

# Hawaii Airbnb Smart Pricing and Market Segmentation System

Jingyi Chen, Zhanyuan Jiang, Jiaqi Li

December 17, 2025

## Abstract

This project analyzes Airbnb pricing in Oahu using a sequence of predictive models and market segmentation techniques. We first establish baseline performance using standard linear and tree-based models, which serve as reference points for later comparisons. Building on these benchmarks, we investigate whether incorporating market segmentation can further improve predictive accuracy. We compare two segmentation-based strategies: cluster-then-predict and cluster-as-features and find that the latter delivers more reliable performance. By incorporating cluster membership directly into a global nonlinear model, we achieve the best balance between predictive accuracy and stability.

## 1 Background

Airbnb has become one of the most widely used platforms for short-term rentals in the world, allowing hosts to list properties and set nightly prices while guests search for accommodations that match their preferences and budget. In this process, pricing plays a central role in host revenue and guest decision-making.

Airbnb does provide a default Smart Pricing tool to assist hosts in setting competitive prices. The Hawaii Airbnb market, however, has highly differentiated pricing patterns due to geographic location and seasonal tourism. Ocean-front views, access to beaches, amenities such as pools and AC, and proximity to tourist areas significantly influence price dispersion. Unfortunately, these localized factors are not fully captured by Airbnb’s generic Smart Pricing algorithm.

Hence, to better serve Hawaii hosts, this research aims to build a localized pricing system that identifies data-driven sub-markets and provides more transparent price recommendations. As part of this research, we envision a dashboard where hosts can see their market segment and obtain a data-driven nightly price estimate tailored to Hawaii’s unique market conditions.

## 2 Data Description

### 2.1 Data Resource

The dataset used in this project is obtained from [Inside Airbnb](#), a publicly accessible platform that provides detailed, regularly updated snapshots of Airbnb activity across global cities. Inside Airbnb offers quarterly datasets for the past year for each region, allowing time-specific and market-level analyses.

The Hawaii Islands comprise several distinct regions, but this study focuses exclusively on Oahu, which offers the largest concentration of Airbnb listings and sufficient data for reliable analysis.

### 2.2 Feature Engineering

The original dataset obtained from Inside Airbnb contains a wide range of listing-level information, including property characteristics (e.g., room type, accommodates, bathrooms, bedrooms), host attributes, review scores, geographic coordinates, and a textual amenities list.

To better capture pricing drivers in the Oahu market, we engineered additional features using the Google Maps Distance Matrix API, computing each listing’s driving distance and travel time to Honolulu International Airport and Waikiki Beach. We also parsed the amenities text into structured

indicators (e.g., pool, AC, ocean view) and created an amenities count variable. The target variable was transformed using  $\log(\text{price})$ , categorical attributes were one-hot encoded, and missing values were handled through imputation with corresponding missing-value flags.

### 3 Methodology

#### 3.1 Baseline Pricing Model

To establish a clear benchmark for evaluating subsequent models, we first fit a set of standard predictive models directly on the full set of cleaned and engineered features. At this stage, the goal is not aggressive optimization, but rather to assess how well common modeling approaches perform when the Airbnb market is treated as a homogeneous population, without explicitly accounting for market segmentation or structural heterogeneity. These baseline models provide reference points for judging whether additional modeling structure is warranted. An overview of the baseline models and their out-of-sample performance is reported in Table 1.

Model	Type	Tuning Strategy	Test $R^2$	Test $RMSE$
OLS	Linear	No regularization	0.620	0.679
Ridge with CV ( $\alpha=3.16$ )	Linear	5-fold CV to select $\alpha$	0.620	0.679
Random Forest with RandomizedSearchCV	Nonlinear	5-fold CV	0.828	0.457
LightGBM with RandomizedSearchCV	Nonlinear	5-fold CV	0.861	0.410
CatBoost with RandomizedSearchCV	Nonlinear	5-fold CV	0.867	0.402
XGBoost with RandomizedSearchCV	Nonlinear	5-fold CV	0.874	0.392

Table 1: Overview of the Baseline Models

##### 3.1.1 Linear Models: OLS and Ridge Regression

We begin with linear models, including Ordinary Least Squares (OLS) and Ridge regression. Both models are trained on the same feature set, with Ridge introducing  $\ell_2$  regularization to stabilize coefficient estimates in the presence of multicollinearity. In practice, both OLS and Ridge achieve similar out-of-sample performance, with test  $R^2$  values around 0.62. This indicates that regularization improves coefficient stability but does not meaningfully enhance predictive accuracy. Overall, linear models appear insufficient to capture the complexity of Airbnb pricing dynamics.

##### 3.1.2 Tree-Based Models

To capture nonlinear relationships and feature interactions, we next evaluate several tree-based models, including Random Forest, LightGBM, CatBoost, and XGBoost[1]. All tree-based approaches substantially outperform linear models, confirming the presence of strong nonlinear effects in the data. Among them, XGBoost delivers the best performance, achieving a test  $R^2$  of approximately 0.87, representing a large improvement over the linear baseline. The remaining tree-based models serve as useful robustness checks and achieve competitive performance, but do not surpass XGBoost.

Overall, these baseline results suggest that while increased model flexibility yields large gains in predictive accuracy, further improvements are unlikely to come from algorithmic complexity alone. This motivates the exploration of models that explicitly incorporate market structure and segmentation in later sections.

#### 3.2 Market Segmentation

##### 3.2.1 Clustering in Original Feature Space

Although several baseline models (e.g., XGBoost) achieve strong predictive performance, they rely on a single global pricing relationship across all listings. In practice, Airbnb pricing on Oahu is highly location-dependent: beachfront properties near Waikiki, airport-adjacent listings, and inland residential units often follow distinct pricing mechanisms. To capture such heterogeneity, we apply unsupervised clustering to partition the market into data-driven sub-markets.

Clustering is performed on a subset of interpretable features capturing geographic accessibility[2] and listing characteristics, including latitude, longitude, driving distances to Honolulu International Airport and Waikiki Beach, review-based quality measures, and accommodation-related attributes (capacity, bedrooms, bathrooms). All features are standardized prior to clustering. Using Agglomerative Clustering and silhouette scores to select the number of clusters, we find that a coarse segmentation with  $k = 3$  provides a good balance between within-cluster cohesion and between-cluster separation (silhouette score = 0.524). A robustness check using K-means yields broadly consistent partitions, suggesting that the identified sub-markets are not sensitive to the specific clustering algorithm. Figure 1 visualizes the resulting clusters on the latitude–longitude map, confirming clear geographic differentiation.

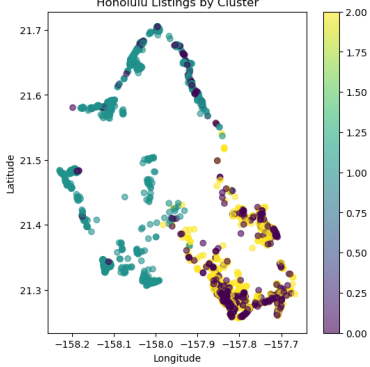


Figure 1: Geographic Clusters

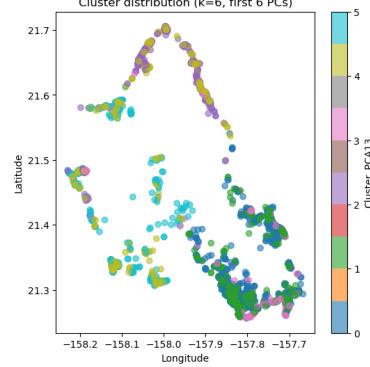


Figure 2: PCA-Based Clusters

### 3.2.2 PCA-Based Analysis and Supplementary Clustering

To explore whether additional latent structure exists beyond these coarse groupings, we further apply Principal Component Analysis (PCA) to the same standardized feature set. The first four principal components explain approximately 90 % of the total variance, indicating that much of the information in the original features can be summarized in a low-dimensional, decorrelated representation.

Clustering in the PCA-transformed space yields a different silhouette[3] profile, with the optimal value attained at  $k = 6$  (silhouette score = 0.61). As illustrated in Figure 2, these clusters exhibit finer spatial separation compared to clustering in the original geographic space. This does not imply that PCA-based clustering is superior; rather, PCA changes the notion of similarity by emphasizing directions of maximal variance and removing feature correlations. Taken together, the results suggest that the Oahu Airbnb market exhibits multi-scale structure: a small number of intuitive, geographically interpretable segments, alongside more nuanced subdivisions revealed in a variance-driven representation.

Motivated by these findings, we subsequently consider incorporating segmentation information into pricing models, either by adding cluster membership as features or by fitting separate models within each identified sub-market.

## 4 Results

### 4.1 Baseline and Segmentation-Aware Models

To evaluate whether explicitly modeling market heterogeneity improves pricing performance, we compare two segmentation-aware strategies: *cluster-then-predict* and *cluster-as-features*. Both approaches build upon the baseline tree-based models introduced earlier, which already outperform linear models.

### 4.2 Cluster-Then-Predict

In the cluster-then-predict approach, we first partition listings into sub-markets using agglomerative clustering based on geographic and listing characteristics. Separate predictive models are then trained within each cluster. Due to computational constraints, we focus on two representative tree-based methods: Random Forest (RF) and LightGBM (LGBM).

Table 2 reports the aggregated out-of-sample performance across clusters. While cluster-specific models achieve strong in-sample fit, overall test performance is comparable to the best global baseline models, with LightGBM achieving a test  $R^2$  of 0.840.

Model	Test $R^2$	Test RMSE
RF (cluster-then-predict)	0.815	0.474
LGBM (cluster-then-predict)	0.840	0.441

Table 2: Overall test performance for cluster-then-predict

### 4.3 Cluster-As-Features

Alternatively, we incorporate cluster assignments directly into the feature space as categorical indicators, allowing a single global model to condition predictions on sub-market membership. This approach avoids training separate models per cluster while still leveraging segmentation information.

Table 3 shows that adding cluster features consistently improves performance across tree-based models. XGBoost with cluster features achieves the best overall performance, with a test  $R^2$  of 0.873 and RMSE of 0.393.

Model	Test $R^2$	Test RMSE
RF + Cluster Features	0.830	0.454
LightGBM + Cluster Features	0.868	0.401
CatBoost + Cluster Features	0.861	0.410
XGBoost + Cluster Features	<b>0.873</b>	<b>0.393</b>

Table 3: Performance of cluster-as-features models

Overall, the cluster-as-features strategy provides a better balance between predictive accuracy and computational efficiency, outperforming both global baselines and cluster-then-predict models.

### 4.4 Interpretation of Price Drivers

To interpret pricing drivers, we analyze both linear and nonlinear models. Ridge regression highlights capacity, host portfolio size, review quality, and location proxies as key linear predictors. To capture nonlinear effects, we further examine SHAP values from tree-based models with cluster features.

SHAP-based feature importance analyses for the tree-based models are reported in the Appendix. Across models, listing capacity (e.g., accommodates and bedrooms), host activity, availability, and review scores consistently dominate model predictions. Cluster indicators contribute to pricing predictions but play a secondary role relative to core listing characteristics.

## 5 Discussion and Future Work

### 5.1 Linear versus Tree-Based Models

Results from the baseline analysis show that linear models (OLS and Ridge) achieve similar out-of-sample performance, with test  $R^2$  values around 0.62. Introducing  $\ell_2$  regularization stabilizes coefficient estimates but does not meaningfully improve predictive accuracy, indicating that multicollinearity is not the primary limitation of the linear specification. Instead, the main constraint is the limited expressive capacity of global linear models when applied to Airbnb pricing.

