# Purifying Interaction Effects with Functional ANOVA: Paper Walkthrough + Two Experiments

Aik Tarkhanyan

https://github.com/HaykTarkhanyan/fanova_purification

# Agenda

**Part I: Paper Walkthrough**

1. The identifiability problem
2. Functional ANOVA target
3. Mass-moving algorithm
4. Convergence & properties
5. Weight estimation ($w$)

**Part II: Experiments**

6. German Credit (real data)
   Main effects & interactions under different $w$

7. Stress-test (synthetic)
   How dependence $\rho$ affects purification

**Conclusions**

Summary & takeaways

## 1) Additive models with interactions

We consider models of the form:

$$Y \approx f_0 + f_1(X_1) + f_2(X_2) + f_3(X_1, X_2)$$

More generally:

$$F(X) = f_0 + \sum_{i=1}^{d} f_i(X_i) + \sum_{i \neq j} f_{ij}(X_i, X_j) + \cdots$$

- Goal: interpret $f_i$ (mains) and $f_{ij}$ (interactions).
- Problem: without constraints, the decomposition is **non-identifiable**.

# 2) Identifiability issue: effects can move across orders

For a fixed predictor, you can shift "mass" between terms while keeping predictions identical.

### Simple intuition

Add any $h(X_1)$ to $f_1$ and subtract the same $h(X_1)$ from $f_{12}(X_1, X_2)$: the sum stays unchanged.

- This allows **contradictory interpretations** for the exact same function.
- So "main effect" and "interaction" are not uniquely defined by the predictor.
- This motivates a **canonical** notion of "pure" interactions.

# 3) Boolean example (OR/AND as XOR + mains)

Let $X_1, X_2 \in \{0, 1\}$.

## Key equivalence

$$X_1 \vee X_2 = 0.25\,(X_1 \oplus X_2) + 0.5\,(X_1 - 0.5) + 0.5\,(X_2 - 0.5) + 0.75$$

Similarly:

$$X_1 \wedge X_2 = -0.25\,(X_1 \oplus X_2) + 0.5\,(X_1 - 0.5) + 0.5\,(X_2 - 0.5) + 0.25$$

- The interaction parts are (centered) XOR up to sign; the rest is absorbed into mains+intercept.
- AND and OR have identical interaction structure (centered XOR up to sign); differences are absorbed into mains/intercept.

# Representational degeneracy



**Takeaway:** Multiple decompositions can yield identical predictions but assign credit differently across mains/interactions; purification selects a canonical one.

# 4) Multiplicative model and the $(\alpha, \beta)$ degree of freedom

A classic interaction model:

$$Y \approx a + bX_1 + cX_2 + dX_1X_2.$$

An equivalent re-parameterization (for any $\alpha, \beta$):

$$Y \approx (a - d\alpha\beta) + (b + d\beta)X_1 + (c + d\alpha)X_2 + d(X_1 - \alpha)(X_2 - \beta).$$

- Identifiable as a function, but coefficients are **not meaningful** without rules for choosing $\alpha, \beta$.
- A canonical choice is the fANOVA one, which yields **minimum variance in higher-order terms**.
- Conceptually, centering constants align with moments like $\mathbb{E}_w[X_1]$ and $\mathbb{E}_w[X_2]$; the decomposition depends on $w$.

# Multiplicative example: varying $(\alpha, \beta)$



**Takeaway:** As $(\alpha, \beta)$ vary, variance attribution shifts between mains and interactions even though predictions do not.

# 5) Functional ANOVA (fANOVA): the canonical target

The functional ANOVA decomposes $F$ into **orthogonal** components under weight $w(X)$:

$$\underbrace{F(X)}_{\text{predictor}} = \underbrace{f_0}_{\text{intercept}} + \underbrace{\sum_i f_i(X_i)}_{\text{main effects}} + \underbrace{\sum_{i<j} f_{ij}(X_i, X_j)}_{\text{pairwise interactions}} + \cdots$$

**Orthogonality:** $\mathbb{E}_w[f_u \cdot f_v] = 0$ for $u \neq v$.

Equivalently, each component satisfies **integrate-to-zero** (marginalizing out any variable gives 0):

$$\boxed{\mathbb{E}_w[f_u(X_u) \mid X_{u \setminus i}] = 0 \quad \forall i \in u}$$

- For main effect $f_i$: $\mathbb{E}_w[f_i(X_i)] = 0$ (centered).
- For interaction $f_{ij}$: $\mathbb{E}_w[f_{ij} \mid X_i] = 0$ and $\mathbb{E}_w[f_{ij} \mid X_j] = 0$ (no marginal structure).

