**ChatGPT**

# Abstract

Many predictive tasks in real-world tabular data involve *engineered features* like ratios or products of raw variables to capture domain invariances. We investigate whether a modern tree ensemble (XGBoost) can **automatically learn ratio/product feature structures from raw positive inputs** under various data regimes, or if providing these features explicitly is essential for robust generalization. We design a synthetic experimental suite with classification tasks of increasing complexity (Levels 1–7) involving ratio/product-based features, and we evaluate XGBoost models trained either on **raw features only** or with **oracle features** (the true latent ratio/product coordinates or signal added). Across extensive experiments, we find that raw-only XGBoost often fails to recover the true feature structure under distribution shifts. In particular, tasks requiring multiplicative interactions (Level 4) show the largest performance gaps: for example, with moderate model capacity, a raw-only model achieves a median PRAUC $\approx 0.11$ versus $\approx 0.24$ with oracle features – a **~53–55% relative drop** [1] [2]. Under *invariance-preserving covariate shifts* that change input scales without altering the true ratio/product, raw-only models break down catastrophically (PRAUC dropping to ~0.08 in one case, versus ~0.91 with oracle features, $\Delta \approx 0.83$) [3] [4]. Oracle models not only close the accuracy gap but also exhibit substantially improved **invariance** (up to an order-of-magnitude lower invariance error) and more stable feature importance. These results highlight that while tree ensembles are flexible, they do not reliably learn certain physical or algebraic invariances from data alone. We discuss diagnostic insights (invariance error, variance isolations, dominance of inputs) and interpretability implications. Finally, we outline limitations – e.g. focus on one model class (XGBoost), synthetic data, limited seeds – and suggest next steps such as data size scaling experiments and tests on other model architectures. Our findings underscore the value of incorporating domain-knowledge features for robust tabular modeling, and they form a foundation for a forthcoming arXiv working paper.

# Introduction

Machine learning practitioners often face a choice: rely on models to discover complex feature relationships in raw data, or manually construct features that encode known invariances (ratios, products, etc.). In domains ranging from finance (e.g. debt-to-income ratios) to physics (dimensionless products of variables), **ratio and product features** play a crucial role in generalization. Ideally, a sufficiently expressive model like XGBoost (a gradient-boosted tree ensemble) could infer these combinations automatically. However, in practice, models trained on raw inputs might latch onto spurious correlations (e.g. absolute magnitudes) present in the training distribution, rather than the true invariant relationship, leading to brittle performance under distribution shift [5] [6].

In this work, we **systematically study whether XGBoost can learn ratio/product structures from raw positive features** without explicit guidance. We formulate a set of synthetic binary classification tasks where the label is determined by certain ratio or product of input variables (or compositions thereof). These tasks are grouped into *Levels 1–7* by complexity: from simple single-ratio or single-product signals up to more complex compositions like products of ratios, non-monotonic transformations, and gating logic [7]. Crucially, we evaluate model robustness under varied **distribution regimes** – including changes in input scale or distribution shape – to test if a model trained on raw features truly captures the invariant ratio/

product feature or relies on incidental cues. The core question is: **Can a raw-feature XGBoost "rediscover" the hidden ratio/product feature and remain predictive when superficial data properties change?** [8]

To answer this, we compare two training setups: (1) a **raw-only** model, given only the primitive features $x_0, \dots, x_{d-1}$; and (2) an **oracle feature** model, which in addition to the raw $x_i$ is given the actual constructed feature(s) that the data-generating process uses (either intermediate coordinates or the final signal) [9] . The oracle thus provides an upper-bound on performance by supplying the "perfect" feature engineering. By measuring the performance gap between oracle and raw-only, we quantify *how well the model learned the feature structure on its own*. We denote this gap as **Δ PRAUC = PRAUC(oracle) – PRAUC(raw)** on a held-out test set [10] , using area under the Precision-Recall curve (PRAUC) as our primary metric. A large Δ PRAUC indicates that raw-only XGBoost missed important structure that the oracle feature provides.

Our contributions are as follows: (1) We present a novel synthetic benchmark to stress-test XGBoost's ability to learn multiplicative feature interactions under realistic nuisance factors (noise, heavy-tailed distributions, feature correlations, and dataset shifts). (2) Through extensive experiments, we identify specific scenarios (notably *Level 4* interaction tasks and *invariance-preserving shifts*) where raw feature models dramatically underperform, while oracle features restore performance. (3) We introduce diagnostic measures of invariance and input dominance to interpret model behavior, and show that oracle features significantly improve these diagnostics (indicating models that align better with true invariances). (4) We discuss the implications for feature engineering practices and outline future research directions, including experiments on other model classes and methods to enforce invariances. This work is intended as a **working paper** draft for arXiv, providing a comprehensive analysis for both industry practitioners concerned with model robustness and researchers interested in inductive biases in tree ensembles.