In contrast, tree-based models substantially outperform linear baselines. Random Forest, LightGBM, CatBoost, and XGBoost all achieve much higher test  $R^2$  values, confirming the presence of strong nonlinear effects and feature interactions in pricing behavior. These gains suggest that modeling flexibility, rather than regularization strength, is the key driver of improved predictive performance.

## 5.2 Comparing Cluster-Then-Predict and Cluster-as-Features

Motivated by the observed market heterogeneity, we explored two strategies for incorporating clustering information into the pricing models: cluster-then-predict and cluster-as-features.

The cluster-then-predict approach fits separate models within each cluster, aiming to capture segment-specific pricing relationships. However, this strategy fragments the data, reducing effective sample size within each cluster and increasing estimation variance. Although cluster-level models often achieve strong in-sample fit, their out-of-sample performance is unstable, particularly for smaller or more heterogeneous clusters.

In contrast, the cluster-as-features approach augments a global model with cluster membership indicators while preserving the full training sample. This method consistently outperforms cluster-then-predict and yields more stable improvements across all tree-based models. By conditioning predictions on sub-market membership without sacrificing data efficiency, cluster-as-features achieves a more favorable bias–variance tradeoff.

## 5.3 Why Segmentation Yields Modest but Consistent Gains

Although incorporating cluster information improves predictive performance, the magnitude of the gain is modest. This suggests that much of the segmentation-related information is already implicitly captured by strong tree-based models through geographic variables, capacity measures, and availability constraints.

Rather than introducing fundamentally new information, cluster indicators act as coarse summaries of patterns that tree models partially learn on their own. As a result, clustering refines predictions rather than transforming the pricing mechanism, leading to incremental but consistent performance improvements.

## 5.4 Interpretability and Feature Effects

Both Ridge coefficients and SHAP analyses highlight a consistent set of price drivers across models. Listing capacity (e.g., accommodates, bedrooms), host portfolio size, availability constraints, and review-related measures emerge as the most influential features.

SHAP results further reveal nonlinear and asymmetric effects that are not captured by linear models. Cluster indicators appear in SHAP rankings but contribute less than core listing and host characteristics, reinforcing the interpretation that segmentation information complements rather than dominates the predictive structure.

## 5.5 Future Work

Several directions remain for future research. First, more flexible spatial methods—such as spatially varying coefficient models or graph-based clustering—may capture finer-grained geographic effects. Second, incorporating temporal dynamics and seasonality could improve pricing accuracy in highly seasonal markets like Hawaii. Finally, integrating demand-side signals, such as booking probability or occupancy rates, may allow pricing models to move beyond static prediction toward revenue optimization.

## References

- [1] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794, 2016.
- [2] Yanfang Liu, Xiang Kong, Xiaoping Liu, and Shixuan Wang. Analyzing spatial variance of airbnb pricing determinants using multiscale geographically weighted regression (mgwr) approach. *Sustainability*, 12(11):4710, 2020.
- [3] Peter J Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987.

## A Additional Figures

Figure A.1–A.4 present SHAP summaries for tree-based models with cluster-as-features. These figures complement the discussion in Section 4.

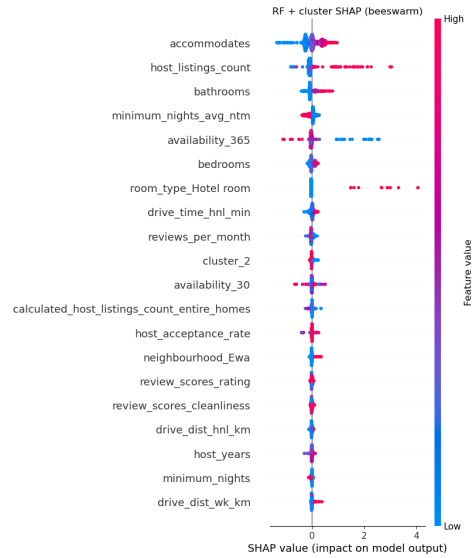


Figure A.1: SHAP summary for Random Forest with cluster-as-features.

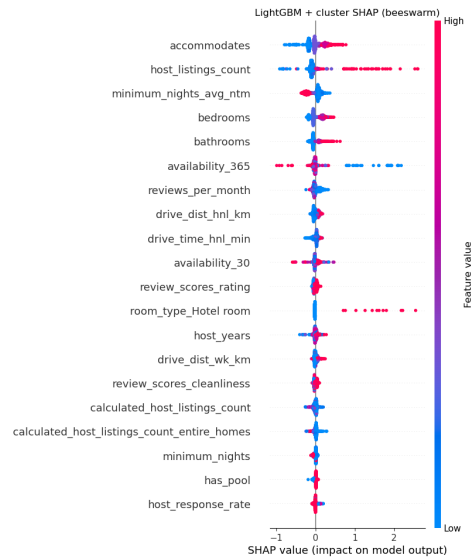


Figure A.2: SHAP summary for LightGBM with cluster-as-features.

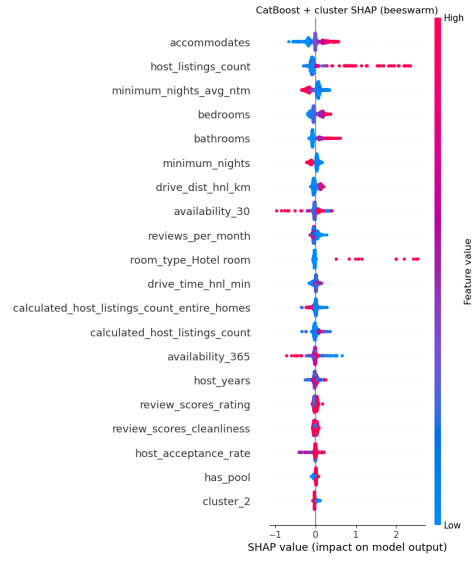


Figure A.3: SHAP summary for CatBoost with cluster-as-features.

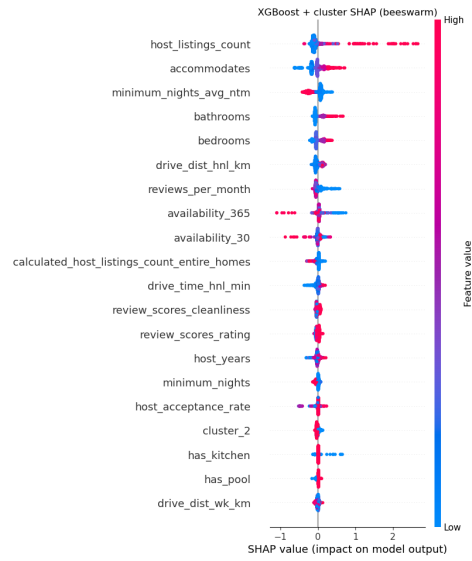


Figure A.4: SHAP summary for XGBoost with cluster-as-features.



## B Additional Tables

Model	Cluster	Test $R^2$	Test RMSE	Train $R^2$	Train RMSE
RF	0	0.681	0.349	0.921	0.204
LGBM	0	0.778	0.291	0.983	0.094
RF	1	0.513	1.083	0.979	0.265
LGBM	1	0.545	1.047	1.000	0.025
RF	2	0.771	0.428	0.944	0.205
LGBM	2	0.799	0.401	1.000	0.010
RF	3	0.885	0.621	0.964	0.356
LGBM	3	0.879	0.638	0.992	0.170
RF	4	0.588	0.418	0.927	0.185
LGBM	4	0.602	0.411	0.975	0.108
RF	5	0.677	0.519	0.919	0.261
LGBM	5	0.776	0.432	0.998	0.039

Table B.1: Cluster-level performance for the cluster-then-predict strategy

## C Code Appendix

```
1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # ### Clustering
5
6  # In[ ]:
7
8
9  import pandas as pd
10 import numpy as np
11
12 df = pd.read_csv("/Users/jiangzhanyuan/Desktop/second_year/IEOR242A/Final Project/
    Listing_Honolulu.csv")
13 df.head()
14
15
16 # In[3]:
17
18
19 features_for_clustering = [
20     "latitude",
21     "longitude",
22     "drive_dist_hnl_km",    # distance to HNL
23     "drive_dist_wk_km",    # distance to Waikiki
24     #"drive_time_hnl_min",
25     #"drive_time_wk_min",
26     "review_scores_rating",
27     "review_scores_cleanliness",
28     # "review_scores_location",
29     # "reviews_per_month",
30     "accommodates",
31     "bedrooms",
32     "bathrooms",
33     #"host_listings_count"
34     #"host_years"
35     #"price"
36
37
38
39     # "accommodates",
40     # "beds",
41     # "bedrooms",
42     # "bathrooms",
43     # "host_years",
44     # "host_listings_count",
45     # "calculated_host_listings_count",
46     # "calculated_host_listings_count_entire_homes",
47     # "host_response_rate",
48     # "host_acceptance_rate",
49     # "availability_365",
50     # "availability_30",
51     # "minimum_nights",
52     # "maximum_nights",
53     # "minimum_nights_avg_ntm",
54     # "maximum_nights_avg_ntm",
55     # "number_of_reviews",
56     # "number_of_reviews_ltm",
57     # "number_of_reviews_l30d",
58     # "reviews_per_month",
59     # "review_scores_rating",
60     # "review_scores_cleanliness",
61     # "review_scores_location",
62     # "review_scores_communication",
63     # "drive_dist_hnl_km",
64     # "drive_time_hnl_min",
65     # "drive_dist_wk_km",
66     # "drive_time_wk_min",
67     # "latitude",
68     # "longitude",
```

```

69 ]
70 ]
71
72 clust_df = df[features_for_clustering].copy()
73 clust_df.describe()
74
75
76 # In[4]:
77
78
79 extra_features = [
80     #"has_air_conditioning",
81     #"has_gym",
82     #"has_pool",
83
84     #"is_room_in_hotel",
85     #"is_entire_home"
86
87     #"host_is_superhost"
88 ]
89
90 features_for_clustering = features_for_clustering + extra_features
91 clust_df = df[features_for_clustering].copy()
92
93
94 # In[5]:
95
96
97 # Standardize the features
98 from sklearn.preprocessing import StandardScaler
99
100 scaler = StandardScaler()
101 X_scaled = scaler.fit_transform(clust_df)
102
103 X_scaled.shape
104
105
106 # In[6]:
107
108
109 from sklearn.cluster import AgglomerativeClustering
110 from sklearn.metrics import silhouette_score
111
112 sil_scores = {}
113
114 for k in range(2, 10):
115     model_k = AgglomerativeClustering(
116         n_clusters=k,
117         linkage="ward" # affinity/metric is fixed to euclidean for ward
118     )
119     labels_k = model_k.fit_predict(X_scaled)
120     score_k = silhouette_score(X_scaled, labels_k)
121     sil_scores[k] = score_k
122     print(f"k = {k}, silhouette_score = {score_k:.3f}")
123
124
125
126 # In[7]:
127
128
129 import matplotlib.pyplot as plt
130
131 ks = list(sil_scores.keys())
132 scores = [sil_scores[k] for k in ks]
133
134 plt.plot(ks, scores, marker="o")
135 plt.xlabel("Number of clusters k")
136 plt.ylabel("Silhouette score")
137 plt.title("Silhouette Score vs k")
138 plt.show()
139

