

Privacy-preserving Meta-Analysis through Low-Rank Basis Hunting

Kosuke Imai

Harvard University

Applied Statistics Seminar

Indian Statistical Institute, Kolkata

March 24, 2026

Joint work with Wenqi Shi (Harvard) and Yi Zhang (Netflix)

Motivation

- **Meta analysis**: statistical analysis to combine the results of multiple independent studies
 - widely used in social and medical sciences for evidence synthesis
 - generalizing results from *sources* to a *target* population
 - common parametric approaches: random effect models
- Key challenges:
 - ① **unknown heterogeneity** across sources and between source and target populations
 - covariate shift $\mathcal{P}(\mathbf{X})$
 - conditional shift $\mathcal{P}(Y | \mathbf{X})$
 - ② **function-valued quantities** of interest, going beyond vector-valued parameters
 - regression functions $\mathbb{E}[Y | \mathbf{X}] \iff$ means $\mathbb{E}[Y]$
 - conditional average treatment effect (CATE) $\mathbb{E}[Y(1) - Y(0) | \mathbf{X}] \iff$ ATE $\mathbb{E}[Y(1) - Y(0)]$
 - ③ **Privacy preservation**: limited direct access to source data
 - ④ **Use of machine learning (ML) models**: flexible nonparametric estimation for each source

Overview of Our Contributions

- **MetaHunt**: privacy-preserving functional meta-analysis
 - estimates function-valued quantities for a *new target* population
 - requires only aggregate information from sources
 - allows for the use of (possibly different and unknown) ML models in each source
 - provides asymptotically valid (pointwise) statistical inference
- Key idea:
 - source and target populations share a **common low-rank structure**
 - all study-level functions lie in the convex hull of a small number of latent basis functions
- **functional Successive Projection Algorithm** (fSPA) to recover latent basis functions
- **Conformal inference** to provide a confidence interval with marginal coverage control

Low-rank Cross-Study Heterogeneity

- Study-specific functions of interest $f^{(i)}(\mathbf{x})$ may differ in complex and unknown ways
- But, we assume that they share a common underlying structure
- All study-specific functions lie within the convex hull of a set of basis functions $\{g_k(\mathbf{x})\}_{k=1}^K$
- Enable dimension reduction, nonparametric modeling, and generalization

Assumption 1 (Low-rank cross-study heterogeneity)

There exists a set of basis functions $\{g_k(\mathbf{x})\}_{k=1}^K$ ($K < m$) such that for all $i = 0, 1, \dots, m$

$$f^{(i)}(\mathbf{x}) = \sum_{k=1}^K \pi_{ik} g_k(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X},$$

where $\boldsymbol{\pi}_i = (\pi_{i1}, \dots, \pi_{iK})^\top \in \Delta_{K-1} = \{\boldsymbol{\pi} \in \mathbb{R}^K : \sum_{k=1}^K \pi_k = 1, \pi_k \geq 0\}$

Weight Model

- The weights π_i determine how study i combines the shared basis functions $\{g_k\}_{k=1}^K$
- Use study-level information \mathbf{W}_i to predict variation in π_i
- Studies with similar \mathbf{W}_i will have similar mixing of the basis functions

Assumption 2 (Weight model)

For all $i \in \{0, 1, \dots, m\}$, the weight vector π_i is drawn independently from a conditional distribution given $\mathbf{W}_i \in \mathcal{W}$

$$\pi_i \mid \mathbf{W}_i \stackrel{\text{ind.}}{\sim} \mathcal{P}_{\pi \mid \mathcal{W}}(\cdot \mid \mathbf{W}_i),$$

where $\mathcal{P}_{\pi \mid \mathcal{W}}$ is an arbitrary distributional map from \mathcal{W} to the simplex Δ_{K-1}

Weight Model Example: Dirichlet Regression

Example 1 (Kernelized Dirichlet regression model)

