Confidently Comparing Estimates with the c-value

Brian Trippe, Sameer K. Deshpande, Tamara Broderick





Learning from Educational Testing Data

 National Center for Education Statistics gathers standardized tests from U.S. high schools







- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances







- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school







- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!







- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!
 Hierarchical Bayesian Approach
 - Share strength across similar schools



- Enrollment size
- Type (Catholic, charter, public)

- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!
 Hierarchical Bayesian Approach
 - Share strength across similar schools
 - Estimate school performances by the posterior mean



- Enrollment size
- Type (Catholic, charter, public)

- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!
 Hierarchical Bayesian Approach
 - Share strength across similar schools
 - Estimate school performances by the posterior mean [Hoff, 2020]
 - 5-50 students tested at 676 schools



- Enrollment size
- Type (Catholic, charter, public)

Learning from Educational Testing Data

- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!

Hierarchical Bayesian Approach [Lindley and

- Share strength across similar schools
- Estimate school performances by the posterior mean [Hoff, 2020]
 - 5-50 students tested at 676 schools



- Enrollment size
- Type (Catholic, charter, public)

Learning from Educational Testing Data

- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!

Hierarchical Bayesian Approach [Lindley and

- Share strength across similar schools
- Estimate school performances by the posterior mean [Hoff, 2020]
 - 5-50 students tested at 676 schools
- Limitation: complexity, subjectivity of the prior



- Enrollment size
- Type (Catholic, charter, public)

Learning from Educational Testing Data

- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!

Hierarchical Bayesian Approach [Lindley and

- Share strength across similar schools
- Estimate school performances by the posterior mean [Hoff, 2020]
 - 5-50 students tested at 676 schools
- Limitation: complexity, subjectivity of the prior
- Question: is this more accurate than simple averages?



- Enrollment size
- Type (Catholic, charter, public)

Learning from Educational Testing Data

- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisy!

Hierarchical Bayesian Approach [Lindley and

Smith, 1972, Rubin, 1981, Gelman et al., 2013]

- Share strength across similar schools
- Estimate school performances by the posterior mean [Hoff, 2020]
 - 5-50 students tested at 676 schools
- Limitation: complexity, subjectivity of the prior
- Question: is this more accurate than simple ____averages?

This question comes up in many new analyses!



- Region
- Enrollment size
- Type (Catholic, charter, public)

Learning from Educational Testing Data

- National Center for Education Statistics gathers standardized tests from U.S. high schools
- Want to know school-level performances
- Standardized tests on small sample of students for each school – Noisv!

Hierarchical Bayesian Approach [Lindley and

- Share strength across similar schools
- Estimate school performances by the posterior mean [Hoff, 2020]
 - 5-50 students tested at 676 schools
- Limitation: complexity, subjectivity of the prior
- Question: is this more accurate than simple averages? Yes (c=99.26%)!This question comes up in many new analyses!



- Enrollment size
- Type (Catholic, charter, public)

Roadmap

- Justifying complexity
- Methods for choosing methods
- The c-value as a measure of confidence (our method)
- How and when we can compute c-values
- Application to educational testing
- Extensions to nonlinear models and estimates

Roadmap

- Justifying complexity
- Methods for choosing methods
- The c-value as a measure of confidence (our method)
- How and when we can compute c-values
- Application to educational testing
- Extensions to nonlinear models and estimates

In Machine Learning

- ► Information criteria (AIC / BIC)
- Cross-validation

In Machine Learning

Information criteria (AIC / BIC)
 Cross-validation
 Apply to prediction, not parameter estimation!

In Machine Learning

- Cross-validation

Decision Theory

 $\begin{array}{l} \text{Model: } y \sim p(\cdot; \theta) \\ \\ \text{Estimates: } \hat{\theta}(y) \text{ vs. } \theta^*(y) \end{array}$

Information criteria (AIC / BIC)
 Cross-validation
 Apply to prediction, not parameter estimation!

In Machine Learning

- Cross-validation

Decision Theory

$$\begin{array}{ll} \text{Model: } y \sim p(\cdot; \theta) & \text{Loss: } L(\theta, \cdot) \\ \text{Estimates: } \hat{\theta}(y) \text{ vs. } \theta^*(y) \end{array}$$

Information criteria (AIC / BIC)
 Cross-validation
 Apply to prediction, not parameter estimation!

In Machine Learning

Decision Theory

Information criteria (AIC / BIC)
Apply to prediction, not parameter estimation!

$$\begin{array}{ll} \mbox{Model: } y \sim p(\cdot; \theta) & \mbox{Loss: } L(\theta, \cdot) \\ \mbox{Estimates: } \hat{\theta}(y) \mbox{ vs. } \theta^*(y) \end{array}$$

 $\text{Standard criterion} - \text{Risk } R(\theta, \hat{\theta}(\cdot)) := \mathbb{E}_{y \sim p(\cdot; \theta)} \left[L(\theta, \hat{\theta}(y)) \right]$

In Machine Learning

Decision Theory

Information criteria (AIC / BIC)
Apply to prediction, not parameter estimation!

$$\begin{array}{ll} \mbox{Model: } y \sim p(\cdot; \theta) & \mbox{Loss: } L(\theta, \cdot) \\ \mbox{Estimates: } \hat{\theta}(y) \mbox{ vs. } \theta^*(y) \end{array}$$

Standard criterion – Risk $R(\theta, \hat{\theta}(\cdot)) := \mathbb{E}_{y \sim p(\cdot; \theta)} \left[L(\theta, \hat{\theta}(y)) \right]$ • Choose $\hat{\theta}(y)$ over $\theta^*(y)$ if $R(\theta, \hat{\theta}(\cdot)) < R(\theta, \theta^*(\cdot))$

