# Homework 1: Using OLS Regression to Predict Median House Values in Philadelphia

Yiming Cao, Sujan Kakumanu, Angel Sanaa Rutherford

2025-10-15

## 1 Introduction

## 2 Methods

### 2.1 Data Cleaning

### 2.2 Exploratory Data Analysis

### 2.3 Multiple Regression Analysis

Multiple regression models a dependent variable as a function of multiple predictors, rather than a single predictor such as in simple regression. These predictors each have a coefficient that represents their effect on a dependent variable, controlling for all other predictors. This approach improves model accuracy in situations where multiple variables better explain outcomes of a dependent variable.

This report regressed log-transformed median house value (LNMEDHVAL) on the proportion of housing units that are vacant (PCTVACANT), percent of housing units that are single family detached homes (PCTSINGLES), proportion of residents with at least a bachelor's degree (PCTBACHMOR), and log-transformed number of households with incomes below 100% poverty level (LNNBELPOV). This regression function can be expressed as follows:

$$\text{LNMEDHVAL} = \beta_0 + \beta_1 \text{PCTVACANT} + \beta_2 \text{PCTSINGLES} + \beta_3 \text{PCTBACHMOR} + \beta_4 \text{LNNBELPOV} + \varepsilon$$

Multiple regression relies on several key assumptions, most of which mirror the assumptions of simple regression. First, linear relationships should exist between the dependent variable and each predictor, assessed through scatterplots or residual plots and addressed via transformations if needed. Second, residuals should be approximately normally distributed, which can be assessed

1

through a histogram. Third, residuals must be random — indicating that observations are not systematically related. Fourth, residuals must be homoscedastic, exhibiting constant variance across all values. Fifth, the dependent variable should be continuous.

A unique assumption for multiple regression is avoiding perfect multicollinearity: no predictor should be strongly correlated with others. Multicollinearity inflates standard errors and produces unstable coefficient estimates. This assumption can be checked by analyzing the correlation coefficients between all dependent variables, with anything greater than 0.9 generally being a cause for concern. Variance Inflation Factor (VIF) can be used to further inspect a suspicion of multicollinearity, with a VIF $< 5$ being generally acceptable and a VIF $< 10$ warranting more inspection. A VIF $> 10$ strongly indicates multicollineariy.

In the above multiple regression function, $\beta_0$ represents the depedent variable when all predictors are zero. The coefficients of the predictors $\beta_1, \beta_2, \beta_3, \beta_4$ each represent the change in the dependent variable with a one unit increase in the predictor, holding all other predictors constant.

These $\beta$ coefficients in multiple regression are simultaneously estimated in order to minimize the Error Sum of Squares (SSE). The general formula and breakdown of what is to be minimized is provided below (with n being the number of observations, and k is the number of predictors):

$$SSE = \sum_{i=1}^{n} \varepsilon^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} \left[ y_i - \left( \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \cdots + \hat{\beta}_k x_{ki} \right) \right]^2$$

This minimization works by finding the $\beta$ coefficients that, when raw predictor $(x_i)$ data is used, will minimize the residuals $(y_i - \hat{y}_i)$. SSE is also used to calculate Mean Squared Error (MSE), noted by the estimated parameter $\hat{\sigma}^2$. This is the estimate of the variance of the error term $\epsilon$. The formula for MSE, in terms of SSE is noted below:

$$MSE = \frac{SSE}{n - (k + 1)}$$

Another term in regression analysis is Total Sum of Squares (SST). It measures the total variation in the dependent variable around it's mean by using the following formula:

$$SST = \sum_{i=1}^{n} (y_i - \bar{y})^2$$

Using this formula for SST, and the previously stated formula for SSE, we can calculate $R^2$ — the coefficient of multiple determination. This is the proportion of variance in the model explained by all k predictors, and is the represented by the following:

$$R^2 = 1 - \frac{\text{SSE}}{\text{SST}}$$

Multiple regression presents a unique dillema in comparison to simple regression, in that adding more predictors will generally increase $R^2$. Adjusting $R^2$, noted below, can account for additional predictors and determine whether or not they are improving the model.

$$R^2_{\text{adj}} = \frac{(n-1)R^2 - k}{n - (k+1)}$$

vii. State the hypotheses you test. Specifically, talk about the F-ratio and the H0 and Ha associated with it, as well as the hypotheses you test about each of the individual i's (again, state H0 and Ha).

