Using Bagged Posteriors for Robust Inference

Jonathan Huggins Harvard University

Joint work with Jeff Miller

• Goal: predict future insurance claims based on (real) historic data

- Goal: predict future insurance claims based on (real) historic data
- Try Bayesian inference (non-trivial) model (data is 10 time series)



- Goal: predict future insurance claims based on (real) historic data
- Try Bayesian inference (non-trivial) model (data is 10 time series)
- **Problem:** uncertainty not well-calibrated because model is wrong



- Goal: predict future insurance claims based on (real) historic data
- Try Bayesian inference (non-trivial) model (data is 10 time series)
- Problem: uncertainty not well-calibrated because model is wrong
- Alternative: the bootstrap \Rightarrow too little data



- Goal: predict future insurance claims based on (real) historic data
- Try Bayesian inference (non-trivial) model (data is 10 time series)
- **Problem:** uncertainty not well-calibrated because model is wrong
- Alternative: the bootstrap \Rightarrow too little data
- Solution: use the bagged posterior (BayesBag)



Agenda

- BayesBag for parameter inference (and prediction)
- BayesBag methodology
- BayesBag for model selection

 Goal: learn about unobserved phenomenon (parameter) of interest θ [e.g. future claims]

- Goal: learn about unobserved phenomenon (parameter) of interest θ [e.g. future claims]
- **Prior** beliefs $\pi_0(\theta)$ about the phenomenon



- Goal: learn about unobserved phenomenon (parameter) of interest θ [e.g. future claims]
- **Prior** beliefs $\pi_0(\theta)$ about the phenomenon
- Observe data Y via **model** $p(Y \mid \theta)$



- Goal: learn about unobserved phenomenon (parameter) of interest θ [e.g. future claims]
- **Prior** beliefs $\pi_0(\theta)$ about the phenomenon
- Observe data Y via **model** $p(Y \mid \theta)$
- Combine prior & likelihood to form **posterior**:

 $\pi(\theta \mid Y) \propto p(Y \mid \theta) \pi_0(\theta)$



- Goal: learn about unobserved phenomenon (parameter) of interest θ [e.g. future claims]
- **Prior** beliefs $\pi_0(\theta)$ about the phenomenon
- Observe data Y via **model** $p(Y \mid \theta)$
- Combine prior & likelihood to form **posterior**:

 $\pi(\theta \mid Y) \propto p(Y \mid \theta) \pi_0(\theta)$

• **Benefits:** coherent belief updates, uncertainty quantification, flexible modeling, and more



- Goal: learn about unobserved phenomenon (parameter) of interest θ [e.g. future claims]
- **Prior** beliefs $\pi_0(\theta)$ about the phenomenon
- Observe data Y via **model** $p(Y \mid \theta)$
- Combine prior & likelihood to form **posterior**:

 $\pi(\theta \mid Y) \propto p(Y \mid \theta) \pi_0(\theta)$

- **Benefits:** coherent belief updates, uncertainty quantification, flexible modeling, and more
- Assumption #1: measurement model correct: observed Y has distribution $p(Y | \theta_{true})$



- Goal: learn about unobserved phenomenon (parameter) of interest θ [e.g. future claims]
- **Prior** beliefs $\pi_0(\theta)$ about the phenomenon
- Observe data Y via **model** $p(Y \mid \theta)$
- Combine prior & likelihood to form **posterior**:

 $\pi(\theta \,|\, Y) \propto p(Y \,|\, \theta) \pi_0(\theta)$

- **Benefits:** coherent belief updates, uncertainty quantification, flexible modeling, and more
- Assumption #1: measurement model correct: observed Y has distribution $p(Y | \theta_{true})$
- Assumption #2: Prior puts sufficient mass on true parameter θ_{true}



• Have data $Y = (Y_1, \ldots, Y_n)$



- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$



- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$
- Bootstrap dataset $Y^* = (Y_1^*, \dots, Y_m^*)$, where Y_i^* i.i.d. $\sim P_n$



- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$ • Bootstrap dataset $Y^* = (Y_1^*, \dots, Y_m^*)$, where Y_i^* i.i.d. $\sim P_n$ Y^* not always equal to n!

- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$
- Bootstrap dataset $Y^* = (Y_1^*, \dots, Y_n^*)$, where Y_i^* i.i.d. $\sim P_n$
- Bagged posterior $\pi^*(\theta \mid Y) = \mathbb{E}\{\pi(\theta \mid Y^*) \mid Y\}$

Y
not always equal to n!

Y

 P_n

- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$
- Bootstrap dataset $Y^* = (Y_1^*, \dots, Y_n^*)$, where Y_i^* i.i.d. $\sim P_n$
- Bagged posterior $\pi^*(\theta \mid Y) = \mathbb{E}\{\pi(\theta \mid Y^*) \mid Y\}$

- We show: when m = n, **conservative** uncertainty even when model wrong

- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$
- Bootstrap dataset $Y^* = (Y_1^*, \dots, Y_n^*)$, where Y_i^* i.i.d. $\sim P_n$
- Bagged posterior $\pi^*(\theta \mid Y) = \mathbb{E}\{\pi(\theta \mid Y^*) \mid Y\}$



- We show: when m = n, **conservative** uncertainty even when model wrong
- In practice, sample B bootstrap datasets: $\pi^*(\theta \mid Y) \approx \frac{1}{B} \sum_{b=1}^{B} \pi(\theta \mid Y_{(b)}^*)$

- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$
- Bootstrap dataset $Y^* = (Y_1^*, \dots, Y_n^*)$, where Y_i^* i.i.d. $\sim P_n$
- Bagged posterior $\pi^*(\theta \mid Y) = \mathbb{E}\{\pi(\theta \mid Y^*) \mid Y\}$



- We show: when m = n, **conservative** uncertainty even when model wrong
- In practice, sample B bootstrap datasets: $\pi^*(\theta \mid Y) \approx \frac{1}{B} \sum_{b=1}^{B} \pi(\theta \mid Y_{(b)}^*)$
- Suffices to take B = 50 or 100

