

Robust Empirical Bayes Confidence Intervals^{*}

Timothy B. Armstrong[†]

Michal Kolesár[‡]

University of Southern California

Princeton University

Mikkel Plagborg-Møller[§]

Princeton University

September 27, 2021

Abstract

We construct robust empirical Bayes confidence intervals (EBCIs) in a normal means problem. The intervals are centered at the usual linear empirical Bayes estimator, but use a critical value accounting for shrinkage. Parametric EBCIs that assume a normal distribution for the means (Morris, 1983b) may substantially undercover when this assumption is violated. In contrast, our EBCIs control coverage regardless of the means distribution, while remaining close in length to the parametric EBCIs when the means are indeed Gaussian. If the means are treated as fixed, our EBCIs have an average coverage guarantee: the coverage probability is at least $1 - \alpha$ on average across the n EBCIs for each of the means. Our empirical application considers the effects of U.S. neighborhoods on intergenerational mobility.

Keywords: average coverage, empirical Bayes, confidence interval, shrinkage

JEL codes: C11, C14, C18

^{*}This paper is dedicated to the memory of Gary Chamberlain, who had a profound influence on our thinking about decision problems in econometrics, and empirical Bayes methods in particular. Luther Yap provided excellent research assistance. We received helpful comments from four anonymous referees, Otávio Bartalotti, Toru Kitagawa, Laura Liu, Ulrich Müller, Stefan Wager, Mark Watson, Martin Weidner, and numerous seminar participants. We are especially indebted to Bruce Hansen and Roger Koenker for inspiring our simulation study. Armstrong acknowledges support by the National Science Foundation Grant SES-2049765. Kolesár acknowledges support by the Sloan Research Fellowship and by the National Science Foundation Grant SES-22049356. Plagborg-Møller acknowledges support by the National Science Foundation Grant SES-1851665.

[†]email: timothy.armstrong@usc.edu

[‡]email: mcolesar@princeton.edu

[§]email: mikkelpm@princeton.edu

1 Introduction

Empirical researchers in economics are often interested in estimating effects for many individuals or units, such as estimating teacher quality for teachers in a given geographic area. In such problems, it has become common to shrink unbiased but noisy preliminary estimates of these effects toward baseline values, say the average fixed effect for teachers with the same experience. In addition to estimating teacher quality (Kane and Staiger, 2008; Jacob and Lefgren, 2008; Chetty et al., 2014), shrinkage techniques have been used recently in a wide range of applications including estimating school quality (Angrist et al., 2017), hospital quality (Hull, 2020), the effects of neighborhoods on intergenerational mobility (Chetty and Hendren, 2018), and patient risk scores across regional health care markets (Finkelstein et al., 2017).

The shrinkage estimators used in these applications can be motivated by an empirical Bayes (EB) approach. One imposes a working assumption that the individual effects are drawn from a normal distribution (or, more generally, a known family of distributions). The mean squared error (MSE) optimal point estimator then has the form of a Bayesian posterior mean, treating this distribution as a prior distribution. Rather than specifying the unknown parameters in the prior distribution *ex ante*, the EB estimator replaces them with consistent estimates, just as in random effects models. This approach is attractive because one does not need to assume that the effects are in fact normally distributed, or even take a “Bayesian” or “random effects” view: the EB estimators have lower MSE (averaged across units) than the unshrunk unbiased estimators, even when the individual effects are treated as nonrandom (James and Stein, 1961).

In spite of the popularity of EB methods, it is currently not known how to provide uncertainty assessments to accompany the point estimates without imposing strong parametric assumptions on the effect distribution. Indeed, Hansen (2016, p. 116) describes inference in shrinkage settings as an open problem in econometrics. The natural EB version of a confidence interval (CI) takes the form of a Bayesian credible interval, again using the postulated effect distribution as a prior (Morris, 1983b). If the distribution is correctly specified, this *parametric empirical Bayes confidence interval (EBCI)* will cover 95%, say, of the true effect parameters, under repeated sampling of the observed data *and* of the effect parameters. We refer to this notion of coverage as “EB coverage”, following the terminology in Morris (1983b, Eq. 3.6). Unfortunately, we show that, in the context of a normal means model, the parametric EBCI with nominal level 95% can have actual EB coverage as low as 74% for certain non-normal effect distributions. The potential undercoverage is increasing in the degree of shrinkage, and we derive a simple “rule of thumb” for gauging the potential coverage

distortion.

To allow easy uncertainty assessment in EB applications that is reliable irrespective of the degree of shrinkage, we construct novel *robust EBCIs* that take a simple form and control EB coverage *regardless* of the true effect distribution. Our baseline model is an (approximate) normal means problem $Y_i \sim N(\theta_i, \sigma_i^2)$, $i = 1, \dots, n$. In applications, Y_i represents a preliminary asymptotically unbiased estimate of the effect θ_i for unit i . Like the parametric EBCI that assumes a normal distribution for θ_i , the robust EBCI we propose is centered at the normality-based EB point estimate $\hat{\theta}_i$, but it uses a larger critical value to take into account the bias due to shrinkage. We provide software implementing our methods. EB coverage is controlled in the class of all distributions for θ_i that satisfy certain moment bounds, which we estimate consistently from the data (similarly to the parametric EBCI, which uses the second moment). We show that the baseline implementation of our robust EBCI is “adaptive” in the sense that its length is close to that of the parametric EBCI when the θ_i ’s are in fact normally distributed. Thus, little efficiency is lost from using the robust EBCI in place of the non-robust parametric one.

In addition to controlling EB coverage, we show that the robust $1 - \alpha$ EBCIs have a frequentist *average coverage* property: If the means $\theta_1, \dots, \theta_n$ are treated as *fixed*, the coverage probability—averaged across the n parameters θ_i —is at least $1 - \alpha$. This average coverage property weakens the usual notion of coverage, which would be imposed separately for each θ_i . As discussed in Remark 2.1 below, the average coverage criterion is motivated by the same idea as the usual EB point estimator (Efron, 2010): we seek CIs that are short and control coverage *on average*, without requiring good performance for every single unit i . Due to the weaker coverage requirement, our robust EBCIs are shorter than the usual CI centered at the unshrunk estimate Y_i , and often substantially so. Intuitively, the average coverage criterion only requires us to guard against the *average* coverage distortion induced by the biases of the individual shrinkage estimators $\hat{\theta}_i$, and the data is quite informative about whether *most* of these biases are large, even though individual biases are difficult to estimate. To complement the frequentist properties, our EBCIs can also be viewed as Bayesian credible sets that are robust to the prior on θ_i , in terms of *ex ante* coverage.

Our underlying ideas extend to other linear and non-linear shrinkage settings with possibly non-Gaussian data. For example, our techniques allow for the construction of robust EBCIs that contain (nonlinear) soft thresholding estimators, as well as average coverage confidence bands for nonparametric regression functions.

We illustrate our results by computing EBCIs for the causal effects of growing up in different U.S. neighborhoods (specifically commuting zones) on intergenerational mobility. We follow Chetty and Hendren (2018), who apply EB shrinkage to initial fixed effects estimates.

Depending on the specification, we find that the robust EBCIs are on average 12–25% as long as the unshrunk CIs.

The average coverage criterion was originally introduced in the literature on nonparametric regression (Wahba, 1983; Nychka, 1988; Wasserman, 2006, Ch. 5.8). Cai et al. (2014) construct rate-optimal adaptive confidence bands that achieve average coverage. These procedures are challenging to implement in our EB setting, and do not have a clear finite-sample justification, unlike our procedure. Liu et al. (2019) construct forecast intervals in a dynamic panel data model that guarantee average coverage in a Bayesian sense (for a fixed prior). We give a detailed discussion of alternative approaches to inference in EB settings in Section 5.

The rest of this paper is organized as follows. Section 2 illustrates our methods in the context of a simple homoskedastic Gaussian model. Section 3 presents our recommended baseline procedure and discusses practical implementation issues. Section 4 presents our main results on the coverage and efficiency of the robust EBCI, and on the coverage distortions of the parametric EBCI; we also verify the finite-sample coverage accuracy of the robust EBCI through extensive simulations. Section 5 compares our EBCI with other inference approaches. Section 6 discusses extensions of the basic framework. Section 7 contains an empirical application to inference on neighborhood effects. Appendices A to C give details on finite-sample corrections, computational details, and formal asymptotic coverage results. The Online Supplement contains proofs as well as further technical results. Applied readers are encouraged to focus on Sections 2, 3 and 7.

2 Simple example

This section illustrates the construction of the robust EBCIs that we propose in a simplified setting with no covariates and with known, homoskedastic errors. Section 3 relaxes these restrictions, and discusses other empirically relevant extensions of the basic framework, as well as implementation issues.

We observe n estimates Y_i of elements of the parameter vector $\theta = (\theta_1, \dots, \theta_n)'$. Each estimate is normally distributed with common, known variance σ^2 ,

$$Y_i \mid \theta \sim N(\theta_i, \sigma^2), \quad i = 1, \dots, n. \quad (1)$$

In many applications, the Y_i 's arise as preliminary least squares estimates of the parameters θ_i . For instance, they may correspond to fixed effect estimates of teacher or school value added, neighborhood effects, or firm and worker effects. In such cases, Y_i will only be *approximately* normal in large samples by the central limit theorem (CLT); we take this

explicitly into account in the theory in Appendix C.

A popular approach to estimation that substantially improves upon the raw estimator $\hat{\theta}_i = Y_i$ under the compound MSE $\sum_{i=1}^n E[(\hat{\theta}_i - \theta_i)^2]$ is based on empirical Bayes (EB) shrinkage. In particular, suppose that the θ_i 's are themselves normally distributed,

$$\theta_i \sim N(0, \mu_2). \quad (2)$$

Our discussion below applies if Eq. (2) is viewed as a subjective Bayesian prior distribution for a single parameter θ_i , but for concreteness we will think of Eq. (2) as a “random effects” sampling distribution for the n mean parameters $\theta_1, \dots, \theta_n$. Under this normal sampling distribution, it is optimal to estimate θ_i using the posterior mean $\hat{\theta}_i = w_{EB}Y_i$, where $w_{EB} = 1 - \sigma^2/(\sigma^2 + \mu_2)$. To avoid having to specify the variance μ_2 , the EB approach treats it as an unknown parameter, and replaces the marginal precision of Y_i , $1/(\sigma^2 + \mu_2)$, with a method of moments estimate $n/\sum_{i=1}^n Y_i^2$, or the degrees-of-freedom adjusted estimate $(n-2)/\sum_{i=1}^n Y_i^2$. The latter leads to $\hat{w}_{EB} = (1 - \sigma^2(n-2)/\sum_{i=1}^n Y_i^2)$, which is the classic estimator of James and Stein (1961).

One can also use Eq. (2) to construct CIs for the θ_i 's. In particular, since the marginal distribution of $w_{EB}Y_i - \theta_i$ is normal with mean zero and variance $(1 - w_{EB})^2\mu_2 + w_{EB}^2\sigma^2 = w_{EB}\sigma^2$, this leads to the $1 - \alpha$ CI

$$w_{EB}Y_i \pm z_{1-\alpha/2}w_{EB}^{1/2}\sigma, \quad (3)$$

where z_α is the α quantile of the standard normal distribution. Since the form of the interval is motivated by the parametric assumption (2), we refer to it as a parametric EBCI. With μ_2 unknown, one can replace w_{EB} by \hat{w}_{EB} .¹ This is asymptotically equivalent to (3) as $n \rightarrow \infty$.

The coverage of the parametric EBCI in (3) is $1 - \alpha$ under repeated sampling of (Y_i, θ_i) according to Eqs. (1) and (2). To distinguish this notion of coverage from the case with fixed θ , we refer to coverage under repeated sampling of (Y_i, θ_i) as “empirical Bayes coverage”. This follows the definition of an empirical Bayes confidence interval (EBCI) in Morris (1983b, Eq. 3.6) and Carlin and Louis (2000, Ch. 3.5). Unfortunately, this coverage property relies heavily on the parametric assumption (2). We show in Section 4.3 that the actual EB coverage of the nominal $1 - \alpha$ parametric EBCI can be as low as $1 - 1/\max\{z_{1-\alpha/2}, 1\}$ for certain non-normal distributions of θ_i with variance μ_2 ; for 95% EBCIs, this evaluates to 74%. This contrasts with existing results on estimation: although the empirical Bayes estimator is motivated by the parametric assumption (2), it performs well even if this assumption is

¹Alternatively, to account for estimation error in \hat{w}_{EB} , Morris (1983b) suggests adjusting the variance estimate $\hat{w}_{EB}\sigma^2$ to $\hat{w}_{EB}\sigma^2 + 2Y_i^2(1 - \hat{w}_{EB})^2/(n-2)$. The adjustment does not matter asymptotically.

dropped, with low MSE even if we treat θ as fixed.

This paper constructs an EBCI with a similar robustness property: the interval will be close in length to the parametric EBCI when Eq. (2) holds, but its EB coverage is at least $1 - \alpha$ without any parametric assumptions on the distribution of θ_i . To describe the construction, suppose that all that is known is that θ_i is sampled from a distribution with second moment given by μ_2 (in practice, we can replace μ_2 by the consistent estimate $n^{-1} \sum_{i=1}^n Y_i^2 - \sigma^2$). Conditional on θ_i , the estimator $w_{EB}Y_i$ has bias $(w_{EB} - 1)\theta_i$ and variance $w_{EB}^2\sigma^2$, so that the t -statistic $(w_{EB}Y_i - \theta_i)/w_{EB}\sigma$ is normally distributed with mean $b_i = (1 - 1/w_{EB})\theta_i/\sigma$ and variance 1. Therefore, if we use a critical value χ , the non-coverage of the CI $w_{EB}Y_i \pm \chi w_{EB}\sigma$, conditional on θ_i , will be given by the probability

$$r(b_i, \chi) = P(|Z - b_i| \geq \chi \mid \theta_i) = \Phi(-\chi - b_i) + \Phi(-\chi + b_i), \quad (4)$$

where Z denotes a standard normal random variable, and Φ denotes the standard normal cdf. Thus, by iterated expectations, under repeated sampling of θ_i , the non-coverage is bounded by

$$\rho(\sigma^2/\mu_2, \chi) = \sup_F E_F[r(b, \chi)] \quad \text{s.t.} \quad E_F[b^2] = \frac{(1 - 1/w_{EB})^2}{\sigma^2} \mu_2 = \frac{\sigma^2}{\mu_2}, \quad (5)$$

where E_F denotes expectation under $b \sim F$. Although this is an infinite-dimensional optimization problem over the space of distributions, it turns out that it admits a simple closed-form solution, which we give in Proposition B.1 in Appendix B. Moreover, because the optimization is a linear program, it can be solved even in the more general settings of applied relevance that we consider in Section 3.

Set $\chi = \text{cva}_\alpha(\sigma^2/\mu_2)$, where $\text{cva}_\alpha(t) = \rho^{-1}(t, \alpha)$, and the inverse is with respect to the second argument. Then the resulting interval

$$w_{EB}Y_i \pm \text{cva}_\alpha(\sigma^2/\mu_2)w_{EB}\sigma \quad (6)$$

will maintain coverage $1 - \alpha$ among all distributions of θ_i with $E[\theta_i^2] = \mu_2$ (recall that we estimate μ_2 consistently from the data). For this reason, we refer to it as a robust EBCI. Figure 1 in Section 3.1 gives a plot of the critical values for $\alpha = 0.05$. We show in Section 4.2 below that by also imposing a constraint on the fourth moment of θ_i , in addition to the second moment constraint, one can construct a robust EBCI that “adapts” to the Gaussian case in the sense that its length will be close to that of the parametric EBCI in Eq. (3) if these moment constraints are compatible with a normal distribution.

Instead of considering EB coverage, one may alternatively wish to assess uncertainty associated with the estimates $\hat{\theta}_i = w_{EB}Y_i$ when θ is treated as fixed. In this case, the EBCI

in Eq. (6) has an average coverage guarantee that

$$\frac{1}{n} \sum_{i=1}^n P(\theta_i \in [w_{EB}Y_i \pm \text{cva}_\alpha(\sigma^2/\mu_2)w_{EB}\sigma] \mid \theta) \geq 1 - \alpha, \quad (7)$$

provided that the moment constraint can be interpreted as a constraint on the empirical second moment on the θ_i 's, $n^{-1} \sum_{i=1}^n \theta_i^2 = \mu_2$. In other words, if we condition on θ , then the coverage is at least $1 - \alpha$ on average across the n EBCIs for $\theta_1, \dots, \theta_n$. To see this, note that the average non-coverage of the intervals is bounded by (5), except that the supremum is only taken over possible empirical distributions for $\theta_1, \dots, \theta_n$ satisfying the moment constraint. Since this supremum is necessarily smaller than $\rho(\sigma^2/\mu_2, \chi)$, it follows that the average coverage is at least $1 - \alpha$.²

The usual CIs $Y_i \pm z_{1-\alpha/2}\sigma$ also of course achieve average coverage $1 - \alpha$. The robust EBCI in Eq. (6) will, however, be shorter, especially when μ_2 is small relative to σ^2 —see Figure 3 below: by weakening the requirement that each CI covers the true parameter $1 - \alpha$ percent of the time to the requirement that the coverage is $1 - \alpha$ on average across the CIs, we can substantially shorten the CI length. It may seem surprising at first that we can achieve this by centering the CI at the shrinkage estimates $w_{EB}Y_i$. The intuition for this is that the shrinkage reduces the variability of the estimates. This comes at the expense of introducing bias in the estimates. However, we can on average control the resulting coverage loss by using the larger critical value $\text{cva}_\alpha(\sigma^2/\mu_2)$. Because under the average coverage criterion we only need to control the bias *on average* across i , rather than for each individual θ_i , this increase in the critical value is smaller than the reduction in the standard error.

