

# Data Preprocessing Pipelines

Nandan

September 25, 2025

Reference: 02.Stage1\_staging.ipynb

## 1 Systematic Missing Value Treatment

### 1.1 Objective

Address missing competitor pricing data using forward-fill imputation followed by group median fallback, following Section 2.3 methodology for handling dirty hotel pricing data.

### 1.2 Implementation

Applied two-stage missing value treatment:

$$\text{Stage 1: Forward Fill : } P_{i,t} = P_{i,t-1} \text{ if } P_{i,t} = \text{NA} \quad (1)$$

$$\text{Stage 2: Group Median : } P_{i,t} = \text{median}(P_i) \text{ if still NA} \quad (2)$$

### 1.3 Results

- Missing competitor prices: 30 observations (1.6% of total)
- Non-random pattern: Aqua Pacific (16 missing), Courtyard Marriott (11 missing)
- Treatment successful: 0 remaining missing values

## 2 Data-Driven Price Validation and Conservative Outlier Treatment

### 2.1 Methodology

Implemented data-driven price bounds using 99th percentile thresholds rather than arbitrary limits:

$$\text{Upper Bound} = 1.5 \times \max(P99_{focal}, P99_{competitors}) \quad (3)$$

$$\text{Calculated Threshold} = \$1498 \text{ (vs arbitrary \$5000)} \quad (4)$$

Applied conservative IQR-based outlier detection:

$$\text{Outlier Threshold} = Q_{0.75} \pm 3.0 \times (Q_{0.75} - Q_{0.25}) \quad (5)$$

### 2.2 Critical Signal Preservation

$$\text{Outliers Removed : } 169 \text{ focal + 100 competitor observations} \quad (6)$$

$$\text{Competitive Correlation : } 0.440 \text{ (strengthened from initial)} \quad (7)$$

Conservative 3.0 multiplier prevented signal destruction that would occur with aggressive 1.5x multiplier.

## 3 Base Rate Extraction and Temporal Alignment

### 3.1 Critical Methodological Change

Following Professor Dan's feedback, implemented base rate focus since focal hotel has room-type pricing while competitors provide only base rates:

$$\text{Base Rate}_{focal,t} = \min_{r \in \text{room types}} P_{focal,r,t} \quad (8)$$

### 3.2 Base Rate Characteristics

Base Rate Observations : 365 daily observations (9)

Base Rate Range : \$209 – \$379 (10)

Temporal Overlap : 364 days (Sep 2025 - Sep 2026) (11)

### 3.3 Alignment Results

- Overlap shapes: Focal (364 observations), Competitors (1720 observations)
- Daily density: Focal 1.0, Competitors 4.7 observations per day
- Perfect base-rate-to-base-rate competitive comparison achieved

## 4 Cross-Hotel Normalization Implementation

### 4.1 Location-Scale Robust Normalization

Following Section 3.1, applied robust normalization using median and median absolute deviation:

Focal Base Rate Parameters : Median = \$279.00, MAD = \$20.00 (12)

Competitor Parameters : Median = \$298.00, MAD = \$36.57 (13)

$$\text{Normalization Transform} : P_{i,t}^{(norm)} = \frac{P_{i,t} - M_i}{\text{MAD}_i} \times \text{MAD}_{ref} + M_{ref} \quad (14)$$

### 4.2 Normalization Impact

Base rate normalization creates proper scale alignment for competitive analysis while preserving relative price variation within each hotel's distribution.

## 5 Feature Engineering for Instrumental Variables

### 5.1 Temporal Instruments

Created cyclical temporal features for Stage 1 modeling:

Day Cycle :  $\sin(2\pi \cdot \text{day\_of\_week}/7), \cos(2\pi \cdot \text{day\_of\_week}/7)$  (15)

Monthly Cycle :  $\sin(2\pi \cdot \text{month}/12), \cos(2\pi \cdot \text{month}/12)$  (16)

Weekend Indicator :  $\mathbf{1}[\text{day\_of\_week} \geq 5]$  (17)

Week of Year : ISO calendar week number (18)

## 5.2 Price Dynamics Features

Constructed lagged price features for dynamic analysis:

$$\text{Base Rate Lags : } P_{focal,t-k} \text{ for } k \in \{1, 2, 3\} \quad (19)$$

$$\text{Price Changes : } \Delta P_{focal,t} = P_{focal,t} - P_{focal,t-1} \quad (20)$$

$$\text{Competitor Lags : } P_{comp_i,t-k} \text{ for each competitor } i \quad (21)$$

# 6 Data Structure Preparation for 2SRI

## 6.1 Base Rate Focus

Created simplified structure focusing on base rate competition:

- Focal base rate daily series: 364 observations
- Competitor price matrix: 364 dates  $\times$  5 hotels
- Eliminated room-type complexity for clean competitive modeling

## 6.2 Stage 1 Preparation

Structured data for flexible Stage 1 modeling:

$$\text{Focal Daily Structure : Single base rate per day} \quad (22)$$

$$\text{Competitor Matrix Structure : Hotel-date pivot with normalized prices} \quad (23)$$

$$\text{Temporal Features : Comprehensive instrument set} \quad (24)$$

# 7 Quality Validation Post-Processing

## 7.1 Competitive Relationship Validation

Final validation of base rate competitive relationships:

$$\text{Base Rate Correlation} = 0.440 \quad (25)$$

$$\text{Sample Sizes : } 364 \text{ focal, } 1720 \text{ competitor observations} \quad (26)$$

$$\text{Warning Threshold : No correlation weakening detected} \quad (27)$$

## 7.2 2SRI Readiness Assessment

Base rate approach successfully creates:

- Meaningful competitive correlation ( $0.440 > 0.2$  threshold)
- Clean temporal alignment without room-type complexity
- Proper base-rate-to-base-rate competitive comparison
- Foundation for Stage 1 flexible modeling implementation

## 8 Export Pipeline for Stage 1 Implementation

### 8.1 Output Structure

Created analysis-ready datasets:

- `focal_hotel_clean.csv`: Base rate time series with features
- `competitors_clean.csv`: Normalized competitor pricing data
- `focal_daily_aggregated.csv`: Daily focal base rates
- `competitor_price_matrix.csv`: Pivoted competitor matrix