuals:
##     Min      1Q  Median      3Q     Max
## -4.0704 -0.7182  0.0362  0.7453  4.9447
##
## Coefficients:
##                         Estimate Std. Error t value Pr(>|t|)
## (Intercept)            4.079e+00  1.313e+00   3.107 0.001929 **
## Shapeoval              5.270e-01  3.790e-01   1.391 0.164549
## Shaperound            -4.849e-01  3.470e-01  -1.397 0.162492
## Shapesqu_rect          5.221e-01  2.585e-01   2.019 0.043641 *
## school_pntgD/FL       -6.293e-01  1.220e+00  -0.516 0.606085
## school_pntgF          -1.319e+00  1.221e+00  -1.080 0.280330
## school_pntgG          -3.098e+00  1.712e+00  -1.809 0.070667 .
## school_pntgI          -1.260e+00  1.224e+00  -1.030 0.303359
## school_pntgS          -7.612e-01  1.514e+00  -0.503 0.615257
## school_pntgX          -1.855e+00  1.235e+00  -1.501 0.133501
## dealerL                2.658e+00  2.168e-01  12.259  < 2e-16 ***
## dealerP                1.347e+00  2.868e-01   4.697 2.89e-06 ***
## dealerR                1.875e+00  1.868e-01  10.034  < 2e-16 ***
## Interm1               -1.137e+00  4.806e-01  -2.365 0.018145 *
## Surface                8.925e-04  1.963e-04   4.546 5.94e-06 ***
## paired1                1.968e-01  1.917e-01   1.027 0.304787
## finished1              1.169e+00  2.219e-01   5.266 1.60e-07 ***
## discauth1              9.803e-01  3.892e-01   2.519 0.011891 *
## diff_origin1          -5.455e-01  1.062e-01  -5.135 3.20e-07 ***
## portrait1             -9.041e-01  2.054e-01  -4.402 1.15e-05 ***
## artistliving1         -1.265e-01  8.358e-01  -0.151 0.879742
## endbuyerC             -3.360e-01  3.735e-01  -0.899 0.368569
## endbuyerD             -3.344e-01  3.693e-01  -0.906 0.365238
## endbuyerE             -1.106e+00  3.822e-01  -2.893 0.003874 **
## endbuyerU             -8.277e-01  3.801e-01  -2.178 0.029576 *
## endbuyerX             -1.583e+00  3.777e-01  -4.192 2.94e-05 ***
## authorstyle1          -1.766e+00  4.213e-01  -4.193 2.92e-05 ***
## lrgfont1               1.779e+00  1.762e-01  10.100  < 2e-16 ***
## still_life1           -1.728e-01  2.422e-01  -0.713 0.475821
## prevcoll1              1.519e+00  1.704e-01   8.912  < 2e-16 ***
## dealerL:Interm1        1.413e+00  8.523e-01   1.658 0.097582 .
