# Welcome to the post-Bayesian seminar!

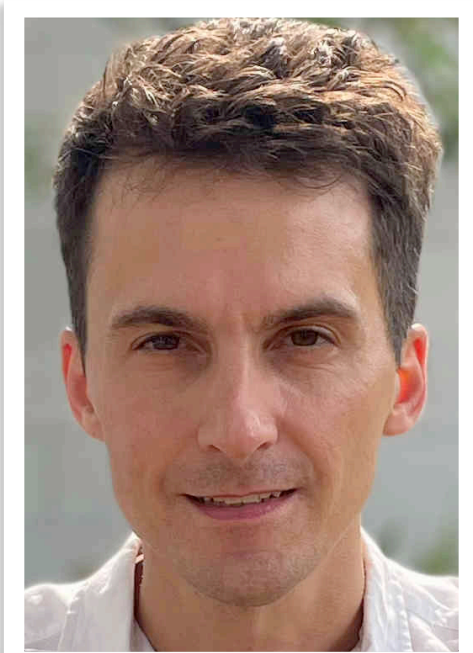**When:** every 2 weeks @ Tuesdays either 9 AM (9:00) or 2 PM (14:00) GMT

**Structure:**
Chapter 1: Generalised Bayes (11/02—22/04)
Chapter 2: Resampling & Martingale Posteriors (06/05—15/07)
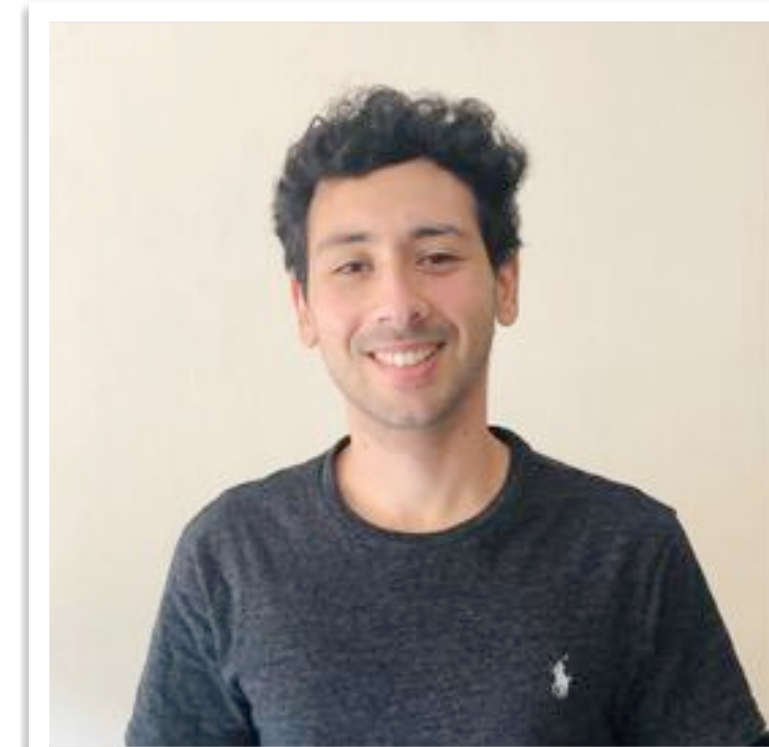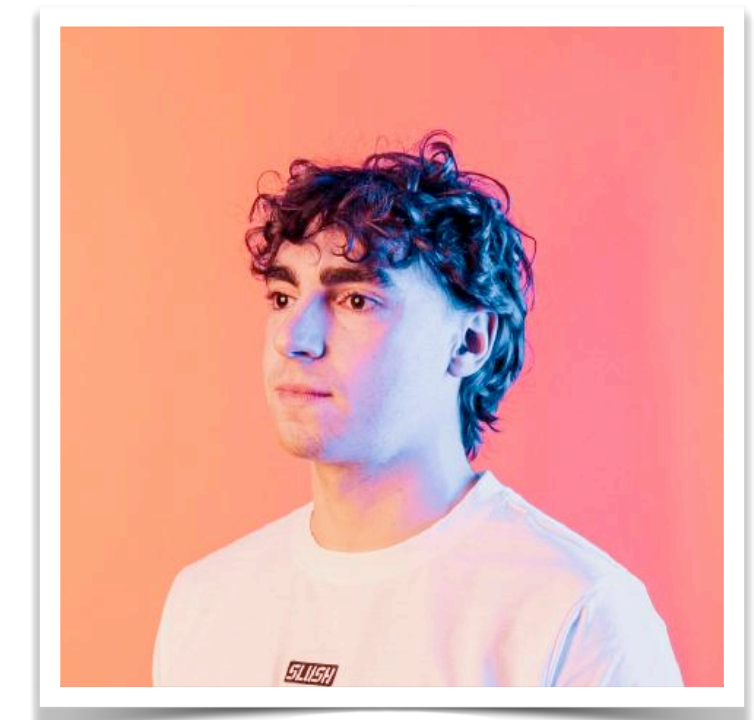Chapter 3: PAC-Bayes (after the summer break)

**Organisers:**



**Prof. Pierre Alquier (ESSEC Singapore)**

**Dr. Edwin Fong (University of Hong Kong)**
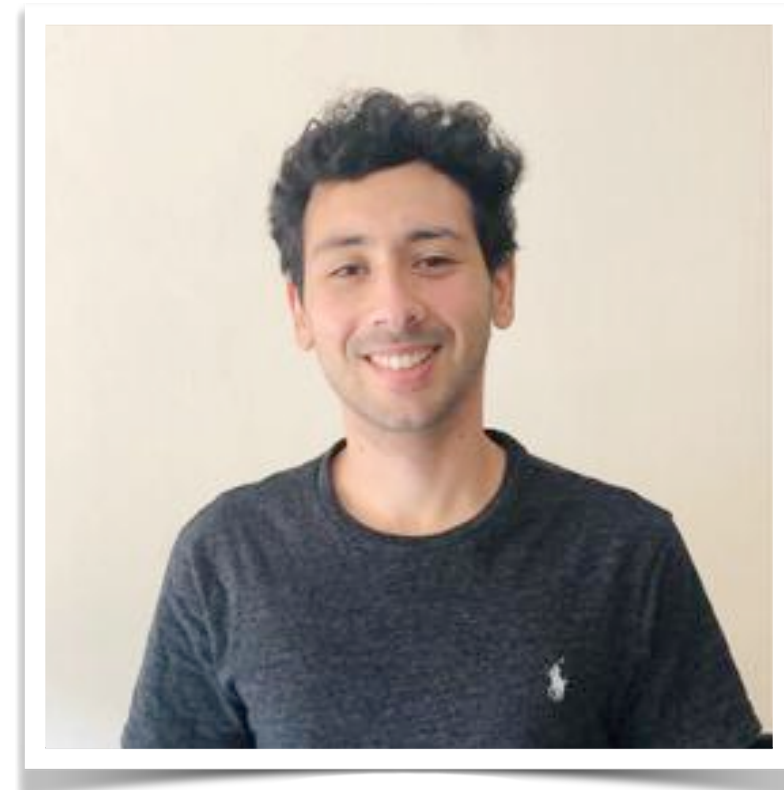
**Matias Altamirano (UCL)**
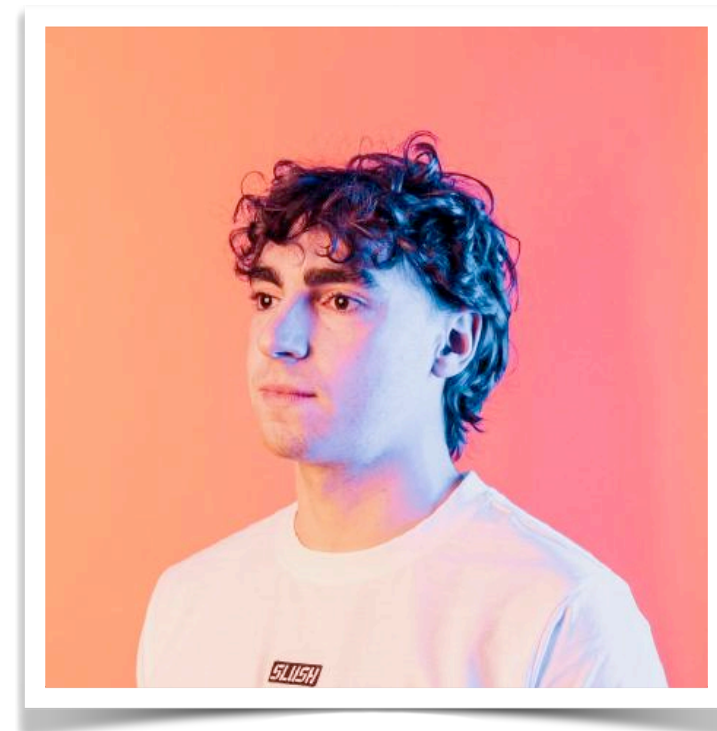
**Yann McLatchie (UCL)**

# Welcome to the post-Bayesian seminar!

**Shameless Plug:**

**Workshop @ UCL on post-Bayesian methods
15. /16. May 2025!!!**



**Matias Altamirano
(UCL)**

**Yann McLatchie
(UCL)**

# Welcome to the post-Bayesian seminar!

## <u>Important Links</u>

At a glance/website:                           <u>https://tinyurl.com/postBayesWebsite</u>
Where to subscribe to mailing list:      <u>https://tinyurl.com/postBayesSubscribe</u>
Where to subscribe to calendar:          <u>https://tinyurl.com/postBayesCalendar</u>
Where to attend the seminars:            <u>https://tinyurl.com/postBayesZoom</u>
Where recorded seminars are stored:    <u>https://tinyurl.com/postBayesYT</u>

## Please share widely! :)

# Welcome to the post-Bayesian seminar!

## Questions / Comments during talks

**During talk:**
- use Q/A function in zoom
- Other questions can be upvoted
- We will try to monitor questions and ask relevant ones in natural breaks

**After talk:**
- Raise your hand in zoom
- We will do our best to decide who gets to ask a question fairly
- We will do our best to resolve remaining questions in Q / A function

# The Bayesian hangover:
## updating beliefs about updating beliefs

Jeremias Knoblauch
Department of Statistical Science
University College London

11/02/25

# Key Questions addressed in talk

**Part I: What's the Bayesian hangover?
And why do we need this seminar?**

**Part II: What is the (post-Bayesian) aspirin?**





## Part III: What will Chapter 1 cover?

Power/Fractional/Cold Posterior

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

Gibbs/Quasi/Pseudo Posterior

$$\pi_n^{L}(\theta \mid x_{1:n}) = \frac{\exp\{-L(x_{1:n}, p_{\theta})\} \cdot \pi(\theta)}{\int \exp\{-L(x_{1:n}, p_{\theta})\} \cdot \pi(\theta) d\theta}$$

Optimisation-centric Posterior

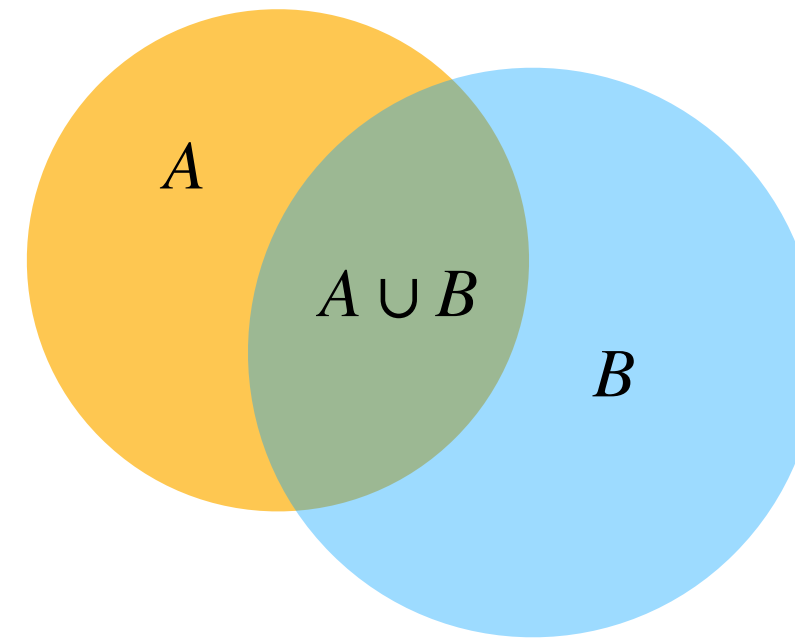$$q_n^*(\theta) = \arg\min_{q \in \mathcal{Q}} \left\{ \mathcal{L}(q, x_{1:n}) + D(q, \pi) \right\}$$

# Part I: The Hangover

# Preamble: Bayesian Data Analysis

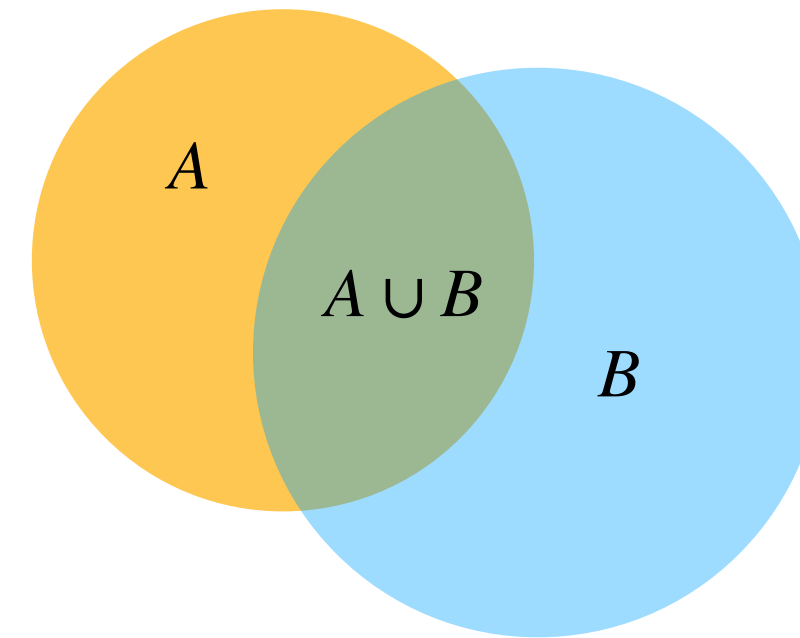**Bayes' Theorem: Inversion of conditionals**

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(A)}{P(B)}$$

# Preamble: Bayesian Data Analysis

**Bayes' Theorem: Inversion of conditionals**

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(A)}{P(B)}$$



Data model: $p(x_{1:n} \mid \theta)$
$x_{1:n} \in \mathcal{X}^n$

Prior probability: $\pi(\theta)$
$\theta \in \Theta$

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

**(Bayes) Posterior**

# Preamble: Bayesian Data Analysis

**Bayes' Theorem: Inversion of conditionals**

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(A)}{P(B)}$$

Data model:  $p(x_{1:n} \mid \theta)$
$x_{1:n} \in \mathcal{X}^n$

Prior probability:  $\pi(\theta)$
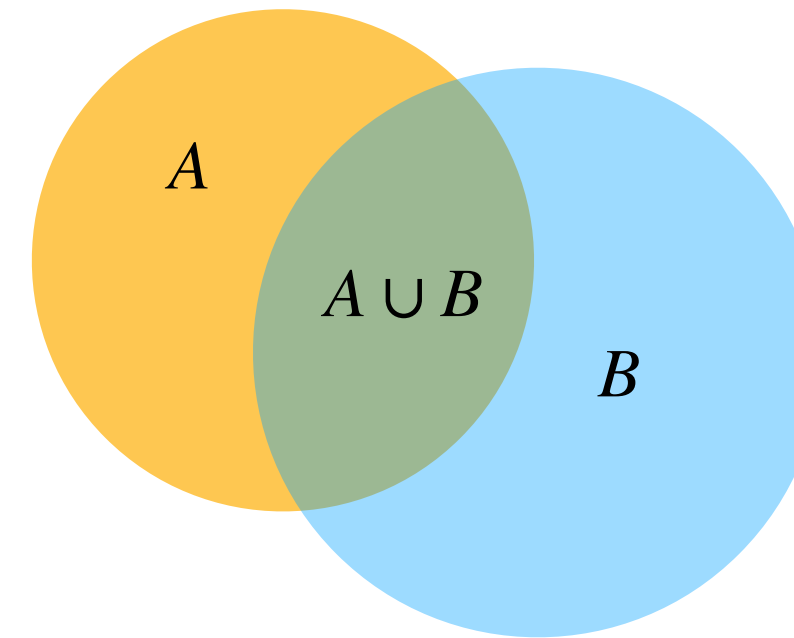$\theta \in \Theta$

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

**(Bayes) Posterior**

$A$

$A \cup B$

$B$

$\oplus$  Averages models (instead of picking only one)

$\oplus$  Quantifies uncertainty about $\theta$ via $\pi_n(\theta \mid x_{1:n})$

$\oplus$  Inclusion of domain expertise via prior $\pi$

# Preamble: Bayesian Data Analysis

**Bayes' Theorem: Inversion of conditionals**

$$P(\textcolor{orange}{A} \mid \textcolor{blue}{B}) = \frac{P(\textcolor{blue}{B} \mid \textcolor{orange}{A}) \cdot P(\textcolor{orange}{A})}{P(\textcolor{blue}{B})}$$

Data model: $p(x_{1:n} \mid \theta)$

$x_{1:n} \in \mathcal{X}^n$

Prior probability: $\pi(\theta)$

$\theta \in \Theta$

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$
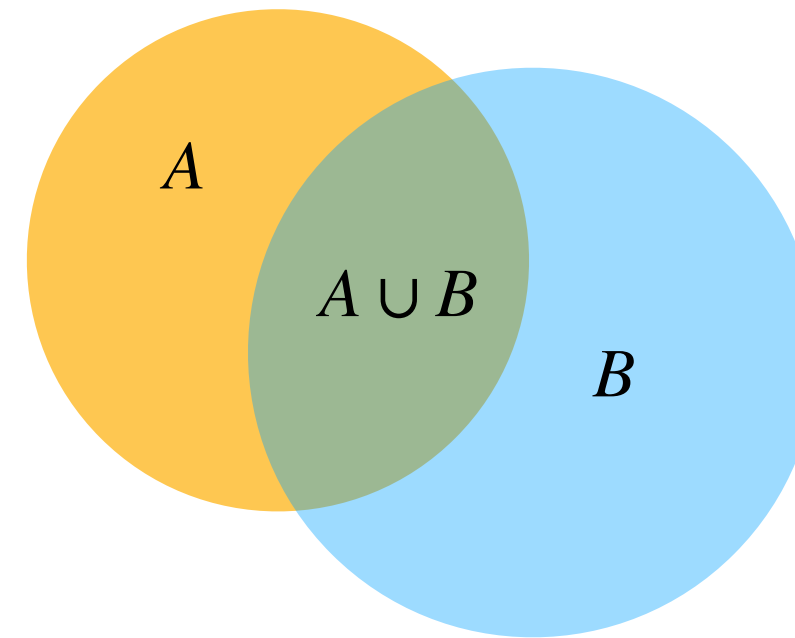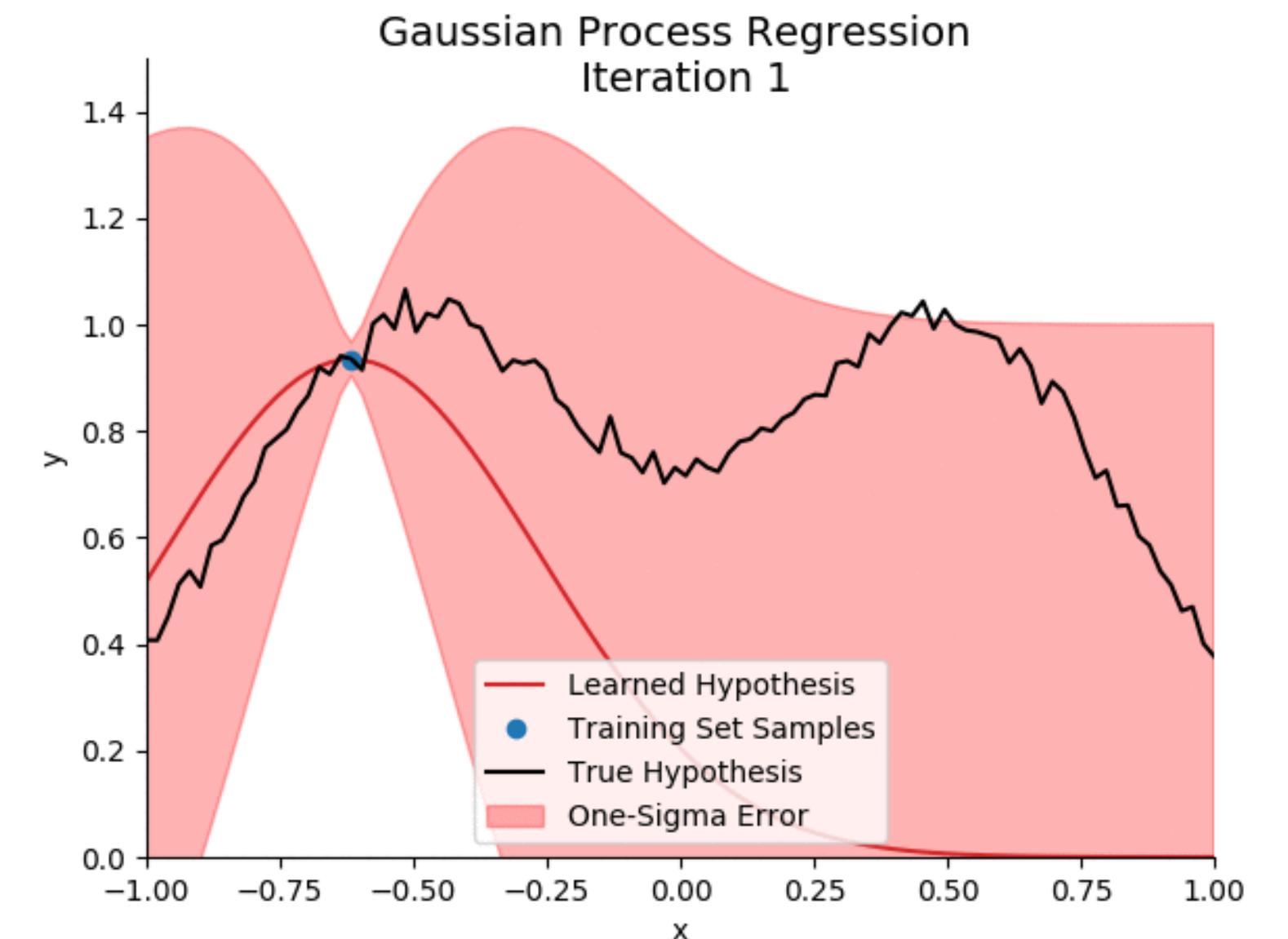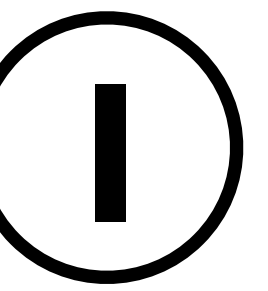
**(Bayes) Posterior**



Gaussian Process Regression
Iteration 1

- Learned Hypothesis
- Training Set Samples
- True Hypothesis
- One-Sigma Error

(+) Averages models (instead of picking only one)

(+) Quantifies uncertainty about $\theta$ via $\pi_n(\theta \mid x_{1:n})$

(+) Inclusion of domain expertise via prior $\pi$

# Problematic Assumptions for Bayesian Analysis  ⓘ
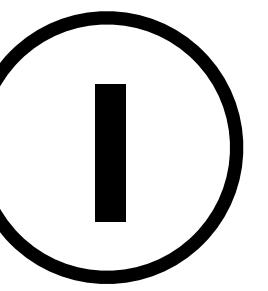
(A1)  $\boxed{x_{1:n} \sim p(x_{1:n} \mid \theta*) \text{ for some } \theta* \in \Theta}$

$\Theta = $ Only relevant State of the world

# Problematic Assumptions for Bayesian Analysis ①

(A1) $\boxed{x_{1:n} \sim p(x_{1:n} \mid \theta^*) \text{ for some } \theta^* \in \Theta}$

$\Theta$ = Only relevant State of the world

(A2) $\boxed{\pi(\theta) = \text{uncertainty about the true State of the world}}$

How rational decision-makers choose the prior

# Problematic Assumptions for Bayesian Analysis ①

| | |
|---|---|
| (A1) | model well-specified |
| (A2) | prior well-specified |
| (A3) | computationally feasible |

**(A1)** $\boxed{x_{1:n} \sim p(x_{1:n} \mid \theta^*) \text{ for some } \theta^* \in \Theta}$

$\Theta$ = Only relevant State of the world

**(A2)** $\boxed{\pi(\theta) = \text{ uncertainty about the true State of the world}}$

How rational decision-makers choose the prior

**(A3)** $\boxed{\pi_n(\theta \mid x_{1:n}) \text{ computable in practice}}$

Guarantees real-world relevance

# Problematic Assumptions for Bayesian Analysis

(I)

(A1)    model well-specified
(A2)    prior well-specified
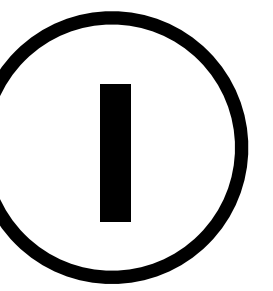(A3)    computationally feasible