Remark 2.1 (Interpretation of average coverage). While the average coverage criterion is weaker than the classical requirement of guaranteed coverage for each parameter, we believe it is useful, particularly in the EB context, for three reasons. First, the EB *point estimator* achieves lower MSE on average across units at the expense of potentially worse performance for some individual units (see, for example, Efron, 2010, Ch. 1.3). Thus, researchers who use EB estimators instead of the unshrunk Y_i 's prioritize favorable group performance over protecting individual performance. It is natural to resolve the trade-off in the same way when it comes to uncertainty assessments. Our average coverage intervals do exactly this: they guarantee coverage and achieve short length on average across units at the expense of giving up on a coverage guarantee for every individual unit. See Section 5 for further

²This link between average risk of separable decision rules (here coverage of CIs, each of which depends only on Y_i) when the parameters $\theta_1, \dots, \theta_n$ are treated as fixed and the risk of a single decision rule when these parameters are random and identically distributed is a special case of what Jiang and Zhang (2009) call the fundamental theorem of compound decisions, which goes back to Robbins (1951).

discussion.

Second, one motivation for the usual notion of coverage is that if one constructs many CIs, and there is not too much dependence between the data used to construct each interval, then by the law of large numbers, at least a $1 - \alpha$ fraction of them will contain the corresponding parameter. As we discuss further in Remark 4.1, average coverage intervals also have this interpretation.

Finally, under the classical requirement of guaranteed coverage for each θ_i , it is not possible to substantially improve upon the usual CI centered at the unshrunk estimate Y_i , regardless of how one forms the CI.³ It is only by relaxing the coverage requirement that we can circumvent this impossibility result and obtain intervals that reflect the efficiency improvement from the empirical Bayes approach.

3 Practical implementation

We now describe how to compute a robust EBCI that allows for heteroskedasticity, shrinks towards more general regression estimates rather than towards zero, and exploits higher moments of the bias to yield a narrower interval. In Section 3.1, we describe the empirical Bayes model that motivates our baseline approach. Section 3.2 describes the practical implementation of our baseline approach.

3.1 Motivating model and robust EBCI

In applied settings, the unshrunk estimates Y_i will typically have heteroskedastic variances. Furthermore, rather than shrinking towards zero, it is common to shrink toward an estimate of θ_i based on some covariates X_i , such as a regression estimate $X_i'\hat{\delta}$. We now describe how to adapt the ideas in Section 2 to such settings.

Consider a generalization of the model in Eq. (1) that allows for heteroskedasticity and covariates,

$$Y_i \mid \theta_i, X_i, \sigma_i \sim N(\theta_i, \sigma_i^2), \quad i = 1, \dots, n. \quad (8)$$

The covariate vector X_i may contain just the intercept, and it may also contain (functions of) σ_i . To construct an EB estimator of θ_i , consider the working assumption that the sampling

³The results in Pratt (1961) imply that for CIs with coverage 95%, one cannot achieve expected length improvements greater than 15% relative to the usual unshrunk CIs, even if one happens to optimize length for the true parameter vector $(\theta_1, \dots, \theta_n)$. See, for example, Corollary 3.3 in Armstrong and Kolesár (2018) and the discussion following it.

distribution of the θ_i 's is conditionally normal:

$$\theta_i \mid X_i, \sigma_i \sim N(\mu_{1,i}, \mu_2), \quad \text{where} \quad \mu_{1,i} = X_i' \delta. \quad (9)$$

The hierarchical model (8)–(9) leads to the Bayes estimate $\hat{\theta}_i = \mu_{1,i} + w_{EB,i}(Y_i - \mu_{1,i})$, where $w_{EB,i} = \frac{\mu_2}{\mu_2 + \sigma_i^2}$. This estimate shrinks the unrestricted estimate Y_i of θ_i toward $\mu_{1,i} = X_i' \delta$. In contrast to (8), the normality assumption (9) typically cannot be justified simply by appealing to the CLT; the linearity of the conditional mean $\mu_{1,i} = X_i' \delta$ may also be suspect. Our robust EBCI will therefore be constructed so that it achieves valid EB coverage even if assumption (9) fails. To obtain a narrow robust EBCI, we augment the second moment restriction used to compute the critical value in Eq. (5) with restrictions on higher moments of the bias of $\hat{\theta}_i$. In our baseline specification, we add a restriction on the fourth moment.

In particular, we replace assumption (9) with the much weaker requirement that the conditional second moment and kurtosis of $\varepsilon_i = \theta_i - X_i' \delta$ do not depend on (X_i, σ_i) :

$$E[(\theta_i - X_i' \delta)^2 \mid X_i, \sigma_i] = \mu_2, \quad E[(\theta_i - X_i' \delta)^4 \mid X_i, \sigma_i] / \mu_2^2 = \kappa, \quad (10)$$

where δ is defined as the probability limit of the regression estimate $\hat{\delta}$.⁴ We discuss this requirement further in Remark 3.1 below, and we relax it in Remark 3.2 below.

We now apply analysis analogous to that in Section 2. Let us suppose for simplicity that δ , μ_2 , κ , and σ_i are known; we relax this assumption in Section 3.2 below, and in the theory in Section 4. Denote the conditional bias of $\hat{\theta}_i$ normalized by the standard error by $b_i = (w_{EB,i} - 1)\varepsilon_i / (w_{EB,i}\sigma_i) = -\sigma_i \varepsilon_i / \mu_2$. Under repeated sampling of θ_i , the non-coverage of the CI $\hat{\theta}_i \pm \chi w_{EB,i} \sigma$, conditional on (X_i, σ_i) , depends on the distribution of the normalized bias b_i , as in Section 2. Given the moments μ_2 and κ , the *maximal* non-coverage is given by

$$\rho(m_{2,i}, \kappa, \chi) = \sup_F E_F[r(b, \chi)] \quad \text{s.t.} \quad E_F[b^2] = m_{2,i}, \quad E_F[b^4] = \kappa m_{2,i}^2, \quad (11)$$

where b is distributed according to the distribution F . Here $m_{2,i} = E[b_i^2 \mid X_i, \sigma_i] = \sigma_i^2 / \mu_2$. Observe that the kurtosis of b_i matches that of ε_i . Appendix B shows that the infinite-dimensional linear program (11) can be reduced to two nested *univariate* optimizations. We also show that the least favorable distribution—the distribution F maximizing (11)—is a discrete distribution with up to 4 support points (see Remark B.1).

Define the critical value $\text{cva}_\alpha(m_{2,i}, \kappa) = \rho^{-1}(m_{2,i}, \kappa, \alpha)$, where the inverse is in the last argument. Figure 1 plots this function for $\alpha = 0.05$ and selected values of κ . This leads to

⁴Our framework can be modified to let (X_i, σ_i) be fixed, in which case δ depends on n . See the discussion following Theorem 4.1 below.

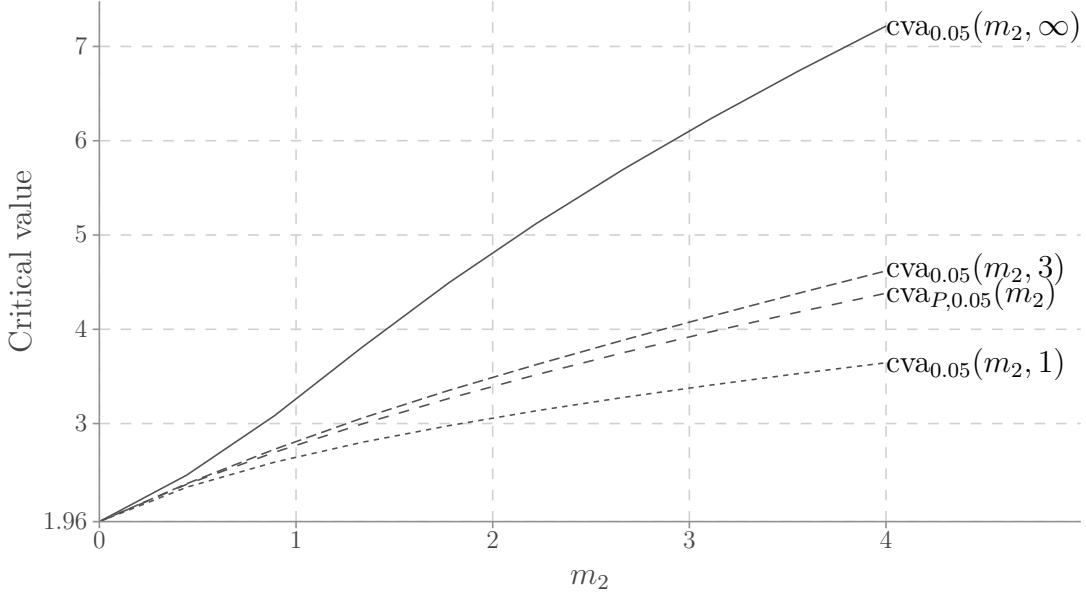


Figure 1: Function $\text{cva}_\alpha(m_2, \kappa)$ for $\alpha = 0.05$ and selected values of κ . The function $\text{cva}_\alpha(m_2)$, defined in Section 2, that only imposes a constraint on the second moment, corresponds to $\text{cva}_\alpha(m_2, \infty)$. The function $\text{cva}_{P,\alpha}(m_2) = z_{1-\alpha/2}\sqrt{1+m_2}$ corresponds to the critical value under the assumption that θ_i is normally distributed.

the robust EBCI

$$\hat{\theta}_i \pm \text{cva}_\alpha(m_{2,i}, \kappa) w_{EB,i} \sigma_i, \quad (12)$$

which, by construction, has coverage at least $1 - \alpha$ under repeated sampling of (Y_i, θ_i) , conditional on (X_i, σ_i) , so long as Eq. (10) holds; it is not required that (9) holds. Note that both the critical value and the CI length are increasing in σ_i .

3.2 Baseline implementation

Our baseline implementation of the robust EBCI plugs in consistent estimates of the unknown quantities in Eq. (12), based on the data $\{Y_i, X_i, \hat{\sigma}_i\}_{i=1}^n$, where $\hat{\sigma}_i$ is a consistent estimate of σ_i (such as the standard error of the preliminary estimate Y_i), and X_i is a vector of covariates that are thought to help predict θ_i .

1. Regress Y_i on X_i to obtain the fitted values $X_i' \hat{\delta}$, with $\hat{\delta} = (\sum_{i=1}^n \omega_i X_i X_i')^{-1} \sum_{i=1}^n \omega_i X_i Y_i$ denoting the weighted least squares estimate with precision weights ω_i . Two natural choices are setting $\omega_i = \hat{\sigma}_i^{-2}$, or setting $\omega_i = 1/n$ for unweighted estimates; see Appendix A.2 for further discussion. Let $\hat{\mu}_2 = \max \left\{ \frac{\sum_{i=1}^n \omega_i (\hat{\epsilon}_i^2 - \hat{\sigma}_i^2)}{\sum_{i=1}^n \omega_i}, \frac{2 \sum_{i=1}^n \omega_i^2 \hat{\sigma}_i^4}{\sum_{i=1}^n \omega_i \cdot \sum_{i=1}^n \omega_i \hat{\sigma}_i^2} \right\}$, and $\hat{\kappa} = \max \left\{ \frac{\sum_{i=1}^n \omega_i (\hat{\epsilon}_i^4 - 6 \hat{\sigma}_i^2 \hat{\epsilon}_i^2 + 3 \hat{\sigma}_i^4)}{\hat{\mu}_2^2 \sum_{i=1}^n \omega_i}, 1 + \frac{32 \sum_{i=1}^n \omega_i^2 \hat{\sigma}_i^8}{\hat{\mu}_2^2 \sum_{i=1}^n \omega_i \cdot \sum_{i=1}^n \omega_i \hat{\sigma}_i^4} \right\}$, where $\hat{\epsilon}_i = Y_i - X_i' \hat{\delta}$.

2. Form the EB estimate

$$\hat{\theta}_i = X_i' \hat{\delta} + \hat{w}_{EB,i}(Y_i - X_i' \hat{\delta}), \quad \text{where} \quad \hat{w}_{EB,i} = \frac{\hat{\mu}_2}{\hat{\mu}_2 + \hat{\sigma}_i^2}.$$

3. Compute the critical value $\text{cva}_\alpha(\hat{\sigma}_i^2/\hat{\mu}_2, \hat{\kappa})$ defined below Eq. (11).

4. Report the robust EBCI

$$\hat{\theta}_i \pm \text{cva}_\alpha(\hat{\sigma}_i^2/\hat{\mu}_2, \hat{\kappa}) \hat{w}_{EB,i} \hat{\sigma}_i. \quad (13)$$

We provide a fast and stable software package that automates these steps.⁵ We discuss the assumptions needed for validity of the robust EBCI in Remarks 3.1 to 3.3 below.

Remark 3.1 (Conditional EB coverage and moment independence). A potential concern about EB coverage in a heteroskedastic setting is that in order to reduce the length of the CI on average, one could choose to overcover parameters θ_i with small σ_i and undercover parameters θ_i with large σ_i . Our robust EBCI ensures that this does not happen by requiring EB coverage to hold conditional on (X_i, σ_i) . This also avoids analogous coverage concerns as a result of the value of X_i .

The key to ensuring this property is assumption (10) that the conditional second moment and kurtosis of $\varepsilon_i = \theta_i - X_i' \delta$ does not depend on (X_i, σ_i) . Conditional moment independence assumptions of this form are common in the literature. For instance, it is imposed in the analysis of neighborhood effects in Chetty and Hendren (2018) (their approach requires independence of the second moment), which is the basis for our empirical application in Section 7. Nonetheless, such conditions may be strong in some settings, as argued by Xie et al. (2012) in the context of EB point estimation. In Remark 3.2 below, we drop condition (10) entirely by replacing $\hat{\mu}_2$ and $\hat{\kappa}$ with nonparametric estimates of these conditional moments; alternatively, one could relax it by using a flexible parametric specification.⁶

Remark 3.2 (Nonparametric moment estimates). As a robustness check to guard against failure of the moment independence assumption (10), one may replace the critical value in Eq. (13) with $\text{cva}_\alpha((1 - 1/\hat{w}_{EB,i})^2 \hat{\mu}_{2i}/\hat{\sigma}_i^2, \hat{\kappa}_i)$, where $\hat{\mu}_{2i}$ and $\hat{\kappa}_i$ are consistent nonparametric estimates of $\mu_{2i} = E[(\theta_i - X_i' \delta)^2 | X_i, \sigma_i]$ and $\kappa_i = E[(\theta_i - X_i' \delta)^4 | X_i, \sigma_i]/\mu_{2i}^2$. The resulting CI will be asymptotically equivalent to the CI in the baseline implementation if Eq. (13) holds, but it will achieve valid EB coverage even if this assumption fails. In our empirical

⁵Matlab, Stata, and R packages are available at <https://github.com/kolesarm/ebci>

⁶Another way to drop condition (10) is to base shrinkage on the t -statistics Y_i/σ_i , applying the baseline implementation above with $Y_i/\hat{\sigma}_i$ in place of Y_i and 1 in place of $\hat{\sigma}_i$. Then the homoskedastic analysis in Section 2 applies, leading to valid EBCIs without any assumptions about independence of the moments. See Remark 3.8 and Appendix D.1 in Armstrong et al. (2020) for further discussion.

application, we use nearest-neighbor estimates, as described in Appendix A.1. As a simple diagnostic to gauge how the second moment of $\theta_i - X_i'\delta$ varies with (X_i, σ_i) , one can report the R^2 gain in predicting $\hat{\varepsilon}_i^2 - \hat{\sigma}_i^2$ using $\hat{\mu}_{2i}$ rather than the baseline estimate $\hat{\mu}_2$, as we illustrate in our empirical application.

Remark 3.3 (Average coverage and non-independent sampling). We show in Section 4 that the robust EBCI satisfies an average coverage criterion of the form (7) when the parameters $\theta = (\theta_1, \dots, \theta_n)$ are considered fixed, in addition to achieving valid EB coverage when the θ_i 's are viewed as random draws from some underlying distribution. To guarantee average coverage or EB coverage, we do not need to assume that the Y_i 's and θ_i 's are drawn independently across i . This is because the average coverage and EB coverage criteria only depend on the marginal distribution of (Y_i, θ_i) , not the joint distribution. Indeed, in deriving the infeasible CI in Eq. (12), we made no assumptions about the dependence structure of (Y_i, θ_i) across i . Consequently, to guarantee asymptotic coverage of the feasible interval in Eq. (13) as $n \rightarrow \infty$, we only need to ensure that the estimates $\hat{\mu}_2, \hat{\kappa}, \hat{\delta}, \hat{\sigma}_i$ are consistent for $\mu_2, \kappa, \delta, \sigma_i$, which is the case under many forms of weak dependence or clustering. Furthermore, our baseline implementation above does not require the researcher to take an explicit stand on the dependence of the data; for example, in the case of clustering, the researcher does not need to take an explicit stand on how the clusters are defined.

Remark 3.4 (Estimating moments of the distribution of θ_i). The estimators $\hat{\mu}_2$ and $\hat{\kappa}$ in step 1 of our baseline implementation above are based on the moment conditions $E[(Y_i - X_i'\delta)^2 - \sigma_i^2 \mid X_i, \sigma_i] = \mu_2$ and $E[(Y_i - X_i'\delta)^4 + 3\sigma_i^4 - 6\sigma_i^2(Y_i - X_i'\delta)^2 \mid X_i, \sigma_i] = \kappa\mu_2^2$, replacing population expectations by weighted sample averages. In addition, to avoid small-sample coverage issues when μ_2 and κ are near their theoretical lower bounds of 0 and 1, respectively, these estimates incorporate truncation on $\hat{\mu}_2$ and $\hat{\kappa}$, motivated by an approximation to a Bayesian estimate with flat prior on μ_2 and κ as in Morris (1983a,b). We verify the small-sample coverage accuracy of the resulting EBCIs through extensive simulations in Section 4.4. Appendix A.1 discusses the choice of the moment estimates, as well as other ways of performing truncation.

Remark 3.5 (Using higher moments and other forms of shrinkage). In addition to using the second and fourth moment of bias, one may augment (11) with restrictions on higher moments of the bias in order to further tighten the critical value. In Section 4.2, we show that using other moments in addition to the second and fourth moment does not substantially decrease the critical value in the case where θ_i is normally distributed. Thus, the CI in our baseline implementation is robust to failure of the normality assumption (9), while being near-optimal when this assumption does hold. Section 4.2 also shows that further efficiency

gains are possible if one uses the linear estimator $\tilde{\theta}_i = \mu_{1,i} + w_i(Y_i - \mu_{1,i})$ with the shrinkage coefficient w_i chosen to optimize CI length, instead of using the MSE-optimal shrinkage $w_{EB,i}$. For efficiency under a non-normal distribution of θ_i , one needs to consider non-linear shrinkage; we discuss this extension in Section 6.1.