# 6) Piecewise-constant case: tensors + slice mean-zero

For piecewise-constant $F$ on bins (one set per feature), each $f_u$ becomes a **tensor** $T_u$.
**Integrate-to-zero** $\rightarrow$ **weighted slice mean-zero**:

$$\sum_k T_u[\ldots, k, \ldots] \cdot w_i[k] = 0 \quad \text{for each slice along dimension } i \in u$$

## Pairwise intuition ($f_{ij}$ is a matrix)

- Each **row** has weighted mean zero: $\sum_k T_{ij}[r, k] \cdot w_j[k] = 0$
- Each **column** has weighted mean zero: $\sum_k T_{ij}[k, c] \cdot w_i[k] = 0$

**Takeaway:** This constraint enables an exact post-hoc algorithm for tree ensembles.

# Slice Mean-Zero: Visual Example

**Before purification**

| | | | |
|---|---|---|---|
| 4 | 2 | 3 | $\longrightarrow \bar{r}_3 = 3$ |
| 1 | 2 | 3 | $\longrightarrow \bar{r}_2 = 2$ |
| 1 | 5 | 3 | $\longrightarrow \bar{r}_1 = 3$ |

$X_2$ (vertical axis label)

$X_1$

Row means $\neq 0$

**After purification**

| | | | |
|---|---|---|---|
| 1 | $-1$ | 0 | $\longrightarrow \bar{r}_3 = 0$ |
| $-1$ | 0 | 1 | $\longrightarrow \bar{r}_2 = 0$ |
| $-2$ | 2 | 0 | $\longrightarrow \bar{r}_1 = 0$ |

$X_2$ (vertical axis label)

$X_1$

Row means $= 0$ ✓

### What happens to the removed mass?

Row means are subtracted from the interaction tensor $\rightarrow$ added to the main effect $f_2(X_2)$.
Column means are subtracted $\rightarrow$ added to $f_1(X_1)$. Predictions unchanged!

# 7) "Pure interaction effects" = fANOVA of $\mathbb{E}[Y \mid X]$

**Pure interactions** = variance that *cannot* be explained by any subset of variables.

## Formal definition

Find $\{f_u\}$ that minimize $\mathbb{E}_w\big[(F(X) - \sum_u f_u(X_u))^2\big]$ subject to:

$$\mathbb{E}_w[f_u(X_u) \mid X_v] = 0 \quad \forall v \subsetneq u$$

## Key insight

This is exactly the fANOVA decomposition with $w(X) = p(X)$ (true data density).

**Takeaway:** Pure = irreducible. If $f_{ij}$ can be absorbed into $f_i$ or $f_j$, it's not a "pure" interaction.

# 8) Purification for trees: mass-moving

For tree ensembles, represent effects as tensors $T_u$. **Idea**: iteratively move weighted slice means from higher-order to lower-order tensors.
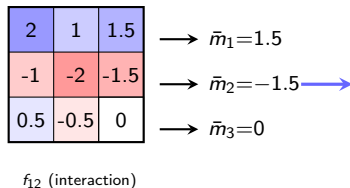
## One step of mass-moving

1. Compute weighted mean of slice: $\bar{m} = \sum_k T_u[\ldots, k, \ldots] \cdot w_i[k]$
2. Subtract from interaction: $T_u[\ldots, k, \ldots] \leftarrow T_u[\ldots, k, \ldots] - \bar{m}$
3. Add to lower-order tensor: $T_{u \setminus i}[\ldots] \leftarrow T_{u \setminus i}[\ldots] + \bar{m}$
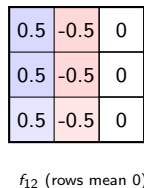
**Algorithm**:

- **Purify-Matrix**: cycle through dimensions until all slice means $\approx 0$.
- **Purify-All**: process tensors from highest-order down to intercept.
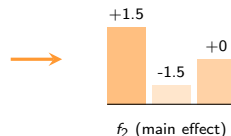
# Mass-moving (Purify-Matrix): Visual Walkthrough

**Step 1: Compute row means**

| 2 | 1 | 1.5 |
|---|---|-----|
| -1 | -2 | -1.5 |
| 0.5 | -0.5 | 0 |

$\longrightarrow \bar{m}_1 = 1.5$
$\longrightarrow \bar{m}_2 = -1.5$
$\longrightarrow \bar{m}_3 = 0$

$f_{12}$ (interaction)

**Step 2: Subtract from rows**

| 0.5 | -0.5 | 0 |
|-----|------|---|
| 0.5 | -0.5 | 0 |
| 0.5 | -0.5 | 0 |

$f_{12}$ (rows mean 0)

**Step 3: Add to main**



+1.5    -1.5    +0

$f_2$ (main effect)

---

## Key insight

Subtracting $\bar{m}_i$ from row $i$ of $f_{12}$ and adding it to $f_2[i]$ preserves predictions:
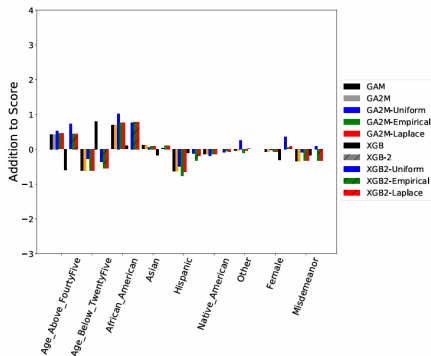
$$f_2(x_2) + f_{12}(x_1, x_2) = \text{unchanged}$$
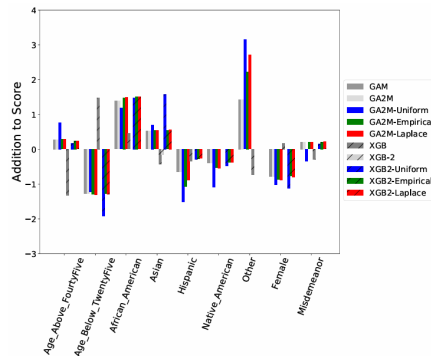
## Result