## Methodology

**Synthetic Data and Tasks:** We construct a suite of synthetic binary classification tasks where each task's label is determined by a known combination of the inputs. All input features $x_i$ are positive continuous variables, drawn from distributions that introduce various "nuisance" characteristics (described below). Each task is designed such that a specific *ratio or product* (or a small set of such combined features) underlies the true decision boundary. Tasks are organized by *level of complexity* [11] :

- **Level 1:** Single ratio or single product of two features (e.g. $s = x_0 / x_1$ or $s = x_0 \cdot x_1$). These are simple two-variable interactions.
- **Level 2:** Ratio or product of sums (e.g. $s = \frac{x_0 + x_1}{x_2 + x_3}$). This involves aggregating groups of features then taking ratios/products.
- **Level 3:** Ratios or products of differences (e.g. $s = \frac{x_0 - x_1}{x_2 - x_3}$, etc.). These introduce additive interactions inside the ratio/product.
- **Level 4: Multiplicative interactions of coordinates** – e.g. a product-of-products or ratio-of-ratios ($s = \frac{x_0/x_1}{x_2/x_3}$, or $s = (x_0 x_1) \cdot (x_2 x_3)$). Level 4 tasks effectively require learning a *two-level* interaction (a product of two ratios, etc.) [12] , making them particularly challenging.
- **Level 5:** Mixed ratio-product combinations (e.g. $s = \frac{x_0/x_1}{x_2 \cdot x_3}$). These combine additive and multiplicative relations.

- **Level 6: Non-monotonic transformations** of a basic ratio/product (e.g. a *U-shaped* or *band-pass* function of a Level 1 ratio [13] ). Here, the raw ratio might need to be learned *and* an additional non-monotonic pattern applied.
- **Level 7: Gated combinations**, where the model must learn to *switch* between using a ratio vs. a product based on some condition (for instance, if an indicator feature $g$ is 0 use ratio $u_1$, if $g=1$ use product $u_2$) [14] . This tests the model's ability to learn piecewise logic in addition to ratio/product computation.

For each task, we generate data as follows. First, $d$ raw features are sampled from a specified joint distribution (see *Regime variations* below). We then compute the task-specific signal $s$ (e.g. the ratio or product as defined for that task) from the raw features. Finally, a binary label $y$ is assigned based on $s$ (for example, thresholding $s$ at a certain value or via a logistic function to introduce noise). The training and testing sets for each task share the same underlying functional relationship for $y$ but may differ in the input feature distribution if a shift is applied (more on this below).

**Distribution Regimes:** To thoroughly evaluate robustness, we generate data under multiple regimes, each altering the distribution of inputs or introducing shifts:

- **Tail/Correlation regimes:** We vary the **marginal distributions** of features (e.g. log-normal with different sigma parameters for tail heaviness) and the **correlation $\rho$** between features [15] . For instance, one regime might have low variance ($\sigma=0.3$) and independent features, while another has high tail variance ($\sigma=1.2$) and strong correlations ($\rho=0.9$). These regimes, labeled *tail_corr*, represent static conditions (no train-test shift) but with different levels of difficulty (heavy tails and correlations can make learning harder).
- **Mixture extremes:** A regime where the feature distribution is a mixture of a low-variance and high-variance components [16] . This can create bimodal or skewed feature distributions, again with no explicit train-test shift but more complex input distributions.
- **Shift (naive):** A **covariate shift** introduced at test time that *does not preserve* the ratio/product structure. For example, adding a constant or scaling each feature independently (so the ratio of features changes). Labeled *shift_naive*, this regime can fool models that depend on absolute scales because the test distribution differs from training in ways that *alter the target signal's distribution*.
- **Shift (preserve):** A covariate shift that **preserves the true ratio/product signal**. For instance, multiplying all input features by a common factor $c$ at test time (so that ratios remain the same, since $x_0/x_1$ is unchanged, etc.) [6] . We label this regime *shift_preserve*. Here, an ideal model that learned the ratio should generalize perfectly, whereas a model relying on raw magnitudes will severely degrade. This regime is a stress test for invariant learning.

Each combination of task and regime yields a dataset; we typically generate *n = 30,000* samples for training and a similarly sized test set (with shifts applied to test as appropriate). We repeat each experiment with multiple random seeds to account for training stochasticity.