In Machine Learning

Decision Theory

Information criteria (AIC / BIC)
Apply to prediction, not parameter estimation!

$$\begin{array}{ll} \text{Model: } y \sim p(\cdot; \theta) & \text{Loss: } L(\theta, \cdot) \\ \\ \text{Estimates: } \hat{\theta}(y) \text{ vs. } \theta^*(y) \end{array}$$

Standard criterion – Risk $R(\theta, \hat{\theta}(\cdot)) := \mathbb{E}_{y \sim p(\cdot; \theta)} \left| L(\theta, \hat{\theta}(y)) \right|$

- Choose $\hat{\theta}(y)$ over $\theta^*(y)$ if $R(\theta, \hat{\theta}(\cdot)) < R(\theta, \theta^*(\cdot))$
- Risk depends on θ !

In Machine Learning

- Information criteria (AIC / BIC)
 - Cross-validation

Decision Theory

$$\begin{array}{ll} \text{Model: } y \sim p(\cdot; \theta) & \text{Loss: } L(\theta, \cdot) \\ \\ \text{Estimates: } \hat{\theta}(y) \text{ vs. } \theta^*(y) \end{array}$$

Standard criterion – Risk $R(\theta, \hat{\theta}(\cdot)) := \mathbb{E}_{y \sim p(\cdot;\theta)} \left[L(\theta, \hat{\theta}(y)) \right]$

- Choose $\hat{\theta}(y)$ over $\theta^*(y)$ if $R(\theta, \hat{\theta}(\cdot)) < R(\theta, \theta^*(\cdot))$
- Risk depends on θ !
- \blacktriangleright I care about **my** y

In Machine Learning

Decision Theory

Information criteria (AIC / BIC)
Apply to prediction, not parameter estimation!

Standard criterion – Risk $R(\theta, \hat{\theta}(\cdot)) := \mathbb{E}_{y \sim p(\cdot;\theta)} \left| L(\theta, \hat{\theta}(y)) \right|$

- Choose $\hat{\theta}(y)$ over $\theta^*(y)$ if $R(\theta, \hat{\theta}(\cdot)) < R(\theta, \theta^*(\cdot))$
- Risk depends on $\theta!$
- I care about my y

Goal: Measure of confidence that $\theta^*(\cdot)$ has smaller loss than $\hat{\theta}(\cdot)$

In Machine Learning

Decision Theory

Information criteria (AIC / BIC)
Apply to prediction, not parameter estimation!

Standard criterion – Risk $R(\theta, \hat{\theta}(\cdot)) := \mathbb{E}_{y \sim p(\cdot;\theta)} \left| L(\theta, \hat{\theta}(y)) \right|$

- Choose $\hat{\theta}(y)$ over $\theta^*(y)$ if $R(\theta, \hat{\theta}(\cdot)) < R(\theta, \theta^*(\cdot))$
- Risk depends on $\theta!$
- I care about my y

Goal: Measure of confidence that $\theta^*(\cdot)$ has smaller loss than $\hat{\theta}(\cdot)$

- On the observed dataset
- 2. Without needing subjective assumptions about θ

$$\label{eq:model:product} \begin{array}{lll} \mbox{Model:} & y \sim p(\cdot; \theta) & \mbox{Loss:} & L(\theta, \cdot) \\ \mbox{Estimates:} & \hat{\theta}(y) & \mbox{vs.} & \theta^*(y) \end{array}$$







Challenge: identify if $W(\theta, y) > 0$



Challenge: identify if $W(\theta, y) > 0$

Our approach

Introduce high probability lower bound on the win, $b(y, \alpha)$



Challenge: identify if $W(\theta, y) > 0$

Our approach

Introduce high probability lower bound on the win, $b(y, \alpha)$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{y \sim p(\cdot;\theta)} \left[W(\theta, y) > b(y, \alpha) \right] \geq \alpha$



Challenge: identify if $W(\theta, y) > 0$

Our approach

Introduce high probability lower bound on the win, $b(y, \alpha)$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{y \sim p(\cdot;\theta)} \left[W(\theta, y) > b(y, \alpha) \right] \geq \alpha$

▶ 95% confidence in $\theta^*(y)$ if b(y, 0.95) > 0



Our approach

Introduce high probability lower bound on the win, $b(y, \alpha)$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{y \sim p(\cdot;\theta)} \left[W(\theta, y) > b(y, \alpha) \right] \geq \alpha$

▶ 95% confidence in $\theta^*(y)$ if b(y, 0.95) > 0



Our approach

Introduce high probability lower bound on the win, $b(y, \alpha)$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{y \sim p(\cdot;\theta)} \left[W(\theta, y) > b(y, \alpha) \right] \geq \alpha$

▶ 95% confidence in $\theta^*(y)$ if b(y, 0.95) > 0



Our approach

Introduce high probability lower bound on the win, $b(y, \alpha)$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{y \sim p(\cdot;\theta)} \left[W(\theta, y) > b(y, \alpha) \right] \geq \alpha$

▶ 95% confidence in $\theta^*(y)$ if b(y, 0.95) > 0

Define
$$\mathbf{c}(\mathbf{y}) \coloneqq \inf_{\alpha \in [\mathbf{0}, \mathbf{1}]} \left\{ \alpha | \mathbf{b}(\mathbf{y}, \alpha) \leq \mathbf{0} \right\}$$

• Loosely, largest level α below which $b(y, \alpha) > 0$
Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta} \left[W(\theta, y) > b(y, \alpha) \right] \geq \alpha$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta}\left[W(\theta, y) > b(y, \alpha)\right] \geq \alpha$