```

```

140
141 # In[8]:
142
143
144 # choose k based on the silhouette plot
145 best_k = 3
146
147 cluster_model = AgglomerativeClustering(
148     n_clusters=best_k,
149     linkage="ward"
150 )
151
152 cluster_labels = cluster_model.fit_predict(X_scaled)
153
154 # Attach to original df
155 df["Cluster_ID"] = cluster_labels
156 df["Cluster_ID"].value_counts()
157
158
159 # ### K-means robustness check
160
161 # In[9]:
162
163
164 clust_df.describe()
165 clust_df.nunique()
166
167
168 # In[10]:
169
170
171 from sklearn.cluster import KMeans
172
173 kmeans = KMeans(n_clusters=best_k, random_state=42, n_init="auto")
174 km_labels = kmeans.fit_predict(X_scaled)
175
176 df["Cluster_ID_kmeans"] = km_labels
177 df[["Cluster_ID", "Cluster_ID_kmeans"]].head()
178
179
180 # In[11]:
181
182
183 clust_df.nunique()
184
185
186
187 # ## Cluster Profiling and Interpretation
188
189 # In[12]:
190
191
192 cluster_profile = df.groupby("Cluster_ID")[features_for_clustering + ["price"]].agg(
193     ["mean", "median"]
194 )
195
196 cluster_profile
197
198
199 #
200 # - Cluster 0 is the closest to waikiki
201 # - Cluster 1 has the highest average price
202
203 # In[ ]:
204
205
206
207
208
209 # In[13]:
210

```

```

211
212 import matplotlib.pyplot as plt
213
214 plt.figure(figsize=(6, 6))
215 scatter = plt.scatter(
216     df["longitude"],
217     df["latitude"],
218     c=df["Cluster_ID"],
219     alpha=0.6
220 )
221 plt.xlabel("Longitude")
222 plt.ylabel("Latitude")
223 plt.title("Honolulu Listings by Cluster")
224 plt.colorbar(scatter, label="Cluster_ID")
225 plt.show()
226
227
228 # In[14]:
229
230
231 import seaborn as sns
232
233 plt.figure(figsize=(16, 8))
234 sns.boxplot(
235     data=df,
236     x="Cluster_ID",
237     y="price"
238 )
239 plt.yscale("log") # price outliers log-scale
240 plt.title("Price Distribution by Cluster")
241 plt.show()
242
243
244 # In[15]:
245
246
247 output_cols = df.columns
248
249 #
250 df.to_csv("listing_honolulu_clustered.csv", index=False)
251
252
253 # ## PCA
254 #
255
256 # In[ ]:
257
258
259 from sklearn.decomposition import PCA
260 from sklearn.preprocessing import StandardScaler
261
262
263 num_cols = [
264
265     "latitude",
266     "longitude",
267     "drive_dist_hnl_km", # distance to HNL
268     "drive_dist_wk_km", # distance to Waikiki
269     "drive_time_hnl_min",
270     "drive_time_wk_min",
271     "review_scores_rating",
272     "review_scores_cleanliness",
273     # "review_scores_location",
274     # "reviews_per_month",
275     "accommodates",
276     "bedrooms",
277     "bathrooms",
278 ]
279
280
281 pca_df = df[num_cols].copy()

```

```

282 scaler_pca = StandardScaler()
283 X_pca = scaler_pca.fit_transform(pca_df)
284
285
286 pca_full = PCA(random_state=42)
287 pca_full.fit(X_pca)
288
289 explained = pca_full.explained_variance_ratio_
290 cum_explained = explained.cumsum()
291
292 k_85 = int(np.argmax(cum_explained >= 0.85) + 1)
293 k_90 = int(np.argmax(cum_explained >= 0.90) + 1)
294 k_95 = int(np.argmax(cum_explained >= 0.95) + 1)
295
296 explained_df = pd.DataFrame(
297     {
298         "component": np.arange(1, len(explained) + 1),
299         "explained_variance_ratio": explained,
300         "cumulative_variance_ratio": cum_explained,
301     }
302 )
303 explained_df.head()
304
305
306 # In[ ]:
307
308
309 fig, axes = plt.subplots(1, 2, figsize=(12, 4))
310
311
312 axes[0].bar(explained_df["component"], explained_df["explained_variance_ratio"], color
313             = "#4c72b0")
314 axes[0].set_xlabel("Principal Component")
315 axes[0].set_ylabel("Explained Variance Ratio")
316 axes[0].set_title("Each Principal Component's Contribution")
317 axes[0].set_xticks(explained_df["component"])
318 axes[0].tick_params(axis="x", rotation=45)
319
320 axes[1].plot(explained_df["component"], explained_df["cumulative_variance_ratio"],
321             marker="o")
322 axes[1].axhline(0.85, color="gray", linestyle="--", label="85%")
323 axes[1].axhline(0.90, color="red", linestyle="--", label="90%")
324 axes[1].axhline(0.95, color="green", linestyle="--", label="95%")
325 axes[1].axvline(k_85, color="gray", linestyle=":", alpha=0.7)
326 axes[1].axvline(k_90, color="red", linestyle=":", alpha=0.7)
327 axes[1].axvline(k_95, color="green", linestyle=":", alpha=0.7)
328 axes[1].annotate(f"k={k_85}", (k_85, cum_explained[k_85-1]), textcoords="offset points",
329                 xytext=(0,10), ha="center", color="gray")
330 axes[1].annotate(f"k={k_90}", (k_90, cum_explained[k_90-1]), textcoords="offset points",
331                 xytext=(0,10), ha="center", color="red")
332 axes[1].annotate(f"k={k_95}", (k_95, cum_explained[k_95-1]), textcoords="offset points",
333                 xytext=(0,10), ha="center", color="green")
334 axes[1].set_xlabel("Number of Components")
335 axes[1].set_ylabel("Cumulative Explained Variance")
336 axes[1].set_ylim(0, 1.05)
337 axes[1].set_title("Cumulated Contribution for first k principal components")
338 axes[1].legend()
339
340 plt.tight_layout()
341 plt.show()
342
343
344 print(f"85% of the total variance: {k_85}")
345 print(f"90% of the total variance: {k_90}")
346 print(f"95% of the total variance: {k_95}")
347
348
349 # ## Use first 4 principles to do hierachical clustering
350 #
351
352 # In[ ]:

```

```

348
349
350 from sklearn.cluster import AgglomerativeClustering
351 from sklearn.metrics import silhouette_score
352
353
354 pc_scores = pca_full.transform(X_pca)
355 X_pca_13 = pc_scores[:, :4]
356
357 k_range = range(2, 9)
358 sil_records = []
359 for k in k_range:
360     model = AgglomerativeClustering(n_clusters=k, linkage="ward")
361     labels = model.fit_predict(X_pca_13)
362     sil = silhouette_score(X_pca_13, labels)
363     sil_records.append({"k": k, "silhouette": sil})
364
365 sil_df = pd.DataFrame(sil_records)
366 sil_df
367
368
369 # In[ ]:
370
371
372 plt.figure(figsize=(6, 4))
373 plt.plot(sil_df["k"], sil_df["silhouette"], marker="o")
374 plt.xlabel("Number of clusters k")
375 plt.ylabel("Silhouette score")
376 plt.title("Silhouette vs k (Agglomerative on first 13 PCs)")
377 plt.grid(alpha=0.3)
378 plt.show()
379
380 best_k_13 = int(sil_df.sort_values("silhouette", ascending=False).iloc[0]["k"])
381 print(f"Silhouette best k: {best_k_13}")
382
383 final_model_13 = AgglomerativeClustering(n_clusters=best_k_13, linkage="ward")
384 cluster_labels_13 = final_model_13.fit_predict(X_pca_13)
385
386 df["Cluster_PCA13"] = cluster_labels_13
387 df["Cluster_PCA13"].value_counts().sort_index()
388
389
390 # In[ ]:
391
392
393 plt.figure(figsize=(6, 6))
394 scatter = plt.scatter(
395     df["longitude"],
396     df["latitude"],
397     c=df["Cluster_PCA13"],
398     alpha=0.6,
399     cmap="tab10",
400 )
401 plt.xlabel("Longitude")
402 plt.ylabel("Latitude")
403 plt.title(f"Cluster distribution (k={best_k_13}, first 6 PCs)")
404 plt.colorbar(scatter, label="Cluster_PCA13")
405 plt.show()
406
407
408 # In[ ]:
409
410
411 price_series = pd.to_numeric(df["price"], errors="coerce")
412
413 plt.figure(figsize=(10, 5))
414 sns.boxplot(data=df.assign(price_num=price_series), x="Cluster_PCA13", y="price_num")
415 plt.yscale("log")
416 plt.title("Price Distribution by Cluster (PCA13 Agglomerative)")
417 plt.xlabel("Cluster_PCA13")
418 plt.ylabel("Price (log scale)")

```

```

419 plt.show()
420
421
422
423 # In[ ]:

1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # ## Baseline 1. OLS (log price)
5  # - Y 'log1p(price)'
6  # - drop 'id', 'neighbourhood_group', 'latitude', 'longitude'
7  # - One-hot 'room_type', 'neighbourhood', 'host_response_time' drop_first
8
9  # In[16]:
10
11
12 import pandas as pd
13 import numpy as np
14 from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV
15 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
16 from sklearn.preprocessing import StandardScaler
17 from sklearn.linear_model import Ridge
18 from sklearn.pipeline import make_pipeline
19 from sklearn.ensemble import RandomForestRegressor
20 from lightgbm import LGBMRegressor
21 from xgboost import XGBRegressor
22 import statsmodels.api as sm
23
24 # Load dataset
25 DATA_PATH = "/Users/jiangzhanyuan/Desktop/second year/IEOR242A/Final Project/
    Listing_Honolulu.csv"
26 df = pd.read_csv(DATA_PATH)
27
28 print(f"Loaded dataset shape: {df.shape}")
29 print(df[['room_type', 'neighbourhood', 'host_response_time']].dtypes)
30
31
32
33 # In[7]:
34
35
36 # Price distribution by buckets
37 price = df["price"].astype(float)
38 print("Basic stats:")
39 print(price.describe(percentiles=[0.01, 0.05, 0.1, 0.25, 0.5, 0.75, 0.9, 0.95, 0.99]))
40
41 bins = [0, 50, 100, 150, 200, 300, 400, 500, 750, 1000, 1500, 2000, 3000, 5000, 10000,
    20000, 50000, np.inf]
42 price_bucket = pd.cut(price, bins=bins, right=False)
43 print("\nFrequency by price buckets (right-open intervals):")
44 print(price_bucket.value_counts().sort_index())
45
46
47
48 # In[8]:
49
50
51 # Baseline 1: OLS on log(price)
52 target = "price"
53 drop_cols = ["id", "neighbourhood_group", "latitude", "longitude"]
54 categorical_cols = ["room_type", "neighbourhood", "host_response_time"]
55
56 model_df = df.copy()
57 cols_to_drop = [c for c in drop_cols if c in model_df.columns]
58 X = model_df.drop(columns=cols_to_drop + [target])
59 y_price = model_df[target].astype(float)
60 y = np.log1p(y_price)
61
62 X_enc = pd.get_dummies(X, columns=categorical_cols, drop_first=True, dtype=float)
63 X_enc = X_enc.astype(float)