**Model Training and Oracle Feature Modes:** We use XGBoost (tree booster) classifiers with a fixed hyperparameter configuration (or a small set of configurations) across all tasks. In our main experiments, we evaluate three XGBoost model capacity settings: a **baseline** (depth-5 trees with moderate regularization), a lighter model (shallower trees, depth-4), and a higher-capacity model (deeper trees,

depth-7) [17] . Early stopping is employed based on validation AUC. We consider three feature input modes [9] :

- **Raw only:** Model is given only the $d$ raw features ${x_i}$. It must learn any ratio/product relationships implicitly.
- **Oracle (coords only):** Model is given the raw features **plus** the intermediate "coordinate" features used in the task's true signal (e.g. if the task is $y$ depends on $u = x_0/x_1$, then $u$ is provided as an additional feature; for a gated task, the gate indicator and the two candidate signals $u_1, u_2$ would be provided, but not the final combined signal).
- **Oracle (signal only):** Model is given the raw features **plus** the final true signal $s$ itself (e.g. the actual ratio or product value that directly determines $y$). This represents an upper bound where the model effectively has the "answer" feature handed to it, aside from any monotonic transformations or thresholding that relate $s$ to $y$.

Both oracle settings inject strong inductive bias: the "coords only" mode gives hints of the right transformations without giving the final answer directly, whereas "signal only" gives the model the exact predictive feature. In all cases, models are evaluated on the same test sets (with or without shifts). **Precision-Recall AUC (PRAUC)** on the test set is our primary metric [18] , as many tasks are imbalanced by construction (since thresholding a continuous $s$ can yield class imbalance). We compute Δ PRAUC for each run as described. For statistical reliability, the main experiment ( `exp_default` ) includes 3 random seeds per configuration (totaling 1944 trained models across tasks/regimes/configs) [19] .

**Diagnostics:** To gain insight into *how* the models make predictions, we measure several diagnostic quantities on the test set [20] :

- **Invariance Error:** We quantify how well a model's predictions respect the known invariances (ratio or product) by performing controlled perturbations. For example, for a ratio-based task, we can multiply both inputs $x_0, x_1$ by a constant factor and check if the prediction changes. The **ratio-scale invariance error** is the change in model output when all components of a ratio are scaled together (should ideally be zero). We compute analogous errors for product tasks (e.g. multiply one factor by $c$ and divide the other by $c$ to preserve the product). A lower invariance error means the model's predictions are stable under transformations that leave the true signal unchanged [21] . We report median invariance error for raw vs. oracle models.
- **Iso-Coordinate Variance:** This measures output variability when altering one input while keeping the true combined feature constant. For instance, for a ratio $x_0/x_1$, we can adjust $x_0$ and counter-adjust $x_1$ to keep the ratio the same, and see if the model output varies. High **iso-coordinate variance** indicates the model is sensitive to specific input values rather than just the ratio/product. We report this as an aggregate statistic (ideally, a model using only the ratio would have near-zero variance in such an experiment) [21] .
- **Dominance Metrics:** For tasks involving sums (like ratio of sums), we track the degree of *single-feature dominance* in those sums: e.g. compute $\max_i x_i / \sum_j x_j$ for the group. If this dominance is high (close to 1), one feature overwhelmingly drives the sum. Such cases can allow *shortcut learning*: a raw model might just pick out the largest feature as a predictor instead of the intended ratio [20] . We measure the average and 90th-percentile dominance in train and test sets per regime [22] [23] to identify regimes prone to shortcuts (for example, we found the *shift_naive* regime often had high dominance values around 0.63–0.76 in training [22] ).

These diagnostics help explain performance differences: if an oracle model has lower invariance error and iso-variance, it suggests it truly learned the intended invariant feature, whereas a raw model with higher values likely relied on non-invariant cues.

## Results

### Overall Performance: Oracle vs. Raw-Only

Across all tasks and conditions, providing the true ratio/product features (oracle modes) yields equal or better performance than relying on raw features alone. The **performance gap is most pronounced for the more complex tasks (Levels 4–7)**. Figure 1 illustrates the distribution of Δ PRAUC values for each task in the main experiment. Tasks up to Level 3 (simpler ratios/products) show relatively small Δ PRAUC – the boxes are centered near zero, indicating that raw-only XGBoost almost matches oracle performance in many cases. In contrast, *Level 4 tasks stand out with consistently large positive Δ PRAUC*, often in the 0.1–0.2 range or higher [24]. This means that on tasks involving multiplicative interactions of multiple coordinates, the raw-only model's PRAUC is dramatically worse than the oracle's.