After convergence: every row/column of interaction has weighted mean = 0

**Takeaway:** Mass moves from interaction to mains; predictions stay identical; interaction becomes "pure" (only irreducible structure remains).

# Example: COMPAS



(a) Prediction of Recidivism

(b) Prediction of COMPAS Score

Figure 4: Main effects of additive models with interactions trained to predict the (a) ground-truth recidivism and (b) COMPAS risk score. The implications of the main effects depend on the model class, the use of purification, and the distribution used for purification.
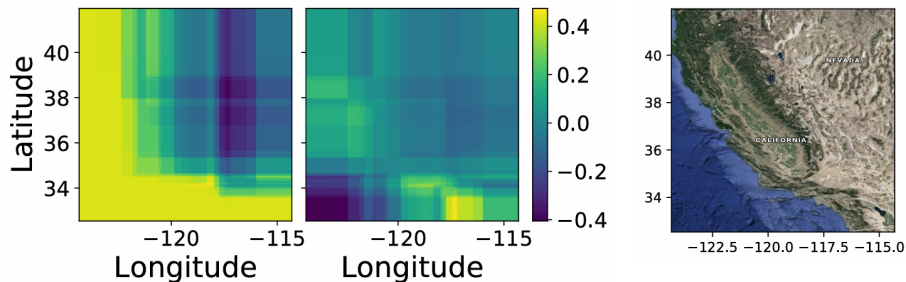
# Example: California housing



Figure 5: Interaction of the Latitude/Longitude features in an XGB2 model trained on the California housing data. The left pane is the unpurified interaction, the middle pane is the purified interaction, and the right is the map of California from which samples were drawn. Purification sorts out the influence from the Los Angeles and the San Francisco metro areas.

# 10) Convergence + correctness

The mass-moving procedure converges to tensors satisfying slice mean-zero constraints. In practice, a small number of sweeps per interaction tensor is enough (see next slide).
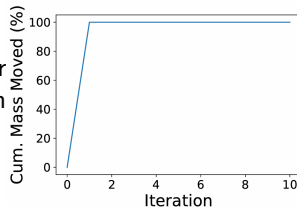
## Theorem 1 (special distributions)

For many simple weights (e.g., uniform along row/column dimensions), purification converges in a single pass.
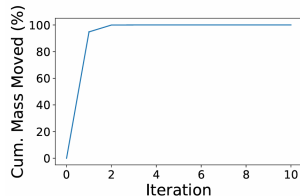
## Theorem 2 (generic non-degenerate $w$)

Converges to tolerance $\varepsilon$ in $O(\log(M_0) - \log(\varepsilon))$ iterations per interaction tensor.

# Empirical convergence: mass moves quickly

- Generate random tensors $T \sim \mathcal{N}(0, \sigma I)$.

- Use weights either: (i) uniform ($w \propto 1$) or (ii) Gaussian ($w \sim \mathcal{N}(0, \sigma I)$) in dimension $P$.

- In practice, almost all mass is moved in the first iteration.

- With uniform weights, convergence occurs in a single pass (per row/column).



(a) Uniform density, $\sigma = 1$, $P = 100$



(b) Density drawn from multivariate normal $\sigma = 10$, $P = 100$

**Takeaway:** Most of the mass moves on the first sweep; uniform weights converge in one pass, enabling purification at scale.

# 11) Useful properties

Because the purified decomposition is the unique fANOVA form:

- **Permutation invariance:** reordering categorical codes does not change purified interactions.
- **Linearity:** purification commutes with averaging / bagging (purify ensembles without changing results).
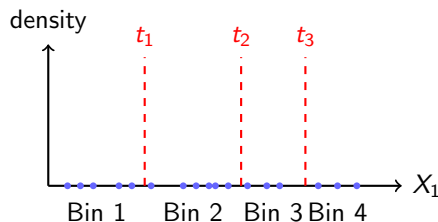
# 12) Estimating $w$ is part of interpretation

Effects are only meaningful *relative to* a distribution $w(X)$. The correct target is the true data density $p(X)$, but it must be estimated.