The c-value as a quantification of confidence

Thm 1: For any θ and $\alpha \in [0, 1]$, $\mathbb{P}_{\theta} \left[c(y) > \alpha \text{ and } W(\theta, y) \leq 0 \right] \leq 1 - \alpha$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta}\left[W(\theta, y) > b(y, \alpha)\right] \geq \alpha$

The c-value as a quantification of confidence

Thm 1: For any θ and $\alpha \in [0, 1]$, $\mathbb{P}_{\theta} \left[c(y) > \alpha \text{ and } W(\theta, y) \leq 0 \right] \leq 1 - \alpha$

Unlikely that both (A.) c-value is close to 1 and (B.) θ*(y) is not more accurate than θ(y)

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta}\left[W(\theta, y) > b(y, \alpha)\right] \geq \alpha$

The c-value as a quantification of confidence

Thm 1: For any θ and $\alpha \in [0, 1]$, $\mathbb{P}_{\theta} \left[c(y) > \alpha \text{ and } W(\theta, y) \leq 0 \right] \leq 1 - \alpha$

Unlikely that both (A.) c-value is close to 1 and (B.) θ*(y) is not more accurate than θ(y)

Using c(y) to choose between $\hat{\theta}(y)$ and $\theta^*(y)$

Define two-stage estimator

$$\theta^{\dagger}(y,\alpha) := \mathbb{1}[c(y) > \alpha]\theta^{*}(y) + \mathbb{1}[c(y) \le \alpha]\hat{\theta}(y)$$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta}\left[W(\theta, y) > b(y, \alpha)\right] \geq \alpha$

The c-value as a quantification of confidence

Thm 1: For any θ and $\alpha \in [0, 1]$, $\mathbb{P}_{\theta} \left[c(y) > \alpha \text{ and } W(\theta, y) \leq 0 \right] \leq 1 - \alpha$

Unlikely that both (A.) c-value is close to 1 and (B.) θ*(y) is not more accurate than θ̂(y)

Using c(y) to choose between $\hat{\theta}(y)$ and $\theta^*(y)$

Define two-stage estimator

 $\boldsymbol{\theta}^{\dagger}(\boldsymbol{y},\boldsymbol{\alpha}) \mathrel{\mathop:}= \mathbbm{1}[\boldsymbol{c}(\boldsymbol{y}) > \boldsymbol{\alpha}] \boldsymbol{\theta}^{*}(\boldsymbol{y}) + \mathbbm{1}[\boldsymbol{c}(\boldsymbol{y}) \leq \boldsymbol{\alpha}] \hat{\boldsymbol{\theta}}(\boldsymbol{y})$

Deviates from default only when confident in alternative

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta}\left[W(\theta, y) > b(y, \alpha)\right] \geq \alpha$

The c-value as a quantification of confidence

Thm 1: For any θ and $\alpha \in [0, 1]$, $\mathbb{P}_{\theta} \left[c(y) > \alpha \text{ and } W(\theta, y) \leq 0 \right] \leq 1 - \alpha$

Unlikely that both (A.) c-value is close to 1 and (B.) θ*(y) is not more accurate than θ(y)

Using c(y) to choose between $\hat{\theta}(y)$ and $\theta^*(y)$

Define two-stage estimator

$$\theta^{\dagger}(y,\alpha) := \mathbb{1}[c(y) > \alpha]\theta^{*}(y) + \mathbb{1}[c(y) \le \alpha]\hat{\theta}(y)$$

Deviates from default only when confident in alternative

Thm 2: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta} \left[L(\theta, \theta^{\dagger}(y, \alpha)) > L(\theta, \hat{\theta}(y)) \right] \leq 1 - \alpha$

Condition: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta}\left[W(\theta, y) > b(y, \alpha)\right] \geq \alpha$

The c-value as a quantification of confidence

Thm 1: For any θ and $\alpha \in [0, 1]$, $\mathbb{P}_{\theta} \left[c(y) > \alpha \text{ and } W(\theta, y) \leq 0 \right] \leq 1 - \alpha$

Unlikely that both (A.) c-value is close to 1 and (B.) θ*(y) is not more accurate than θ(y)

Using c(y) to choose between $\hat{\theta}(y)$ and $\theta^*(y)$

Define two-stage estimator

$$\theta^{\dagger}(y,\alpha) := \mathbb{1}[c(y) > \alpha]\theta^{*}(y) + \mathbb{1}[c(y) \le \alpha]\hat{\theta}(y)$$

Deviates from default only when confident in alternative

Thm 2: For any θ and $\alpha \in [0,1]$, $\mathbb{P}_{\theta} \left[L(\theta, \theta^{\dagger}(y, \alpha)) > L(\theta, \hat{\theta}(y)) \right] \leq 1 - \alpha$

▶ If we report $\theta^*(y)$ only when c(y) > 0.95, we do worse than $\hat{\theta}(\cdot)$ at $_{6/15}$

Roadmap

- Justifying complexity
- Methods for choosing methods
- The c-value as a measure of confidence
- How and when we can compute c-values
- Application to educational testing
- Extensions to nonlinear models and estimates

Roadmap

- Justifying complexity
- Methods for choosing methods
- The c-value as a measure of confidence
- How and when we can compute c-values
- Application to educational testing
- Extensions to nonlinear models and estimates







Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$

$$W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$$



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$

▶ Relies on unknown θ , and $L(\theta, \hat{\theta})$ and $L(\theta, \theta^*)$ are dependent

$$W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$$



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$

▶ Relies on unknown θ , and $L(\theta, \hat{\theta})$ and $L(\theta, \theta^*)$ are dependent

$$W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$$

= $\|\hat{\theta}(y) - y\|^2 - \|\theta^*(y) - y\|^2 + 2\langle \hat{\theta}(y) - \theta^*(y), y - \theta \rangle$



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$ \blacktriangleright Relies on unknown θ , and $L(\theta, \hat{\theta})$ and $L(\theta, \theta^*)$ are dependent $W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$ $= \underbrace{\|\hat{\theta}(y) - y\|^2 - \|\theta^*(y) - y\|^2}_{\text{observed and computable}} + \underbrace{2\langle \hat{\theta}(y) - \theta^*(y), y - \theta \rangle}_{\text{unobserved}(\ddagger)}$



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$

▶ Relies on unknown θ , and $L(\theta, \hat{\theta})$ and $L(\theta, \theta^*)$ are dependent

$$W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$$

=
$$\underbrace{\|\hat{\theta}(y) - y\|^2 - \|\theta^*(y) - y\|^2}_{\text{observed and computable}} + \underbrace{2\langle\hat{\theta}(y) - \theta^*(y), y - \theta\rangle}_{\text{unobserved}(\ddagger)}$$

$$\geq \underbrace{\text{observed}}_{\text{berved}} + F_{\ddagger}^{-1}(1 - \alpha; \theta), \text{ with prob. } \alpha$$



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$ \blacktriangleright Relies on unknown θ , and $L(\theta, \hat{\theta})$ and $L(\theta, \theta^*)$ are dependent $W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$ $= \underbrace{\|\hat{\theta}(y) - y\|^2 - \|\theta^*(y) - y\|^2}_{\text{observed and computable}} + \underbrace{2\langle \hat{\theta}(y) - \theta^*(y), y - \theta \rangle}_{\text{unobserved}(\ddagger)}$ $\geq \text{observed} + F_{\ddagger}^{-1}(1 - \alpha; g(\theta)), \text{ with prob. } \alpha$

• Key observation: F_{\ddagger} depends on θ only through a scalar $g(\theta)$



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$ **•** Relies on unknown θ , and $L(\theta, \hat{\theta})$ and $L(\theta, \theta^*)$ are dependent $W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$ $= \|\hat{\theta}(y) - y\|^2 - \|\theta^*(y) - y\|^2 + 2\langle \hat{\theta}(y) - \theta^*(y), y - \theta \rangle$ observed and computable unobserved(t) \geq observed + $F_{\dagger}^{-1}(1-\alpha;g(\theta))$, with prob. α \geq observed + $\inf_{\lambda \in C(u, \frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$, with prob. α

Key observation: F[‡] depends on θ only through a scalar g(θ)
Split excess α across interval and quantile (union bound)



Idea: Use $1 - \alpha$ lower quantile of $W(\theta, y)$ **•** Relies on unknown θ , and $L(\theta, \hat{\theta})$ and $L(\theta, \theta^*)$ are dependent $W(\theta, y) = \|\hat{\theta}(y) - \theta\|^2 - \|\theta^*(y) - \theta\|^2$ $= \|\hat{\theta}(y) - y\|^2 - \|\theta^*(y) - y\|^2 + 2\langle \hat{\theta}(y) - \theta^*(y), y - \theta \rangle$ observed and computable unobserved(t) \geq observed + $F_{\dagger}^{-1}(1-\alpha;g(\theta))$, with prob. α $\geq \text{observed} + \inf_{\lambda \in C(y, \frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right), \text{ with prob. } \alpha$ $b(y,\alpha)$

Key observation: F_{\ddagger} depends on θ only through a scalar $g(\theta)$

8 / 15

Model: $\theta, y \in \mathbb{R}^N$ $L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2$

 $\mathsf{Model:} \ \theta, y \in \mathbb{R}^N \ \mathsf{with} \ y \sim \mathcal{N}(\theta, I_N) \quad \ L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2$





• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$ Filling out the details of the bound

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^{N} y_n$ Filling out the details of the bound

$$b(y,\alpha) = \text{ observed} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger =$ unobserved

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$ Filling out the details of the bound

$$b(y,\alpha) = \|\hat{\theta}(y) - y\|^2 - \|\theta^*(y) - y\|^2 + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger =$ unobserved

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

θ*(y) is a classic Bayes estimate, shrinks towards mean
τ > 0 and ȳ := N⁻¹ Σ^N_{n=1} y_n
Filling out the details of the bound

$$b(y,\alpha) = -\|\frac{y+\tau^{-2}\mathbf{1}_N\bar{y}}{1+\tau^{-2}} - y\|^2 + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger =$ unobserved

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

θ*(y) is a classic Bayes estimate, shrinks towards mean
τ > 0 and ȳ := N⁻¹ Σ^N_{n=1} y_n
Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger =$ unobserved

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

θ*(y) is a classic Bayes estimate, shrinks towards mean
τ > 0 and ȳ := N⁻¹ Σ^N_{n=1} y_n
Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger =$ **unobserved**

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

θ^{*}(*y*) is a classic Bayes estimate, shrinks towards mean
τ > 0 and *y* := N⁻¹ ∑_{n=1}^N y_n
Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger = 2 \langle \hat{\theta}(y) - \theta^*(y), y - \theta \rangle$

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

θ^{*}(*y*) is a classic Bayes estimate, shrinks towards mean
τ > 0 and *y* := N⁻¹ ∑_{n=1}^N y_n
Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger = \frac{2}{1+\tau^2} \langle y - \bar{y} \mathbf{1}_N, y - \theta \rangle$