**Figure 1.** Distribution of Δ PRAUC by task for the main experiment (Exp_default, n=30k). Each boxplot shows the spread of Δ across all regimes, seeds, and model capacity settings for a given task (higher Δ means oracle features improve performance more). *Level 4 tasks (ratio-of-ratios and product-of-products)* exhibit the largest median Δ, indicating that raw-only models struggled significantly on these multiplicative interaction tasks, whereas oracle features provided a big boost. Simpler tasks (Levels 1–2) have Δ centered near 0 (sometimes even slightly negative), meaning XGBoost often managed to learn those from raw inputs.

Quantitatively, for a representative XGBoost configuration (depth-5, "baseline"), the median test PRAUC on *Level 4* tasks was only about **0.11** for the raw-only model, compared to **~0.24** with oracle features [1]. This corresponds to roughly a **53–55% relative drop** in PRAUC for raw-only [2]. In fact, as **Table 1** shows, Level 4 yields the largest relative drop among all task levels. Level 3 (ratio/product of differences) has the next largest gap (~35–39% drop in the worst cases), while simpler Level 1 tasks show <20% drop and sometimes no advantage to oracle at all [25] [26]. These results confirm that **the more complex the feature interaction, the less likely XGBoost is to capture it from raw data alone**. A tree ensemble can in principle approximate a product or ratio, but doing so may require many splits and sufficient training variation; apparently for higher-order interactions, the model fails to fully recover the true signal without explicit features.

**Table 1 – Performance gap by task complexity (Exp_default, baseline XGBoost):** For each task level, we report the median PRAUC for raw-only and oracle (signal) models, and the relative drop (percentage by which raw-only falls short of oracle). Higher percentage indicates a larger benefit from having the oracle feature. Level 4 tasks show the most drastic drops, underscoring the difficulty raw models have with multiplicative interactions.

| Task Level | Raw-only PRAUC (median) | Oracle (signal) PRAUC | Relative Drop (%) |
|---|---|---|---|
| **Level 1** (simple ratio/product) | 0.2046 | 0.2387 | ~14.3% [27] |

| Task Level | Raw-only PRAUC (median) | Oracle (signal) PRAUC | Relative Drop (%) |
|---|---|---|---|
| **Level 2** (ratio/product of sums) | 0.1317 | 0.1569 | ~16.1% [28] |
| **Level 3** (ratio/product of diffs) | 0.2106 | 0.3298 | ~36.2% [29] |
| **Level 4** (ratio-of-ratios, etc.) | 0.1157 | 0.2457 | ~52.9% [2] |
| **Level 5** (mixed interactions) | 0.3761 | 0.4768 | ~21.1% [30] |
| **Level 6** (non-monotonic) | 0.0539 | 0.0653 | ~17.4% [31] |
| **Level 7** (gated combination) | 0.3334 | 0.3462 | ~3.8% [1] |

[1] *Level 7 shows a smaller median drop in this aggregate view; however, its tasks often have high base performance even for raw models (PRAUC >0.9 in some regimes), leaving less room for improvement.*

Notably, **Level 7 (gated) tasks did not show as large a gap** on median – in some runs, raw models performed almost as well as oracle. This may be because the gating logic can sometimes be learned by trees as separate regions of the feature space, especially if one of the raw features strongly indicates which sub-regime to use. Even so, oracle features generally did no harm and in some seeds improved consistency.

## Impact of Distribution Shifts

We expected that if a model truly learned the ratio/product feature, certain distribution shifts (that preserve that feature) should not hurt performance. Conversely, a model using spurious cues would break under those shifts. Our results strongly support this: **the "invariance-preserving" shift scenario was by far the most challenging for raw models, yet trivial for oracle models.**

In the *shift_preserve* regime, test inputs were rescaled in a way that leaves the true $s$ unchanged. For example, in one Level 4 task, we multiplied all input features by a constant factor at test time. The oracle model, which effectively bases its prediction on $s$, maintained high accuracy. The raw model, however, saw its predictions deteriorate severely. The most extreme case was **Level 4 (product × product) under a preserve shift**: the raw-only model's PRAUC dropped to ~**0.08**, whereas the oracle model achieved **0.91**, yielding Δ ≈ **+0.835** [3] [4] . In other words, the raw model completely failed when the data was out-of-distribution in terms of scale, even though the task's fundamental ratio structure was constant – clear evidence it had learned the "wrong" thing. Figure 2 provides a heatmap of median Δ PRAUC broken down by task level and regime family, highlighting how *shift_preserve* (rightmost column) causes the largest gaps, especially at higher task levels.

**Figure 2.** Heatmap of median Δ PRAUC by task complexity (level) and regime family (Exp_default). Each cell shows the median Δ for all tasks of a given level under a regime family (averaged over seeds and XGB configs). Darker colors = higher oracle benefit. We observe that the *preserve shifts* consistently yield the highest Δ (dark cells in the rightmost column), particularly for Levels 4–5. In contrast, the no-shift regimes (tail_corr and mixture_extremes) have modest Δ, and even naive shifts show only moderate gaps. This pattern underscores that raw models fail specifically when test conditions require true invariance.