## dealerP:Interm1        1.803e+00  1.345e+00   1.341 0.180099
## dealerR:Interm1        2.410e+00  4.942e-01   4.876 1.20e-06 ***
## dealerL:Surface       -6.209e-04  2.062e-04  -3.012 0.002644 **
## dealerP:Surface       -4.019e-04  3.426e-04  -1.173 0.240852
## dealerR:Surface       -6.054e-04  1.984e-04  -3.051 0.002325 **
## dealerL:paired1       -1.202e+00  2.553e-01  -4.711 2.71e-06 ***
## dealerP:paired1       -6.826e-01  3.558e-01  -1.918 0.055273 .
## dealerR:paired1       -3.008e-01  2.101e-01  -1.432 0.152505
## dealerL:finished1     -6.222e-01  4.984e-01  -1.248 0.212108
## dealerP:finished1     -7.051e-01  4.080e-01  -1.728 0.084140 .
## dealerR:finished1     -5.962e-01  2.421e-01  -2.463 0.013901 *
## dealerL:discauth1      2.511e-02  9.195e-01   0.027 0.978216
## dealerP:discauth1     -1.755e+00  9.475e-01  -1.852 0.064254 .
## dealerR:discauth1     -9.142e-01  3.611e-01  -2.532 0.011458 *
## Surface:diff_origin1  -1.506e-04  9.739e-05  -1.546 0.122308
## diff_origin1:portrait1 8.810e-01  4.096e-01   2.151 0.031628 *
## artistliving1:endbuyerC 1.038e+00 8.097e-01   1.282 0.200074
```
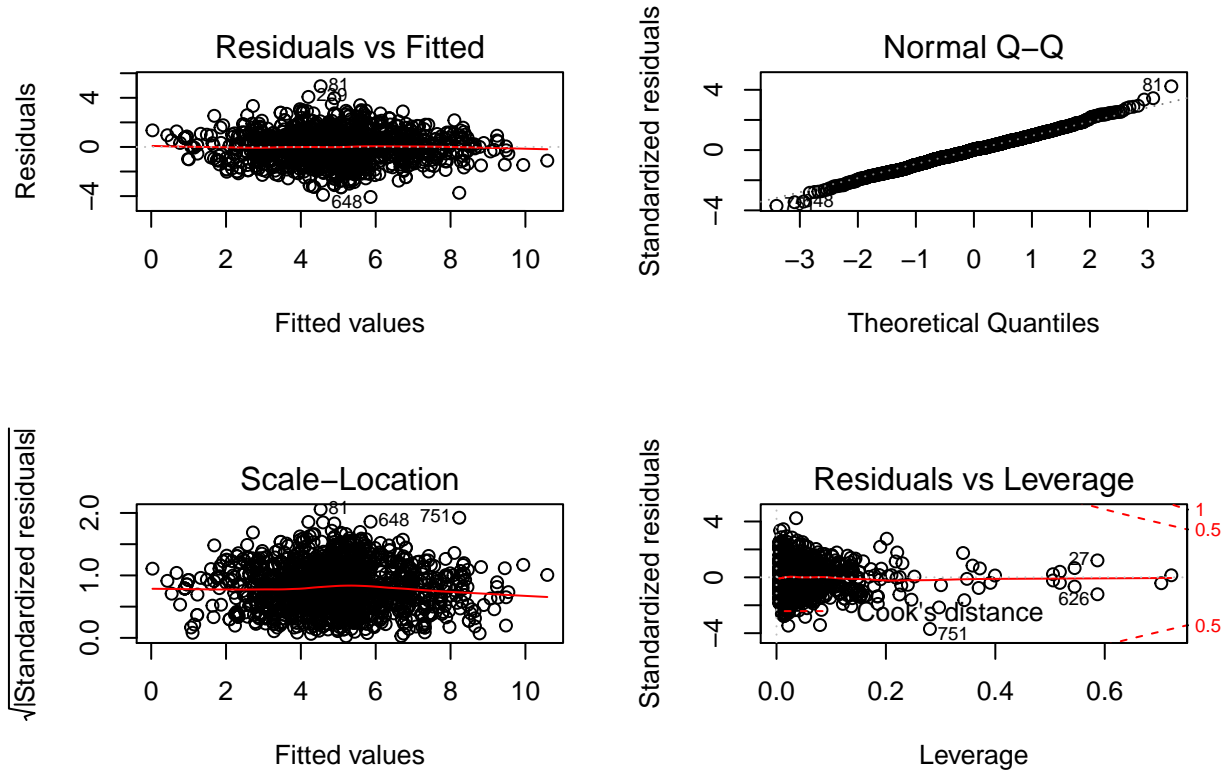
```
## artistliving1:endbuyerD      1.209e+00  8.013e-01   1.509 0.131472
## artistliving1:endbuyerE      2.127e+00  8.594e-01   2.475 0.013444 *
## artistliving1:endbuyerU      1.246e+00  8.264e-01   1.508 0.131811
## artistliving1:endbuyerX      1.455e+00  8.221e-01   1.770 0.076949 .
## artistliving1:authorstyle1   1.432e+00  8.946e-01   1.601 0.109673
## Interm1:Surface              4.456e-04  1.559e-04   2.858 0.004327 **
## Interm1:lrgfont1            -8.787e-01  2.693e-01  -3.263 0.001127 **
## Surface:lrgfont1            -2.213e-04  1.148e-04  -1.928 0.054113 .
## Surface:still_life1         -5.516e-04  2.449e-04  -2.252 0.024460 *
## Surface:discauth1           -1.473e-04  2.244e-04  -0.656 0.511886
## finished1:prevcoll1         -1.077e+00  3.250e-01  -3.313 0.000946 ***
## paired1:lrgfont1            -7.749e-01  2.542e-01  -3.049 0.002341 **
## paired1:discauth1           -4.412e-01  3.435e-01  -1.284 0.199193
## diff_origin1:authorstyle1    9.431e-01  4.509e-01   2.092 0.036640 *
## diff_origin1:still_life1    -7.855e-01  3.475e-01  -2.261 0.023932 *
## finished1:discauth1          5.814e-01  3.099e-01   1.876 0.060877 .
## discauth1:lrgfont1          -5.129e-01  4.376e-01  -1.172 0.241333
## finished1:artistliving1     -4.352e-01  2.836e-01  -1.535 0.125069
## Interm1:portrait1           -8.756e-01  5.850e-01  -1.497 0.134661
## dealerL:artistliving1       -9.966e-01  3.591e-01  -2.775 0.005592 **
## dealerP:artistliving1       -6.857e-01  4.341e-01  -1.579 0.114451
## dealerR:artistliving1       -5.102e-01  2.976e-01  -1.714 0.086724 .
## authorstyle1:prevcoll1      -1.618e+00  1.270e+00  -1.274 0.202911
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.189 on 1429 degrees of freedom