- Have data $Y = (Y_1, \ldots, Y_n)$
- Empirical data distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$
- Bootstrap dataset $Y^* = (Y_1^*, \dots, Y_n^*)$, where Y_i^* i.i.d. $\sim P_n$
- Bagged posterior $\pi^*(\theta \mid Y) = \mathbb{E}\{\pi(\theta \mid Y^*) \mid Y\}$



- We show: when m = n, **conservative** uncertainty even when model wrong
- In practice, sample B bootstrap datasets: $\pi^*(\theta \mid Y) \approx \frac{1}{B} \sum_{b=1}^{B} \pi(\theta \mid Y_{(b)}^*)$
- Suffices to take B = 50 or 100
- Benefits: easy to use, can parallelize across B

 P_n

• Assume $\theta \in \mathbb{R}^d$

- Assume $\theta \in \mathbb{R}^d$
- Assume independent observations, so posterior is $\pi(\theta \mid Y) \propto \pi_0(\theta) \prod_{i=1}^n p_{\theta}(Y_i)$

- Assume $\theta \in \mathbb{R}^d$
- Assume independent observations, so posterior is $\pi(\theta \mid Y) \propto \pi_0(\theta) \prod_{i=1}^n p_{\theta}(Y_i)$
- Assume data Y_1, \ldots, Y_n i.i.d. $\sim P_{\mathsf{true}}$

- Assume $\theta \in \mathbb{R}^d$
- Assume independent observations, so posterior is $\pi(\theta \mid Y) \propto \pi_0(\theta) \prod_{i=1}^n p_{\theta}(Y_i)$
- Assume data Y_1, \ldots, Y_n i.i.d. $\sim P_{\mathsf{true}}$
- From now on, we consider the large $n \, \operatorname{limit}$

- Assume $\theta \in \mathbb{R}^d$
- Assume independent observations, so posterior is $\pi(\theta \mid Y) \propto \pi_0(\theta) \prod_{i=1}^n p_{\theta}(Y_i)$
- Assume data Y_1, \ldots, Y_n i.i.d. $\sim P_{\mathsf{true}}$
- From now on, we consider the large $n \ {\rm limit}$
- We will be ignoring the finite-sample benefits of Bayes

- Assume $\theta \in \mathbb{R}^d$
- Assume independent observations, so posterior is $\pi(\theta \mid Y) \propto \pi_0(\theta) \prod_{i=1}^n p_{\theta}(Y_i)$
- Assume data Y_1, \ldots, Y_n i.i.d. $\sim P_{\mathsf{true}}$
- From now on, we consider the large $n \ {\rm limit}$
- We will be ignoring the finite-sample benefits of Bayes
- But, we still get many useful insights!

• Interested in quantifying uncertainty about **KL-optimal parameter**

$$\theta_{\mathsf{opt}} = \arg\max_{\theta} \mathbb{E}\{\log p_{\theta}(Y_1)\}$$

• Interested in quantifying uncertainty about KL-optimal parameter

$$\theta_{\mathsf{opt}} = \arg\max_{\theta} \mathbb{E}\{\log p_{\theta}(Y_1)\}$$

• Why? $p_{\theta_{opt}}$ best explains the data amongst the models $\{p_{\theta} : \theta \in \mathbb{R}^d\}$

• Interested in quantifying uncertainty about KL-optimal parameter

$$\theta_{\mathsf{opt}} = \arg\max_{\theta} \mathbb{E}\{\log p_{\theta}(Y_1)\}$$

- Why? $p_{\theta_{opt}}$ best explains the data amongst the models $\{p_{\theta} : \theta \in \mathbb{R}^d\}$
- When model correctly specified, $p_{\theta_{opt}} = P_{true}$
Uncertainty about the optimal parameter

• Interested in quantifying uncertainty about KL-optimal parameter

$$\theta_{\mathsf{opt}} = \arg\max_{\theta} \mathbb{E}\{\log p_{\theta}(Y_1)\}$$

- Why? $p_{\theta_{opt}}$ best explains the data amongst the models $\{p_{\theta} : \theta \in \mathbb{R}^d\}$
- When model correctly specified, $p_{\theta_{opt}} = P_{true}$
- Define the maximum likelihood estimator

Uncertainty about the optimal parameter

• Interested in quantifying uncertainty about KL-optimal parameter

$$\theta_{\mathsf{opt}} = \arg\max_{\theta} \mathbb{E}\{\log p_{\theta}(Y_1)\}$$

- Why? $p_{\theta_{opt}}$ best explains the data amongst the models $\{p_{\theta}: \theta \in \mathbb{R}^d\}$
- When model correctly specified, $p_{\theta_{\rm opt}} = P_{\rm true}$
- Define the maximum likelihood estimator

$$\hat{\theta}(Y) = \arg\max_{\theta} \prod_{i=1}^{n} p_{\theta}(Y_i)$$

Uncertainty about the optimal parameter

• Interested in quantifying uncertainty about KL-optimal parameter

$$\theta_{\mathsf{opt}} = \arg\max_{\theta} \mathbb{E}\{\log p_{\theta}(Y_1)\}$$

- Why? $p_{\theta_{opt}}$ best explains the data amongst the models $\{p_{\theta} : \theta \in \mathbb{R}^d\}$
- When model correctly specified, $p_{\theta_{\text{opt}}} = P_{\text{true}}$
- Define the maximum likelihood estimator

$$\hat{\theta}(Y) = \arg\max_{\theta} \prod_{i=1}^{n} p_{\theta}(Y_i)$$

• Under mild conditions, $\hat{\theta}(Y) \to \theta_{\mathsf{opt}}$ and $\pi(\cdot \mid Y) \stackrel{\mathcal{D}}{\Longrightarrow} \delta_{\theta_{\mathsf{opt}}}$ as $n \to \infty$

[Kleijn & van der Vaart 2012] For $\vartheta \sim \pi(\cdot \,|\, Y)$,

 $\vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \qquad \left(\mathsf{that is, } n^{1/2} \{\vartheta - \hat{\theta}(Y)\} \stackrel{\mathcal{D}}{\Longrightarrow} \mathcal{N}(0, \Sigma_{\mathsf{M}})\right)$