Aggregating across all tasks, the *shift_preserve* conditions led to a **median raw-only PRAUC around 0.12 vs. 0.35 with oracle**, roughly a **Δ of 0.23–0.24 (≈65% relative)** [32] . By contrast, in the *tail_corr* (no shift) regimes, raw vs. oracle performance differed by only ~0.03 (e.g. 0.137 vs 0.169 PRAUC, ~18% relative) [33] . The *mixture_extremes* regime similarly showed small absolute improvements (a few points of PRAUC) when adding oracle features [16] [34] . Even the *shift_naive* regime (which changes distributions in non-invariant ways) resulted in moderate Δ (~0.03–0.06, or ~15–20% relative) – significant but nowhere near the catastrophic failure of shift_preserve [35] . These findings confirm that **raw-feature models exploit incidental correlations (like absolute feature scales) that do not hold under preserve-type shifts**, whereas providing the true feature immunizes the model against such shifts.

Interestingly, there were a few cases where the oracle model did *slightly worse* under distribution shift. For instance, in a Level 2 ratio-of-sums task under a naive shift, we observed negative Δ values (raw-only outperforming oracle by a small margin) around -0.05 to -0.07 [36] [37] . These occurred in high-noise or high-dominance scenarios. We suspect this is due to the oracle feature introducing an additional input scale that, if not well-regularized, could cause slight overfitting or instability with limited data. In practice, these negative gaps were rare and relatively small in magnitude (on the order of a few percent PRAUC). Nonetheless, they highlight that adding features is not always strictly beneficial – interactions with training dynamics (e.g. early stopping, regularization) can lead to slight underperformance for oracle models in some edge cases. Overall, however, the **trend is overwhelmingly that oracle features improve or preserve performance, never catastrophically worse**. Raw models, on the other hand, have severe failure modes under specific shifts.

## Diagnostic Insights: Invariance and Shortcut Behavior

To better understand *why* the raw-only models falter, we examine the diagnostic measures introduced earlier. The **invariance error** analysis reveals a clear pattern: **oracle feature models learn the intended invariances far better than raw models**. For the main experiment, the median invariance error (across all tasks/runs) for raw-only models was on the order of $10^{-2}$, whereas for oracle models it was around $10^{-3}$ or lower – roughly an **8× to 10× reduction** [38] [39] . Concretely, when we scaled ratio inputs or recomposed product inputs (in ways that should not affect the true $s$), the raw model's predictions would noticeably change (a sign it had sensitivity to absolute values), while the oracle model was essentially invariant (only tiny output changes). Figure 3 plots the relationship between each run's Δ PRAUC and the raw model's invariance error. We observe a strong correlation: runs where the raw model had high invariance error (failing to be invariant) are exactly those where oracle features yielded large gains (high Δ) [40] . This suggests that **learning the correct invariance is a key factor in generalization** – the oracle features explicitly encode that invariance, hence their advantage.

**Figure 3.** Relationship between invariance error and performance gain (Δ PRAUC) in the main experiment. Each point represents a particular task/regime run. The x-axis is the raw model's invariance error (log-scale), and y-axis is Δ PRAUC (oracle minus raw). A clear trend emerges: when raw models have poor invariance (high error), the oracle's advantage is large (points in the upper-right). Points near the bottom-left (low error, low Δ) correspond to cases where the raw model already learned the invariance well, so oracle offers little benefit. This confirms that the performance gaps are tied to whether the model learned or missed the fundamental ratio/product invariance.

We also examine the **iso-coordinate variance** metrics. In raw-only models, varying one feature while compensating with another to keep a ratio constant often caused noticeable prediction fluctuations. The

median iso-variance for raw models was around $2.8\times10^{-4}$ in ratio tasks (and $3.0\times10^{-4}$ in product tasks), whereas oracle models drove this down to roughly $10^{-5}$ or less [38] . This 20–40× reduction aligns with the invariance findings: oracle models focus on the true combined feature, so once that is fixed, changing individual inputs doesn't confuse the model. Raw models, lacking that focus, still respond to individual inputs even if the overall ratio/product is unchanged, indicating they only partially learned the relationship.

