Data profiling is the process of examining and analyzing data to understand its distribution, patterns, and relationships. Here's a data profiling report for the provided dataset:

**Overview**

- The dataset contains 100 rows and 13 columns.
- The columns are:
    1. brokered_by
    2. status
    3. price
    4. bed
    5. bath
    6. acre_lot
    7. street
    8. city
    9. state
    10. zip_code
    11. house_size
    12. prev_sold_date

**Data Types**

- **brokered_by**: integer (missing values: 5)
- **status**: categorical (missing values: 0)
    - Unique values: 'sold', 'for_sale', 'f', 's', '-', 'Unknown'
- **price**: numeric (missing values: 2)
- **bed**: integer (missing values: 3)
- **bath**: integer (missing values: 2)
- **acre_lot**: numeric (missing values: 2)
- **street**: string (missing values: 4)
- **city**: string (missing values: 4)
- **state**: string (missing values: 4)
- **zip_code**: integer (missing values: 6)
- **house_size**: numeric (missing values: 4)
- **prev_sold_date**: datetime (missing values: 4)

**Distribution**

- **price**:
    - Min: 11,400
    - Max: 2,359,000
    - Mean: 434,911
    - Median: 349,000
- **bed**:
    - Min: 1
    - Max: 12
    - Mean: 3.5
    - Median: 3
- **bath**:
    - Min: 1
    - Max: 12
    - Mean: 2.7
    - Median: 2
- **acre_lot**:
    - Min: -1
    - Max: 27.59
    - Mean: 0.24
    - Median: 0.15
- **house_size**:
    - Min: 85
    - Max: 6,828

– Mean: 2,044
– Median: 1,874

**Relationships**

- **status** and **price**:
  – 'sold' prices range from 11,400 to 2,359,000
  – 'for_sale' prices range from 105,000 to 2,199,000
- **bed** and **bath**:
  – Most common combination: 3 beds, 2 baths
- **acre_lot** and **house_size**:
  – No clear correlation

**Missing Values**

- **brokered_by**: 5 missing values
- **price**: 2 missing values
- **bed**: 3 missing values
- **bath**: 2 missing values
- **acre_lot**: 2 missing values
- **street**: 4 missing values
- **city**: 4 missing values
- **state**: 4 missing values
- **zip_code**: 6 missing values
- **house_size**: 4 missing values
- **prev_sold_date**: 4 missing values

**Outliers**

- **price**: 2,359,000 (likely an outlier)
- **acre_lot**: 27.59 (likely an outlier)
- **house_size**: 6,828 (likely an outlier)

**Data Quality Issues**

- Inconsistent formatting in **prev_sold_date** column
- Missing values in multiple columns
- Outliers in **price**, **acre_lot**, and **house_size** columns
- Inconsistent values in **status** column (e.g., 'f', 's', '-')

Overall, the dataset appears to be a collection of real estate listings with various attributes. However, there are several data quality issues that need to be addressed, including missing values, outliers, and inconsistent formatting.