[Kleijn & van der Vaart 2012] For $\vartheta \sim \pi(\cdot \,|\, Y)$,

$$\vartheta \mid Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \qquad \left(\text{that is, } n^{1/2} \{ \vartheta - \hat{\theta}(Y) \} \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_{\mathsf{M}}) \right)$$

"model" covariance

[Kleijn & van der Vaart 2012] For $\vartheta \sim \pi(\cdot \,|\, Y)$,

$$\vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \sum_{\mathsf{M}}/n) \qquad \left(\text{that is, } n^{1/2} \{ \vartheta - \hat{\theta}(Y) \} \stackrel{\mathcal{D}}{\Longrightarrow} \mathcal{N}(0, \Sigma_{\mathsf{M}}) \right)$$

"model" covariance

[Huber 1967, White 1982] Sampling distribution of the MLE satisfies

 $\hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n)$

[Kleijn & van der Vaart 2012] For $\vartheta \sim \pi(\cdot \,|\, Y)$,

$$\vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \sum_{\mathsf{M}}/n) \qquad \left(\mathsf{that is, } n^{1/2} \{ \vartheta - \hat{\theta}(Y) \} \stackrel{\mathcal{D}}{\Longrightarrow} \mathcal{N}(0, \Sigma_{\mathsf{M}}) \right)$$

"model" covariance

[Huber 1967, White 1982] Sampling distribution of the MLE satisfies

$$\hat{\theta}(Y) \approx \mathcal{N}(\theta_{\text{opt}}, \sum_{S}/n)$$

[Kleijn & van der Vaart 2012] For $\vartheta \sim \pi(\cdot \,|\, Y)$,

$$\vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \sum_{\mathsf{M}} / n) \qquad \left(\mathsf{that is, } n^{1/2} \{ \vartheta - \hat{\theta}(Y) \} \stackrel{\mathcal{D}}{\Longrightarrow} \mathcal{N}(0, \Sigma_{\mathsf{M}}) \right)$$

"model" covariance

[Huber 1967, White 1982] Sampling distribution of the MLE satisfies

$$\hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \underline{\Sigma_{\mathsf{S}}}/n)$$

$$sandwich \ covariance$$

• If $P_{true} = p_{\theta_{opt}}$, then $\Sigma_{S} = \Sigma_{M}$ and posterior uncertainty about θ_{opt} correct

[Kleijn & van der Vaart 2012] For $\vartheta \sim \pi(\cdot \,|\, Y)$,

$$\vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \underbrace{\Sigma_{\mathsf{M}}/n}) \qquad \left(\mathsf{that is, } n^{1/2} \{ \vartheta - \hat{\theta}(Y) \} \stackrel{\mathcal{D}}{\Longrightarrow} \mathcal{N}(0, \Sigma_{\mathsf{M}}) \right)$$
 "model" covariance

[Huber 1967, White 1982] Sampling distribution of the MLE satisfies

$$\hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \sum_{\mathsf{S}}/n)$$
sandwich covariance

• If $P_{\text{true}} = p_{\theta_{\text{opt}}}$, then $\Sigma_{\text{S}} = \Sigma_{\text{M}}$ and posterior uncertainty about θ_{opt} correct

• If model misspecified, then in general $\Sigma_{S} \neq \Sigma_{M}$ and posterior uncertainty about θ_{opt} could be very wrong

MLE: Posterior:
$$\begin{split} \hat{\theta}(Y) &\approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n) \\ \vartheta \,|\, Y &\approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \end{split}$$

 $\begin{array}{ll} \mathsf{MLE:} & \hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n) \\ \\ \mathsf{Posterior:} & \vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \end{array}$

[van der Vaart & Wellner 1996] Bootstrap distribution of the MLE satisfies

 $\hat{\theta}(Y^*) \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{S}}/m)$

 $\begin{array}{ll} \mathsf{MLE:} & \hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n) \\ \\ \mathsf{Posterior:} & \vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \end{array} \end{array}$

[van der Vaart & Wellner 1996] Bootstrap distribution of the MLE satisfies

$$\hat{\theta}(Y^*) \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{S}}/m)$$

- With m = n, the bootstrap
 - quantifies sampling variability

 $\begin{array}{ll} \mathsf{MLE:} & \hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n) \\ \\ \mathsf{Posterior:} & \vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \end{array} \end{array}$

[van der Vaart & Wellner 1996] Bootstrap distribution of the MLE satisfies

$$\hat{\theta}(Y^*) \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{S}}/m)$$

- With m = n, the bootstrap
 - quantifies sampling variability
 - \circ provides correct frequentist confidence intervals

MLE: Posterior: Bootstrap:
$$\begin{split} \hat{\theta}(Y) &\approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n) \\ \vartheta \,|\, Y &\approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \\ \hat{\theta}(Y^*) \,|\, Y &\approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{S}}/m) \end{split}$$



[H & Miller 2019] BayesBag includes sampling and model-based uncertainty: for $\vartheta^* \sim \pi^*(\cdot \,|\, Y)$,

 $\vartheta^* \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m)$



[H & Miller 2019] BayesBag includes sampling and model-based uncertainty: for $\vartheta^* \sim \pi^*(\cdot \,|\, Y)$,

 $\vartheta^* \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m)$

• If model correct, m = 2n yields correct uncertainty:

 $\operatorname{Cov}(\vartheta^* \,|\, Y) \approx \Sigma_{\mathsf{S}}/n$

$$\begin{array}{ll} \mathsf{MLE:} & \hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n) \\ \\ \mathsf{Posterior:} & \vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \\ \\ \\ \mathsf{Bootstrap:} & \hat{\theta}(Y^*) \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{S}}/m) \end{array}$$

