



Covariate-Powered Empirical Bayes Estimation

Nikolaos Ignatiadis¹ and Stefan Wager²

¹ Department of Statistics, Stanford University (ignat@stanford.edu)

² Graduate School of Business, Stanford University (swager@stanford.edu)

Setup

Model I: We observe $(i = 1, \dots, n)$

$$\begin{aligned} X_i &\stackrel{\text{iid}}{\sim} \mathbb{P}^X \\ \mu_i | X_i &\sim \mathcal{N}(m(X_i), A) \\ Z_i | \mu_i &\sim \mathcal{N}(\mu_i, \sigma^2), \end{aligned}$$

where $m(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$ and $A > 0$ unknown, $\sigma^2 > 0$ known.

Goal: Estimate μ_i by $\hat{\mu}_i$ s.t.

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\mu_i - \hat{\mu}_i)^2] \text{ is small}$$

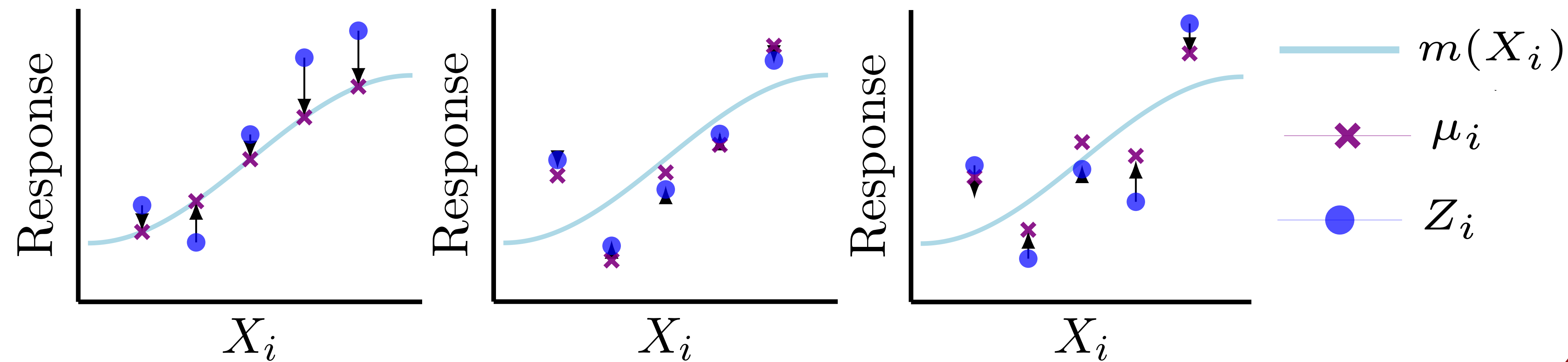
Benchmark: The Bayes rule

$$t_{m,A}^*(x, z) = \mathbb{E}_{m,A}[\mu_i | X_i = x, Z_i = z] = \frac{A}{\sigma^2 + A} z + \frac{\sigma^2}{\sigma^2 + A} m(x)$$

$$\begin{aligned} A = 0: \\ t_{m,A}^*(x, z) &= m(x) \end{aligned}$$

$$\begin{aligned} A \gg \sigma^2: \\ t_{m,A}^*(x, z) &\approx z \end{aligned}$$

$$\begin{aligned} A \approx \sigma^2: \\ \text{Convex combination} \end{aligned}$$



Empirical Bayes with Cross-Fitting

1. Form a partition of $\{1, \dots, n\}$ into two folds l_1 and l_2 .
2. Use observations in l_1 , to estimate the regression $m(x) = \mathbb{E}[Z_i | X_i = x]$ by $\hat{m}_{l_1}(\cdot)$.
3. Use observations in l_2 , to estimate A , through the formula

