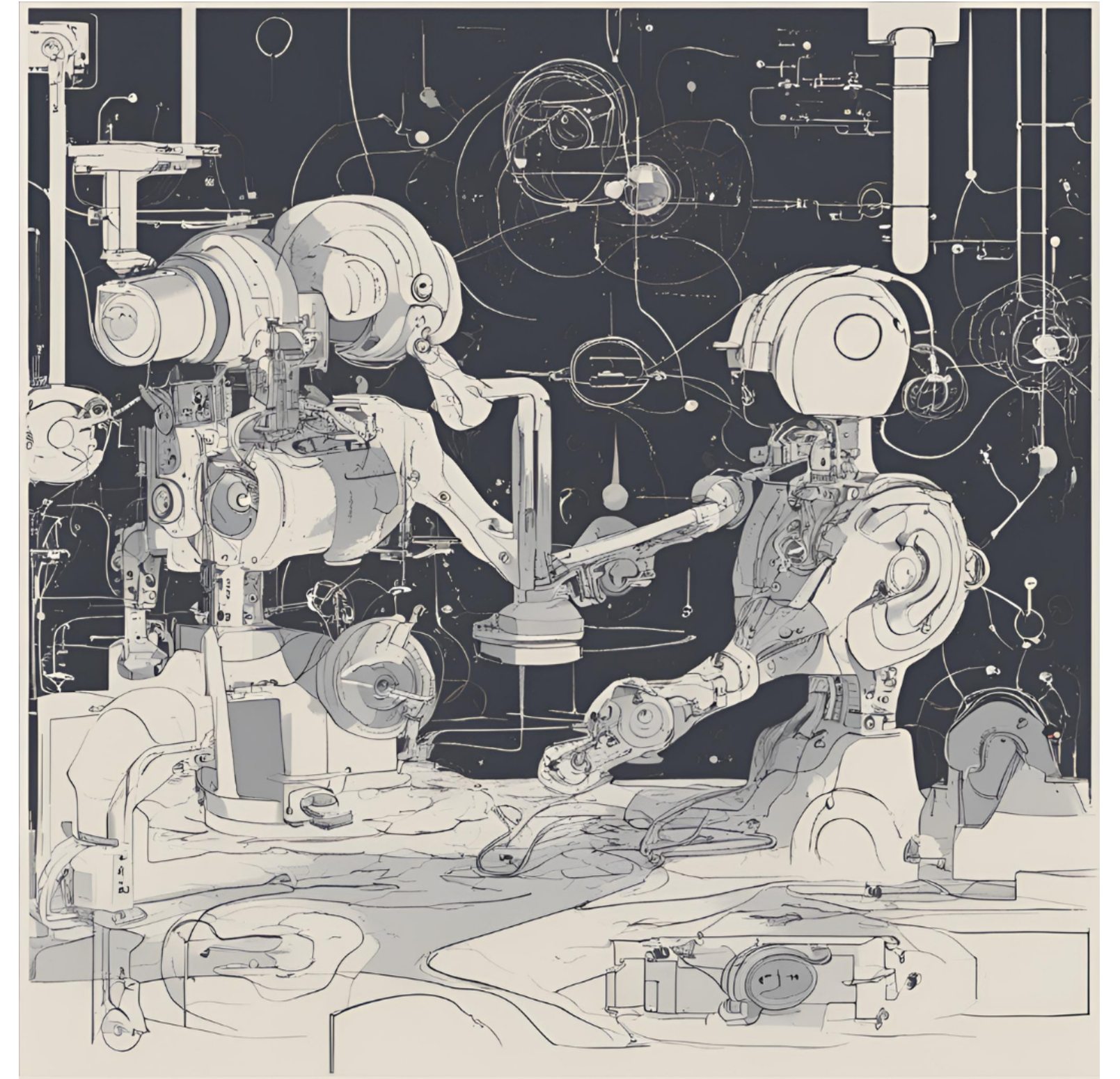


Post-Bayesian Machine Learning

Jeremias Knoblauch; Lecturer & EPSRC Fellow @ UCL Stats

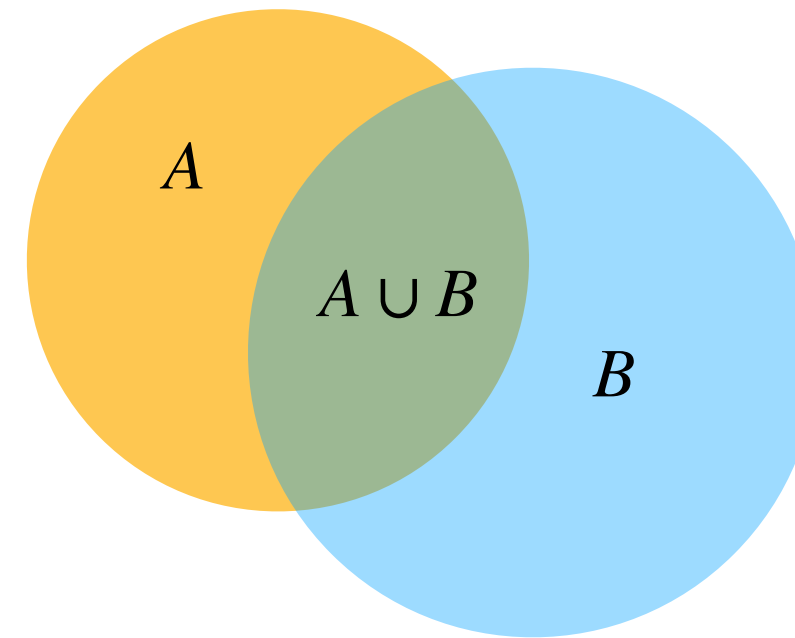


25/07/24

Preamble: Bayesian ML

Bayes' Theorem: Inversion of conditionals

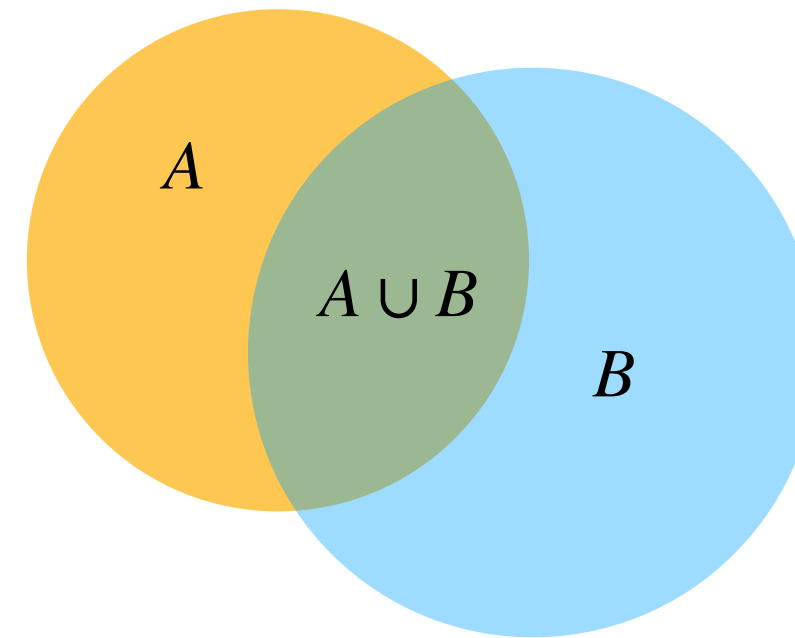
$$P(A | B) = \frac{P(B | A) \cdot P(A)}{P(B)}$$



Preamble: Bayesian ML