**(A1)** $x_{1:n} \sim p(x_{1:n} \mid \theta*)$ for some $\theta* \in \Theta$

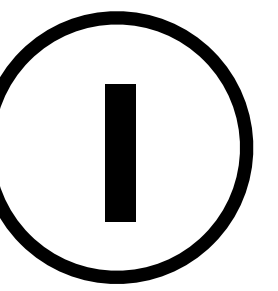$\Theta$  =  Only relevant State of the world

**(A2)** $\pi(\theta) =$ uncertainty about the true State of the world

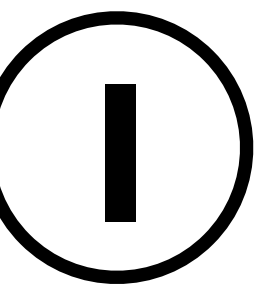How rational decision-makers choose the prior

**(A3)** $\pi_n(\theta \mid x_{1:n})$ computable in practice

Guarantees real-world relevance

FRAGILE

# Case Study: Bayesian ML & Boston Housing Data

Traditional Bayesian analysis in science

Expert with
research question

$x_{1:n}$

| | |
|---|---|
| (A1) | model well-specified |
| (A2) | prior well-specified |
| (A3) | computationally feasible |

# Case Study: Bayesian ML & Boston Housing Data  Ⓘ

Traditional Bayesian analysis in science

Expert with
research question

Statistical modelling &
expert knowledge

$x_{1:n}$

$p(x_{1:n} \mid \theta), \pi(\theta)$

(A1)  model well-specified

(A2)  prior well-specified

(A3)  computationally feasible

# Case Study: Bayesian ML & Boston Housing Data  (I)

Traditional Bayesian analysis in science

Expert with research question $\longrightarrow$ Statistical modelling & expert knowledge $\longrightarrow$ Inference
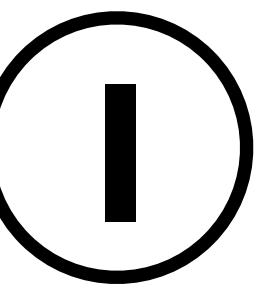
$x_{1:n}$ 　　　　$p(x_{1:n} \mid \theta), \pi(\theta)$ 　　　 $\pi_n(\theta \mid x_{1:n})$
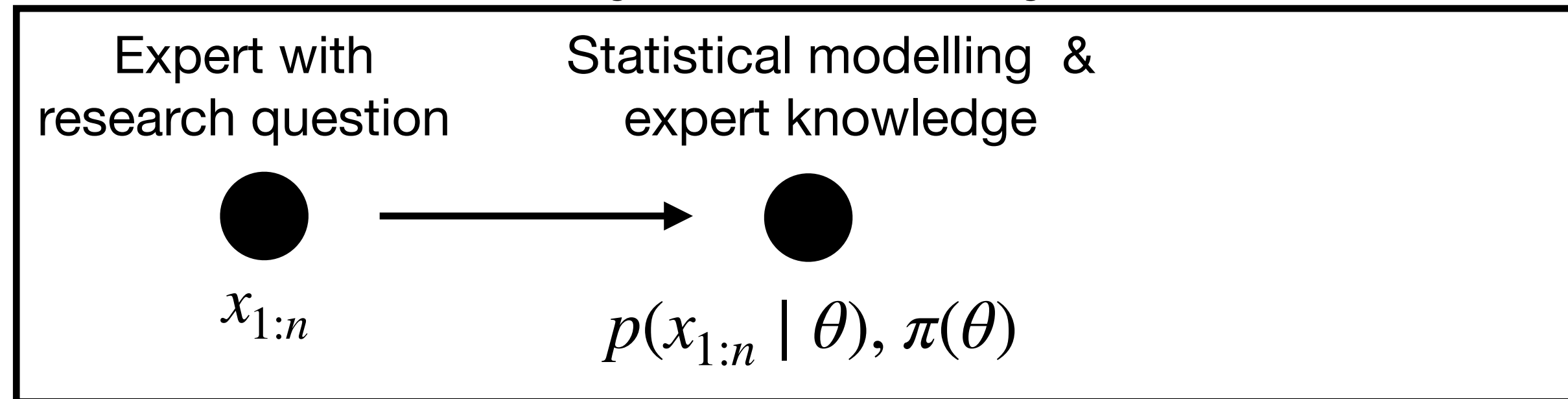
(A1)　model well-specified

(A2)　prior well-specified

(A3)　computationally feasible

# Case Study: Bayesian ML & Boston Housing Data  (I)

## Traditional Bayesian analysis in science

Expert with research question

Statistical modelling & expert knowledge

Inference

$x_{1:n}$

$p(x_{1:n} \mid \theta), \pi(\theta)$

$\pi_n(\theta \mid x_{1:n})$

**Harrison & Rubinfeld (1978)**
**Research Question:** influence of air pollution on house prices?

(A1) model well-specified

(A2) prior well-specified

(A3) computationally feasible
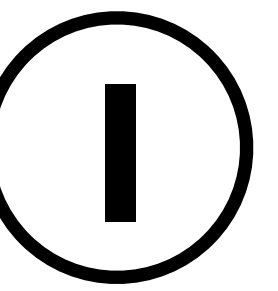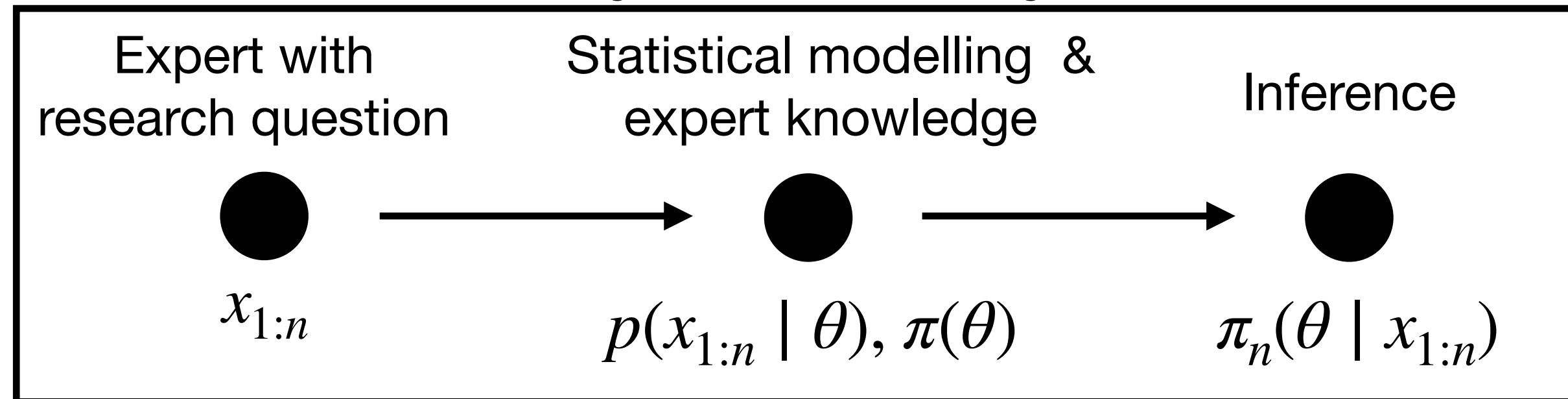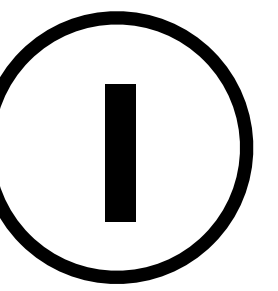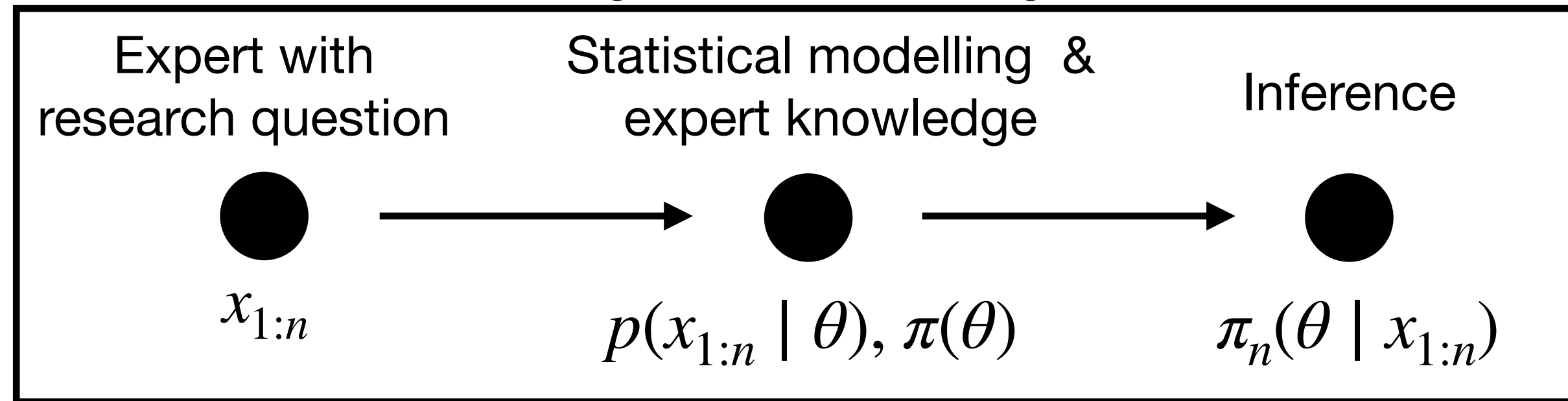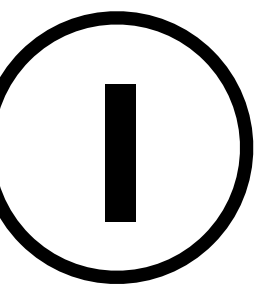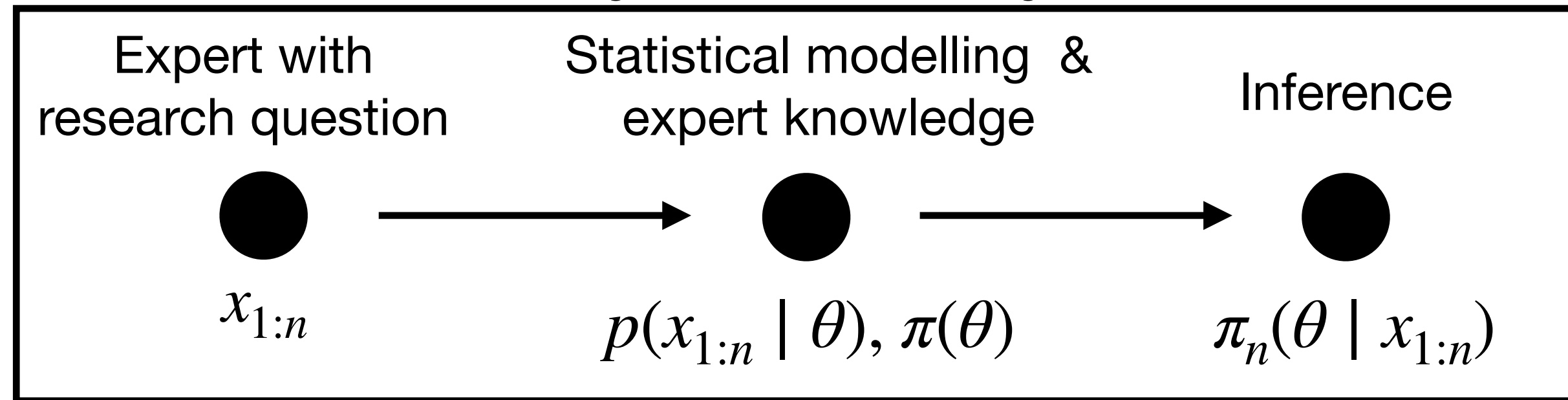
# Case Study: Bayesian ML & Boston Housing Data  (I)

## Traditional Bayesian analysis in science

| Expert with research question | Statistical modelling & expert knowledge | Inference |
|---|---|---|
| ● | ● | ● |
| $x_{1:n}$ | $p(x_{1:n} \mid \theta), \pi(\theta)$ | $\pi_n(\theta \mid x_{1:n})$ |

**Harrison & Rubinfeld (1978)**

**Research Question:** influence of air pollution on house prices?

(A1) ✔

parameters of interest

incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay ⟶

pollutants ⟶

rooms, sqm, ... ⟶

measurement error ⟶

$$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$$

| (A1) | model well-specified |
|---|---|
| (A2) | prior well-specified |
| (A3) | computationally feasible |

# Case Study: Bayesian ML & Boston Housing Data ⓘ

Traditional Bayesian analysis in science

Expert with research question → Statistical modelling & expert knowledge → Inference

$x_{1:n}$  $\qquad$  $p(x_{1:n} \mid \theta), \pi(\theta)$  $\qquad$  $\pi_n(\theta \mid x_{1:n})$

**Harrison & Rubinfeld (1978)**

**Research Question:** influence of air pollution on house prices?

(A1) ✓  parameters of interest  incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

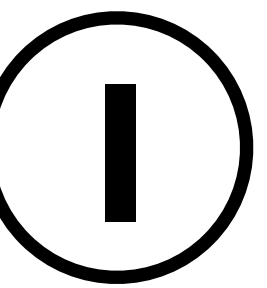willingness to pay  $\qquad$  pollutants  $\qquad$  rooms, sqm, ...

measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$

$\pi(\theta) \sim$ hand-crafted by experts   (A2) ✓
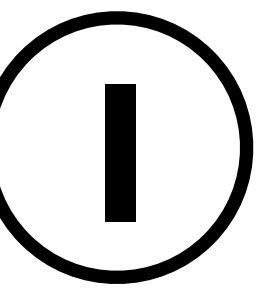
(A1)  model well-specified
(A2)  prior well-specified
(A3)  computationally feasible

Traditional Bayesian analysis in science

| Expert with research question | Statistical modelling & expert knowledge | Inference |
|---|---|---|
| ● | ● | ● |
| $x_{1:n}$ | $p(x_{1:n} \mid \theta),\ \pi(\theta)$ | $\pi_n(\theta \mid x_{1:n})$ |

**Harrison & Rubinfeld (1978)**

**Research Question:** influence of air pollution on house prices?

(A1) ✔

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

parameters of interest — $p_j$

incidental parameters — $c_0$, $c_j$

willingness to pay ↑

pollutants ↑

rooms, sqm, ... ↑

measurement error ↑

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$
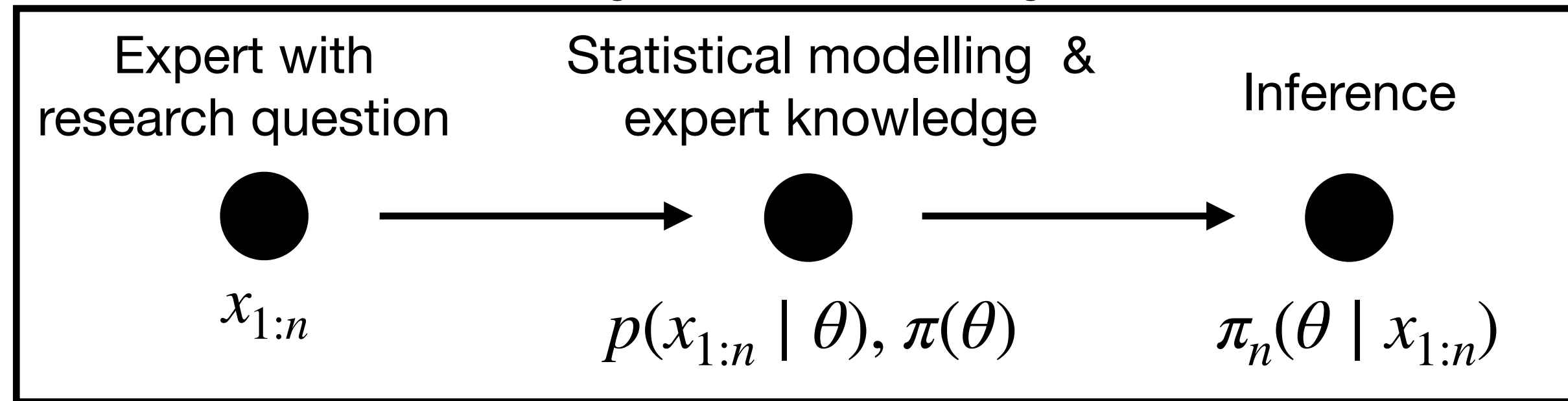
$\pi(\theta) \sim$ hand-crafted by experts   (A2) ✔

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly   (A3) ✔

| | |
|---|---|
| (A1) | model well-specified |
| (A2) | prior well-specified |
| (A3) | computationally feasible |

# Case Study: Bayesian ML & Boston Housing Data  (I)

## Traditional Bayesian analysis in science

Expert with research question
$\quad\quad$ Statistical modelling & expert knowledge
$\quad\quad$ Inference

$x_{1:n}$
$\quad\quad\quad$ $p(x_{1:n} \mid \theta), \pi(\theta)$
$\quad\quad\quad$ $\pi_n(\theta \mid x_{1:n})$

## Modern Bayesian ML

Flexible model

$p(x_{1:n} \mid \theta), \pi(\theta)$

**Harrison & Rubinfeld (1978)**
**Research Question:** influence of air pollution on house prices?

(A1) ✔

parameters of interest $\quad$ incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay $\quad$ pollutants $\quad$ rooms, sqm, ...

measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$
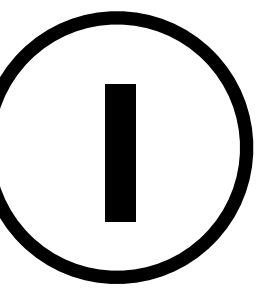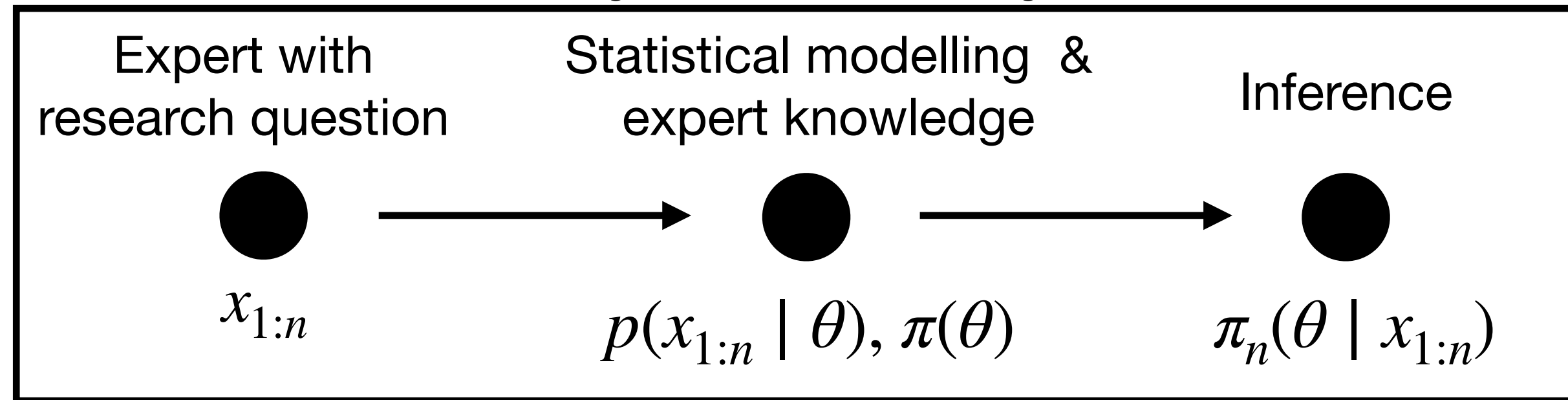
$\pi(\theta) \sim$ hand-crafted by experts $\quad$ (A2) ✔

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly $\quad$ (A3) ✔

(A1) $\quad$ model well-specified
(A2) $\quad$ prior well-specified
(A3) $\quad$ computationally feasible

## Traditional Bayesian analysis in science

Expert with research question

Statistical modelling & expert knowledge

Inference

$x_{1:n}$

$p(x_{1:n} \mid \theta), \pi(\theta)$

$\pi_n(\theta \mid x_{1:n})$

## Modern Bayesian ML

Flexible model

Different data/problems

$p(x_{1:n} \mid \theta), \pi(\theta)$

$x_{1:n}$

**Harrison & Rubinfeld (1978)**

**Research Question:** influence of air pollution on house prices?

(A1) ✔

parameters of interest

incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay

pollutants

rooms, sqm, ...

measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$
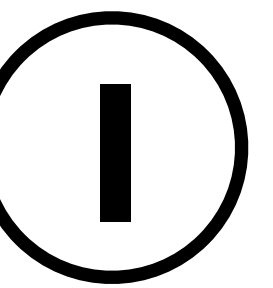
$\pi(\theta) \sim$ hand-crafted by experts    (A2) ✔

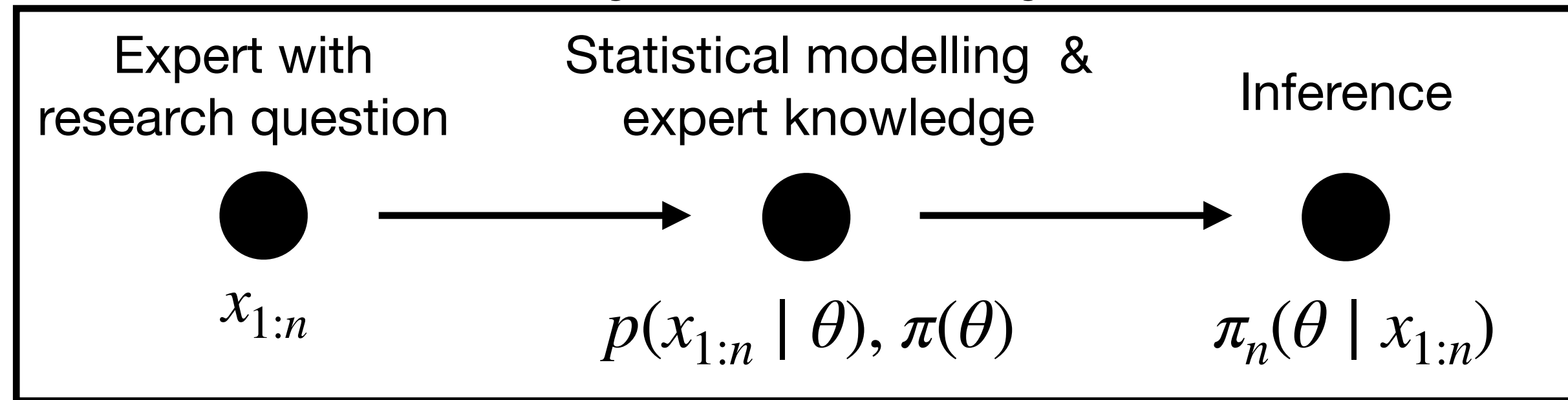$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly    (A3) ✔

(A1)  model well-specified
(A2)  prior well-specified
(A3)  computationally feasible