4 Main results

This section provides formal statements of the coverage properties of the CIs presented in Sections 2 and 3. Furthermore, we show that the CIs presented in Sections 2 and 3 are highly efficient when the mean parameters are in fact normally distributed. Next, we calculate the maximal coverage distortion of the parametric EBCI, and derive a rule of thumb for gauging the potential coverage distortion. Finally, we present a comprehensive simulation study of the finite-sample performance of the robust EBCI. Applied readers interested primarily in implementation issues may skip ahead to the empirical application in Section 7.

4.1 Coverage under baseline implementation

In order to state the formal result, let us first carefully define the notions of coverage that we consider. Consider intervals CI_1, \dots, CI_n for elements of the parameter vector $\theta = (\theta_1, \dots, \theta_n)'$. The probability measure P denotes the joint distribution of θ and CI_1, \dots, CI_n . Following Morris (1983b, Eq. 3.6) and Carlin and Louis (2000, Ch. 3.5), we say that the interval CI_i is an (asymptotic) $1 - \alpha$ empirical Bayes confidence interval (EBCI) if

$$\liminf_{n \rightarrow \infty} P(\theta_i \in CI_i) \geq 1 - \alpha. \quad (14)$$

We say that the intervals CI_i are (asymptotic) $1 - \alpha$ average coverage intervals (ACIs) under the parameter sequence $\theta_1, \dots, \theta_n$ if

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n P(\theta_i \in CI_i \mid \theta) \geq 1 - \alpha. \quad (15)$$

The average coverage property (15) is a property of the distribution of the data conditional on θ and therefore does not require that we view the θ_i 's as random (as in a Bayesian or “random effects” analysis). To maintain consistent notation, we nonetheless use the conditional notation $P(\cdot \mid \theta)$ when considering average coverage. See Appendix C for a formulation with θ treated as nonrandom.

Observe that under the exchangeability condition that $P(\theta_i \in CI_i) = P(\theta_j \in CI_j)$ for all

i, j , if the ACI property (15) holds almost surely, then the EBCI property (14) holds, since then

$$P(\theta_i \in CI_i) = \frac{1}{n} \sum_{j=1}^n P(\theta_j \in CI_j) \geq 1 - \alpha + o(1) \quad \text{for all } i.$$

We now provide coverage results for the baseline implementation described in Section 3.2. To keep the statements in the main text as simple as possible, we (i) maintain the assumption that the unshrunk estimates Y_i follow an exact normal distribution conditional on the parameter θ_i , (ii) state the results only for the homoskedastic case where the variance σ_i of the unshrunk estimate Y_i does not vary across i , and (iii) consider only unconditional coverage statements of the form (14) and (15). In Appendix C, we allow the estimates Y_i to be only approximately normally distributed and allow σ_i to vary, and we verify that our assumptions hold in a linear fixed effects panel data model. We also formalize the statements about conditional coverage made in Remark 3.1.

Theorem 4.1. *Suppose $Y_i \mid \theta \sim N(\theta_i, \sigma^2)$. Let $\mu_{j,n} = \frac{1}{n} \sum_{i=1}^n (\theta_i - X_i' \delta)^j$ and let $\kappa_n = \mu_{4,n} / \mu_{2,n}^2$. Suppose the sequence $\theta = \theta_1, \dots, \theta_n$ and the conditional distribution $P(\cdot \mid \theta)$ satisfy the following conditions with probability one:*

1. $\mu_{2,n} \rightarrow \mu_2$ and $\mu_{4,n} / \mu_{2,n}^2 \rightarrow \kappa$ for some $\mu_2 \in (0, \infty)$ and $\kappa \in (1, \infty)$.
2. Conditional on θ , $(\hat{\delta}, \hat{\sigma}, \hat{\mu}_2, \hat{\kappa})$ converges in probability to $(\delta, \sigma, \mu_2, \kappa)$.

Then the CIs in Eq. (13) with $\hat{\sigma}_i = \hat{\sigma}$ satisfy the ACI property (15) with probability one. Furthermore, if $\theta_1, \dots, \theta_n$ follow an exchangeable distribution and the estimators $\hat{\delta}$, $\hat{\sigma}$, $\hat{\mu}_2$ and $\hat{\kappa}$ are exchangeable functions of the data $(X_1', Y_1)', \dots, (X_n', Y_n)'$, then these CIs satisfy the EB coverage property (14).

Theorem 4.1 follows immediately from Theorem C.2 in Appendix C. In order to cover both the EB coverage condition (14) and the average coverage condition (15), Theorem 4.1 considers a random sequence of parameters $\theta_1, \dots, \theta_n$, and shows average coverage conditional on these parameters. See Appendix C for a formulation with θ treated as nonrandom.

The condition on the moments μ_2 and κ avoids degenerate cases such as when $\mu_2 = 0$, in which case the EB point estimator $\hat{\theta}_i$ shrinks each preliminary estimate Y_i all the way to $X_i' \hat{\delta}$. Note also that the theorem does not require that $\hat{\delta}$ be the ordinary least squares (OLS) estimate in a regression of Y_i onto X_i , and that δ be the population analog; one can define δ in other ways, the theorem only requires that $\hat{\delta}$ be a consistent estimate of it. The definition of δ does, however, affect the plausibility of the moment independence assumption in Eq. (10) needed for conditional coverage results stated in Appendix C.⁷

⁷In addition, specification of $\mu_{1i} = X_i' \delta$ affects the width of the resulting EBCIs through its effect on μ_2 and κ .

Remark 4.1. As shown in Appendix C, if CIs satisfy the average coverage condition (15) given $\theta_1, \dots, \theta_n$, they will typically also satisfy the stronger condition

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\theta_i \in CI_i\} \geq 1 - \alpha + o_{P(\cdot|\theta)}(1), \quad (16)$$

where $o_{P(\cdot|\theta)}(1)$ denotes a sequence that converges in probability to zero conditional on θ (Eq. (16) implies Eq. (15) since the left-hand side is uniformly bounded). That is, at least a fraction $1 - \alpha$ of the n CIs contain their respective true parameters, asymptotically. This is analogous to the result that for estimation, the difference between the squared error $\frac{1}{n} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)^2$ and the MSE $\frac{1}{n} \sum_{i=1}^n E[(\hat{\theta}_i - \theta_i)^2 | \theta]$ typically converges to zero.

4.2 Relative efficiency

The robust EBCI in Eq. (12), unlike the parametric EBCI $\hat{\theta}_i \pm z_{1-\alpha/2} \sigma_i \sqrt{w_{EB,i}}$, does not rely on the normality assumption in Eq. (9) for its validity. We now show that this robustness does not come at a high cost in terms of efficiency when in fact the normality assumption (9) holds: the inefficiency is small unless the signal-to-noise ratio μ_2/σ_i^2 is very small.

There are two reasons for the inefficiency relative to this normal benchmark. First, the robust EBCI only makes use of the second and fourth moment of the conditional distribution of $\theta_i - X_i' \delta$, rather than its full distribution. Second, if we only have knowledge of these two moments, it is no longer optimal to center the EBCI at the estimator $\hat{\theta}_i$: one may need to consider other, perhaps non-linear, shrinkage estimators, as we do below in Section 6.1.

We decompose the sources of inefficiency by studying the relative length of the robust EBCI relative to the EBCI that picks the amount of shrinkage optimally. For the latter, we maintain assumption (10), and consider a more general class of estimators $\tilde{\theta}(w_i) = \mu_{1,i} + w_i(Y_i - \mu_{1,i})$: we impose the requirement that the shrinkage is linear for tractability, but allow the amount of shrinkage w_i to be optimally determined. The normalized bias of $\tilde{\theta}(w_i)$ is given by $b_i = (1/w_i - 1)\varepsilon_i/\sigma_i$, which leads to the EBCI

$$\mu_{1,i} + w_i(Y_i - \mu_{1,i}) \pm \text{cva}_\alpha((1 - 1/w_i)^2 \mu_2/\sigma_i^2, \kappa) w_i \sigma_i.$$

The half-length of this EBCI, $\text{cva}_\alpha((1 - 1/w_i)^2 \mu_2/\sigma_i^2, \kappa) w_i \sigma_i$, can be numerically minimized as a function of w_i to find the EBCI length-optimal shrinkage. Denote the minimizer by $w_{\text{opt}}(\mu_2/\sigma_i^2, \kappa, \alpha)$. Like $w_{EB,i}$, the optimal shrinkage depends on μ_2 and σ_i^2 only through the signal-to-noise ratio μ_2/σ_i^2 . Numerically evaluating the minimizer shows that $w_{\text{opt}}(\cdot, \kappa, \alpha) \geq w_{EB,i}$ for $\kappa \geq 3$ and $\alpha \in \{0.05, 0.1\}$. The resulting EBCI is optimal among all fixed-length

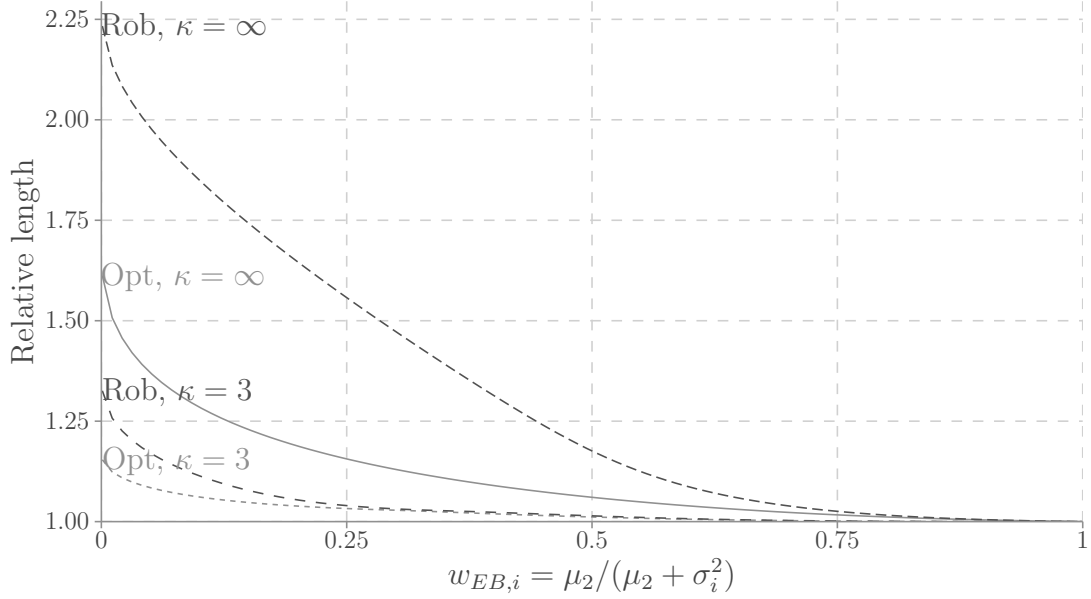


Figure 2: Relative efficiency of robust EBCI (Rob) and optimal robust EBCI (Opt) relative to the normal benchmark, for $\alpha = 0.05$. The figure plots ratios of Rob length, $2 \text{cva}_\alpha(\sigma_i^2/\mu_2, \kappa) \cdot \sigma_i \mu_2 / (\mu_2 + \sigma_i^2)$, and Opt length $2 \text{cva}_\alpha((1 - 1/w_{\text{opt}}(\mu_2/\sigma_i^2, \kappa, \alpha))^2 \mu_2/\sigma_i^2, \kappa) \cdot \sigma_i w_{\text{opt}}(\mu_2/\sigma_i^2, \kappa, \alpha)$, relative to the parametric EBCI length $2z_{1-\alpha/2} \sqrt{\mu_2/(\mu_2 + \sigma_i^2)} \sigma_i$ as a function of the shrinkage factor $w_{EB,i} = \mu_2/(\mu_2 + \sigma_i^2)$, which maps the signal-to-noise ratio μ_2/σ_i^2 to the interval $[0, 1]$.

EBCIs centered at linear estimators under (10), and we call it the optimal robust EBCI.

Figure 2 plots the ratio of lengths of the optimal robust EBCI and robust EBCI relative to the parametric EBCI, for $\alpha = 0.05$. The figure shows that for efficiency relative to the normal benchmark, it is relatively more important to impose the fourth moment constraint than to use the optimal amount of shrinkage (and only impose the second moment constraint). It also shows that the efficiency loss of the robust EBCI is modest unless the signal-to-noise ratio is very small: if $w_{EB,i} \geq 0.1$ (which is equivalent to $\mu_2/\sigma_i^2 \geq 1/9$), the efficiency loss is at most 11.4% for $\alpha = 0.05$; up to half of the efficiency loss is due to not using the optimal shrinkage. For $\alpha = 0.1$ (not plotted), the results are very similar; in particular, if $w_{EB,i} \geq 0.1$, the efficiency loss is at most 12.9%.

When the signal-to-noise ratio is very small, so that $w_{EB,i} < 0.1$, the efficiency loss of the robust EBCI is higher (up to 39% for $\alpha = 0.05$ or 0.1). Using the optimal robust EBCI ensures that the efficiency loss is below 20%, irrespective of the signal-to-noise ratio. On the other hand, when the signal-to-noise ratio is small, any of these CIs will be significantly tighter than the unshrunk CI $Y_i \pm z_{1-\alpha/2} \sigma_i$. To illustrate this point, Figure 3 plots the efficiency of the robust EBCI that imposes the second moment constraint only, relative to this unshrunk CI. It can be seen from the figure that shrinkage methods allow us to tighten

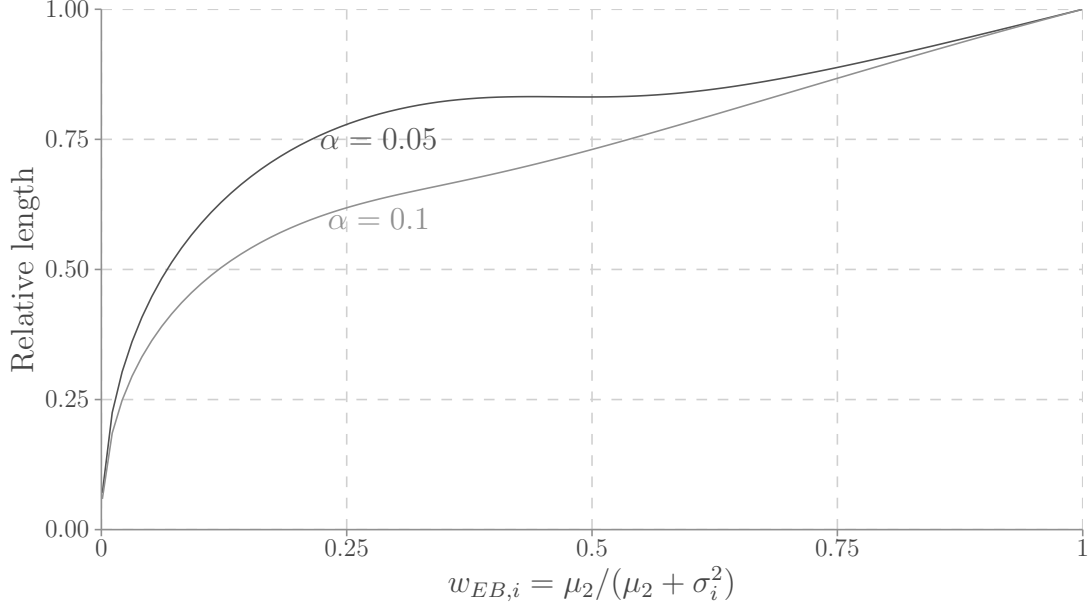


Figure 3: Relative efficiency of robust EBCI $\hat{\theta}_i \pm \text{cva}_\alpha(\sigma_i^2/\mu_2, \kappa = \infty) \cdot \sigma\mu_2/(\mu_2 + \sigma_i^2)$ relative to the unshrunk CI $Y_i \pm z_{1-\alpha/2}\sigma_i$. The figure plots the ratio of the length of the robust EBCI relative to the unshrunk CI as a function of the shrinkage factor $w_{EB,i} = \mu_2/(\mu_2 + \sigma_i^2)$.

the CI by 44% or more when $\mu_2/\sigma_i^2 \leq 0.1$.

4.3 Undercoverage of parametric EBCI

The parametric EBCI $\hat{\theta}_i \pm z_{1-\alpha/2}w_{EB,i}^{1/2}\sigma_i$ is an EB version of a Bayesian credible interval that treats (9) as a prior. We now assess its potential undercoverage when Eq. (9) is violated.

Given knowledge of only the second moment μ_2 of $\varepsilon_i = Y_i - X_i'\delta$, the maximal undercoverage of this interval is given by

$$\rho(1/w_{EB,i} - 1, z_{1-\alpha/2}/\sqrt{w_{EB,i}}), \quad (17)$$

since $w_{EB,i} = \mu_2/(\mu_2 + \sigma_i^2)$. Here ρ is the non-coverage function defined in Eq. (5). Figure 4 plots the maximal non-coverage probability as a function of $w_{EB,i}$, for significance levels $\alpha = 0.05$ and $\alpha = 0.10$. The figure suggests a simple “rule of thumb”: if $w_{EB,i} \geq 0.3$, the maximal coverage distortion is less than 5 percentage points for these values of α .

The following lemma confirms that the maximal non-coverage is decreasing in $w_{EB,i}$, as suggested by the figure. It also gives an expression for the maximal non-coverage across all values of $w_{EB,i}$ (which is achieved in the limit $w_{EB,i} \rightarrow 0$).

