

Bayesian Inference with Sets of Conjugate Priors

Gero Walter

Department of Statistics
Ludwig-Maximilians-Universität München (LMU)

June 29th, 2012







Introduction

- ▶ Bernoulli observations: 0/1 observations (team wins no/yes)



Introduction

- ▶ Bernoulli observations: 0/1 observations (team wins no/yes)
- ▶ given: a set of observations (team won 12 out of 16 matches)



Introduction

- ▶ Bernoulli observations: 0/1 observations (team wins no/yes)
- ▶ given: a set of observations (team won 12 out of 16 matches)
- ▶ additional to observations, we have strong prior information (we are convinced that $P(\text{win})$ should be around 0.75)



Introduction

- ▶ Bernoulli observations: 0/1 observations (team wins no/yes)
- ▶ given: a set of observations (team won 12 out of 16 matches)
- ▶ additional to observations, we have strong prior information (we are convinced that $P(\text{win})$ should be around 0.75)
- ▶ we are, e.g., interested in (predictive) probability P that team wins in the next match



Introduction

- ▶ Bernoulli observations: 0/1 observations (team wins no/yes)
- ▶ given: a set of observations (team won 12 out of 16 matches)
- ▶ additional to observations, we have strong prior information (we are convinced that $P(\text{win})$ should be around 0.75)
- ▶ we are, e.g., interested in (predictive) probability P that team wins in the next match
- ▶ standard statistical model for this situation:
Beta-Bernoulli/Binomial Model



Beta-Bernoulli/Binomial Model (BBM)

- ▶ Beta prior on $p = P(\text{win})$
- ▶ here in parameterization used, e.g., by Walley (1991):



Beta-Bernoulli/Binomial Model (BBM)

- ▶ Beta prior on $p = P(\text{win})$
- ▶ here in parameterization used, e.g., by Walley (1991):

Data :	s	\sim	$\text{Binom}(p, n)$
conjugate prior:	p	\sim	$\text{Beta}(n^{(0)}, y^{(0)})$
posterior:	$p \mid s$	\sim	$\text{Beta}(n^{(n)}, y^{(n)})$

$$y^{(n)} = \frac{n^{(0)}}{n^{(0)} + n} \cdot y^{(0)} + \frac{n}{n^{(0)} + n} \cdot \frac{s}{n}, \quad n^{(n)} = n^{(0)} + n$$



Beta-Bernoulli/Binomial Model (BBM)

- ▶ Beta prior on $p = P(\text{win})$
- ▶ here in parameterization used, e.g., by Walley (1991):

Data :	s	\sim	$\text{Binom}(p, n)$
conjugate prior:	p	\sim	$\text{Beta}(n^{(0)}, y^{(0)})$
posterior:	$p \mid s$	\sim	$\text{Beta}(n^{(n)}, y^{(n)})$

$$y^{(n)} = \frac{n^{(0)}}{n^{(0)} + n} \cdot y^{(0)} + \frac{n}{n^{(0)} + n} \cdot \frac{s}{n}, \quad n^{(n)} = n^{(0)} + n$$

$$y^{(n)} = E[p \mid s] \quad \text{Var}(p \mid s) = \frac{y^{(n)}(1 - y^{(n)})}{n^{(n)} + 1}$$



Beta-Bernoulli/Binomial Model (BBM)

- ▶ Beta prior on $p = P(\text{win})$
- ▶ here in parameterization used, e.g., by Walley (1991):

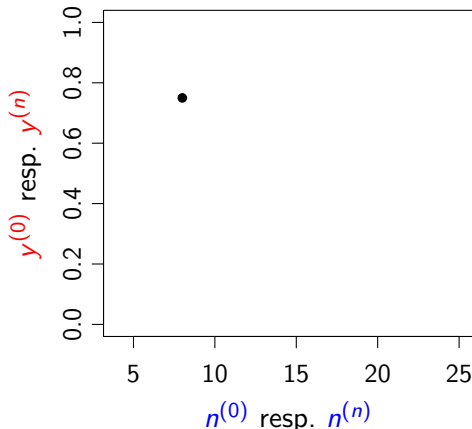
Data :	s	\sim	$\text{Binom}(p, n)$
conjugate prior:	p	\sim	$\text{Beta}(n^{(0)}, y^{(0)})$
posterior:	$p \mid s$	\sim	$\text{Beta}(n^{(n)}, y^{(n)})$

$$y^{(n)} = \frac{n^{(0)}}{n^{(0)} + n} \cdot y^{(0)} + \frac{n}{n^{(0)} + n} \cdot \frac{s}{n}, \quad n^{(n)} = n^{(0)} + n$$

$$y^{(n)} = E[p \mid s] = P \quad \text{Var}(p \mid s) = \frac{y^{(n)}(1 - y^{(n)})}{n^{(n)} + 1}$$