[H & Miller 2019] BayesBag includes sampling and model-based uncertainty: for $\vartheta^* \sim \pi^*(\cdot \,|\, Y)$,

$$\vartheta^* \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m)$$

• If model correct, m = 2n yields correct uncertainty:

$$\operatorname{Cov}(\vartheta^* | Y) \approx \Sigma_{\mathsf{S}}/n$$

• No matter what, m = n ensures **conservative uncertainty**:

 $\operatorname{Cov}(\vartheta^* \,|\, Y) \ge \Sigma_{\mathsf{S}}/n$

$$\begin{array}{ll} \mathsf{MLE:} & \hat{\theta}(Y) \approx \mathcal{N}(\theta_{\mathsf{opt}}, \Sigma_{\mathsf{S}}/n) \\ \\ \mathsf{Posterior:} & \vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \\ \\ \\ \mathsf{Bootstrap:} & \hat{\theta}(Y^*) \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{S}}/m) \end{array}$$

[H & Miller 2019] BayesBag includes sampling and model-based uncertainty: for $\vartheta^* \sim \pi^*(\cdot \,|\, Y)$,

$$\vartheta^* \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m)$$

• If model correct, m = 2n yields correct uncertainty:

$$\operatorname{Cov}(\vartheta^* | Y) \approx \Sigma_{\mathsf{S}}/n$$

• No matter what, m = n ensures **conservative uncertainty**:

$$\operatorname{Cov}(\vartheta^* | Y) \ge \Sigma_{\mathsf{S}}/n$$

• Remember: we are ignoring finite-sample benefits of Bayes

Agenda

- BayesBag for parameter inference (and prediction)
- BayesBag methodology
- BayesBag for model selection

 $\begin{array}{ll} \text{Posterior:} & \vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \\ \text{BayesBag:} & \vartheta^* \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m) \end{array}$

- Posterior: $\vartheta \mid Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n)$ BayesBag: $\vartheta^* \mid Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m)$
- $\operatorname{Var}(\vartheta) = \operatorname{posterior variance} \approx \Sigma_{\mathsf{M}}/n$

- $\begin{array}{ll} \text{Posterior:} & \vartheta \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n) \\ \text{BayesBag:} & \vartheta^* \,|\, Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m) \end{array}$
- $\operatorname{Var}(\vartheta) = \operatorname{posterior variance} \approx \Sigma_{\mathsf{M}}/n$
- $\operatorname{Var}(\vartheta^*) = \operatorname{BayesBag} \operatorname{variance} \approx (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/n$





• optimal bootstrap sample size m_{opt} estimate

$$\hat{m}_{\mathsf{opt}} = \frac{\operatorname{Var}(\vartheta^*)}{\operatorname{Var}(\vartheta^*) - \operatorname{Var}(\vartheta)} \times n$$

Posterior:
$$\vartheta \mid Y \approx \mathcal{N}(\hat{\theta}(Y), \Sigma_{\mathsf{M}}/n)$$
BayesBag: $\vartheta^* \mid Y \approx \mathcal{N}(\hat{\theta}(Y), (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/m)$ • $\operatorname{Var}(\vartheta) = \operatorname{posterior variance} \approx \Sigma_{\mathsf{M}}/n$ Computable• $\operatorname{Var}(\vartheta^*) = \operatorname{BayesBag variance} \approx (\Sigma_{\mathsf{M}} + \Sigma_{\mathsf{S}})/n$

• optimal bootstrap sample size m_{opt} estimate

$$\hat{m}_{\mathsf{opt}} = \frac{\operatorname{Var}(\vartheta^*)}{\operatorname{Var}(\vartheta^*) - \operatorname{Var}(\vartheta)} \times n$$

 $\rightarrow~$ finite-sample version also available using a Gaussian model approximation

1) compute standard posterior $\pi(\cdot \mid Y)$

1) compute standard posterior $\pi(\cdot \,|\, Y)$

– e.g., use MCMC to get approximate samples $\theta_{(1)},\ldots,\theta_{(T)}$ from $\pi(\cdot\,|\,Y)$

1) compute standard posterior $\pi(\cdot \,|\, Y)$

– e.g., use MCMC to get approximate samples $heta_{(1)},\ldots, heta_{(T)}$ from $\pi(\cdot\,|\,Y)$

2) compute bagged posterior $\pi^*(\cdot \,|\, Y)$ with m=n using $B\approx 50$ bootstrap datasets

1) compute standard posterior $\pi(\cdot \,|\, Y)$

– e.g., use MCMC to get approximate samples $heta_{(1)},\ldots, heta_{(T)}$ from $\pi(\cdot\,|\,Y)$

2) compute bagged posterior $\pi^*(\cdot \,|\, Y)$ with m=n using $B\approx 50$ bootstrap datasets

- e.g., use MCMC to get approximate samples $\theta^*_{(b,1)}, \ldots, \theta^*_{(b,T)}$ from $\pi(\cdot \mid Y^*_{(b)})$ for $b = 1, \ldots, B$

1) compute standard posterior $\pi(\cdot \,|\, Y)$

– e.g., use MCMC to get approximate samples $heta_{(1)},\ldots, heta_{(T)}$ from $\pi(\cdot\,|\,Y)$

2) compute bagged posterior $\pi^*(\cdot \,|\, Y)$ with m=n using $B\approx 50$ bootstrap datasets

- e.g., use MCMC to get approximate samples $\theta^*_{(b,1)},\ldots,\theta^*_{(b,T)}$ from $\pi(\cdot\,|\,Y^*_{(b)})$ for $b=1,\ldots,B$

if Gaussian approximation to standard and bagged posteriors decent then

1) compute standard posterior $\pi(\cdot \,|\, Y)$

– e.g., use MCMC to get approximate samples $heta_{(1)},\ldots, heta_{(T)}$ from $\pi(\cdot\,|\,Y)$

2) compute bagged posterior $\pi^*(\cdot \,|\, Y)$ with m=n using $B\approx 50$ bootstrap datasets