Lemma 4.1. *The non-coverage probability (17) of the parametric EBCI is weakly decreasing*

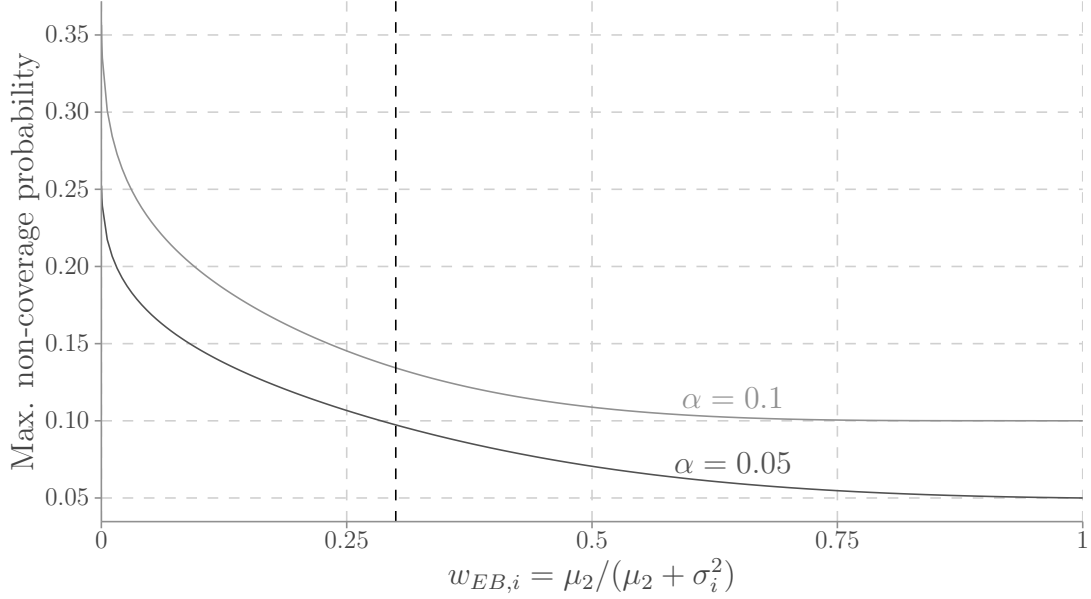


Figure 4: Maximal non-coverage probability of parametric EBCI, $\alpha \in \{0.05, 0.10\}$. The vertical line marks the “rule of thumb” value $w_{EB,i} = 0.3$, above which the maximal coverage distortion is less than 5 percentage points for these two values of α .

as a function of $w_{EB,i}$, with the supremum given by $1/\max\{z_{1-\alpha}^2, 1\}$.

The maximal non-coverage probability $1/\max\{z_{1-\alpha/2}^2, 1\}$ equals 0.260 for $\alpha = 0.05$ and 0.370 for $\alpha = 0.10$. For $\alpha > 2\Phi(-1) \approx 0.317$, the maximal non-coverage probability is 1.

If we additionally impose knowledge of the kurtosis of ε_i , the maximal non-coverage of the parametric EBCI can be similarly computed using Eq. (11), as illustrated in the application in Section 7.

4.4 Monte Carlo simulations

Here we show through simulations that the robust EBCI achieves accurate average coverage in finite samples.

4.4.1 Design

The data generating process (DGP) is a simple linear fixed effects panel data model. We first draw θ_i , $i = 1, \dots, n$, i.i.d. from a random effects distribution specified below. Then we simulate panel data from the model

$$W_{it} = \theta_i + U_{it}, \quad i = 1, \dots, n, \quad t = 1, \dots, T,$$

where the errors U_{it} are mean zero and i.i.d. across (i, t) and independent of the θ_i 's. The unshrunk estimator of θ_i is the sample average of W_{it} for unit i , with standard error obtained from the usual unbiased variance estimator:

$$Y_i = \frac{1}{T} \sum_{t=1}^T W_{it}, \quad \hat{\sigma}_i = \sqrt{\frac{1}{T(T-1)} \sum_{t=1}^T (W_{it} - Y_i)^2}.$$

We draw U_{it} from one of two distributions: (1) a normal distribution and (2) a (shifted) chi-squared distribution with 3 degrees of freedom. In case (1), Y_i is exactly normal conditional on θ_i , but $\hat{\sigma}_i^2$ does not exactly equal $\text{var}(Y_i \mid \theta_i)$ for finite T . In case (2), Y_i is non-normal and positively skewed (conditional on θ_i) for finite T .

We consider six random effects distributions for θ_i (see Supplemental Appendix E.1 for detailed definitions): (i) normal (kurtosis $\kappa = 3$); (ii) scaled chi-squared with 1 degree of freedom ($\kappa = 15$); (iii) two-point distribution ($\kappa \approx 8.11$); (iv) three-point distribution ($\kappa = 2$); (v) the least favorable distribution for the robust EBCI that exploits only second moments (κ depends on μ_2 , see Appendix B); and (vi) the least favorable distribution for the parametric EBCI.

Given T , we scale the θ_i distribution to match one of four signal-to-noise ratios $\mu_2 / \text{var}(Y_i \mid \theta_i) \in \{0.1, 0.5, 1, 2\}$, for a total of $6 \times 4 = 24$ DGPs for each distribution of U_{it} . We shrink towards the grand mean ($X_i = 1$ for all i). We construct the robust EBCIs following the baseline implementation in Section 3.2 (with $\omega_i = 1/n$), as well as a version that does not impose constraints on the kurtosis.

As $T \rightarrow \infty$, we recover the idealized setting in Section 2, with $(Y_i - \theta_i) / \sqrt{\text{var}(Y_i \mid \theta_i)}$ converging in distribution to a standard normal (conditional on θ_i), and $\hat{\sigma}_i^2 / \text{var}(Y_i \mid \theta_i)$ converging in probability to 1, for each i .

4.4.2 Results

Table 1 shows that the 95% robust EBCIs achieve good average coverage when the panel errors U_{it} are normally distributed. This is true for all DGPs, panel dimensions n and T , and whether we exploit one or both of the (estimated) moments μ_2 and κ . When the time dimension T equals 10, the maximal coverage distortion across all DGPs and all cross-sectional dimensions $n \in \{100, 200, 500\}$ is 3.2 percentage points. For $T \geq 20$, the coverage distortion of the robust EBCIs is always below 2.1 percentage points.

Table 2 shows that coverage distortions are somewhat larger when the panel errors U_{it} are chi-squared distributed and T is small. The robust EBCIs undercover by up to 7.2 percentage points when $T = 10$ due to the pronounced non-normality of Y_i given θ_i . However, the

Table 1: Monte Carlo simulation results, panel data with normal errors.

T	Robust, μ_2 only				Robust, μ_2 & κ				Parametric			
	10	20	∞	ora	10	20	∞	ora	10	20	∞	ora
Panel A: Average coverage (%), minimum across 24 DGPs												
$n = 100$	92.1	93.7	94.0	95.0	91.8	93.2	93.2	94.6	79.2	79.7	79.3	86.9
$n = 200$	91.9	93.4	92.9	95.0	91.8	93.3	92.9	94.8	80.7	80.3	81.0	86.3
$n = 500$	91.9	93.6	94.8	95.0	91.9	93.5	94.3	94.9	84.2	85.1	85.1	85.6
Panel B: Relative average length, average across 24 DGPs												
$n = 100$	1.09	1.10	1.11	1.16	1.03	1.02	1.02	1.00	0.81	0.82	0.83	0.86
$n = 200$	1.09	1.10	1.12	1.16	1.02	1.02	1.01	1.00	0.81	0.82	0.84	0.86
$n = 500$	1.10	1.11	1.13	1.16	1.04	1.03	1.01	1.00	0.82	0.83	0.84	0.86

Notes: Normally distributed errors. Nominal average confidence level $1 - \alpha = 95\%$. All EBCI procedures use baseline estimate of $\hat{\mu}_2$ and (if applicable) $\hat{\kappa}$, except columns labeled “ora”, which use oracle values of μ_2 and κ . Columns $T = \infty$ and “ora” use oracle standard errors σ_i . For each DGP, “average coverage” and “average length” refer to averages across units $i = 1, \dots, n$ and across 2,000 Monte Carlo repetitions. Average CI length is measured relative to the robust EBCI that exploits the oracle values of μ_2 , κ , and σ_i (but not of the grand mean $\delta = E[\theta]$).

Table 2: Monte Carlo simulation results, panel data with chi-squared errors.

T	Robust, μ_2 only				Robust, μ_2 & κ				Parametric			
	10	20	50	ora	10	20	50	ora	10	20	50	ora
Panel A: Average coverage (%), minimum across 24 DGPs												
$n = 100$	87.9	90.9	93.1	95.0	87.8	90.8	92.6	94.7	79.9	79.3	79.3	87.0
$n = 200$	87.9	90.8	93.0	94.9	87.8	90.8	92.8	94.8	77.8	79.8	80.3	86.2
$n = 500$	87.8	90.8	93.0	95.0	87.8	90.7	92.9	94.9	82.0	84.1	84.8	85.6
Panel B: Relative average length, average across 24 DGPs												
$n = 100$	1.05	1.08	1.10	1.16	1.01	1.02	1.02	1.00	0.79	0.81	0.82	0.86
$n = 200$	1.04	1.08	1.10	1.16	0.99	1.00	1.00	1.00	0.78	0.81	0.82	0.86
$n = 500$	1.05	1.09	1.11	1.16	0.99	1.00	1.00	1.00	0.79	0.82	0.83	0.86

Notes: Chi-squared distributed errors. See caption for Table 1. Results for $T = \infty$ are by definition the same as in Table 1.

distortion is at most 4.3 percentage points when $T = 20$, and at most 2.4 percentage points when $T \geq 50$. The coverage distortion due to non-normality when T is small is similar to the coverage distortion of the usual unshrunk CI (not reported).

Importantly, in all cases considered in Tables 1 and 2, the worst-case coverage distortion of the parametric EBCI substantially exceeds that of the corresponding robust EBCIs, sometimes by more than 10 percentage points. Nevertheless, the cost of robustness in terms of extra CI length is modest and consistent with the theoretical results in Section 4.2.

Both the estimation of the standard errors σ_i and the estimation of the moments μ_2 and κ contribute to the finite-sample coverage distortions. The “ora” columns in Table 1 exploit the oracle (true) values of μ_2 , κ , and $\sigma_i = \sqrt{\text{var}(Y_i | \theta_i)}$, while the $T = \infty$ columns use oracle standard errors but not oracle moments. By comparing these columns, we see that estimation of μ_2 and κ is responsible for modest coverage distortions when $n = 100$ or 200 . However, estimation of the standard errors σ_i also contributes to the distortions, as can be seen by comparing the $T = 10$ and $T = \infty$ columns.

In Supplemental Appendix E.2 we show that the robust EBCI also has good coverage in a heteroskedastic design calibrated to the empirical application in Section 7 below.

5 Comparison with other approaches

Here we compare our EBCI procedure with other approaches to confidence interval construction in the normal means model. We also discuss other related inference problems.

5.1 Average coverage vs. alternative coverage concepts

The average coverage requirement in Eq. (15) is less stringent than the usual (pointwise) notion of frequentist coverage that $P(\theta_i \in CI_i | \theta) \geq 1 - \alpha$ for all i . An even stronger coverage requirement is that of simultaneous coverage: $P(\forall i: \theta_i \in CI_i | \theta) \geq 1 - \alpha$. As discussed in Remark 2.1, under the pointwise coverage criterion, one cannot achieve substantial reductions in length relative to the unshrunk CI. Under the simultaneous coverage criterion, it is likewise impossible to substantially improve upon the usual sup- t confidence band based on the unshrunk estimates (Cai et al., 2014). Thus, undercoverage for some θ_i ’s must be tolerated if one wants to use shrinkage to improve CI length.

The fact that our EBCIs achieve improvements in average length at the expense of undercovering for certain units i is analogous to well-known properties of EB point estimators, as discussed in Remark 2.1. We now show that the units i for which our EBCI undercovers are quantitatively similar to the units for which the shrinkage estimator $\hat{\theta}_i$ has higher MSE

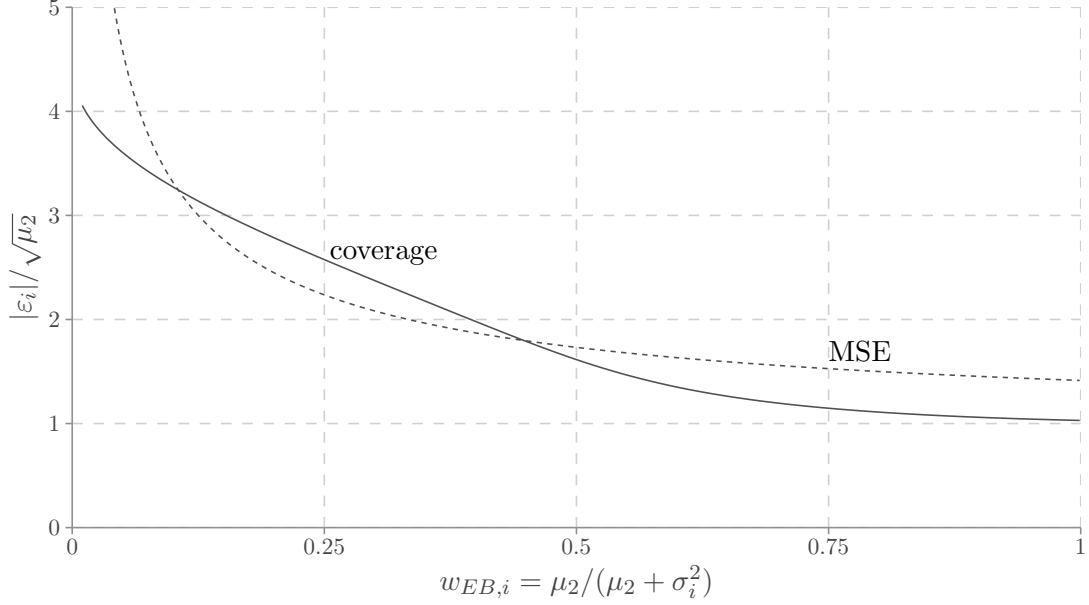


Figure 5: Value of $|\varepsilon_i|/\sqrt{\mu_2}$, as a function of $w_{EB,i}$, such that the MSE of the shrinkage point estimator equals that of the unshrunk estimator (MSE), and such that the coverage of the robust EBCI with $\kappa = \infty$ equals the nominal average coverage $1 - \alpha$ (coverage), for $\alpha = 0.05$.

than the unshrunk estimator Y_i . The pointwise coverage of our EBCI (conditional on X_i) is given by $1 - r(\sqrt{1/w_{EB,i} - 1} \cdot |\varepsilon_i|/\sqrt{\mu_2}, \text{cva}_\alpha(1/w_{EB,i} - 1, \kappa))$, with r defined in Eq. (4), $w_{EB,i} = \mu_2/(\mu_2 + \sigma_i^2)$, and $\varepsilon_i = \theta_i - X_i'\delta$ is the “shrinkage error” defined in Section 3.1. Since r is increasing in its first argument, for a fixed signal-to-noise ratio μ_2/σ_i^2 (and hence fixed amount of shrinkage $w_{EB,i}$), the coverage is decreasing in the normalized shrinkage error $|\varepsilon_i|/\sqrt{\mu_2}$: the units i for which our EBCI undercovers are those whose covariate-predicted value $X_i'\delta$ fails to approximate their true effect θ_i well. The MSE of the shrinkage estimator (for an individual unit i), normalized by the MSE σ_i^2 of the unshrunk estimator, equals $E[(\hat{\theta}_i - \theta_i)^2 | \theta_i, X_i]/\sigma_i^2 = w_{EB,i}^2 + (1 - w_{EB,i})w_{EB,i} \cdot |\varepsilon_i|/\sqrt{\mu_2}$; it is also increasing in $|\varepsilon_i|/\sqrt{\mu_2}$.

Figure 5 shows that the knife-edge value of $|\varepsilon_i|/\sqrt{\mu_2}$ for which the pointwise coverage of our EBCI equals $1 - \alpha$ is quantitatively close to the value of $|\varepsilon_i|/\sqrt{\mu_2}$ for which the MSE of the shrinkage estimator equals that of the unshrunk estimator. In other words, to the extent that one worries about undercoverage for certain types of θ_i values, one should simultaneously worry about the relative performance of the shrinkage point estimator for those same values.

We stress that the pointwise coverage depends on the unobservable shrinkage error ε_i , which cannot be gauged directly from the observables (Y_i, X_i) . If one wishes to avoid systematic differences in coverage across units i with different genders, say (i.e., one is worried that

ε_i correlates with gender) one can simply add gender to the set of covariates X_i : the baseline procedure in Section 3.2 ensures control of average coverage conditional on the covariates X_i . In Section 6.2, we show how to adapt our EBCIs to settings where one focuses the analysis on a subset of units i based on the values of their unshrunk estimates Y_i (e.g., keeping only the largest estimates).

From a Bayesian point of view, our robust EBCI can be viewed as an uncertainty interval that is robust to the choice of prior distribution in the *unconditional* gamma-minimax sense: the coverage probability of this CI is at least $1 - \alpha$ when averaged over the distribution of the data and over the prior distribution for θ_i , for any prior distribution that satisfies the moment bounds. This follows directly from the derivations in Section 2, reinterpreting the random effects distribution for θ_i as a prior distribution. In contrast, *conditional* gamma-minimax credible intervals, discussed recently by [Giacomini et al. \(2019, p. 6\)](#), are too stringent in our setting. This notion requires that the posterior credibility of the interval be at least $1 - \alpha$ regardless of the choice of prior, in any data sample, which would require reporting the entire parameter space (up to the moment bounds).