```



```

64
65 # 80/20 shuffle split train vs test
66 X_train, X_test, y_train, y_test = train_test_split(
67     X_enc, y, test_size=0.2, shuffle=True, random_state=42
68 )
69
70 X_train_const = sm.add_constant(X_train, has_constant="add")
71 X_test_const = sm.add_constant(X_test, has_constant="add")
72
73 ols_model = sm.OLS(y_train, X_train_const).fit()
74
75 print(f"Train rows: {len(X_train)}, Test rows: {len(X_test)}")
76 print(f"Encoded feature count: {X_enc.shape[1]}")
77 print("Target: log(price)")
78
79
80
81 # In[ ]:
82
83
84 # Evaluate (log space + original price space)
85
86 # Train predictions (log space)
87 train_pred_log = ols_model.predict(X_train_const)
88 train_mse_log = mean_squared_error(y_train, train_pred_log, squared=True)
89 train_rmse_log = mean_squared_error(y_train, train_pred_log, squared=False)
90 train_mae_log = mean_absolute_error(y_train, train_pred_log)
91 train_r2 = ols_model.rsquared
92 train_r2_adj = ols_model.rsquared_adj
93
94 # Test predictions (log space)
95 test_pred_log = ols_model.predict(X_test_const)
96 test_mse_log = mean_squared_error(y_test, test_pred_log, squared=True)
97 test_rmse_log = mean_squared_error(y_test, test_pred_log, squared=False)
98 test_mae_log = mean_absolute_error(y_test, test_pred_log)
99 test_r2_log = r2_score(y_test, test_pred_log)
100
101 test_osr2_log = 1 - ((y_test - test_pred_log) ** 2).sum() / ((y_test - y_train.mean())
102     ** 2).sum()
103
104 # Back-transform to price space
105 train_pred_price = np.exp1(train_pred_log)
106 test_pred_price = np.exp1(test_pred_log)
107 train_price = np.exp1(y_train)
108 test_price = np.exp1(y_test)
109
110 train_mse_price = mean_squared_error(train_price, train_pred_price, squared=True)
111 train_rmse_price = mean_squared_error(train_price, train_pred_price, squared=False)
112 train_mae_price = mean_absolute_error(train_price, train_pred_price)
113
114 test_mse_price = mean_squared_error(test_price, test_pred_price, squared=True)
115 test_rmse_price = mean_squared_error(test_price, test_pred_price, squared=False)
116 test_mae_price = mean_absolute_error(test_price, test_pred_price)
117 test_r2_price = r2_score(test_price, test_pred_price)
118
119 test_osr2_price = 1 - ((test_price - test_pred_price) ** 2).sum() / ((test_price -
120     train_price.mean()) ** 2).sum()
121
122 print("Log-space metrics:")
123 print(
124     f"Train MSE: {train_mse_log:0.4f} | RMSE: {train_rmse_log:0.4f} | MAE: {
125         train_mae_log:0.4f} | R^2: {train_r2:0.3f} | Adj R^2: {train_r2_adj:0.3f}"
126 )
127 print(
128     f"Test MSE: {test_mse_log:0.4f} | RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log
129         :0.4f} | R^2: {test_r2_log:0.3f} | OSR^2: {test_osr2_log:0.3f}"
130 )
131
132 print("\nOriginal-price-space metrics:")
133 print(
134     f"Train MSE: {train_mse_price:0.2f} | RMSE: {train_rmse_price:0.2f} | MAE: {

```

```

        train_mae_price:0.2f}"
131 )
132 print(
133     f"Test MSE: {test_mse_price:0.2f} | RMSE: {test_rmse_price:0.2f} | MAE: {
        test_mae_price:0.2f} | R^2: {test_r2_price:0.3f} | OSR^2: {test_osr2_price:0.3f
        }"
134 )
135
136 # Top coefficients by absolute value (log-price model)
137 coef_table = (
138     ols_model.params.rename("coef")
139     .reset_index()
140     .rename(columns={"index": "feature"})
141     .assign(abs_coef=lambda d: d["coef"].abs())
142     .sort_values("abs_coef", ascending=False)
143 )
144 print("\nTop 15 coefficients by |coef|:")
145 print(coef_table.head(15))
146
147 print("\nOLS summary (truncated):")
148 print(ols_model.summary())
149
150
151
152 # In[76]:
153
154
155 # Export OLS report and coefficients (log-price model)
156 report_path = "/Users/jiangzhanyuan/Desktop/second year/IEOR242A/Final Project/
    ols_log_report.txt"
157 coef_path = "/Users/jiangzhanyuan/Desktop/second year/IEOR242A/Final Project/
    ols_log_coefs.csv"
158 metrics_path = "/Users/jiangzhanyuan/Desktop/second year/IEOR242A/Final Project/
    ols_log_metrics.json"
159
160 with open(report_path, "w") as f:
161     f.write(ols_model.summary().as_text())
162
163 coef_table.to_csv(coef_path, index=False)
164
165 metrics = {
166     # log-space
167     "train_mse_log": float(train_mse_log),
168     "train_rmse_log": float(train_rmse_log),
169     "train_mae_log": float(train_mae_log),
170     "train_r2_log": float(train_r2),
171     "train_r2_adj_log": float(train_r2_adj),
172     "test_mse_log": float(test_mse_log),
173     "test_rmse_log": float(test_rmse_log),
174     "test_mae_log": float(test_mae_log),
175     "test_r2_log": float(test_r2_log),
176     "test_osr2_log": float(test_osr2_log),
177     # price-space (back-transformed)
178     "train_mse_price": float(train_mse_price),
179     "train_rmse_price": float(train_rmse_price),
180     "train_mae_price": float(train_mae_price),
181     "test_mse_price": float(test_mse_price),
182     "test_rmse_price": float(test_rmse_price),
183     "test_mae_price": float(test_mae_price),
184     "test_r2_price": float(test_r2_price),
185     "test_osr2_price": float(test_osr2_price),
186     # data shape
187     "n_train": int(len(X_train)),
188     "n_test": int(len(X_test)),
189     "n_features_encoded": int(X_enc.shape[1]),
190 }
191 import json
192 with open(metrics_path, "w") as f:
193     json.dump(metrics, f, indent=2)
194
195 print(f"Saved summary to: {report_path}")

```

```

196 print(f"Saved coefficients to: {coef_path}")
197 print(f"Saved metrics to: {metrics_path}")
198
199
200
201 # ## Baseline 2. Ridge (log price)
202 # - Y 'log1p(price)'
203 # - drop 'id', 'neighbourhood_group', 'latitude', 'longitude'
204 # - One-hot 'room_type', 'neighbourhood', 'host_response_time' drop_first
205 #
206 #
207 #
208
209 # In[79]:
210
211
212 # Ridge (fixed alpha)
213 ridge_drop_cols = ["id", "neighbourhood_group", "latitude", "longitude"]
214 ridge_cat_cols = ["room_type", "neighbourhood", "host_response_time"]
215
216 ridge_df = df.copy()
217 ridge_cols_to_drop = [c for c in ridge_drop_cols if c in ridge_df.columns]
218 X_ridge = ridge_df.drop(columns=ridge_cols_to_drop + ["price"])
219 y_price_ridge = ridge_df["price"].astype(float)
220 y_log_ridge = np.log1p(y_price_ridge)
221
222 X_ridge_enc = pd.get_dummies(X_ridge, columns=ridge_cat_cols, drop_first=True, dtype=
float).astype(float)
223
224 X_train_r, X_test_r, y_train_r, y_test_r = train_test_split(
225     X_ridge_enc, y_log_ridge, test_size=0.2, shuffle=True, random_state=42
226 )
227
228 ridge_alpha = 10.0
229 ridge_model = make_pipeline(StandardScaler(with_mean=False), Ridge(alpha=ridge_alpha))
230 ridge_model.fit(X_train_r, y_train_r)
231
232 y_train_pred_log = ridge_model.predict(X_train_r)
233 y_test_pred_log = ridge_model.predict(X_test_r)
234
235 ridge_train_rmse_log = mean_squared_error(y_train_r, y_train_pred_log, squared=False)
236 ridge_test_rmse_log = mean_squared_error(y_test_r, y_test_pred_log, squared=False)
237 ridge_train_mae_log = mean_absolute_error(y_train_r, y_train_pred_log)
238 ridge_test_mae_log = mean_absolute_error(y_test_r, y_test_pred_log)
239 ridge_train_r2_log = r2_score(y_train_r, y_train_pred_log)
240 ridge_test_r2_log = r2_score(y_test_r, y_test_pred_log)
241
242 print("Ridge (alpha=10) log-space metrics:")
243 print(
244     f"Train RMSE: {ridge_train_rmse_log:0.4f} | MAE: {ridge_train_mae_log:0.4f} | R^2:
{ridge_train_r2_log:0.3f}"
245 )
246 print(
247     f"Test RMSE: {ridge_test_rmse_log:0.4f} | MAE: {ridge_test_mae_log:0.4f} | R^2: {
ridge_test_r2_log:0.3f}"
248 )
249
250 ridge_coefs = pd.DataFrame({
251     "feature": X_ridge_enc.columns,
252     "coef": ridge_model.named_steps["ridge"].coef_,
253 }).assign(abs_coef=lambda d: d.coef.abs()).sort_values("abs_coef", ascending=False)
254 print("\nTop 15 Ridge coefficients by |coef| (log-price space):")
255 print(ridge_coefs.head(15))
256
257
258
259 # ## Baseline 2a. Ridge CV (log price)
260 # - Same features as above 80 / 20 Train/Test split
261 # - Train set 5-fold CV Search alpha
262 #
263 #

```

```

264
265 # In[80]:
266
267
268 # Ridge CV
269 ridge_df_cv = df.copy()
270 ridge_cols_to_drop_cv = [c for c in ["id", "neighbourhood_group", "latitude", "
    longitude"] if c in ridge_df_cv.columns]
271 X_ridge_cv = ridge_df_cv.drop(columns=ridge_cols_to_drop_cv + ["price"])
272 y_log_ridge_cv = np.log1p(ridge_df_cv["price"].astype(float))
273
274 X_ridge_enc_cv = pd.get_dummies(
275     X_ridge_cv,
276     columns=["room_type", "neighbourhood", "host_response_time"],
277     drop_first=True,
278     dtype=float,
279 ).astype(float)
280
281 X_train_cv, X_test_cv, y_train_cv, y_test_cv = train_test_split(
282     X_ridge_enc_cv, y_log_ridge_cv, test_size=0.2, shuffle=True, random_state=42
283 )
284
285 alphas = np.logspace(-3, 3, 13)
286 ridge_cv_model = make_pipeline(StandardScaler(with_mean=False), Ridge())
287 param_grid = {"ridge__alpha": alphas}
288 search = GridSearchCV(ridge_cv_model, param_grid, cv=5, scoring="
    neg_mean_squared_error", n_jobs=1)
289 search.fit(X_train_cv, y_train_cv)
290
291 best_alpha = search.best_params_["ridge__alpha"]
292 best_model = search.best_estimator_
293 print(f"Best alpha: {best_alpha}")
294
295 y_train_pred_log = best_model.predict(X_train_cv)
296 y_test_pred_log = best_model.predict(X_test_cv)
297
298 train_rmse_log = mean_squared_error(y_train_cv, y_train_pred_log, squared=False)
299 test_rmse_log = mean_squared_error(y_test_cv, y_test_pred_log, squared=False)
300 train_mae_log = mean_absolute_error(y_train_cv, y_train_pred_log)
301 test_mae_log = mean_absolute_error(y_test_cv, y_test_pred_log)
302 train_r2_log = r2_score(y_train_cv, y_train_pred_log)
303 test_r2_log = r2_score(y_test_cv, y_test_pred_log)
304
305 print("Ridge CV log-space metrics:")
306 print(
307     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
        train_r2_log:0.3f}"
308 )
309 print(
310     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
        :0.3f}"
311 )
312
313 ridge_cv_coefs = pd.DataFrame({
314     "feature": X_ridge_enc_cv.columns,
315     "coef": best_model.named_steps["ridge"].coef_,
316 }).assign(abs_coef=lambda d: d.coef.abs()).sort_values("abs_coef", ascending=False)
317 print("\nTop 15 Ridge-CV coefficients by |coef| (log-price space):")
318 print(ridge_cv_coefs.head(15))
319
320
321
322 # ## Baseline 3. Random Forest (log price)
323 # - Y 'log1p(price)'
324 # - drop 'id', 'neighbourhood_group', 'latitude', 'longitude'
325 # - One-hot 'room_type', 'neighbourhood', 'host_response_time' drop_first
326 # - 80/20 random split GridSearchCV
327 #
328 #
329
330 # In[11]:

```

```

331
332
333 # Random Forest (fixed hyperparams)
334 rf_drop_cols = ["id", "neighbourhood_group", "latitude", "longitude"]
335 rf_cat_cols = ["room_type", "neighbourhood", "host_response_time"]
336
337 rf_df = df.copy()
338 rf_cols_to_drop = [c for c in rf_drop_cols if c in rf_df.columns]
339 X_rf = rf_df.drop(columns=rf_cols_to_drop + ["price"])
340 y_rf_log = np.log1p(rf_df["price"].astype(float))
341
342 X_rf_enc = pd.get_dummies(X_rf, columns=rf_cat_cols, drop_first=True, dtype=float).
    astype(float)
343
344 X_train_rf, X_test_rf, y_train_rf, y_test_rf = train_test_split(
345     X_rf_enc, y_rf_log, test_size=0.2, shuffle=True, random_state=42
346 )
347
348 rf_model = RandomForestRegressor(
349     n_estimators=50,
350     max_depth=5,
351     min_samples_split=10,
352     min_samples_leaf=4,
353     bootstrap=True,
354     random_state=42,
355     n_jobs=-1,
356 )
357 rf_model.fit(X_train_rf, y_train_rf)
358
359 train_pred_log = rf_model.predict(X_train_rf)
360 test_pred_log = rf_model.predict(X_test_rf)
361
362 train_mse_log = mean_squared_error(y_train_rf, train_pred_log)
363 test_mse_log = mean_squared_error(y_test_rf, test_pred_log)
364 train_rmse_log = np.sqrt(train_mse_log)
365 test_rmse_log = np.sqrt(test_mse_log)
366 train_mae_log = mean_absolute_error(y_train_rf, train_pred_log)
367 test_mae_log = mean_absolute_error(y_test_rf, test_pred_log)
368 train_r2_log = r2_score(y_train_rf, train_pred_log)
369 test_r2_log = r2_score(y_test_rf, test_pred_log)
370
371 print("RandomForest (fixed) log-space metrics:")
372 print(
373     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
        train_r2_log:0.3f}"
374 )
375 print(
376     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
        :0.3f}"
377 )
378
379
380
381 # In[ ]:
382
383
384 # Random Forest GridSearchCV (log price)
385 param_grid = {
386     "n_estimators": [50, 100, 200],
387     "max_depth": [10, 20, 40],
388     "min_samples_split": [2, 5, 10],
389     "min_samples_leaf": [2, 4, 8],
390 }
391
392 rf_base = RandomForestRegressor(
393     bootstrap=True,
394     random_state=42,
395     n_jobs=-1,
396 )
397
398 rf_search = GridSearchCV(

```

```

399     rf_base,
400     param_grid,
401     cv=5,
402     scoring="neg_mean_squared_error",
403     n_jobs=-1,
404     return_train_score=False,
405 )
406 rf_search.fit(X_train_rf, y_train_rf)
407
408 best_params = rf_search.best_params_
409 best_rf = rf_search.best_estimator_
410 print("Best params:", best_params)
411
412 y_train_pred_log = best_rf.predict(X_train_rf)
413 y_test_pred_log = best_rf.predict(X_test_rf)
414
415 train_mse_log = mean_squared_error(y_train_rf, y_train_pred_log)
416 test_mse_log = mean_squared_error(y_test_rf, y_test_pred_log)
417 train_rmse_log = np.sqrt(train_mse_log)
418 test_rmse_log = np.sqrt(test_mse_log)
419 train_mae_log = mean_absolute_error(y_train_rf, y_train_pred_log)
420 test_mae_log = mean_absolute_error(y_test_rf, y_test_pred_log)
421 train_r2_log = r2_score(y_train_rf, y_train_pred_log)
422 test_r2_log = r2_score(y_test_rf, y_test_pred_log)
423
424 print("RandomForest (GridSearchCV) log-space metrics:")
425 print(
426     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
427         train_r2_log:0.3f}"
428 )
429 print(
430     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
431         :0.3f}"
432 )
433
434 # ## Baseline 4. LightGBM (log price) + GridSearchCV
435 # - Y 'log1p(price)'
436 # - Same features as above drop 'id', 'neighbourhood_group', 'latitude', 'longitude'
437 #   one -hot 'room_type/neighbourhood/host_response_time'
438 # - 80/20 Train/ Test 5-fold GridSearchCV
439 #
440
441 # In[ ]:
442
443
444 from lightgbm import LGBMRegressor
445 from sklearn.model_selection import train_test_split, GridSearchCV
446 # LightGBM Regressor with GridSearchCV (log price)
447 param_grid_lgbm = {
448     "n_estimators": [200, 500, 800],
449     "num_leaves": [31, 63, 127],
450     "max_depth": [-1, 10, 20],
451     "learning_rate": [0.05, 0.1],
452     "subsample": [0.8, 1.0],
453     "colsample_bytree": [0.8, 1.0],
454     "min_child_samples": [10, 20, 40],
455 }
456
457 lgbm = LGBMRegressor(
458     objective="regression",
459     random_state=42,
460     n_jobs=-1,
461 )
462
463 gs_lgbm = GridSearchCV(
464     lgbm,
465     param_grid_lgbm,
466     cv=5,

```

```

467     scoring="neg_mean_squared_error",
468     n_jobs=-1,
469     return_train_score=False,
470 )
471 gs_lgbm.fit(X_train_rf, y_train_rf)
472
473 print("Best params (LightGBM):", gs_lgbm.best_params_)
474 best_lgbm = gs_lgbm.best_estimator_
475
476 y_train_pred_log = best_lgbm.predict(X_train_rf)
477 y_test_pred_log = best_lgbm.predict(X_test_rf)
478
479 train_rmse_log = mean_squared_error(y_train_rf, y_train_pred_log, squared=False)
480 test_rmse_log = mean_squared_error(y_test_rf, y_test_pred_log, squared=False)
481 train_mae_log = mean_absolute_error(y_train_rf, y_train_pred_log)
482 test_mae_log = mean_absolute_error(y_test_rf, y_test_pred_log)
483 train_r2_log = r2_score(y_train_rf, y_train_pred_log)
484 test_r2_log = r2_score(y_test_rf, y_test_pred_log)
485
486 print("LightGBM (GridSearchCV) log-space metrics:")
487 print(
488     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
489         train_r2_log:0.3f}"
490 )
491 print(
492     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
493         :0.3f}"
494 )
495
496 # ## Random Forest (log price)      RandomizedSearchCV (faster)
497 # - Same features as above and same Y 80 /20
498 # - RandomizedSearchCV to save time
499 #
500 #
501
502 # In[ ]:
503
504
505 from sklearn.model_selection import RandomizedSearchCV
506
507 # Randomized search for RF (faster)
508 rand_param_dist = {
509     "n_estimators": [50, 100, 200],
510     "max_depth": [10, 20, 40],
511     "min_samples_split": [2, 5, 10],
512     "min_samples_leaf": [2, 4, 8],
513 }
514
515 rf_base_rand = RandomForestRegressor(
516     bootstrap=True,
517     random_state=42,
518     n_jobs=-1,
519 )
520
521 rf_rand = RandomizedSearchCV(
522     rf_base_rand,
523     rand_param_dist,
524     n_iter=20,
525     cv=5,
526     scoring="neg_mean_squared_error",
527     n_jobs=-1,
528     return_train_score=False,
529     random_state=42,
530 )
531 rf_rand.fit(X_train_rf, y_train_rf)
532
533 best_params_rand = rf_rand.best_params_
534 best_rf_rand = rf_rand.best_estimator_
535 print("Best params (RandomizedSearchCV):", best_params_rand)

```

```

536
537 y_train_pred_log = best_rf_rand.predict(X_train_rf)
538 y_test_pred_log = best_rf_rand.predict(X_test_rf)
539
540 train_mse_log = mean_squared_error(y_train_rf, y_train_pred_log)
541 test_mse_log = mean_squared_error(y_test_rf, y_test_pred_log)
542 train_rmse_log = np.sqrt(train_mse_log)
543 test_rmse_log = np.sqrt(test_mse_log)
544 train_mae_log = mean_absolute_error(y_train_rf, y_train_pred_log)
545 test_mae_log = mean_absolute_error(y_test_rf, y_test_pred_log)
546 train_r2_log = r2_score(y_train_rf, y_train_pred_log)
547 test_r2_log = r2_score(y_test_rf, y_test_pred_log)
548
549 print("RandomForest (RandomizedSearchCV) log-space metrics:")
550 print(
551     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
552         train_r2_log:0.3f}"
553 )
554 print(
555     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
556         :0.3f}"
557 )
558
559 # ## LightGBM (log price) RandomizedSearchCV (faster)
560 #
561 #
562 #
563
564 # In[ ]:
565
566
567 # LightGBM with RandomizedSearchCV (log price)
568 param_dist_lgbm = {
569     "n_estimators": [200, 400, 800],
570     "num_leaves": [31, 63, 127],
571     "max_depth": [-1, 10, 20],
572     "learning_rate": [0.05, 0.1],
573     "subsample": [0.8, 1.0],
574     "colsample_bytree": [0.8, 1.0],
575     "min_child_samples": [10, 20, 40],
576 }
577
578 lgbm_base = LGBMRegressor(
579     objective="regression",
580     random_state=42,
581     n_jobs=-1,
582 )
583
584 lgbm_rand = RandomizedSearchCV(
585     lgbm_base,
586     param_distributions=param_dist_lgbm,
587     n_iter=25,
588     cv=5,
589     scoring="neg_mean_squared_error",
590     n_jobs=-1,
591     return_train_score=False,
592     random_state=42,
593 )
594 lgbm_rand.fit(X_train_rf, y_train_rf)
595
596 print("Best params (LGBM RandomizedSearchCV):", lgbm_rand.best_params_)
597 best_lgbm_rand = lgbm_rand.best_estimator_
598
599 y_train_pred_log = best_lgbm_rand.predict(X_train_rf)
600 y_test_pred_log = best_lgbm_rand.predict(X_test_rf)
601
602 train_mse_log = mean_squared_error(y_train_rf, y_train_pred_log)
603 test_mse_log = mean_squared_error(y_test_rf, y_test_pred_log)
604 train_rmse_log = np.sqrt(train_mse_log)