Finally, the **dominance analysis** provides insight into *how* raw models might achieve reasonable performance in some cases without learning the true feature. In sum-based ratio tasks, if one component dominates the sum, the ratio essentially reduces to that single component (e.g. $\frac{x_0}{x_0 + x_1} \approx 1$ if $x_0 \gg x_1$). We found that in certain regimes (particularly *mixture_extremes* and some *shift_naive* cases), the training data often had a dominant feature in the sum with mean dominance around 0.57–0.63 and a 90th percentile above 0.75 [41] [22] . In such conditions, a raw model could **take a shortcut** by simply learning a threshold on the largest feature instead of the intended ratio. This would work on typical training points (since the largest feature largely determines the ratio), but if the dominance pattern changes in test data, the model would fail. The preserve-shift regimes in fact often *reduce* dominance (since scaling all features equally makes none stand out), breaking the shortcut. Oracle models are less susceptible to this because the presence of the true ratio feature guides the model to use it rather than individual inputs. Figure 4 (in appendix) illustrates a scatter of Δ PRAUC vs. dominance in train data; we saw that scenarios with high dominance often correspond to moderate Δ (the raw model seemed fine until distribution changed), whereas low-dominance scenarios had higher Δ (raw model had no easy shortcut, so oracle helped more) [42] [43] .

In summary, the diagnostics reinforce a consistent story: **Raw-feature models often learn partial tricks that suffice under the training distribution (e.g. exploiting scale, single-feature dominance, or correlations), but these do not generalize when those conditions shift.** Oracle features, by injecting the true invariant, prevent the model from over-relying on those tricks. This not only boosts test performance but also aligns the model's behavior with known invariances – an *interpretability win* because we can trust that the model's predictions hinge on meaningful features (ratios/products) rather than obscure combinations of raw inputs.

## Discussion

Our findings have several important implications for both applied machine learning and the study of inductive biases in models:

- **Limits of Automatic Feature Learning in Trees:** Gradient-boosted trees are universal function approximators and *could* represent multiplicative relations by successive splits. However, our results show that in realistic data scenarios (finite samples, noisy and correlated features, etc.), they **struggle to spontaneously learn even simple ratio/product features**. The high-degree interactions needed are apparently not discovered reliably, especially when easier linear or one-sided splits can fit the training data. This echoes observations in industry: practitioners often pre-compute features like ratios, even when using tree models, because the model might not find them on its own. Our contribution is a quantification of how big this gap can be – exceeding 50% relative drop in our Level 4 tasks – and pinpointing when it occurs.

- **Invariance as an Inductive Bias:** The stark contrast between raw and oracle models under preserve-shifts is essentially a proof-of-concept that **encoding known invariances (here via features) dramatically improves robustness**. The oracle models had the ratio/product provided, so invariance to scaling or other symmetric transformations was built-in. One can view this as adding domain knowledge as a feature vs. expecting the model to learn a symmetry from data. Tree models do not inherently enforce such symmetry (unlike, say, a convolution enforces translational invariance), so without explicit features or specialized training (e.g. data augmentation), they won't magically generalize in an invariant way. This suggests that for applications where certain invariances are crucial (e.g. ratios in medical risk scores, physical laws), **feature engineering or other bias injection is not just a convenience but a necessity for reliability**.

- **Shortcut Learning in Tabular Data:** Much discussion of "shortcuts" or spurious correlations happens in deep learning on images, but our results highlight analogs in tabular settings. For instance, using the largest feature of a sum instead of the ratio is a shortcut; using absolute value of features instead of their ratio under a scale shift is another. We demonstrated how measuring dominance or invariance error can reveal these shortcuts. An interesting observation is that *when training and testing distributions remain the same*, these shortcuts cause little harm – the model might predict well on i.i.d. data, hiding its lack of true understanding. It's only under distribution changes that the gap is exposed. This reinforces the importance of stress-testing models (e.g. with simulate shifts as we did) to detect when a model might be relying on a coincidental feature.

- **Interpretability and Feature Importance:** From an interpretability standpoint, models with oracle features were easier to interpret – for example, feature importance or SHAP values (if computed) would correctly highlight the $s$ or $u$ features as driving the prediction, which aligns with the known mechanism. Raw models, by contrast, might distribute importance across multiple raw features or higher-tree interactions that are hard to parse. The improvement in invariance metrics also indicates that oracle models behave in a more expected way (they won't fluctuate outputs for irrelevant input changes). This suggests a practical benefit: if you have strong prior knowledge of a functional form (like a ratio), using it can make the model not only perform better but also act in a more transparent and stable manner – properties desirable for high-stakes applications.