Beta-Bernoulli/Binomial Model (BBM)



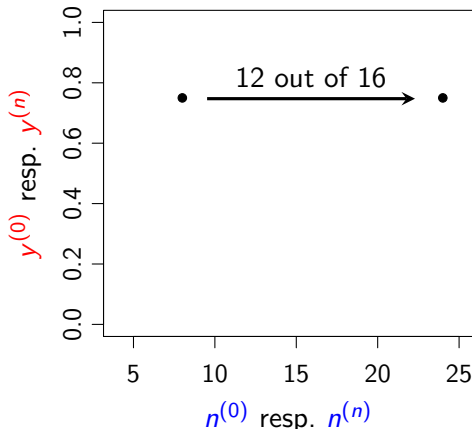
no conflict:

prior $n^{(0)} = 8$, $y^{(0)} = 0.75$

data $s/n = 12/16 = 0.75$



Beta-Bernoulli/Binomial Model (BBM)



no conflict:

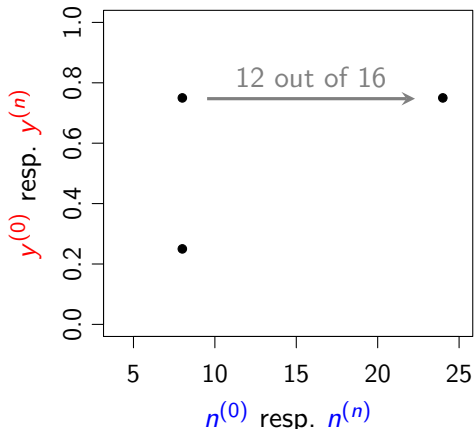
prior $n^{(0)} = 8$, $y^{(0)} = 0.75$

data $s/n = 12/16 = 0.75$

$n^{(n)} = 24$, $y^{(n)} = 0.75$



Beta-Bernoulli/Binomial Model (BBM)



no conflict:

prior $n^{(0)} = 8$, $y^{(0)} = 0.75$

data $s/n = 12/16 = 0.75$



$n^{(n)} = 24$, $y^{(n)} = 0.75$

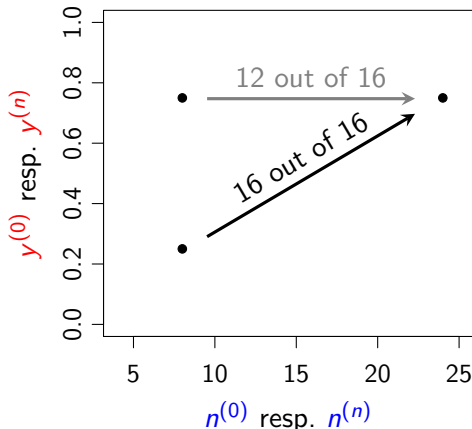
prior-data conflict:

prior $n^{(0)} = 8$, $y^{(0)} = 0.25$

data $s/n = 16/16 = 1$



Beta-Bernoulli/Binomial Model (BBM)



no conflict:

prior $n^{(0)} = 8$, $y^{(0)} = 0.75$
data $s/n = 12/16 = 0.75$



$n^{(n)} = 24$, $y^{(n)} = 0.75$



prior-data conflict:

prior $n^{(0)} = 8$, $y^{(0)} = 0.25$
data $s/n = 16/16 = 1$



Prior-Data Conflict $\hat{=}$ situation in which...

- ▶ ... *informative prior beliefs* and *trusted data* (sampling model correct, no outliers, etc.) are in conflict.
- ▶ "... the prior [places] its mass primarily on distributions in the sampling model for which the observed data is surprising." (Evans & Moshonov, 2006)
- ▶ ... there are not enough data to overrule the prior.

We should notice prior-data conflict in the posterior.



Prior-Data Conflict $\hat{=}$ situation in which...

- ▶ ... *informative prior beliefs* and *trusted data* (sampling model correct, no outliers, etc.) are in conflict.
- ▶ "... the prior [places] its mass primarily on distributions in the sampling model for which the observed data is surprising."
(Evans & Moshonov, 2006)
- ▶ ... there are not enough data to overrule the prior.

