

2SRI Hotel Pricing Implementation

Notebook 1: Data Exploration and Quality Assessment

Nandan

September 25, 2025

1 Executive Summary

This report documents the data exploration phase of the 2SRI hotel pricing implementation. The analysis reveals a well-structured dataset with 1 focal hotel containing 14 room types and 5 competitor hotels. The data shows strong temporal overlap (364 days) and meaningful competitive relationships, providing a solid foundation for the subsequent 2SRI modeling approach.

2 Data Loading and Basic Validation

2.1 Dataset Structure

The implementation uses two primary datasets:

- **Focal Hotel Data:** 5,110 observations across 1 hotel with 14 distinct room types
- **Competitor Data:** 1,820 observations across 5 competitor hotels

2.2 Temporal Coverage

Both datasets cover a full year period with proper datetime formatting. The focal hotel provides daily pricing across multiple room categories, while competitors provide base rate pricing only.

2.3 Data Types Validation

All critical fields show appropriate data types:

- Hotel IDs: Object type (string identifiers)
- Stay dates: Properly converted to datetime64
- Prices: Numeric (int64/float64)
- Room types: Categorical (int64) for focal hotel

3 Data Quality Assessment

3.1 Missing Data Analysis

- **Focal Hotel:** 0 missing values (complete dataset)
- **Competitor Hotels:** 30 missing price values requiring treatment

Missing data patterns by competitor hotel:

- Aqua Pacific Monarch: 16 missing (4.4 percent missing rate)
- Castle Kamaole Sands: 0 missing (complete)
- Courtyard Marriott: 11 missing (3.0 percent missing rate)
- Koheia Kai Resort: 2 missing (0.5 percent missing rate)
- Ohana Waikiki Malia: 1 missing (0.3 percent missing rate)

3.2 Price Validation

Invalid Price Detection: No prices less than or equal to 0 detected in either dataset.

Price Range Analysis:

- Focal hotel range: \$219.00 - \$999.00
- Competitor range: \$207.17 - \$904.12

3.3 Room Type Analysis (Focal Hotel)

The focal hotel operates 14 distinct room types with consistent pricing patterns:

- Room type price range: \$286.15 - \$541.66 (mean prices)
- Standard deviation: \$33.33 - \$139.29 across room types
- All room types have complete 365-day coverage

3.4 Availability Metrics

High availability rates across all competitor hotels:

- Aqua Pacific Monarch: 95.3 percent
- Castle Kamaole Sands: 100.0 percent
- Courtyard Marriott: 97.0 percent
- Koheia Kai Resort: 99.5 percent
- Ohana Waikiki Malia: 99.7 percent

4 Temporal Alignment and Overlap Analysis

4.1 Overlap Assessment

Temporal Overlap: 364 days of common coverage between focal and competitor datasets.

Sample Sizes in Overlap Period:

- Focal observations: 5,096
- Competitor observations: 1,820

This extensive overlap provides sufficient data for robust 2SRI estimation.

5 Price Normalization

5.1 Normalization Statistics

Robust normalization using median and median absolute deviation (MAD):

- **Focal Hotel:** Median = \$319.00, MAD = \$40.00
- **Competitors:** Median = \$298.90, MAD = \$49.03

Normalized Price Ranges:

- Focal normalized range: \$164.07 - \$1,438.85
- Competitor range: \$179.00 - \$904.12

6 Competitive Relationship Analysis

6.1 Price Correlation Analysis

Key Findings:

- Common dates for analysis: 364
- Price level correlation: 0.288
- Price change correlation: 0.133
- Normalized price correlation: 0.288
- Strongest correlation at lag 1: 0.297

The correlation of 0.288 indicates moderate competitive interdependence, sufficient for 2SRI identification but not so strong as to suggest perfect competition.

6.2 Lead-Lag Analysis

The lead-lag analysis reveals temporal pricing relationships:

- Peak correlation at 1-day lag: 0.297
- Suggests focal hotel pricing may lead competitor adjustments
- Provides validation for using lagged competitor prices as instruments

7 Instrument Validation for 2SRI

7.1 Available Instruments

The analysis confirms availability of key exogenous instruments:

- **Day of week effects:** Available (temporal variation)
- **Monthly/seasonal effects:** Available (calendar patterns)
- **Availability-Focal correlation:** 0.104
- **Availability-Competitor correlation:** -0.440

7.2 Room Type Segmentation Analysis

Room type correlations with competitor pricing range from 0.274 to 0.529, indicating:

- Consistent competitive relationships across room categories
- Stronger correlations for premium room types (0.5+ range)
- Sufficient variation for identification across segments

8 2SRI Implementation Readiness Assessment

8.1 Data Quality Checklist

- **Sufficient overlap:** PASS (364 days)
- **Clean price data:** PASS (no invalid prices)
- **Missing data manageable:** PASS (1.6 percent missing)
- **Price variation adequate:** PASS (sufficient variance)
- **Meaningful correlation:** PASS (0.288 correlation)
- **Multiple competitors:** PASS (5 hotels)
- **Instruments available:** PASS (temporal + availability)

Overall Readiness Score: 100.0 percent

Status: READY FOR 2SRI IMPLEMENTATION

9 Conclusions and Next Steps

9.1 Key Findings

1. Dataset provides excellent foundation for 2SRI implementation with complete temporal coverage and minimal missing data
2. Moderate competitive correlation (0.288) offers ideal balance for identification - strong enough for meaningful relationships but not perfect substitution
3. Room type analysis reveals potential for segmented modeling approaches
4. Availability data provides additional exogenous variation for instrument strength

9.2 Methodological Implications

1. Base rate extraction from focal hotel (daily minimum across room types) will be required for proper competitive comparison
2. Missing data treatment using forward-fill and median imputation is appropriate given low missing rates
3. Temporal instruments (day-of-week, seasonality) combined with availability measures should provide sufficient identification strength

9.3 Recommended Next Steps

1. Proceed to Notebook 2: Data preprocessing with base rate extraction
2. Implement robust price validation using data-driven bounds rather than arbitrary thresholds
3. Maintain focus on temporal alignment for Stage 1 and Stage 2 2SRI implementation

This comprehensive data exploration confirms that the dataset meets all requirements for a robust 2SRI competitive pricing analysis.