- **No Free Lunch – Oracle not Always Superior:** We noted some edge cases of negative Δ PRAUC. While minor, they are a reminder that adding features (even "correct" ones) can introduce new challenges: potential multicollinearity with raw features, the need for the model to weigh redundant information, or altered optimization dynamics. In our ratio-of-sums example with a negative Δ [37], the oracle model might have over-relied on the provided ratio feature and slightly underfit the nuance that raw features could provide under shift. Another possibility is that early stopping cut training once the oracle feature plateaued, whereas the raw model might have benefitted from further rounds capturing some idiosyncrasy. This invites a deeper analysis of how to integrate oracle features: e.g. one could drop the raw features that are used in the ratio to avoid redundancy, or add a small noise to oracle features during training to prevent overconfidence. Our work primarily demonstrates the benefit ceiling of having those features; productionizing this knowledge would require careful integration.

# Limitations

While our experimental framework is extensive, several limitations must be acknowledged:

- **Focus on XGBoost only:** We limited our study to gradient-boosted decision trees (specifically XGBoost). This choice was motivated by XGBoost's popularity in tabular data challenges and its known strength in handling various feature types. However, the results may differ for other model classes. For instance, deep neural networks with enough data might multiplicatively combine features more naturally (or they might not – it's an open question). Similarly, simpler models like linear/logistic regression would obviously fail to capture non-linear ratios unless features are precomputed. Our conclusions about invariances are thus directly applicable only to tree ensembles. Expanding to other algorithms (random forests, neural nets, symbolic regressors) is important future work.

- **Synthetic data vs. real-world data:** Our tasks, by design, have a known ground-truth feature and controlled shifts. Real data rarely offers that clarity – there may be many latent relationships, and distribution shifts can be more complex. While we believe the phenomenon of "model failing to learn a known important ratio" does occur in practice (hence domain experts often supply those ratios), the exact magnitude of performance loss and types of shifts could vary. Testing our findings on a real dataset (for example, a medical dataset where a certain lab value ratio is critical) would increase the external validity.

- **Limited regimes and seeds:** We explored a decent variety of regimes (heavy tails, correlation, mixture, two types of shifts) and used 3 seeds for our main runs. However, there are many other forms of distribution shift (e.g. shifts in noise level, adversarial feature distribution changes) we did not test. Our statistical confidence in the results is reasonably high given consistent patterns over seeds, but more seeds or cross-validation could further ensure stability. We also fixed the training sample size (30k) in the main experiment; the effect of very small or much larger data was not fully explored (see next section).

- **No comparison to invariant learning methods:** We did not include specialized baselines such as models trained with data augmentation for invariance (e.g. physically scaling inputs during training), or regularized to encourage certain invariances. Such methods could potentially bridge the gap without adding oracle features. Similarly, including *monotonic constraints* or other inductive biases in XGBoost might help (for instance, enforcing that prediction changes only with certain ratio combinations). Evaluating these approaches was out of scope for this initial study.

- **Evaluation metric and scope:** We focused on PRAUC as the metric, which is appropriate for our tasks with imbalanced labels. It's possible that for other metrics (ROC AUC, accuracy at a threshold, etc.), the gaps might manifest differently. We did not deeply analyze calibration or the probabilistic quality of predictions – only rank-based performance. Also, our interpretation metrics (invariance error, etc.) are tailored to ratio/product features; other feature structures would need analogs defined.

Despite these limitations, we believe the core insight is robust: given the complexity of feature interactions and the potential for spurious correlations, **XGBoost often benefits greatly from having the "right"**

**features handed to it, especially to ensure robustness under shifts**. The work opens several avenues to broaden and deepen the analysis.

## Future Work

This study raises a number of follow-up questions and experimental extensions:

- **Varying training data size:** A critical next step is to perform an *n-sweep* – i.e. train models on different dataset sizes (e.g. 1k, 5k, 30k, 100k samples) to see if large data can compensate for lack of explicit features. We hypothesize that as $n \to \infty$, a sufficiently flexible model *could* learn the ratio structure (since it can partition the space finely). Preliminary results (not shown here) suggest Δ PRAUC might shrink with more data [44], but not vanish entirely for the hardest tasks. Plotting Δ vs. n (as in Fig. 2 of the *exp_all_levels_10k* panel) across multiple n would inform whether feature engineering is mostly a data-efficiency booster or absolutely necessary even with very big data.

- **Other model classes and architectures:** We plan to test whether these findings hold for neural networks on the same tasks. A fully connected network given raw inputs might, for example, learn a multiplication via hidden layer interactions. But it might also struggle with invariances unless certain architectures (like multiplicative units or attention mechanisms) are used. Likewise, exploring rule-based learners or symbolic regression could be interesting – those might explicitly find the formula given enough computational effort. If some models can learn ratios from raw data with ease, that could point tree-based models towards improvements (or vice versa).