- e.g., use MCMC to get approximate samples $\theta^*_{(b,1)},\ldots,\theta^*_{(b,T)}$ from $\pi(\cdot\,|\,Y^*_{(b)})$ for $b=1,\ldots,B$

if Gaussian approximation to standard and bagged posteriors decent then

3a) compute optimal bootstrap sample size estimate \hat{m}_{opt}

3b) compute bagged posterior $\pi^*(\cdot \mid Y)$ with $m = \hat{m}_{opt}$ and output

1) compute standard posterior $\pi(\cdot \,|\, Y)$

– e.g., use MCMC to get approximate samples $heta_{(1)},\ldots, heta_{(T)}$ from $\pi(\cdot\,|\,Y)$

2) compute bagged posterior $\pi^*(\cdot \,|\, Y)$ with m=n using $B\approx 50$ bootstrap datasets

- e.g., use MCMC to get approximate samples $\theta^*_{(b,1)}, \ldots, \theta^*_{(b,T)}$ from $\pi(\cdot \mid Y^*_{(b)})$ for $b = 1, \ldots, B$

if Gaussian approximation to standard and bagged posteriors decent then

3a) compute optimal bootstrap sample size estimate \hat{m}_{opt}

3b) compute bagged posterior $\pi^*(\cdot \mid Y)$ with $m = \hat{m}_{\rm opt}$ and output else

4) output bagged posterior with m = n computed in step 2

Diagnosing model-data mismatch

Diagnosing model-data mismatch

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
- Recall:
 - \circ -expect optimal bootstrap dataset size to be $m_{\rm opt}=2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty
- Define model-data mismatch index $\mathcal{I} = 2n/\hat{m}_{opt} 1$

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty
- Define model-data mismatch index $\mathcal{I} = 2n/\hat{m}_{opt} 1$
- Provides interpretable diagnostic for **model criticism**

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty
- Define model-data mismatch index $\mathcal{I} = 2n/\hat{m}_{opt} 1$
- Provides interpretable diagnostic for **model criticism**

 $\circ \ -1 < \mathcal{I} < 1$

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty
- Define model-data mismatch index $\mathcal{I} = 2n/\hat{m}_{opt} 1$
- Provides interpretable diagnostic for **model criticism**

 $\circ \quad -1 < \mathcal{I} < 1$

 $\circ~\mathcal{I} \approx 0$: no disagreement

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty
- Define model-data mismatch index $\mathcal{I} = 2n/\hat{m}_{opt} 1$
- Provides interpretable diagnostic for **model criticism**

 $\circ \quad -1 < \mathcal{I} < 1$

- $\circ~\mathcal{I} \approx 0$: no disagreement
- $\circ \mathcal{I} > 0$: posterior overconfident

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty
- Define model-data mismatch index $\mathcal{I} = 2n/\hat{m}_{opt} 1$
- Provides interpretable diagnostic for **model criticism**

 $\circ \quad -1 < \mathcal{I} < 1$

- $\circ~\mathcal{I} \approx 0$: no disagreement
- $\circ \ \mathcal{I} > 0$: posterior overconfident
- $\circ \ \mathcal{I} < 0$: posterior under-confident

- Recall:
 - \circ expect optimal bootstrap dataset size to be $m_{opt} = 2n$ if model is correct
 - $\circ~$ expect optimal bootstrap dataset size to be $m_{\rm opt}\approx n$ if posterior is drastically underestimating uncertainty
- Define model-data mismatch index $\mathcal{I} = 2n/\hat{m}_{opt} 1$
- Provides interpretable diagnostic for model criticism

 $\circ \quad -1 < \mathcal{I} < 1$

- $\circ~\mathcal{I} \approx 0$: no disagreement
- $\circ \mathcal{I} > 0$: posterior overconfident
- $\circ \quad \mathcal{I} < 0: \text{ posterior under-confident}$



• Assumed model: linear regression with a conjugate prior on $\theta = (\sigma^2, \beta)$

$$Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2) \qquad Y_i = (X_i, Z_i)$$

• Assumed model: linear regression with a conjugate prior on $\theta = (\sigma^2, \beta)$

$$Z_{i} \mid X_{i}, \beta, \sigma^{2} \sim \mathcal{N}(X_{i}^{\top}\beta, \sigma^{2}) \qquad Y_{i} = (X_{i}, Z_{i})$$
outcome

• Assumed model: linear regression with a conjugate prior on $\theta = (\sigma^2, \beta)$

$$\begin{array}{c} Z_{i} \mid X_{i}, \beta, \sigma^{2} \sim \mathcal{N}(X_{i}^{\intercal}\beta, \sigma^{2}) \\ \uparrow \\ outcome \end{array} \quad \begin{array}{c} Y_{i} = (X_{i}, Z_{i}) \\ \hline \\ covariates \end{array}$$

• Assumed model: linear regression with a conjugate prior on $\theta = (\sigma^2, \beta)$

$$Z_{i} | X_{i}, \beta, \sigma^{2} \sim \mathcal{N}(X_{i}^{\mathsf{T}}\beta, \sigma^{2}) \qquad Y_{i} = (X_{i}, Z_{i})$$
outcome
$$C_{i} = (X_{i}, Z_{i})$$

• True model P_{true} (that we simulate from):

$$Z_i \mid X_i \sim \mathcal{N}(f(X_i)^\top \beta_{\mathsf{gen}}, 1) \qquad \qquad X_i \sim G$$

• Assumed model: linear regression with a conjugate prior on $\theta = (\sigma^2, \beta)$



• True model P_{true} (that we simulate from):

$$Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{\mathsf{gen}}, 1) \qquad \qquad X_i \sim G$$

• n = 50, d = 10, and $m = \hat{m}_{opt}$

• Assumed model: linear regression with a conjugate prior on $\theta = (\sigma^2, \beta)$



• True model P_{true} (that we simulate from):

$$Z_i \mid X_i \sim \mathcal{N}(f(X_i)^\top \beta_{\mathsf{gen}}, 1) \qquad \qquad X_i \sim G$$

- n = 50, d = 10, and $m = \hat{m}_{opt}$
- Important: when model misspecified, $\theta_{opt} \neq (1, \beta_{gen})$ (in general)

Assumed model: $Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$ True model: $Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \quad X_i \sim G$