We should notice prior-data conflict in the posterior.

$$E[p | s] = y^{(n)} = \frac{n^{(0)}}{n^{(0)} + n} \cdot y^{(0)} + \frac{n}{n^{(0)} + n} \cdot \frac{s}{n}$$

➔ Conflict between prior and data is just averaged out!



Prior-Data Conflict $\hat{=}$ situation in which...

- ▶ ... *informative prior beliefs* and *trusted data* (sampling model correct, no outliers, etc.) are in conflict.
- ▶ "... the prior [places] its mass primarily on distributions in the sampling model for which the observed data is surprising."
(Evans & Moshonov, 2006)
- ▶ ... there are not enough data to overrule the prior.

We should notice prior-data conflict in the posterior.

$$E[p | s] = y^{(n)} = \frac{n^{(0)}}{n^{(0)} + n} \cdot y^{(0)} + \frac{n}{n^{(0)} + n} \cdot \frac{s}{n}$$

➔ Conflict between prior and data is just averaged out!

$$\text{Var}(p | s) = \frac{y^{(n)}(1 - y^{(n)})}{n^{(n)} + 1}, \quad n^{(n)} = n^{(0)} + n$$

➔ does not change systematically with prior-data conflict!



Prior-Data Conflict & Conjugate Priors

Weighted average structure is underneath all common conjugate priors for exponential family sampling distributions!

$X \stackrel{iid}{\sim}$ linear, canonical exponential family, i.e.

$$p(x | \theta) \propto \exp \{ \langle \psi, \tau(x) \rangle - n\mathbf{b}(\psi) \} \quad \left[\psi \text{ transformation of } \theta \right]$$

→ conjugate prior: $p(\psi) \propto \exp \{ n^{(0)} [\langle \psi, \mathbf{y}^{(0)} \rangle - \mathbf{b}(\psi)] \}$

→ (conjugate) posterior: $p(\psi | x) \propto \exp \{ n^{(n)} [\langle \psi, \mathbf{y}^{(n)} \rangle - \mathbf{b}(\psi)] \}$,

$$\text{where } \mathbf{y}^{(n)} = \frac{n^{(0)}}{n^{(0)} + n} \cdot \mathbf{y}^{(0)} + \frac{n}{n^{(0)} + n} \cdot \frac{\tau(x)}{n} \quad \text{and} \quad n^{(n)} = n^{(0)} + n.$$



Why Generalize Bayesian Inference?

Assigning a certain prior distribution on p

↔ Defining a conglomerate of probability statements (on p).

Bayesian theory lacks the ability to specify the degree of uncertainty in these probability statements.



Why Generalize Bayesian Inference?

Assigning a certain prior distribution on p

↔ Defining a conglomerate of probability statements (on p).

Bayesian theory lacks the ability to specify the degree of uncertainty in these probability statements.

Variance or stretch of a distribution for describing uncertainty?



Why Generalize Bayesian Inference?

Assigning a certain prior distribution on p

↔ Defining a conglomerate of probability statements (on p).

Bayesian theory lacks the ability to specify the degree of uncertainty in these probability statements.

Variance or stretch of a distribution for describing uncertainty?

→ Does not work in the case of prior-data conflict:

In conjugate updating, the posterior variance does not depend on the degree of prior-data conflict in most cases.



Why Generalize Bayesian Inference?

Assigning a certain prior distribution on p

↔ Defining a conglomerate of probability statements (on p).

Bayesian theory lacks the ability to specify the degree of uncertainty in these probability statements.

Variance or stretch of a distribution for describing uncertainty?

→ Does not work in the case of prior-data conflict:

In conjugate updating, the posterior variance does not depend on the degree of prior-data conflict in most cases.

→ How to express the precision of a probability statement?



Generalized Bayesian Inference — Basic Idea

Use **set** of priors \rightarrow base inferences on **set** of posteriors
obtained by element-wise updating
 \rightarrow numbers become intervals:

$$\begin{aligned} E[p] &\rightarrow [\underline{E}[p], \bar{E}[p]] \\ P(p \in A) &\rightarrow [\underline{P}(p \in A), \bar{P}(p \in A)] \end{aligned}$$



Generalized Bayesian Inference — Basic Idea

Use **set** of priors \rightarrow base inferences on **set** of posteriors
 obtained by element-wise updating
 \rightarrow numbers become intervals:

$$E[p] \quad \rightarrow \quad [\underline{E}[p], \bar{E}[p]]$$

$$P(p \in A) \quad \rightarrow \quad [\underline{P}(p \in A), \bar{P}(p \in A)]$$