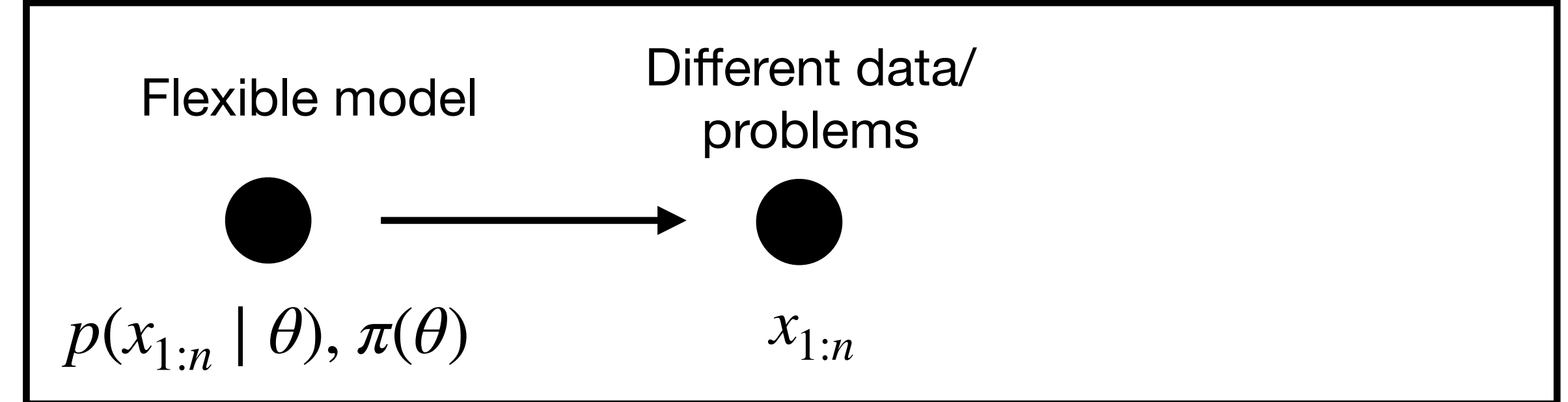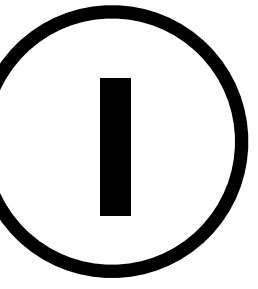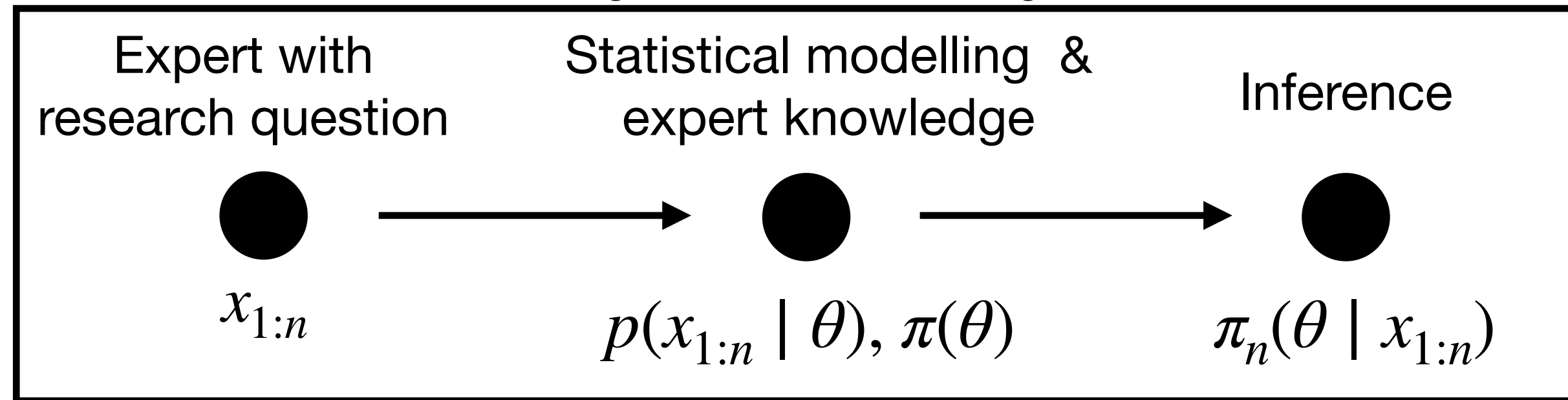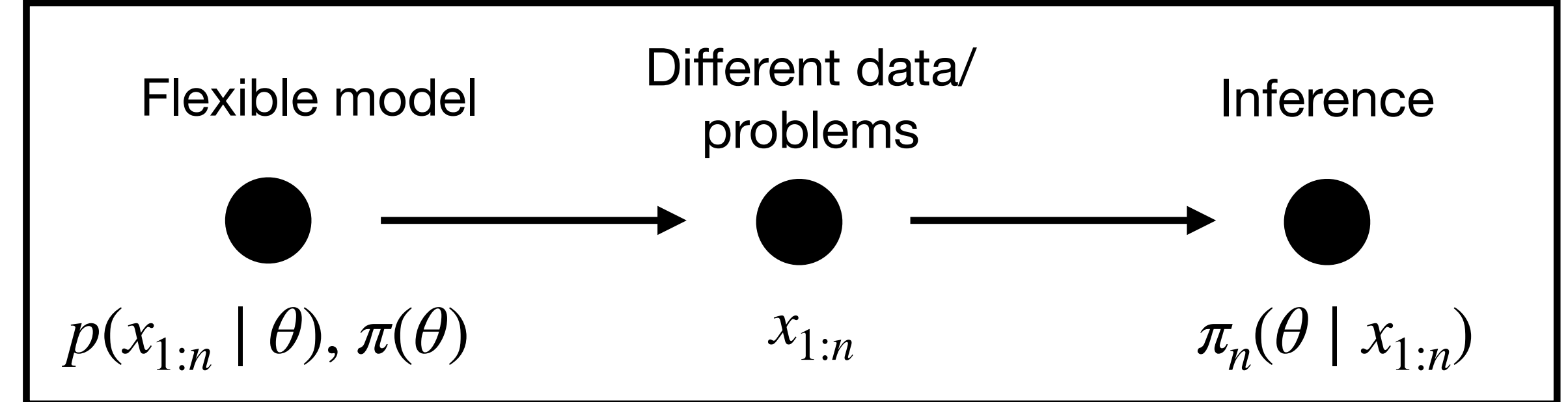# Case Study: Bayesian ML & Boston Housing Data Ⓘ

## Traditional Bayesian analysis in science

Expert with research question → Statistical modelling & expert knowledge → Inference

$x_{1:n}$     $p(x_{1:n} \mid \theta), \pi(\theta)$     $\pi_n(\theta \mid x_{1:n})$

## Modern Bayesian ML

Flexible model → Different data/ problems → Inference

$p(x_{1:n} \mid \theta), \pi(\theta)$     $x_{1:n}$     $\pi_n(\theta \mid x_{1:n})$

**Harrison & Rubinfeld (1978)**

**Research Question:** influence of air pollution on house prices?

(A1) ✓

parameters of interest        incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay        pollutants        rooms, sqm, ...

measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^{\top}$

$\pi(\theta) \sim$ hand-crafted by experts     (A2) ✓

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly     (A3) ✓

| (A1) | model well-specified |
| --- | --- |
| (A2) | prior well-specified |
| (A3) | computationally feasible |

# Case Study: Bayesian ML & Boston Housing Data ⓘ

## Traditional Bayesian analysis in science



Expert with research question → Statistical modelling & expert knowledge → Inference

$x_{1:n}$ → $p(x_{1:n} \mid \theta), \pi(\theta)$ → $\pi_n(\theta \mid x_{1:n})$

## Modern Bayesian ML



Flexible model → Different data/problems → Inference

$p(x_{1:n} \mid \theta), \pi(\theta)$ → $x_{1:n}$ → $\pi_n(\theta \mid x_{1:n})$

**Harrison & Rubinfeld (1978)**
**Research Question:** influence of air pollution on house prices?

**(A1)** ✔

parameters of interest     incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay     pollutants     rooms, sqm, ...

measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$
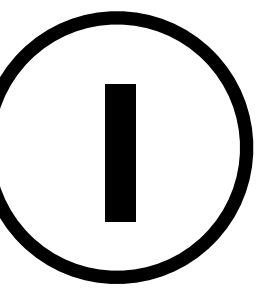
$\pi(\theta) \sim$ hand-crafted by experts     **(A2)** ✔

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly     **(A3)** ✔

**Pearce et al. (2020) [AISTATS]**
**Research Question:** Does my algorithm improve prediction on regression tasks like Boston UCI data?
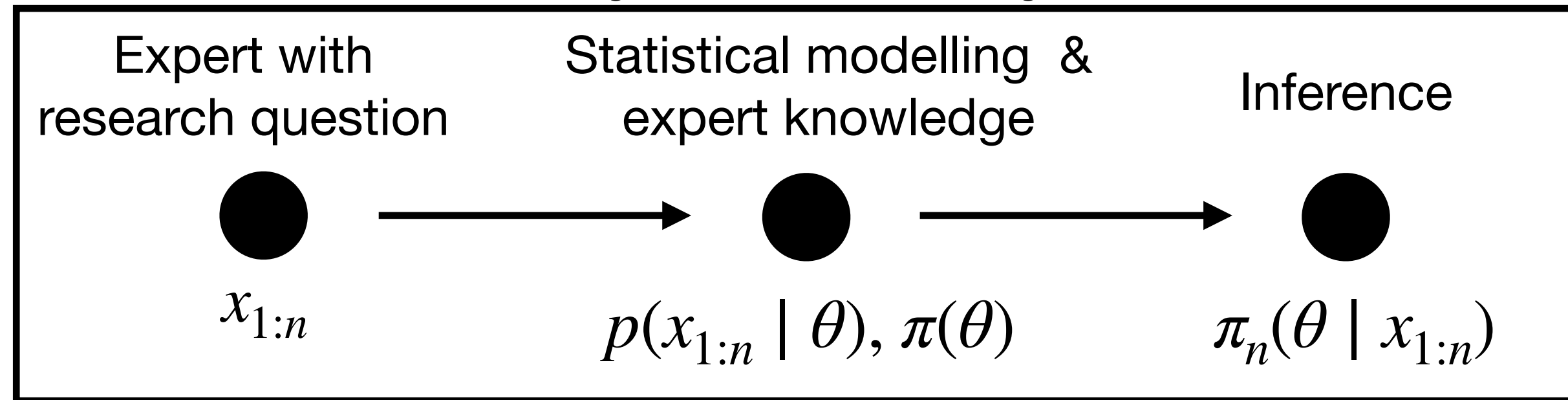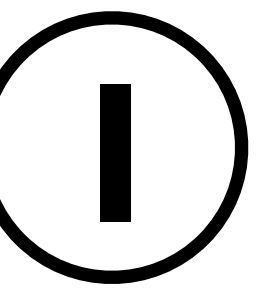
| | |
|---|---|
| (A1) | model well-specified |
| (A2) | prior well-specified |
| (A3) | computationally feasible |

# Case Study: Bayesian ML & Boston Housing Data  Ⓘ

## Traditional Bayesian analysis in science

Expert with research question → Statistical modelling & expert knowledge → Inference

$x_{1:n}$    $p(x_{1:n} \mid \theta), \pi(\theta)$    $\pi_n(\theta \mid x_{1:n})$

## Modern Bayesian ML

Flexible model → Different data/ problems → Inference

$p(x_{1:n} \mid \theta), \pi(\theta)$    $x_{1:n}$    $\pi_n(\theta \mid x_{1:n})$

### Harrison & Rubinfeld (1978)
**Research Question:** influence of air pollution on house prices?

(A1) ✓

parameters of interest

incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay    pollutants    rooms, sqm, ...

measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$

$\pi(\theta) \sim$ hand-crafted by experts    (A2) ✓

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly    (A3) ✓
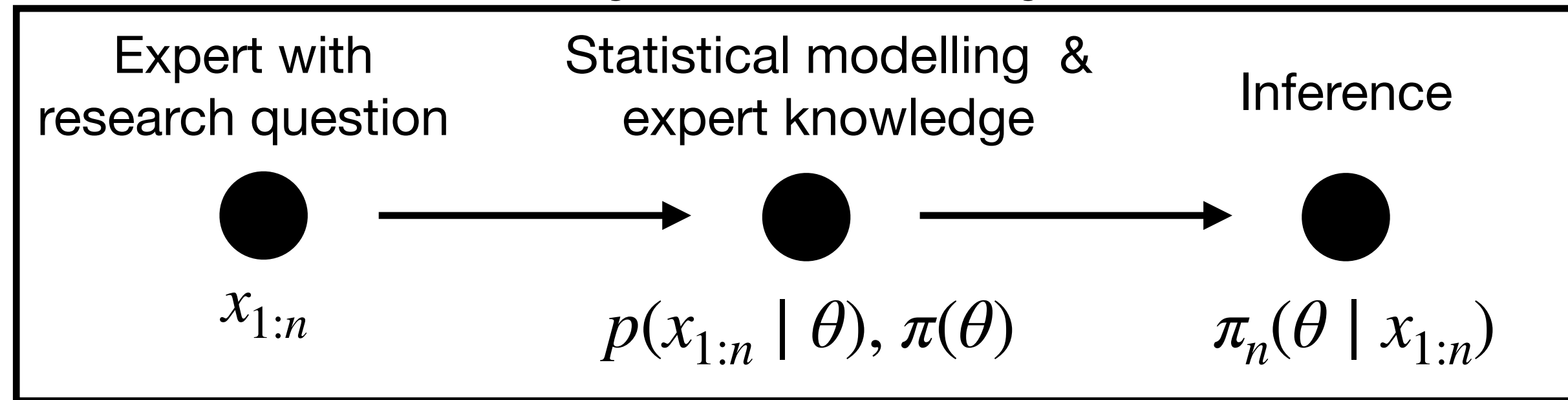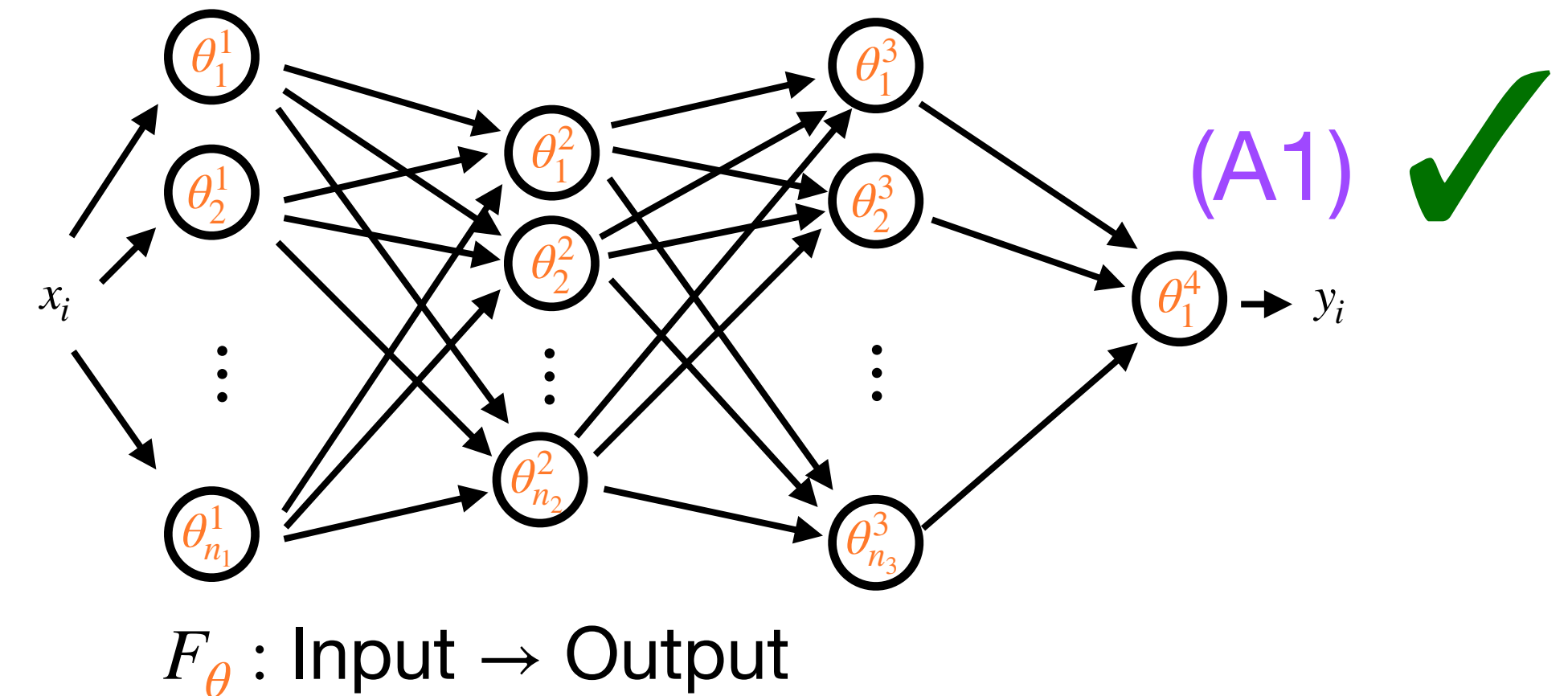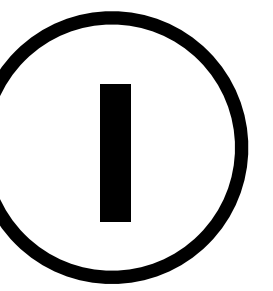
### Pearce et al. (2020) [AISTATS]
**Research Question:** Does my algorithm improve prediction on regression tasks like Boston UCI data?

(A1) ✓

$x_i$    $\theta_1^1, \theta_2^1, \ldots, \theta_{n_1}^1$    $\theta_1^2, \theta_2^2, \ldots, \theta_{n_2}^2$    $\theta_1^3, \theta_2^3, \ldots, \theta_{n_3}^3$    $\theta_1^4 \to y_i$

$F_\theta$ : Input → Output

# Case Study: Bayesian ML & Boston Housing Data  (I)

## Traditional Bayesian analysis in science

Expert with research question → Statistical modelling & expert knowledge → Inference

$x_{1:n}$ → $p(x_{1:n} \mid \theta), \pi(\theta)$ → $\pi_n(\theta \mid x_{1:n})$

## Modern Bayesian ML

Flexible model → Different data/problems → Inference

$p(x_{1:n} \mid \theta), \pi(\theta)$ → $x_{1:n}$ → $\pi_n(\theta \mid x_{1:n})$

**Harrison & Rubinfeld (1978)**
**Research Question:** influence of air pollution on house prices?

(A1) ✓

parameters of interest — incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay — pollutants — rooms, sqm, ... — measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$

$\pi(\theta) \sim$ hand-crafted by experts (A2) ✓

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly (A3) ✓
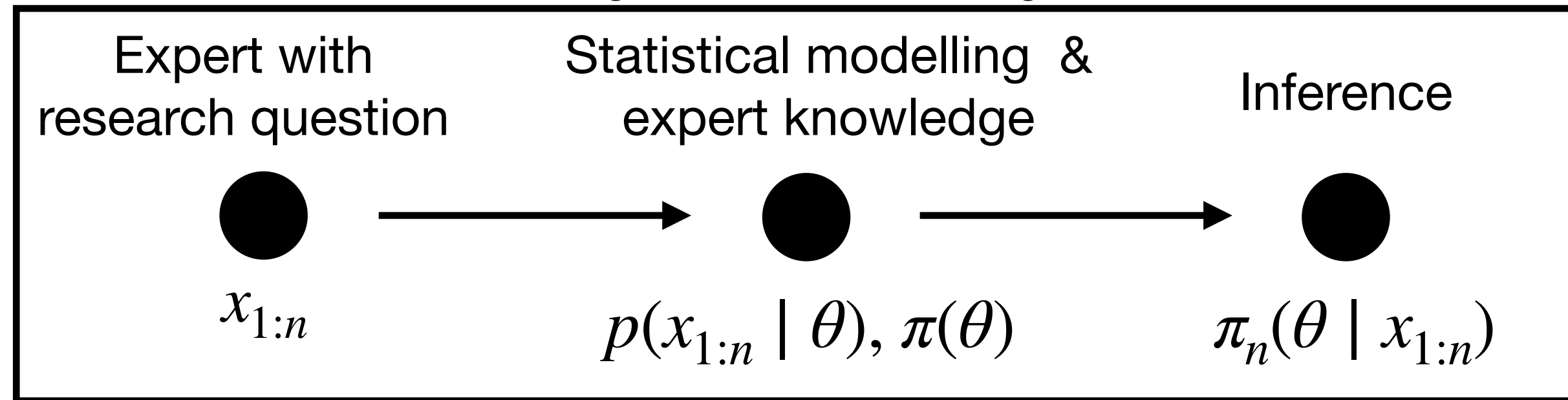
**Pearce et al. (2020) [AISTATS]**
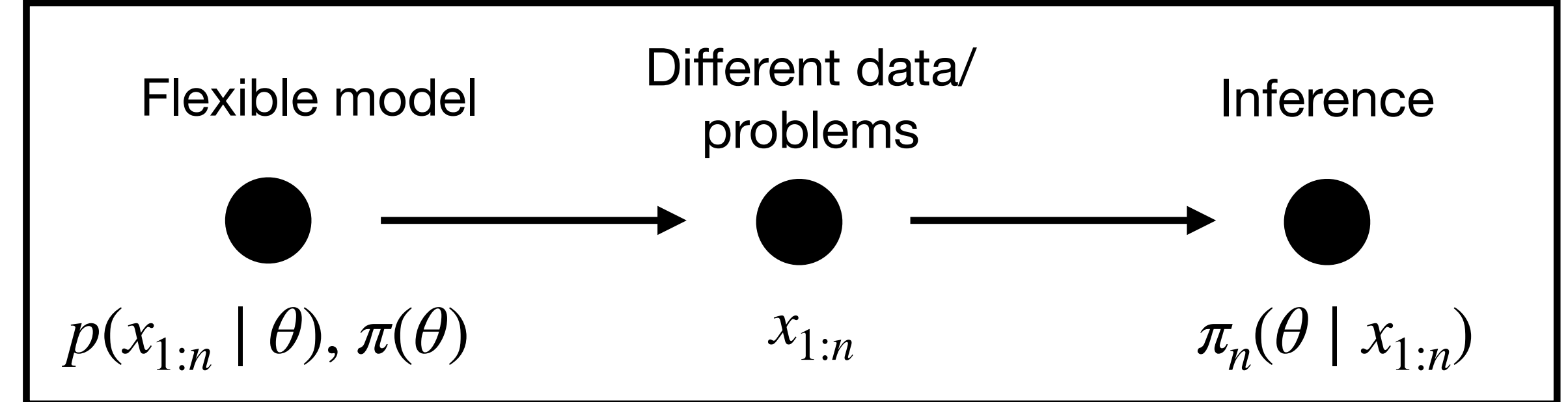**Research Question:** Does my algorithm improve prediction on regression tasks like Boston UCI data?

(A1) ✓



$F_\theta$ : Input → Output

$\pi(\theta) \sim$ Normal (A2) ✗

# Case Study: Bayesian ML & Boston Housing Data  Ⓘ

## Traditional Bayesian analysis in science

Expert with research question → Statistical modelling & expert knowledge → Inference

$x_{1:n}$  →  $p(x_{1:n} \mid \theta), \pi(\theta)$  →  $\pi_n(\theta \mid x_{1:n})$

## Modern Bayesian ML

Flexible model → Different data/ problems → Inference

$p(x_{1:n} \mid \theta), \pi(\theta)$  →  $x_{1:n}$  →  $\pi_n(\theta \mid x_{1:n})$

**Harrison & Rubinfeld (1978)**
**Research Question:** influence of air pollution on house prices?

(A1) ✔

parameters of interest    incidental parameters

$$\log y_i = \sum_{j=1}^{J_1} p_j \log(x_{j,i}) + c_0 + \sum_{j=J_1}^{J_2} c_j \log(x_{j,i}) + \varepsilon_i$$

willingness to pay      pollutants       rooms, sqm, ...

measurement error

$\theta = (c_0, c_2, \ldots, c_{J_1}, p_1, p_2 \ldots p_{J_2})^\top$
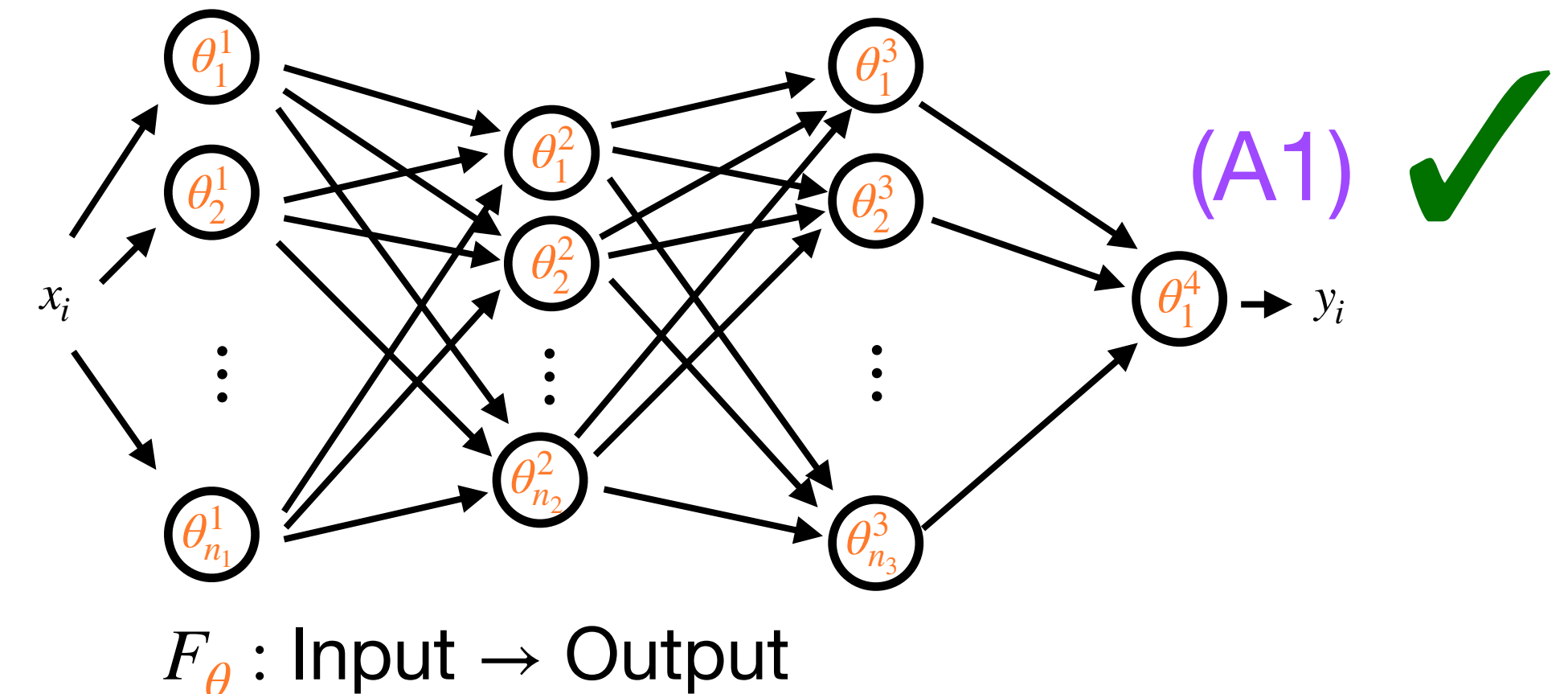
$\pi(\theta) \sim$ hand-crafted by experts  (A2) ✔