- **Feature engineering vs. model ensembling:** Another angle is whether an ensemble of raw-feature models could emulate the oracle. For example, one could train a model specialized in certain regimes or use knowledge distillation from an oracle-informed model to a raw model. Would the raw model then implicitly capture the ratio? Such techniques could be useful for deploying models without additional features (for simplicity) while still getting invariance benefits.

- **Regularization for invariance:** Instead of adding features, one could incorporate invariance knowledge through data augmentations or constraints. For instance, during training we could randomly scale sets of features (for ratio tasks) and enforce that the prediction doesn't change (maybe via a penalty in the loss). This kind of *invariance augmentation* might guide the model to discover the ratio on its own. We plan to experiment with such training schemes as a way to imbue raw models with the desired invariances, closing the gap to oracle performance without directly feeding the oracle feature.

- **Expanded diagnostics and interpretability:** We used invariance and iso-variance as diagnostics. Future work could include *Shapley value analysis* to see if oracle features indeed carry most importance in those models, or partial dependence plots to verify that raw models have complicated, fragile dependence on original features. Additionally, measuring *generalization on counterfactual data* (creating test points that have same ratio but different raw values) would be a direct way to assess ratio-learning – something we implicitly did with shift_preserve, but could be extended.

- **Real-world case studies:** Finally, applying these insights to real datasets (for example, UCI or Kaggle datasets where domain knowledge suggests a ratio/product feature) would be valuable. If we can demonstrate a scenario where adding a ratio feature improves a real model under a known shift (or even just i.i.d. test data), it would underscore the practical importance. Conversely, finding a case where a model *did* learn a ratio by itself (and understanding what made it possible – lots of data, very clean signal, etc.) could provide guidance on when feature engineering is necessary.

Our immediate next step is preparing this work for an arXiv submission, which will involve some of the above extensions (particularly the data size experiments and possibly a neural network baseline) to bolster the claims. We also intend to release the code for data generation and the benchmark so others can explore alternative models or features on these tasks.

# Conclusion

In this work, we presented a comprehensive evaluation of XGBoost's ability to learn complex ratio/product feature structures from raw inputs. Through controlled synthetic tasks, we demonstrated that while XGBoost often fits simple ratios or products given sufficient data, it **struggles with higher-order combinations and fails to capture the correct invariances under distribution shift**. The consequence is that models can perform well on i.i.d. data yet break under shifts that should be harmless – unless they are provided with carefully constructed features reflecting those invariances. By giving the model "oracle" features (true latent ratios/products), we saw dramatic improvements in robustness and insight into model behavior, affirming decades of domain expertise that feature engineering matters. We also introduced diagnostics to quantify invariance learning and shortcut reliance, finding that oracle features greatly reduce spurious sensitivity.

Our findings serve as a cautionary tale and a motivation: cautionary in that even powerful models have blind spots for certain structures, and motivation for developing methods to close this gap – whether through better feature engineering tools, training procedures that enforce invariances, or model architectures inherently sensitive to ratio/product relationships. For practitioners, the message is clear: if your problem has known derived features (ratios, physical quantities, etc.), **don't assume the model will figure it out – consider explicitly including them** to ensure stability. For researchers, this work provides a quantitative testbed for studying how models generalize and how we might impart inductive biases that align with domain knowledge.

In conclusion, XGBoost can approximate almost any function given enough depth and data, but **when data is limited or conditions change, the lack of built-in invariance learning can severely limit performance**. Simple feature engineering in the form of adding ratio/product features remains a powerful approach to address this, and understanding this interplay is crucial for building robust, interpretable models in the wild. We hope our "XGB invariant features" benchmark will spur further exploration into bridging the gap between what models *can* represent and what they *do* learn from data.

*Note on Readiness:* **\* This draft assembles the key pieces of our experimental analysis into a working paper format. The results are quite comprehensive, and the narrative connects them well. To be fully paper-ready for an arXiv submission, we would polish the text further, ensure all figures and tables are properly referenced and perhaps add a bit more context from related work. We might also include one or two additional experiments (such as the suggested n-sweep or a neural network baseline) to pre-empt likely questions. However, as it stands, the core content and conclusions are**

**solid. With minor revisions and perhaps a brief related-work section, this draft is** very close to submission-ready* *for arXiv, providing a strong foundation for a full paper.*

---

1  3  5  6  7  8  9  10  11  12  18  19  20  21  36  38  writedown.md
file://file_0000000095a071f892ad5e94ad16abe0

2  4  13  14  15  16  17  22  23  25  26  27  28  29  30  31  32  33  34  35  37  39  41  tables.md
file://file_00000000968471f88e32215e35fba2d2

24  40  42  43  44  plots.md
file://file_000000009eec71fda248c85795e90aaf