Shorter intervals \leftrightarrow more precise probability statements



Generalized Bayesian Inference — Basic Idea

Use **set** of priors \rightarrow base inferences on **set** of posteriors
 obtained by element-wise updating
 \rightarrow numbers become intervals:

$$E[p] \quad \rightarrow \quad [\underline{E}[p], \bar{E}[p]]$$

$$P(p \in A) \quad \rightarrow \quad [\underline{P}(p \in A), \bar{P}(p \in A)]$$

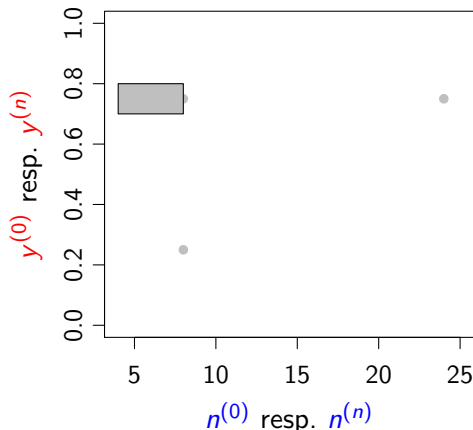
Shorter intervals \leftrightarrow more precise probability statements

\rightarrow differentiate between

- ▶ stochastic uncertainty (“risk”) vs.
- ▶ non-stochastic uncertainty (“ambiguity”)



pdc-Imprecise BBM (pdc-IBBM): Walley 1991, Ch.5.4.3

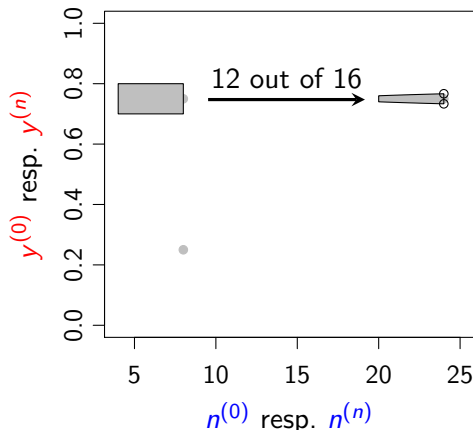


no conflict:

prior $n^{(0)} \in [4, 8]$, $y^{(0)} \in [0.7, 0.8]$
data $s/n = 12/16 = 0.75$



pdc-Imprecise BBM (pdc-IBBM): Walley 1991, Ch.5.4.3



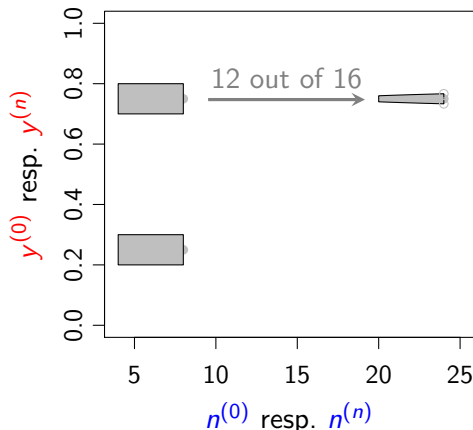
no conflict:

prior $n^{(0)} \in [4, 8]$, $y^{(0)} \in [0.7, 0.8]$
data $s/n = 12/16 = 0.75$

▼
"spotlight" shape



pdc-Imprecise BBM (pdc-IBBM): Walley 1991, Ch.5.4.3



no conflict:

prior $n^{(0)} \in [4, 8]$, $y^{(0)} \in [0.7, 0.8]$
data $s/n = 12/16 = 0.75$

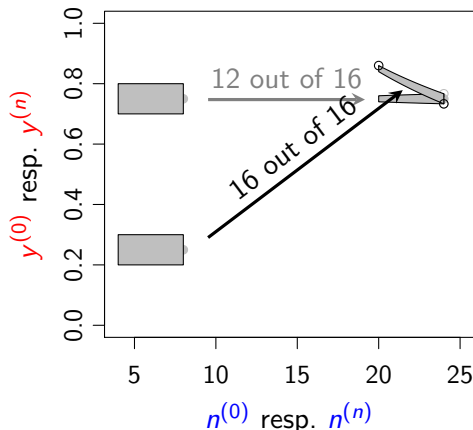
▼
“spotlight” shape

prior-data conflict:

prior $n^{(0)} \in [4, 8]$, $y^{(0)} \in [0.2, 0.3]$
data $s/n = 16/16 = 1$



pdc-Imprecise BBM (pdc-IBBM): Walley 1991, Ch.5.4.3



no conflict:

prior $n^{(0)} \in [4, 8]$, $y^{(0)} \in [0.7, 0.8]$
data $s/n = 12/16 = 0.75$

“spotlight” shape

prior-data conflict:

prior $n^{(0)} \in [4, 8]$, $y^{(0)} \in [0.2, 0.3]$
data $s/n = 16/16 = 1$

“banana” shape



Inner Workings

- ▶ *convex* sets of distributions (“credal sets”)