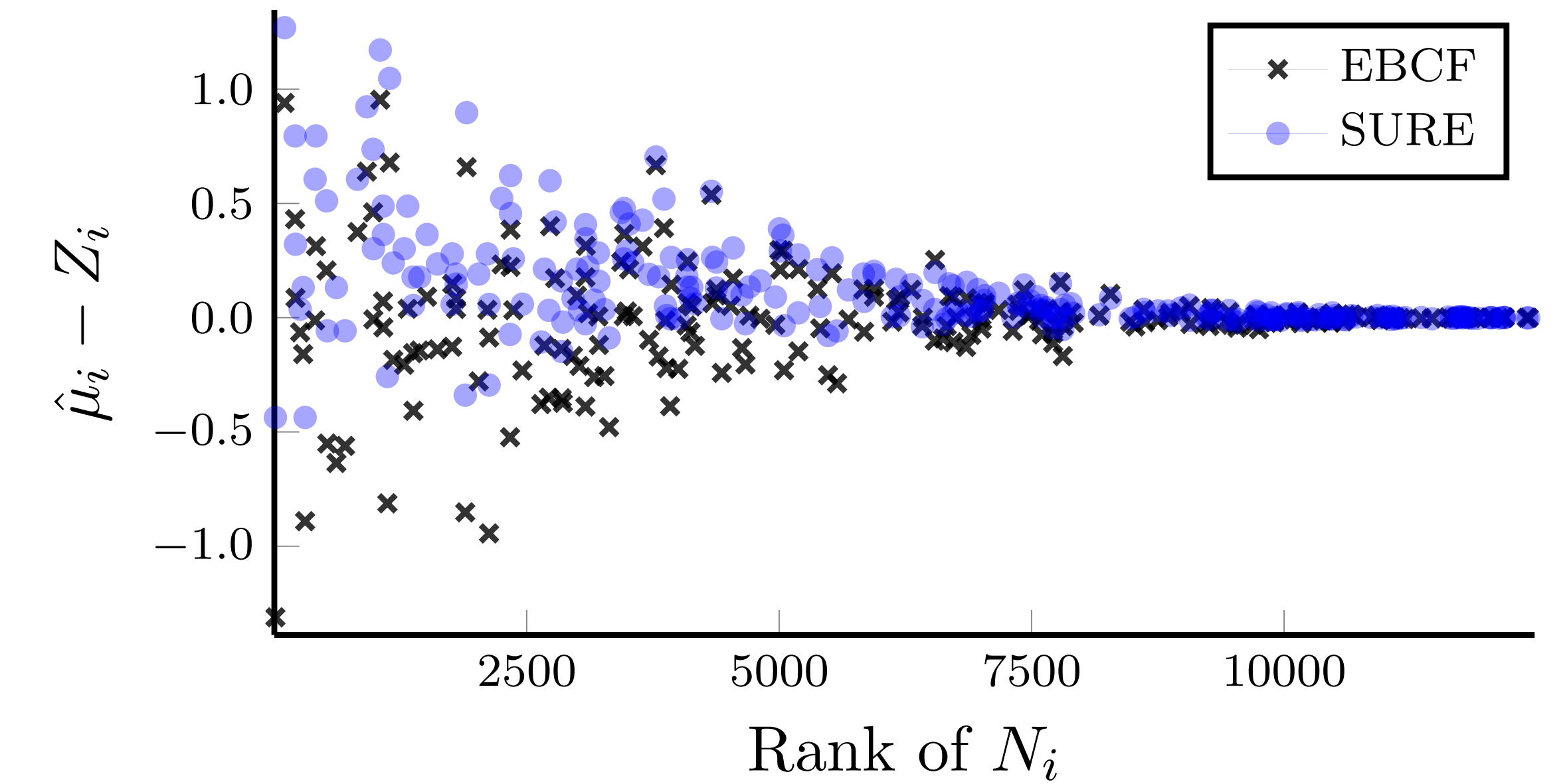
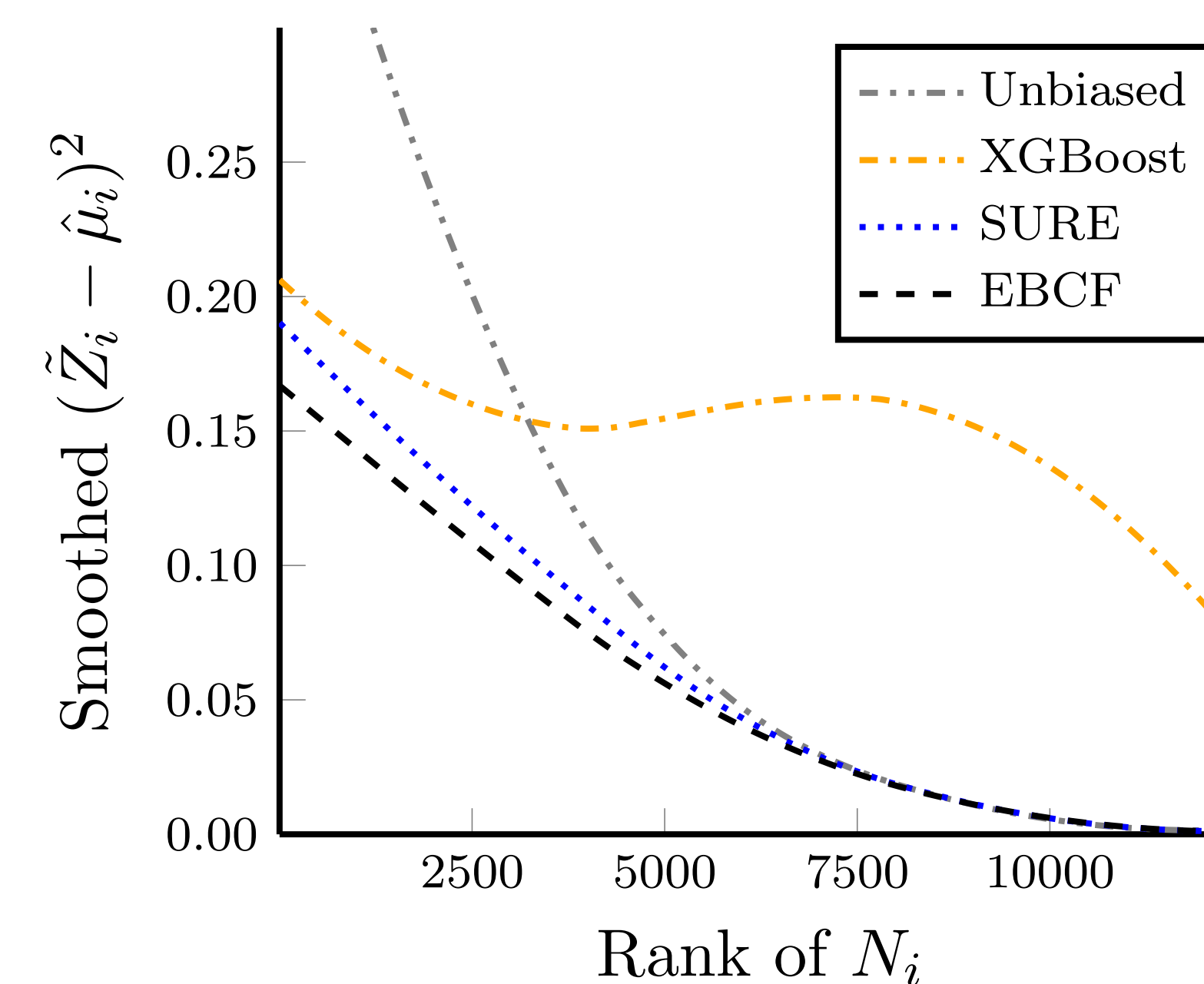
$$\hat{A}_{l_2} = \left(\frac{1}{|l_2|} \sum_{i \in l_2} (\hat{m}_{l_1}(X_i) - Z_i)^2 - \sigma^2 \right)_+$$