5.2 Finite-sample vs. asymptotic coverage

Our procedures are asymptotically valid as $n \rightarrow \infty$, as proved in Section 4.1. These asymptotics do not capture the impact of estimation error in the “hyper-parameters” $\hat{\sigma}_i$, $\hat{\delta}$, $\hat{\mu}_2$, and $\hat{\kappa}$, or the impact of lack of exact normality of the Y_i ’s, on the finite-sample performance of the EBCIs. As detailed in Section 3.2 and Appendix A, we do apply a finite-sample adjustment to the moments $\hat{\mu}_2$ and $\hat{\kappa}$, which is motivated by the same heuristic arguments that [Morris \(1983a,b\)](#) uses to motivate finite-sample adjustments to the parametric EBCI.⁸ The promising simulation results in Section 4.4 notwithstanding, these adjustments do not ensure exact average coverage control in finite samples.⁹

Our results are thus analogous to the results on coverage of Eicker-Huber-White CIs in cross-sectional OLS: asymptotic validity follows by consistency of the OLS variance estimate and asymptotic normality of the outcomes, while adjustments to account for finite-sample issues (such as the HC2 or HC3 variance estimators studied in [MacKinnon and White, 1985](#)) are justified heuristically. Deriving EBCIs with finite-sample coverage guarantees is an interesting problem that we leave for future research; the problem appears to be challenging

⁸An alternative approach would be to adapt the bootstrap adjustment proposed by [Carlin and Louis \(2000, Ch. 3.5.3\)](#) in the context of parametric EBCI construction (see also [Efron, 2019](#)). As with the [Morris \(1983a,b\)](#) adjustment, we are not aware of a formal result justifying it.

⁹Alternatively, one could account for hyperparameter uncertainty by computing the critical value $\sup_{\tilde{\sigma}_i, \tilde{\mu}_2, \tilde{\kappa} \in \hat{C}_i} \text{cva}_\alpha(\tilde{\sigma}_i^2 / \tilde{\mu}_2, \tilde{\kappa})$ over an initial confidence set \hat{C}_i for the hyper-parameters, coupled with a Bonferroni adjustment of the confidence level $1 - \alpha$. This approach appears to be highly conservative in practice.

even in the context of constructing parametric EBCIs.

5.3 Local vs. global optimality

Our EBCIs are designed to provide uncertainty assessments to accompany linear shrinkage estimates that, as the Introduction argues, have been popular in applied work. Our procedure’s global validity, as well as local near-optimality when the θ_i ’s are normal (cf. Section 4.2), is analogous to Eicker-Huber-White CIs for OLS estimators: these CIs are optimal under normal homoskedastic regression errors, but remain valid when this assumption is dropped.

Our EBCIs are, however, not globally efficient: when the θ_i ’s are not normally distributed, one may construct shorter CIs using non-linear shrinkage. In Section 6.1, we show how our method can be adapted to construct EBCIs that are locally near-optimal under non-normal baseline priors using non-linear shrinkage, such as soft thresholding. Since the distribution of θ_i is nonparametrically identified under the normal sampling model (8), it is in principle possible to construct EBCIs that are globally efficient using nonparametric methods. In the context of the homoskedastic model with no covariates in Eq. (1), various approaches to non-parametric point estimation of the θ_i ’s have been proposed, including kernels (Brown and Greenshtein, 2009), splines (Efron, 2019), or nonparametric maximum likelihood (Kiefer and Wolfowitz, 1956; Jiang and Zhang, 2009; Koenker and Mizera, 2014). It is an interesting problem for future research to adapt these methods to EBCI construction, while ensuring asymptotic validity, good finite-sample performance, and allowing for covariates, heteroskedasticity, and possible dependence across i .

5.4 Other inference problems

While we have sought to construct CIs for each of the n effects $\theta_1, \dots, \theta_n$, other papers have considered alternative inference problems in the normal means model. Efron (2015) develops a formula for the frequentist standard error of EB estimators, but this cannot be used to construct CIs without a corresponding estimate of the bias. Bonhomme and Weidner (2021) and Ignatiadis and Wager (2021) consider robust estimation and inference on functionals of the random effects *distribution*, rather than on the effects themselves. Finally, there is a substantial literature on shrinkage confidence balls, i.e., confidence sets of the form $\{\theta: \sum_{i=1}^n (\theta_i - \hat{\theta}_i)^2 \leq \hat{c}\}$ (see Casella and Hwang, 2012, for a review). While interesting from a theoretical perspective, these sets can be difficult to visualize and report in practice.¹⁰

¹⁰Confidence balls can be translated into intervals satisfying the average coverage criterion using Chebyshev’s inequality (see Wasserman, 2006, Ch. 5.8). However, the resulting intervals are very conservative compared to the ones we construct.

Finally, while we focus on CI length in our relative efficiency comparisons, our approach can be fruitfully applied when the goal of CI construction is to discern non-null effects, rather than to construct short CIs. In particular, suppose one forms a test of the null hypothesis $H_{0,i} : \theta_i = \theta_0$ for some null value θ_0 by rejecting when $\theta_0 \notin CI_i$, where CI_i is our robust EBCI given in (12). In Supplemental Appendix F, we show that the test based on our EBCI has higher average power than the usual z -test based on the unshrunk estimate when $X_i'\delta$ (the regression line towards which we shrink) is far enough from the null value θ_0 , and that these power gains can be substantial. Furthermore, such tests can be combined with corrections from the multiple testing literature to form procedures that asymptotically control the false discovery rate (FDR), a commonly used criterion for multiple testing.¹¹

6 Extensions

We now discuss two extensions of our method: adapting our intervals to general, possibly non-linear shrinkage, and constructing intervals that achieve coverage conditional on Y_i falling into a pre-specified interval.

6.1 General shrinkage

Our method can be generalized to cover general, possibly non-linear shrinkage based on possibly non-Gaussian data. Let $\mathcal{S}(y; \chi, \tilde{X}_i) \subseteq \mathbb{R}$ be a family of candidate confidence sets for a parameter θ_i , which depends on the data $Y_i = y$, a tuning parameter $\chi \in \mathbb{R}$ to be selected below, and covariates \tilde{X}_i (that include any known nuisance parameters) that we treat as fixed. We assume that \mathcal{S} is increasing in χ , in the sense of set containment, and that the non-coverage probability conditional on θ satisfies

$$P(\theta_i \notin \mathcal{S}(Y_i; \chi, \tilde{X}_i) \mid \theta, \tilde{X}^{(n)}) = \tilde{r}(a_i, \chi), \quad (18)$$

where a_i is some function of θ_i , $\tilde{X}^{(n)} = (\tilde{X}_1, \dots, \tilde{X}_n)$, and \tilde{r} is a known function (perhaps computed numerically or through simulation). Similarly to linear shrinkage in the normal means model, Eq. (18) may only hold approximately if the set \mathcal{S} depends on estimated

¹¹In particular, Storey (2002) shows that the Benjamini and Hochberg (1995) procedure asymptotically controls the FDR so long as the p -values do not exhibit too much statistical dependence and the proportion of rejected null hypotheses does not converge too quickly to zero. While Storey (2002) assumes that the uncorrected tests control size in the classical sense, the argument goes through essentially unchanged so long as the tests invert CIs that satisfy (16), which holds so long as the CIs do not exhibit too much statistical dependence, as discussed in Remark 4.1. We note, however, that this does not hold for modifications of the Benjamini and Hochberg (1995) procedure that use initial estimates of the proportion of true null hypotheses.

parameters (such as standard error estimates or tuning parameters), or if we use a large-sample approximation to the distribution of Y_i . We assume that a_i satisfies the moment constraints $E_F[g(a_i) \mid \tilde{X}^{(n)}] = m$, where g is a p -vector of moment functions,¹² and the expectation is over the conditional distribution F of a_i conditional on $\tilde{X}^{(n)}$. To guarantee EB coverage, we compute the maximal non-coverage

$$\rho_g(m, \chi) = \sup_F E_F[\tilde{r}(a, \chi)], \quad E_F[g(a)] = m, \quad (19)$$

analogously to Eq. (11). This is a linear program, which can be computed numerically to a high degree of precision even with several constraints; see Appendix B for details. Given an estimate \hat{m} of the moment vector m , we form a robust EBCI as

$$\mathcal{S}(Y_i; \hat{\chi}, \tilde{X}_i), \quad \text{where} \quad \hat{\chi} = \inf\{\chi: \rho_g(\hat{m}, \chi) \leq \alpha\}. \quad (20)$$

Example 6.1 (Linear shrinkage in the normal model). The setting in Section 3.1 obtains if we set $\tilde{X}_i = (X_i, \sigma_i)$ and $\mathcal{S}(y; \chi, \tilde{X}_i) = \{(1 - w_{EB,i})X_i'\delta + w_{EB,i}Y_i \pm \chi w_{EB,i}\sigma_i\}$. Here a_i is given by the normalized bias $b_i = (1/w_{EB,i} - 1)(\theta_i - X_i'\delta)/\sigma_i$, and the function \tilde{r} is given by the function $r(b, \chi)$ defined in (4). Our baseline implementation uses constraints on the second and fourth moments, $g(a_i) = (a_i^2, a_i^4)$.

Example 6.2 (Nonlinear soft thresholding). Consider for simplicity the homoskedastic normal model $Y_i \mid \theta_i \sim N(\theta_i, \sigma^2)$ without covariates. A popular alternative to linear estimators is the soft thresholding estimator $\hat{\theta}_{ST,i} = \text{sign}(Y_i) \max\{|Y_i| - \sqrt{2/\mu_2}, 0\}$ (e.g. Abadie and Kasy, 2019). It equals the posterior mode corresponding to a baseline Laplace prior with second moment μ_2 , which has density $\pi_0(\theta) = \frac{1}{\sqrt{2\mu_2}} \exp(-|\theta|\sqrt{2/\mu_2})$ (Johnstone, 2019, Example 2.5). To construct a robust EBCI that always contains the soft thresholding estimator, we calibrate the corresponding highest posterior density set:

$$\mathcal{S}(Y_i; \chi) = \left\{ t \in \mathbb{R}: \log \frac{\sigma^{-1}\phi((Y_i - t)/\sigma)\pi_0(t)}{\int_{-\infty}^{\infty} \sigma^{-1}\phi((Y_i - \tilde{\theta})/\sigma)\pi_0(\tilde{\theta}) d\tilde{\theta}} + \chi \geq 0 \right\}, \quad (21)$$

where ϕ is the standard normal density. This set is available in closed form and takes the form of an interval (see Supplemental Appendix G.1). Here $a_i = \theta_i$, and the function $\tilde{r}(a, \chi)$ in (18) can be computed via numerical integration.

In Supplemental Appendix G.1, we show that the resulting robust EBCI that imposes

¹²The moment functions g need not be simple moments, and could incorporate constraints used for selection of hyper-parameters, such as constraints on the marginal data distribution or, if an unbiased risk criterion is used, the constraint that the derivative of the risk equals zero at the selected prior hyper-parameters.

the constraint $E[\theta_i^2] = \mu_2$ not only has robust EB coverage (by definition), it also achieves substantial expected length improvements when the θ_i 's are in fact Laplace distributed. For $\alpha = 0.05$ and $\mu_2/\sigma^2 \leq 0.2$, the expected length under the Laplace distribution of the soft thresholding EBCI is at least 49% smaller than the length of the unshrunk CI. This exceeds the length reduction achieved by the linear robust EBCI shown in Figure 3.

Example 6.3 (Poisson shrinkage). Supplemental Appendix G.2 constructs a robust EBCI for the rate parameter θ_i in a Poisson model $Y_i \mid \theta_i \sim \text{Poisson}(\theta_i)$. This example demonstrates that our general approach does not require normality of the data.

Example 6.4 (Linear estimators in other settings). While our focus has been on EB shrinkage, our approach applies to other settings in which an estimator $\hat{\theta}_i$ is approximately normally distributed with non-negligible bias. In particular, suppose $(\hat{\theta}_i - \theta_i)/\text{se}_i$ is distributed $N(a_i, 1)$, where se_i is the standard deviation of the estimate $\hat{\theta}_i$, which for simplicity we take to be known. This holds whenever $\hat{\theta}_i$ is a linear function of jointly normal observations W_1, \dots, W_N , i.e., $\hat{\theta}_i = \sum_{j=1}^N k_{ij} W_j$ for some deterministic weights k_{ij} . Examples include series, kernel, or local polynomial estimators in a nonparametric regression with fixed covariates and normal errors. We can construct a confidence interval for θ_i as $\hat{\theta}_i \pm \chi \cdot \text{se}_i$, in which case Eq. (18) holds with $\tilde{r} = r$ given in Eq. (4). It follows from Theorem C.1 in Appendix C that if the moment constraints m on the normalized bias in Eq. (19) are replaced by consistent estimates, the resulting robust EBCI will satisfy the average coverage property (15) in large samples. We leave a full treatment of these applications for future research.

6.2 Coverage after selection

In some applications, researchers may be primarily interested in the θ_i parameters corresponding to those units i whose initial Y_i estimates fall in a given interval $[\iota_1, \iota_2]$, where $-\infty \leq \iota_1 < \iota_2 \leq \infty$. For example, in a teacher value added application, we may only be interested in the ability θ_i of those teachers i whose fixed effect estimates Y_i are positive, corresponding to setting $\iota_1 = 0$ and $\iota_2 = \infty$. Because of the selection on outcomes, naively applying our baseline EBCI procedure to the selected sample $\{i: Y_i \in [\iota_1, \iota_2]\}$ does not yield the desired average coverage across the selected units i (the same issue arises with classical CIs; see Benjamini and Yekutieli, 2005; Lee et al., 2016; Andrews et al., 2021). We now show how to correct for the selection bias in the simple homoskedastic model $Y_i \mid \theta_i \sim N(\theta_i, \sigma^2)$ without covariates from Section 2 (reintroducing the extra model features in Section 3.1 only complicates notation).

We seek a critical value χ such that the average coverage of the CI $[\hat{\theta}_i \pm \chi w_{EB} \sigma]$ is at

least $1 - \alpha$ *conditional* on the sample selection, i.e.,

$$P(\theta_i \in \hat{\theta}_i \pm \chi w_{EB} \sigma \mid Y_i \in [\iota_1, \iota_2]) \geq 1 - \alpha$$

under repeated sampling of (Y_i, θ_i) , regardless of the distribution for θ_i . Straightforward calculations show that the conditional (on θ_i and on selection) non-coverage equals

$$\begin{aligned} \tilde{r}_{\iota_1, \iota_2}(\theta_i, \chi) &= P(\theta_i \notin \hat{\theta}_i \pm \chi w_{EB} \sigma \mid Y_i \in [\iota_1, \iota_2], \theta_i) \\ &= \min \left\{ 1 - \frac{\Phi(\min\{\chi - b_i, (\iota_2 - \theta_i)/\sigma\}) - \Phi(\max\{-\chi - b_i, (\iota_1 - \theta_i)/\sigma\})}{\Phi((\iota_2 - \theta_i)/\sigma) - \Phi((\iota_1 - \theta_i)/\sigma)}, 1 \right\}, \end{aligned}$$

where $b_i = (1 - 1/w_{EB})\theta_i/\sigma$ as in Section 2. Among all distributions for θ_i consistent with the conditional moment $\tilde{\mu}_{2, \iota_1, \iota_2} = E[\theta_i^2 \mid Y_i \in [\iota_1, \iota_2]]$, the worst-case non-coverage probability, conditional on selection, is given by

$$\tilde{\rho}_{\iota_1, \iota_2}(\tilde{\mu}_{2, \iota_1, \iota_2}, \chi) \equiv \sup_F E_F[\tilde{r}_{\iota_1, \iota_2}(\theta_i, \chi)] \quad \text{s.t.} \quad E_F[\theta_i^2] = \tilde{\mu}_{2, \iota_1, \iota_2},$$

where E_F denotes expectation under $\theta_i \sim F$. This is an infinite-dimensional linear program that can be solved numerically to a high degree of accuracy, cf. Appendix B. To achieve robust conditional coverage, we solve numerically for the χ such that $\tilde{\rho}_{\iota_1, \iota_2}(\tilde{\mu}_{2, \iota_1, \iota_2}, \chi) = \alpha$.

We can estimate the conditional second moment $\tilde{\mu}_{2, \iota_1, \iota_2}$ as follows. Denote the log marginal density of Y_i by $\ell(y) \equiv \log \int \phi(y - \theta) d\Gamma_0(\theta)$, where Γ_0 is the true distribution of θ_i . Tweedie's formulas (e.g. Efron, 2019, Eq. (26)) imply

$$\tilde{\mu}_{2, \iota_1, \iota_2} = E[\theta_i^2 \mid Y_i \in [\iota_1, \iota_2]] = 1 + E[(Y_i + \ell'(Y_i))^2 + \ell''(Y_i) \mid Y_i \in [\iota_1, \iota_2]]. \quad (22)$$

Let $\hat{\ell}(y)$ be a kernel estimate of the log marginal density function of the data Y_1, \dots, Y_n . Then the estimate

$$\hat{\mu}_{2, \iota_1, \iota_2} \equiv 1 + \frac{\sum_{i: Y_i \in [\iota_1, \iota_2]} \{(Y_i + \hat{\ell}'(Y_i))^2 + \hat{\ell}''(Y_i)\}}{\#\{i: Y_i \in [\iota_1, \iota_2]\}}$$

will be consistent as $n \rightarrow \infty$ for $\tilde{\mu}_{2, \iota_1, \iota_2}$ in (22) under mild regularity conditions.

7 Empirical application

We illustrate our methods using the data and model in Chetty and Hendren (2018), who are interested in the effect of neighborhoods on intergenerational mobility.

7.1 Framework

We adopt the main specification of [Chetty and Hendren \(2018\)](#), which focuses on two definitions of a “neighborhood effect” θ_i . The first defines it as the effect of spending an additional year of childhood in commuting zone (CZ) i on children’s rank in the income distribution at age 26, for children with parents at the 25th percentile of the national income distribution. The second definition is analogous, but for children with parents at the 75th percentile. Using de-identified tax returns for all children born between 1980 and 1986 who move across CZs exactly once as children, [Chetty and Hendren \(2018\)](#) exploit variation in the age at which children move between CZs to obtain preliminary fixed effect estimates Y_i of θ_i .

Since these preliminary estimates are measured with noise, to predict θ_i , [Chetty and Hendren \(2018\)](#) shrink Y_i towards average outcomes of permanent residents of CZ i (children with parents at the same percentile of the income distribution who spent all of their childhood in the CZ). To give a sense of the accuracy of these forecasts, [Chetty and Hendren \(2018\)](#) report estimates of their unconditional MSE (i.e. treating θ_i as random), under the implicit assumption that the moment independence assumption in Eq. (10) holds. Here we complement their analysis by constructing robust EBCIs associated with these forecasts.