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ computed exactly  (A3) ✔
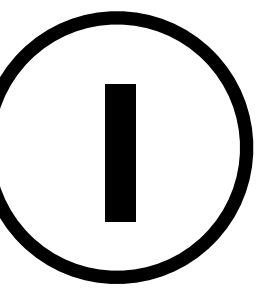
**Pearce et al. (2020) [AISTATS]**
**Research Question:** Does my algorithm improve prediction on regression tasks like Boston UCI data?
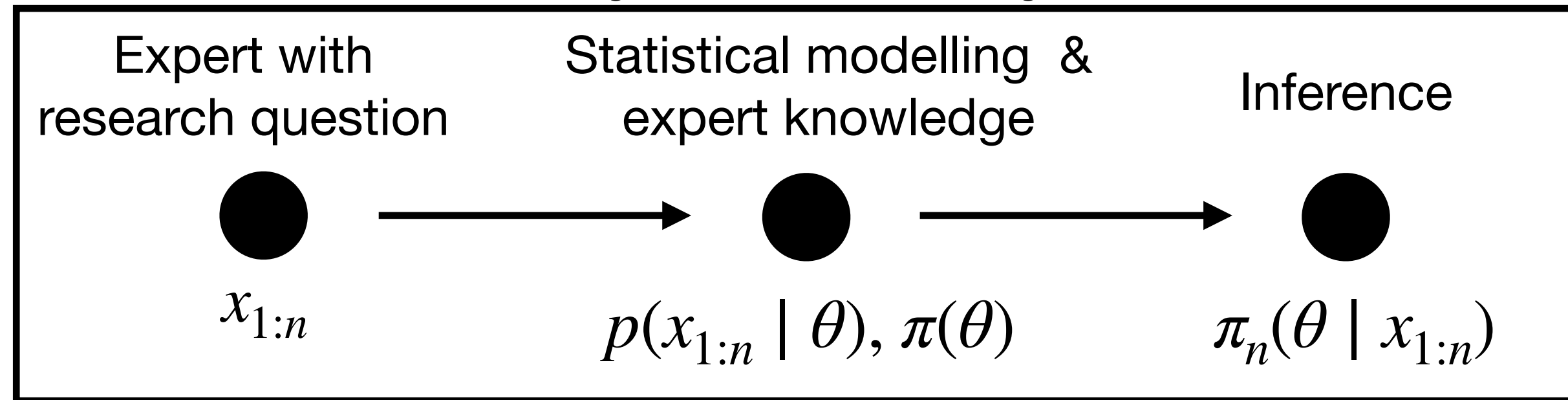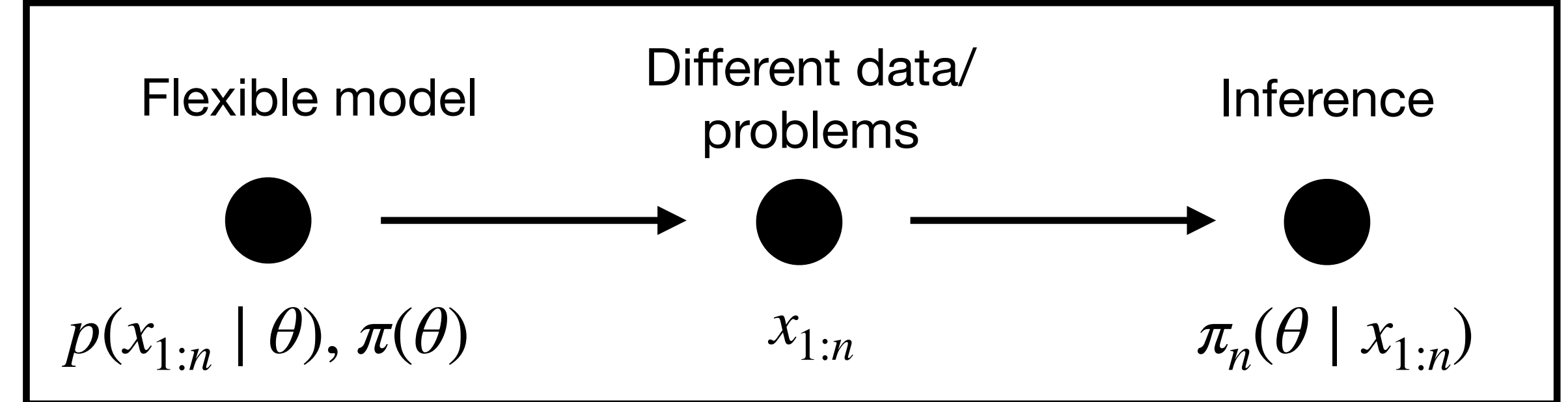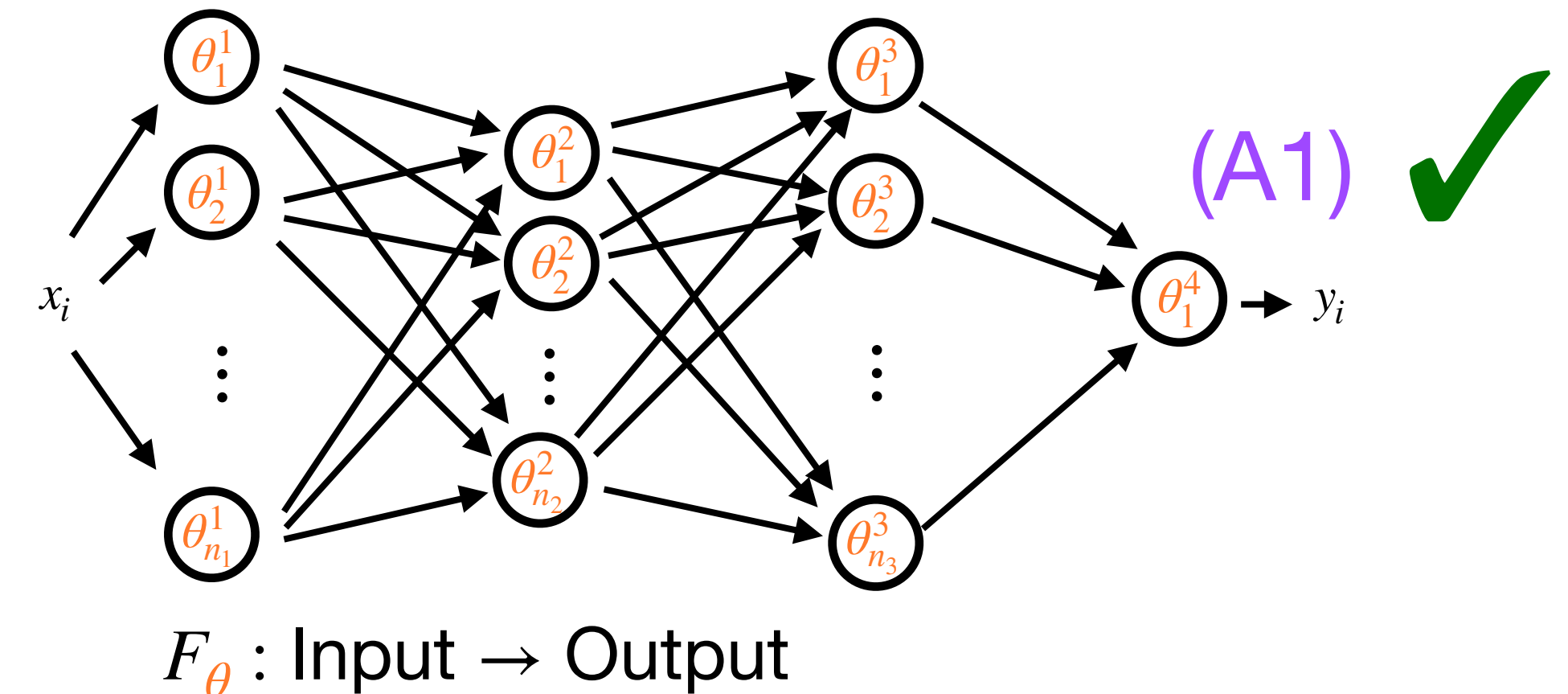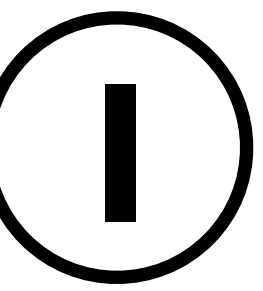
(A1) ✔

$F_\theta :$ Input → Output

$\pi(\theta) \sim$ Normal  (A2)

$\pi_n(\theta \mid x_{1:n}) \longrightarrow$ coarse approximation  (A3)

# Assumptions & Foundations

## Traditional Bayesian analysis in science

Expert with research question

Statistical modelling & expert knowledge

Inference

$x_{1:n}$

$p(x_{1:n} \mid \theta), \pi(\theta)$

$\pi_n(\theta \mid x_{1:n})$

(A1)  model well-specified

(A2)  prior well-specified

(A3)  computationally feasible

## Modern Bayesian ML

Flexible model

Different data/ problems

Inference

$p(x_{1:n} \mid \theta), \pi(\theta)$

$x_{1:n}$

$\pi_n(\theta \mid x_{1:n})$

(A1)  model well-specified

(A2)  prior well-specified

(A3)  computationally feasible

FRAGILE

# Assumptions & Foundations

Ⓘ

**Traditional Bayesian analysis in science**

Expert with research question

Statistical modelling & expert knowledge

Inference

$x_{1:n}$

$p(x_{1:n} \mid \theta), \pi(\theta)$

$\pi_n(\theta \mid x_{1:n})$

**Modern Bayesian ML**

Flexible model

Different data/problems

Inference

$p(x_{1:n} \mid \theta), \pi(\theta)$

$x_{1:n}$

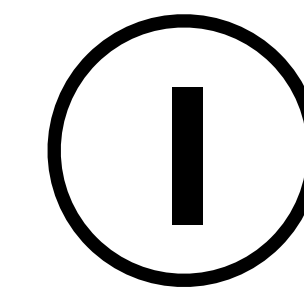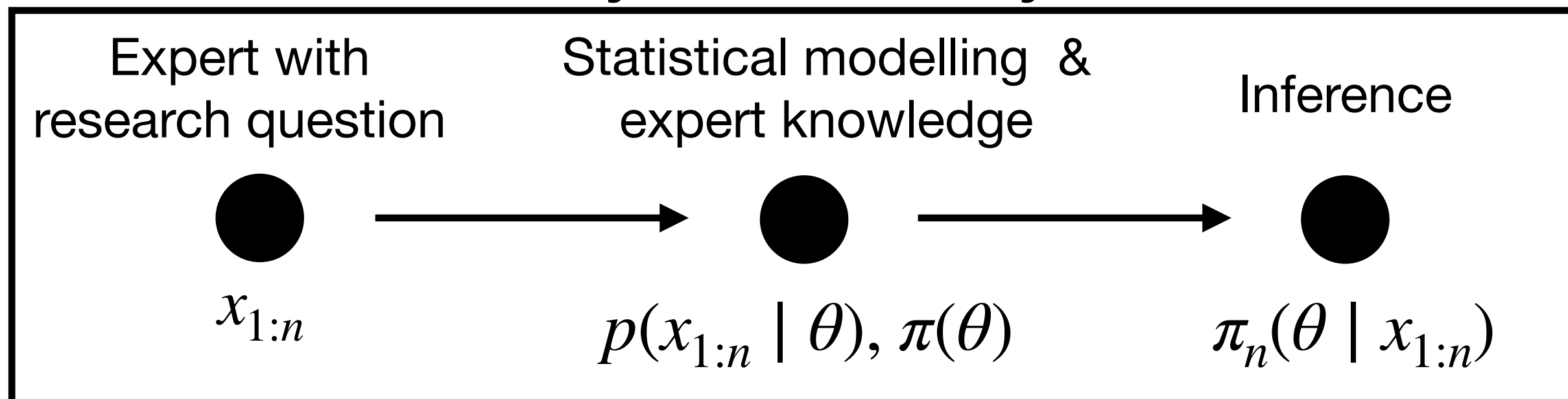$\pi_n(\theta \mid x_{1:n})$

(A1)  model well-specified

(A2)  prior well-specified

(A3)  computationally feasible

(A1)  model well-specified

(A2)  prior well-specified

(A3)  computationally feasible

FRAGILE

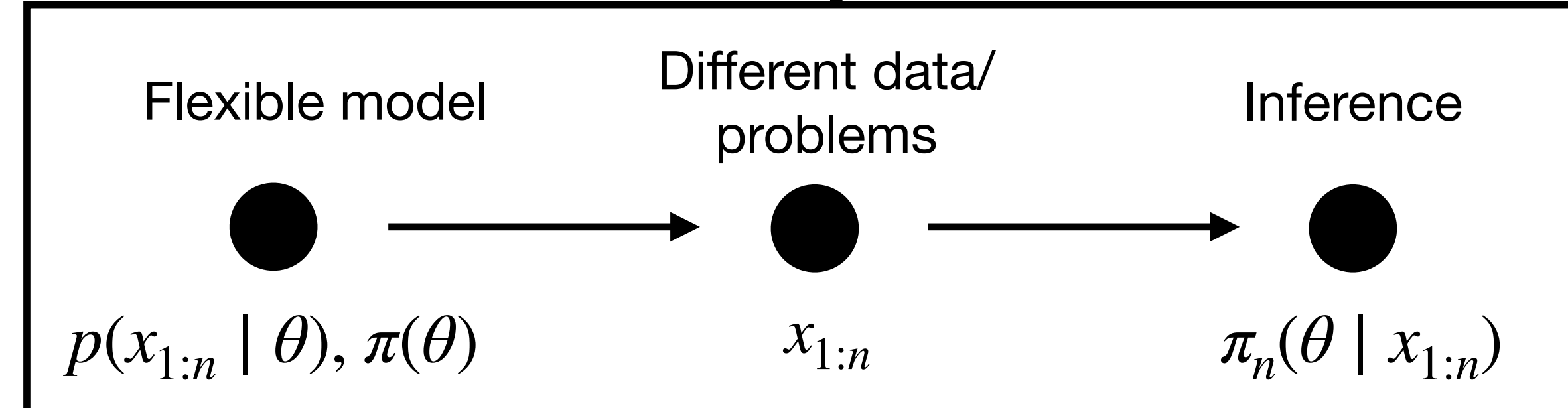## Post-Bayesian Approaches ask:
Can we keep benefits of Bayesianism without these assumptions???

**This seminar is an attempt to organise ourselves under a common banner!!!**

# Part II: What is the (post-Bayesian) Aspirin?

# Post-Bayesian Inference

**Possible belief updates**

×

Bayes' Posterior  (A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

(A1)   model well-specified
(A2)   prior well-specified
(A3)   computationally feasible

# Post-Bayesian Inference

space of priors

space of hyperparameters

Belief updates $= \left\{ \mathscr{B} : \mathscr{P}(\Theta) \times \mathscr{X}^n \times \mathscr{H} \longrightarrow \mathscr{P}(\Theta) \right\}$

data space

space of posteriors

$\times$

Possible belief updates

Bayes' Posterior    (A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

(A1)    model well-specified
(A2)    prior well-specified
(A3)    computationally feasible

# Post-Bayesian Inference

space of priors

space of hyperparameters

Belief updates $= \left\{ \mathscr{B} : \mathscr{P}(\Theta) \times \mathscr{X}^n \times \mathscr{H} \longrightarrow \mathscr{P}(\Theta) \right\}$

data space

space of posteriors

$$p(x_{1:n} \mid \theta) \longrightarrow p(x_{1:n} \mid \theta)^{\lambda}, \ \lambda > 0$$
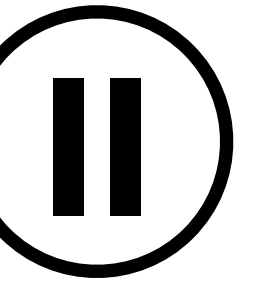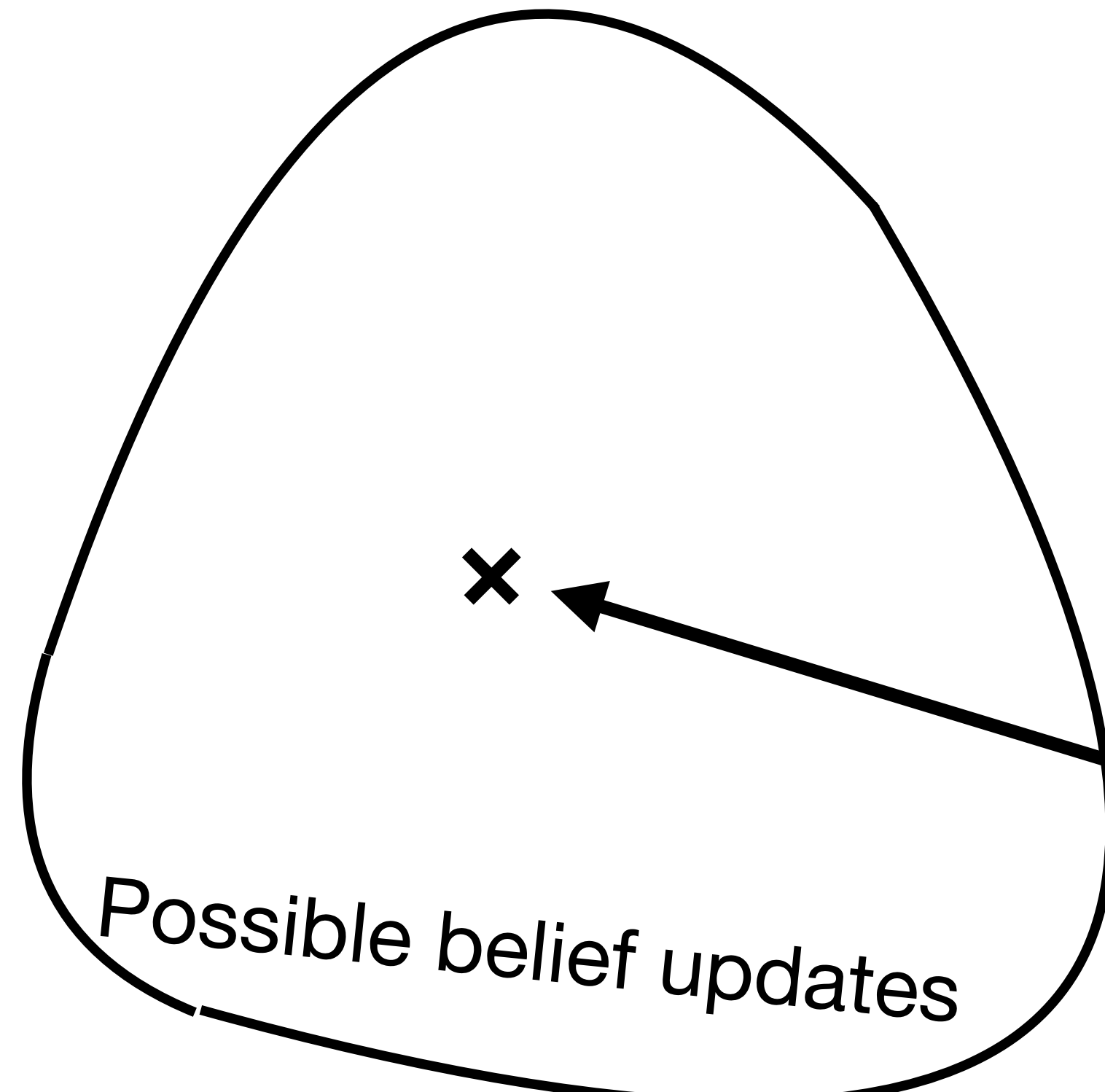


Possible belief updates

(A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

(A1)  model well-specified
(A2)  prior well-specified
(A3)  computationally feasible

# Post-Bayesian Inference

space of priors

space of hyperparameters

$$\text{Belief updates} \; = \; \Big\{ \mathscr{B} : \mathscr{P}(\Theta) \; \times \; \mathscr{X}^n \; \times \; \mathscr{H} \; \longrightarrow \; \mathscr{P}(\Theta) \Big\}$$

data space

space of posteriors

*[See Grünwald (2011)]*

Power/Fractional/
Cold Posterior

(A1), (A2), (A3)

$$p(x_{1:n} \mid \theta) \longrightarrow p(x_{1:n} \mid \theta)^{\lambda}, \; \lambda > 0$$

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

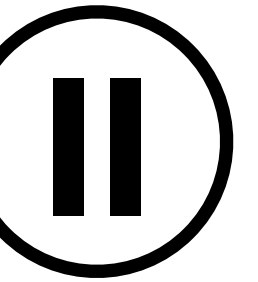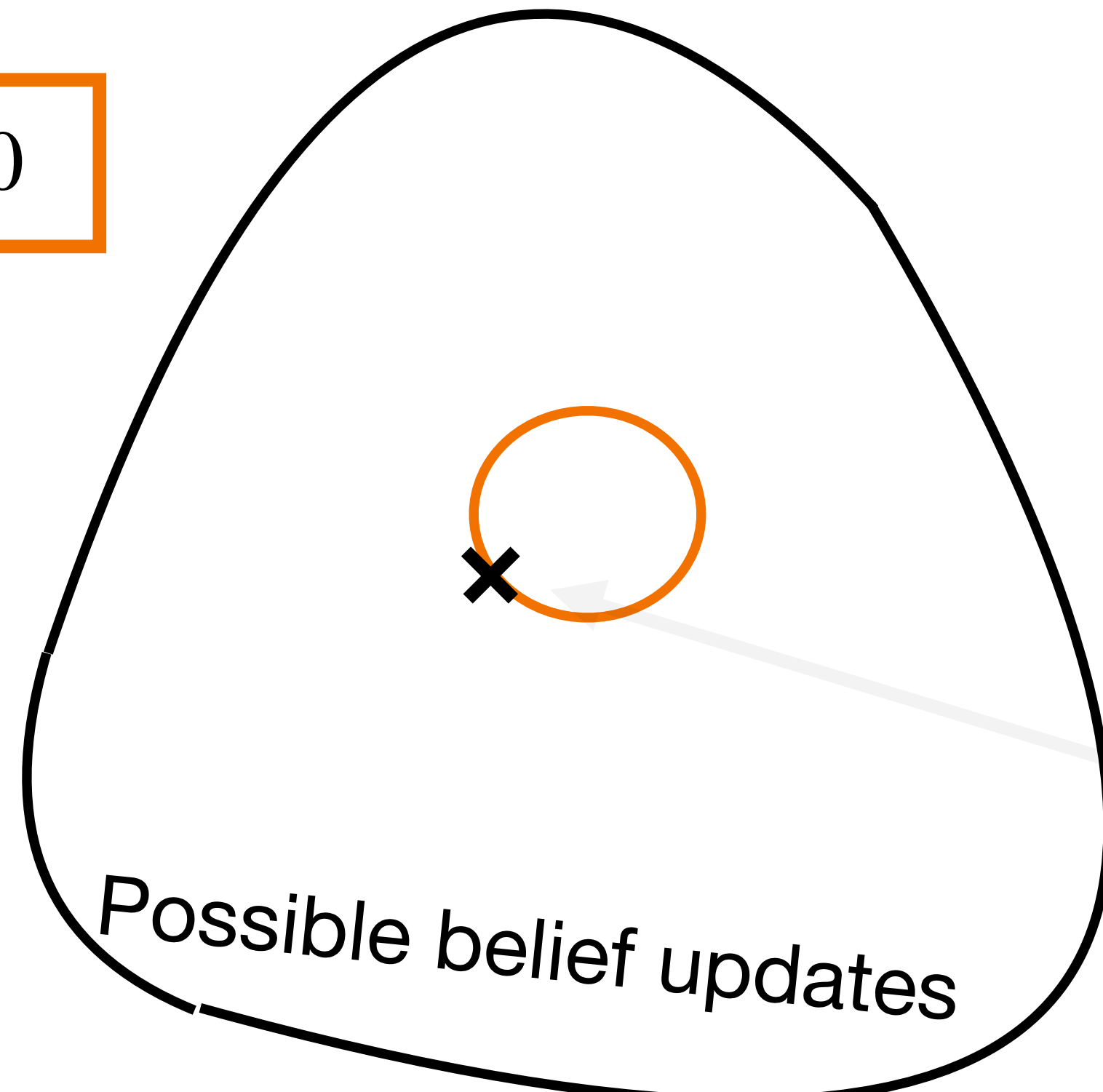Possible belief updates

(A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

(A1) model well-specified
(A2) prior well-specified
(A3) computationally feasible

# Post-Bayesian Inference

II

$$p(x_{1:n} \mid \theta) \longrightarrow p(x_{1:n} \mid \theta)^{\lambda}, \ \lambda > 0$$

$$p(x_{1:n} \mid \theta) \longrightarrow \exp\{-L(x_{1:n}, p_{\theta})\}, \ \text{loss } L$$

(A1)   model well-specified
(A2)   prior well-specified
(A3)   computationally feasible

Possible belief updates

(A1), (A2), (A3)

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$
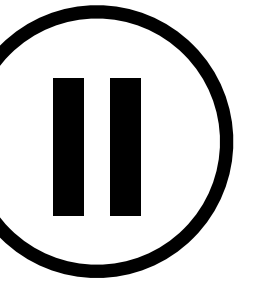
(A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