Bayes' Theorem: Inversion of conditionals

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



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

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

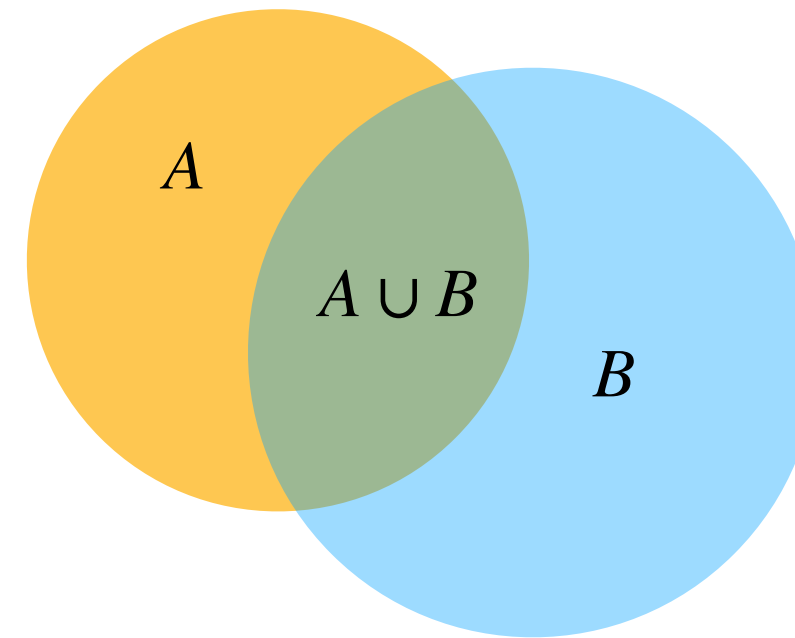
$$\pi_n(\theta | x_{1:n}) = \frac{p(x_{1:n} | \theta) \cdot \pi(\theta)}{\int p(x_{1:n} | \theta) \cdot \pi(\theta) d\theta}$$