Our sample consists of 595 U.S. CZs, with population over 25,000 in the 2000 census, which is the set of CZs for which [Chetty and Hendren \(2018\)](#) report baseline fixed effect estimates Y_i of the effects θ_i . These baseline estimates are normalized so that their population-weighted mean is zero. Thus, we may interpret the effects θ_i as being relative to an “average” CZ. We follow the baseline implementation from Section 3.2 with standard errors $\hat{\sigma}_i$ reported by [Chetty and Hendren \(2018\)](#), and covariates X_i corresponding to a constant and the average outcomes for permanent residents. In line with the original analysis, we use precision weights $\omega_i = 1/\hat{\sigma}_i^2$ when constructing the estimates $\hat{\delta}$, $\hat{\mu}_2$ and $\hat{\kappa}$.

7.2 Results

Columns (1) and (2) in Table 3 summarize the main estimation and efficiency results. The shrinkage magnitude and relative efficiency results are similar for children with parents at the 25th and 75th percentiles of the income distribution. In both columns, the estimate of the kurtosis κ is large enough so that it does not affect the critical values or the form of the optimal shrinkage: specifications that only impose constraints on the second moment yield identical results.¹³ In line with this finding, a density plot of the t -statistics (reported as Figure S2 in [Armstrong et al. \(2020\)](#)) exhibits a fat lower tail. As a robustness check,

¹³The truncation in the $\hat{\kappa}$ formula in our baseline algorithm in Section 3.2 binds in columns (1) and (2), although the non-truncated estimates 345.3 and 5024.9 are similarly large; using these non-truncated estimates yields identical results.

Table 3: Statistics for 90% EBCIs for neighborhood effects.

Percentile	Baseline		Nonparametric	
	(1)	(2)	(3)	(4)
	25th	75th	25th	75th
Panel A: Summary statistics				
$E[\sqrt{\mu_{2,i}}]$	0.079	0.044	0.076	0.042
$E[\kappa_i]$	778.5	5948.6	1624.9	43009.9
$E[\mu_{2i}/\sigma_i^2]$	0.142	0.040	0.139	0.072
$\hat{\delta}_{\text{intercept}}$	-1.441	-2.162	-1.441	-2.162
$\hat{\delta}_{\text{perm. resident}}$	0.032	0.038	0.032	0.038
$E[w_{EB,i}]$	0.093	0.033	0.093	0.033
$E[w_{opt,i}]$	0.191	0.100	0.191	0.100
$E[\text{non-cov of parametric EBCI}_i]$	0.227	0.278	0.210	0.292
Panel B: $E[\text{half-length}_i]$				
Robust EBCI	0.195	0.122	0.186	0.116
Optimal robust EBCI	0.149	0.090	0.145	0.094
Parametric EBCI	0.123	0.070	0.123	0.070
Unshrunk CI	0.786	0.993	0.786	0.993
Panel C: Efficiency relative to robust EBCI				
Optimal robust EBCI	1.312	1.352	1.289	1.238
Parametric EBCI	1.582	1.731	1.509	1.648
Unshrunk CI	0.248	0.123	0.237	0.117

Notes: Columns (1) and (2) correspond to shrinking Y_i as in the baseline implementation that imposes Eq. (10), so that $\mu_{2i} = E[(\theta_i - X_i'\delta)^2 \mid X_i, \sigma_i]$ and $\kappa_i = E[(\theta_i - X_i'\delta)^4 \mid X_i, \sigma_i]/\mu_{2i}^2$ do not vary with i . Columns (3) and (4) use nonparametric estimates of μ_{2i} and κ_i , using the nearest neighbor estimator described in Appendix A.1. The number of nearest neighbors $J = 422$ (Column (3)) and $J = 525$ (Column (4)) is selected using cross-validation. For all columns, $\hat{\delta} = (\hat{\delta}_{\text{intercept}}, \hat{\delta}_{\text{perm. resident}})$ is computed by regressing Y_i onto a constant and outcomes for permanent residents. “Optimal Robust EBCI” refers to a robust EBCI based on length-optimal shrinkage $w_{opt,i}$, described in Section 4.2. “ $E[\text{non-cov of parametric EBCI}_i]$ ”: average of maximal non-coverage probability of parametric EBCI, given the estimated moments.

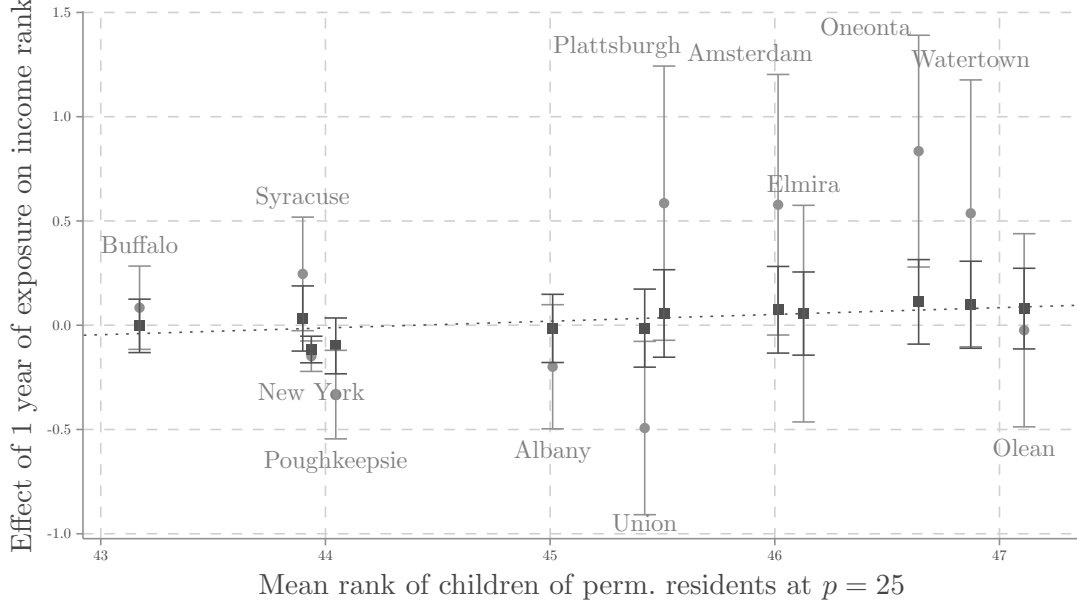


Figure 6: Neighborhood effects for New York and 90% robust EBCIs for children with parents at the $p = 25$ percentile of the national income distribution, plotted against mean outcomes of permanent residents. Gray lines correspond to CIs based on unshrunk estimates represented by circles, and black lines correspond to robust EBCIs based on EB estimates represented by squares that shrink towards a dotted regression line based on permanent residents’ outcomes. Baseline implementation as in Section 3.2.

columns (3) and (4) report results based on nonparametric moment estimates (see Remark 3.2 and Appendix A.1): the results are very similar. Indeed, the R^2 gain in predicting $\hat{\varepsilon}_i^2 - \hat{\sigma}_i^2$ using $\hat{\mu}_{2i}$ is less than 0.001 in both specifications, indicating that there is little evidence in the data against the moment independence assumption.

The baseline robust 90% EBCIs are 75.2–87.7% shorter than the usual unshrunk CIs $Y_i \pm z_{1-\alpha/2}\hat{\sigma}_i$. To interpret these gains in dollar terms, for children with parents at the 25th percentile of the income distribution, a percentile gain corresponds to an annual income gain of \$818 (Chetty and Hendren, 2018, p. 1183). Thus, the average half-length of the baseline robust EBCIs in column (1) implies CIs of the form $\pm\$160$ on average, while the unshrunk CIs are of the form $\pm\$643$ on average. These large gains are a consequence of a low signal-to-noise ratio μ_2/σ_i^2 in this application. Because the shrinkage magnitude is so large on average, the tail behavior of the bias matters, and since the kurtosis estimates suggests these tails are fat, it is important to use the robust critical value: the parametric EBCI exhibits average potential size distortions of 12.7–17.8 percentage points. Indeed, for over 90% of the CIs in the specifications in columns (1) and (2), the shrinkage coefficient $w_{EB,i}$ falls below the “rule of thumb” threshold of 0.3 derived in Section 4.3.

To visualize these results, Figure 6 plots the unshrunk 90% CIs based on the preliminary estimates, as well as robust EBCIs based on EB estimates for cities in the state of New York for children with parents at the 25th percentile. While the EBCIs for large CZs like New York City or Buffalo are similar to the unshrunk CIs, they are much tighter for smaller CZs like Plattsburgh or Watertown, with point estimates that shrink the preliminary estimates Y_i most of the way toward the regression line $X_i'\hat{\delta}$.

In summary, using shrinkage allows us to considerably tighten the CIs based on preliminary estimates. This is true even though the CIs only effectively use second moment constraints—imposing constraints on the kurtosis does not affect the critical values in this application.

Appendix A Moment estimates

The EBCI in our baseline implementation has valid EB coverage asymptotically as $n \rightarrow \infty$, so long as the estimates $\hat{\mu}_2$ and $\hat{\kappa}$ are consistent. While the particular choice of the estimates $\hat{\mu}_2$ and $\hat{\kappa}$ does not affect the CI asymptotically, finite sample considerations can be important for small to moderate values of n . In particular, unrestricted moment-based estimates of μ_2 and κ may fall below their theoretical lower bounds of 0 and 1, in which case it is not clear how to define the EBCI.¹⁴ To address this issue, in analogy to finite-sample corrections to parametric EBCIs proposed in Morris (1983a,b), Appendix A.1 derives two finite-sample corrections to the unrestricted estimates that approximate a Bayesian estimate under a flat hyperprior on (μ_2, κ) . We verify that these corrections give good coverage in an extensive set of Monte Carlo designs in Section 4.4. We also discuss implementation of nonparametric moment estimates. Appendix A.2 discusses the choice of weights ω_i .

A.1 Finite n corrections and nonparametric moment estimates

To derive our estimates of μ_2 and κ , we first consider unrestricted estimation under the moment independence condition (10). For μ_2 , this condition implies the moment condition $E[(Y_i - X_i'\delta)^2 - \sigma_i^2 \mid X_i, \sigma_i] = \mu_2$. Replacing $Y_i - X_i'\delta$ with the residual $\hat{\varepsilon}_i = Y_i - X_i'\hat{\delta}$ yields the estimate

$$\hat{\mu}_{2,UC} = \frac{\sum_{i=1}^n \omega_i \mathcal{W}_{2i}}{\sum_{i=1}^n \omega_i}, \quad \mathcal{W}_{2i} = \hat{\varepsilon}_i^2 - \hat{\sigma}_i^2, \quad (23)$$

for any weights $\omega_i = \omega_i(X_i, \hat{\sigma}_i)$. Here, UC stands for “unconstrained,” since the estimate $\hat{\mu}_{2,UC}$ can be negative. To incorporate the constraint $\mu_2 > 0$, we use an approximation

¹⁴Formally, our results are asymptotic and require $\mu_2 > 0$ and $\kappa > 1$, so that these issues do not occur when n is large enough. We discuss the difficulty of providing finite-sample coverage guarantees in Section 5.

to a Bayesian approach with a flat prior on the set $[0, \infty)$. A full Bayesian approach to estimating μ_2 would place a hyperprior on possible joint distributions of X_i, σ_i, θ_i , which could potentially lead to using complicated functions of the data to estimate μ_2 . For simplicity, we compute the posterior mean given $\hat{\mu}_{2,\text{UC}}$, and we use a normal approximation to the likelihood. Since the posterior distribution only uses knowledge of $\hat{\mu}_{2,\text{UC}}$, we refer to this as a flat prior limited information Bayes (FPLIB) approach.

To derive this formula, first note that, if \hat{m} is an estimate of a parameter m with $\hat{m} \mid m \sim N(m, V)$, then under a flat prior for m on $[0, \infty)$, the posterior mean of m is given by

$$b(\hat{m}, V) = \hat{m} + \sqrt{V} \phi(\hat{m}/\sqrt{V}) / \Phi(\hat{m}/\sqrt{V}),$$

where ϕ and Φ are the standard normal pdf and cdf respectively. Furthermore, if $\hat{m} = \sum_{i=1}^n \omega_i Z_i / \sum_{i=1}^n \omega_i$ where the Z_i 's are independent with mean m conditional on the weights $\omega = (\omega_1, \dots, \omega_n)'$, then an unbiased estimate of the variance of \hat{m} given ω is given by

$$V(Z, \omega) = \frac{\sum_{i=1}^n \omega_i^2 (Z_i^2 - \hat{m}^2)}{(\sum_{i=1}^n \omega_i)^2 - \sum_{i=1}^n \omega_i^2}.$$

Conditioning on the X_i 's and σ_i 's (and ignoring sampling variation in $\hat{\delta}$ and the $\hat{\sigma}_i$'s), we can then apply this formula to $\hat{\mu}_{2,\text{UC}}$, with $Z_i = \mathcal{W}_{2i}$, where \mathcal{W}_{2i} is given in (23). This gives the FPLIB estimate for μ_2 :

$$\hat{\mu}_{2,\text{FPLIB}} = b(\hat{\mu}_{2,\text{UC}}, V(\mathcal{W}_2, \omega)).$$

To derive the FPLIB estimate for κ , we begin with an unconstrained estimate of $\mu_4 = E[(\theta_i - X_i' \delta)^4]$. The moment independence condition (10) delivers the moment condition $\mu_4 = E[(Y_i - X_i' \delta)^4 + 3\sigma_i^4 - 6\sigma_i^2(Y_i - X_i' \delta)^2 \mid X_i, \sigma_i]$, which leads to the unconstrained estimate

$$\hat{\mu}_{4,\text{UC}} = \frac{\sum_{i=1}^n \omega_i \mathcal{W}_{4i}}{\sum_{i=1}^n \omega_i}, \quad \mathcal{W}_{4i} = \hat{\varepsilon}_i^4 - 6\hat{\sigma}_i^2 \hat{\varepsilon}_i^2 + 3\hat{\sigma}_i^4.$$

To avoid issues with small values of estimates of μ_2 in the denominator, we apply the FPLIB approach to an estimate of $\mu_4 - \mu_2^2$, using a flat prior on the parameter space $[0, \infty)$. Using the delta method leads to approximating the variance of $\hat{\mu}_{4,\text{UC}} - \hat{\mu}_{2,\text{UC}}^2$ with the variance of $\sum_{i=1}^n \omega_i (\mathcal{W}_{4i} - 2\mu_2 \mathcal{W}_{2i}) / \sum_{i=1}^n \omega_i$, so that the FPLIB estimate of $\mu_4 - \mu_2^2$ is $b(\hat{\mu}_{4,\text{UC}} - \hat{\mu}_{2,\text{UC}}^2, V(\mathcal{W}_4 - 2\hat{\mu}_{2,\text{FPLIB}} \mathcal{W}_2, \omega))$, and the FPLIB estimate of κ is

$$\hat{\kappa}_{\text{FPLIB}} = 1 + \frac{b(\hat{\mu}_{4,\text{UC}} - \hat{\mu}_{2,\text{UC}}^2, V(\mathcal{W}_4 - 2\hat{\mu}_{2,\text{FPLIB}} \mathcal{W}_2, \omega))}{\hat{\mu}_{2,\text{FPLIB}}^2}.$$

As a further simplification, we derive approximations in which the posterior mean formula $b(\hat{m}, V)$ is replaced by a simple truncation formula. We refer to this approach as posterior mean trimming (PMT). In particular, suppose we apply the formula $b(\hat{m}, V)$ to an estimator \hat{m} such that $\hat{m} \geq m_0$ and $V \geq V_0$ by construction, where $m_0 < 0$. Then the posterior mean satisfies $b(\hat{m}, V) \geq b(m_0, V_0)$ (Pinelis, 2002, Proposition 1.2). Thus, a simple approximation to the FPLIB estimator is to truncate \hat{m} from below at $b(m_0, V_0)$. To obtain an even simpler formula, we use the approximation $b(m_0, V_0) = -V_0/m_0 + O(V_0^{3/2})$ (Pinelis, 2002, Proposition 1.3), which holds as $V_0 \rightarrow 0$ (or, equivalently, as $n \rightarrow \infty$, provided the estimator \hat{m} is consistent). The variance of $\hat{\mu}_{2,UC}$ conditional on (X_i, σ_i) is bounded below by $2 \sum_{i=1}^n \omega_i^2 \sigma_i^4 / (\sum_{i=1}^n \omega_i)^2$, and $\hat{\mu}_{2,UC} \geq -\sum_{i=1}^n \omega_i \sigma_i^2 / \sum_{i=1}^n \omega_i$, so we can use $V_0/m_0 = -\frac{2 \sum_{i=1}^n \omega_i^2 \sigma_i^4}{\sum_{i=1}^n \omega_i \sigma_i^2 \cdot \sum_{i=1}^n \omega_i}$, which gives the PMT estimator

$$\hat{\mu}_{2,PMT} = \max \left\{ \hat{\mu}_{2,UC}, \frac{2 \sum_{i=1}^n \omega_i^2 \sigma_i^4}{\sum_{i=1}^n \omega_i \sigma_i^2 \cdot \sum_{i=1}^n \omega_i} \right\}.$$

For κ , we simplify our approach to deriving a trimming rule by treating μ_2 as known, and considering the variance of the infeasible estimate $\hat{\kappa}_{UC}^* = \frac{\sum_{i=1}^n \omega_i (\hat{\epsilon}_i^4 - 6\hat{\sigma}_i^2 \mu_2 - 3\hat{\sigma}_i^4)}{\mu_2^2 \sum_{i=1}^n \omega_i}$. Using the above truncation formula for $\hat{\kappa}_{UC}^* - 1$ along with the fact that $\hat{\kappa}_{UC}^* \geq \frac{\sum_{i=1}^n \omega_i (-6\hat{\sigma}_i^2 \mu_2 - 3\hat{\sigma}_i^4)}{\mu_2^2 \sum_{i=1}^n \omega_i}$ and the lower bound $8 \sum_i \omega_i^2 (2\mu_2^3 \sigma_i^2 + 21\mu_2^2 \sigma_i^4 + 48\mu_2 \sigma_i^6 + 12\sigma_i^8) / \mu_2^4 (\sum_i \omega_i)^2$ on the variance yields $V_0/m_0 = -\frac{8 \sum_i \omega_i^2 (2\mu_2^3 \sigma_i^2 + 21\mu_2^2 \sigma_i^4 + 48\mu_2 \sigma_i^6 + 12\sigma_i^8)}{\mu_2^2 (\sum_i \omega_i) \sum_{i=1}^n \omega_i (\mu_2^2 + 6\hat{\sigma}_i^2 \mu_2 + 3\hat{\sigma}_i^4)}$. To simplify the trimming rule even further, we only use the leading term of V_0/m_0 as $\mu_2 \rightarrow 0$, $V_0/m_0 = -\frac{32 \sum_i \omega_i^2 \sigma_i^8}{\mu_2^2 (\sum_i \omega_i) \sum_{i=1}^n \omega_i \hat{\sigma}_i^4} + o(1/\mu_2^2)$. Plugging in $\hat{\mu}_{2,PMT}$ in place of the unknown μ_2 then gives the PMT estimator

$$\hat{\kappa}_{PMT} = \max \left\{ \frac{\hat{\mu}_{4,UC}}{\hat{\mu}_{2,PMT}^2}, 1 + \frac{32 \sum_{i=1}^n \omega_i^2 \hat{\sigma}_i^8}{\hat{\mu}_{2,PMT}^2 \sum_{i=1}^n \omega_i \cdot \sum_{i=1}^n \omega_i \hat{\sigma}_i^4} \right\}.$$

The estimators in step 1 of our baseline implementation in Section 3.2 correspond to $\hat{\mu}_{2,PMT}$ and $\hat{\kappa}_{PMT}$, due to their slightly simpler form relative to the FPLIB estimators. In unreported simulations based on the designs described in Section 4.4 and Supplemental Appendix E.2, we find that EBCIs based on FPLIB lead to even smaller finite-sample coverage distortions than those based on the baseline implementation that uses PMT, at the expense of slightly longer average length.