Let $\kappa : \mathcal{W} \times \mathcal{W} \rightarrow \mathbb{R}$ be a positive-definite kernel with the associated reproducing kernel Hilbert space (RKHS) \mathcal{H}_κ and feature map $\psi : \mathcal{W} \rightarrow \mathcal{H}_\kappa$. For all $i \in \{0, 1, \dots, m\}$, assume

$$\boldsymbol{\pi}_i := (\pi_{i1}, \pi_{i2}, \dots, \pi_{iK}) \stackrel{ind.}{\sim} \text{Dirichlet}(\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{iK}),$$

with $\alpha_{ik} = \exp \{ \langle \boldsymbol{\beta}_k, \psi(\mathbf{w}_i) \rangle_{\mathcal{H}_\kappa} \}$ and $\boldsymbol{\beta}_k \in \mathcal{H}_\kappa$

- Alternative modeling approaches:
 - Log-ratio regression
 - Neural networks with softmax

Exchangeability of Study-level Covariates

Assumption 3 (Study-level covariate exchangeability)

The study-level covariates $\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_m$ are exchangeable

- Exchangeability required for conformal prediction
- No exchangeability assumption at the individual level
- Under Assumptions 1–3, $(\mathbf{W}_i, \pi_i, f^{(i)})$ are jointly exchangeable across $i = 0, 1, \dots, m$

Overview of MetaHunt

Estimation stage

Input: $\{\mathbf{W}_i, \hat{f}^{(i)}(\mathbf{x})\}_{i=1}^m$

Output: $\{\hat{g}_k\}_{k=1}^K, \widehat{\mathcal{M}} : \mathcal{W} \rightarrow \Delta_{K-1}$

- 1 Basis hunting: obtain $\{\hat{g}_k\}_{k=1}^K$
- 2 Estimate weights $\hat{\pi}_i$ by projection
- 3 Fit the weight model $\widehat{\mathcal{M}}$

Prediction stage

Input: $\{\hat{g}_k\}_{k=1}^K, \widehat{\mathcal{M}}, \mathbf{W}_0$

Output: $\tilde{f}^{(0)}(\mathbf{x})$ and a confidence interval

- 1 Predict the weights: $\tilde{\pi}_0 = \widehat{\mathcal{M}}(\mathbf{W}_0)$
- 2 Predict the target function:

$$\tilde{f}^{(0)}(\mathbf{x}) = \sum_{k=1}^K \tilde{\pi}_{0k} \hat{g}_k(\mathbf{x})$$

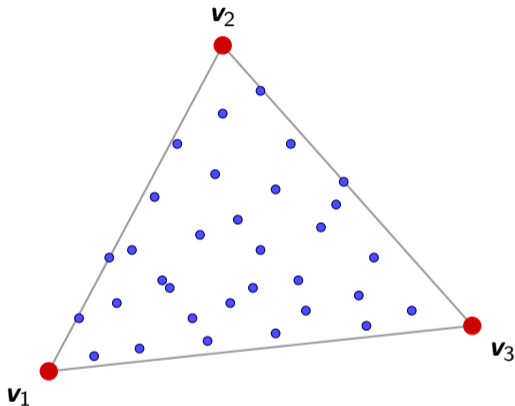
- 3 Construct a conformal prediction interval

Basis Hunting vs. Vertex Hunting

- **Vertex hunting:**

Find the vertices \mathbf{v}_k such that for all \mathbf{x}_i we have

$$\mathbf{x}_i = \sum_{k=1}^K \pi_{ik} \mathbf{v}_k, \quad \sum_{k=1}^K \pi_{ik} = 1, \quad \pi_{ik} \geq 0$$

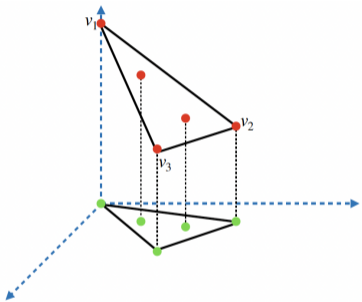


- **Basis hunting:**

Recover the latent basis functions $\{\mathbf{g}_k\}_{k=1}^K$ from the observed functions $\{\hat{f}^{(i)}\}_{i=1}^m$