Inner Workings

- ▶ *convex* sets of distributions (“credal sets”)
- ▶ convexity needed for consistency properties (“coherence”)



Inner Workings

- ▶ *convex* sets of distributions (“credal sets”)
- ▶ convexity needed for consistency properties (“coherence”)
- ▶ sets of distributions induced by sets of parameters:
not necessarily convex



Inner Workings

- ▶ *convex* sets of distributions (“credal sets”)
- ▶ convexity needed for consistency properties (“coherence”)
- ▶ sets of distributions induced by sets of parameters:
not necessarily convex
- ▶ take convex hull of these parametric distributions:
credal set = finite convex mixtures of parametric distributions



Inner Workings

- ▶ *convex* sets of distributions (“credal sets”)
- ▶ convexity needed for consistency properties (“coherence”)
- ▶ sets of distributions induced by sets of parameters:
not necessarily convex
- ▶ take convex hull of these parametric distributions:
credal set = finite convex mixtures of parametric distributions
- ▶ prior/posterior credal set: convex hull of distributions
induced by prior/posterior parameter set



Inner Workings

- ▶ *convex* sets of distributions (“credal sets”)
- ▶ convexity needed for consistency properties (“coherence”)
- ▶ sets of distributions induced by sets of parameters:
not necessarily convex
- ▶ take convex hull of these parametric distributions:
credal set = finite convex mixtures of parametric distributions
- ▶ prior/posterior credal set: convex hull of distributions
induced by prior/posterior parameter set
- ▶ pictures show parameter sets (that need not be convex)



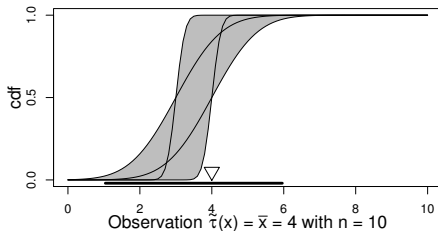
Properties

Works for any canonical exponential family sampling distribution!
→ *generalized iLUCK models*, Walter & Augustin (2009)

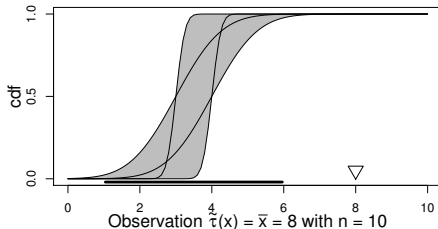


$$X \sim N(\mu, 1) \rightarrow \mu \sim N(y^{(0)}, \frac{1}{n^{(0)}})$$

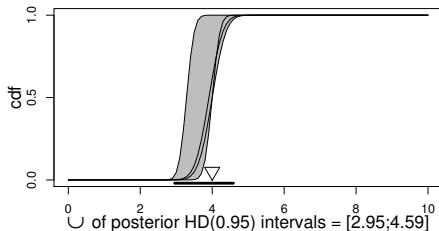
Set of priors: $y^{(0)} \in [3;4]$ and $n^{(0)} \in [1;25]$



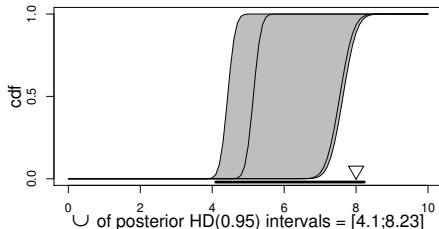
Set of priors: $y^{(0)} \in [3;4]$ and $n^{(0)} \in [1;25]$



Set of posteriors: $y^{(1)} \in [3.29;4]$ and $n^{(1)} \in [11;35]$



Set of posteriors: $y^{(1)} \in [4.43;7.64]$ and $n^{(1)} \in [11;35]$





Properties

Works for any canonical exponential family sampling distribution!