Assumed model: $Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$ True model: $Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \quad X_i \sim G$

• Performance metric is difference in log posterior density at θ_{opt} :

 $\log \pi^*(\theta_{\mathsf{opt}} \,|\, Y) - \log \pi(\theta_{\mathsf{opt}} \,|\, Y)$

Assumed model:
$$Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$$

True model: $Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \quad X_i \sim G$



Assumed model:
$$Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$$

True model: $Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \quad X_i \sim G$



Assumed model:
$$Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$$

True model: $Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \quad X_i \sim G$



Assumed model:
$$Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$$

True model: $Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \quad X_i \sim G$



Assumed model: $Z_i | X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$ True model: $Z_i | X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \quad X_i \sim G$

Assumed model: True model:

$$\begin{aligned} Z_i \,|\, X_i, \beta, \sigma^2 &\sim \mathcal{N}(X_i^\top \beta, \sigma^2) \\ Z_i \,|\, X_i &\sim \mathcal{N}(f(X_i)^\top \beta_{\mathsf{gen}}, 1), \qquad X_i \sim G \end{aligned}$$

Assumed model correct



Assumed model: True model:

$$\begin{split} Z_i \,|\, X_i, \beta, \sigma^2 &\sim \mathcal{N}(X_i^\top \beta, \sigma^2) \\ Z_i \,|\, X_i &\sim \mathcal{N}(f(X_i)^\top \beta_{\mathsf{gen}}, 1), \qquad X_i \sim G \end{split}$$



[**H** & Miller 2019]

Assumed model: True model:

 $Z_i \mid X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^{\top}\beta, \sigma^2)$ $Z_i \mid X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1), \qquad X_i \sim G$



Assumed model: True model:

is very wrong $Z_i \mid X_i, \beta, \sigma^2 \sim \mathcal{N}(X_i^\top \beta, \sigma^2)$ $Z_i \mid X_i \sim \mathcal{N}(f(X_i)^\top \beta_{gen}, 1),$ $X_i \sim G$



something

Agenda

- BayesBag for parameter inference (and prediction)
- BayesBag methodology
- BayesBag for model selection

 Goal: based on data *Y*, select between a (finite or countable) set of models *M* = {*m*₁, *m*₂, ...}

- Goal: based on data *Y*, select between a (finite or countable) set of models *M* = {*m*₁, *m*₂, ...}
- Example: systematics

- Goal: based on data *Y*, select between a (finite or countable) set of models *M* = {*m*₁, *m*₂, ...}
- Example: systematics
 - Goal: learn about evolutionary history of a set of species [e.g. whales]





- Goal: based on data *Y*, select between a (finite or countable) set of models *M* = {*m*₁, *m*₂, ...}
- Example: systematics
 - Goal: learn about evolutionary history of a set of species [e.g. whales]
 - Approach: infer which phylogenetic trees are consistent with observed species characteristics Y
 [e.g. genetic data, physical features such as coloring]



- Goal: based on data *Y*, select between a (finite or countable) set of models *M* = {*m*₁, *m*₂, ...}
- Example: systematics
 - Goal: learn about evolutionary history of a set of species [e.g. whales]
 - Approach: infer which phylogenetic trees are consistent with observed species characteristics Y
 [e.g. genetic data, physical features such as coloring]
- **Problem:** Bayesian model selection still assumes some model in *M* is correct



BayesBag stabilizes Bayesian model selection

• Assume two models m_1 and m_2

BayesBag stabilizes Bayesian model selection

- Assume two models m_1 and m_2
- Assume m_1 and m_2 explain the data-generating distribution equally well: $\mathbb{E}\{\log p(Y \mid m_1)\} = \mathbb{E}\{\log p(Y \mid m_2)\}$

BayesBag stabilizes Bayesian model selection

- Assume two models m_1 and m_2
- Assume m_1 and m_2 explain the data-generating distribution equally well: $\mathbb{E}\{\log p(Y \mid m_1)\} = \mathbb{E}\{\log p(Y \mid m_2)\}$
- Then, we hope models have equal posterior probability $(n \to \infty)$: $\pi(m_1 | Y) = \pi(m_2 | Y) = 1/2$
- Assume two models m_1 and m_2
- Assume m_1 and m_2 explain the data-generating distribution equally well: $\mathbb{E}\{\log p(Y \mid m_1)\} = \mathbb{E}\{\log p(Y \mid m_2)\}$
- Then, we hope models have equal posterior probability $(n \to \infty)$: $\pi(m_1 | Y) = \pi(m_2 | Y) = 1/2$
- However...

[Yang & Zhu 2018, H & Miller 2019] $\pi(m_1 \mid Y) = 0$ or 1 with equal probability.

- Assume two models m_1 and m_2
- Assume m_1 and m_2 explain the data-generating distribution equally well: $\mathbb{E}\{\log p(Y \mid m_1)\} = \mathbb{E}\{\log p(Y \mid m_2)\}$
- Then, we hope models have equal posterior probability $(n \to \infty)$: $\pi(m_1 | Y) = \pi(m_2 | Y) = 1/2$
- However...

[Yang & Zhu 2018, H & Miller 2019] $\pi(m_1 \mid Y) = 0$ or 1 with equal probability.

all posterior mass on a single, arbitrary model

- Assume two models m_1 and m_2
- Assume m_1 and m_2 explain the data-generating distribution equally well: $\mathbb{E}\{\log p(Y \mid m_1)\} = \mathbb{E}\{\log p(Y \mid m_2)\}$
- Then, we hope models have equal posterior probability $(n \to \infty)$: $\pi(m_1 | Y) = \pi(m_2 | Y) = 1/2$
- However...

[Yang & Zhu 2018, H & Miller 2019] $\pi(m_1 | Y) = 0$ or 1 with equal probability. [H & Miller 2019] A) If m = n, then $\pi^*(m_1 | Y) \sim \text{Uniform}(0, 1)$ B) If $m/n \rightarrow 0$, then $\pi^*(m_1 | Y) \rightarrow 1/2$.

