## 1) Data plumbing, config, and guards

## **Data loading & caching**

- \*\*What it does:\*\* `load\_df(None)` pulls `cleaned\_large\_bgg\_dataset.parquet`
   (preferred) or `cleaned\_large\_bgg\_dataset.csv` from disk.
   `@st.cache\_data(ttl=3600)` keeps it hot for an hour.
- \*\*How it works:\*\* Tries parquet → csv, otherwise halts with an error. All exceptions are surfaced.
- \*\*Why it matters:\*\* Eliminates UI state jitter and prevents slow reloads midanalysis.
- \*\*Read it like this:\*\* If you see "No dataset found...", the app never reached analytics; fix paths first.

## **Column partitioning**

- \*\*What it does:\*\* `split\_features(df)` separates modelable numeric features
   (`X`) from metadata (`meta`) and removes excluded fields
   (`EXCLUDE\_FOR\_CLUSTERING`).
- \*\*How it works:\*\* Keeps numeric columns, fills NA with 0, preserves `Name` in metadata.
- \*\*Why it matters:\*\* Prevents leakage (ratings, rank, owners) into clustering and keeps numeric shape stable.

## **Global constants & theming**

 Color palettes (`PALETTE`, `CHART\_COLORS`) and CSS are set early to keep visuals consistent and readable (light backgrounds, muted text, clear hover states).

# 2) Sidebar controls (the "contract" for every downstream view)

#### wny this matters

Advanced Barch Gars/tables are \*\* post-filter actified a chart looks empty or odd, check the filter state first.

• \*\*Mechanic/theme matching pitfall:\*\* "All" can starve the result set fast; use "Any" when exploring.

## 3) Clustering & embedding (market segmentation)

#### K-means + PCA scaffold

- \*\*What it does:\*\* `fit\_clusterer(X\_all, k)` standardizes numeric features, chooses `k\_eff` safely (handles low-var cases), runs K-means (n\_init=10), and produces 2-D coordinates with PCA for plotting.
- \*\*How to read:\*\* Clusters are \*\*behavioral/structural\*\* segments across mechanics, players, time, etc. They're not "themes"; they're adjacent neighborhoods in feature space.
- \*\*Labeling:\*\* `generate\_cluster\_labels(view\_f)` names clusters using percentiles of complexity/playtime, mean rating ("weakly/solid/strongly"), and \*\*over-indexed\*\* mechanics/themes (`\_top\_diff\_features`). If insufficient signal, falls back to "Segment X: similar to <best example>".

#### What to do with it

- Use segment labels to anchor discussion: \*"Mid-weight, strongly rated, with worker-placement bias (medium play)"\* is more precise than "like mid Euros."
- If a label looks wrong, it's a data artifact; drill into the top mechanics/themes list inside Segment Explorer (below).

## 4) Derived fields & guardrails

- \*\*`Play Time Hours`\*\* = minutes/60, clipped to 10h (for axes sanity).
- \*\*\* Success Score \*\*\* = `(AvgRating 5) \* Owned Users / 1000` quick, transparent signal combining quality & reach.
- \*\*`Market Penetration`\*\* = `Owned Users / max(Owned Users)` normalizes reach to [0,1].

\*\*Interpretation:\*\* `Success Score` is not a profit estimate. It's a coarse "did this travel?" indicator relative to rating.

# 5) Home KPIs (top metrics row)

- \*\*Total Games (filtered)\*\* and new this year (delta).
- \*\*Market Segments:\*\* count of unique cluster labels in view.
- \*\*Median Rating / Median Owners / Avg Complexity / Success Rate (≥7.0).\*\*

\*\*How to read:\*\* These are \*\*context baselines\*\* for all subsequent charts. If median rating is 6.2, don't set a 7.8 target without a differentiator strategy.

Advanced Board Game Developer Console — Feature Reference

## 7) Design Wizard tab

## 7.1 Archetypes (presets + grounded insights)

- \*\*What it does:\*\* 3×3 chooser with defaults for cats/mechs/weight/time/players and a \*\*market insight\*\* that's terse, numeric, and specific.
- \*\*How to use:\*\* Pick the closest lane, then customize in the form.

## 7.2 Design form → Profile → Segmenting → Neighbors

- \*\*Mapping:\*\* Mechanics/themes toggles map to the dataset's one-hot columns (with guard for "Cooperative Game" and "Solo Mode").
- \*\*Cluster assignment:\*\* Input is aligned to `X\_all`, scaled with the clustering scaler, and assigned a `cluster\_id`.
- \*\*Nearest neighbors:\*\* Within that cluster, we compute pairwise distances in standardized space and take the `topn` closest.

\*\*Interpretation:\*\* Your "similar games" are \*\*feature-space\*\* nearest neighbors, not semantic clones. They are the right baseline for predictions and pricing.

## 7.3 Predictions (rating, owners, risk)

- \*\*Models:\*\* Tries to load `rating\_xgb` and `sales\_xgb`. If missing, falls back to \*\*neighbor averages\*\* with small noise for rating and median for owners.
- \*\*Confidence:\*\* Based on neighbor distance dispersion (tighter = higher). Clipped [40,95] to avoid false precision.
- \*\*Percentile:\*\* Where your predicted rating sits in the filtered market.

#### \*\*How to read:\*\*

- \*\*Predicted Rating:\*\* Directional. If it's high but owners low → likely niche; if owners high but rating mid → fun toy, tighten decisions.
- \*\*Confidence:\*\* A proxy for \*\*model agreement\*\* given your feature profile; low confidence → explore a second preset or adjust features (players/time) and recompute.