# Post-Bayesian Inference

*[See Bissiri, Holmes & Walker (2016)]*

Gibbs/Generalised/
Pseudo Posterior

~~(A1)~~, (A2), (A3)

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

$$p(x_{1:n} \mid \theta) \longrightarrow p(x_{1:n} \mid \theta)^\lambda, \ \lambda > 0$$

$$p(x_{1:n} \mid \theta) \longrightarrow \exp\{-\mathsf{L}(x_{1:n}, p_\theta)\}, \ \text{loss } \mathsf{L}$$

Possible belief updates

~~(A1)~~, (A2), (A3)

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta) d\theta}$$

(A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

(A1)   model well-specified
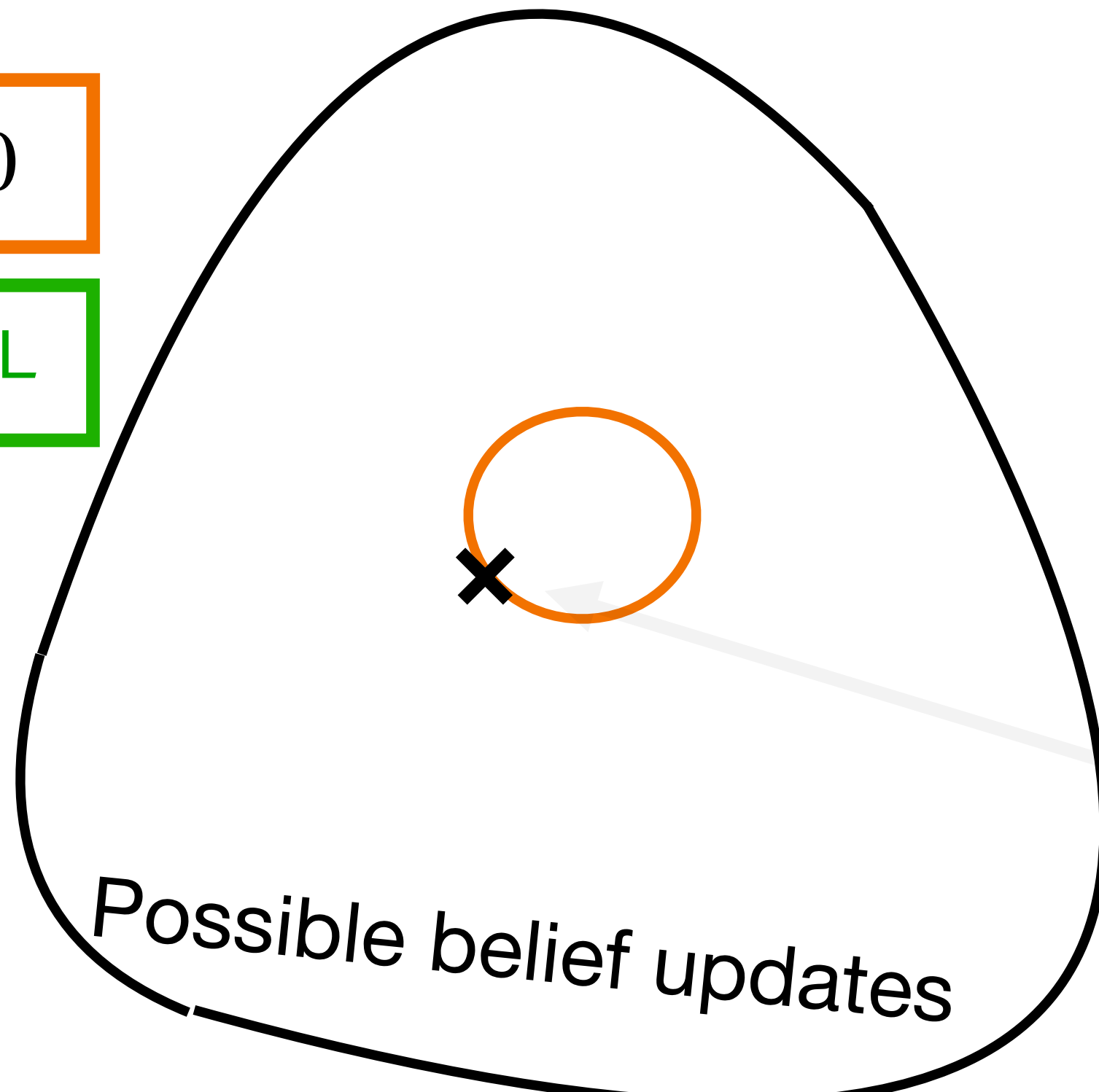(A2)   prior well-specified
(A3)   computationally feasible

# Post-Bayesian Inference

Optimisation-centric posteriors /
Generalised Variational Inference

~~(A1)~~, ~~(A2)~~, ~~(A3)~~

$$q_n^*(\theta) = \arg\min_{q \in \mathbb{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + \mathrm{D}(q, \pi) \right\};$$

$$\mathbb{Q} \subseteq \mathscr{P}(\Theta)$$

Data-fitting    Prior regularisation

*[See Knoblauch, Jewson, & Damoulas (2019/2022)]*

$$p(x_{1:n} \mid \theta) \longrightarrow p(x_{1:n} \mid \theta)^\lambda, \ \lambda > 0$$

$$\begin{aligned} \mathrm{KL} &\longrightarrow \mathrm{D} \\ \mathscr{P}(\Theta) &\longrightarrow \mathbb{Q} \end{aligned}$$

(A1)   model well-specified
(A2)   prior well-specified
(A3)   computationally feasible
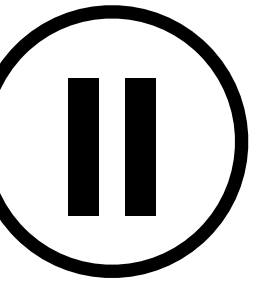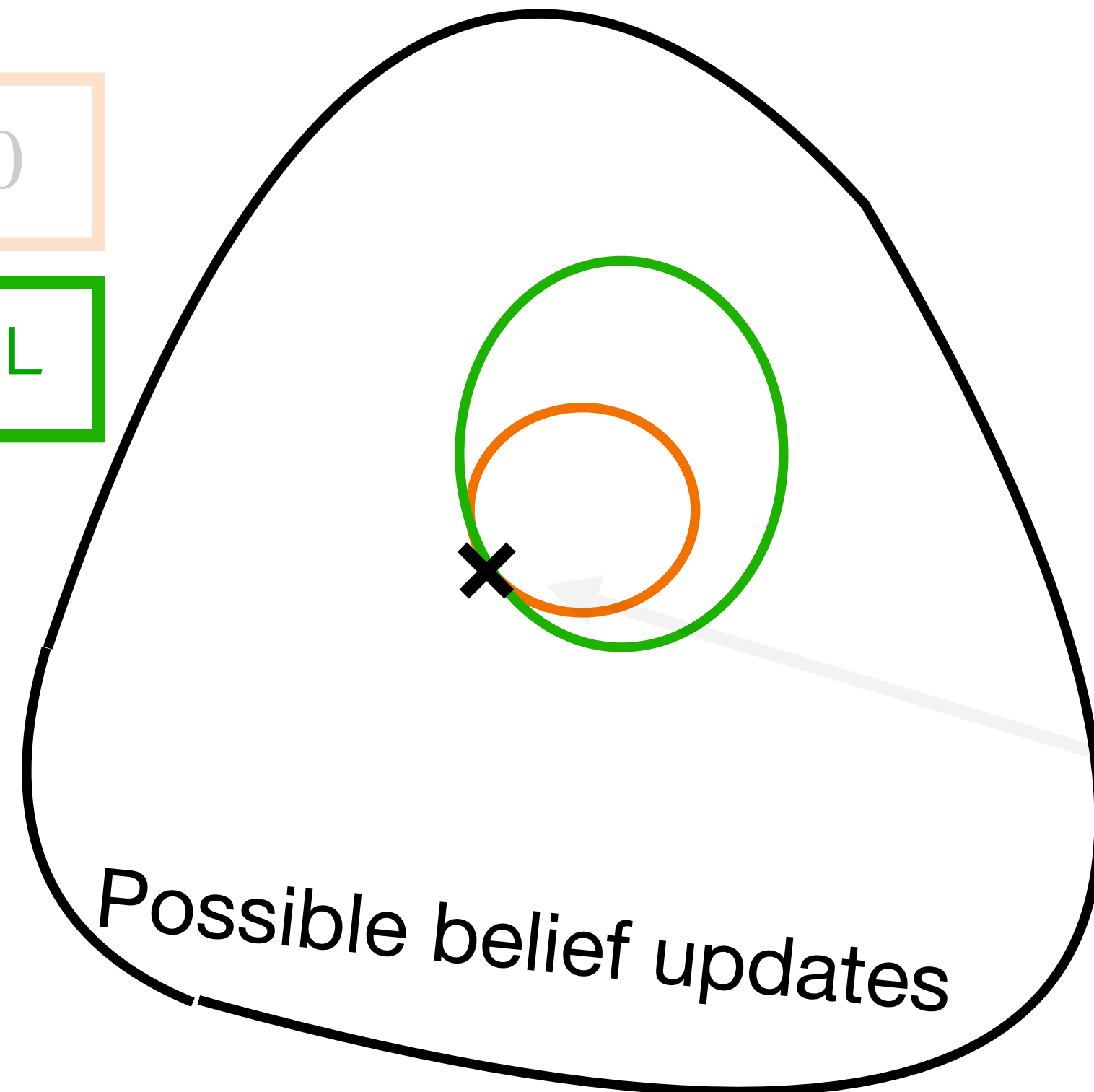
$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta) d\theta}$$
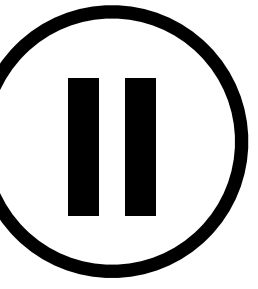
(A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$
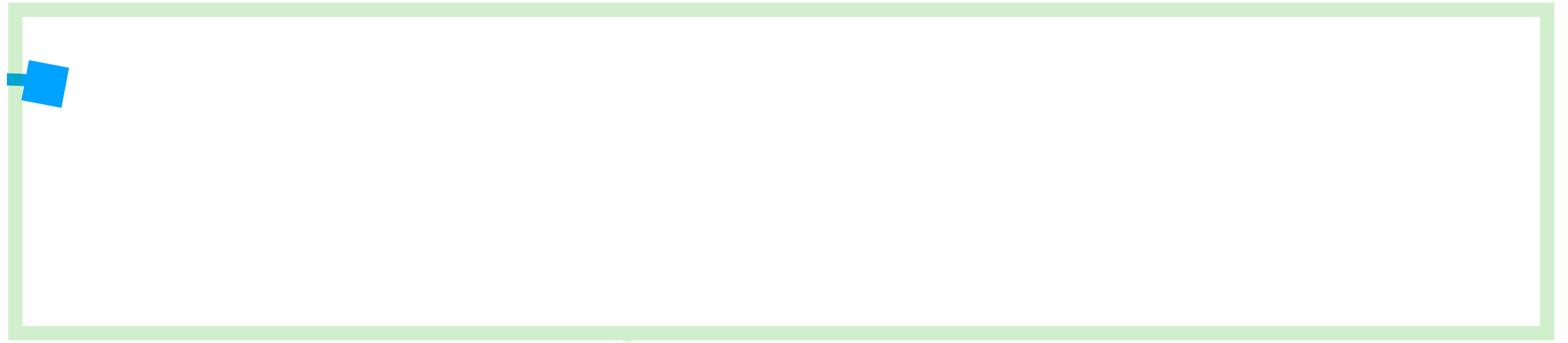
Possible belief updates

# Post-Bayesian Inference

Ⓘ**II**

**Optimisation-centric posteriors /
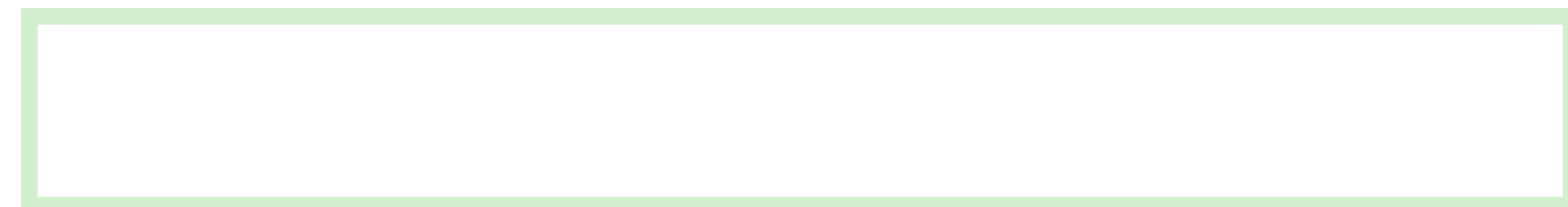Generalised Variational Inference** ~~(A1), (A2), (A3)~~

$$q_n^*(\theta) = \arg\min_{q \in \mathcal{Q}} \left\{ \underbrace{\mathscr{L}(q, x_{1:n})}_{\text{Data-fitting}} + \underbrace{D(q, \pi)}_{\text{Prior regularisation}} \right\};$$
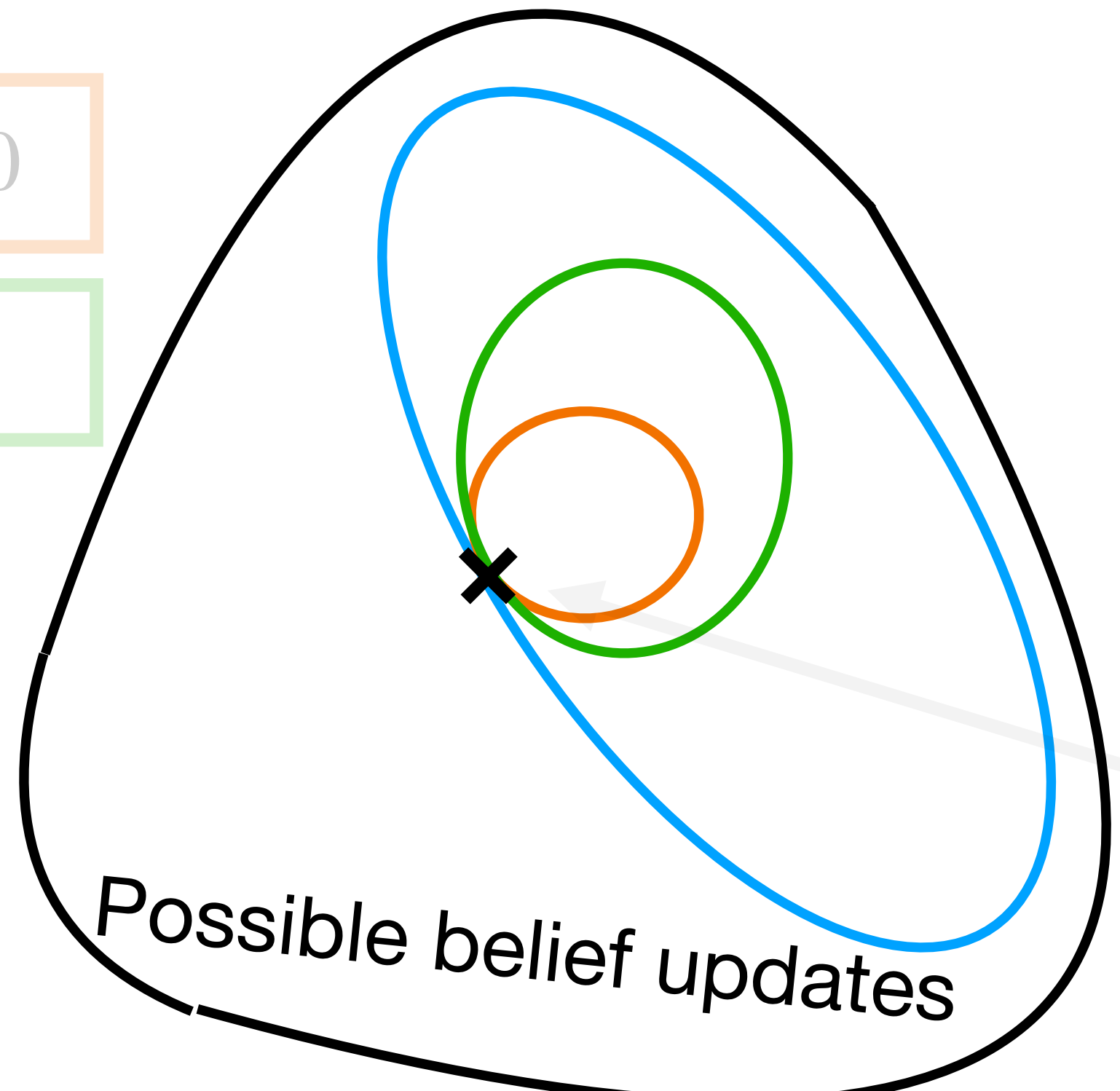
$\mathcal{Q} \subseteq \mathscr{P}(\Theta)$

**Gibbs/Generalised/
Pseudo Posterior** ~~(A1)~~, (A2), (A3)

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Martingale posteriors &
resampling-based
approaches** *[See Fong, Holmes, & walker (2023)]* ~~(A1), (A2), (A3)~~

For $i = 1, 2, \dots$

$X_{n+i+1} \sim p(X_{n+i} \mid x_{1:n}, X_{n+1:n+i})$

$\theta^\infty = \mathrm{argmin}_{\theta \in \Theta} \mathsf{L}\left([x_{1:n}, X_{n+1:\infty}], \theta\right)$

**Power/Fractional/
Cold Posterior** ~~(A1)~~, (A2), (A3)

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta) d\theta}$$

| (A1) | model well-specified |
| (A2) | prior well-specified |
| (A3) | computationally feasible |

Possible belief updates

**Bayes' Posterior** (A1), (A2), (A3)

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

# Chapter 2

## Resampling & Martingale Posteriors
## (06/05—15/07)

Martingale posteriors & resampling-based approaches

~~(A1)~~, ~~(A2)~~, ~~(A3)~~

For $i = 1, 2, \ldots$

$$X_{n+i+1} \sim p(X_{n+i} \mid x_{1:n}, X_{n+1:n+i})$$

$$\theta^\infty = \operatorname{argmin}_{\theta \in \Theta} \mathsf{L}\left([x_{1:n}, X_{n+1:\infty}], \theta\right)$$

**Dr. Edwin Fong
(University of
Hong Kong)**

Optimisation-centric posteriors /
Generalised Variational Inference     (A1), (A2), (A3)

$$q_n^*(\theta) = \arg\min_{q \in \mathcal{Q}} \left\{ \underbrace{\mathcal{L}(q, x_{1:n})}_{\text{Data-fitting}} + \underbrace{D(q, \pi)}_{\text{Prior regularisation}} \right\};$$

$\mathcal{Q} \subseteq \mathcal{P}(\Theta)$

Data-fitting     Prior regularisation

Gibbs/Generalised/
Pseudo Posterior          (A1), (A2), (A3)

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

# Chapter 3
## PAC-Bayes
## (after summer break)

## Prof. Pierre Alquier
## (ESSEC Singapore)

For $i = 1, 2, \ldots$

$X_{n+i+1} \sim p(X_{n+i} \mid x_{1:n}, X_{n+1:n+i})$

$(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$
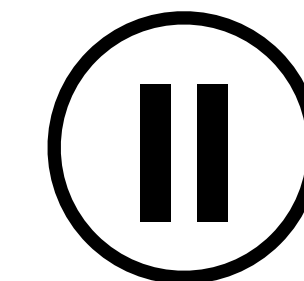
(A1), (A2), (A3)

$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$

(A3)   computationally feasible

Possible belief updates

Optimisation-centric posteriors /
Generalised Variational Inference        (A1), (A2), (A3)

$$q_n^*(\theta) = \arg\min_{q \in \mathcal{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + D(q, \pi) \right\};$$

$\mathcal{Q} \subseteq \mathscr{P}(\Theta)$     Data-fitting     Prior regularisation

Gibbs/Generalised/
Pseudo Posterior        (A1), (A2), (A3)

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

Power/Fractional/
Cold Posterior        (A1), (A2), (A3)

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta) d\theta}$$

# Chapter 1
## Generalised Bayes (11/02—22/04)

For $i = 1, 2, \ldots$
$X_{n+i+1} \sim p(X_{n+i} \mid x_{1:n}, X_{n+1:n+i})$

(A1), (A2), (A3)

Possible belief updates

(A3)  computationally feasible

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

# Structure of Chapter 1

**Today: Overview of post-Bayesian ideas & generalised Bayes**

**25/02:  Theoretical foundations
 (Prof. David Frazier)**

**11/03:  Learning rate selection & the power posterior
(Prof. Ryan Martin)**

**25/03:  Prediction-centric approaches
(Prof. Chris Oates)**

**08/04:  Coarsened Bayes & applications for biomedical problems
(Prof. David Dunson)**

**22/04:  From generalised Bayes to Martingale Posteriors
(Prof. Chris Holmes)**

# Part III: Basics of generalised Bayes

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$
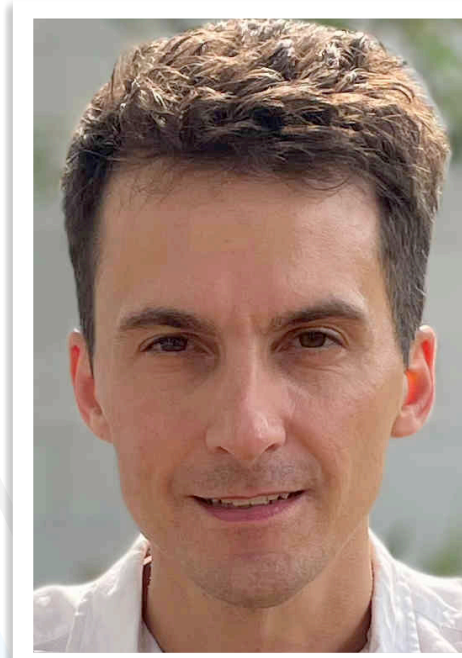
Gibbs/Generalised/
Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

Optimisation-centric posteriors /
Generalised Variational Inference

$$q_n^*(\theta) = \arg\min_{q \in \mathcal{Q}} \left\{ \mathscr{L}(q, x_{1:n}) \quad + \quad \mathsf{D}(q, \pi) \right\};$$

$$\mathcal{Q} \subseteq \mathscr{P}(\Theta)$$

Data-fitting    Prior regularisation

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

**Q: What does it do?**

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

**Q: What does it do?**
**A: Trades off prior vs data**

posterior  $\pi_n^{(\lambda)}(\theta \mid x_{1:n})$

prior  $\pi(\theta)$

likelihood  $p(x_{1:n} \mid \theta)^{\lambda}$



$\lambda = 0.5$

$\lambda = 1$

$\lambda = 2$

$\theta$

Picture from Kallionen, Paananen, Bürkner, & Vehtari (2023)

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

**Q: Why do it do?**

Regression model (misspecified):

$$p(y_i \mid \theta, x_i) = \mathcal{N}\left( y_i; \sum_{d=1}^{\overset{50}{D}} \theta_i \, x_{i,d} \, , \sigma^2 \right)$$



Squared Risk vs. Number of observations $n$

$\frac{n}{D} < 6$    $\frac{n}{D} > 6$

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)d\theta}$$

**Q: Why do it do?**

Regression model (misspecified):

$$p(y_i \mid \theta, x_i) = \mathcal{N}\left(y_i; \sum_{d=1}^{\overset{50}{D}} \theta_i\, x_{i,d}\, , \sigma^2\right)$$



Bayes Posterior & Maximum A Posteriori

Squared Risk

$\frac{n}{D} < 6$  $\frac{n}{D} > 6$

True Model   Number of observations $n$

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

**Q: Why do it do?**

Regression model (misspecified):

$$p(y_i \mid \theta, x_i) = \mathcal{N}\left(y_i; \sum_{d=1}^{\overset{50}{D}} \theta_i \, x_{i,d} \, , \sigma^2\right)$$

The **'Safe Bayes' effect** (see Grünwald, 2012)
(picture from Grünwald & van Ommen, 2017)

Bayes Posterior & Maximum A Posteriori

Power posteriors
($\lambda$ chosen with
'Safe Bayes' approach)



Squared Risk

$\frac{n}{D} < 6$          $\frac{n}{D} > 6$

True Model          Number of observations $n$

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta) d\theta}$$
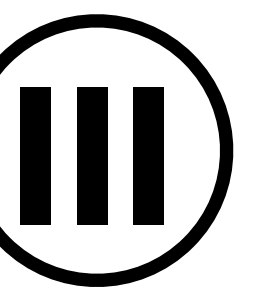
**Q: Why do it do?**
**A: better risk properties/predictions**

*if $\dfrac{n}{D}$ is small*

Regression model (misspecified):

$$p(y_i \mid \theta, x_i) = \mathcal{N}\left(y_i; \sum_{d=1}^{D} \theta_i \, x_{i,d}, \sigma^2\right)$$

(with 50 circled over the sum)

The **'Safe Bayes' effect** (see Grünwald, 2012)
(picture from Grünwald & van Ommen, 2017)

Bayes Posterior & Maximum A Posteriori

Power posteriors
($\lambda$ chosen with
'Safe Bayes' approach)