```



```

605 test_rmse_log = np.sqrt(test_mse_log)
606 train_mae_log = mean_absolute_error(y_train_rf, y_train_pred_log)
607 test_mae_log = mean_absolute_error(y_test_rf, y_test_pred_log)
608 train_r2_log = r2_score(y_train_rf, y_train_pred_log)
609 test_r2_log = r2_score(y_test_rf, y_test_pred_log)
610
611 print("LightGBM (RandomizedSearchCV) log-space metrics:")
612 print(
613     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
        train_r2_log:0.3f}"
614 )
615 print(
616     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
        :0.3f}"
617 )
618
619
620
621 # ## Baseline 5. XGBoost (log price)      RandomizedSearchCV
622 #
623 #
624 #
625
626 # In[ ]:
627
628
629 # XGBoost with RandomizedSearchCV (log price)
630 param_dist_xgb = {
631     "n_estimators": [100, 200, 400, 800],
632     "max_depth": [2, 3, 5, 7, 10],
633     "learning_rate": [0.03, 0.05, 0.1, 0.2],
634     "subsample": [0.7, 0.85, 1.0],
635     "colsample_bytree": [0.7, 0.85, 1.0],
636     "min_child_weight": [1, 3, 5],
637     "gamma": [0, 0.1, 0.3],
638 }
639
640 xgb_base = XGBRegressor(
641     objective="reg:squarederror",
642     random_state=42,
643     n_jobs=-1,
644     tree_method="hist",
645 )
646
647 xgb_rand = RandomizedSearchCV(
648     estimator=xgb_base,
649     param_distributions=param_dist_xgb,
650     n_iter=25,
651     cv=5,
652     scoring="neg_mean_squared_error",
653     n_jobs=-1,
654     return_train_score=False,
655     random_state=42,
656 )
657
658 xgb_rand.fit(X_train_rf, y_train_rf)
659
660 print("Best params (XGBoost RandomizedSearchCV):", xgb_rand.best_params_)
661 best_xgb = xgb_rand.best_estimator_
662
663 y_train_pred_log = best_xgb.predict(X_train_rf)
664 y_test_pred_log = best_xgb.predict(X_test_rf)
665
666 train_mse_log = mean_squared_error(y_train_rf, y_train_pred_log)
667 test_mse_log = mean_squared_error(y_test_rf, y_test_pred_log)
668 train_rmse_log = np.sqrt(train_mse_log)
669 test_rmse_log = np.sqrt(test_mse_log)
670 train_mae_log = mean_absolute_error(y_train_rf, y_train_pred_log)
671 test_mae_log = mean_absolute_error(y_test_rf, y_test_pred_log)
672 train_r2_log = r2_score(y_train_rf, y_train_pred_log)
673 test_r2_log = r2_score(y_test_rf, y_test_pred_log)

```

```

674
675 print("XGBoost (RandomizedSearchCV) log-space metrics:")
676 print(
677     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
        train_r2_log:0.3f}"
678 )
679 print(
680     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
        :0.3f}"
681 )
682
683
684
685 # ## Baseline 6. CatBoost (log price) RandomizedSearchCV
686 # - 'id', 'neighbourhood_group', 'latitude', 'longitude'
        one -hot 'room_type/neighbourhood/host_response_time'
687 # - 'log1p(price)' 80 /20
688 # -
689
690 # In[ ]:
691
692
693 from catboost import CatBoostRegressor
694 from sklearn.model_selection import RandomizedSearchCV
695
696 # CatBoost with RandomizedSearchCV (log price)
697 cat_param_dist = {
698     "n_estimators": [300, 500, 800, 1200],
699     "depth": [4, 6, 8, 10],
700     "learning_rate": [0.03, 0.05, 0.1],
701     "subsample": [0.7, 0.85, 1.0],
702     "colsample_bylevel": [0.7, 0.85, 1.0],
703     "l2_leaf_reg": [1, 3, 5, 7, 10],
704 }
705
706 cat_base = CatBoostRegressor(
707     loss_function="RMSE",
708     eval_metric="RMSE",
709     random_state=42,
710     verbose=False,
711     thread_count=-1,
712 )
713
714 cat_rand = RandomizedSearchCV(
715     estimator=cat_base,
716     param_distributions=cat_param_dist,
717     n_iter=25,
718     cv=5,
719     scoring="neg_mean_squared_error",
720     n_jobs=1,
721     return_train_score=False,
722     random_state=42,
723 )
724
725 cat_rand.fit(X_train_rf, y_train_rf)
726
727 print("Best params (CatBoost RandomizedSearchCV):", cat_rand.best_params_)
728 best_cat = cat_rand.best_estimator_
729
730 y_train_pred_log = best_cat.predict(X_train_rf)
731 y_test_pred_log = best_cat.predict(X_test_rf)
732
733 train_mse_log = mean_squared_error(y_train_rf, y_train_pred_log)
734 test_mse_log = mean_squared_error(y_test_rf, y_test_pred_log)
735 train_rmse_log = np.sqrt(train_mse_log)
736 test_rmse_log = np.sqrt(test_mse_log)
737 train_mae_log = mean_absolute_error(y_train_rf, y_train_pred_log)
738 test_mae_log = mean_absolute_error(y_test_rf, y_test_pred_log)
739 train_r2_log = r2_score(y_train_rf, y_train_pred_log)
740 test_r2_log = r2_score(y_test_rf, y_test_pred_log)
741

```

```

742 print("CatBoost (RandomizedSearchCV) log-space metrics:")
743 print(
744     f"Train RMSE: {train_rmse_log:0.4f} | MAE: {train_mae_log:0.4f} | R^2: {
          train_r2_log:0.3f}"
745 )
746 print(
747     f"Test RMSE: {test_rmse_log:0.4f} | MAE: {test_mae_log:0.4f} | R^2: {test_r2_log
          :0.3f}"
748 )

1  import pandas as pd
2  import numpy as np
3
4  from sklearn.model_selection import train_test_split, RandomizedSearchCV
5  from sklearn.preprocessing import StandardScaler
6  from sklearn.cluster import AgglomerativeClustering
7  from sklearn.metrics import silhouette_score, mean_squared_error, mean_absolute_error,
   r2_score
8  from sklearn.ensemble import RandomForestRegressor
9  from lightgbm import LGBMRegressor
10
11 from scipy.spatial.distance import cdist
12
13
14 # Paths and feature definitions (mirrors baseline + clustering notebooks)
15 DATA_PATH = "/Users/jiangzhanyuan/Desktop/second_year/IEOR242A/Final Project/
   Listing_Honolulu.csv"
16
17 cluster_features = [
18     "latitude",
19     "longitude",
20     "drive_dist_hnl_km",
21     "drive_dist_wk_km",
22     "review_scores_rating",
23     "review_scores_cleanliness",
24     "accommodates",
25     "bedrooms",
26     "bathrooms",
27 ]
28
29 # Baseline feature space (69 columns after one-hot with drop_first)
30 drop_cols = ["id", "neighbourhood_group", "latitude", "longitude"]
31 cat_cols = ["room_type", "neighbourhood", "host_response_time"]
32
33 # Target
34 TARGET = "price"
35
36 df = pd.read_csv(DATA_PATH)
37 print(f"Loaded dataset: {df.shape}")
38
39
40 def rmse(y_true, y_pred):
41     # Older sklearn may not support squared=False
42     return np.sqrt(mean_squared_error(y_true, y_pred))
43
44
45 def mae(y_true, y_pred):
46     return mean_absolute_error(y_true, y_pred)
47
48
49 def r2(y_true, y_pred):
50     return r2_score(y_true, y_pred)
51
52
53 def assign_to_centroid(X_scaled, centroids):
54     """Assign each row in X_scaled to nearest centroid (Euclidean)."""
55     distances = cdist(X_scaled, centroids, metric="euclidean")
56     return distances.argmin(axis=1)
57
58
59 # Prepare baseline feature matrix (69 cols) and log-price target

```

```

60 model_df = df.copy()
61 cols_to_drop = [c for c in drop_cols if c in model_df.columns]
62 X_raw = model_df.drop(columns=cols_to_drop + [TARGET])
63 y_log = np.log1p(model_df[TARGET].astype(float))
64
65 X_enc = pd.get_dummies(X_raw, columns=cat_cols, drop_first=True, dtype=float).astype(
    float)
66 print(f"Encoded feature count: {X_enc.shape[1]}")
67
68 # Train/test split must mirror baseline (80/20, shuffle=True, random_state=42)
69 X_train, X_test, y_train, y_test = train_test_split(
70     X_enc, y_log, test_size=0.2, shuffle=True, random_state=42
71 )
72 print(f"Train rows: {len(X_train)}, Test rows: {len(X_test)}")
73
74
75 # ---- Clustering (train-only fit) ----
76 clust_df = df[cluster_features].copy()
77 clust_train = clust_df.loc[X_train.index]
78 clust_test = clust_df.loc[X_test.index]
79
80 scaler = StandardScaler()
81 X_clust_train_scaled = scaler.fit_transform(clust_train)
82 X_clust_test_scaled = scaler.transform(clust_test)
83
84 # Silhouette check on train (k=2..9) to mirror notebook
85 sil_scores = {}
86 for k in range(2, 10):
87     model_k = AgglomerativeClustering(n_clusters=k, linkage="ward")
88     labels_k = model_k.fit_predict(X_clust_train_scaled)
89     sil_scores[k] = silhouette_score(X_clust_train_scaled, labels_k)
90
91 best_k = max(sil_scores, key=sil_scores.get)
92 print("Silhouette scores (train):", sil_scores)
93 print(f"Best k (train silhouette): {best_k}")
94
95 cluster_model = AgglomerativeClustering(n_clusters=best_k, linkage="ward")
96 train_cluster_labels = cluster_model.fit_predict(X_clust_train_scaled)
97
98 # Compute centroids in scaled space and assign test by nearest centroid
99 centroids = np.vstack([
100     X_clust_train_scaled[train_cluster_labels == c].mean(axis=0)
101     for c in range(best_k)
102 ])
103 test_cluster_labels = assign_to_centroid(X_clust_test_scaled, centroids)
104
105 print("Train cluster counts:", pd.Series(train_cluster_labels).value_counts().to_dict())
106 print("Test cluster counts:", pd.Series(test_cluster_labels).value_counts().to_dict())
107
108
109 # Hyperparameter search spaces (align with baseline RandomizedSearchCV choices)
110 rf_param_dist = {
111     "n_estimators": [50, 100, 200],
112     "max_depth": [10, 20, 40],
113     "min_samples_split": [2, 5, 10],
114     "min_samples_leaf": [2, 4, 8],
115 }
116
117 lgbm_param_dist = {
118     "n_estimators": [200, 400, 800],
119     "num_leaves": [31, 63, 127],
120     "max_depth": [-1, 10, 20],
121     "learning_rate": [0.05, 0.1],
122     "subsample": [0.8, 1.0],
123     "colsample_bytree": [0.8, 1.0],
124     "min_child_samples": [10, 20, 40],
125 }
126
127
128 # Containers