# 7.4 Pricing & Unit Economics

\*\*Inputs:\*\* MSRP, channel fees, unit COGS, shipping, returns, fixed costs, elasticity, sales window.

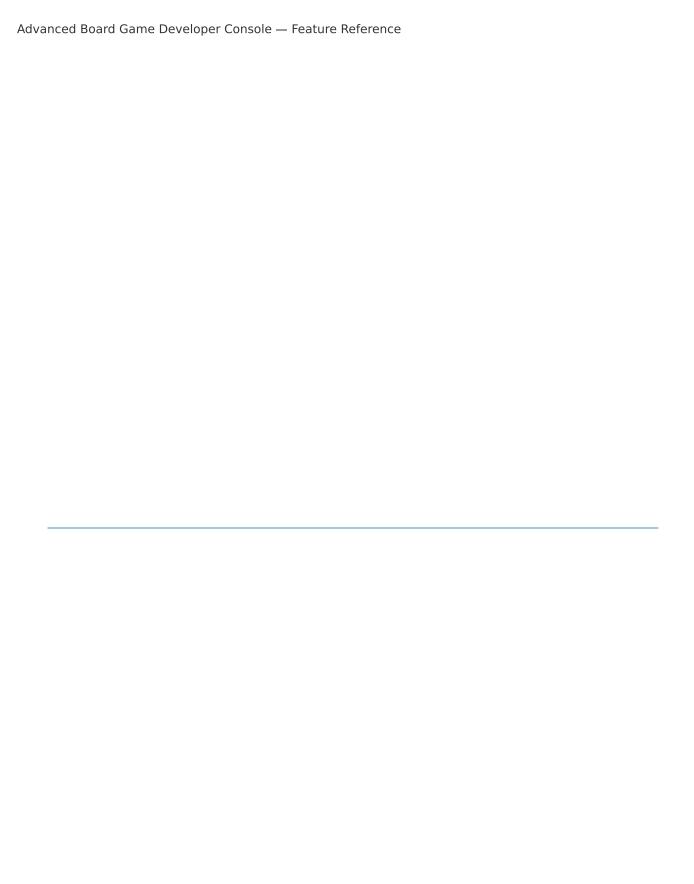
#### \*\*Core formulas:\*\*

• Anchor price (heuristic):

`base =  $35 + (complexity - 2.5)*6 + component/premium adders + players/playtime adders <math>\rightarrow$  clipped [15,150]`

- Net per unit: `msrp \* (1 channel fee pct)`
- Gross profit/unit: `net per unit (unit cogs + shipping)`
- Owners adjustment (elasticity, optional): `owners\_adj = owners\_base \* (msrp/anchor\_price)^elasticity` (clipped to ±40%)
- Effective units: `owners adj \* (1 returns pct)`
- Total GP: `gross profit per unit \* max(effective units, 0)`
- Net profit: `Total GP fixed\_costs`
- ROI multiple: `Net profit / fixed costs` (∞ if fixed=0)
- Break-even units: `fixed costs / gross profit per unit`
- Payback months: `ceil(breakeven units / (effective units / sales window))`
- Gross margin % (unit): `gross profit per unit / net per unit`

Advanced Board Game Developer Console —	Feature Reference



selected segment.
Advanced Board Game Developer Console — Feature Reference
<ul> <li>**What to do:** Prototype the top 1-2 changes; don't stack five.</li> </ul>

# 10) Mechanic Synergies tab

• \*\*How to use:\*\* If you've already pinned your complexity target, this narrows pairing options fast.

## 11) Narrative system

- \*\*What it does:\*\* `narr(md)` blockquotes explanatory text when "Show narrative insights" is on (sidebar).
- \*\*Why it matters:\*\* Keeps prescriptive copy out of code paths and lets you toggle it for focused data views (e.g., demos vs deep work).

## 12) Model and input alignment (robustness)

## Aligning inputs to training (`align\_profile\_to\_training`)

- \*\*Type coercion:\*\* Numeric cols coerced; bad inputs fall back to 0 (harsh but safe).
- \*\*Consistency:\*\* If `Min Players > Max Players`, we correct.
- \*\*One-hot groups:\*\* Ensures at least one `Cat:` or `Mechanic\_` feature is on if a group is present (prevents all-zeros traps).
- \*\*Clamping:\*\* If an input scaler is available, numeric inputs are clipped to the model's training min/max.

\*\*Why it matters:\*\* Keeps predictions from undefined territory and prevents "nans exploding the tree."

## 13) Performance & UX

- \*\*Caching:\*\* Data (`@st.cache\_data`) and models (`@st.cache\_resource`) minimize compute churn.
- \*\*Plot ergonomics:\*\* Unified hover modes where helpful; colorbars added explicitly when Plotly won't auto-explain encodings.
- \*\*Legend gaps:\*\* The mechanics network builds a manual legend (dummy traces) because every edge is its own trace.

## 14) How to interpret the entire console (meta-guidance)

# 15) Known edge cases / failure modes

- \*\*Empty views:\*\* Over-tight filters → "No games match..." Fix filters first.
- \*\*Skewed predictions:\*\* If neighbors are sparse (tiny segment), confidence will drop. Use a nearby preset to triangulate.
- $\bullet$  \*\*Elasticity sensitivity:\*\* If ROI flips sign with  $\pm$ \$5 price changes, you're bandwidth-limited by assumptions (fees, returns, or COGS). Stress-test those.
- \*\*Mechanic columns mismatch:\*\* If your dataset uses different prefixes or spellings, discovery helpers (`discover\_mechanics`, `discover\_themes`) will include them but presets may not light them up—toggle manually.