(Bayes) Posterior

Preamble: Bayesian ML

Bayes' Theorem: Inversion of conditionals

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



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

$$p(x_{1:n} | \theta)$$

Prior probability:
 $\theta \in \Theta$

$$\pi(\theta)$$

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

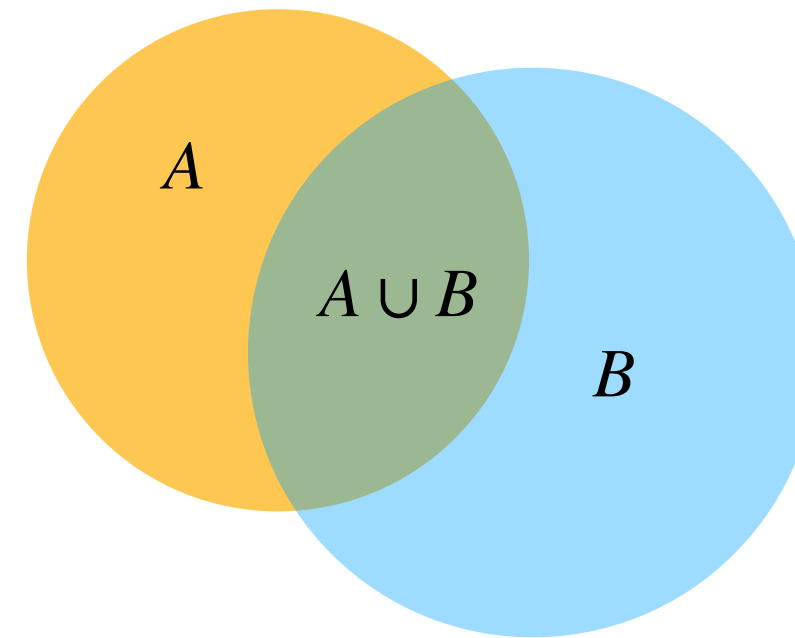
(Bayes) Posterior

- ⊕ Averages models (instead of picking only one)
- ⊕ Quantifies uncertainty about θ via $\pi_n(\theta | x_{1:n})$
- ⊕ Inclusion of domain expertise via prior π

Preamble: Bayesian ML

Bayes' Theorem: Inversion of conditionals

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

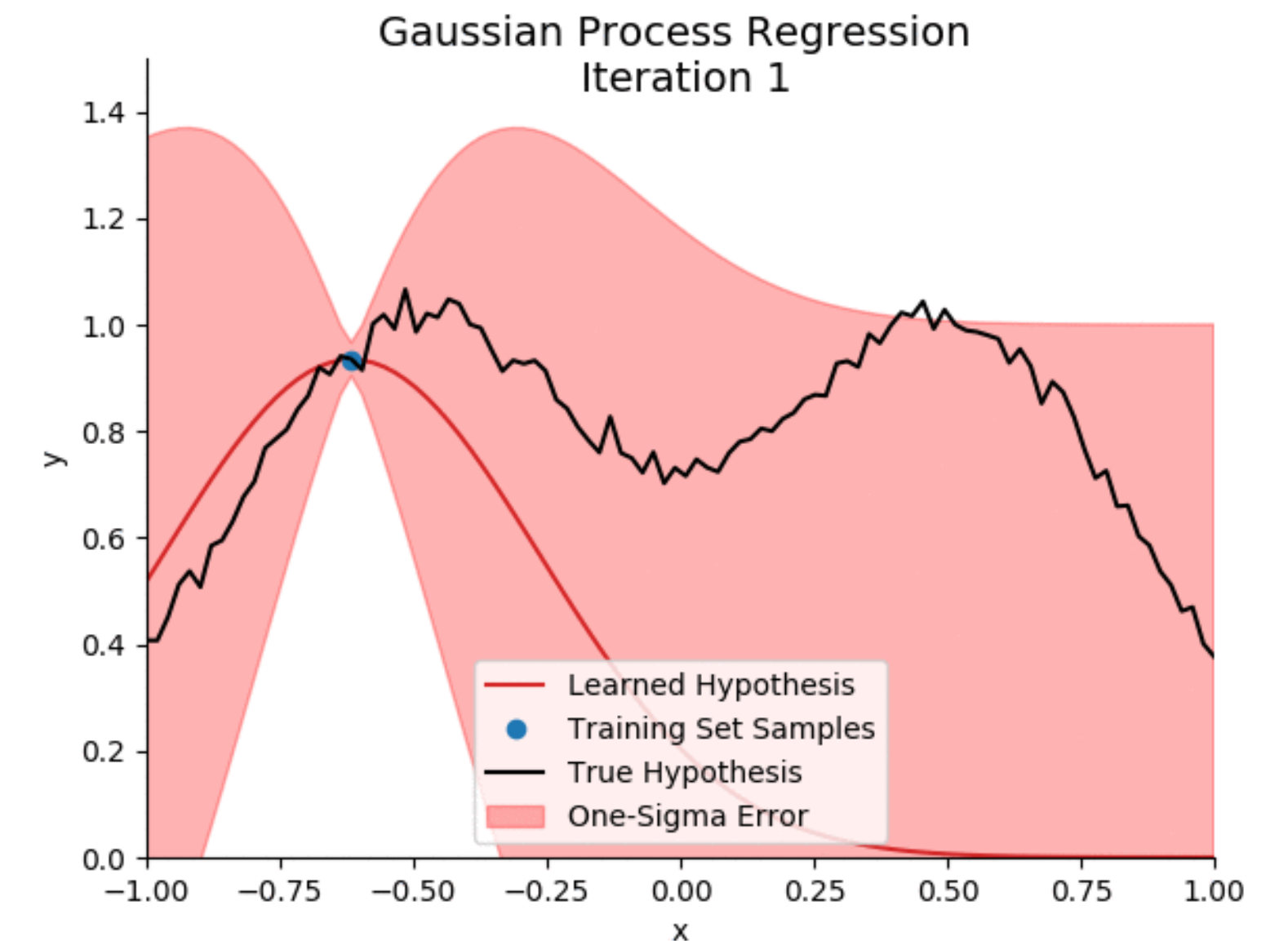


Data model:
 $x_{1:n} \in \mathcal{X}^n$
 $p(x_{1:n} | \theta)$

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

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

(Bayes) Posterior



- ⊕ Averages models (instead of picking only one)
- ⊕ Quantifies uncertainty about θ via $\pi_n(\theta | x_{1:n})$
- ⊕ Inclusion of domain expertise via prior π