## Multiple R-squared:  0.6335, Adjusted R-squared:  0.6155
## F-statistic: 35.28 on 70 and 1429 DF,  p-value: < 2.2e-16
```

## Residuals & Diagnostics Analysis



After fitting the model, we created the four model diagnostic plots. The overall appearances of all four plots seem acceptable, with no obvious outlier or highly influential points shown. The model also does not violate the normality assumption. The constant variance of residuals assumption seems to be satisfied. However, there are 2 cases that are dropped from the plots because they both have leverage of 1, indicating that they could potentially be the outlying cases of underpriced/overpriced paintings that we will later on investigate in, or have extreme price values. It is worth our attention to specifically look at these cases.

## Variables

|  | Coefficient | 2.5% | 97.5% |
|---|---|---|---|
| (Intercept) | 59.076 | 4.497 | 776.049 |
| Shapeoval | 1.694 | 0.805 | 3.562 |
| Shaperound | 0.616 | 0.312 | 1.216 |
| Shapesqu_rect | 1.686 | 1.015 | 2.799 |
| school_pntgD/FL | 0.533 | 0.049 | 5.836 |
| school_pntgF | 0.267 | 0.024 | 2.935 |
| school_pntgG | 0.045 | 0.002 | 1.299 |
| school_pntgI | 0.284 | 0.026 | 3.128 |
| school_pntgS | 0.467 | 0.024 | 9.108 |
| school_pntgX | 0.157 | 0.014 | 1.766 |
| dealerL | 14.272 | 9.327 | 21.838 |
| dealerP | 3.846 | 2.191 | 6.750 |
| dealerR | 6.519 | 4.519 | 9.404 |
| Interm1 | 0.321 | 0.125 | 0.824 |
| Surface | 1.001 | 1.001 | 1.001 |

|  | Coefficient | 2.5% | 97.5% |
|---|---|---|---|
| paired1 | 1.217 | 0.836 | 1.773 |
| finished1 | 3.217 | 2.082 | 4.972 |
| discauth1 | 2.665 | 1.242 | 5.719 |
| diff_origin1 | 0.580 | 0.471 | 0.714 |
| portrait1 | 0.405 | 0.271 | 0.606 |
| artistliving1 | 0.881 | 0.171 | 4.540 |
| endbuyerC | 0.715 | 0.343 | 1.487 |
| endbuyerD | 0.716 | 0.347 | 1.477 |
| endbuyerE | 0.331 | 0.156 | 0.701 |
| endbuyerU | 0.437 | 0.207 | 0.921 |
| endbuyerX | 0.205 | 0.098 | 0.431 |
| authorstyle1 | 0.171 | 0.075 | 0.391 |
| lrgfont1 | 5.925 | 4.194 | 8.371 |
| still_life1 | 0.841 | 0.523 | 1.353 |
| prevcoll1 | 4.565 | 3.268 | 6.377 |
| dealerL:Interm1 | 4.108 | 0.772 | 21.866 |
| dealerP:Interm1 | 6.070 | 0.434 | 84.864 |
| dealerR:Interm1 | 11.132 | 4.222 | 29.347 |