This report will conduct two tests to evaluate the model. First, there is the F-ratio — a model utility test. F-ratio tests the following null hypothesis $H_0$ and alternative hypothesis $H_a$:

$$H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$$

$$H_a : \text{At least one } \beta_i \neq 0$$

In essence, the null hypothesis states that all of the model $\beta$ parameters are zero, and the alternative states that at least one of those parameters is not zero. Failure to reject the null hypothesis suggests that the model is incredibly weak, and should be reevaluated. If the null hypothesis is rejected, the second test can be conducted with the following hypotheses.

$$H_0 : \beta_i = 0$$

$$H_a : \beta_i \neq 0$$

In this test, we evaluate the performance of each predictor i (in the case of this report, the 4 predictors stated earlier). A t-test can be conducted, where the t-statistic for each predictor is calculated as the estimated coefficient divided by its standard error:

$$t_i = \frac{\hat{\beta}_i - \beta_i}{s_{\hat{\beta}_i}}$$

Each predictor has its own p-value calculated using the above t-statistic. If the p-value is less than 0.05, we reject the null hypothesis for that predictor and conclude that it is a statistically significant predictor of the dependent variable. If the p-value is greater than or equal to 0.05, we fail to reject the null hypothesis and conclude that the predictor is not statistically significant.

## 2.4 Additional Analysis

## 2.5 Software Used

# 3 Results

## 3.1 Exploratory Results

## 3.2 Regression Results

   i. Present the regression output from R. Be sure that your output presents the parameter estimates (and associated standard errors, t-statistics and p-values), as well as the R2, the adjusted R2, and the relevant F-ratio and associated p-value.

```
Call:
lm(formula = LNMEDHVAL ~ PCTVACANT + PCTSINGLES + PCTBACHMOR +
    LNNBELPOV, data = Regression_shpData)

Residuals:
     Min       1Q   Median       3Q      Max
-2.25817 -0.20391  0.03822  0.21743  2.24345

Coefficients:
              Estimate Std. Error t value            Pr(>|t|)
(Intercept) 11.1137781  0.0465318 238.843 < 0.0000000000000002 ***
PCTVACANT   -0.0191563  0.0009779 -19.590 < 0.0000000000000002 ***
PCTSINGLES   0.0029770  0.0007032   4.234            0.0000242 ***
PCTBACHMOR   0.0209095  0.0005432  38.494 < 0.0000000000000002 ***
LNNBELPOV   -0.0789035  0.0084567  -9.330 < 0.0000000000000002 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3665 on 1715 degrees of freedom
Multiple R-squared:  0.6623,    Adjusted R-squared:  0.6615
F-statistic: 840.9 on 4 and 1715 DF,  p-value: < 0.00000000000000022


Analysis of Variance Table

Response: LNMEDHVAL
            Df  Sum Sq Mean Sq  F value               Pr(>F)
```
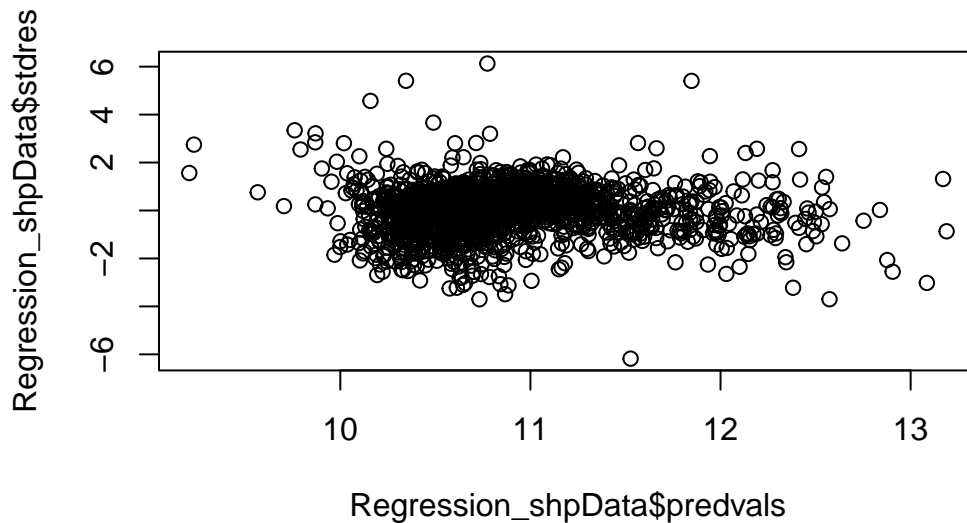
```
PCTVACANT     1 180.383 180.383 1343.093 < 0.00000000000000022 ***
PCTSINGLES    1  24.543  24.543  182.741 < 0.00000000000000022 ***
PCTBACHMOR    1 235.111 235.111 1750.586 < 0.00000000000000022 ***
LNNBELPOV     1  11.692  11.692   87.054 < 0.00000000000000022 ***
Residuals  1715 230.332   0.134
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



ii. Referencing the regression output in (i) above, interpret the results as in the example included above this report outline. NOTE: YOUR DEPENDENT VARIABLE (AND SOME PREDICTORS) WOULD BE LOG-TRANSFORMED, UNLIKE IN THE EXAMPLE HERE. LOOK AT THE SLIDES FOR EXAMPLES OF INTERPRETING REGRESSION OUTPUT WITH LOG-TRANSFORMED VARIABLES.

## 3.3 Regression Assumption Checks

## 3.4 Additional Models

# 4 Discussion & Limitations