Preamble: Bayesian ML

Mathematical
Foundations

1764—1786: **Bayes' Theorem**

1930: **DeFinetti's Representation Theorem**

1950s/60s: **Savage Axioms, Birnbaum's Likelihood Principle, ...**

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

Preamble: Bayesian ML

Mathematical Foundations

1764—1786: **Bayes' Theorem**

1930: **DeFinetti's Representation Theorem**

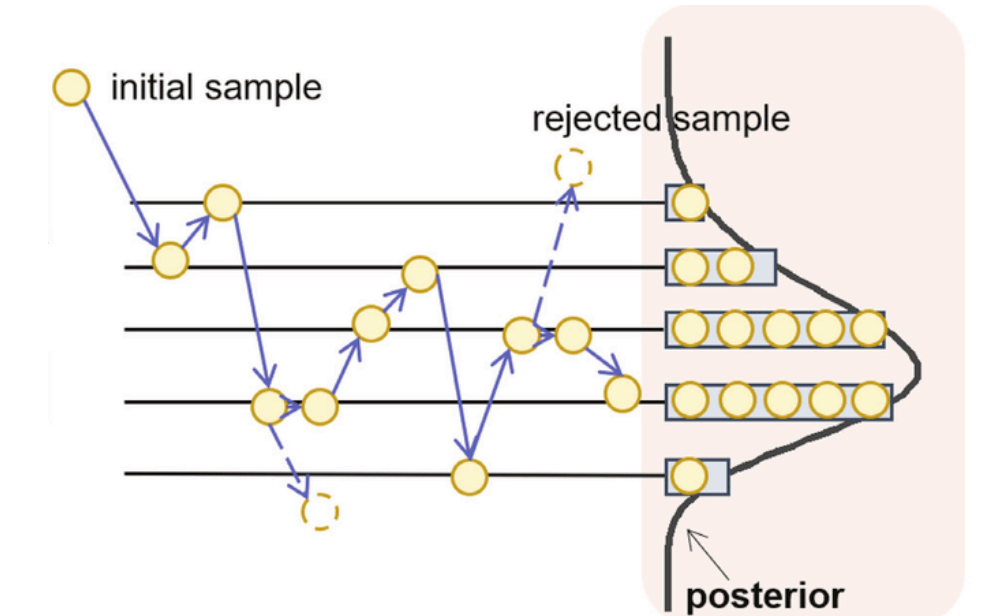
1950s/60s: **Savage Axioms, Birnbaum's Likelihood Principle, ...**

Computation

1960s/70s: **Computational Principles** for Bayesian computation

1980/90s onwards: computation becomes **feasible**

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



Preamble: Bayesian ML

Mathematical Foundations

1764—1786: **Bayes' Theorem**

1930: **DeFinetti's Representation Theorem**

1950s/60s: **Savage Axioms, Birnbaum's Likelihood Principle, ...**

Computation

1960s/70s: **Computational Principles** for Bayesian computation

1980/90s onwards: computation becomes **feasible**

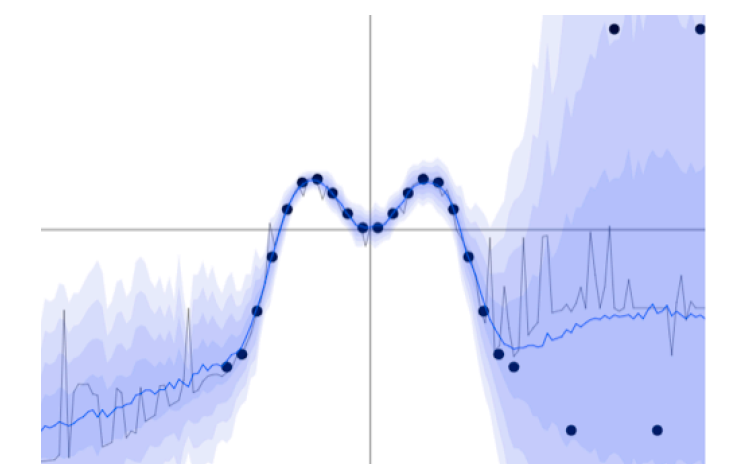
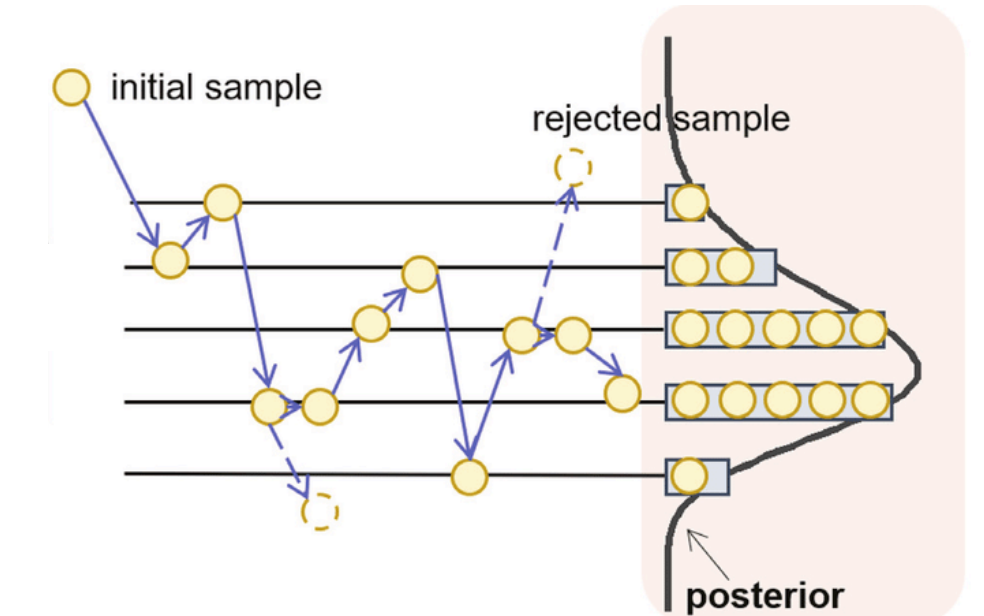
Bayesian ML