To implement the nonparametric estimates $\hat{\kappa}_i$ and $\hat{\mu}_{2i}$ in Remark 3.2, we use the nearest-neighbor estimator that for each i computes the PMT estimates $\hat{\mu}_{2,PMT}$ and $\hat{\kappa}_{PMT}$ described above, using only the J observations closest to i , rather than the full sample of n observations. We define distance as a Euclidean distance on (X_i, σ_i) , after scaling elements of this vector by their standard deviations. Under regularity conditions, the resulting estimates will be

consistent for μ_{2i} and κ_i , so long as $J \rightarrow \infty$ and $J/n \rightarrow 0$. We select J using leave-one-out cross-validation, using the squared prediction error in predicting \mathcal{W}_{2i} as the criterion. For simplicity, we use the same J for estimating the kurtosis as that used for estimating the second moment.

A.2 Choice of weighting

Under condition (10), the weights ω_i used to estimate μ_2 and κ can be any function of X_i, σ_i . Furthermore, while $\hat{\delta}$ can be essentially arbitrary as long as it converges in probability to some δ such that Eq. (10) holds, that equation will often be motivated by the assumption that the conditional mean of θ_i is linear in X_i ,

$$E[\theta_i - X_i' \delta \mid X_i, \sigma_i] = 0. \quad (24)$$

Under this condition, the weights ω_i used to estimate δ can also be any function of X_i, σ_i .

Thus, under conditions (10) and (24), the choice of weighting can be guided by efficiency considerations. In general, the optimal weights are different for each of the three estimates of δ, μ_2 , and κ , and implementing them requires first stage estimates of the variances of Y_i, \mathcal{W}_{2i} and \mathcal{W}_{4i} , conditional on (X_i, σ_i) (with \mathcal{W}_{2i} and \mathcal{W}_{4i} defined in Appendix A.1). To avoid estimation of these variances, consider the limiting case where the signal-to-noise ratio goes to 0, i.e. $\mu_2 / \min_i \sigma_i^2 \rightarrow 0$. The resulting weights will be near-optimal under a low signal-to-noise ratio, when precise estimation of these parameters is relatively more important for accurate coverage (under a high signal-to-noise ratio, shrinkage is limited, and estimation error in these parameters has little effect on coverage). Let us also ignore estimation error in δ for simplicity, and suppose that the Y_i 's are independent conditional on $(\theta_i, X_i, \sigma_i)$. Then, as $\mu_2 / \min_i \sigma_i^2 \rightarrow 0$, the weights $\hat{\sigma}_i^{-2}, \hat{\sigma}_i^{-4}$, and $\hat{\sigma}_i^{-8}$, for estimating δ, μ_2 , and μ_4 , respectively, become optimal. For simplicity, the baseline implementation in Section 3.2 uses the same weights ω_i for each of the estimates; the choice $\omega_i = \hat{\sigma}_i^{-2}$ targets optimal estimation of δ . However, one could relax this constraint, and use the weights $\hat{\sigma}_i^{-4}$, and $\hat{\sigma}_i^{-8}$ for estimating μ_2 and μ_4 instead. The choice $\omega_i = 1/n$ has the advantage of simplicity; one may also motivate it by robustness concerns when Eq. (10) fails, though our preferred robustness check is to use nonparametric moment estimates, as outlined in Remark 3.2.

Appendix B Computational details

To simplify the statement of the results below, let $r_0(b, \chi) = r(\sqrt{b}, \chi)$, and put $m_2 = \sigma^2 / \mu_2$. The next proposition shows that, if only a second moment constraint is imposed, the maximal

non-coverage probability $\rho(m_2, \chi)$, defined in Eq. (5), has a simple solution:

Proposition B.1. *Consider the problem in Eq. (5). The solution is given by*

$$\rho(m_2, \chi) = \begin{cases} r_0(0, \chi) + \frac{m_2}{t_0}(r_0(t_0, \chi) - r_0(0, \chi)) & \text{if } m_2 < t_0, \\ r_0(m_2, \chi) & \text{otherwise.} \end{cases}$$

Here $t_0 = 0$ if $\chi < \sqrt{3}$, otherwise t_0 is the unique solution to $r_0(t, \chi) + u \frac{\partial}{\partial u} r_0(u, \chi) = r_0(u, \chi)$.

The proof of Proposition B.1 shows that $\rho(m_2, \chi)$ corresponds to the least concave majorant of the function r_0 .

The next result shows that, if in addition to a second moment constraint, we impose a constraint on the kurtosis, the maximal non-coverage probability can be computed as a solution to two nested univariate optimizations:

Proposition B.2. *Suppose $\kappa > 1$ and $m_2 > 0$. Then the solution to the problem*

$$\rho(m_2, \kappa, \chi) = \sup_F E_F[r(b, \chi)] \quad \text{s.t.} \quad E_F[b^2] = m_2, E_F[b^4] = \kappa m_2^2,$$

is given by $\rho(m_2, \kappa, \chi) = r_0(m_2, \chi)$ if $m_2 \geq t_0$, with t_0 defined in Proposition B.1. If $m_2 < t_0$, then the solution is given by

$$\inf_{0 < x_0 \leq t_0} \left\{ r_0(x_0, \chi) + (m_2 - x_0)r'_0(x_0, \chi) + ((x_0 - m_2)^2 + (\kappa - 1)m_2^2) \sup_{0 \leq x \leq t_0} \delta(x; x_0) \right\}, \quad (25)$$

where $r'_0(x_0, \chi) = \partial r_0(x_0, \chi) / \partial x_0$, $\delta(x; x_0) = \frac{r_0(x, \chi) - r_0(x_0, \chi) - (x - x_0)r'_0(x_0, \chi)}{(x - x_0)^2}$ if $x \neq x_0$, and $\delta(x_0; x_0) = \lim_{x \rightarrow x_0} \delta(x; x_0) = \frac{1}{2} \frac{\partial^2}{\partial x_0^2} r_0(x_0, \chi)$.

If $m_2 \geq t_0$, then imposing a constraint on the kurtosis does not help to reduce the maximal non-coverage probability, and $\rho(m_2, \kappa, \chi) = \rho(m_2, \chi)$.

Remark B.1 (Least favorable distributions). It follows from the proof of these propositions that distributions maximizing Eq. (5)—the least favorable distributions for the normalized bias b —have two support points if $m_2 \geq t_0$, namely $-\sqrt{m_2}$ and $\sqrt{m_2}$ (since the rejection probability $r(b, \chi)$ depends on b only through its absolute value, any distribution with these two support points maximizes Eq. (5)). If $m_2 < t_0$, there are three support points, $b = 0$, with probability $1 - m_2/t_0$ and $b = \pm\sqrt{t_0}$ with total probability m_2/t_0 (again, only the sum of the probabilities is uniquely determined). If the kurtosis constraint is also imposed, then there are four support points, $\pm\sqrt{x_0}$ and $\pm\sqrt{x}$, where x and x_0 optimize Eq. (25).

Finally, the characterization of the solution to the general program in Eq. (19) depends on the form of the constraint function g . To solve the program numerically, discretize the support of F to turn the problem into a finite-dimensional linear program, which can be solved using a standard linear solver. In particular, we solve the problem

$$\rho_g(m, \chi) = \sup_{p_1, \dots, p_K} \sum_{k=1}^K p_k r(x_k, \chi) \quad \text{s.t.} \quad \sum_{k=1}^K p_k g(x_k) = m, \quad \sum_{k=1}^K p_k = 1, \quad p_k \geq 0.$$

Here x_1, \dots, x_K denote the support points of b , with p_k denoting the associated probabilities.

Appendix C Coverage results

This Appendix provides coverage results that generalize Theorem 4.1. Appendix C.1 introduces the general setup. Appendix C.2 provides results for general shrinkage estimators that satisfy an approximate normality assumption. Appendix C.3 considers a generalization of our baseline specification in the EB setting, and states a generalization of Theorem 4.1.

C.1 General setup and notation

Let $\hat{\theta}_1, \dots, \hat{\theta}_n$ be estimates of parameters $\theta_1, \dots, \theta_n$, with standard errors $\text{se}_1, \dots, \text{se}_n$. The standard errors may be random variables that depend on the data. We are interested in coverage properties of the intervals $CI_i = \{\hat{\theta}_i \pm \text{se}_i \cdot \chi_i\}$ for some χ_1, \dots, χ_n , which may be chosen based on the data. In some cases, we will condition on a variable \tilde{X}_i when defining EB coverage or average coverage. Let $\tilde{X}^{(n)} = (\tilde{X}_1, \dots, \tilde{X}_n)'$ and let $\chi^{(n)} = (\chi_1, \dots, \chi_n)'$.

As discussed in Section 4.1, the average coverage criterion does not require thinking of θ as random. To save on notation, we will state most of our average coverage results and conditions in terms of a general sequence of probability measures $\tilde{P} = \tilde{P}^{(n)}$ and triangular arrays θ and $\tilde{X}^{(n)}$. We will use $E_{\tilde{P}}$ to denote expectation under the measure \tilde{P} . We can then obtain EB coverage statements by considering a distribution P for the data and $\theta, \tilde{X}^{(n)}$ and an additional variable ν such that these conditions hold for the measure $\tilde{P}(\cdot) = P(\cdot \mid \theta, \nu, \tilde{X}^{(n)})$ for $\theta, \nu, \tilde{X}^{(n)}$ in a probability one set. The variable ν is allowed to depend on n , and can include nuisance parameters as well as additional variables.

It will be useful to formulate a conditional version of the average coverage criterion (15), to complement the conditional version of EB coverage discussed in the main text. Due to discreteness of the empirical measure of the \tilde{X}_i 's, we consider coverage conditional on each set in some family \mathcal{A} of sets. To formalize this, let $\mathcal{I}_{\mathcal{X}, n} = \{i \in \{1, \dots, n\} : \tilde{X}_i \in \mathcal{X}\}$, and let

$N_{\mathcal{X},n} = \#\mathcal{I}_{\mathcal{X},n}$. The sample average non-coverage on the set \mathcal{X} is then given by

$$ANC_n(\chi^{(n)}; \mathcal{X}) = \frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} \mathbf{I}\{\theta_i \notin \{\hat{\theta} \pm \text{se}_i \cdot \chi_i\}\} = \frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} \mathbf{I}\{|Z_i| > \chi_i\},$$

where $Z_i = (\hat{\theta}_i - \theta_i)/\text{se}_i$. We consider the following notions of average coverage control, conditional on the set $\mathcal{X} \in \mathcal{A}$:

$$ANC_n(\chi; \mathcal{X}) \leq \alpha + o_{\tilde{P}}(1), \quad (26)$$

and

$$\limsup_n E_{\tilde{P}}[ANC_n(\chi; \mathcal{X})] = \limsup_n \frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} \tilde{P}(|Z_i| > \chi_i) \leq \alpha. \quad (27)$$

Note that (26) implies (27), since $ANC_n(\chi; \mathcal{X})$ is uniformly bounded. Furthermore, if we integrate with respect to some distribution on $\nu, \tilde{X}^{(n)}$ such that (27) holds with $\tilde{P}(\cdot) = P(\cdot \mid \theta, \nu, \tilde{X}^{(n)})$ almost surely, we get (again by uniform boundedness)

$$\limsup_n E[ANC_n(\chi; \mathcal{X}) \mid \theta] \leq \alpha,$$

which, in the case where \mathcal{X} contains all \tilde{X}_i 's with probability one, is condition (15) from the main text.

Now consider EB coverage, as defined in Eq. (14) in the main text, but conditioning on \tilde{X}_i . We consider EB coverage under a distribution P for the data, $\tilde{X}^{(n)}$, θ and ν , where ν includes additional nuisance parameters and covariates, and where the average coverage condition (27) holds with $P(\cdot \mid \theta, \nu, \tilde{X}^{(n)})$ playing the role of \tilde{P} with probability one. Consider the case where \tilde{X}_i is discretely distributed under P . Suppose that the exchangeability condition

$$P(\theta_i \in CI_i \mid \mathcal{I}_{\{\tilde{x}\},n}) = P(\theta_j \in CI_j \mid \mathcal{I}_{\{\tilde{x}\},n}) \text{ for all } i, j \in \mathcal{I}_{\{\tilde{x}\},n} \quad (28)$$

holds with probability one. Then, for each j ,

$$\begin{aligned} P(\theta_j \in CI_j \mid \tilde{X}_j = \tilde{x}) &= P(\theta_j \in CI_j \mid j \in \mathcal{I}_{\{\tilde{x}\},n}) = E[P(\theta_j \in CI_j \mid \mathcal{I}_{\{\tilde{x}\},n}) \mid j \in \mathcal{I}_{\{\tilde{x}\},n}] \\ &= E\left[\frac{1}{N_{\{\tilde{x}\},n}} \sum_{i \in \mathcal{I}_{\{\tilde{x}\}}} P(\theta_i \in CI_i \mid \mathcal{I}_{\{\tilde{x}\}}) \mid j \in \mathcal{I}_{\{\tilde{x}\},n}\right]. \end{aligned}$$

Plugging in $P(\cdot \mid \theta, \nu, \tilde{X}^{(n)})$ for \tilde{P} in the coverage condition (27), taking the expectation conditional on $\mathcal{I}_{\{\tilde{x}\},n}$ and using uniform boundedness, it follows that the \liminf of the term

in the conditional expectation is no less than $1 - \alpha$. Then, by uniform boundedness of this term,

$$\liminf_{n \rightarrow \infty} P(\theta_j \in CI_j \mid \tilde{X}_j = \tilde{x}) \geq 1 - \alpha. \quad (29)$$

This is a conditional version of the EB coverage condition (14) from the main text.

C.2 Results for general shrinkage estimators

We assume that $Z_i = (\hat{\theta}_i - \theta_i)/\text{se}_i$ is approximately normal with variance one and mean b_i under the sequence of probability measures $\tilde{P} = \tilde{P}^{(n)}$. To formalize this, we consider a triangular array of distributions satisfying the following conditions.

Assumption C.1. *For some random variables \tilde{b}_i and constants $b_{i,n}$, $Z_i - \tilde{b}_i$ satisfies*

$$\lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} \left| \tilde{P}(Z_i - \tilde{b}_i \leq t) - \Phi(t) \right| = 0$$

for all $t \in \mathbb{R}$ and, for all $\mathcal{X} \in \mathcal{A}$ and any $\varepsilon > 0$, $\frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} \tilde{P}(|\tilde{b}_i - b_{i,n}| \geq \varepsilon) \rightarrow 0$.

Note that, when applying the results with $\tilde{P}(\cdot)$ given by the sequence of measures $P(\cdot \mid \theta, \nu, \tilde{X}^{(n)})$, the constants $b_{i,n}$ will be allowed to depend on $\theta, \nu, \tilde{X}^{(n)}$.

Let $g: \mathbb{R} \rightarrow \mathbb{R}^p$ be a vector of moment functions. We consider critical values $\hat{\chi}^{(n)} = (\hat{\chi}_1, \dots, \hat{\chi}_n)$ based on an estimate of the conditional expectation of $g(b_{i,n})$ given \tilde{X}_i , where the expectation is taken with respect to the empirical distribution of $\tilde{X}_i, b_{i,n}$. Due to the discreteness of this measure, we consider the behavior of this estimate on average over sets $\mathcal{X} \in \mathcal{A}$. We assume that there exists a function $m: \mathcal{X} \rightarrow \mathbb{R}^p$ that plays the role of the conditional expectation of $g(b_{i,n})$ given \tilde{X}_i , along with estimates \hat{m}_i of $m(\tilde{X}_i)$, which satisfy the following assumptions.