- Random functions vs random vectors

Successive Projection Algorithm (SPA; Arajuño et al. 2001) for Vertex Hunting



- 1 Select the data point with the largest length as the first vertex v_1
- 2 Project all data points onto the space orthogonal to the selected vertex
- 3 Among the projected points, choose the one with the largest length as the next vertex v_2
- 4 Repeat the projection and selection steps until K vertices are chosen

- **Denoising step** (Jin et al. 2024): For each (projected) point,
 - 1 If there are fewer than N points in the neighborhood of x_i , remove x_i
 - 2 Otherwise, replace x_i with the average of N neighborhood points

d-fSPA: Denoised Basis Hunting Algorithm

Algorithm 1: d-fSPA

Input: Estimated functions $\{\hat{f}^{(i)}(\cdot)\}_{i=1}^m$; number of basis K ; tuning parameters (N, Δ)

Denoising step:

- 1 If there are fewer than N functions in $B_\Delta(\hat{f}^{(i)})$, remove this function
- 2 Otherwise, replace $\hat{f}^{(i)}$ by the average of all functions in $B_\Delta(\hat{f}^{(i)})$

where $B_\Delta(\hat{f}^{(i)}) = \{f : \|f - \hat{f}^{(i)}\| \leq \Delta\}$ denotes the Δ -neighborhood of $\hat{f}^{(i)}$

Initialize $S = \emptyset$, $h_i(\cdot) = \hat{f}^{(i)}(\cdot)$ for $1 \leq i \leq m$

for $k = 1, \dots, K$ **do**

- 1 Project $h_i = \hat{f}^{(i)}$ onto the orthogonal space of $\text{span}(S)$: $h_i = h_i - P_{\text{span}(S)} h_i$
- 2 Find s_k such that h_{s_k} has the largest norm $s_k = \text{argmax}_{1 \leq i \leq m} \|h_i\|$
- 3 Update the set $S = S \cup \{\hat{f}^{(s_k)}\}$

Output: Estimated basis functions $\hat{g}_k(\cdot) = \hat{f}^{(s_k)}$ for $1 \leq k \leq K$

d-fSPA Implementation Details

- Functional norm:
 - L^2 norm under the target covariate distribution $\mathcal{P}_{0,\mathbf{X}}$ (use sampled covariates for evaluation)

$$\|f - \tilde{f}\| = \left\{ \mathbb{E}_{\mathcal{P}_{0,\mathbf{X}}} [(f(\mathbf{X}) - \tilde{f}(\mathbf{X}))^2] \right\}^{1/2}$$

- If the target distribution is unavailable, use proxies (e.g., pooled source covariate distribution)
- Denoising parameters: heuristic choice (Jin et al. 2024)
 $N = 0.5 \log m$ and $\Delta = \max_{ij} \|\hat{f}^{(i)} - \hat{f}^{(j)}\|/10$
- Study-specific functions $\hat{f}^{(i)}$ may be estimated using (possibly different and unknown) black-box ML models

Selecting the Number of Basis Functions K

1 Elbow plot:

- Plot K against the reconstruction error

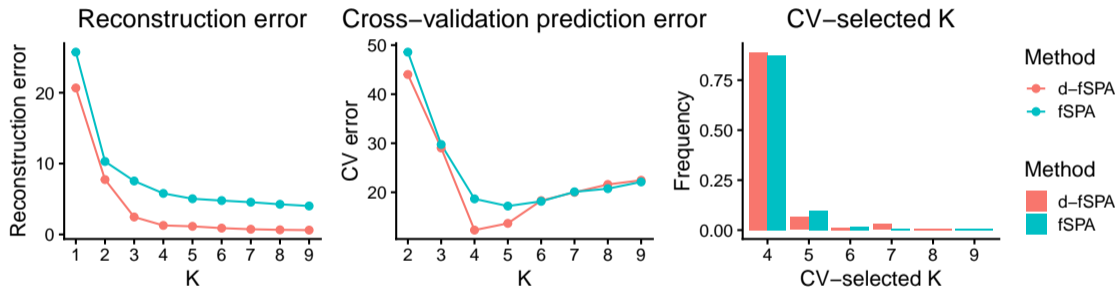
$$\mathcal{E}(K) := \frac{1}{m} \sum_{i=1}^m \min_{\pi_i} \left\| \hat{f}^{(i)} - \sum_{k=1}^K \pi_{ik} \hat{g}_k \right\|$$