1990s: **ML** becomes **data-driven**

Mid 1990s—mid 2000s: **Bayesian ML** emerges

Classic perspective on Bayesian ML

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



Preamble: Bayesian ML

Mathematical Foundations

1764—1786: **Bayes' Theorem**

1930: **DeFinetti's Representation Theorem**

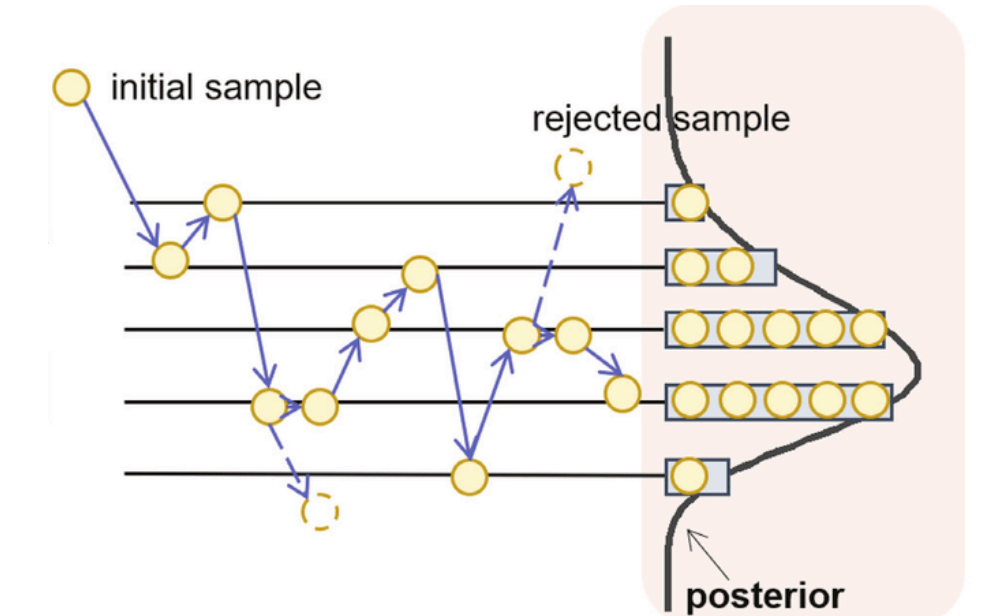
1950s/60s: **Savage Axioms, Birnbaum's Likelihood Principle, ...**

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

Computation

1960s/70s: **Computational Principles** for Bayesian computation

1980/90s onwards: computation becomes **feasible**



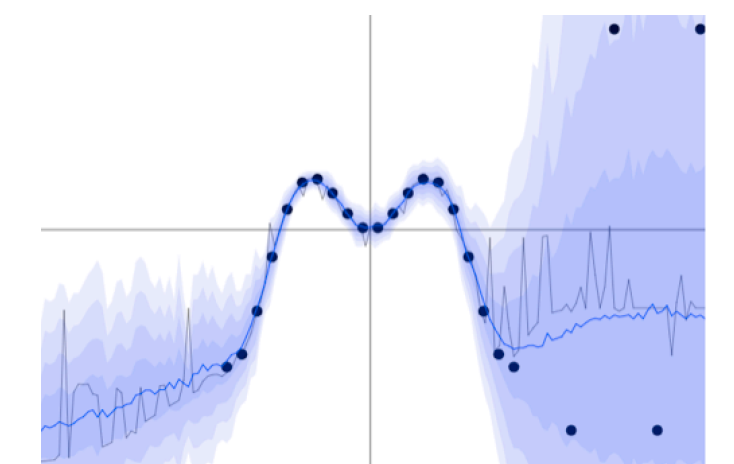
Bayesian ML

1990s: **ML** becomes **data-driven**

Classic perspective

Mid 1990s—mid 2000s: **Bayesian ML** emerges **on Bayesian ML**

2000s—present: **Adaptations of Bayesian principles** for ML/DL



Post-Bayesian ML and Today's Talk

Today's Talk in a Nutshell



Pitch for ML via Bayes' Rule:

- ⊕ Uncertainty Quantification
- ⊕ Inclusion of Prior Knowledge
- ⊕ Model averaging

Today's Talk in a Nutshell



Pitch for ML via Bayes' Rule:

- ⊕ Uncertainty Quantification
- ⊕ Inclusion of Prior Knowledge
- ⊕ Model averaging

Problem:

- ⊖ ML violates underlying assumptions
- ⇒ Unreliable & not robust

Today's Talk in a Nutshell

Fix: Post-Bayesian ML

Algorithms with 'Bayesian characteristics' that are **not using Bayes' Rule**



Pitch for ML via Bayes' Rule:

- ⊕ Uncertainty Quantification
- ⊕ Inclusion of Prior Knowledge
- ⊕ Model averaging

Problem:

- ⊖ ML violates underlying assumptions
- ⇒ Unreliable & not robust

Foundations of Bayesian ML

Reviewer 2: *“The Bayes posterior is principled, and grounded in foundational work of DeFinetti, Savage, Birnbaum, and others. This work gives up these guarantees [...] without clear gain in return.”*

Mathematical Foundations

1764—1786: Bayes' Theorem

1930: DeFinetti's Representation Theorem

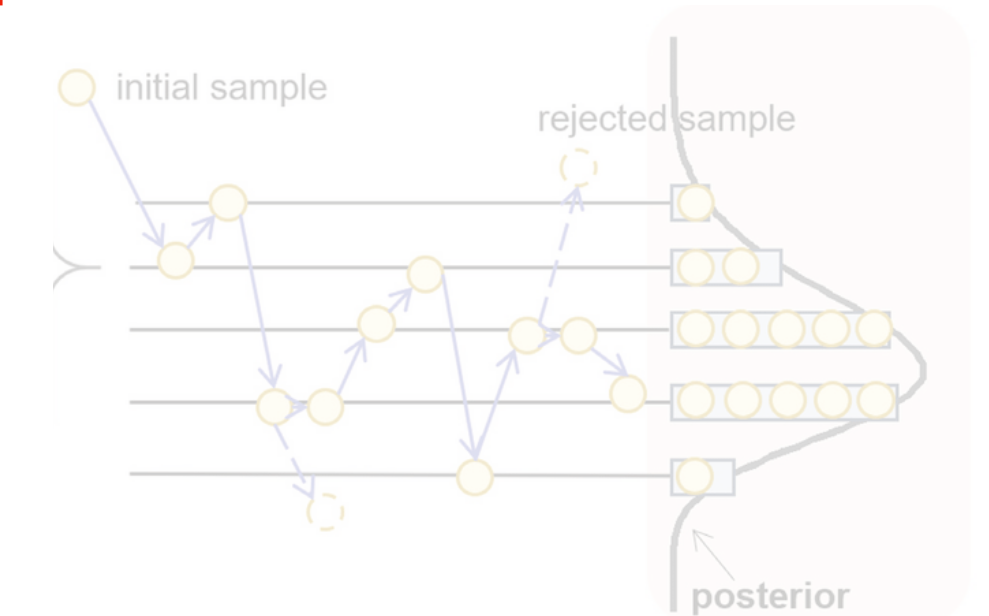
1950s/60s: Savage Axioms, Birnbaum's Likelihood Principle, ...

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

Computation

1960s/70s: Computational Principles (Markov chain Monte Carlo)

1980s onwards: computation becomes feasible



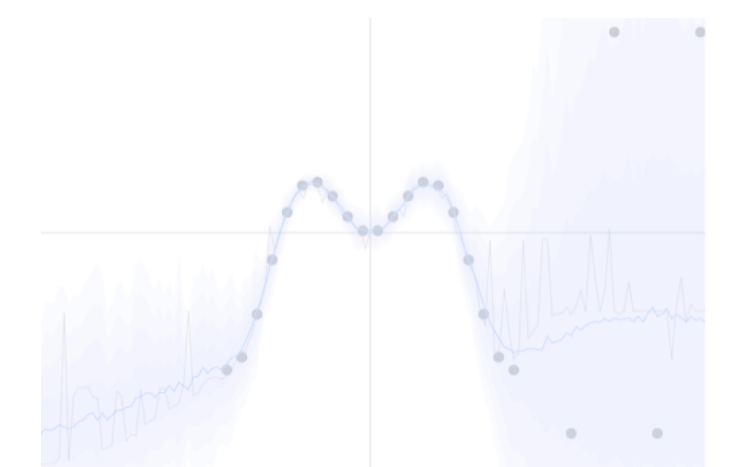
Bayesian ML

1990s: ML becomes data-driven

Mid 1990s—mid 2000s: Bayesian ML emerges

2000s—present: Adaptations of Bayesian principles for ML/DL

Classic perspective on Bayesian ML



Post-Bayesian ML and Today's Talk

Assumptions needed for Foundations

(A1)

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

Θ = Only relevant State of the world

Assumptions needed for Foundations

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

Θ = Only relevant State of the world

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

How rational decision-makers choose the prior

Assumptions needed for Foundations

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

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

Θ = 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) $\pi_n(\theta | x_{1:n})$ computable in practice

Guarantees real-world relevance

Assumptions needed for Foundations

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

(A1)

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

Θ = 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)