Assumption C.2. *For all $\mathcal{X} \in \mathcal{A}$, $N_{\mathcal{X},n} \rightarrow \infty$, $\frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} (g(b_{n,i}) - m(\tilde{X}_i)) \rightarrow 0$, and, for all $\varepsilon > 0$, $\frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} \tilde{P}(\|\hat{m}_i - m(\tilde{X}_i)\| \geq \varepsilon) \rightarrow 0$.*

Assumption C.3. *For every $\mathcal{X} \in \mathcal{A}$ and every $\varepsilon > 0$, there is a partition $\mathcal{X}_1, \dots, \mathcal{X}_J \in \mathcal{A}$ of \mathcal{X} and m_1, \dots, m_J such that, for each j and all $x \in \mathcal{X}_j$, $m(x) \in B_\varepsilon(m_j)$, where $B_\varepsilon(m) = \{\tilde{m} : \|\tilde{m} - m\| \leq \varepsilon\}$.*

Assumption C.4. *For some compact set M in the interior of the set of values of $\int g(b) dF(b)$ where F ranges over all probability measures on \mathbb{R} , we have $m(x) \in M$ for all x .*

Let $\rho_g(m, \chi)$ and $\text{cva}_{\alpha,g}(m)$ be defined as in Section 6,

$$\text{cva}_{\alpha,g}(m) = \inf\{\chi : \rho_g(m, \chi) \leq \alpha\} \quad \text{where} \quad \rho_g(m, \chi) = \sup_F E_F[r(b, \chi)] \text{ s.t. } E_F[g(b)] = m.$$

Let $\hat{\chi}_i = \text{cva}_{\alpha,g}(\hat{m}_i)$. We will consider the average non-coverage $ANC_n(\hat{\chi}^{(n)}; \mathcal{X})$ of the collection of intervals $\{\hat{\theta}_i \pm \text{se}_i \cdot \hat{\chi}_i\}$.

Theorem C.1. *Suppose that Assumptions C.1, C.2, C.3 and C.4 hold, and that, for some j , $\lim_{b \rightarrow \infty} g_j(b) = \lim_{b \rightarrow -\infty} g_j(b) = \infty$ and $\inf_b g_j(b) \geq 0$. Then, for all $\mathcal{X} \in \mathcal{A}$,*

$$E_{\tilde{P}} ANC_n(\hat{\chi}^{(n)}; \mathcal{X}) \leq \alpha + o(1).$$

If, in addition, $Z_i - \tilde{b}_i$ is independent over i under \tilde{P} , then $ANC_n(\hat{\chi}^{(n)}; \mathcal{X}) \leq \alpha + o_{\tilde{P}}(1)$.

C.3 Empirical Bayes shrinkage toward regression estimate

We now apply the general results in Appendix C.2 to the EB setting. As in Section 3, we consider unshrunk estimates Y_1, \dots, Y_n of parameters $\theta = (\theta_1, \dots, \theta_n)'$, along with regressors $X^{(n)} = (X_1, \dots, X_n)$ and variables $\tilde{X}^{(n)} = (\tilde{X}_1, \dots, \tilde{X}_n)'$, which include σ_i and which play the role of the conditioning variables. (While Section 3 uses X_i, σ_i as the conditioning variable \tilde{X}_i , here we generalize the results by allowing the conditioning variables to differ from X_i .) The initial estimate Y_i has standard deviation σ_i , and we observe an estimate $\hat{\sigma}_i$. We obtain average coverage results by considering a triangular array of probability distributions $\tilde{P} = \tilde{P}^{(n)}$, in which the X_i 's, σ_i 's and θ_i 's are fixed. EB coverage can then be obtained for a distribution P of the data, θ and some nuisance parameter $\tilde{\nu}$ such that these conditions hold almost surely with $P(\cdot \mid \theta, \tilde{\nu}, \tilde{X}^{(n)}, X^{(n)})$ playing the role of \tilde{P} .

We consider the following generalization of the baseline specification considered in the main text. Let

$$\hat{\theta}_i = \hat{X}_i' \hat{\delta} + w(\hat{\gamma}, \hat{\sigma}_i)(Y_i - \hat{X}_i' \hat{\delta})$$

where \hat{X}_i is an estimate of X_i (we allow for the possibility that some elements of X_i are estimated rather than observed directly, which will be the case, for example, when σ_i is included in X_i), $\hat{\delta}$ is any random vector that depends on the data (such as the OLS estimator in a regression of Y_i on X_i), and $\hat{\gamma}$ is a tuning parameter that determines shrinkage and may depend on the data. This leads to the standard error $\text{se}_i = w(\hat{\gamma}, \hat{\sigma}_i) \hat{\sigma}_i$ so that the t -statistic is

$$Z_i = \frac{\hat{\theta}_i - \theta_i}{\text{se}_i} = \frac{\hat{X}_i' \hat{\delta} + w(\hat{\gamma}, \hat{\sigma}_i)(Y_i - \hat{X}_i' \hat{\delta}) - \theta_i}{w(\hat{\gamma}, \hat{\sigma}_i) \hat{\sigma}_i} = \frac{Y_i - \theta_i}{\hat{\sigma}_i} + \frac{[w(\hat{\gamma}, \hat{\sigma}_i) - 1](\theta_i - \hat{X}_i' \hat{\delta})}{w(\hat{\gamma}, \hat{\sigma}_i) \hat{\sigma}_i}.$$

We use estimates of moments of order $\ell_1 < \dots < \ell_p$ of the bias, where $\ell_1 < \dots < \ell_p$ are positive integers. Let $\hat{\mu}_\ell$ be an estimate of the ℓ th moment of $(\theta_i - X_i' \delta)$, and suppose that this moment is independent of σ_i in a sense formalized below. Then an estimate of the ℓ_j th

moment of the bias is $\hat{m}_{i,j} = \frac{[w(\hat{\gamma}, \hat{\sigma}_i) - 1]^{\ell_j} \hat{\mu}_{\ell_j}}{w(\hat{\gamma}, \hat{\sigma}_i)^{\ell_j} \hat{\sigma}_i^{\ell_j}}$. Let $\hat{m}_i = (\hat{m}_1, \dots, \hat{m}_p)'$. The EBCI is then given by $\hat{\theta}_i \pm w(\hat{\gamma}, \hat{\sigma}_i) \hat{\sigma}_i \cdot \text{cva}_{\alpha,g}(\hat{m}_i)$ where $g_j(b) = b^{\ell_j}$. We obtain the baseline specification in Section 3.2 when $p = 2$, $\ell_1 = 2$, $\ell_2 = 4$, $\hat{\gamma} = \hat{\mu}_2$ and $w(\hat{\mu}_2, \hat{\sigma}_i) = \hat{\mu}_2/(\hat{\mu}_2 + \hat{\sigma}_i^2)$.

We make the following assumptions.

Assumption C.5.

$$\lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} \left| \tilde{P} \left(\frac{Y_i - \theta_i}{\hat{\sigma}_i} \leq t \right) - \Phi(t) \right| = 0.$$

We give primitive conditions for Assumption C.5 in Supplemental Appendix D.1, and verify them in a linear fixed effects panel data model. The primitive conditions involve considering a triangular array of parameter values such that sampling error and empirical moments of the parameter value sequence are of the same order of magnitude, and defining θ_i to be a scaled version of the corresponding parameter.

Assumption C.6. *The standard deviations σ_i are bounded away from zero. In addition, for some δ and γ , $\hat{\delta}$ and $\hat{\gamma}$ converge to δ and γ under \tilde{P} , and, for any $\varepsilon > 0$,*

$$\lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} \tilde{P}(|\hat{\sigma}_i - \sigma_i| \geq \varepsilon) = 0 \text{ and } \lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} \tilde{P}(|\hat{X}_i - X_i| \geq \varepsilon) = 0.$$

Assumption C.7. *The variable \tilde{X}_i takes values in $\mathcal{S}_1 \times \dots \times \mathcal{S}_s$ where, for each k , either $\mathcal{S}_k = [\underline{x}_k, \bar{x}_k]$ (with $-\infty < \underline{x}_k < \bar{x}_k < \infty$) or \mathcal{S}_k is a finitely discrete set with minimum element \underline{x}_k and maximum element \bar{x}_k . In addition, $\tilde{X}_{i1} = \sigma_i$ (the first element of \tilde{X}_i is given by σ_i). Furthermore, for some μ_0 such that $(\mu_{0,\ell_1}, \dots, \mu_{0,\ell_p})$ is in the interior of the set of values of $\int g(b) dF(b)$ where F ranges over probability measures on \mathbb{R} where $g_j(b) = b^{\ell_j}$ and some constant K , the following holds. Let \mathcal{A} denote the collection of sets $\tilde{\mathcal{S}}_1 \times \dots \times \tilde{\mathcal{S}}_s$ where $\tilde{\mathcal{S}}_k$ is a positive Lebesgue measure interval contained in $[\underline{x}_k, \bar{x}_k]$ in the case where $\mathcal{S}_k = [\underline{x}_k, \bar{x}_k]$, and $\tilde{\mathcal{S}}_k$ is a nonempty subset of \mathcal{S}_k in the case where \mathcal{S}_k is finitely discrete. For any $\mathcal{X} \in \mathcal{A}$, $N_{\mathcal{X},n} \rightarrow \infty$ and*

$$\frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} (\theta_i - X_i' \delta)^{\ell_j} \rightarrow \mu_{0,\ell_j}, \quad \frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} |\theta_i|^{\ell_j} \leq K, \quad \text{and} \quad \frac{1}{N_{\mathcal{X},n}} \sum_{i \in \mathcal{I}_{\mathcal{X},n}} \|X_i\|^{\ell_j} \leq K.$$

In addition, the estimate $\hat{\mu}_{\ell_j}$ converges in probability to μ_{0,ℓ_j} under \tilde{P} for each j .

Theorem C.2. *Let $\hat{\theta}_i$ and se_i be given above and let $\hat{\chi}_i = \text{cva}_{\alpha,g}(\hat{m}_i)$ where \hat{m}_i is given above and $g(b) = (b^{\ell_1}, \dots, b^{\ell_p})$ for some positive integers ℓ_1, \dots, ℓ_p , at least one of which is even. Suppose that Assumptions C.5, C.6 and C.7 hold, and that $w(\cdot)$ is continuous in an open set containing $\{\gamma\} \times \mathcal{S}_1$ and is bounded away from zero on this set. Let \mathcal{A} be as given*

in Assumption C.7. Then, for all $\mathcal{X} \in \mathcal{A}$, $E_{\tilde{P}} \text{ANC}_n(\hat{\chi}^{(n)}; \mathcal{X}) \leq \alpha + o(1)$. If, in addition, $(Y_i, \hat{\sigma}_i)$ is independent over i under \tilde{P} , then $\text{ANC}_n(\hat{\chi}^{(n)}; \mathcal{X}) \leq \alpha + o_{\tilde{P}}(1)$.

As a consequence of Theorem C.2, we obtain, under the exchangeability condition (28), conditional EB coverage, as defined in Eq. (29), for any distribution P of the data and $\theta, \tilde{\nu}$ such that the conditions of Theorem C.2 hold with probability one with the sequence of probability measures $P(\cdot \mid \theta, \tilde{\nu}, X^{(n)}, \tilde{X}^{(n)})$ playing the role of \tilde{P} . This follows from the arguments in Appendix C.1.

Corollary C.1. *Let $\theta, \nu, X^{(n)}, \tilde{X}^{(n)}, Y_i$ follow a sequence of distributions P such that the conditions of Theorem C.2 hold with \tilde{X}_i taking on finitely many values, and $P(\cdot \mid \theta, \nu, X^{(n)}, \tilde{X}^{(n)})$ playing the role of \tilde{P} with probability one, and such that the exchangeability condition (28) holds. Then the intervals $CI_i = \{\hat{\theta}_i \pm w(\hat{\gamma}, \hat{\sigma}_i) \hat{\sigma}_i \cdot \text{cva}_{\alpha, g}(\hat{m}_i)\}$ satisfy the conditional EB coverage condition (29).*

The first part of Theorem 4.1 (average coverage) follows by applying Theorem C.2 with the conditional distribution $P(\cdot \mid \theta)$ playing the role of \tilde{P} . The second part (EB coverage) follows immediately from Corollary C.1.

References

- Abadie, A. and Kasy, M. (2019). Choosing among regularized estimators in empirical economics: The risk of machine learning. *The Review of Economics and Statistics*, 101(5):743–762.
- Andrews, I., McCloskey, A., and Kitagawa, T. (2021). Inference on winners. Unpublished manuscript, Harvard University.
- Angrist, J. D., Hull, P. D., Pathak, P. A., and Walters, C. R. (2017). Leveraging lotteries for school value-added: Testing and estimation. *The Quarterly Journal of Economics*, 132(2):871–919.
- Armstrong, T. B., Kolesár, M., and Plagborg-Møller, M. (2020). Robust empirical Bayes confidence intervals. arXiv: 2004.03448v2.
- Armstrong, T. B. and Kolesár, M. (2018). Optimal inference in a class of regression models. *Econometrica*, 86(2):655–683.
- Benjamini, Y. and Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):289–300.

- Benjamini, Y. and Yekutieli, D. (2005). False discovery rate-adjusted multiple confidence intervals for selected parameters. *Journal of the American Statistical Association*, 100(469):71–81.
- Bonhomme, S. and Weidner, M. (2021). Posterior average effects. arXiv: 1906.06360.
- Brown, L. D. and Greenshtein, E. (2009). Nonparametric empirical Bayes and compound decision approaches to estimation of a high-dimensional vector of normal means. *The Annals of Statistics*, 37(4):1685–1704.
- Cai, T. T., Low, M., and Ma, Z. (2014). Adaptive confidence bands for nonparametric regression functions. *Journal of the American Statistical Association*, 109(507):1054–1070.
- Carlin, B. P. and Louis, T. A. (2000). *Bayes and Empirical Bayes Methods for Data Analysis*. Chapman & Hall/CRC, New York, NY, 2nd edition.
- Casella, G. and Hwang, J. T. G. (2012). Shrinkage confidence procedures. *Statistical Science*, 27(1):51–60.
- Chetty, R., Friedman, J. N., and Rockoff, J. E. (2014). Measuring the impacts of teachers I: Evaluating bias in teacher value-added estimates. *American Economic Review*, 104(9):2593–2632.
- Chetty, R. and Hendren, N. (2018). The impacts of neighborhoods on intergenerational mobility II: County-level estimates. *The Quarterly Journal of Economics*, 133(3):1163–1228.
- Efron, B. (2010). *Large-Scale Inference: Empirical Bayes Methods for Estimation, Testing, and Prediction*. Cambridge University Press, New York, NY.
- Efron, B. (2015). Frequentist accuracy of Bayesian estimates. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77(3):617–646.
- Efron, B. (2019). Bayes, oracle Bayes and empirical Bayes. *Statistical Science*, 34(2):177–201.
- Finkelstein, A., Gentzkow, M., Hull, P., and Williams, H. (2017). Adjusting risk adjustment—accounting for variation in diagnostic intensity. *New England Journal of Medicine*, 376(7):608–610.
- Giacomini, R., Kitagawa, T., and Uhlig, H. (2019). Estimation under ambiguity. Cemmap Working Paper 24/19.

- Hansen, B. E. (2016). Efficient shrinkage in parametric models. *Journal of Econometrics*, 190(1):115–132.
- Hull, P. (2020). Estimating hospital quality with quasi-experimental data. Unpublished manuscript, University of Chicago.
- Ignatiadis, N. and Wager, S. (2021). Confidence intervals for nonparametric empirical Bayes analysis. arXiv: 1902.02774.
- Jacob, B. A. and Lefgren, L. (2008). Can principals identify effective teachers? Evidence on subjective performance evaluation in education. *Journal of Labor Economics*, 26(1):101–136.
- James, W. and Stein, C. M. (1961). Estimation with quadratic loss. In Neyman, J., editor, *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pages 361–379, Berkeley, CA. University of California Press.
- Jiang, W. and Zhang, C.-H. (2009). General maximum likelihood empirical Bayes estimation of normal means. *The Annals of Statistics*, 37(4):1647–1684.
- Johnstone, I. M. (2019). *Gaussian Estimation: Sequence and Multiresolution Models*. Book draft.
- Kane, T. and Staiger, D. (2008). Estimating teacher impacts on student achievement: An experimental evaluation. Technical Report 14607, National Bureau of Economic Research, Cambridge, MA.
- Kiefer, J. and Wolfowitz, J. (1956). Consistency of the maximum likelihood estimator in the presence of infinitely many incidental parameters. *The Annals of Mathematical Statistics*, 27(4):887–906.
- Koenker, R. and Mizera, I. (2014). Convex optimization, shape constraints, compound decisions, and empirical Bayes rules. *Journal of the American Statistical Association*, 109(506):674–685.
- Lee, J. D., Sun, D. L., Sun, Y., and Taylor, J. E. (2016). Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907–927.
- Liu, L., Moon, H. R., and Schorfheide, F. (2019). Forecasting with a panel tobit model. Unpublished manuscript, University of Pennsylvania.

- MacKinnon, J. G. and White, H. (1985). Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties. *Journal of Econometrics*, 29(3):305–325.
- Morris, C. N. (1983a). Parametric empirical Bayes confidence intervals. In Box, G. E. P., Leonard, T., and Wu, C.-F., editors, *Scientific Inference, Data Analysis, and Robustness*, pages 25–50, New York, NY. Academic Press.
- Morris, C. N. (1983b). Parametric empirical Bayes inference: Theory and applications. *Journal of the American Statistical Association*, 78(381):47–55.
- Nychka, D. (1988). Bayesian confidence intervals for smoothing splines. *Journal of the American Statistical Association*, 83(404):1134–1143.
- Pinelis, I. (2002). Monotonicity properties of the relative error of a Padé approximation for Mills’ ratio. *Journal of Inequalities in Pure & Applied Mathematics*, 3(2).
- Pratt, J. W. (1961). Length of confidence intervals. *Journal of the American Statistical Association*, 56(295):549–567.
- Robbins, H. (1951). Asymptotically subminimax solutions of compound statistical decision problems. In Neyman, J., editor, *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, pages 131–149. University of California Press, Berkeley, California.
- Storey, J. D. (2002). A direct approach to false discovery rates. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):479–498.
- Wahba, G. (1983). Bayesian “confidence intervals” for the cross-validated smoothing spline. *Journal of the Royal Statistical Society: Series B (Methodological)*, 45(1):133–150.
- Wasserman, L. (2006). *All of Nonparametric Statistics*. Springer, New York, NY.
- Xie, X., Kou, S. C., and Brown, L. D. (2012). SURE estimates for a heteroscedastic hierarchical model. *Journal of the American Statistical Association*, 107(500):1465–1479.