- The “elbow” point indicates where adding more bases yields small improvement
- Requires only the estimated functions $\{\hat{f}^{(i)}(\mathbf{x})\}$

2 Cross-validation:

- Select K by minimizing cross-validation prediction error
- Requires both the estimation and prediction stages
- Requires site-level covariates \mathbf{W}_i

Selection of K : Simulation Illustration



Exact Recovery of Basis Functions by fSPA under No Estimation Error

Proposition 1 (Exact recovery by fSPA under no estimation error)

Suppose the basis functions are non-degenerate, Assumption 1 holds, and the following conditions are satisfied:

- 1 No estimation error: $\|f^{(i)} - \hat{f}^{(i)}\| = 0$ for all $i = 1, \dots, m$
- 2 Pure node condition: For each $k = 1, \dots, K$, there exists a study i such that $\pi_{ik} = 1$

Then, up to permutation of $\{1, \dots, K\}$,

$$\max_{1 \leq k \leq K} \|\hat{g}_k - g_k\| = 0$$

Dealing with Estimation Error and Degree of Purity

- Estimation error:

$$\epsilon^{(i)}(\mathbf{x}) := \hat{f}^{(i)}(\mathbf{x}) - f^{(i)}(\mathbf{x})$$

- Degree of purity:

$$\delta_m := \max_{1 \leq k \leq K} \min_{1 \leq i \leq m} (1 - \pi_{ik})$$

Assumption 4 (Error control)

- 1 $\mathbb{E}[\epsilon^{(i)}(\mathbf{x})^2] = O(n_i^{-r})$ for some $r > 0$, for all i and \mathbf{x}
- 2 There exists $0 < a < r$ such that $m = o(\inf_i n_i^a)$
- 3 $\delta_m \max_{1 \leq k \leq K} \|g_k\| = o(1)$

- 1 Parametric models $r = 1$; ML models $r < 1$
- 2 The number of studies m remains moderate relative to study sample sizes
- 3 As m increases, we expect some studies to approach pure nodes $\delta_m \rightarrow 0$

Theorem 2 (Asymptotic recovery by d-fSPA)

Under Assumptions 1 and 4, if the basis functions are non-degenerate, then, with an appropriate choice of the denoising parameters (N, Δ) , we have, up to permutation,

$$\max_{1 \leq k \leq K} \|\hat{g}_k - g_k\| = o_P(1)$$

- Exact conditions for the denoising parameters are currently being finalized...

Fitting the Weight Model

- 1 Estimate study-level weights:

$$\hat{\boldsymbol{\pi}}_i = \operatorname{argmin}_{\boldsymbol{\pi}_i \in \Delta_{K-1}} \left\| \hat{f}^{(i)} - \sum_{k=1}^K \pi_{ik} \hat{\boldsymbol{g}}_k \right\|$$

- 2 Fit the weight model: Use $\{(\mathbf{W}_i, \hat{\boldsymbol{\pi}}_i)\}_{i=1}^m$ to estimate

$$\hat{\mathcal{M}} : \mathcal{W} \rightarrow \Delta_{K-1}$$

Predicting for a Target Site

- 1 Predict mixing weights:

$$\tilde{\boldsymbol{\pi}}_0 = \widehat{\mathcal{M}}(\mathbf{W}_0).$$

- 2 Predict the target function:

$$\tilde{f}^{(0)}(\mathbf{x}) = \sum_{k=1}^K \tilde{\pi}_{0k} \hat{g}_k(\mathbf{x}).$$