4. For $i \in l_2$, estimate μ_i by $\hat{\mu}_i^{\text{EBCF}} = t_{\hat{m}_{l_1}, \hat{A}_{l_2}}^*(X_i, Z_i)$.
5. Repeat with folds l_1 and l_2 flipped.

MovieLens [2] data analysis

- ▶ 20 million ratings in $\{0, 0.5, \dots, 5\}$ from 138,000 users applied to 27,000 movies ($n \geq 10,000$ after filtering).
- ▶ Keep 10% of users, calculate average rating Z_i for each movie based on N_i users.
- ▶ X_i includes N_i , year of release, genres.
- ▶ μ_i is "true" movie rating.
- ▶ Posit that $Z_i | \mu_i, X_i \sim (\mu_i, \sigma^2/N_i)$.
- ▶ "Ground-truth": \tilde{Z}_i , the average movie rating based on other 90% of users.
- ▶ Evaluation by mean-squared error: $\sum_{i=1}^n (\tilde{Z}_i - \hat{\mu}_i)^2/n$

		All	Sci-Fi & Horror
$Z_i \rightarrow$	Unbiased	0.098 (± 0.005)	0.098 (± 0.032)
[1] \rightarrow	XGBoost	0.150 (± 0.005)	0.210 (± 0.036)
[4] \rightarrow	SURE	0.061 (± 0.002)	0.064 (± 0.018)
This work \rightarrow	EBCF (with XGBoost)	0.055 (± 0.002)	0.051 (± 0.012)



EBCF is minimax optimal

Model I: $\implies X_i \stackrel{\text{iid}}{\sim} \mathbb{P}^X, Z_i | X_i \sim \mathcal{N}(m(X_i), A + \sigma^2)$

Minimax regression error over $\mathcal{C} \subset \{f : \mathcal{X} \rightarrow \mathbb{R}\}$

$$\mathfrak{M}_n^{\text{Reg}}(\mathcal{C}; A + \sigma^2) := \inf_{\hat{m}_n} \max_{m \in \mathcal{C}} \mathbb{E}_{m,A} \left[\int (\hat{m}_n(x) - m(x))^2 d\mathbb{P}^X(x) \right]$$

Minimax empirical Bayes excess risk [3] over \mathcal{C} , with $A > 0$ fixed (but unknown)

$$\mathfrak{M}_n^{\text{EB}}(\mathcal{C}; A, \sigma^2) \widehat{=} \inf_{\hat{t}_n} \max_{m \in \mathcal{C}} \{ \text{Expected risk of } \hat{t}_n - \text{Bayes risk} \}$$

Theorem: For many \mathcal{C} , e.g., Lipschitz class in \mathbb{R}^d

$$\mathfrak{M}_n^{\text{EB}}(\mathcal{C}; A, \sigma^2) \asymp \frac{\sigma^4}{(\sigma^2 + A)^2} \mathfrak{M}_n^{\text{Reg}}(\mathcal{C}; A + \sigma^2)$$

EBCF is robust to misspecification

Model II: Non-Gaussian, equal variances

$$(X_i, \mu_i, Z_i) \sim \mathbb{P}^{(X_i, \mu_i, Z_i)}, \mathbb{E}[Z_i | \mu_i, X_i] = \mu_i, \text{Var}[Z_i | \mu_i, X_i] = \sigma^2$$

Guarantees for EBCF in fold l_2 (under bounded $\mathbb{E}[Z_i^4 | \mu_i, X_i], \mu_i$):

$$\frac{1}{|l_2|} \sum_{i \in l_2} \mathbb{E}[(\mu_i - \hat{\mu}_i^{\text{EBCF}})^2] \leq \left\{ \frac{1}{|l_2|} \sum_{i \in l_2} \mathbb{E}[(\mu_i - \hat{m}_{l_1}(X_i))^2] \right\} + o\left(\frac{1}{\sqrt{|l_2|}}\right)$$

EBCF can be extended (with similar guarantees) to

Model III: Non-Gaussian, unequal variances

$$(X_i, \mu_i, Z_i) \sim \mathbb{P}^{(X_i, \mu_i, Z_i)}, \mathbb{E}[Z_i | \mu_i, X_i] = \mu_i, \text{Var}[Z_i | \mu_i, X_i] = \sigma_i^2$$

Resources

Code availability

Software: <https://github.com/ignatiadis/EBayes.jl>
Reproducibility: <https://github.com/ignatiadis/EBCrossFitPaper>

References

- [1] Tianqi Chen and Carlos Guestrin. *ACM SIGKDD* 22:785–794, 2016.
- [2] F Maxwell Harper and Joseph A Konstan. *ACM TIS*, 5(4):19, 2016.
- [3] Herbert Robbins. *Annals of Mathematical Statistics*, 35:1–20, 1964.
- [4] Xianchao Xie, SC Kou, and Lawrence D Brown. *JASA*, 107(500):1465–1479, 2012.

Acknowledgments

This research was funded by a gift from Google.