# Predicting Housing Prices with Course Methods: A Trimmed, Reproducible Study

Daniel Phelps — 10 September 2025

Email: dphelps9693@floridapoly.edu | Project Webpage | Repository

**Abstract**

This proposal narrows scope to fit course timelines while keeping the project meaningful. Using the Ames Housing dataset, I will: (1) clean and explore the data; (2) reduce dimensions with PCA; (3) form simple market segments with k-means; and (4) mine association rules that link feature combinations to Low/Medium/High price bands. Deliverables include code, figures, and a concise write-up on a GitHub Pages site.

## 1. Motivation

Accurate pricing helps buyers, sellers, and planners. Numbers such as living area, overall quality, and neighborhood are known drivers. This project uses only techniques covered in the syllabus to explain patterns in a plain, visual way—no heavy modeling—so results are reproducible and easy to grade. The focus is on clarity and alignment with course topics (EDA/visualization, PCA, clustering, frequent pattern mining, and simple anomaly checks).

## 2. Related Work / Literature Review

Hedonic pricing models express home prices as a function of characteristics (Rosen, 1974). The Ames dataset is a widely used alternative to the Boston Housing data for benchmarking feature effects and predictive workflows (De Cock, 2011). For unsupervised structure, clustering with k-means and silhouette analysis is common in real-estate segmentation studies (Rousseeuw, 1987). Association rule mining (Agrawal & Srikant, 1994; Han et al., 2000) summarizes frequent co-occurring attributes; several housing papers use rules to describe high/low value patterns. For basic text, TF-IDF and lexicon-based sentiment (VADER/AFINN) are standard tools to turn short descriptions into features. This project combines these well-established methods in a compact, transparent pipeline.

**Key references (short list)**

• Rosen, S. (1974). Hedonic prices and implicit markets. Journal of Political Economy.

• De Cock, D. (2011). Ames, Iowa: Alternative to the Boston Housing Data Set. Journal of Statistics Education.

• Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation of cluster analysis. J. Comput. Appl. Math.

• Agrawal, R., & Srikant, R. (1994). Fast algorithms for mining association rules. VLDB.

• Han, J., Pei, J., & Yin, Y. (2000). Mining frequent patterns without candidate generation (FP-Growth). SIGMOD.

• Hutto, C., & Gilbert, E. (2014). VADER: A parsimonious rule-based sentiment model. ICWSM.

• Nielsen, F. Å. (2011). AFINN word list for sentiment analysis.

## 3. Data & Preprocessing (Ames Housing)

Data split: 70/15/15 (train/val/test) with a fixed seed. Cleaning includes: removing invalid prices; imputing numeric medians and categorical modes; grouping rare categories into "Other"; capping extreme outliers for stable visuals. For rules, discretize key numerics (e.g., living area → small/medium/large) and define price bands (Low/Medium/High) using quantiles. Standardize numeric features for PCA and k-means. Keep transformations simple and well-documented.

## 4. Methods (from the course)

### 4.1 Exploratory Data Analysis (EDA) & Visualization

I will begin with a structured EDA pass to understand distributions, spot quality issues, and identify variables that are plausibly related to price.

**Variables in scope:** Numerics: SalePrice, GrLivArea (above-ground living area), TotalBsmtSF, GarageArea, GarageCars, YearBuilt, OverallQual, OverallCond, FullBath, LotArea.

Categoricals: Neighborhood, HouseStyle, BldgType, MSZoning, Exterior1st, KitchenQual, CentralAir.

**Cleaning/transform hints:** Because SalePrice is right-skewed, I will inspect both the raw scale and log10(SalePrice). I will flag extreme outliers (e.g., top/bottom 0.5–1%) and show results **with** and **without** them to make patterns robust. Missing values in numerics will be median-imputed; rare categories will be merged into "Other" when appropriate.

**Plots and tables (with one-line takeaways).**

- Distribution of SalePrice (linear and log) to justify log scale and outlier handling.

- Scatter of GrLivArea vs SalePrice (log-y), highlighting large-area outliers; report a simple correlation.
- Box plots of SalePrice by Neighborhood and by OverallQual to visualize major shifts in central tendency.
- Heatmap of correlations among numeric predictors to spot redundancy and guide PCA/feature selection.
- A compact summary table (N, mean, median, IQR) for SalePrice and 3–5 top predictors.