$\pi_n(\theta | x_{1:n})$ computable in practice

Guarantees real-world relevance

FRAGILE

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science

Expert with
research question

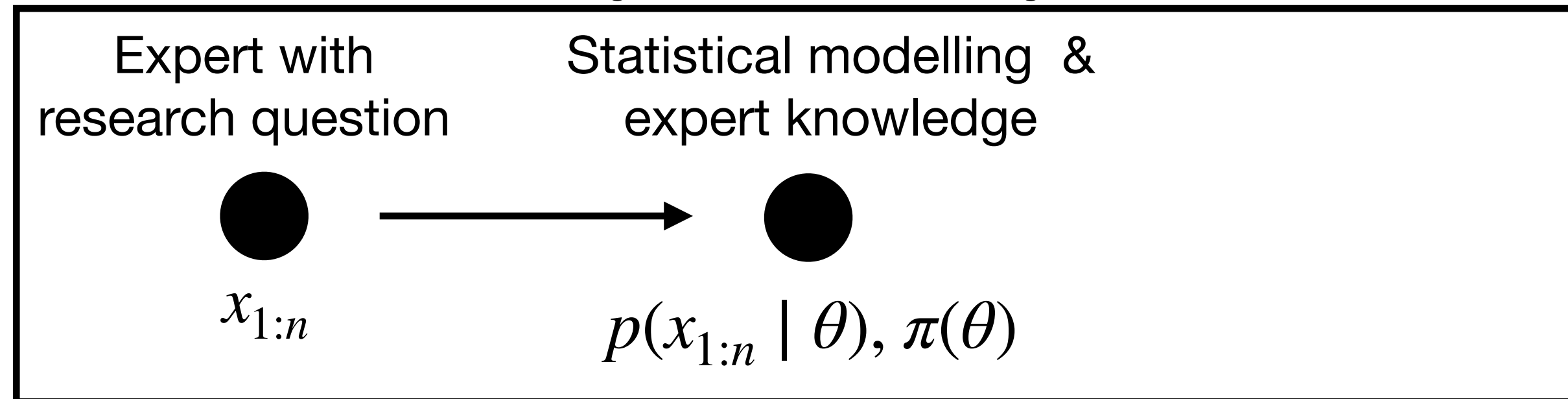


$x_{1:n}$

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

Case Study: Regression with Boston Housing Data

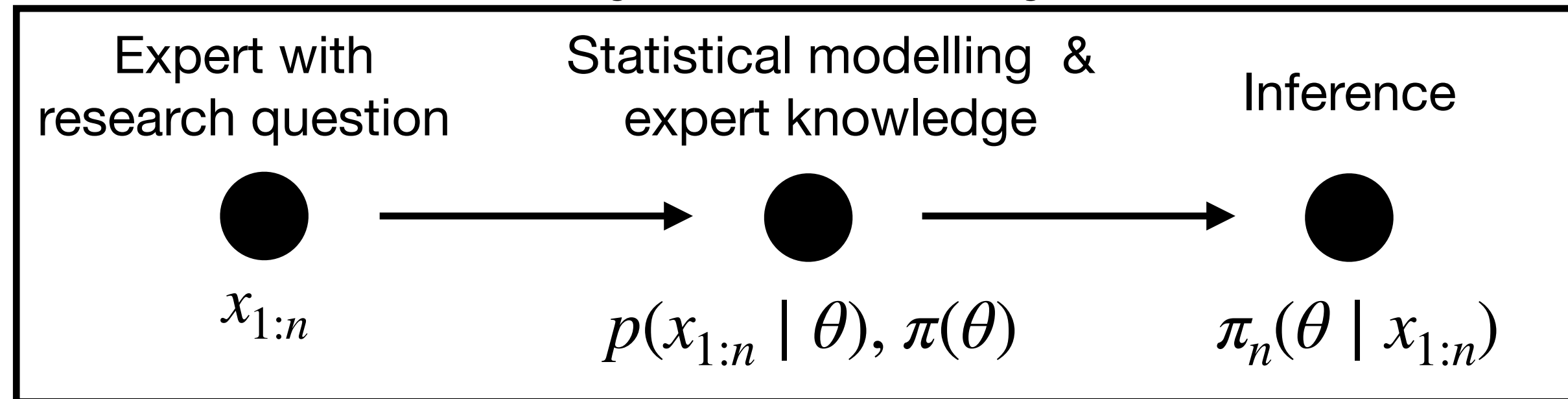
Traditional Bayesian analysis in science



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

Case Study: Regression with Boston Housing Data

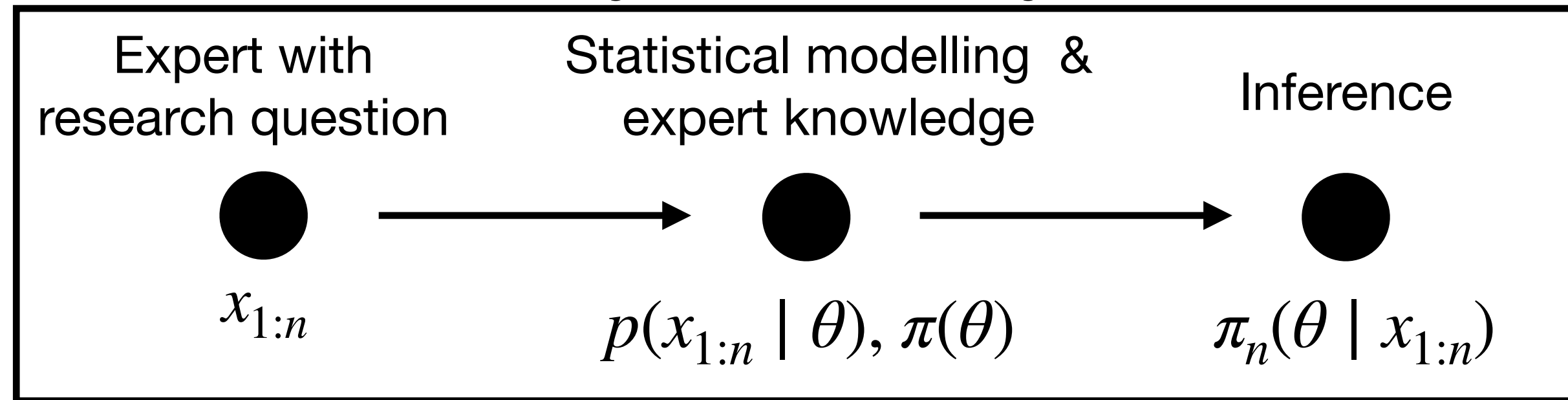
Traditional Bayesian analysis in science



