I'll help you profile this dataset. Here's a comprehensive analysis:

1. Dataset Overview:

- Total rows: 89
- Total columns: 11

2. Column Analysis:

brokered_by: - Contains numeric IDs and missing values - Many "-1" and "nan" values - Mixed data quality

status: - Main categories: "sold", "for_sale", "f", "s", "Unknown" - Contains missing values - Inconsistent formatting

price: - Mix of numeric values and string formats (some with "$" prefix) - Range appears to be from ~$13,400 to $2,365,000 - Contains missing values and "-1"

bed: - Mix of numeric and text representations ("three", "four", etc.) - Range appears to be 1-11 bedrooms - Contains missing values

bath: - Similar to bed, mixed numeric and text representations - Range appears to be 1-10 bathrooms - Contains missing values

acre_lot: - Numeric values representing property size - Contains many missing values and "-1" - Range varies significantly (0.0017 to 48.0)

street: - Mixed format (some with numbers separated by commas) - Contains "Unknown" entries - Many missing values

city: - Various city names - Contains "Unknown" entries - Many missing values

state: - Mix of full state names and abbreviations - Contains "Unknown" entries - Many missing values

zip_code: - Numeric postal codes - Many missing values and "-1" - Some invalid entries

house_size: - Numeric values with varying scales - Contains very small decimal values and larger whole numbers - Many missing values and "-1"

prev_sold_date: - Mixed date formats (YYYY-MM-DD, MM/DD/YY) - Contains "Unknown" entries - Many missing values

3. Data Quality Issues:

- Inconsistent formatting across multiple columns
- Multiple representations for missing values ("-1", "nan", empty strings)
- Mixed data types within same columns
- Inconsistent date formats
- Inconsistent state name formatting
- Various ways to represent the same status (e.g., "sold" vs "s")

4. Recommendations for Cleaning:
5. Standardize missing value representation
6. Convert price to numeric format, removing "$" symbols
7. Standardize bed/bath counts to numeric values
8. Normalize state names to either all full names or all abbreviations
9. Standardize date format
10. Standardize status values
11. Validate and clean zip codes
12. Convert house_size to consistent units

This dataset would require significant cleaning before being used for analysis, particularly in handling the various inconsistencies and standardizing the formats across columns.