- Assume two models m_1 and m_2
- Assume m_1 and m_2 explain the data-generating distribution equally well: $\mathbb{E}\{\log p(Y \mid m_1)\} = \mathbb{E}\{\log p(Y \mid m_2)\}$
- Then, we hope models have equal posterior probability $(n \to \infty)$: $\pi(m_1 | Y) = \pi(m_2 | Y) = 1/2$
- However...

[Yang & Zhu 2018, H & Miller 2019] $\pi(m_1 | Y) = 0$ or 1 with equal probability. [H & Miller 2019] A) If m = n, then $\pi^*(m_1 | Y) \sim \text{Uniform}(0, 1)$ B) If $m/n \rightarrow 0$, then $\pi^*(m_1 | Y) \rightarrow 1/2$. all posterior mass on a single, arbitrary model bagged posterior mass more evenly distributed

• Models are $m_1 = \mathcal{N}(-1, 1)$ and $m_2 = \mathcal{N}(1, 1)$

- Models are $m_1 = \mathcal{N}(-1, 1)$ and $m_2 = \mathcal{N}(1, 1)$
- True distribution is $P_{\text{true}} = \mathcal{N}(0, 1)$

- Models are $m_1 = \mathcal{N}(-1, 1)$ and $m_2 = \mathcal{N}(1, 1)$
- True distribution is $P_{\text{true}} = \mathcal{N}(0, 1)$



- Models are $m_1 = \mathcal{N}(-1, 1)$ and $m_2 = \mathcal{N}(1, 1)$
- True distribution is $P_{\text{true}} = \mathcal{N}(0, 1)$
- Generate datasets $Y^{(1)},Y^{(2)},\ldots$ of size n=1000, where $Y^{(i)}_j\sim \mathcal{N}(0,1).$



- Models are $m_1 = \mathcal{N}(-1, 1)$ and $m_2 = \mathcal{N}(1, 1)$
- True distribution is $P_{\text{true}} = \mathcal{N}(0, 1)$
- Generate datasets $Y^{(1)},Y^{(2)},\ldots$ of size n=1000 , where $Y^{(i)}_j\sim\mathcal{N}(0,1).$



$$\pi(m_1 | Y^{(1)}) = 1$$

$$\pi^*(m_1 | Y^{(1)}) = 0.82$$



- Models are $m_1 = \mathcal{N}(-1, 1)$ and $m_2 = \mathcal{N}(1, 1)$
- True distribution is $P_{\text{true}} = \mathcal{N}(0, 1)$
- Generate datasets $Y^{(1)},Y^{(2)},\ldots$ of size n=1000 , where $Y^{(i)}_j\sim \mathcal{N}(0,1).$





- Models are $m_1 = \mathcal{N}(-1, 1)$ and $m_2 = \mathcal{N}(1, 1)$
- True distribution is $P_{\text{true}} = \mathcal{N}(0, 1)$

2

- Generate datasets $Y^{(1)},Y^{(2)},\ldots$ of size n=1000 , where $Y^{(i)}_j\sim\mathcal{N}(0,1).$

$$\pi(m_1 | Y^{(1)}) = 1 \qquad \pi(m_1 | Y^{(2)}) = 10^{-5}$$

$$\pi^*(m_1 | Y^{(1)}) = 0.82 \qquad \pi^*(m_1 | Y^{(2)}) = 0.38$$

4

-4

-2

2



-4

-2

• **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA

all

Minke GACCCGAACGTAATAA...ATCCGTTCCCATACTC Blue CACCCCCCGTACTAT...TGAGTCCGAATTGGAA Fin TGTCTTCTACACTCCA...ACAGGTTGTACGTCAC Grey GGGTCGCTGTAGACCA...GATACCGCTCTCACAT

• **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA

all	1st half	
Minke	GACCCGAACGTAATAA	ATCCGTTCCCATACTC
Blue	CACCCCCCCGTACTAT	.TGAGTCCGAATTGGAA
Fin	TGTCTTCTACACTCCA.	.ACAGGTTGTACGTCAC
Grey	GGGTCGCTGTAGACCA	.GATACCGCTCTCACAT

• **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA

all	1st half	2nd half
Minke	GACCCGAACGTAATAA	ATCCGTTCCCATACTC
Blue	CACCCCCCGTACTAT	TGAGTCCGAATTGGAA
Fin	TGTCTTCTACACTCCA	ACAGGTTGTACGTCAC
Grey	GGGTCGCTGTAGACCA	GATACCGCTCTCACAT

• **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA





- **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA
- Compute posterior tree probabilities based on all,
 1st half, and 2nd half



- **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA
- Compute posterior tree probabilities based on all, 1st half, and 2nd half
- Compute overlap of 99%
 high probability regions



- **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA
- Compute posterior tree probabilities based on all, 1st half, and 2nd half
- Compute overlap of 99%
 high probability regions





- **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA
- Compute posterior tree probabilities based on all, 1st half, and 2nd half
- Compute overlap of 99%
 high probability regions





- **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA
- Compute posterior tree probabilities based on all, 1st half, and 2nd half
- Compute overlap of 99% high probability regions





- **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA
- Compute posterior tree probabilities based on all, 1st half, and 2nd half
- Compute overlap of 99%
 high probability regions





tree 3

- **Goal:** infer phylogeny of 13 whale species from mitochondrial coding DNA
- Compute posterior tree probabilities based on all, 1st half, and 2nd half
- Compute overlap of 99%
 high probability regions
- 0% overlap = contradiction

all

.2

tree 1

tree 2



.0

tree 4

0

1

.5

0.75

0.5

0.25

total

• Goal: infer phylogeny of 13 whale species

- Goal: infer phylogeny of 13 whale species
- For some evolutionary models, little to no overlap





- Goal: infer phylogeny of 13 whale species
- For some evolutionary models, little to no overlap





- Goal: infer phylogeny of 13 whale species
- For some evolutionary models, little to no overlap





- Goal: infer phylogeny of 13 whale species
- For some evolutionary models, little to no overlap
- Bayesian model selection is unstable and not reproducible [Wilcox et al. 2002, Alfaro et al. 2003, Douady et al. 2003, ...]