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$ Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y-\bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger \sim \frac{2}{1+\tau^2} \left[\chi_{N-1}^2(\frac{1}{4}\|\theta - \bar{\theta}\mathbf{1}_N\|^2) - \frac{1}{4}\|\theta - \bar{\theta}\mathbf{1}_N\|^2\right]$
 $\searrow \chi_{N-1}^2(\lambda)$ is non-central χ^2 with $N-1$ degrees of freedom, non-centrality λ

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$ Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger \sim \frac{2}{1+\tau^2} \left[\chi^2_{N-1}(g(\theta)) - g(\theta) \right], \ g(\theta) = \frac{1}{4} \|\theta - \bar{\theta} \mathbf{1}_N\|^2$ $\blacktriangleright \chi^2_{N-1}(\lambda)$ is non-central χ^2 with N-1 degrees of freedom, non-centrality λ

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$ Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger \sim \frac{2}{1+\tau^2} \left[\chi^2_{N-1}(g(\theta)) - g(\theta) \right], \ g(\theta) = \frac{1}{4} \|\theta - \bar{\theta} \mathbf{1}_N\|^2$

▶ $\chi^2_{N-1}(\lambda)$ is non-central χ^2 with N-1 degrees of freedom, non-centrality λ

► Interval $C(y, 1 - \alpha)$ for $g(\theta)$ from $||y - \bar{y}\mathbf{1}_N||^2 \sim \chi^2_{N-1}(4g(\theta))$
Example Bound – The Lindley and Smith [1972] Estimator

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$ Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger \sim \frac{2}{1+\tau^2} \left[\chi^2_{N-1}(g(\theta)) - g(\theta) \right], \ g(\theta) = \frac{1}{4} \|\theta - \bar{\theta} \mathbf{1}_N\|^2$

▶ $\chi^2_{N-1}(\lambda)$ is non-central χ^2 with N-1 degrees of freedom, non-centrality λ

► Interval $C(y, 1 - \alpha)$ for $g(\theta)$ from $||y - \bar{y}\mathbf{1}_N||^2 \sim \chi^2_{N-1}(4g(\theta))$ Take-aways: Correct coverage by construction

Example Bound – The Lindley and Smith [1972] Estimator

 $\begin{array}{ll} \text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, I_N) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \text{Estimates: } \underbrace{\hat{\theta}(y) = y}_{\text{Default (MLE)}} & \text{vs.} & \underbrace{\theta^*(y) = \frac{y + \tau^{-2} \bar{y} \mathbf{1}_N}{1 + \tau^{-2}}}_{\text{Alternative (Lindley and Smith)}} \end{array}$

• $\theta^*(y)$ is a classic Bayes estimate, shrinks towards mean • $\tau > 0$ and $\bar{y} := N^{-1} \sum_{n=1}^N y_n$ Filling out the details of the bound

$$b(y,\alpha) = \frac{-\|y - \bar{y}\mathbf{1}_N\|^2}{(1+\tau^2)^2} + \inf_{\lambda \in C(y,\frac{1-\alpha}{2})} F_{\ddagger}^{-1}\left(\frac{1-\alpha}{2}; g(\theta) = \lambda\right)$$

where $\ddagger \sim \frac{2}{1+\tau^2} \left[\chi^2_{N-1}(g(\theta)) - g(\theta) \right], \ g(\theta) = \frac{1}{4} \|\theta - \bar{\theta} \mathbf{1}_N\|^2$

▶ $\chi^2_{N-1}(\lambda)$ is non-central χ^2 with N-1 degrees of freedom, non-centrality λ

► Interval $C(y, 1 - \alpha)$ for $g(\theta)$ from $||y - \bar{y}\mathbf{1}_N||^2 \sim \chi^2_{N-1}(4g(\theta))$ Take-aways: Correct coverage by construction, Computable



Use simulated data for calibration, power, and risk



• $b(y, \alpha)$ is conservative across levels α and θ

Use simulated data for calibration, power, and risk



10 / 15

Use simulated data for calibration, power, and risk



 Coverage has little θ dependence

Use simulated data for calibration, power, and risk



dependence











Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **•** Goal: estimate school-specific means, $\theta \in \mathbb{R}^N$
- ▶ Default: $\hat{\theta}(y) = y$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **•** Goal: estimate school-specific means, $\theta \in \mathbb{R}^N$
- ▶ Default: $\hat{\theta}(y) = y$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

- ▶ D = 8 features of each school (region, school type, enrollment, ...)
- ▶ Prior: $\theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2)$ for $\beta \in \mathbb{R}^D$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

• Prior:
$$\theta_n \stackrel{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2)$$
 for $\beta \in \mathbb{R}^D$

• Likelihood: $y_n \overset{indep}{\sim} \mathcal{N}(\theta_n, \sigma^2/\text{size}_n)$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **•** Goal: estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

▶ Prior:
$$\theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2)$$
 for $\beta \in \mathbb{R}^D$

► Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \operatorname{diag}(\sigma^2/\operatorname{size}_1, \dots, \sigma^2/\operatorname{size}_N)$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

- ▶ D = 8 features of each school (region, school type, enrollment, ...)
- $\blacktriangleright \text{ Prior: } \theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2) \text{ for } \beta \in \mathbb{R}^D$
- Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \operatorname{diag}(\sigma^2/\operatorname{size}_1, \ldots, \sigma^2/\operatorname{size}_N)$
 - Estimate τ, β, σ by empirical Bayes (lme4)

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

$$\blacktriangleright \text{ Prior: } \theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2) \text{ for } \beta \in \mathbb{R}^D$$

• Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \operatorname{diag}(\sigma^2/\operatorname{size}_1, \ldots, \sigma^2/\operatorname{size}_N)$

Estimate τ, β, σ by empirical Bayes (lme4)

• Alternative $\theta^*(y) = [I_N + \tau^{-2}\Sigma]^{-1}y + [I_N + \tau^2\Sigma^{-1}]^{-1}X\beta$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

$$\blacktriangleright \text{ Prior: } \theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2) \text{ for } \beta \in \mathbb{R}^D$$

► Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \operatorname{diag}(\sigma^2/\operatorname{size}_1, \ldots, \sigma^2/\operatorname{size}_N)$

Estimate τ, β, σ by empirical Bayes (lme4)

• Alternative $\theta^*(y) = [I_N + \tau^{-2}\Sigma]^{-1}y + [I_N + \tau^2\Sigma^{-1}]^{-1}X\beta$

• $\theta^*(y)$ is an affine transformation of y

 $\label{eq:model} \begin{tabular}{ll} \begin{$

 $\begin{cases} \mathsf{Model:} \ \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, \Sigma) & L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2 \\ \mathsf{Estimates:} \ \hat{\theta}(y) = \underline{Ay + k} & \mathsf{vs.} \ \theta^*(y) = \underline{Cy + \ell} \text{ for } A, C \in \mathbb{R}^{N \times N} & k, \ell \in \mathbb{R}^N \end{cases}$

 $\text{Model: } \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, \Sigma) \quad \ L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2$

 $\text{Estimates: } \hat{\theta}(y) \!=\! Ay + k \quad \text{vs. } \theta^*(y) \!=\! Cy + \ell \text{ for } A, C \!\in\! \mathbb{R}^{N \times N} \hspace{0.1cm} k, \ell \!\in\! \mathbb{R}^N$

 Applications: Gaussian process kernel selection, shrinkage estimation, linear regression

$$\mathsf{Model:} \ \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, \Sigma) \quad \ L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2$$

 $\text{Estimates:} \ \hat{\theta}(y) \!=\! Ay + k \quad \text{vs.} \ \ \theta^*(y) \!=\! Cy + \ell \ \text{for} \ A, C \!\in\! \mathbb{R}^{N \times N} \ \ k, \ell \!\in\! \mathbb{R}^N$

 Applications: Gaussian process kernel selection, shrinkage estimation, linear regression

$$\begin{split} b(y,\alpha) &:= \|\hat{\theta} - y\|^2 - \|\theta^* - y\|^2 + 2\mathrm{tr}[(A - C)\Sigma] + \\ & 2z_{\frac{1-\alpha}{2}} \sqrt{U\left(\frac{1-\alpha}{2}\right) + \frac{1}{2}\|\Sigma^{\frac{1}{2}}(A + A^\top - C - C^\top)\Sigma^{\frac{1}{2}}\|_{F}^2} \end{split}$$

where

$$\begin{split} U(1-\alpha) &:= \inf_{\delta > 0} \left\{ \delta \left| \| \hat{\theta}(y) - \theta^*(y) \|_{\Sigma}^2 \leq (\delta + \| \Sigma^{\frac{1}{2}} (A - C) \Sigma^{\frac{1}{2}} \|_F^2) + \right. \\ z_{1-\alpha} \sqrt{2 \| \Sigma^{\frac{1}{2}} (A - C) \Sigma (A - C)^\top \Sigma^{\frac{1}{2}} \|_F^2 + 4 \| \Sigma^{\frac{1}{2}} (A - C) \Sigma^{\frac{1}{2}} \|_{\mathsf{Op}}^2 \delta} \right\} \\ \\ \text{is a high confidence upper bound on } g(\theta) &:= \left\| (A - C) \theta + (k - \ell) \right\|_{\Sigma}^2 \end{split}$$

$$\mathsf{Model:} \ \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, \Sigma) \quad \ L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2$$

 $\text{Estimates:} \ \hat{\theta}(y) \!=\! Ay + k \quad \text{vs.} \ \ \theta^*(y) \!=\! Cy + \ell \ \text{for} \ A, C \!\in\! \mathbb{R}^{N \times N} \ \ k, \ell \!\in\! \mathbb{R}^N$

 Applications: Gaussian process kernel selection, shrinkage estimation, linear regression

$$\begin{split} b(y,\alpha) &:= \|\hat{\theta} - y\|^2 - \|\theta^* - y\|^2 + 2\mathrm{tr}[(A-C)\Sigma] + \\ & 2z_{\frac{1-\alpha}{2}} \sqrt{U\left(\frac{1-\alpha}{2}\right) + \frac{1}{2} \|\Sigma^{\frac{1}{2}}(A+A^\top - C - C^\top)\Sigma^{\frac{1}{2}}\|_F^2} \end{split}$$

where

$$U(1-\alpha) := \inf_{\delta>0} \left\{ \delta \left\| \|\hat{\theta}(y) - \theta^*(y)\|_{\Sigma}^2 \le (\delta + \|\Sigma^{\frac{1}{2}}(A-C)\Sigma^{\frac{1}{2}}\|_F^2) + z_{1-\alpha} \sqrt{2\|\Sigma^{\frac{1}{2}}(A-C)\Sigma(A-C)^\top \Sigma^{\frac{1}{2}}\|_F^2 + 4\|\Sigma^{\frac{1}{2}}(A-C)\Sigma^{\frac{1}{2}}\|_{\mathsf{OP}}^2 \delta} \right\}$$

is a high confidence upper bound on $g(\theta) := \left\| (A - C)\theta + (k - \ell) \right\|_{\Sigma}^2$