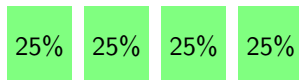## Three practical estimators (piecewise-constant)

- **Uniform:** $\hat{w}_{\text{unif}}(x_{-u}) \propto 1$
- **Empirical:** $\hat{w}_{\text{emp}}(x_{-u}) \propto \sum_{x \in X_{\text{train}}} \mathbf{1}\{x_{-u} = x'_{-u}\}$
- **Laplace:** $\hat{w}_{\text{lap}} \propto \hat{w}_{\text{unif}} + \hat{w}_{\text{emp}}$ (avoids zero-count bins; stabilizes slice means when support is sparse)
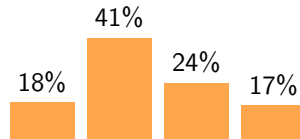
# Binning & Weight Distributions
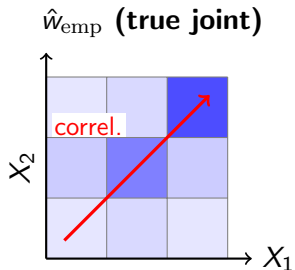
## Tree splits → bins



## Uniform weights

| 25% | 25% | 25% | 25% |

## Weight options

| **uniform** | All bins equal |
| **empirical** | Counts from data |
| **laplace** | empirical + smooth |
| **indep** | $w_{12}=w_1 \cdot w_2$ |

## Empirical weights

18%   41%   24%   17%

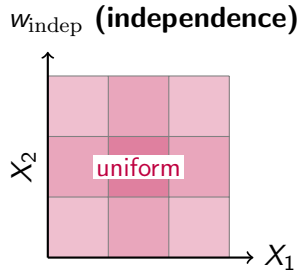**Takeaway:** Bins derived from tree splits. Weight choice affects the "slice mean zero" constraint.

# 2D Joint Weights: Correlated Features

$\hat{w}_{\mathrm{emp}}$ **(true joint)**



Mass concentrated on diagonal

$w_{\mathrm{indep}}$ **(independence)**



$w_{\mathrm{indep}} = w(x_1) \cdot w(x_2)$

## Key difference

$\hat{w}_{\mathrm{emp}}$ captures where data lives (diagonal). $w_{\mathrm{indep}}$ spreads mass uniformly, ignoring correlation.

# Laplace Smoothing: Handling Empty Bins

**Problem:** Some bins may have zero training samples.

- Weighted mean becomes undefined (0/0)
- Rare regions get ignored entirely
- Sensitive to sampling noise
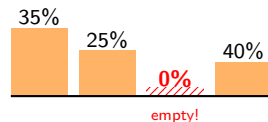
**Solution: Laplace (add-one) smoothing**

$$\hat{w}_{\mathrm{lap}} = \hat{w}_{\mathrm{emp}} + \hat{w}_{\mathrm{unif}}$$

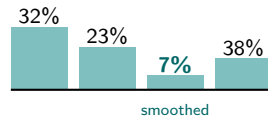Equivalently: add a "pseudo-count" of 1 to each bin before computing weights.

**Effect:**

- Empty bins get small but non-zero weight
- Popular bins still dominate

**Before: Empirical only**

35%  25%  **0%**  40%
empty!

**After: Laplace smoothed**

32%  23%  **7%**  38%
smoothed

$$\hat{w}_{\mathrm{lap}} = \frac{\mathrm{count}+1}{\mathrm{total}+\#\mathrm{bins}}$$

**Part I: Paper Walkthrough**

**Part II: Experiments**

**Conclusions**

# Our additions: two experiments

## Experiment 1 (Real data): German Credit

- Fit depth-2 XGBoost (piecewise-constant $\Rightarrow$ tensorized effects)
- Purify same tensors under $\hat{w}_{\mathrm{emp}}$, $\hat{w}_{\mathrm{unif}}$, and $w_{\mathrm{indep}}$
- Hypothesis: ignoring dependence in $w$ reallocates mass between mains and interactions.

## Experiment 2 (Synthetic): stress-test under strong dependence

- Control dependence with $\rho \in \{0, 0.3, 0.6, 0.9, 0.97\}$
- Compare purification under $\hat{w}_{\mathrm{emp}}$ vs misspecified $w$ ($\hat{w}_{\mathrm{unif}}/w_{\mathrm{indep}}/w_{\mathrm{perm}}$)
- Track divergence metrics and visualize surfaces
- Hypothesis: sensitivity to misspecified $w$ increases with $\rho$.

# Experiment 1: German Credit (setup)

## ML hygiene

- Target: `Credit amount` (regression). Split: 75/25 train/test.
- Preprocessing: fill missing with `unknown`, one-hot encode categoricals.

## Model + purification

- Model: depth-2 XGBoost (n_estimators=120, lr=0.1).
- Binning: thresholds from tree splits (fixed across $w$).
- Purification: compare $\hat{w}_{\mathrm{emp}}$, $\hat{w}_{\mathrm{unif}}$, and $w_{\mathrm{indep}}$.

## What to look for

Does the qualitative explanation (sign, monotonicity, quadrant pattern) change when dependence is ignored in $w$?

# German Credit: Summary Statistics

## Prediction invariance confirmed

RMSE = 2043.3 across all weight modes ($\hat{w}_{\mathrm{emp}}$, $\hat{w}_{\mathrm{unif}}$, $w_{\mathrm{indep}}$) — predictions are **identical**.