Squared Risk

$\dfrac{n}{D} < 6$        $\dfrac{n}{D} > 6$
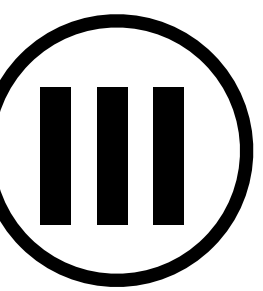
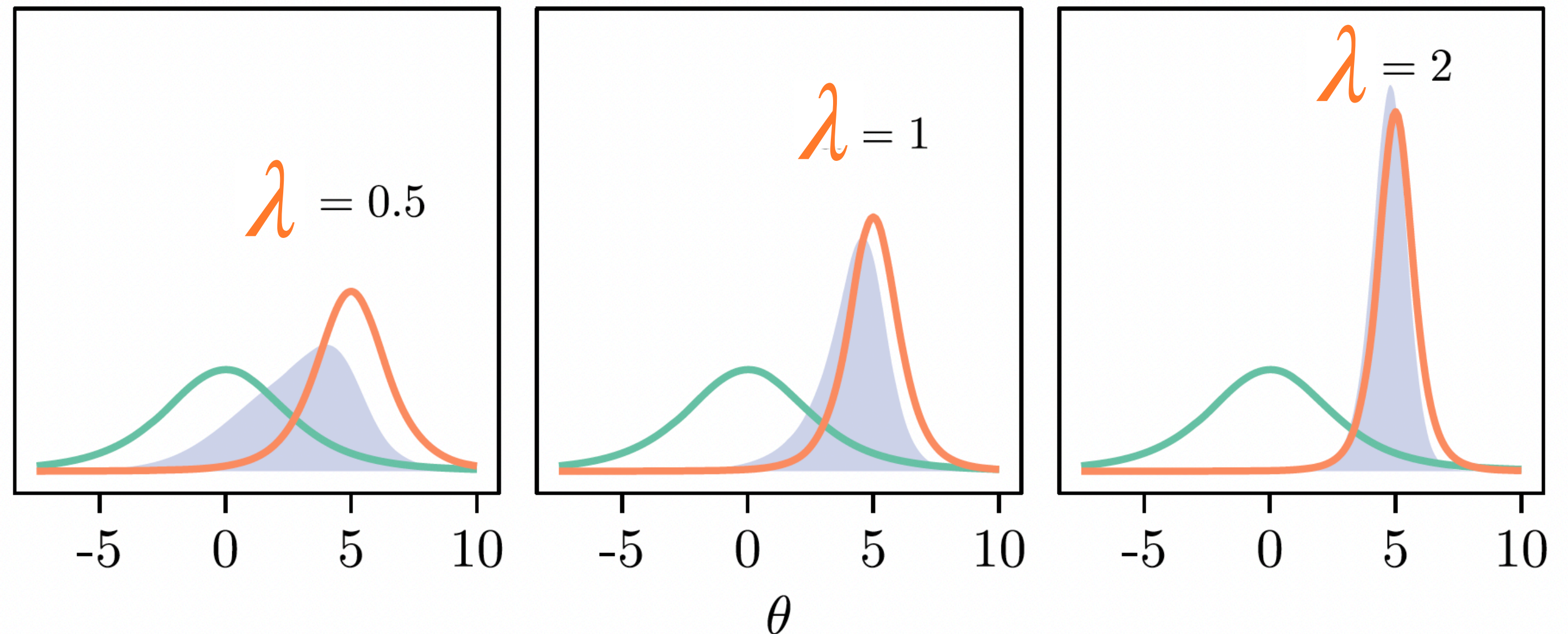True Model        Number of observations $n$

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

**Q: What else have we discovered?**

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$
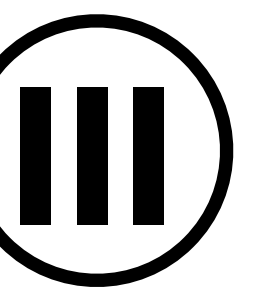
**25/02**

**Q: What else have we discovered?**

**Power posterior & their variational approximations concentrate in situations where standard Bayes wouldn't**

*Bhattacharya, Pati, & Yang (2019)*
*Alquier & Ridgeway (2020)*
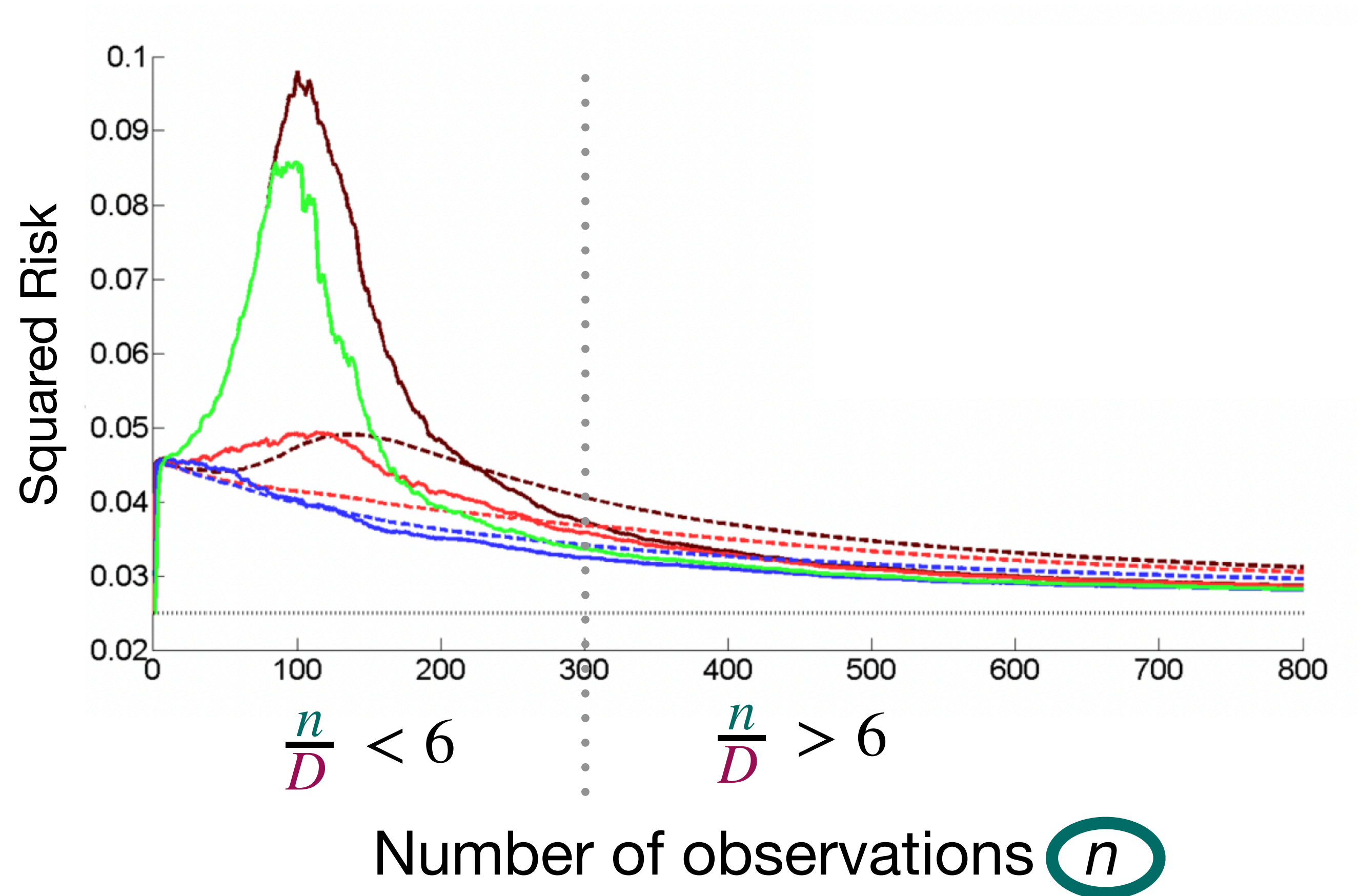*Yang, Pati, & Bhattacharya (2020)*

# Basics: power posteriors
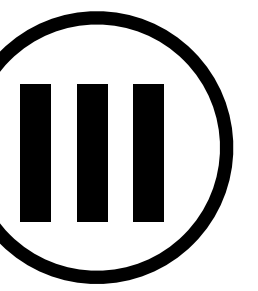
Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

25/02

**It is often surprisingly difficult to choose $\lambda$**

*Grünwald (2012)*
*Lyddon, Holmes, & Walker (2019)*
*Wu & Martin (2023)*

**Q: What else have we discovered?**

**Power posterior & their variational approximations concentrate in situations where standard Bayes wouldn't**

*Bhattacharya, Pati, & Yang (2019)*
*Alquier & Ridgeway (2020)*
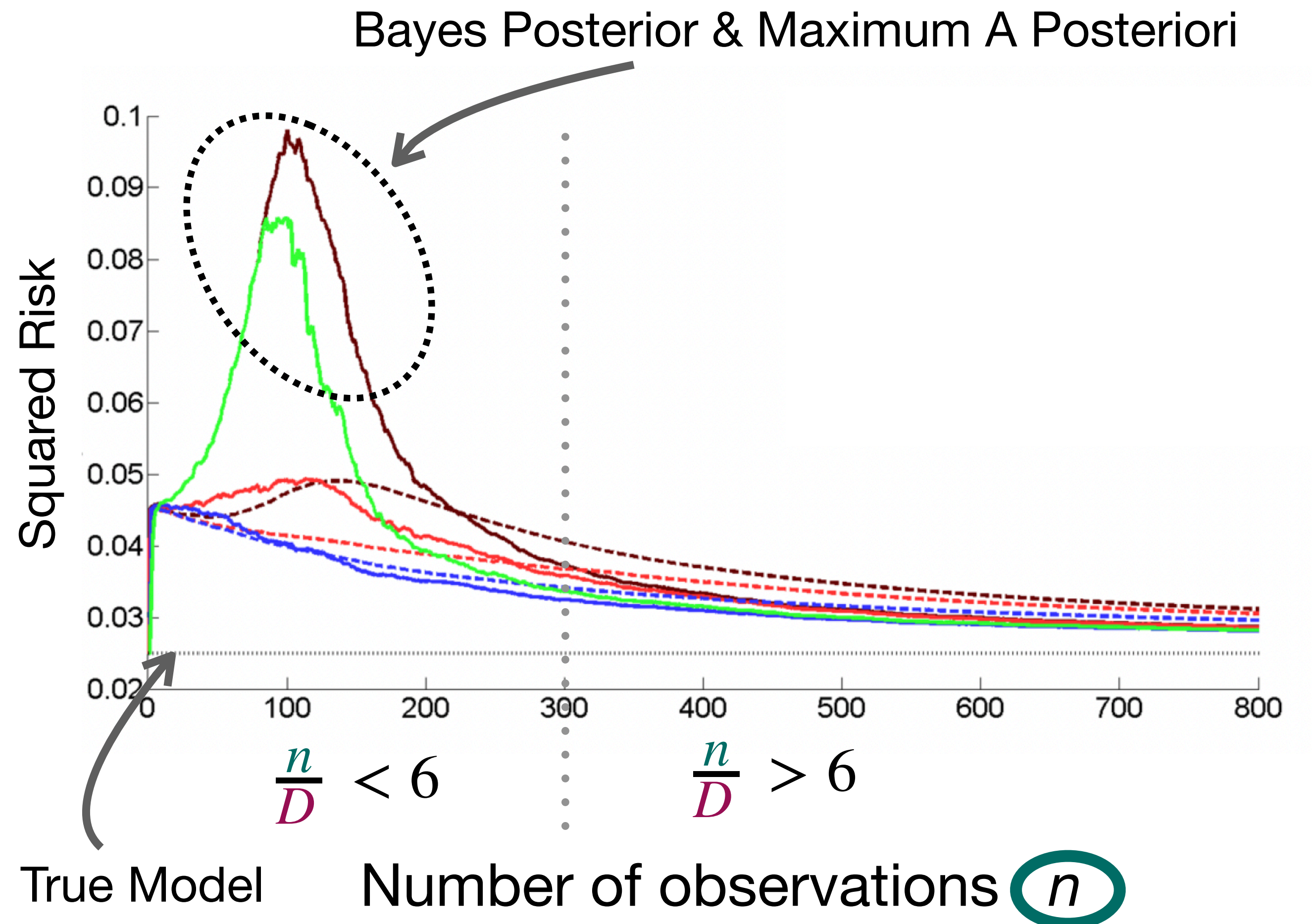*Yang, Pati, & Bhattacharya (2020)*

# Basics: power posteriors
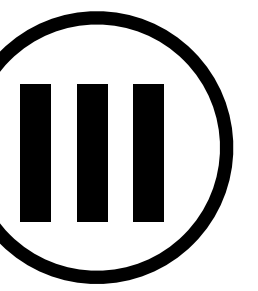
Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^\lambda \cdot \pi(\theta) d\theta}$$

**Q: What else have we discovered?**

**Power posterior & their variational approximations concentrate in situations where standard Bayes wouldn't**

*Bhattacharya, Pati, & Yang (2019)*
*Alquier & Ridgeway (2020)*
*Yang, Pati, & Bhattacharya (2020)*

**It is often surprisingly difficult to choose $\lambda$**

*Grünwald (2012)*
*Lyddon, Holmes, & Walker (2019)*
*Wu & Martin (2023)*

**Predictive/robustness gains vanish provably & very quickly for even moderate $\dfrac{n}{D}$**

*Medina, Olea, Rush, & Velez (2022)*
*McLatchie, Fong, Frazier, & Knoblauch (2024)*

11/03

25/02

25/03

# Basics: power posteriors
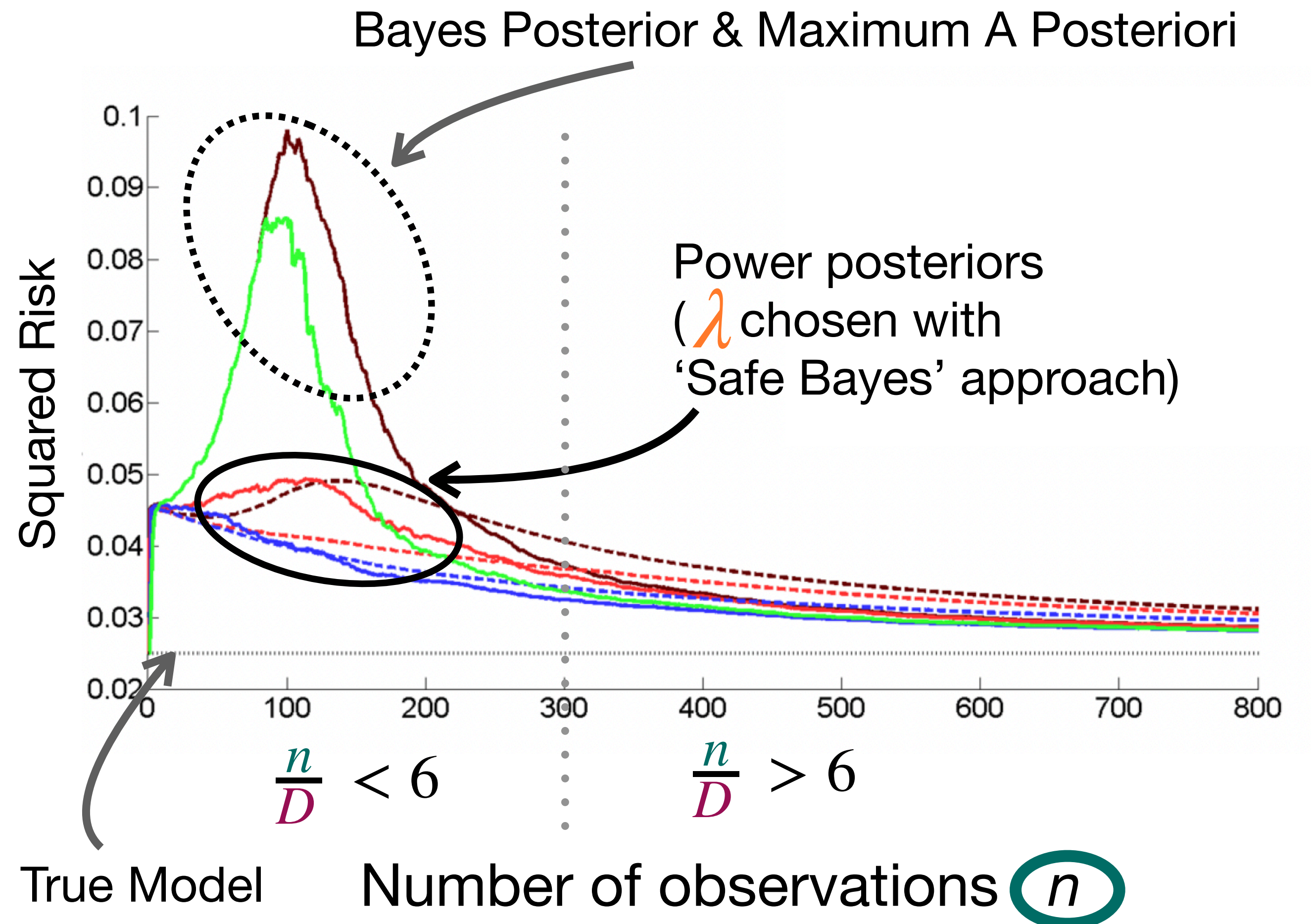
Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$
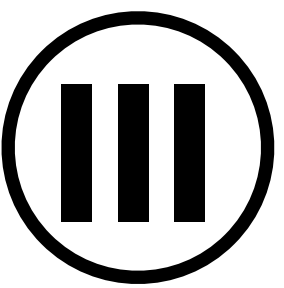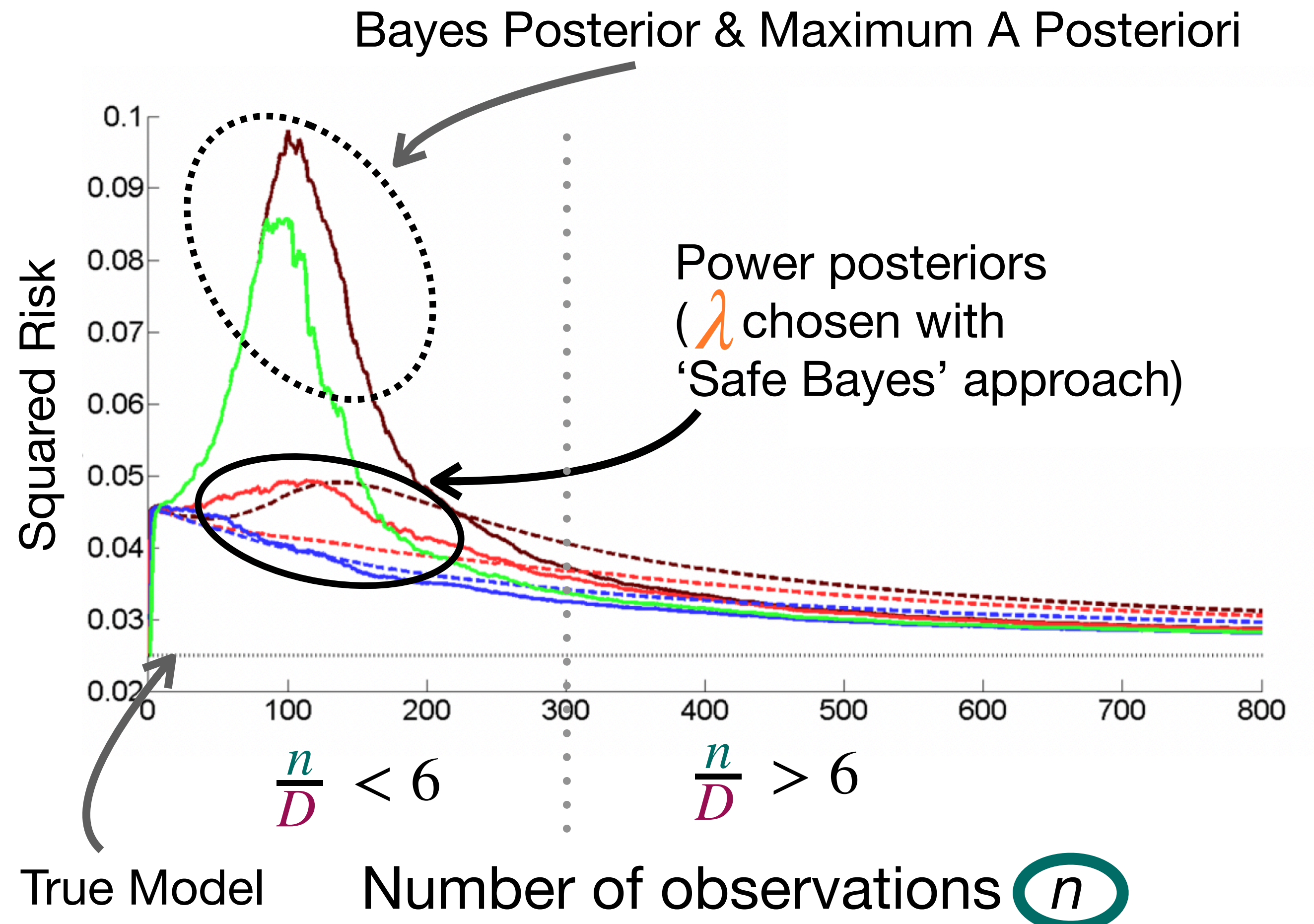
**25/02**

**11/03**

**Q: What else have we discovered?**

**Power posterior & their variational approximations concentrate in situations where standard Bayes wouldn't**

*Bhattacharya, Pati, & Yang (2019)*
*Alquier & Ridgeway (2020)*
*Yang, Pati, & Bhattacharya (2020)*

**Power posterior $\approx$ 'Coarsened' posterior (conditioning on a neighbourhood of observed data)**

*Miller & Dunson (2019)*

**08/04**

**It is often surprisingly difficult to choose $\lambda$**

*Grünwald (2012)*
*Lyddon, Holmes, & Walker (2019)*
*Wu & Martin (2023)*

**25/03**

**Predictive/robustness gains vanish provably & very quickly for even moderate $\dfrac{n}{D}$**

*Medina, Olea, Rush, & Velez (2022)*
*McLatchie, Fong, Frazier, & Knoblauch (2024)*

# Basics: power posteriors

Power/Fractional/Cold Posteriors

$$\pi_n^{(\lambda)}(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta)^{\lambda} \cdot \pi(\theta) d\theta}$$

**Q: What else have we discovered?**

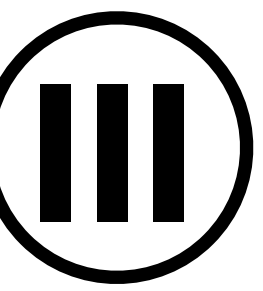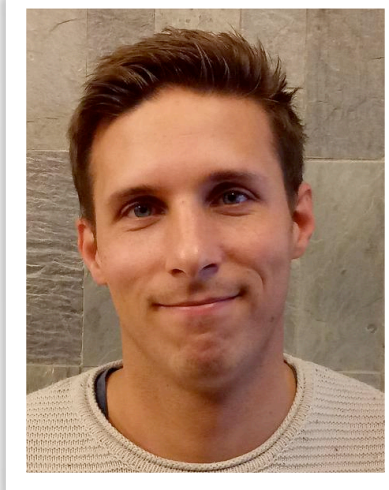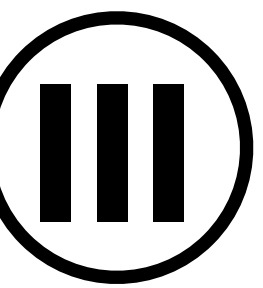**Power posterior & their variational approximations concentrate in situations where standard Bayes wouldn't**

*Bhattacharya, Pati, & Yang (2019)*
*Alquier & Ridgeway (2020)*
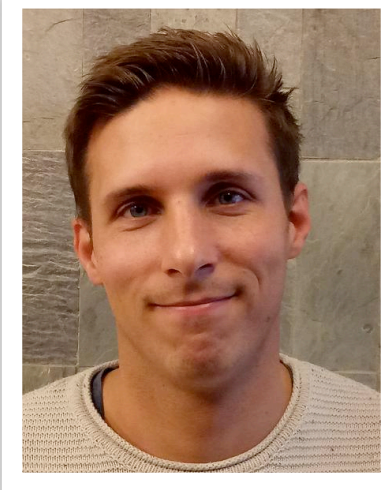*Yang, Pati, & Bhattacharya (2020)*

**Power posterior $\approx$ 'Coarsened' posterior (conditioning on a neighbourhood of observed data)**

*Miller & Dunson (2019)*

**It is often surprisingly difficult to choose $\lambda$**

*Grünwald (2012)*
*Lyddon, Holmes, & Walker (2019)*
*Wu & Martin (2023)*

**Predictive/robustness gains vanish provably & very quickly for even moderate $\frac{n}{D}$**

*Medina, Olea, Rush, & Velez (2022)*
*McLatchie, Fong, Frazier, & Knoblauch (2024)*

$\Longrightarrow$ **Motivates** $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \dfrac{\exp\{-\mathsf{L}(x_{1:n}, p_{\theta})\} \cdot \pi(\theta)}{\int \exp\{-\mathsf{L}(x_{1:n}, p_{\theta})\} \cdot \pi(\theta) d\theta}$