```

```

129 cluster_metrics = []
130 rf_preds_test = pd.Series(index=X_test.index, dtype=float)
131 lgbm_preds_test = pd.Series(index=X_test.index, dtype=float)
132
133 for cluster_id in range(best_k):
134     train_mask = pd.Series(train_cluster_labels, index=X_train.index) == cluster_id
135     test_mask = pd.Series(test_cluster_labels, index=X_test.index) == cluster_id
136
137     X_train_c = X_train.loc[train_mask]
138     y_train_c = y_train.loc[train_mask]
139     X_test_c = X_test.loc[test_mask]
140     y_test_c = y_test.loc[test_mask]
141
142     print(f"\nCluster {cluster_id}: train {len(X_train_c)}, test {len(X_test_c)}")
143
144     # Random Forest
145     rf_base = RandomForestRegressor(bootstrap=True, random_state=42, n_jobs=-1)
146     rf_rand = RandomizedSearchCV(
147         rf_base,
148         rf_param_dist,
149         n_iter=20,
150         cv=5,
151         scoring="neg_mean_squared_error",
152         n_jobs=-1,
153         random_state=42,
154         return_train_score=False,
155     )
156     rf_rand.fit(X_train_c, y_train_c)
157     rf_train_pred = rf_rand.predict(X_train_c)
158     rf_test_pred = rf_rand.predict(X_test_c) if len(X_test_c) else np.array([])
159
160     rf_metrics = {
161         "cluster": cluster_id,
162         "model": "RF",
163         "train_rmse": rmse(y_train_c, rf_train_pred),
164         "train_mae": mae(y_train_c, rf_train_pred),
165         "train_r2": r2(y_train_c, rf_train_pred),
166         "test_rmse": rmse(y_test_c, rf_test_pred) if len(X_test_c) else np.nan,
167         "test_mae": mae(y_test_c, rf_test_pred) if len(X_test_c) else np.nan,
168         "test_r2": r2(y_test_c, rf_test_pred) if len(X_test_c) else np.nan,
169     }
170     cluster_metrics.append(rf_metrics)
171
172     # LightGBM
173     lgbm_base = LGBMRegressor(objective="regression", random_state=42, n_jobs=-1)
174     lgbm_rand = RandomizedSearchCV(
175         lgbm_base,
176         param_distributions=lgbm_param_dist,
177         n_iter=25,
178         cv=3,
179         scoring="neg_mean_squared_error",
180         n_jobs=-1,
181         random_state=42,
182         return_train_score=False,
183     )
184     lgbm_rand.fit(X_train_c, y_train_c)
185     lgbm_train_pred = lgbm_rand.predict(X_train_c)
186     lgbm_test_pred = lgbm_rand.predict(X_test_c) if len(X_test_c) else np.array([])
187
188     lgbm_metrics = {
189         "cluster": cluster_id,
190         "model": "LGBM",
191         "train_rmse": rmse(y_train_c, lgbm_train_pred),
192         "train_mae": mae(y_train_c, lgbm_train_pred),
193         "train_r2": r2(y_train_c, lgbm_train_pred),
194         "test_rmse": rmse(y_test_c, lgbm_test_pred) if len(X_test_c) else np.nan,
195         "test_mae": mae(y_test_c, lgbm_test_pred) if len(X_test_c) else np.nan,
196         "test_r2": r2(y_test_c, lgbm_test_pred) if len(X_test_c) else np.nan,
197     }
198     cluster_metrics.append(lgbm_metrics)
199

```

```

200     # Store test predictions for overall metrics
201     rf_preds_test.loc[test_mask] = rf_test_pred
202     lgbm_preds_test.loc[test_mask] = lgbm_test_pred
203
204     cluster_metrics_df = pd.DataFrame(cluster_metrics)
205     cluster_metrics_df
206
207
208     # ---- Overall metrics on test (log-price space) ----
209     rf_overall = {
210         "model": "RF",
211         "rmse": rmse(y_test, rf_preds_test),
212         "mae": mae(y_test, rf_preds_test),
213         "r2": r2(y_test, rf_preds_test),
214     }
215
216     lgbm_overall = {
217         "model": "LGBM",
218         "rmse": rmse(y_test, lgbm_preds_test),
219         "mae": mae(y_test, lgbm_preds_test),
220         "r2": r2(y_test, lgbm_preds_test),
221     }
222
223     overall_metrics_df = pd.DataFrame([rf_overall, lgbm_overall])
224     print("Cluster-level metrics (log-price):")
225     display(cluster_metrics_df)
226     print("\nOverall test metrics (log-price):")
227     display(overall_metrics_df)

```

```

1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # ## Cluster-as-Feature for RF & LightGBM
5  #
6  # Following the first 14 cells of the \texttt{Honolulu\_clustering} workflow, we
   perform market segmentation using the original feature set \texttt{latitude}, \
   \texttt{longitude}, \texttt{drive\_dist\_hnl\_km}, \texttt{drive\_dist\_wk\_km}, \
   \texttt{review\_scores\_rating}, \texttt{review\_scores\_cleanliness}, \texttt{
   accommodates}, \texttt{bedrooms}, and \texttt{bathrooms}. These features are
   standardized using \texttt{StandardScaler}, and agglomerative clustering with Ward
   linkage is applied. Based on the silhouette score, the optimal number of clusters
   is selected as $k=3$.
7  #
8  # To avoid information leakage, clustering is fitted exclusively on the training set.
   The same scaler and the resulting cluster centroids are then used to assign cluster
   labels to the test set via nearest-centroid matching.
9  #
10 # For downstream predictive modeling, we preserve the same feature preprocessing
   pipeline as in the \texttt{Baseline\_Model}. Specifically, we drop \texttt{id}, \
   \texttt{neighbourhood\_group}, \texttt{latitude}, and \texttt{longitude}, and apply
   one-hot encoding to \texttt{room\_type}, \texttt{neighbourhood}, and \texttt{host\
   _response\_time}.
11 #
12 # Using this setup, we compare the performance of Random Forest and LightGBM models
   under two feature configurations: the baseline feature set and the augmented
   feature set that includes cluster membership indicators (cluster-as-features).
13
14 # In[40]:
15
16
17 import pandas as pd
18 import numpy as np
19 from pathlib import Path
20 import matplotlib.pyplot as plt
21
22 from sklearn.model_selection import train_test_split, RandomizedSearchCV
23 from sklearn.metrics import silhouette_score
24 from catboost import CatBoostRegressor
25 from xgboost import XGBRegressor
26 import shap
27 from sklearn.preprocessing import StandardScaler

```

```

28 from sklearn.cluster import AgglomerativeClustering
29 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
30 from sklearn.ensemble import RandomForestRegressor
31 from lightgbm import LGBMRegressor
32
33 # Repro
34 RANDOM_STATE = 42
35 DATA_PATH = Path("/Users/jiangzhanyuan/Desktop/second year/IEOR242A/Final Project/
    Listing_Honolulu.csv")
36
37 pd.set_option("display.max_columns", None)
38
39
40 # In[41]:
41
42
43 # Load data
44
45 df = pd.read_csv(DATA_PATH)
46 print(f"Loaded dataset: {df.shape}")
47
48
49 # In[42]:
50
51
52 # Baseline preprocessing config (must mirror Baseline_Model.ipynb)
53
54 target_col = "price"
55 base_drop_cols = ["id", "neighbourhood_group", "latitude", "longitude"]
56 base_cat_cols = ["room_type", "neighbourhood", "host_response_time"]
57
58 # Clustering config (matching Honolulu_clustering first 14 cells)
59 cluster_cols = [
60     "latitude",
61     "longitude",
62     "drive_dist_hnl_km",
63     "drive_dist_wk_km",
64     "review_scores_rating",
65     "review_scores_cleanliness",
66     "accommodates",
67     "bedrooms",
68     "bathrooms",
69 ]
70
71 best_k = 4
72
73
74 # In[43]:
75
76
77 def prepare_baseline_features(df_in: pd.DataFrame):
78     df_model = df_in.copy()
79     cols_to_drop = [c for c in base_drop_cols if c in df_model.columns]
80     X = df_model.drop(columns=cols_to_drop + [target_col])
81     y = np.log1p(df_model[target_col].astype(float))
82     X_enc = pd.get_dummies(X, columns=base_cat_cols, drop_first=True, dtype=float).
        astype(float)
83     return X_enc, y
84
85
86 def train_test_split_baseline(df_in: pd.DataFrame):
87     X_enc, y_log = prepare_baseline_features(df_in)
88     X_train, X_test, y_train, y_test = train_test_split(
89         X_enc, y_log, test_size=0.2, shuffle=True, random_state=RANDOM_STATE
90     )
91     return X_train, X_test, y_train, y_test, X_enc
92
93
94 def eval_regression(y_true, y_pred):
95     # squared=False not available in older sklearn; compute RMSE manually
96     mse_val = mean_squared_error(y_true, y_pred)

```

```

97     rmse_val = np.sqrt(mse_val)
98     return {
99         "rmse": rmse_val,
100         "mae": mean_absolute_error(y_true, y_pred),
101         "r2": r2_score(y_true, y_pred),
102     }
103
104
105 # In[37]:
106
107
108 def fit_cluster_labels(train_raw: pd.DataFrame, test_raw: pd.DataFrame):
109     train_clust = train_raw[cluster_cols].copy()
110     test_clust = test_raw[cluster_cols].copy()
111
112     scaler = StandardScaler()
113     X_train_scaled = scaler.fit_transform(train_clust)
114     X_test_scaled = scaler.transform(test_clust)
115
116     agg = AgglomerativeClustering(n_clusters=best_k, linkage="ward")
117     train_labels = agg.fit_predict(X_train_scaled)
118
119     # Derive test labels via nearest centroid in scaled space (no refit on test)
120     centroids = np.vstack([
121         X_train_scaled[train_labels == k].mean(axis=0)
122         for k in range(best_k)
123     ])
124     dists = ((X_test_scaled[:, None, :] - centroids[None, :, :]) ** 2).sum(axis=2) **
125         0.5
126     test_labels = dists.argmin(axis=1)
127
128     return pd.Series(train_labels, index=train_raw.index, name="cluster"), pd.Series(
129         test_labels, index=test_raw.index, name="cluster")
130
131
132 def build_cluster_feature_matrices(X_train_base, X_test_base, df_full: pd.DataFrame):
133     # Align raw rows with split indices
134     train_raw = df_full.loc[X_train_base.index]
135     test_raw = df_full.loc[X_test_base.index]
136
137     train_labels, test_labels = fit_cluster_labels(train_raw, test_raw)
138
139     cluster_cols_full = [f"cluster_{k}" for k in range(best_k)]
140     train_dummies = pd.get_dummies(train_labels, prefix="cluster").reindex(columns=
141         cluster_cols_full, fill_value=0)
142     test_dummies = pd.get_dummies(test_labels, prefix="cluster").reindex(columns=
143         cluster_cols_full, fill_value=0)
144
145     X_train_cluster = pd.concat([X_train_base.copy(), train_dummies], axis=1)
146     X_test_cluster = pd.concat([X_test_base.copy(), test_dummies], axis=1)
147     return X_train_cluster, X_test_cluster, train_labels, test_labels
148
149
150 # In[44]:
151
152
153 rf_param_grid = {
154     "n_estimators": [50, 100, 200],
155     "max_depth": [10, 20, 40],
156     "min_samples_split": [2, 5, 10],
157     "min_samples_leaf": [2, 4, 8],
158 }
159
160
161 lgbm_param_grid = {
162     "n_estimators": [200, 500, 800],
163     "num_leaves": [31, 63, 127],
164     "max_depth": [-1, 10, 20],
165     "learning_rate": [0.05, 0.1],
166     "subsample": [0.8, 1.0],
167     "colsample_bytree": [0.8, 1.0],
168     "min_child_samples": [10, 20, 40],