### Top main effects (by variance)

| Feature | Var |
|---|---|
| Duration | 2,906,619 |
| Job_3 | 247,365 |
| Purpose_vacation | 93,542 |
| Age | 79,656 |
| Purpose_radio/TV | 43,529 |

### Top interactions (by variance)

| Pair | Var |
|---|---|
| Duration $\times$ Purpose_vac | 29,522 |
| Age $\times$ Duration | 17,139 |
| Duration $\times$ Job_3 | 7,498 |
| Saving_mod $\times$ Purp_edu | 5,888 |

**Takeaway: Duration** dominates mains; its effect depends on loan purpose (interaction). Different $w$ reallocates variance but predictions stay identical.

# German Credit: main effects (top features)
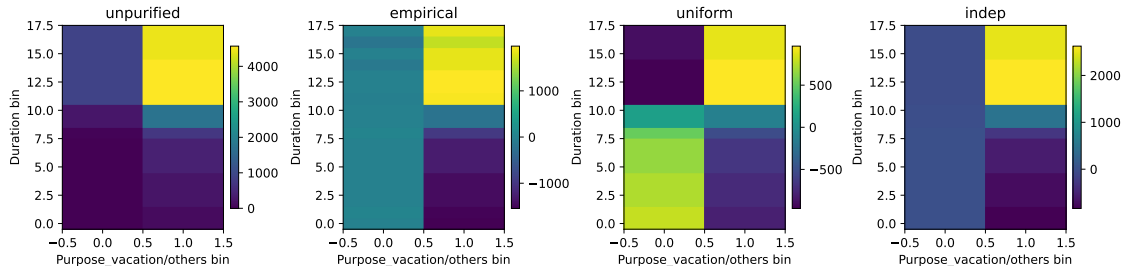


**Takeaway:** Same model and binning; differences across $w$ are reallocations between mains and interactions.

German credit: top interaction Duration × Purpose_vacation/others

**Takeaway:** Interaction residual after purification. If $w_{\mathrm{indep}}$ changes the qualitative pattern relative to $\hat{w}_{\mathrm{emp}}$, the interaction story is sensitive to ignoring dependence.

**Part I: Paper Walkthrough**

1. The identifiability problem
2. Functional ANOVA target
3. Mass-moving algorithm
4. Convergence & properties
5. Weight estimation ($w$)

**Part II: Experiments**

6. German Credit (real data)
   Main effects & interactions under different $w$

▶7. **Stress-test (synthetic)**
   How dependence $\rho$ affects purification

**Conclusions**

   Summary & takeaways

# Experiment 2: core question (weighting distribution $w$)

## Question

Under strong feature dependence, how much do purified mains/interactions change when we purify using:

- the **correct joint** (empirical $\hat{w}_{\mathrm{emp}} \approx p(x)$) vs
- a **misspecified** weighting (uniform or independence-assumed)?

- Same learned piecewise-constant model $\Rightarrow$ same tensors.
- Only $w$ changes.

# Stress-test: data generation

## Linear-Gaussian dependence

$$X_1 \sim \mathcal{N}(0, 1), \quad \varepsilon \sim \mathcal{N}(0, 1) \text{ independent}$$

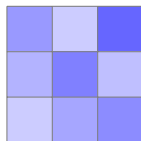$$X_2 = \rho X_1 + \sqrt{1 - \rho^2}\, \varepsilon, \quad \rho \in \{0.0, 0.3, 0.6, 0.9, 0.97\}$$

## Ground-truth functions

- (F1) $F(X_1, X_2) = \beta_1 X_1 + \beta_2 X_2 + \gamma(X_1 X_2)$
- (F2) $F(X_1, X_2) = \gamma \, \text{sign}(X_1)\text{sign}(X_2) + \beta_1 X_1 + \beta_2 X_2$
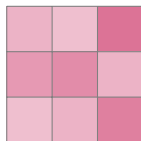
# Evaluation Metrics: Visual Guide

## Weighted $L_2$ Distance

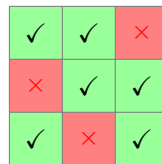$f_{12}^{\text{emp}}$ (reference)      $f_{12}^{\text{other}}$



$$d_{L_2} = \sqrt{\sum_{i,j} w_{ij} \left( f_{ij}^{\text{emp}} - f_{ij}^{\text{other}} \right)^2}$$

Measures overall surface difference

## Sign Flip Rate

**Sign($f^{\text{emp}}$) vs Sign($f^{\text{other}}$)**



$$\text{flip rate} = \frac{\#\{\text{sign differs}\}}{\#\{\text{bins}\}} = \frac{3}{9}$$