→ *generalized iLUCK-models*, Walter & Augustin (2009)

- ▶ $n^{(0)}$ governs precision of posterior:
 $n^{(0)} \uparrow \leftrightarrow \text{precision} \downarrow$



Properties

Works for any canonical exponential family sampling distribution!

→ *generalized iLUCK-models*, Walter & Augustin (2009)

▶ $n^{(0)}$ governs precision of posterior:

$n^{(0)} \uparrow \longleftrightarrow \text{precision} \downarrow$

▶ $n \rightarrow \infty$: consistency ($y^{(n)}$ set reduces to a point at $\tau(x)/n$)



Properties

Works for any canonical exponential family sampling distribution!

→ *generalized iLUCK-models*, Walter & Augustin (2009)

- ▶ $n^{(0)}$ governs precision of posterior:
 $n^{(0)} \uparrow \longleftrightarrow \text{precision} \downarrow$
- ▶ $n \rightarrow \infty$: consistency ($y^{(n)}$ set reduces to a point at $\tau(x)/n$)
- ▶ $y^{(0)}$ stretch $\uparrow \longleftrightarrow y^{(n)}$ stretch \uparrow



Properties

Works for any canonical exponential family sampling distribution!

→ *generalized iLUCK-models*, Walter & Augustin (2009)

- ▶ $n^{(0)}$ governs precision of posterior:
 $n^{(0)} \uparrow \longleftrightarrow \text{precision} \downarrow$
- ▶ $n \rightarrow \infty$: consistency ($y^{(n)}$ set reduces to a point at $\tau(x)/n$)
- ▶ $y^{(0)}$ stretch $\uparrow \longleftrightarrow y^{(n)}$ stretch \uparrow
- ▶ inferences should be linear in posterior distributions:
 then min/max are attained at the parametric distributions
 (these are the extreme points of the credal set);
 E, Var are linear in the parametric distributions.



Properties

Works for any canonical exponential family sampling distribution!

→ *generalized iLUCK-models*, Walter & Augustin (2009)

- ▶ $n^{(0)}$ governs precision of posterior:
 $n^{(0)} \uparrow \longleftrightarrow \text{precision} \downarrow$
- ▶ $n \rightarrow \infty$: consistency ($y^{(n)}$ set reduces to a point at $\tau(x)/n$)
- ▶ $y^{(0)}$ stretch $\uparrow \longleftrightarrow y^{(n)}$ stretch \uparrow
- ▶ inferences should be linear in posterior distributions:
 then min/max are attained at the parametric distributions
 (these are the extreme points of the credal set);
 E, Var are linear in the parametric distributions.
- ▶ reaction to prior-data conflict due to different 'updating speeds' depending on $n^{(0)}$: $y^{(n)}$ moves "faster" for low $n^{(0)}$



Open Ends/Challenges

- ▶ rectangular prior set (two-dimensional interval) seems natural, but generally any shape possible



Open Ends/Challenges

- ▶ rectangular prior set (two-dimensional interval) seems natural, but generally any shape possible
- ▶ posterior parameter sets are not rectangular anyway



Open Ends/Challenges

- ▶ rectangular prior set (two-dimensional interval) seems natural, but generally any shape possible
- ▶ posterior parameter sets are not rectangular anyway
- ▶ prior shape influences the posterior inferences



Open Ends/Challenges

- ▶ rectangular prior set (two-dimensional interval) seems natural, but generally any shape possible
- ▶ posterior parameter sets are not rectangular anyway
- ▶ prior shape influences the posterior inferences
- ▶ shape can be tailored to enable desired inference properties (e.g. bonus precision if prior and data agree especially well)



Open Ends/Challenges

- ▶ rectangular prior set (two-dimensional interval) seems natural, but generally any shape possible
- ▶ posterior parameter sets are not rectangular anyway
- ▶ prior shape influences the posterior inferences
- ▶ shape can be tailored to enable desired inference properties (e.g. bonus precision if prior and data agree especially well)
- ▶ for more complex shapes, elicitation becomes more difficult



Open Ends/Challenges

- ▶ rectangular prior set (two-dimensional interval) seems natural, but generally any shape possible
- ▶ posterior parameter sets are not rectangular anyway
- ▶ prior shape influences the posterior inferences
- ▶ shape can be tailored to enable desired inference properties (e.g. bonus precision if prior and data agree especially well)
- ▶ for more complex shapes, elicitation becomes more difficult
- ▶ take two $y^{(0)}$ intervals at two different $n^{(0)}$ values?