# Kalman Filter Example: generalised / Gibbs posteriors

Bayes' Posterior

$$\pi_n(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \cdot \pi(\theta)}{\int p(x_{1:n} \mid \theta) \cdot \pi(\theta) d\theta}$$

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$



Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch, Briol, & Murphy (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

$$\left( \textbf{NOT necessary for } \theta \textbf{ to come from a model } p_\theta \right)$$

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

## Q: How to think about this?

**Perspective 1: 'General Bayes Updates'**

**(Conditional) independence:**
$$p_\theta(x_{1:n}) = \prod_{i=1}^{n} p_\theta(x_i)$$

$$\Downarrow$$

$$\pi_n(\theta \mid x_{1:n}) \propto \pi_{n-1}(\theta \mid x_{1:(n-1)}) \cdot p_\theta(x_n)$$

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) \propto \pi_n^{\mathsf{L}}(\theta \mid x_{1:(n-1)}) \cdot \exp\{-\lambda \cdot \ell(x_n, p_\theta)\}$$

$$\Uparrow$$

**Summable Losses:**
$$\mathsf{L}(x_{1:n}, p_\theta) = \sum_{i=1}^{n} \ell(x_i, p_\theta)$$

*e.g., Bissiri, Holmes, & Walker (2016)*

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: How to think about this?**
**A: (1) General Bayes Updates**

**Perspective 1: 'General Bayes Updates'**

**(Conditional) independence:**

$$p_\theta(x_{1:n}) = \prod_{i=1}^{n} p_\theta(x_i)$$

$$\Downarrow$$

$$\pi_n(\theta \mid x_{1:n}) \propto \pi_{n-1}(\theta \mid x_{1:(n-1)}) \cdot p_\theta(x_n)$$

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) \propto \pi_n^{\mathsf{L}}(\theta \mid x_{1:(n-1)}) \cdot \exp\{-\lambda \cdot \ell(x_n, p_\theta)\}$$

$$\Updownarrow$$

**Summable Losses:**

$$\mathsf{L}(x_{1:n}, p_\theta) = \sum_{i=1}^{n} \ell(x_i, p_\theta)$$

*e.g., Bissiri, Holmes, & Walker (2016)*

**Main restriction of this interpretation**

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: How to think about this?**
**A: (1) General Bayes Updates**
    **(2) Optimisation-centric view**

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \underbrace{\mathbb{E}_{\theta \sim q}\left[\mathsf{L}(x_{1:n}, p_\theta)\right]}_{\text{Data-fitting}} + \underbrace{\frac{1}{\lambda} \mathrm{KL}(q, \pi)}_{\text{Prior regularisation}} \right\};$$

All probability distributions
over parameter space $\Theta$

*e.g., Knoblauch, Jewson, & Damoulas (2022)*

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: How to think about this?**
**A: (1) General Bayes Updates**
**    (2) Optimisation-centric view**

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q} \left[ \mathsf{L}(x_{1:n}, p_\theta) \right] + \frac{1}{\lambda} \mathrm{KL}(q, \pi) \right\};$$

All probability distributions over parameter space $\Theta$

Data-fitting

Prior regularisation

**Need not be summable**

*e.g., Knoblauch, Jewson, & Damoulas (2022)*

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

**Q: How to think about this?**
**A: (1) General Bayes Updates**
    **(2) Optimisation-centric view**

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \underbrace{\mathbb{E}_{\theta \sim q}\left[\mathsf{L}(x_{1:n}, p_\theta)\right]}_{\text{Data-fitting}} + \underbrace{\frac{1}{\lambda}\mathrm{KL}(q, \pi)}_{\substack{\text{Prior} \\ \text{regularisation}}} \right\};$$

All probability distributions
over parameter space Θ

**Role of PAC-Bayes:**
**what choices lead to what** $\}$ Chapter 3
**generalisation guarantees?**

*e.g., Knoblauch, Jewson, & Damoulas (2022)*

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$
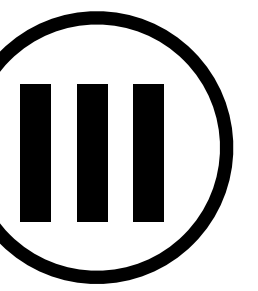
**Q: How to think about this?**
**A: (1) General Bayes Updates**
**(2) Optimisation-centric view**

**NOT asked by PAC-Bayes:**

When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?

How should we design $\mathsf{L}(x_{1:n}, p_\theta)$?

What happens asymptotically?

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q}\left[ \mathsf{L}(x_{1:n}, p_\theta) \right] + \frac{1}{\lambda} \mathrm{KL}\left(q, \pi\right) \right\};$$

Data-fitting

Prior regularisation

All probability distributions over parameter space $\Theta$

**Role of PAC-Bayes:**
**what choices lead to what**
**generalisation guarantees?**

Chapter 3

25/02

*e.g., Knoblauch, Jewson, & Damoulas (2022)*

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: When is** $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ **robust?**

**Setting:** for some small $\varepsilon \geq 0$,

**Data-generating probability distribution**

$\varepsilon$ **-contamination distribution**

$$q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$$

**Part of distribution our model captures**

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$
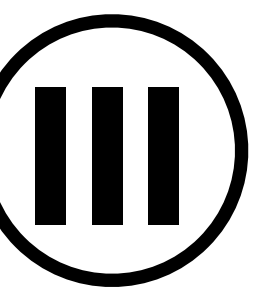
**Q: When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?**

**Setting:** for some small $\varepsilon \geq 0$,

**Data-generating probability distribution**

**$\varepsilon$-contamination distribution**

$$q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$$

**Part of distribution our model captures**

**What we want**

$q_0$

**contamination**

**Standard Bayes**

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

**Q: When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?**

**What we want:**
$$\begin{cases} x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) \\ \qquad\qquad\qquad\qquad \approx \\ z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \end{cases}$$

**Setting:** for some small $\varepsilon \geq 0$,

**Data-generating probability distribution**

**$\varepsilon$-contamination distribution**

$$q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$$

**Part of distribution our model captures**



**What we want**

$q_0$

**contamination**

**Standard Bayes**

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\llcorner}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: When is $\pi_n^{\llcorner}(\theta \mid x_{1:n})$ robust?**

**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$

$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\llcorner}(\theta \mid x_{1:n})$$

$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\llcorner}(\theta \mid z_{1:n})$$

**Robustness:** distance $\left\{ \pi_n^{\llcorner}(\theta \mid x_{1:n}), \pi_n^{\llcorner}(\theta \mid z_{1:n}) \right\} \leq \mathrm{constant}(c) \cdot \varepsilon$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

**Q: When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?**

**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$

$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$$

$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n})$$

**Robustness:** $\sup_{c \in \mathcal{S}} \left\{ \text{distance} \left\{ \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}), \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right\} \right\} \leq \text{constant}(\mathcal{S}) \cdot \varepsilon$

$$= \sup_{\theta \in \Theta} \left| \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) - \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right|$$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?**
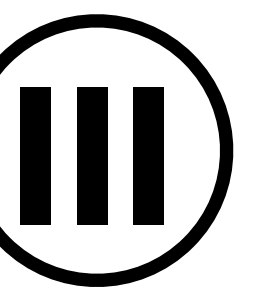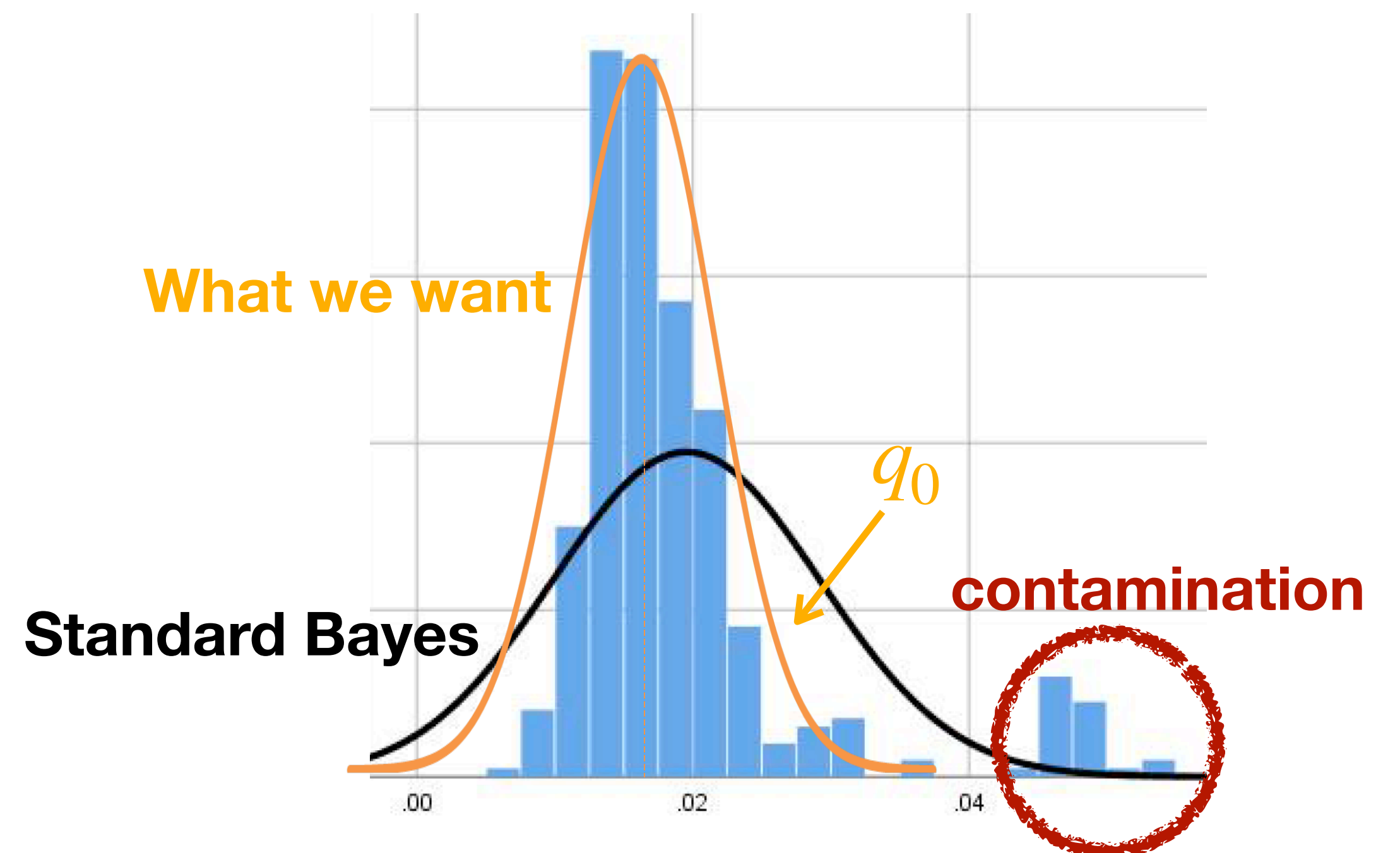
**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$

$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$$

$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n})$$

**Robustness:** $\sup\limits_{c \in \mathcal{S}} \left\{ \text{distance} \left\{ \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}), \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right\} \right\} \leq \text{constant}(\mathcal{S}) \cdot \varepsilon$

$$= \sup\limits_{\theta \in \Theta} \left| \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) - \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right|$$

**Key quantity:** $\dfrac{1}{\varepsilon} \left[ \mathsf{L}(p_\theta, x_{1:n}) - \mathsf{L}(p_\theta, z_{1:n}) \right]$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?**

**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$

$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$$

$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n})$$

**Robustness:** $\sup\limits_{c \in \mathcal{S}} \left\{ \text{distance} \left\{ \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}), \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right\} \right\} \leq \text{constant}(\mathcal{S}) \cdot \varepsilon$

$$= \sup\limits_{\theta \in \Theta} \left| \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) - \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right|$$

**Key quantity:** $\dfrac{\partial}{\partial \varepsilon} \mathsf{L}(p_\theta, x_{1:n}) \Big|_{\varepsilon=0} \approx \dfrac{1}{\varepsilon} \left[ \mathsf{L}(p_\theta, x_{1:n}) - \mathsf{L}(p_\theta, z_{1:n}) \right]$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?**

**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$

$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$$

$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n})$$

**Robustness:** $\sup_{c \in \mathcal{S}} \left\{ \text{distance} \left\{ \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}), \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right\} \right\} \leq \text{constant}(\mathcal{S}) \cdot \varepsilon$

$$= \sup_{\theta \in \Theta} \left| \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) - \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right| \qquad \textbf{\color{green}{Loss robust!}}$$

**Key quantity:** $\boxed{\sup_{\theta \in \Theta} \left| \frac{\partial}{\partial \varepsilon} \mathsf{L}(p_\theta, x_{1:n}) \Big|_{\varepsilon=0} \right| < \infty}$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

**Q: When is** $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ **robust?**

**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$
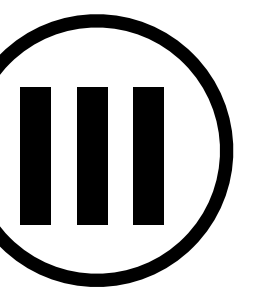
$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$$
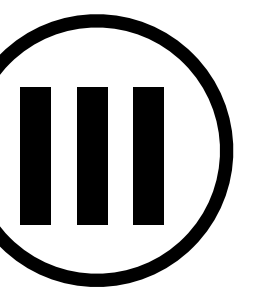
$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n})$$

**Theorem:** $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ is robust over all $c \in \mathcal{S}$ if $\mathsf{L}$ is.

**Robustness:** $\sup_{c \in \mathcal{S}} \left\{ \text{distance} \left\{ \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}), \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right\} \right\} \leq \text{constant}(\mathcal{S}) \cdot \varepsilon$

$$= \sup_{\theta \in \Theta} \left| \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) - \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right|$$

**Loss robust!**

**Key quantity:** $\sup_{\theta \in \Theta} \left| \frac{\partial}{\partial \varepsilon} \mathsf{L}(p_\theta, x_{1:n}) \Big|_{\varepsilon=0} \right| < \infty$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: When is $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ robust?**

**A: Whenever $\mathsf{L}(x_{1:n}, p_\theta)$ is!**

**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$

$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$$

$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n})$$

**Theorem:** $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ is robust over all $c \in \mathcal{S}$ if $\mathsf{L}$ is.

**Robustness:** $\sup_{c \in \mathcal{S}} \left\{ \text{distance} \left\{ \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}), \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right\} \right\} \leq \text{constant}(\mathcal{S}) \cdot \varepsilon$

$$= \sup_{\theta \in \Theta} \left| \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) - \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right|$$

**Loss robust!**

**Key quantity:** $\sup_{\theta \in \Theta} \left| \left. \frac{\partial}{\partial \varepsilon} \mathsf{L}(p_\theta, x_{1:n}) \right|_{\varepsilon=0} \right| < \infty$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: When is** $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ **robust?**

**A: Whenever** $\mathsf{L}(x_{1:n}, p_\theta)$ **is!**

**Setting:** $q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$

$$x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$$

$$z_{1:n} \sim q_0 \longrightarrow \pi_n^{\mathsf{L}}(\theta \mid z_{1:n})$$

**Theorem:** $\pi_n^{\mathsf{L}}(\theta \mid x_{1:n})$ is robust over all $c \in \mathcal{S}$ if $\mathsf{L}$ is.

**Robustness:** $\sup_{c \in \mathcal{S}} \left\{ \text{distance} \left\{ \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}), \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right\} \right\} \leq \text{constant}(\mathcal{S}) \cdot \varepsilon$

$$= \sup_{\theta \in \Theta} \left| \pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) - \pi_n^{\mathsf{L}}(\theta \mid z_{1:n}) \right|$$

**Generally untrue for log likelihoods!**

**Key quantity:** $\sup_{\theta \in \Theta} \left| \left. \frac{\partial}{\partial \varepsilon} \mathsf{L}(p_\theta, x_{1:n}) \right|_{\varepsilon=0} \right| < \infty$

Ghosh & Basu (2015); AISM
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

**Q: How to design robust** $\mathsf{L}(x_{1:n}, p_\theta)$ **?**

Hooker & Vidyashankar (2014); Test
Ghosh & Basu (2016); Statistica Sinica
Jewson, Smith, & Holmes (2018); Entropy
Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Cherieff-Abdellatif & Alquier (2020); AABI
Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & Knoblauch (2023); ICML
Altamirano, Briol, & Knoblauch (2024); ICML
Matsubara, Knoblauch, Briol, & Oates (2023); JASA
Knoblauch*, Frazier*, & Drovandi (2024); preprint
Pacchiardi, Dhoo, & Dutta (2024); EJS
…

**NOT robust to**
**model misspecification**

$$n \cdot \mathrm{KL}(q_\varepsilon, p(\,\cdot\mid\theta))$$

$x_i \sim q_\varepsilon$

$$\approx$$

**Standard Bayes**

$$\mathsf{L}(x_{1:n}, p_\theta) = \sum_{i=1}^{n} -\log p(x_i \mid \theta)$$

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: How to design robust** $\mathsf{L}(x_{1:n}, p_\theta)$ **?**

Hooker & Vidyashankar (2014); Test
Ghosh & Basu (2016); Statistica Sinica
Jewson, Smith, & Holmes (2018); Entropy
Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Cherieff-Abdellatif & Alquier (2020); AABI
Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & Knoblauch (2023); ICML
Altamirano, Briol, & Knoblauch (2024); ICML
Matsubara, Knoblauch, Briol, & Oates (2023); JASA
Knoblauch*, Frazier*, & Drovandi (2024); preprint
Pacchiardi, Dhoo, & Dutta (2024); EJS
…

**NOT robust to
model misspecification**

$$n \cdot \mathrm{KL}(q_\varepsilon, p(\cdot \mid \theta))$$

$x_i \sim q_\varepsilon$

$\approx$

**Standard Bayes**

$$\mathsf{L}(x_{1:n}, p_\theta) = \sum_{i=1}^{n} -\log p(x_i \mid \theta)$$

$$n \cdot \mathrm{D}(q_\varepsilon, p(\cdot \mid \theta))$$

**Robust discrepancy**

$$\mathrm{D}(q_\varepsilon, p(\cdot \mid \theta)) \quad \approx \quad \mathrm{D}(q_0, p(\cdot \mid \theta))$$

# Basics: generalised / Gibbs posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

Hooker & Vidyashankar (2014); Test
Ghosh & Basu (2016); Statistica Sinica
Jewson, Smith, & Holmes (2018); Entropy
Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Cherieff-Abdellatif & Alquier (2020); AABI
Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & Knoblauch (2023); ICML
Altamirano, Briol, & Knoblauch (2024); ICML
Matsubara, Knoblauch, Briol, & Oates (2023); JASA
Knoblauch*, Frazier*, & Drovandi (2024); preprint
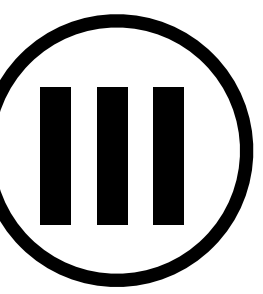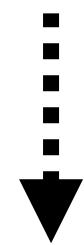Pacchiardi, Dhoo, & Dutta (2024); EJS
…

**Q: How to design robust $\mathsf{L}(x_{1:n}, p_\theta)$ ?**
**A: estimate robust divergence**

**NOT robust to model misspecification**

**Standard Bayes**

$$n \cdot \mathrm{KL}(q_\varepsilon, p(\cdot \mid \theta)) \qquad \overset{x_i \sim q_\varepsilon}{\approx} \qquad \mathsf{L}(x_{1:n}, p_\theta) = \sum_{i=1}^{n} -\log p(x_i \mid \theta)$$

$$n \cdot \mathrm{D}(q_\varepsilon, p(\cdot \mid \theta)) \qquad \overset{x_i \sim q_\varepsilon}{\approx} \qquad \mathsf{L}(x_{1:n}, p_\theta)$$

**Robust discrepancy**

**Robust loss**

$$\mathrm{D}(q_\varepsilon, p(\cdot \mid \theta)) \quad \approx \quad \mathrm{D}(q_0, p(\cdot \mid \theta)) \quad \cdots\!\!\rightarrow \quad \mathsf{L} \text{ is robust over all } c \in \mathcal{S}$$

# Basics: generalised / Gibbs posteriors

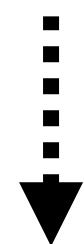Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

**Q: How to design robust** $\mathsf{L}(x_{1:n}, p_\theta)$ **?**
**A: estimate robust divergence**

Hooker & Vidyashankar (2014); Test
Ghosh & Basu (2016); Statistica Sinica
Jewson, Smith, & Holmes (2018); Entropy
Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Cherieff-Abdellatif & Alquier (2020); AABI
Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
Altamirano, Briol, & Knoblauch (2023); ICML
Altamirano, Briol, & Knoblauch (2024); ICML
Matsubara, Knoblauch, Briol, & Oates (2023); JASA
Knoblauch*, Frazier*, & Drovandi (2024); preprint
Pacchiardi, Dhoo, & Dutta (2024); EJS
...

**NOT robust to**
**model misspecification**

$$n \cdot \mathrm{KL}(q_\varepsilon, p(\cdot \mid \theta))$$

$x_i \sim q_\varepsilon$
$\approx$

**Standard Bayes**

$$\mathsf{L}(x_{1:n}, p_\theta) = \sum_{i=1}^n - \log p(x_i \mid \theta)$$

$$n \cdot \mathrm{D}(q_\varepsilon, p(\cdot \mid \theta))$$

$x_i \sim q_\varepsilon$
$\approx$

$$\mathsf{L}(x_{1:n}, p_\theta)$$

**Examples:**

**Robust discrepancy**

$$\mathrm{D}(q_\varepsilon, p(\cdot \mid \theta)) \quad \approx \quad \mathrm{D}(q_0, p(\cdot \mid \theta)) \quad \cdots\cdots\blacktriangleright \quad \mathsf{L} \text{ is robust over all } c \in \mathcal{S}$$

**Robust loss**

MMD
$\alpha/\beta/\gamma$-divergences
Stein discrepancies
...

# Basics: optimisation-centric posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$

## Q: Can we generalise this further?

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q}\left[\mathsf{L}(x_{1:n}, p_\theta)\right] + \frac{1}{\lambda}\mathrm{KL}(q, \pi) \right\}$$

All probability distributions over parameter space $\Theta$

Data-fitting

Prior regularisation

Variational Inference

Robustness to misspecified prior

$$q_n^*(\theta) = \arg \min_{q \in \mathbb{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + \frac{1}{\lambda}\mathrm{D}(q, \pi) \right\}$$

$\mathbb{Q} \subseteq \mathscr{P}(\Theta)$

*Knoblauch, Jewson, & Damoulas (2022)*

# Basics: optimisation-centric posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

## Q: Can we generalise this further?
## A: Yes! In many, many ways

Jankowiak, Pleiss, & Gardner (2020); ICML
Jankowiak, Pleiss, & Gardner (2020); UAI
Alquier (2020); ICML
Wild, Wu, & Sejdinovic (2022); NeurIPS
Javeed, Kouri, & Surowiec (2023); arXiv preprint
Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2023); NeurIPS
Shen, Knoblauch, Power, & Oates (2024); AISTATS
...