Computable:
$$c(y) = c_value(y, \Sigma, A, k, C, 1)$$

$$\mathsf{Model:} \ \theta, y \in \mathbb{R}^N \text{ with } y \sim \mathcal{N}(\theta, \Sigma) \quad \ L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2$$

 $\text{Estimates:} \ \hat{\theta}(y) \!=\! Ay + k \quad \text{vs.} \ \ \theta^*(y) \!=\! Cy + \ell \ \text{for} \ A, C \!\in\! \mathbb{R}^{N \times N} \ \ k, \ell \!\in\! \mathbb{R}^N$

 Applications: Gaussian process kernel selection, shrinkage estimation, linear regression

$$\begin{split} b(y,\alpha) &:= \|\hat{\theta} - y\|^2 - \|\theta^* - y\|^2 + 2\mathrm{tr}[(A-C)\Sigma] + \\ & 2z_{\frac{1-\alpha}{2}} \sqrt{U\left(\frac{1-\alpha}{2}\right) + \frac{1}{2} \|\Sigma^{\frac{1}{2}}(A+A^\top - C - C^\top)\Sigma^{\frac{1}{2}}\|_{1}^{2}} \end{split}$$

where

$$U(1-\alpha) := \inf_{\delta > 0} \left\{ \delta \left| \| \hat{\theta}(y) - \theta^*(y) \|_{\Sigma}^2 \le (\delta + \| \Sigma^{\frac{1}{2}} (A - C) \Sigma^{\frac{1}{2}} \|_F^2 \right\} + \frac{1}{2} \left(\int_{\Sigma}^{1} \frac{1}{$$

$$z_{1-\alpha} \sqrt{2 \|\Sigma^{\frac{1}{2}} (A-C) \Sigma (A-C)^{\top} \Sigma^{\frac{1}{2}} \|_{F}^{2}} + 4 \|\Sigma^{\frac{1}{2}} (A-C) \Sigma^{\frac{1}{2}} \|_{\mathsf{OP}}^{2} \delta \bigg\}$$

is a high confidence upper bound on $g(\theta) \mathrel{\mathop:}= \big\| (A-C)\theta + (k-\ell) \big\|_{\Sigma}^2$

- Computable : c(y) = c_value(y, Σ, A, k, C, 1)
 Some analytical challenges:
 - Non-asymptotic error control [Berry, 1941]

$$\mathsf{Model:} \ \theta, y \in \mathbb{R}^N \ \mathsf{with} \ y \sim \mathcal{N}(\theta, \Sigma) \quad \ L(\theta, \theta'(y)) = \|\theta'(y) - \theta\|^2$$

Estimates: $\hat{\theta}(y) = Ay + k$ vs. $\theta^*(y) = Cy + \ell$ for $A, C \in \mathbb{R}^{N \times N}$ $k, \ell \in \mathbb{R}^N$

 Applications: Gaussian process kernel selection, shrinkage estimation, linear regression

$$\begin{split} b(y,\alpha) &:= \|\hat{\theta} - y\|^2 - \|\theta^* - y\|^2 + 2\mathrm{tr}[(A - C)\Sigma] + \\ & 2z_{\frac{1-\alpha}{2}}\sqrt{U\left(\frac{1-\alpha}{2}\right) + \frac{1}{2}\|\Sigma^{\frac{1}{2}}(A + A^\top - C - C^\top)\Sigma^{\frac{1}{2}}\|^2} \end{split}$$

where

$$U(1-\alpha) := \inf_{\delta > 0} \left\{ \delta \left\| \| \hat{\theta}(y) - \theta^*(y) \|_{\Sigma}^2 \le (\delta + \| \Sigma^{\frac{1}{2}} (A - C) \Sigma^{\frac{1}{2}} \|_F^2 \right\} + \frac{1}{2} \left\| \nabla^2 (A - C) \nabla^2 (A$$

$$z_{1-\alpha} \sqrt{2 \|\Sigma^{\frac{1}{2}} (A-C) \Sigma (A-C)^{\top} \Sigma^{\frac{1}{2}} \|_{F}^{2}} + 4 \|\Sigma^{\frac{1}{2}} (A-C) \Sigma^{\frac{1}{2}} \|_{\mathsf{OP}}^{2} \delta \bigg\}$$

is a high confidence upper bound on $g(\theta) \mathrel{\mathop:}= \big\| (A-C)\theta + (k-\ell) \big\|_{\Sigma}^2$

- Computable : c(y) = c_value(y, Σ, A, k, C, 1)
 Some analytical challenges:
 - Non-asymptotic error control [Berry, 1941]
 - Conservatism

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

• Prior:
$$\theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2)$$
 for $\beta \in \mathbb{R}^D$

- ► Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \operatorname{diag}(\sigma^2/\mathsf{size}_1, \dots, \sigma^2/\mathsf{size}_N)$
- Alternative $\theta^*(y) = [I_N + \tau^{-2}\Sigma]^{-1}y + [I_N + \tau^2\Sigma^{-1}]^{-1}X\beta$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$

• Default:
$$\hat{\theta}(y) = y$$

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

• Prior:
$$\theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2)$$
 for $\beta \in \mathbb{R}^D$

- ► Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \text{diag}(\sigma^2/\text{size}_1, \dots, \sigma^2/\text{size}_N)$
- Alternative $\theta^*(y) = [I_N + \tau^{-2}\Sigma]^{-1}y + [I_N + \tau^2\Sigma^{-1}]^{-1}X\beta$

Compute $c_value(y, \Sigma, A, k, C, 1)$

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$
- Default: $\hat{\theta}(y) = y$ [$A = I_N, k = 0$]

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

• Prior:
$$\theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2)$$
 for $\beta \in \mathbb{R}^D$

► Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \operatorname{diag}(\sigma^2/\operatorname{size}_1, \dots, \sigma^2/\operatorname{size}_N)$

Alternative
$$\theta^*(y) = \underbrace{[I_N + \tau^{-2}\Sigma]^{-1}}_C y + \underbrace{[I_N + \tau^2\Sigma^{-1}]^{-1}X\beta}_\ell$$

Compute c_value(y, Σ , A, k, C, 1)

Educational Longitudinal Study (2002–2012)

- Standardized test of reading ability in 10th grade students
- Sample of 5-50 students at N = 676 schools $(y = [y_1, \dots, y_N])$
- **• Goal:** estimate school-specific means, $\theta \in \mathbb{R}^N$
- ▶ Default: $\hat{\theta}(y) = y$ [$A = I_N, k = 0$]

Small area inference [Fay and Herriot, 1979, Hoff, 2020]

▶ D = 8 features of each school (region, school type, enrollment, ...)

• Prior:
$$\theta_n \overset{indep}{\sim} \mathcal{N}(x_n^\top \beta, \tau^2)$$
 for $\beta \in \mathbb{R}^D$

► Likelihood: $y \sim \mathcal{N}(\theta, \Sigma)$ for $\Sigma = \operatorname{diag}(\sigma^2/\operatorname{size}_1, \dots, \sigma^2/\operatorname{size}_N)$

Alternative
$$\theta^*(y) = \underbrace{[I_N + \tau^{-2}\Sigma]^{-1}}_C y + \underbrace{[I_N + \tau^2\Sigma^{-1}]^{-1}X\beta}_\ell$$

Compute c_value(y, Σ , A, k, C, 1)=0.9926

Beyond Affine Estimates & Gaussian Noise

Beyond Affine Estimates & Gaussian Noise

Many likelihoods are approximately Gaussian
Many likelihoods are approximately Gaussian

► E.g. Logistic Regression

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- Asymptotic normality of MLE

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- \blacktriangleright Asymptotic normality of MLE \rightarrow Gaussian approximation to likelihood

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- Asymptotic normality of MLE \rightarrow Gaussian approximation to likelihood
- \blacktriangleright We show: our bounds provide nominal coverage as sample size $ightarrow\infty$

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- \blacktriangleright Asymptotic normality of MLE \rightarrow Gaussian approximation to likelihood
- \blacktriangleright We show: our bounds provide nominal coverage as sample size $ightarrow\infty$

Many estimates are approximately affine

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- \blacktriangleright Asymptotic normality of MLE \rightarrow Gaussian approximation to likelihood
- \blacktriangleright We show: our bounds provide nominal coverage as sample size $ightarrow\infty$

Many estimates are approximately affine

Empirical Bayes

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- \blacktriangleright Asymptotic normality of MLE \rightarrow Gaussian approximation to likelihood
- \blacktriangleright We show: our bounds provide nominal coverage as sample size $ightarrow\infty$

Many estimates are approximately affine

Empirical Bayes

► E.g. James–Stein estimator:
$$\theta_{\mathsf{JS}}^*(y) = \left(1 - \frac{N-2}{\|y\|^2}\right) y$$

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- Asymptotic normality of MLE \rightarrow Gaussian approximation to likelihood
- \blacktriangleright We show: our bounds provide nominal coverage as sample size $ightarrow\infty$

Many estimates are approximately affine

- Empirical Bayes
- E.g. James–Stein estimator: $\theta^*_{JS}(y) = \left(1 \frac{N-2}{\|y\|^2}\right)y$
- **We show:** our bounds provide nominal coverage as dimension $N \to \infty$

Many likelihoods are approximately Gaussian

- E.g. Logistic Regression
- Asymptotic normality of MLE \rightarrow Gaussian approximation to likelihood
- \blacktriangleright We show: our bounds provide nominal coverage as sample size $ightarrow\infty$

Many estimates are approximately affine

- Empirical Bayes
- E.g. James–Stein estimator: $\theta^*_{JS}(y) = \left(1 \frac{N-2}{\|y\|^2}\right)y$
- We show: our bounds provide nominal coverage as dimension $N \to \infty$

Open Directions:

- 1. Different losses L1, zero-one
- 2. Different models sparse regression
- 3. Tighter bounds overly conservative

Summary

- We proposed c-values to frequentist confidence in new estimates
 - on the observed dataset
 - without assumptions on θ
- Our bounds cover a range of models & estimates for squared error
- We demonstrate conclusive evaluations on real problems

Further Information

Trippe, Brian L., Sameer K. Deshpande, and Tamara Broderick. "Confidently Comparing Estimators with the c-value." *Journal of the American Statistical Association* (2023).

Code Available: github.com/blt2114/c_values

Contact me: btrippe@mit.edu

References

- Andrew C Berry. The accuracy of the Gaussian approximation to the sum of independent variates. *Transactions of the American Mathematical Society*, 49(1):122–136, 1941.
- Robert E Fay and Roger A Herriot. Estimates of income for small places: an application of James-Stein procedures to census data. *Journal of the American Statistical Association*, 74(366a):269–277, 1979.
- Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2013.
- Peter D Hoff. Smaller *p*-values via indirect information. *Journal of the American Statistical Association*, pages 1–35, 2020.
- Dennis V Lindley and Adrian FM Smith. Bayes estimates for the linear model. Journal of the Royal Statistical Society: Series B, 34(1):1–18, 1972.
- Donald B Rubin. Estimation in parallel randomized experiments. *Journal of Educational Statistics*, 6(4):377–401, 1981.

Coverage of Empirical Bayes