Challenges in Constructing Prediction Interval

- 1 Our target $f^{(0)}$ is a function
 - Building uniform confidence bands is hard, i.e., $\Pr(l(\mathbf{x}) \leq f^{(0)}(\mathbf{x}) \leq u(\mathbf{x})) \geq 1 - \alpha$ for all \mathbf{x}
 - We focus on a point-wise interval for $f^{(0)}(\mathbf{x})$ given \mathbf{x}
- 2 Multiple sources of error
 - Possibly correlated errors in basis hunting, weight estimation, weight model
 - We focus on conformal prediction with valid marginal coverage
- 3 Estimation error in $\hat{f}^{(i)}$
 - We aim to predict $f^{(0)}(\mathbf{x})$ but only observe $\hat{f}^{(i)}(\mathbf{x})$ as a proxy
 - Assumptions are needed to control the estimation error $|\hat{f}^{(i)}(\mathbf{x}) - f^{(i)}(\mathbf{x})|$

Conformal Prediction Interval

Algorithm 2: Split conformal prediction for $f^{(0)}(\mathbf{x})$

Input: Estimated functions $\{\hat{f}^{(i)}\}_{i=1}^m$; target-study covariates \mathbf{W}_0 ; evaluation point \mathbf{x} ; miscoverage level $\alpha \in (0, 1)$; estimation pipeline

- 1 **Data splitting:** randomly split the index set $\{1, \dots, m\}$ into \mathcal{I}_{tr} and \mathcal{I}_{cal}
- 2 **Training:** using only training data, run the estimation pipeline to obtain a prediction rule
- 3 **Calibration:** for each $i \in \mathcal{I}_{\text{cal}}$, calculate the conformity score $r_i := |\hat{f}^{(i)}(\mathbf{x}) - \tilde{f}^{(i)}(\mathbf{x})|$
- 4 **Conformal quantile:** set $q_{1-\alpha}(\mathbf{x}) := r_{(\lceil(1-\alpha)(m_{\text{cal}}+1)\rceil)}$
- 5 **Prediction:** compute the point prediction for the new study, $\tilde{f}^{(0)}(\mathbf{x})$

Output: Split conformal prediction interval

$$C_\alpha(\mathbf{x}) = [\tilde{f}^{(0)}(\mathbf{x}) - q_{1-\alpha}(\mathbf{x}), \tilde{f}^{(0)}(\mathbf{x}) + q_{1-\alpha}(\mathbf{x})]$$

Theorem 3 (Asymptotic marginal coverage guarantee)

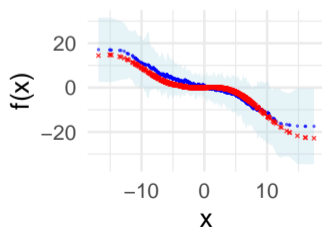
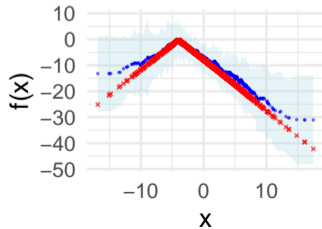
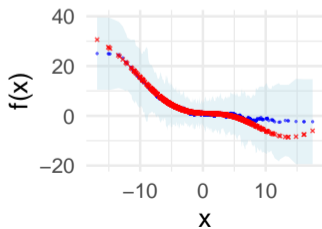
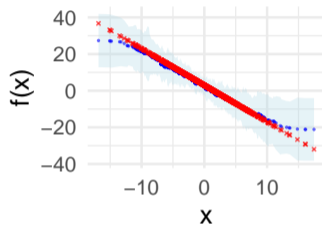
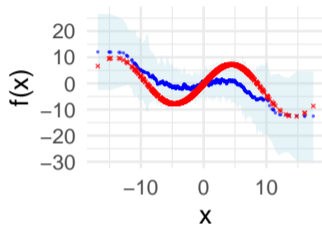
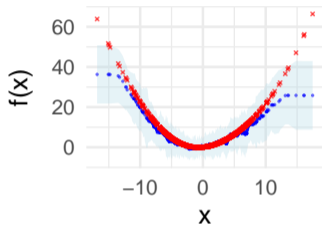
Under Assumptions 1–4, we have asymptotic marginal coverage for any \mathbf{x} ,