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg\min_{q \in \mathscr{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q}\left[ \mathsf{L}(x_{1:n}, p_\theta) \right] + \frac{1}{\lambda}\mathrm{KL}\big(q, \pi\big) \right\}$$

All probability distributions over parameter space Θ

Data-fitting

Prior regularisation

Variational Inference

Robustness to misspecified prior

$$q_n^*(\theta) = \arg\min_{q \in \mathscr{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + \frac{1}{\lambda}\mathrm{D}\big(q, \pi\big) \right\}$$

$$\mathscr{Q} \subseteq \mathscr{P}(\Theta)$$
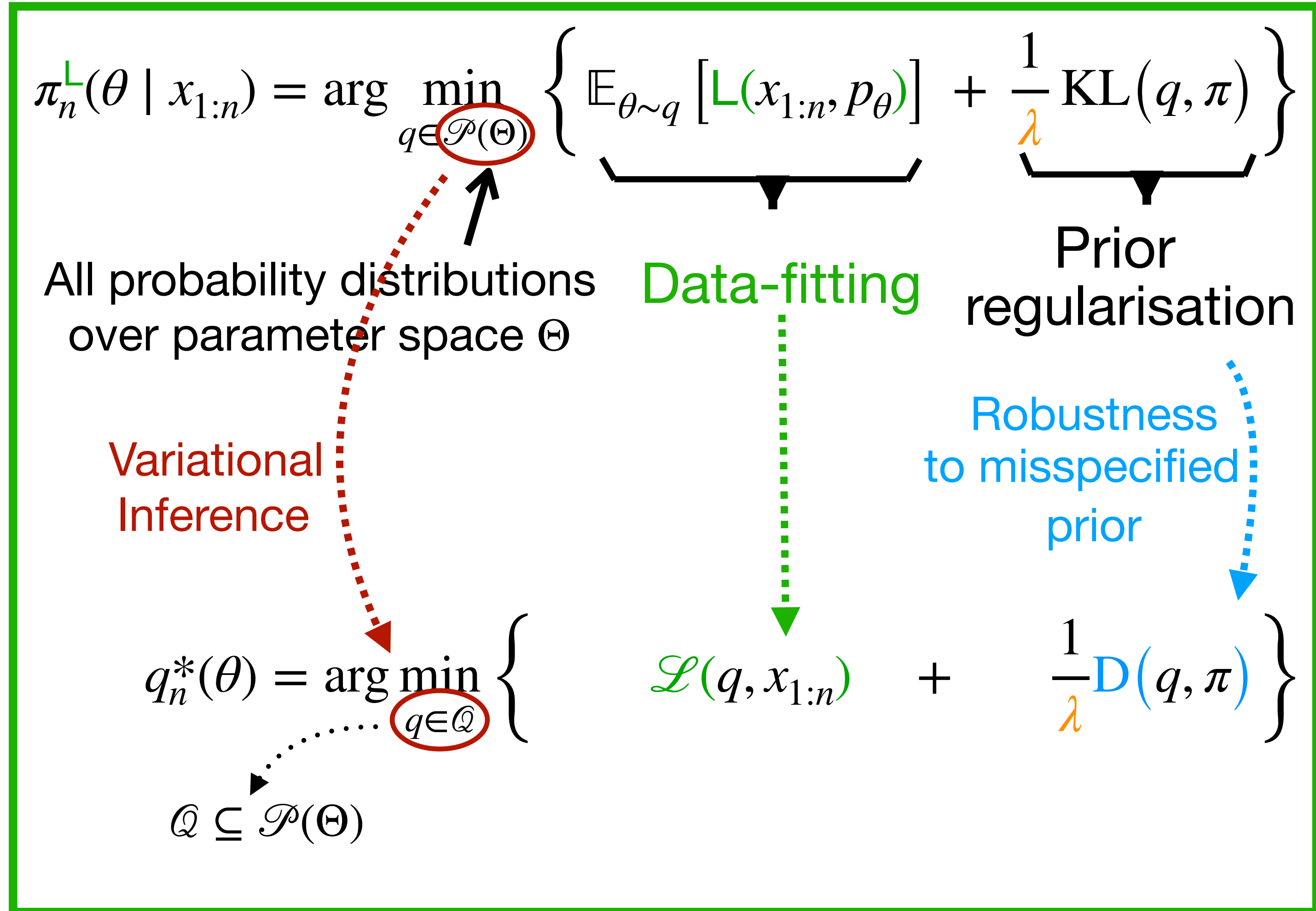
*Knoblauch, Jewson, & Damoulas (2022)*

# Basics: optimisation-centric posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta) d\theta}$$
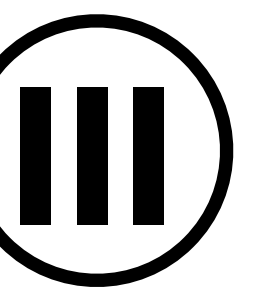
## Q: Can we generalise this further?
## A: Yes! In many, many ways

Jankowiak, Pleiss, & Gardner (2020); ICML
Jankowiak, Pleiss, & Gardner (2020); UAI
Alquier (2020); ICML
Wild, Wu, & Sejdinovic (2022); NeurIPS
Javeed, Kouri, & Surowiec (2023); arXiv preprint
Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2023); NeurIPS
Shen, Knoblauch, Power, & Oates (2024); AISTATS
...

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q}\left[\mathsf{L}(x_{1:n}, p_\theta)\right] + \frac{1}{\lambda}\mathrm{KL}(q, \pi) \right\}$$

All probability distributions over parameter space Θ

Data-fitting

Prior regularisation

**25/03**

Variational Inference

Robustness to misspecified prior

$$q_n^*(\theta) = \arg \min_{q \in \mathscr{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + \frac{1}{\lambda}\mathrm{D}(q, \pi) \right\}$$

$\mathscr{Q} \subseteq \mathscr{P}(\Theta)$

e.g. $\mathscr{L}\left( \int p_\theta(\cdot) \, \mathrm{d}q(\theta), x_{1:n} \right)$

*Knoblauch, Jewson, & Damoulas (2022)*

# Basics: optimisation-centric posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

**Q: Can we generalise this further?**
**A: Yes! In many, many ways**

Jankowiak, Pleiss, & Gardner (2020); ICML
Jankowiak, Pleiss, & Gardner (2020); UAI
Alquier (2020); ICML
Wild, Wu, & Sejdinovic (2022); NeurIPS
Javeed, Kouri, & Surowiec (2023); arXiv preprint
Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2023); NeurIPS
Shen, Knoblauch, Power, & Oates (2024); AISTATS
...

Computation:

$\mathcal{Q} = $ parametric $\implies$ generalised VI

$\mathcal{Q} = \mathscr{P}_2(\Theta) \implies$ Wasserstein Gradient Flow

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \underbrace{\mathbb{E}_{\theta \sim q}\left[\mathsf{L}(x_{1:n}, p_\theta)\right]}_{} + \underbrace{\frac{1}{\lambda}\mathrm{KL}(q, \pi)}_{} \right\}$$

All probability distributions over parameter space Θ

Data-fitting

Prior regularisation

**25/03**

Robustness to misspecified prior

Variational Inference

$$q_n^*(\theta) = \arg \min_{q \in \mathcal{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + \frac{1}{\lambda}\mathrm{D}(q, \pi) \right\}$$

$\mathcal{Q} \subseteq \mathscr{P}(\Theta)$      e.g. $\mathscr{L}\left(\int p_\theta(\cdot)\, dq(\theta), x_{1:n}\right)$

*Knoblauch, Jewson, & Damoulas (2022)*

# Basics: optimisation-centric posteriors

Gibbs/Generalised/Pseudo Posterior

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \frac{\exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)}{\int \exp\{-\lambda \cdot \mathsf{L}(x_{1:n}, p_\theta)\} \cdot \pi(\theta)d\theta}$$

**Q: Can we generalise this further?**
**A: Yes! In many, many ways**

Jankowiak, Pleiss, & Gardner (2020); ICML
Jankowiak, Pleiss, & Gardner (2020); UAI
Alquier (2020); ICML
Wild, Wu, & Sejdinovic (2022); NeurIPS
Javeed, Kouri, & Surowiec (2023); arXiv preprint
Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2023); NeurIPS
Shen, Knoblauch, Power, & Oates (2024); AISTATS
...

Computation:

$$\mathcal{Q} = \text{parametric} \implies \text{generalised VI}$$

$$\mathcal{Q} = \mathscr{P}_2(\Theta) \implies \text{Wasserstein Gradient Flow}$$

First Sampler of its kind!   + 'Morality Tale'

**Perspective 2: 'Optimisation-centric'**

$$\pi_n^{\mathsf{L}}(\theta \mid x_{1:n}) = \arg \min_{q \in \mathscr{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q}\left[ \mathsf{L}(x_{1:n}, p_\theta) \right] + \frac{1}{\lambda}\text{KL}\left(q, \pi\right) \right\}$$

All probability distributions over parameter space $\Theta$
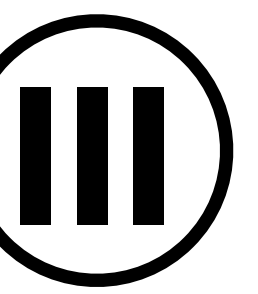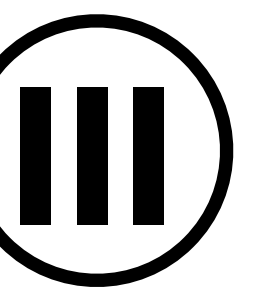
Data-fitting

Prior regularisation

**25/03**

Variational Inference

Robustness to misspecified prior

$$q_n^*(\theta) = \arg \min_{q \in \mathcal{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + \frac{1}{\lambda}\text{D}\left(q, \pi\right) \right\}$$

$$\mathcal{Q} \subseteq \mathscr{P}(\Theta) \qquad \text{e.g. } \mathscr{L}\left( \int p_\theta(\cdot) \, dq(\theta), x_{1:n} \right)$$

*Knoblauch, Jewson, & Damoulas (2022)*

# Morality Tale: Why Post-Bayesian thinking is needed

$$q_n^*(\theta) \quad = \arg\min_{q \in \mathcal{Q}} \left\{ \quad \mathscr{L}(q, x_{1:n}) \; + \; \frac{1}{\lambda} D(q, \pi) \right\}$$

Objective: $q \mapsto \mathbb{E}_{\theta \sim q} \left[ -\log p(x_{1:n} \mid \theta) \right] \; + \; \frac{1}{\lambda} KL(q, \pi)$

Target: **Cold Posterior** $(\lambda \gg 1)$
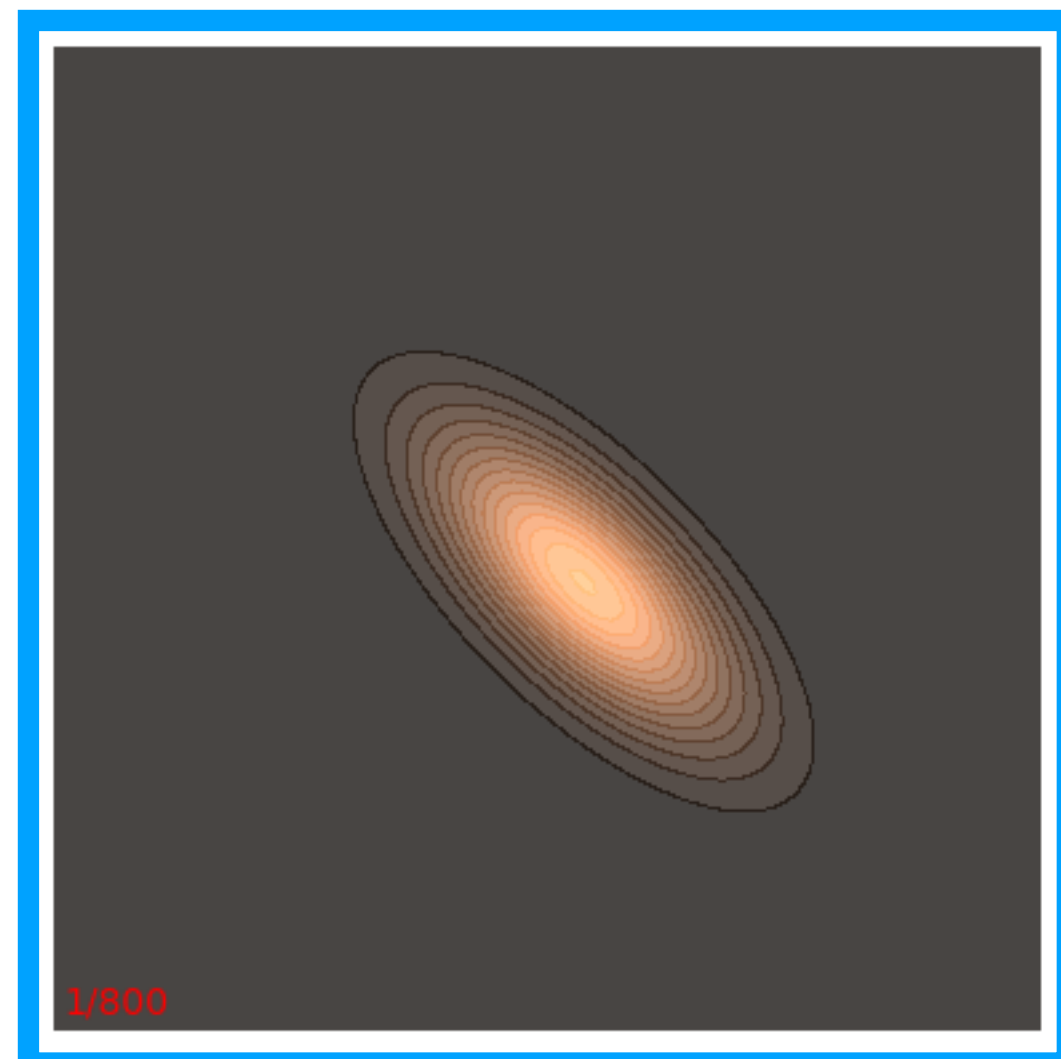
**/ Bayes Posterior** $(\lambda = 1)$

Wild, Ghalebikesabi, Sejdinovic, & **Knoblauch** (2023); NeurIPS Oral

# Morality Tale: Why Post-Bayesian thinking is needed

$$q_n^*(\theta) \quad = \arg\min_{q \in \mathcal{Q}} \left\{ \; \mathscr{L}(q, x_{1:n}) \; + \; \frac{1}{\lambda} D(q, \pi) \right\}$$

Objective: $\quad q \mapsto \mathbb{E}_{\theta \sim q} \left[ -\log p(x_{1:n} \mid \theta) \right] \; + \; \frac{1}{\lambda} KL(q, \pi)$

Target: **Cold Posterior** $(\lambda \gg 1)$
**/ Bayes Posterior** $(\lambda = 1)$

Wasserstein Gradient Flow = Langevin Diffusion



**Converges to well-defined density**

$$q_n^*(\theta) = \pi_n^{(\lambda)}(\theta \mid x_{1:n})$$

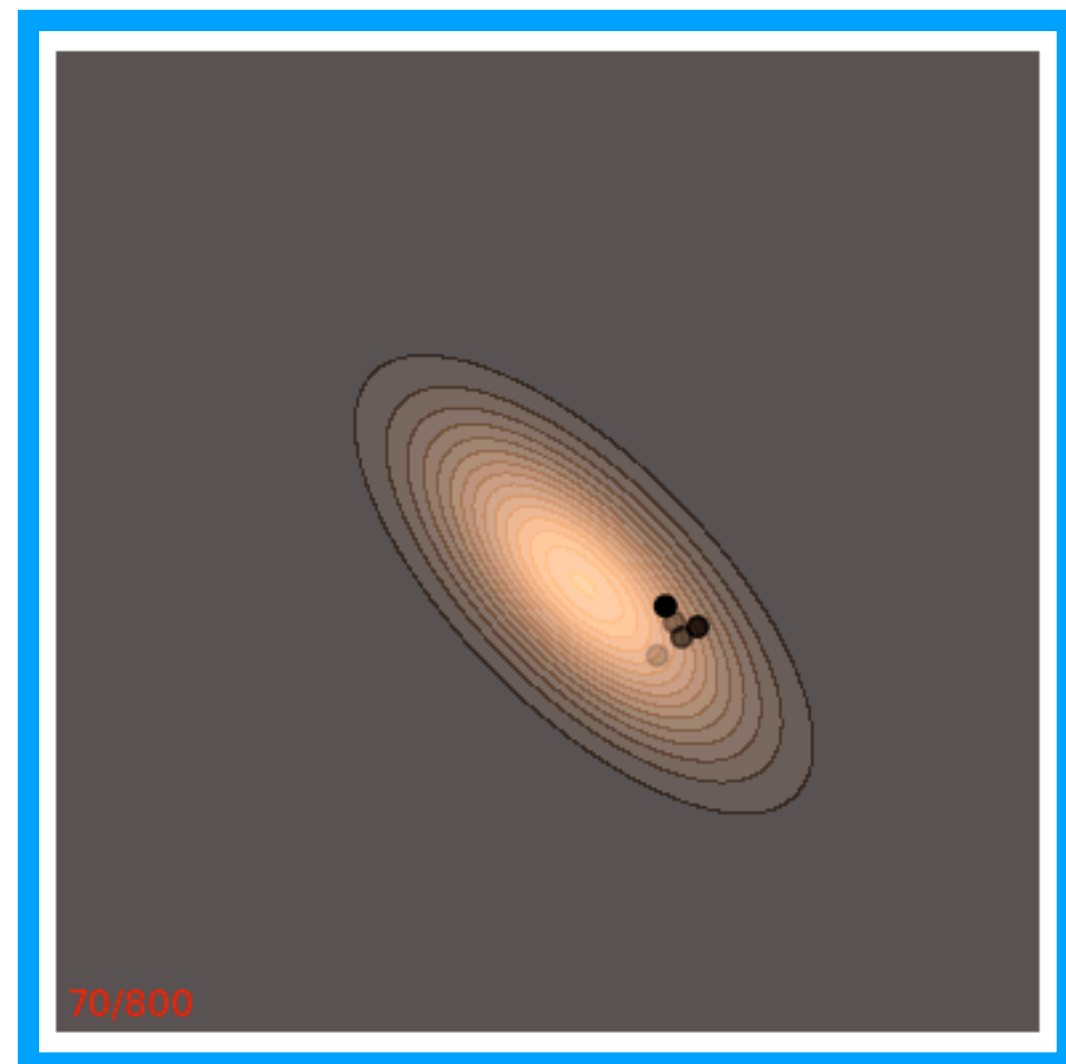Wild, Ghalebikesabi, Sejdinovic, & **Knoblauch** (2023); NeurIPS Oral

# Morality Tale: Why Post-Bayesian thinking is needed

$$q_n^*(\theta) = \arg\min_{q \in \mathcal{Q}} \left\{ \mathscr{L}(q, x_{1:n}) + \frac{1}{\lambda} \mathrm{D}(q, \pi) \right\}$$

Objective: $\quad q \mapsto \mathbb{E}_{\theta \sim q}\left[-\log p(x_{1:n} \mid \theta)\right] + \frac{1}{\lambda}\mathrm{KL}(q, \pi) \quad \xrightarrow{\lambda \to \infty} \quad q \mapsto \mathbb{E}_{\theta \sim q}\left[-\log p(x_{1:n} \mid \theta)\right]$

Target: **Cold Posterior** $(\lambda \gg 1)$ **Deep Ensemble (DE)** $(\lambda \to \infty)$
**/ Bayes Posterior** $(\lambda = 1)$
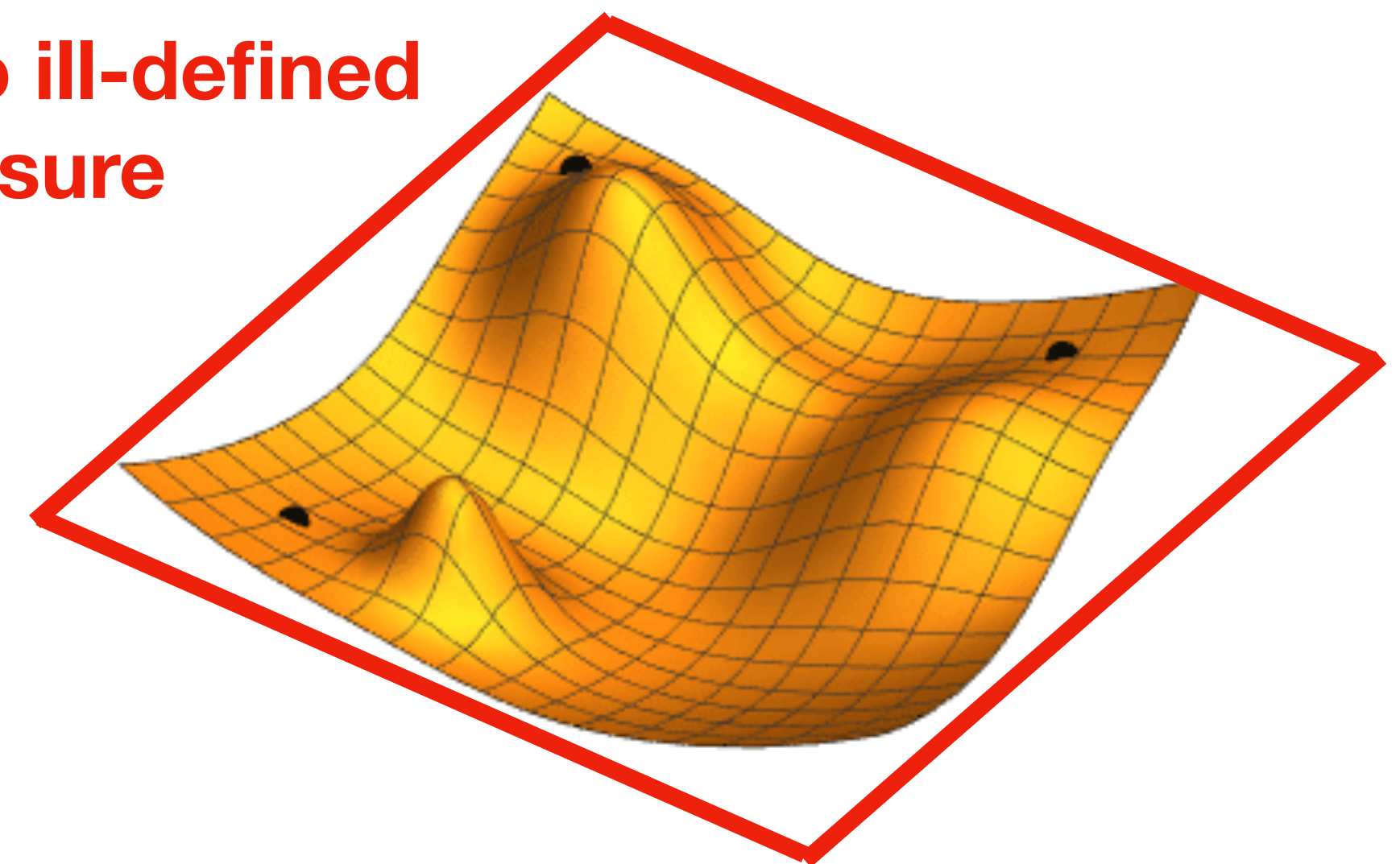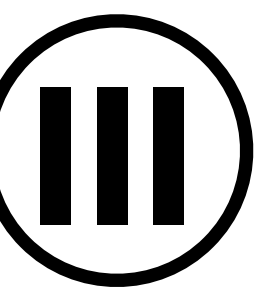
Wasserstein Gradient Flow = Langevin Diffusion

Wasserstein Gradient Flow = DE

**Converges to well-defined density**

$$q_n^*(\theta) = \pi_n^{(\lambda)}(\theta \mid x_{1:n})$$

**Converges to ill-defined discrete measure**



70/800

Wild, Ghalebikesabi, Sejdinovic, & **Knoblauch** (2023); NeurIPS Oral

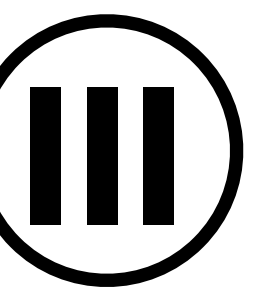# Morality Tale: Why Post-Bayesian thinking is needed

**Claim: 'Deep Ensembles = Bayesian Inference'**

*[…] Deep ensembles (Lakshminarayanan et al., 2017) are not a competing approach to Bayesian inference, but […] a compelling mechanism for Bayesian marginalization.*

**Published 2020 @ NeurIPS**
**(cited ≈ 800 times according to Google scholar)**

# Morality Tale: Why Post-Bayesian thinking is needed
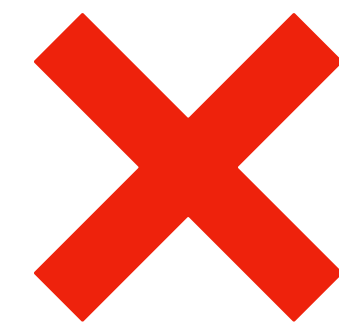
**Claim: 'Deep Ensembles = Bayesian Inference'**

*[…] Deep ensembles (Lakshminarayanan et al., 2017) are not a competing approach to Bayesian inference, but […] a compelling mechanism for Bayesian marginalization.*
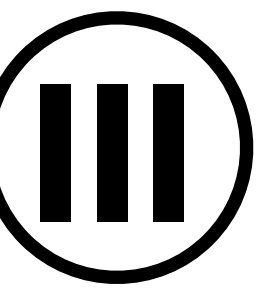
**Published 2020 @ NeurIPS
(cited ≈ 800 times according to Google scholar)**

**Unfortunately, this is not correct.**
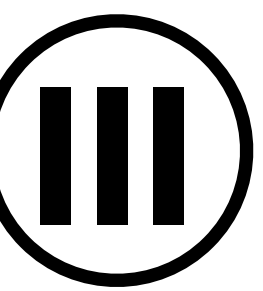
# Morality Tale: Why Post-Bayesian thinking is needed

**I.** In practice, orthodox Bayesianism has already been abandoned
(Bayes posterior: <span style="color:red">prior regulariser, densities</span> ;
Deep Ensembles: <span style="color:red">no prior regulariser, discrete measures</span>)
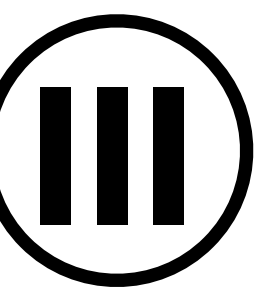
# Morality Tale: Why Post-Bayesian thinking is needed

**I.** In practice, orthodox Bayesianism has already been abandoned
(Bayes posterior: <span style="color:red">prior regulariser, densities</span> ;
Deep Ensembles: <span style="color:red">no prior regulariser, discrete measures</span>)

**II.** But practicioners often don't realise this / pay attention to the ramifications, which in turns leads to incorrect claims and conclusions.
('<span style="color:red">Deep Ensembles are Bayesian</span>')

# Morality Tale: Why Post-Bayesian thinking is needed

**I.** In practice, orthodox Bayesianism has already been abandoned
   (Bayes posterior: prior regulariser, densities ;
   Deep Ensembles: no prior regulariser, discrete measures)

**II.** But practicioners often don't realise this / pay attention to the ramifications, which in turns leads to incorrect claims and conclusions.
   ('Deep Ensembles are Bayesian')

**III.** It's on us to help them! We need to:
   1. acknowledge that post-Bayesian methods are already in use; and
   2. develop clear formalisms and design principles for them.

   Ideas either adapt, or die.