**Outcome.** A short narrative (2–3 sentences per figure) that says what the data show (e.g., "Price rises monotonically with OverallQual; some neighborhoods shift the median price by >$X."). These insights feed the PCA feature set and the clustering choices.

## 4.2 Dimensionality Reduction (PCA)

PCA is used to summarize structure and reduce redundancy among the standardized numeric variables; it is *not* a predictive model in this project.

Inputs. Standardized versions of key numerics from EDA: (GrLivArea, TotalBsmtSF, GarageArea, GarageCars, YearBuilt, OverallQual, FullBath, possibly LotArea). Categorical effects (e.g., Neighborhood) will be kept for later profiling rather than included in PCA.

Procedure.

- Compute PCA on the correlation matrix (i.e., after centering and scaling each feature).
- Report a scree plot and cumulative variance. Retain the first k components that together explain ~70–85% of the variance (anticipated k ≈ 2–4).
- Inspect component loadings to interpret axes (e.g., "size/amenities axis" vs. "age/condition axis").
- Produce a 2-D PCA scatter colored by the price band (Low/Medium/High, defined by tertiles of SalePrice). Optionally mark ellipses for each band.

What I will conclude.

- Which combinations of measurements form the dominant axes of variation.

- Whether price bands separate in the PCA plane (and therefore whether clustering on these standardized variables is sensible).
- Which original features load strongly on the kept components; these will guide the subset used for k-means.

Deliverables: Scree plot, loading table (top absolute loadings per component), and a PCA scatter

## 4.3 Clustering



Fig. 2. Silhouettes of a clustering with *k* = 2 of the twelve countries data of Table 1.
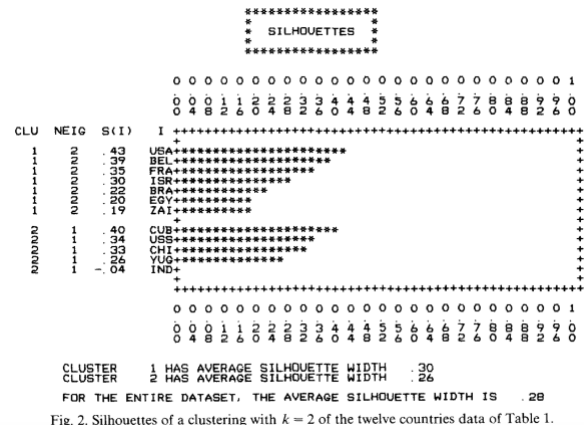
**Figure 1: De Cock, D. (2011). Ames, Iowa: Alternative to the Boston Housing Data Set. Journal of Statistics Education.**

Goal: find market segments in the standardized feature space and describe them in plain English.

Feature set.
The standardized numerics used in PCA (and possibly 1–2 binary indicators such as CentralAir if they add clear separation). Standardization ensures variables contribute equally.

Choosing K.

- Compute k-means for K = 2…8 with nstart = 25 random initializations.
- For each K, compute average silhouette width (Euclidean distance on standardized features).
- Select the K with the highest silhouette, preferring the simplest K if scores are tied or nearly tied. (We will reference Fig. 1 when discussing this choice.)

Cluster assignment and profiling.

- Fit final k-means at the chosen K; record cluster labels.

- Produce a profile table per cluster: size (N), median SalePrice (and log), typical GrLivArea, median OverallQual, top 3 neighborhoods by share, and a short sentence describing the segment (e.g., "Large, higher-quality homes, mostly in Neighborhoods A/B").

- Create a bar plot comparing medians across clusters to communicate differences quickly.

Quality and stability checks.

- Inspect within-cluster variation of price and area to ensure clusters are not dominated by a few outliers.

- If time permits, run a quick bootstrap re-fit on 80% subsamples to see if profiles are stable (not required but informative).

Outcome.
Clear, interpretable segments that can be discussed in the final presentation (e.g., "Entry-level small homes," "Mid-tier family homes," "Large premium homes"). These segments also contextualize the association rules.