Measures qualitative disagreement

**Takeaway:** $L_2$ captures magnitude differences; sign flips capture qualitative interpretation changes (e.g., "positive effect" $\rightarrow$ "negative effect").
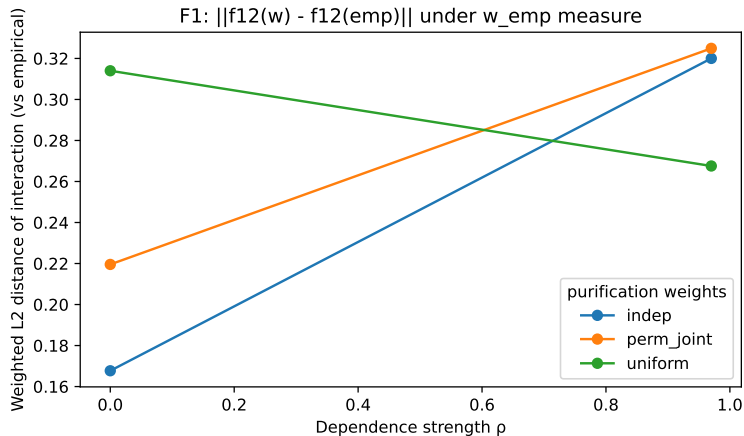
# Stress-test: weighting ablation + metrics

We purify the same tensors under:

$$w \in \{\hat{w}_{\mathrm{emp}}, \ \hat{w}_{\mathrm{unif}}, \ w_{\mathrm{indep}}, \ w_{\mathrm{perm}}\}.$$
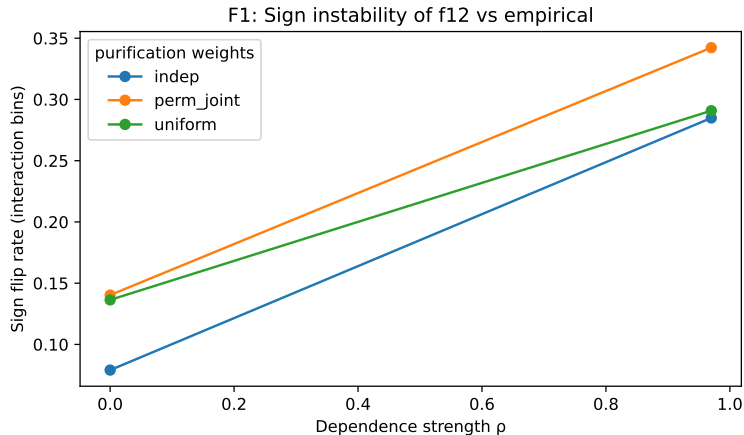
### Metrics vs $\rho$

- Weighted $L_2$ distance between purified interaction surfaces (reference measure: $\hat{w}_{\mathrm{emp}}$)
- Sign flip rate vs $\hat{w}_{\mathrm{emp}}$ (fraction of bins with sign changes)
- Interaction variance share: $\mathrm{Var}_w(f_{12})/\mathrm{Var}_w(F)$

F1: ||f12(w) - f12(emp)|| under w_emp measure

**Takeaway:** Distance from the $\hat{w}_{\mathrm{emp}}$-purified interaction increases with dependence; misspecified $w$ yields divergent explanations.

F1: Sign instability of f12 vs empirical

**Takeaway:** As dependence grows, qualitative interaction conclusions (sign) become sensitive to the choice of $w$.

F1: Interaction variance share depends on w

**Takeaway:** Variance attribution to the interaction changes with $w$; under dependence, misspecifying $w$ reshuffles main vs interaction credit.

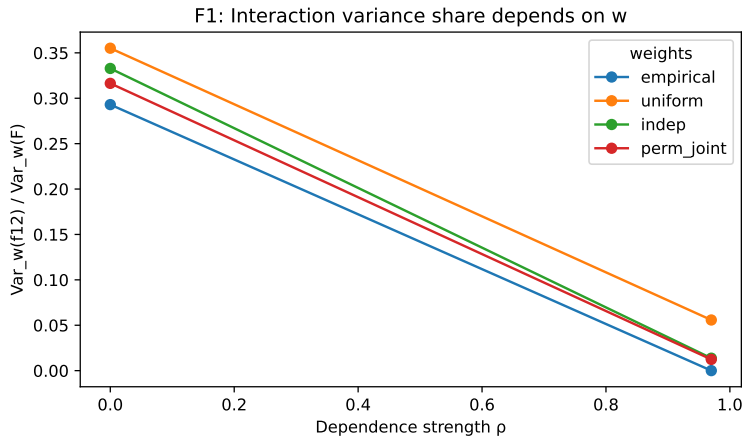F2: ||f12(w) - f12(emp)|| under w_emp measure

**Takeaway:** Distance from the $\hat{w}_{\mathrm{emp}}$-purified interaction increases with dependence; misspecified $w$ yields divergent explanations.