$$\lim \Pr(f^{(0)}(\mathbf{x}) \in C_\alpha(\mathbf{x})) \geq 1 - \alpha$$

Conformal Prediction: Simulation Illustration

Prediction with intervals

Blue circle = prediction, shaded region = prediction interval, red cross = expected values



Extension to a Functional of the Target Function

- $\mathcal{H} : \mathcal{F} \rightarrow \mathbb{R}$: a functional mapping a function to scalar, i.e., $C^{(0)} = \mathcal{H}(f^{(0)})$
- Example: ATE $\tau := \mathbb{E}[Y(1) - Y(0)]$
 - CATE may be smoother than the conditional outcome function
 - shared low-rank structure of the CATE $f(\mathbf{x}) = \mathbb{E}[Y(1) - Y(0) \mid \mathbf{X} = \mathbf{x}]$

$$\tau^{(0)} = \mathcal{H}(f^{(0)}) = \mathbb{E}_{\mathcal{P}_{\mathbf{0}, \mathbf{x}}} [f^{(0)}(\mathbf{X})]$$

where $\mathcal{P}_{\mathbf{0}, \mathbf{x}}$ is the covariate distribution of the target population

- Plug-in point prediction:

$$\tilde{\tau}^{(0)} = \mathcal{H}(\tilde{f}^{(0)})$$

- Conformal interval: For each calibration study i , define the conformity score

$$r_i = \left| \mathcal{H}(\hat{f}^{(i)}) - \mathcal{H}(\tilde{f}^{(i)}) \right|$$

Empirical Application: Many Labs 1 (Klein et al. 2014)

- Large-scale multisite replication project in psychology
- Evaluates the replicability of 13 classic experimental findings, including:
 - anchoring (Jacowitz & Kahneman, 1995)
whether exposure to a high or low numerical anchor influences participants' estimates of the population of Chicago
 - gain vs. loss framing effects (Tversky & Kahneman, 1981)
whether participants make different choices when outcomes are framed in terms of lives saved (gains) versus lives lost (losses)
- Experiments across 36 independent data collection sites ($m = 36$)
- Each hypothesis is tested under a common experimental protocol across studies
- Total sample size: approximately 6,000 participants per hypothesis on average

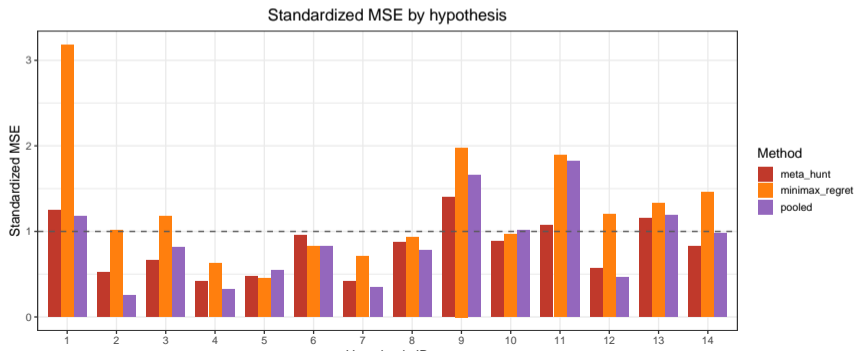
Setup

- Site-level covariates \mathbf{W} :
 - Online vs. laboratory setting
 - US vs. international site
 - Average age
 - Gender ratio
 - Average political ideology
 - Measured covariate shift relative to the target site (varies by target site and hypothesis)
- Individual-level covariates \mathbf{X} : gender, age, race, political ideology, American identity
- **Prediction task**: regression function estimated using random forest
- **Causal task**: site-level CATE functions estimated using Causal Forest
- Selection of K based on cross-validation for each hypothesis
- **Leave-one-site-out prediction** for validation:
 - for each hypothesis, treat one site as the target and use the remaining sites for training
 - repeat so that every site serves as the target
- Empirical benchmark: Target-site ATE estimated directly from the experiment