- Goal: infer phylogeny of 13 whale species
- For some evolutionary models, little to no overlap
- Bayesian model selection is **unstable** and **not reproducible** [Wilcox et al. 2002, Alfaro et al. 2003, Douady et al. 2003, ...]
- Same problem comparing evolutionary models with data fixed





- Goal: infer phylogeny of 13 whale species
- For some evolutionary models, little to no overlap
- Bayesian model selection is unstable and not reproducible [Wilcox et al. 2002, Alfaro et al. 2003, Douady et al. 2003, ...]
- Same problem comparing evolutionary models with data fixed
- Bagged posterior model probabilities more stable and reproducible



- Goal: infer phylogeny of 13 whale species
- For some evolutionary models, little to no overlap
- Bayesian model selection is unstable and not reproducible [Wilcox et al. 2002, Alfaro et al. 2003, Douady et al. 2003, ...]
- Same problem comparing evolutionary models with data fixed
- Bagged posterior model probabilities more stable and reproducible



• We show that BayesBag...

- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)

- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)
 - 2. more **stable** model probabilities (model selection)

- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)
 - 2. more **stable** model probabilities (model selection)
 - 3. demonstrates excellent empirical performance
- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)
 - 2. more **stable** model probabilities (model selection)
 - 3. demonstrates excellent empirical performance
 - 4. is easy to use and widely applicable

- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)
 - 2. more **stable** model probabilities (model selection)
 - 3. demonstrates excellent empirical performance
 - 4. is easy to use and widely applicable
 - 5. combines the **flexible modeling** features of Bayes with the **distributional robustness** of frequentist methods

- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)
 - 2. more **stable** model probabilities (model selection)
 - 3. demonstrates excellent empirical performance
 - 4. is easy to use and widely applicable
 - 5. combines the **flexible modeling** features of Bayes with the **distributional robustness** of frequentist methods
- Future work: finite-sample properties of BayesBag and mismatch index, computation, non-independent data (e.g. time series & spatial models)

- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)
 - 2. more **stable** model probabilities (model selection)
 - 3. demonstrates excellent empirical performance
 - 4. is easy to use and widely applicable
 - 5. combines the **flexible modeling** features of Bayes with the **distributional robustness** of frequentist methods
- Future work: finite-sample properties of BayesBag and mismatch index, computation, non-independent data (e.g. time series & spatial models)
- **Conclusion:** you should give BayesBag a try!

- We show that BayesBag...
 - 1. provides **better-calibrated uncertainty** that is conservative by default (parameter inference)
 - 2. more **stable** model probabilities (model selection)
 - 3. demonstrates excellent empirical performance
 - 4. is easy to use and widely applicable
 - 5. combines the **flexible modeling** features of Bayes with the **distributional robustness** of frequentist methods
- Future work: finite-sample properties of BayesBag and mismatch index, computation, non-independent data (e.g. time series & spatial models)
- **Conclusion:** you should give BayesBag a try!
- On arXiv soon (if you want a heads up: jhuggins@hsph.harvard.edu)

• Let $\Delta_n \triangleq \sum_{i=1}^n \log p(Y_i \mid m_1) - \log p(Y_i \mid m_2)$

• Let
$$\Delta_n \triangleq \sum_{i=1}^n \log p(Y_i \mid m_1) - \log p(Y_i \mid m_2)$$

 δ_i

• Let
$$\Delta_n \triangleq \sum_{i=1}^n \log p(Y_i \mid m_1) - \log p(Y_i \mid m_2)$$

 δ_i

• Then
$$\pi(m_1 | Y) = (1 + \exp\{-\Delta_n\})^{-1}$$



• Let
$$\Delta_n \triangleq \sum_{i=1}^n \log p(Y_i \mid m_1) - \log p(Y_i \mid m_2)$$

 δ_i

• Then
$$\pi(m_1 | Y) = (1 + \exp\{-\Delta_n\})^{-1}$$

• By assumption, $\mathbb{E}[\delta_i] = 0$ but $\sigma^2 = \operatorname{Var}(\delta_i) > 0$



• Let
$$\Delta_n \triangleq \sum_{i=1}^n \log p(Y_i \mid m_1) - \log p(Y_i \mid m_2)$$

 δ_i

• Then
$$\pi(m_1 | Y) = (1 + \exp\{-\Delta_n\})^{-1}$$

- By assumption, $\mathbb{E}[\delta_i] = 0$ but $\sigma^2 = \operatorname{Var}(\delta_i) > 0$
- Hence, Δ_n is a random walk with $\mathbb{E}[\Delta_n^2] = \sigma^2 n$



• Let
$$\Delta_n \triangleq \sum_{i=1}^n \log p(Y_i \mid m_1) - \log p(Y_i \mid m_2)$$

 δ_i

• Then
$$\pi(m_1 | Y) = (1 + \exp\{-\Delta_n\})^{-1}$$

- By assumption, $\mathbb{E}[\delta_i] = 0$ but $\sigma^2 = \operatorname{Var}(\delta_i) > 0$
- Hence, Δ_n is a random walk with $\mathbb{E}[\Delta_n^2] = \sigma^2 n$
- In other words, with very high probability, $|\Delta_n| = \Theta(n^{1/2})$



• Let
$$\Delta_n \triangleq \sum_{i=1}^n \log p(Y_i \mid m_1) - \log p(Y_i \mid m_2)$$

 δ_i

• Then
$$\pi(m_1 | Y) = (1 + \exp\{-\Delta_n\})^{-1}$$

- By assumption, $\mathbb{E}[\delta_i] = 0$ but $\sigma^2 = \operatorname{Var}(\delta_i) > 0$
- Hence, Δ_n is a random walk with $\mathbb{E}[\Delta_n^2] = \sigma^2 n$
- In other words, with very high probability, $|\Delta_n| = \Theta(n^{1/2})$



• Therefore, there is overwhelming evidence of order $n^{1/2}$ for either m_1 or m_2