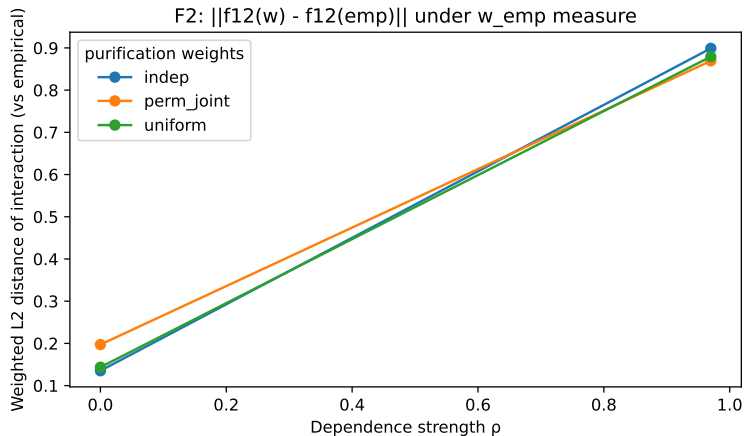F2: Sign instability of f12 vs empirical

**Takeaway:** As dependence grows, qualitative interaction conclusions (sign) become sensitive to the choice of $w$.
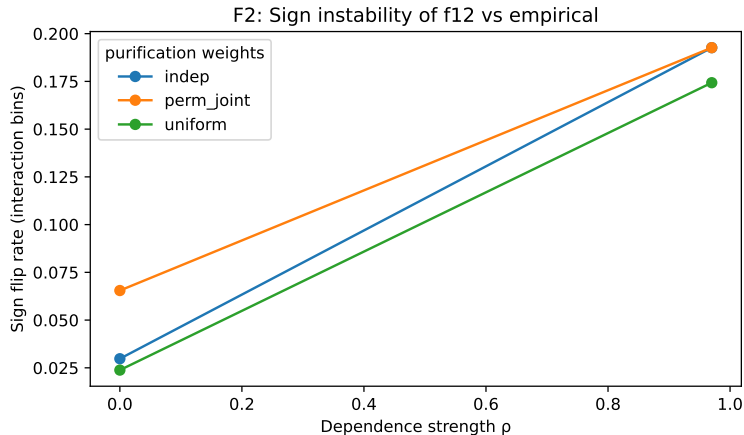
# Stress-test (F2): interaction variance share vs $\rho$



F2: Interaction variance share depends on w

**Takeaway:** Variance attribution to the interaction changes with $w$; under dependence, misspecifying $w$ reshuffles main vs interaction credit.
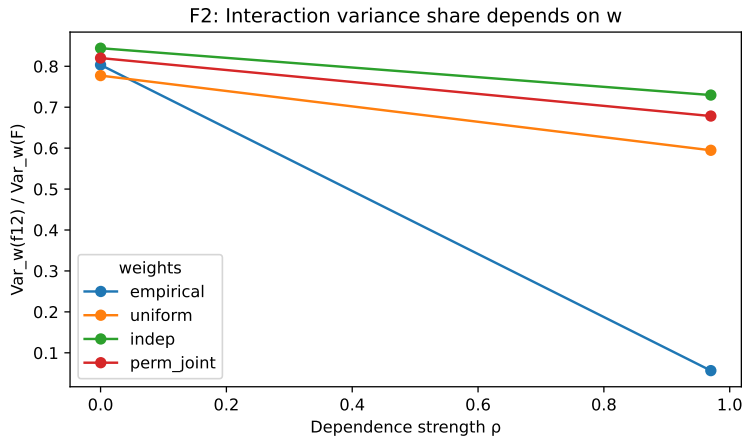
# Stress-test: Summary — $L_2$ Distance & Sign Flips

## $L_2$ distance (interaction surface vs $\hat{w}_{\mathrm{emp}}$ reference)

| Mode | F1 (multiplicative) | | F2 (XOR-like) | |
|---|---|---|---|---|
| | $\rho=0$ | $\rho=0.97$ | $\rho=0$ | $\rho=0.97$ |
| uniform | 0.31 | 0.27 | 0.14 | **0.88** |
| indep | 0.17 | 0.32 | 0.14 | **0.90** |
| perm_joint | 0.22 | 0.32 | 0.20 | **0.87** |

## Sign flip rate at $\rho = 0.97$

| Mode | F1 | F2 |
|---|---|---|
| uniform | 29% | 17% |
| indep | 28% | 19% |

**Takeaway:** At high $\rho$, F2 shows $L_2 \approx 0.9$ (surface almost unrecognizable). 20–30% of bins flip sign!

# Stress-test: Summary — Variance Attribution

## Interaction variance share: $\text{Var}_w(f_{12})/\text{Var}_w(F)$

| Mode | F1 | | F2 | |
|---|---|---|---|---|
| | $\rho=0$ | $\rho=0.97$ | $\rho=0$ | $\rho=0.97$ |
| $\hat{w}_{\text{emp}}$ | 29% | **0.002%** | 80% | **5.6%** |
| uniform | 36% | 5.6% | 78% | 59% |
| indep | 33% | 1.4% | 84% | **73%** |

## Key insight

For F1 at $\rho=0.97$: $\hat{w}_{\text{emp}}$ shows **0% interaction** (absorbed into mains).
But $w_{\text{indep}}$ reports 1.4%, $\hat{w}_{\text{unif}}$ reports 5.6% — **qualitatively different stories!**

