Barcelona Real Estate Price Prediction Model

Applied Probability and Statistics - Section C, Team 44

Shifei Ruan(sr678), Seokho Shin(ss1722), Arayansh Vaish(av323), Lucy Zheng(zz421)

- 1 Introduction: This report outlines the process of developing a predictive model for property prices in Barcelona. The objective was to construct a robust and interpretable statistical model using a dataset of 413 properties and subsequently apply it to predict the prices of 200 new properties.
- **2 Data Preparation and Initial Model:** We began with a multiple linear regression model containing all available variables (model1). This initial full model served as a baseline, yielding an R^2 of 0.7418, indicating that the variables explained a substantial portion of the price variance. However, several variables (Kitchen, Type) were highly insignificant (p-values > 0.9), which suggested model overspecification.

3 Model Refinement Process

The model was refined through an iterative process.

3.1 Removing Insignificant Variables: The first step was to remove the clearly insignificant variables Kitchen and Type to create a more parsimonious model (model2). This simplification did not harm the model's explanatory power (R-squared was still around 0.74) but enhanced its focus on impactful features.

Model 2

3.2 Addressing Non-Linearity and Interaction Effect

- **3.2.1 Interaction Term:** We tested an interaction between m2 and Rooms (model3). The results showed that this interaction term (m2:Rooms) was highly statistically significant (p = 0.000354). This indicates that the effect of size on price is not constant but depends on the number of rooms.
- **3.2.2 Logarithmic Transformation:** A major breakthrough was applying a logarithmic transformation to both the dependent variable (Price) and the key independent variable (m2) to create model4. This transformation successfully linearized the relationship.

```
Estimate Std. Error t Value Pr(>|t|)

(Intercept) 50905.37 27905.19 1.824 0.068858 .

m2 2057.83 352.94 5.830 1.13e-08 ***

Rooms -36637.52 8254.95 -4.438 1.17e-05 ***

Bathrooms 22296.65 7365.08 3.027 0.002626 **

Elevator 7722.08 6730.35 1.147 0.251917

Atico 16062.29 11791.32 1.362 0.173891

Terrasse 8410.14 8222.26 1.023 0.306992

Parking 47076.18 11955.67 3.938 9.70e-05 ***

Yard 31677.87 16182.12 1.958 0.050969 .

m2:Rooms 304.64 84.56 3.603 0.000354 ***

---

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 63180 on 403 degrees of freedom

Multiple R-squared: 0.7498, Adjusted R-squared: 0.7442

F-statistic: 134.2 on 9 and 403 DF, p-value: < 2.2e-16
```

```
Coefficients:

Estimate Std. Error t value Pr(>|t|)

(Intercept) 8.414823 0.222013 37.902 < 2e-16 ***
log(m2) 0.873776 0.054837 15.934 < 2e-16 ***
Rooms -0.100428 0.050077 -2.005 0.045581 *

I(Rooms^2) 0.015093 0.008338 1.810 0.071014 .

Bathrooms 0.115981 0.026116 4.441 1.16e-05 ***

Elevator 0.060001 0.023869 2.514 0.012332 *

Atico 0.088991 0.038536 2.309 0.021431 *

Parking 0.164928 0.043035 3.832 0.000147 ***

Yard 0.088315 0.058294 1.515 0.130559

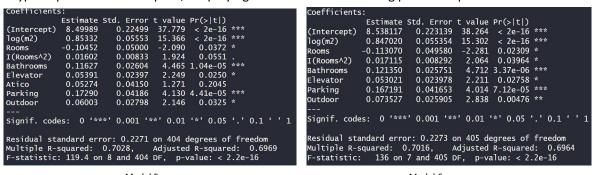
---

Signif. codes: 0 '***' 0.001 '**' 0.01 '* 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2277 on 404 degrees of freedom Multiple R-squared: 0.7011, Adjusted R-squared: 0.6952 F-statistic: 118.4 on 8 and 404 DF, p-value: < 2.2e-16
```

Model 3 Model 4

- **3.2.3 Non-linear Effect of Rooms:** We identified a potential quadratic relationship for the number of Rooms by adding I(Rooms^2). The negative linear and positive quadratic coefficients in later models suggested a complex, U-shaped relationship where price might decrease slightly from 1 to 2-3 rooms and then increase for larger dwellings.
- **3.3 Feature Engineering:** The binary variables Yard and Terrasse were merged into a new variable Outdoor (model5). This was a logical step to create a more generalized feature for any type of private outdoor space, simplifying the model without losing predictive power.



Model 5 Model 6

- **3.4 Final Variable Selection:** The variable Atico remained statistically insignificant (p-value = 0.204) even in the log-transformed model, so it was removed to yield a cleaner, more efficient model (model6).
- **4 Geographic Fixed Effects:** The most significant improvement came from incorporating City Zone as a fixed effect (model7). Location is a paramount factor in real estate valuation. By adding factor('City Zone'), we controlled for all unobserved, time-invariant characteristics of each district. Adjusted R-squared jumped dramatically and coefficients for geographic areas were highly significant.

Crucially, after controlling for location, the variables Rooms and I(Rooms^2) became statistically insignificant. This suggests that the apparent nonlinear effect of rooms was actually a proxy for location-specific housing characteristics.

```
0.045861
                                                        148346
                                                                    0.035705
outdoor
                                                        086857
                                                                    0.022350
                                                                                  3 886 0 000119
                                                        127883
factor(`City Zone`)Eixample
factor(`City Zone
factor(`City Zone
factor(`City Zone
                       )Gràcia
                                                        106433
                                                                    0.046160
                                                                                    306 0.021642
                                                        .043036
                                                                    0.044127
                      )Horta - Guinardó
                                                                                 -0.975 0.330022
                                                                    0.064810
0.042873
                                                        .070028
factor(`City Zone
factor(`City Zone
                                                                                         6.74e-07
                      )Nou Barris
                                                     -0.216522
                                                                                    .050
                                                        142134
                                                                                  3.251 0.001249
                       )Sant Andreu
                      )Sant Marti
factor(`City Zone
                                                                    0.046708
                                                                                     368 0.000831
                       )Sants - Montjuïc
         City Zone
factor(`City Zone`)Sarria - Sant Gervasi 0.326622
                                                                   0.061043
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.0<mark>5 '.' 0.1 ' ' 1</mark>
Residual standard error: 0.193 on 396 degrees of freedom
Multiple R-squared: 0.7894. Adjusted R-squared: 0.78
```

Model 7

5 Final Model Selection: Therefore, the final model (final_simple) was simplified by removing the room-related variables, resulting in a powerful and interpretable specification:

Final Model:

log(Price) ~ log(m2) + Bathrooms + Elevator + Parking + Outdoor + factor('City Zone')