

```
1  * Set the working directory
2  cd "~/Desktop/Replication_Package"
3
4  * Step 1: Load the public school dataset and perform the
   aggregation by ZIP code
5  use "public_school_ny.dta", clear
6
7  * Define school categories
8  gen elementary_school = SCHOOL_LEVEL == "Elementary"
9  gen middle_school = SCHOOL_LEVEL == "Middle"
10 gen high_school = SCHOOL_LEVEL == "High"
11 gen other_school = inlist(SCHOOL_LEVEL, "Other", "Prekindergarten",
   "Secondary", "Ungraded", "Not reported")
12
13 * Aggregate by ZIP code
14 collapse (sum) elementary_schools=elementary_school ///
15           (sum) middle_schools=middle_school ///
16           (sum) high_schools=high_school ///
17           (sum) other_schools=other_school ///
18           (mean) avg_student_teacher_ratio=STUTERATIO ///
19           (sum) total_free_lunch=TOTFRL, by(LZIP)
20
21 * Calculate total schools and rename ZIP code
22 gen total_schools = elementary_schools + middle_schools +
   high_schools + other_schools
23 rename LZIP zip_code
24 save "Public_School_Aggregated.dta", replace
25
26 * Step 2: Merge with the housing data
27 use "NY_Housing.dta", clear
28 merge m:1 zip_code using "Public_School_Aggregated.dta"
29
30 * Keep only matched observations
31 drop if _merge != 3
32 drop _merge
33
34 * Save the merged dataset
35 save "Housing_Schools_Merged.dta", replace
36
37 * Step 3: Load the income dataset
38 use "Income.dta", clear
39 duplicates drop zip_code, force
40 save "Income.dta", replace
41
42 * Reload the housing dataset
43 use "Housing_Schools_Merged.dta", clear
44 merge m:1 zip_code using "Income.dta"
45
```

```
46 * Check merge results
47 drop if _merge != 3
48 drop _merge
49
50 * Save the final merged dataset
51 save "Final_Merged_Data.dta", replace
52
53 * Step 4: Load the final merged dataset
54 use "Final_Merged_Data.dta", clear
55
56 * Step 5: Clean and Filter the Dataset
57 keep price bed bath acre_lot house_size zip_code total_free_lunch
58 total_schools costofliving medianincome avg_student_teacher_ratio
59 drop if price < 60000 | price > 20000000
60 drop if bed > 15
61 drop if bath > 12
62 drop if house_size > 10000 | house_size < 100
63 drop if acre_lot > 60
64 drop if avg_student_teacher_ratio > 20
65 drop if total_schools > 35
66 drop if missing(acre_lot)
67 drop if missing(price) | missing(bed) | missing(bath) | missing(
68 house_size) | missing(total_free_lunch) | missing(
69 avg_student_teacher_ratio)
70
71 * Step 6: Prepare Variables for Analysis
72
73 * Drop unused variables, and log-transform key variables
74 gen log_price = log(price)
75 gen log_house_size = log(house_size)
76 gen log_acre_lot = log(acre_lot)
77 drop if missing(costofliving) | missing(medianincome) | missing(
78 log_acre_lot)
79 gen log_total_free_lunch = log(total_free_lunch)
80 drop if missing(log_total_free_lunch)
81
82 * Step 7: Convert String Variables to Numeric
83
84 * Convert costofliving from string to numeric
85 generate costofliving_num = real(subinstr(subinstr(costofliving,
86 "$", "", .), ",", "", .))
87 drop costofliving
88 rename costofliving_num costofliving
89 format costofliving %12.2f
90
91 * Convert medianincome from string to numeric
92 generate medianincome_num = real(subinstr(subinstr(medianincome,
93 "$", "", .), ",", "", .))