**Takeaway:** Variance attribution shifts dramatically with $w$; misspecifying $w$ yields misleading importance rankings.
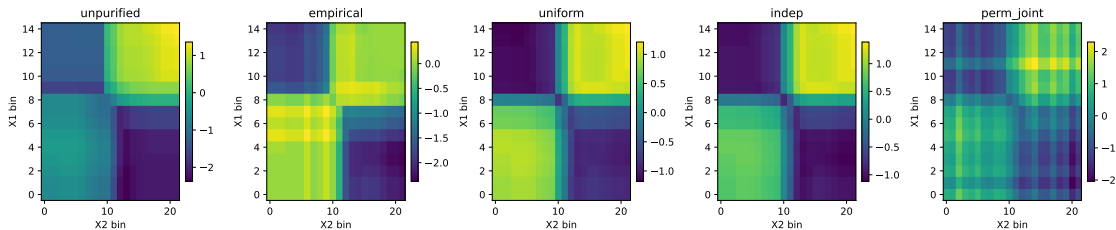
F1: interaction surface at ρ=0.97 under different w

**Takeaway:** At strong dependence ($\rho = 0.97$), the interaction surface can change noticeably across $w$; $\hat{w}_{\mathrm{emp}}$ emphasizes where samples occur.
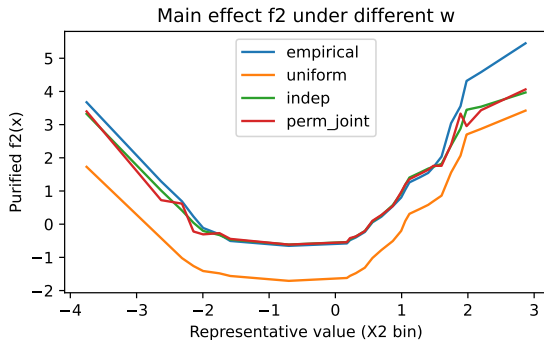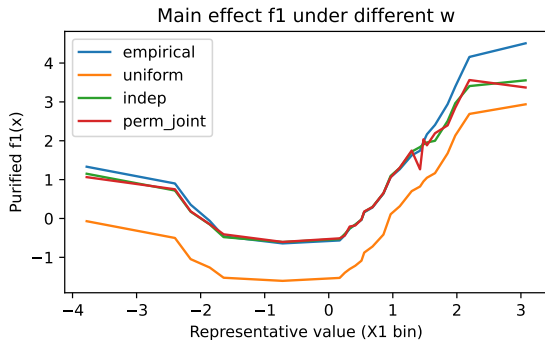
F2: interaction surface at ρ=0.97 under different w

**Takeaway:** At strong dependence ($\rho = 0.97$), the interaction story depends on $w$; misspecification can change qualitative patterns.

F1: main effects at ρ=0.97



**Takeaway:** Main-effect curves can shift across *w* even when the model is fixed; changes are compensating reallocations with interactions (prediction stays the same).

# Conclusions

## Theoretical contribution (paper)

- Interactions are **not identifiable** without constraints; fANOVA yields a canonical decomposition under $w$.
- Purification: exact post-hoc mass-moving recovers fANOVA for piecewise-constant models.

## Our experimental findings

- **Predictions are invariant** to weight choice (confirmed: identical RMSE across $w$).
- **Interpretations are NOT invariant**: under strong dependence ($\rho = 0.97$):
  - Variance attribution shifts dramatically (0% vs 73% interaction share)
  - 20–30% of bins flip sign
  - $L_2$ distance up to 0.9 (interaction surface almost unrecognizable)

**Takeaway:** Always report which $w$ you used. Treat $\hat{w}_{\mathrm{emp}}$ vs $w_{\mathrm{indep}}$ as a **sensitivity analysis** when feature dependence is plausible.

## References & Resources

### Original Paper

Lengerich, B., Tan, S., Chang, C.-H., Hooker, G., & Caruana, R. (2020).

*Purifying Interaction Effects with the Functional ANOVA: An Efficient Algorithm for Recovering Identifiable Additive Models.*

AISTATS 2020, PMLR 108:2402–2412.

`https://proceedings.mlr.press/v108/lengerich20a.html`

### Code & Related Projects

- Original paper implementation: `https://github.com/blengerich/gam_purification`
- InterpretML (Microsoft): `https://github.com/interpretml/interpret`

### Dataset

German Credit Data — UCI Machine Learning Repository

`https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data)`

# AI Usage Disclosure

## AI Tools Used

The following AI tools were used to support development of this project:

- **NotebookLM** (Google) — understanding the paper
- **ChatGPT 5.2 Pro** (OpenAI) — code generation, debugging, and slide drafting
- **Claude Opus / Sonnet 4.5** (Anthropic) — code review, documentation, and slide refinement

All AI-generated content was reviewed, validated, and edited by the author.