| dealerL:Surface | 0.999 | 0.999 | 1.000 |
| dealerP:Surface | 1.000 | 0.999 | 1.000 |
| dealerR:Surface | 0.999 | 0.999 | 1.000 |
| dealerL:paired1 | 0.300 | 0.182 | 0.496 |
| dealerP:paired1 | 0.505 | 0.251 | 1.016 |
| dealerR:paired1 | 0.740 | 0.490 | 1.118 |
| dealerL:finished1 | 0.537 | 0.202 | 1.427 |
| dealerP:finished1 | 0.494 | 0.222 | 1.100 |
| dealerR:finished1 | 0.551 | 0.343 | 0.886 |
| dealerL:discauth1 | 1.025 | 0.169 | 6.226 |
| dealerP:discauth1 | 0.173 | 0.027 | 1.110 |
| dealerR:discauth1 | 0.401 | 0.197 | 0.814 |
| Surface:diff_origin1 | 1.000 | 1.000 | 1.000 |
| diff_origin1:portrait1 | 2.413 | 1.081 | 5.389 |
| artistliving1:endbuyerC | 2.823 | 0.577 | 13.822 |
| artistliving1:endbuyerD | 3.351 | 0.696 | 16.136 |
| artistliving1:endbuyerE | 8.390 | 1.554 | 45.282 |
| artistliving1:endbuyerU | 3.477 | 0.687 | 17.585 |
| artistliving1:endbuyerX | 4.285 | 0.854 | 21.496 |
| artistliving1:authorstyle1 | 4.187 | 0.724 | 24.214 |
| Interm1:Surface | 1.000 | 1.000 | 1.001 |
| Interm1:lrgfont1 | 0.415 | 0.245 | 0.704 |
| Surface:lrgfont1 | 1.000 | 1.000 | 1.000 |
| Surface:still_life1 | 0.999 | 0.999 | 1.000 |
| Surface:discauth1 | 1.000 | 0.999 | 1.000 |
| finished1:prevcoll1 | 0.341 | 0.180 | 0.645 |
| paired1:lrgfont1 | 0.461 | 0.280 | 0.759 |
| paired1:discauth1 | 0.643 | 0.328 | 1.262 |
| diff_origin1:authorstyle1 | 2.568 | 1.060 | 6.218 |
| diff_origin1:still_life1 | 0.456 | 0.231 | 0.901 |
| finished1:discauth1 | 1.789 | 0.974 | 3.285 |
| discauth1:lrgfont1 | 0.599 | 0.254 | 1.413 |
| finished1:artistliving1 | 0.647 | 0.371 | 1.129 |
| Interm1:portrait1 | 0.417 | 0.132 | 1.312 |

|  | Coefficient | 2.5% | 97.5% |
|---|---|---|---|
| dealerL:artistliving1 | 0.369 | 0.182 | 0.747 |
| dealerP:artistliving1 | 0.504 | 0.215 | 1.180 |
| dealerR:artistliving1 | 0.600 | 0.335 | 1.076 |
| authorstyle1:prevcoll1 | 0.198 | 0.016 | 2.396 |

In the linear model we selected, we included `Shape`, `school_pntg`, `dealer`, `Interm`, `Surface`, `paired`, `finished`, `discauth`, `diff_origin`, `portrait`, `artistliving`, `endbuyer`, `authorstyle`, `lrgfont`, `still_life`, and `prevcoll` as our base predictors. Interactions selected by the model selection process and, for the sake of interpretation, those that are reasonable and interpretable are kept in the model as well. Since the response variable was orginally log-transformed, the exponentiated coefficients and confidence intervals are shown in the table.