```

```
88 drop medianincome
89 rename medianincome_num medianincome
90 format medianincome %12.2f
91
92 * Step 8: Standardize Variables
93
94 * Standardize house_size
95 capture drop std_house_size
96 summarize house_size, detail
97 local house_size_mean = r(mean)
98 local house_size_sd = r(sd)
99 gen std_house_size = (house_size - `house_size_mean') /
`house_size_sd'
100
101 * Standardize avg_student_teacher_ratio
102 capture drop std_avg_student_teacher_ratio
103 summarize avg_student_teacher_ratio, detail
104 local avg_student_teacher_ratio_mean = r(mean)
105 local avg_student_teacher_ratio_sd = r(sd)
106 gen std_avg_student_teacher_ratio = (avg_student_teacher_ratio -
`avg_student_teacher_ratio_mean') / `avg_student_teacher_ratio_sd'
107
108 * Standardize acre_lot
109 capture drop std_acre_lot
110 summarize acre_lot, detail
111 local acre_lot_mean = r(mean)
112 local acre_lot_sd = r(sd)
113 gen std_acre_lot = (acre_lot - `acre_lot_mean') / `acre_lot_sd'
114
115 * Standardize medianincome
116 capture drop std_medianincome
117 summarize medianincome, detail
118 local medianincome_mean = r(mean)
119 local medianincome_sd = r(sd)
120 gen std_medianincome = (medianincome - `medianincome_mean') /
`medianincome_sd'
121
122 * Standardize costofliving
123 capture drop std_costofliving
124 summarize costofliving, detail
125 local costofliving_mean = r(mean)
126 local costofliving_sd = r(sd)
127 gen std_costofliving = (costofliving - `costofliving_mean') /
`costofliving_sd'
128
129 * Step 9: Save the cleaned dataset
130 save "Final_Merged_Data_Cleaned.dta", replace
131
```

```
132 * Step 10: Summary Statistics
133 summarize price bed bath acre_lot house_size total_free_lunch
    total_schools costofliving medianincome avg_student_teacher_ratio
134
135 * Step 11: Figures
136
137 * Figure 1: Distribution of Log-Transformed Housing Prices
138 histogram log_price, normal
139 graph export "figure1.png", as(png) replace
140
141 * Figure 2: Residuals for Model 1
142 regress log_price bed
143 predict residuals1, resid
144 histogram residuals1, normal
145 graph export "figure2.png", as(png) replace
146
147 * Figure 3: Residuals for Model 7
148 regress log_price bed bath c.log_house_size##c.log_acre_lot
    total_schools std_avg_student_teacher_ratio log_total_free_lunch
    std_medianincome std_costofliving zip_code
149 predict residuals7, resid
150 histogram residuals7, normal
151 graph export "figure3.png", as(png) replace
152
153 * Step 12: Run Regression Models
154
155 * Model 1: Baseline Model – Bedrooms
156 regress log_price bed
157 est store model1
158
159 * Model 2: Add Housing Characteristics
160 regress log_price bed bath std_house_size std_acre_lot
161 est store model2
162
163 * Model 3: Replace with Log-Transformed House Size and Lot Size
    (r-squared improves)
164 regress log_price bed bath log_house_size log_acre_lot
165 est store model3
166
167 * Model 4: Add Total Schools
168 regress log_price bed bath log_house_size log_acre_lot
    total_schools
169 est store model4
170
171 * Model 5: Add Student-Teacher Ratio and Free Lunch
172 regress log_price bed bath log_house_size log_acre_lot
    total_schools std_avg_student_teacher_ratio log_total_free_lunch
173 est store model5
```

```
175 * Model 6: Add Median Income
176 regress log_price bed bath log_house_size log_acre_lot
    total_schools std_avg_student_teacher_ratio log_total_free_lunch
    std_medianincome
177 est store model6
178
179 * Model 7: Add Cost of Living and Interaction
180 regress log_price bed bath c.log_house_size##c.log_acre_lot
    total_schools std_avg_student_teacher_ratio log_total_free_lunch
    std_medianincome std_costofliving zip_code
181 est store model7
182
183 esttab model1 model2 model3 model4 model5 model6 model7 using
    regression_table.tex, replace ///
184 label b(3) se stats(r2 N) compress
185
```