Prediction Task Results

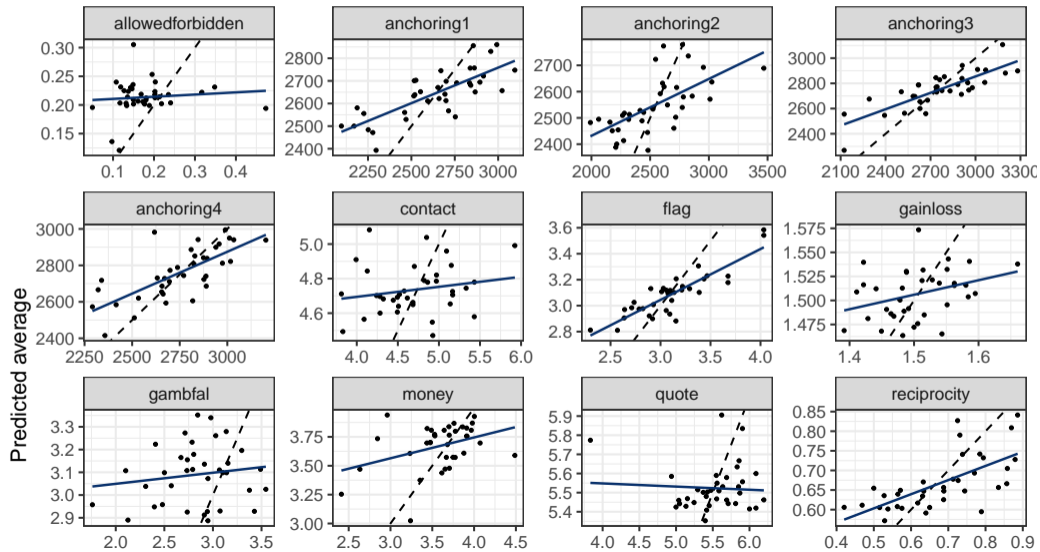
Method	Std. MSE	Coverage	Interval Length (relative)
MetaHunt	0.817	0.988	0.969
Minimax regret (Zhang et al.)	1.266	0.982	0.964
Pooled ML	0.871	0.974	0.806

- Standardized relative to the empirical benchmark
- Pooled ML: ML with conformal inference, using both individual and site-level covariates