```



```

164 }
165
166 cat_param_grid = {
167     "depth": [4, 6, 8, 10],
168     "learning_rate": [0.03, 0.1],
169     "iterations": [200, 500, 800],
170     "l2_leaf_reg": [1, 3, 5],
171     "subsample": [0.7, 0.9, 1.0],
172 }
173
174 xgb_param_grid = {
175     "n_estimators": [200, 400, 800],
176     "max_depth": [4, 6, 10],
177     "learning_rate": [0.05, 0.1, 0.2],
178     "subsample": [0.7, 0.9, 1.0],
179     "colsample_bytree": [0.7, 0.9, 1.0],
180     "min_child_weight": [1, 5],
181 }
182
183
184 def run_rf_grid(X_train, y_train, X_test, y_test, n_iter=20):
185     rf_base = RandomForestRegressor(
186         bootstrap=True,
187         random_state=RANDOM_STATE,
188         n_jobs=-1,
189     )
190     search = RandomizedSearchCV(
191         rf_base,
192         param_distributions=rf_param_grid,
193         n_iter=n_iter,
194         cv=5,
195         scoring="neg_mean_squared_error",
196         n_jobs=-1,
197         random_state=RANDOM_STATE,
198         return_train_score=False,
199     )
200     search.fit(X_train, y_train)
201     best_model = search.best_estimator_
202     train_pred = best_model.predict(X_train)
203     test_pred = best_model.predict(X_test)
204     return {
205         "best_params": search.best_params_,
206         "model": best_model,
207         "train": eval_regression(y_train, train_pred),
208         "test": eval_regression(y_test, test_pred),
209     }
210
211
212 def run_lgbm_grid(X_train, y_train, X_test, y_test, n_iter=20):
213     lgbm = LGBMRegressor(objective="regression", random_state=RANDOM_STATE, n_jobs=-1)
214     search = RandomizedSearchCV(
215         lgbm,
216         param_distributions=lgbm_param_grid,
217         n_iter=n_iter,
218         cv=5,
219         scoring="neg_mean_squared_error",
220         n_jobs=-1,
221         random_state=RANDOM_STATE,
222         return_train_score=False,
223     )
224     search.fit(X_train, y_train)
225     best_model = search.best_estimator_
226     train_pred = best_model.predict(X_train)
227     test_pred = best_model.predict(X_test)
228     return {
229         "best_params": search.best_params_,
230         "model": best_model,
231         "train": eval_regression(y_train, train_pred),
232         "test": eval_regression(y_test, test_pred),
233     }
234

```

```

235
236 def run_cat_grid(X_train, y_train, X_test, y_test, n_iter=20):
237     cat_base = CatBoostRegressor(
238         loss_function="RMSE",
239         random_seed=RANDOM_STATE,
240         verbose=0,
241     )
242     search = RandomizedSearchCV(
243         cat_base,
244         param_distributions=cat_param_grid,
245         n_iter=n_iter,
246         cv=5,
247         scoring="neg_mean_squared_error",
248         n_jobs=-1,
249         random_state=RANDOM_STATE,
250         return_train_score=False,
251     )
252     search.fit(X_train, y_train)
253     best_model = search.best_estimator_
254     train_pred = best_model.predict(X_train)
255     test_pred = best_model.predict(X_test)
256     return {
257         "best_params": search.best_params_,
258         "model": best_model,
259         "train": eval_regression(y_train, train_pred),
260         "test": eval_regression(y_test, test_pred),
261     }
262
263
264 def run_xgb_grid(X_train, y_train, X_test, y_test, n_iter=20):
265     xgb_base = XGBRegressor(
266         objective="reg:squarederror",
267         random_state=RANDOM_STATE,
268         n_jobs=-1,
269         tree_method="hist",
270     )
271     search = RandomizedSearchCV(
272         xgb_base,
273         param_distributions=xgb_param_grid,
274         n_iter=n_iter,
275         cv=5,
276         scoring="neg_mean_squared_error",
277         n_jobs=-1,
278         random_state=RANDOM_STATE,
279         return_train_score=False,
280     )
281     search.fit(X_train, y_train)
282     best_model = search.best_estimator_
283     train_pred = best_model.predict(X_train)
284     test_pred = best_model.predict(X_test)
285     return {
286         "best_params": search.best_params_,
287         "model": best_model,
288         "train": eval_regression(y_train, train_pred),
289         "test": eval_regression(y_test, test_pred),
290     }
291
292
293 # In[ ]:
294
295
296 # Prepare baseline split
297 X_train_base, X_test_base, y_train_log, y_test_log, X_all_base =
298     train_test_split_baseline(df)
299 print(f"Baseline feature count: {X_train_base.shape[1]}")
300
301 # Silhouette check on train-only (same clustering feature set)
302 train_clust_data = df.loc[X_train_base.index, cluster_cols]
303 sil_scaler = StandardScaler()
304 X_train_scaled_for_k = sil_scaler.fit_transform(train_clust_data)
305 sil_scores = {}

```

```

305 for k in range(2, 10):
306     agg_k = AgglomerativeClustering(n_clusters=k, linkage="ward")
307     labels_k = agg_k.fit_predict(X_train_scaled_for_k)
308     sil_scores[k] = silhouette_score(X_train_scaled_for_k, labels_k)
309 print("\nSilhouette scores on train (k=2..9):")
310 for k, s in sil_scores.items():
311     print(f"k={k}: {s:.3f}")
312 print(f"Using predefined best_k={best_k}")
313
314 # Build cluster features (train-only fit, assign test via centroids)
315 X_train_cluster, X_test_cluster, train_clusters, test_clusters =
316     build_cluster_feature_matrices(
317         X_train_base, X_test_base, df
318     )
319 print("Cluster label distribution (train/test):")
320 print(train_clusters.value_counts().sort_index())
321 print(test_clusters.value_counts().sort_index())
322 print(f"Cluster feature count: {X_train_cluster.shape[1]}")
323
324 # Models on baseline features
325 rf_base_res = run_rf_grid(X_train_base, y_train_log, X_test_base, y_test_log)
326 lgb_base_res = run_lgbm_grid(X_train_base, y_train_log, X_test_base, y_test_log)
327 cat_base_res = run_cat_grid(X_train_base, y_train_log, X_test_base, y_test_log)
328 xgb_base_res = run_xgb_grid(X_train_base, y_train_log, X_test_base, y_test_log)
329
330 # Models on cluster-as-feature features
331 rf_cluster_res = run_rf_grid(X_train_cluster, y_train_log, X_test_cluster, y_test_log)
332 lgb_cluster_res = run_lgbm_grid(X_train_cluster, y_train_log, X_test_cluster,
333     y_test_log)
334 cat_cluster_res = run_cat_grid(X_train_cluster, y_train_log, X_test_cluster,
335     y_test_log)
336 xgb_cluster_res = run_xgb_grid(X_train_cluster, y_train_log, X_test_cluster,
337     y_test_log)
338
339 results = []
340 for name, res in [
341     ("RF", rf_base_res),
342     ("RF + cluster", rf_cluster_res),
343     ("LightGBM", lgb_base_res),
344     ("LightGBM + cluster", lgb_cluster_res),
345     ("CatBoost", cat_base_res),
346     ("CatBoost + cluster", cat_cluster_res),
347     ("XGBoost", xgb_base_res),
348     ("XGBoost + cluster", xgb_cluster_res),
349 ]:
350     results.append({
351         "model": name,
352         "set": "train",
353         **res["train"],
354     })
355     results.append({
356         "model": name,
357         "set": "test",
358         **res["test"],
359     })
360
361 metrics_df = pd.DataFrame(results)
362 print("\nMetrics (log-price space):")
363 print(metrics_df)
364
365 # Focused cluster models: both train & test
366 cluster_compare = metrics_df[metrics_df["model"].str.contains("cluster")]
367 print("\nCluster-as-feature metrics (train & test):")
368 print(cluster_compare)
369
370 print("\nBest params:")
371 print("RF baseline:", rf_base_res["best_params"])
372 print("RF + cluster:", rf_cluster_res["best_params"])
373 print("LightGBM baseline:", lgb_base_res["best_params"])
374 print("LightGBM + cluster:", lgb_cluster_res["best_params"])
375 print("CatBoost baseline:", cat_base_res["best_params"])

```

```

372 print("CatBoost + cluster:", cat_cluster_res["best_params"])
373 print("XGBoost baseline:", xgb_base_res["best_params"])
374 print("XGBoost + cluster:", xgb_cluster_res["best_params"])
375
376 # Keep best models for SHAP (RF/LGBM only)
377 rf_cluster_best_model = rf_cluster_res["model"]
378 lgb_cluster_best_model = lgb_cluster_res["model"]
379
380
381 # In[47]:
382
383
384 # Focused comparison: cluster-as-feature models (train & test)
385 cluster_compare = metrics_df[metrics_df["model"].str.contains("cluster")].copy()
386 print("Cluster-as-feature metrics (train & test):")
387 print(cluster_compare)
388
389
390 # In[53]:
391
392
393 # SHAP explanations for cluster-as-feature models (using train split)
394 # To reduce runtime, subsample if needed (set sample_n=None to use all)
395 import matplotlib.pyplot as plt
396 sample_n = 1000
397 X_shap = X_train_cluster.reset_index(drop=True)
398 y_shap = y_train_log.reset_index(drop=True)
399 if sample_n is not None and len(X_shap) > sample_n:
400     X_shap = X_shap.sample(sample_n, random_state=RANDOM_STATE)
401     y_shap = y_shap.loc[X_shap.index]
402
403 # RF SHAP
404 rf_explainer = shap.TreeExplainer(rf_cluster_best_model)
405 rf_shap_values = rf_explainer.shap_values(X_shap)
406 print("RF SHAP computed: shape", np.array(rf_shap_values).shape)
407 plt.figure()
408 shap.summary_plot(rf_shap_values, X_shap, show=False)
409 plt.title("RF + cluster SHAP (beeswarm)")
410 plt.show()
411
412 plt.figure()
413 shap.summary_plot(rf_shap_values, X_shap, plot_type="bar", show=False)
414 plt.title("RF + cluster SHAP (abs bar)")
415 plt.show()
416
417 # LightGBM SHAP
418 lgb_explainer = shap.TreeExplainer(lgb_cluster_best_model)
419 lgb_shap_values = lgb_explainer.shap_values(X_shap)
420 print("LightGBM SHAP computed: shape", np.array(lgb_shap_values).shape)
421 plt.figure()
422 shap.summary_plot(lgb_shap_values, X_shap, show=False)
423 plt.title("LightGBM + cluster SHAP (beeswarm)")
424 plt.show()
425
426 plt.figure()
427 shap.summary_plot(lgb_shap_values, X_shap, plot_type="bar", show=False)
428 plt.title("LightGBM + cluster SHAP (abs bar)")
429 plt.show()
430
431 # CatBoost SHAP
432 cat_explainer = shap.TreeExplainer(cat_cluster_res["model"])
433 cat_shap_values = cat_explainer.shap_values(X_shap)
434 print("CatBoost SHAP computed: shape", np.array(cat_shap_values).shape)
435 plt.figure()
436 shap.summary_plot(cat_shap_values, X_shap, show=False)
437 plt.title("CatBoost + cluster SHAP (beeswarm)")
438 plt.show()
439
440 plt.figure()
441 shap.summary_plot(cat_shap_values, X_shap, plot_type="bar", show=False)
442 plt.title("CatBoost + cluster SHAP (abs bar)")

```

```

443 plt.show()
444
445 # XGBoost SHAP (use permutation explainer on a callable predict)
446 X_shap_np = X_shap.astype(float)
447 masker = shap.maskers.Independent(X_shap_np)
448 xgb_predict = lambda data: xgb_cluster_res["model"].predict(data)
449 xgb_explainer = shap.Explainer(xgb_predict, masker, algorithm="permutation")
450 xgb_shap_values = xgb_explainer(X_shap_np)
451 print("XGBoost SHAP computed: shape", np.array(xgb_shap_values.values).shape)
452 plt.figure()
453 shap.summary_plot(xgb_shap_values, X_shap_np, show=False)
454 plt.title("XGBoost + cluster SHAP (beeswarm)")
455 plt.show()
456
457 plt.figure()
458 shap.summary_plot(xgb_shap_values, X_shap_np, plot_type="bar", show=False)
459 plt.title("XGBoost + cluster SHAP (abs bar)")
460 plt.show()

```