## 4. Summary and Conclusions

a. The median price predicted is `exp(4.401232) = 81.55128` livres. The 95% confidence interval is (6.248, 1064.357) livres. The prediction interval is (2.532, 2626.714) livres.

Table 2: 95% Confidence Interval

| fit | lwr | upr |
|---|---|---|
| 75.173 | 5.754 | 982.155 |

Table 3: 95% Prediction Interval

| fit | lwr | upr |
|---|---|---|
| 75.173 | 2.336 | 2418.59 |

**Interpretation**

From the final model we fitted, we found out that the following variables are statistically significant: `dealer`, `Interm`, `Surface`, `finished`, `discauth`, `diff_origin`, `portrait`, `endbuyer (E,U, X)`, `authorstyle`, `lrgfont`, and `prevcoll`. Some of the interactions are statistically important, such as: `dealer*Interm`, `dealer*paried`, `Interm*lrgfont`, `diff_origin*portrait` etc. We picked the most important ones and interpreted as following:

- dealer: the type of dealer that the auction went through significantly affects the price of the painting. For example, compared with dealer J, the average selling price from dealer L is `exp(2.526) = 12.50339` times higher. (Same interpretation for dealer P and R, with different coefficients)

- Interm: when there is an intermediary involved in the transaction, the selling price is `exp(-1.523) = 0.218` times lower than when there is no intermediary involved.

- Surface: for every one squared inches increase in the painting surface, the selling price is expected to increase 8.779e-4 livres.

- finished: if the painting is finished, the selling price on average is `exp(1.087) = 2.965` times higher than when the painting is not finished.

- portrait: if the painting is described as a portrait, the selling price on average is `exp(-0.9246) = 0.3967` times lower than when the painting is not described as a portrait.

- dealer&Interm interaction: when an intermediary is present, which the price of the auctioned painitngs differs significantly among different dealers. For instance, if the dealer is R and with an intermediary existed, the average selling price is 1.854 times higher when the dealer is J with an intermediary.

## Recommendations

In order to understand the auction prices of 18th century paintings and predict prices of paintings with certain features, we recommend historians focusing on the characteristics mentioned above associated with the painitings (just to mention some, not a complete list). For example, in order to find out highly priced pieces, they might want to look for dealer R, with an intermediary involved; they might want to look for dealer L with the painting sold as a pair with another one; they should look for larger, finished paintings; they should focus on paintings whose authors' names are mentioned during the auction.

## Limitations

1. As mentioned in the data cleaning process, some of the variables have too many levels to be fitted and some levels have few observations that are not sufficient for estimating coefficients. Therefore, we grouped some of the variables, leading to the result that our model cannot investigate the true effects of those combined levels. If we had a larger data set, we could potentially release more levels.

2. As our goal is to find a balance between the prediction accuracy and interpretability, our model will not predict the response variable very accurately (as a sacrifice on interpretability).