Empirical Performance by Hypothesis

Predicted average vs empirical benchmark

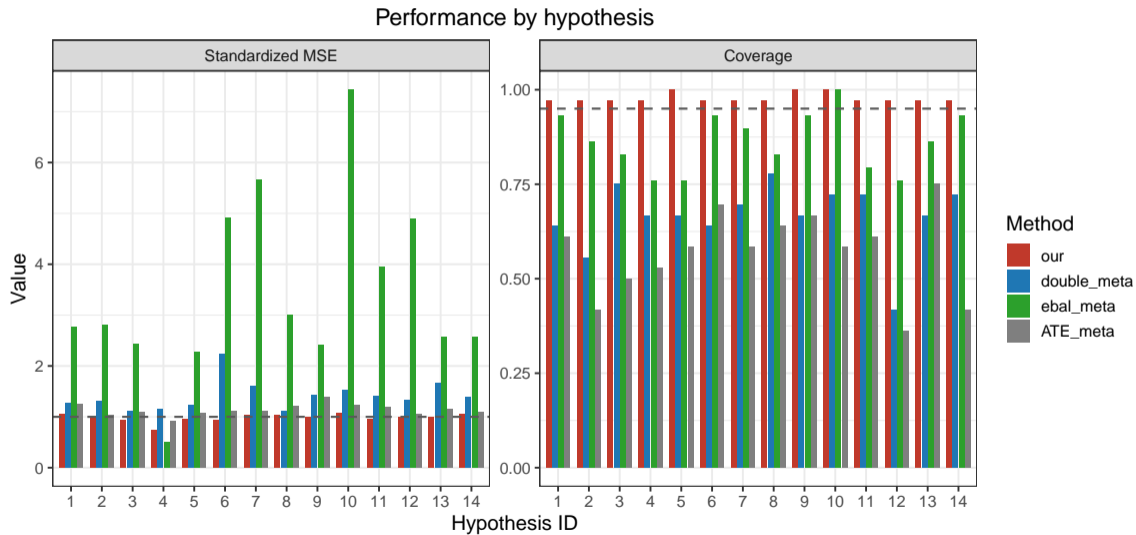


Causal Task Results

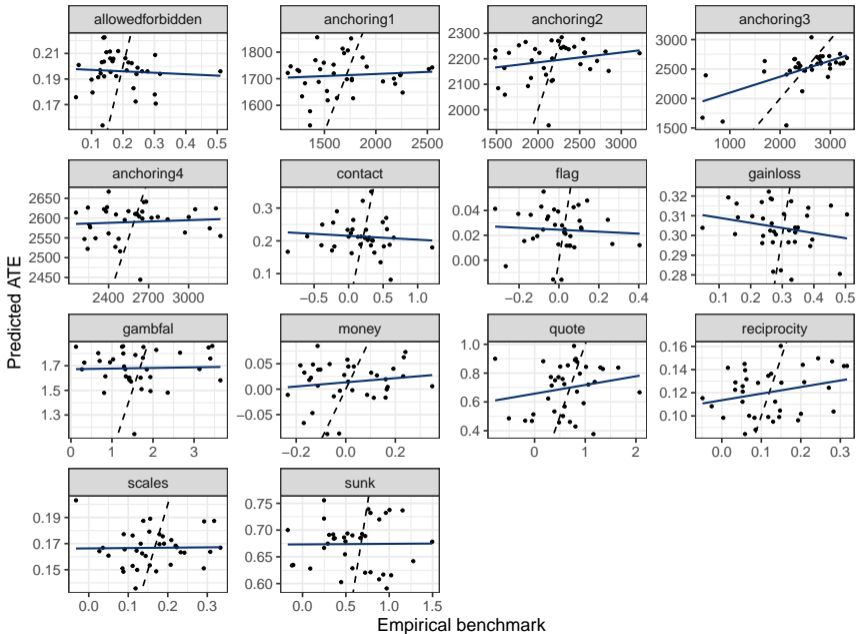
Method	Std. MSE	Coverage	Interval Length (relative)
MetaHunt	0.962	0.976	0.780
DR	1.305	0.700	0.405
ebal	2.532	0.855	0.965
DiffMeans	1.117	0.645	0.316
Pooled ML	0.977	0.958	0.714

- DR: doubly robust estimator at each source site with covariate density ratios
- ebal: entropy balancing estimator (reweighted toward the target) at each source site
- DiffMeans: average of difference-in-means ATE at each source site
- All followed by inverse-variance weighted meta-regression adjusting for site-level covariates

Empirical Performance by Hypotheses



Predicted ATE vs empirical benchmark



Concluding Remarks

- How to generalize function-valued quantities across heterogeneous studies using only aggregate-level information?
- Cross-study heterogeneity can be captured by a low-rank structure
- Methodological contribution:
 - functional basis recovery via denoised functional SPA (d-fSPA)
 - flexible weight modeling linking study-level covariates to mixing proportions
 - distribution-free uncertainty quantification via conformal prediction
- Theoretical guarantees:
 - consistent recovery of latent basis functions
 - asymptotically valid marginal coverage for target predictions
- **MetaHunt** enables privacy-preserving, ML-compatible functional meta-analysis that accommodates both covariate shift and conditional shift