### 4.4 Association Rules (Apriori / FP-Growth)

| TID | Items Bought | (Ordered) Frequent Items |
|---|---|---|
| 100 | $f, a, c, d, g, i, m, p$ | $f, c, a, m, p$ |
| 200 | $a, b, c, f, l, m, o$ | $f, c, a, b, m$ |
| 300 | $b, f, h, j, o$ | $f, b$ |
| 400 | $b, c, k, s, p$ | $c, b, p$ |
| 500 | $a, f, c, e, l, p, m, n$ | $f, c, a, m, p$ |

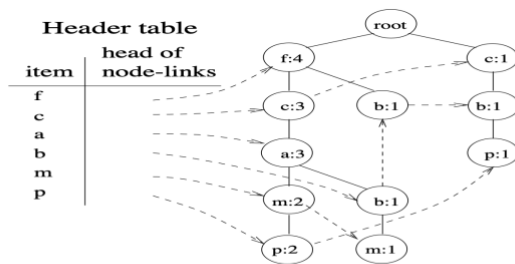Table 1: A transaction database as running example.



Figure 1: The FP-tree in Example 1.

Figure 2: FP-tree structure used by FP-Growth (adapted from Han, Pei, & Yin, 2000).

Objective: describe **combinations of attributes** that frequently occur with **Low** or **High** price bands.

**Preparation.**

- Define price_band as tertiles of SalePrice (Low/Medium/High).

- Discretize selected numerics using quantiles:

GrLivArea → area_band ∈ {small, medium, large}

TotalBsmtSF → {small, medium, large}

YearBuilt → {oldest, mid, newest}

OverallQual → {low, mid, high} (based on practical cut points)

- Keep a small set of cleaned categoricals: Neighborhood (top 6–8 levels + Other), CentralAir (Y/N), KitchenQual (grouped).

**Mining strategy.**

- Use **Apriori** (and/or **FP-Growth**) on a transaction dataset built from the discretized variables.

- **Appearance constraint**: RHS must be price_band=high or price_band=low (we are only interested in rules that imply price band).

- Start with **min support ≈ 2%** and **min confidence ≈ 60%**; adjust to get a manageable set (<50 rules).

- Rank rules by **lift** (how much more likely the RHS is given the LHS); report the **top 10** by lift for each price band.

- Remove **redundant** or **subsumed** rules (where a shorter LHS explains the same RHS with nearly equal metrics).

- Sanity-check for spurious rules due to tiny categories; if needed, raise min support or merge categories.

**Reporting.**

- A compact table with columns: LHS (conditions), RHS (price band), support, confidence, lift.

- Short, human-readable translations, e.g.,"area_band=large + OverallQual=high ⇒ price_band=high (support 9%, conf 78%, lift 1.35)."

- A 2–3 sentence discussion on *why* the strongest rules make sense (linking back to EDA/cluster profiles).

**Outcome:** Descriptive takeaways that connect **combinations** (not just single variables) to observed price levels—useful for storytelling in the final talk. We will reference **Fig. 2** (Apriori workflow) and **Fig. 3** (FP-tree) as small method insets; the results table itself uses our data.

```
1)  L_1 = {large 1-itemsets};
2)  for ( k = 2; L_{k-1} ≠ ∅; k++ ) do begin
3)      C_k = apriori-gen(L_{k-1});  // New candidates
4)      forall transactions t ∈ D do begin
5)          C_t = subset(C_k, t);  // Candidates contained in t
6)          forall candidates c ∈ C_t do
7)              c.count++;
8)      end
9)      L_k = {c ∈ C_k | c.count ≥ minsup}
10) end
11) Answer = ∪_k L_k;
```

Figure 1: Algorithm Apriori

Figure 3: Apriori workflow (adapted from Agrawal & Srikant, 1994).

## 5. Planned Figures

F1: Price distribution; F2: Price vs. living area (log scale); F3: PCA scree + 2-D PCA scatter; F4: Cluster profiles (bars/table); F5: Top association rules table.

## 6. Milestones & What "Done" Looks Like

• Proposal (this document) posted on the site.
• Checkpoint I: EDA complete; PCA plots; initial K selection; draft cluster profiles.
• Checkpoint II: association rules complete; refined clusters.
• Final: short report (9–11 pages), slides, repo with notebooks and figures.

## 7. Risks & Mitigations

• Outliers dominating visuals → cap extremes and show with/without views.
• Sparse categories → merge into "Other"; discretize key numerics.
• Time constraints → prioritize EDA, PCA, clustering, and rules.

## 8. Reproducibility & Artifacts

Public GitHub repo with notebooks (`01_eda.ipynb` → `04_association_rules.ipynb'), figures folder for exported plots, and a GitHub Pages site with links. A README describes how to run the notebooks.