Outlier detection on this dataset requires a multi-pronged approach due to the mixed data types (numerical and categorical). Here's a breakdown of how we can tackle this, along with Python code using Pandas and Scikit-learn:

**1. Data Loading and Preparation:**

```python
import pandas as pd
import numpy as np
from sklearn.ensemble import IsolationForest
from sklearn.preprocessing import StandardScaler

# Load the data (assuming it's in a CSV file named 'real_estate.csv')
df = pd.read_csv('real_estate.csv')

# Convert prev_sold_date to datetime objects (if not already)
df['prev_sold_date'] = pd.to_datetime(df['prev_sold_date'])

# Handle missing values (if any).  For simplicity, this example drops rows with NaNs.  A more robust approach
df.dropna(inplace=True)
```

**2. Numerical Outlier Detection:**

- **Isolation Forest:** Effective for high-dimensional data and doesn't assume a normal distribution.

```python
# Select numerical features for outlier detection
numerical_features = ['price', 'bed', 'bath', 'acre_lot', 'house_size']

# Scale the numerical features (important for Isolation Forest)
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df[numerical_features])

# Apply Isolation Forest
iso_forest = IsolationForest(contamination='auto', random_state=42)  # 'auto' estimates contamination based on
outliers = iso_forest.fit_predict(df_scaled)

# Add outlier predictions back to the DataFrame
df['outlier'] = outliers

# Filter out the outliers (outliers are labeled -1)
outliers_df = df[df['outlier'] == -1]

print("Outliers detected by Isolation Forest:")
print(outliers_df[numerical_features])
```

- **Other Methods (briefly mentioned):**
  - **Z-score:** Assumes a normal distribution, less robust to skewed data.
  - **IQR (Interquartile Range):** Good for skewed data, but can be sensitive to small datasets.
  - **One-Class SVM:** Useful for novelty detection (identifying data points that deviate significantly from the training set).

**3. Categorical Outlier Detection (more nuanced):**

- **Frequency-based analysis:** Identify infrequent categories. This doesn't necessarily mean they are outliers, but they warrant investigation.

```python
categorical_features = ['city', 'state']

for feature in categorical_features:
    counts = df[feature].value_counts()
    threshold = 0.01 * len(df)  # Example threshold: 1% of the data
    infrequent_categories = counts[counts < threshold].index
    print(f"\nInfrequent categories in {feature}:")
    print(infrequent_categories)
```

- **Combination analysis:** Look for unusual combinations of categorical variables. For example, a "house_size" of 100 sq ft in a city known for large properties might be an outlier. This requires domain knowledge.

**4. Handling Outliers:**

- **Removal:** If confident they are errors or truly don't represent the population of interest.
- **Transformation:** Apply logarithmic or other transformations to reduce the impact of extreme values.
- **Imputation:** Replace outliers with more reasonable values (e.g., median or mean).
- **Separate analysis:** Analyze outliers separately to understand their characteristics.

**Key Considerations:**

- **Context is crucial:** Domain expertise is essential to interpret the results of outlier detection. What constitutes an outlier depends on the specific problem.
- **"Contamination" parameter:** In Isolation Forest, this parameter controls the expected proportion of outliers. Tuning this parameter is important.
- **Visualizations:** Scatter plots, box plots, and histograms can help visualize potential outliers.

This comprehensive approach provides a solid foundation for outlier detection in your real estate dataset. Remember to adapt the code and methods based on your specific needs and the insights you gain from exploring the data.