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

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



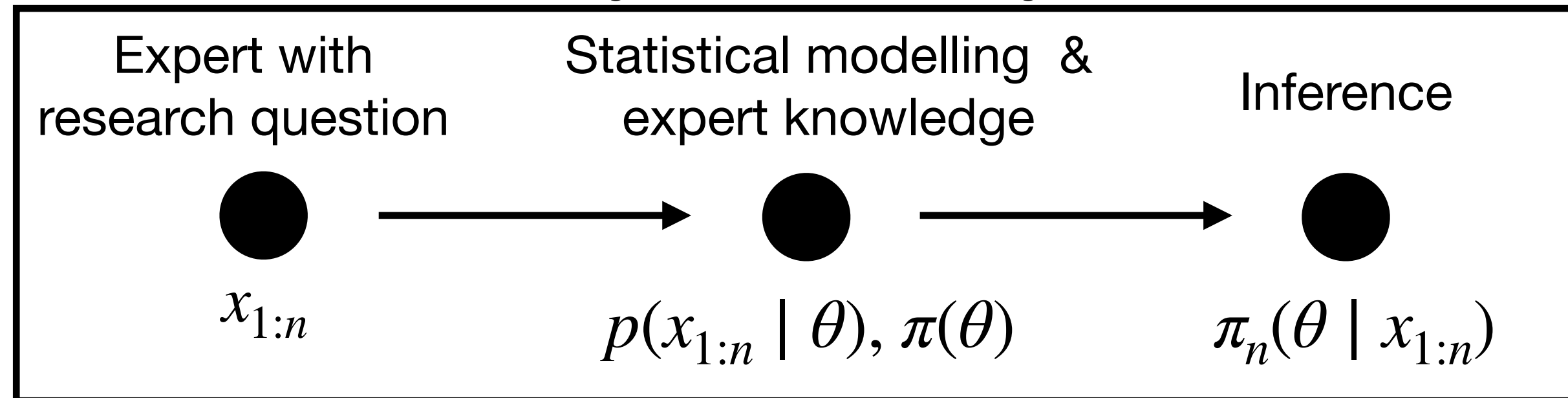
Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

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

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

willingness to pay \uparrow p_j $\log(x_{j,i})$ \uparrow pollutants \uparrow c_j $\log(x_{j,i})$ \uparrow rooms, sqm, ... \uparrow measurement error \uparrow ε_i

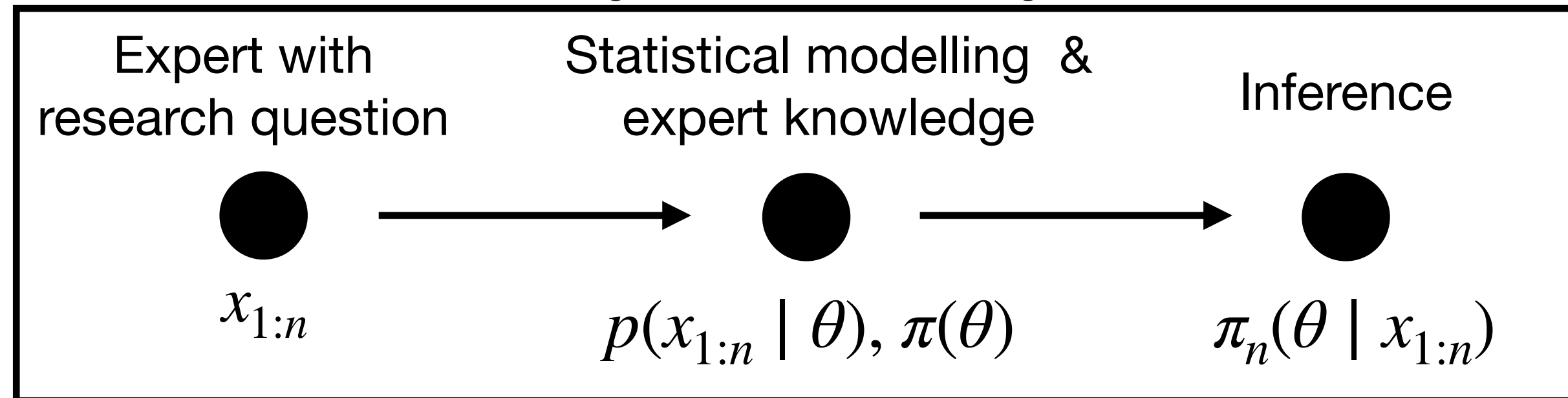
$\theta = (c_0, c_2, \dots, c_{J_1}, p_1, p_2, \dots, p_{J_2})^\top$

parameters of interest incidental parameters

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

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

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

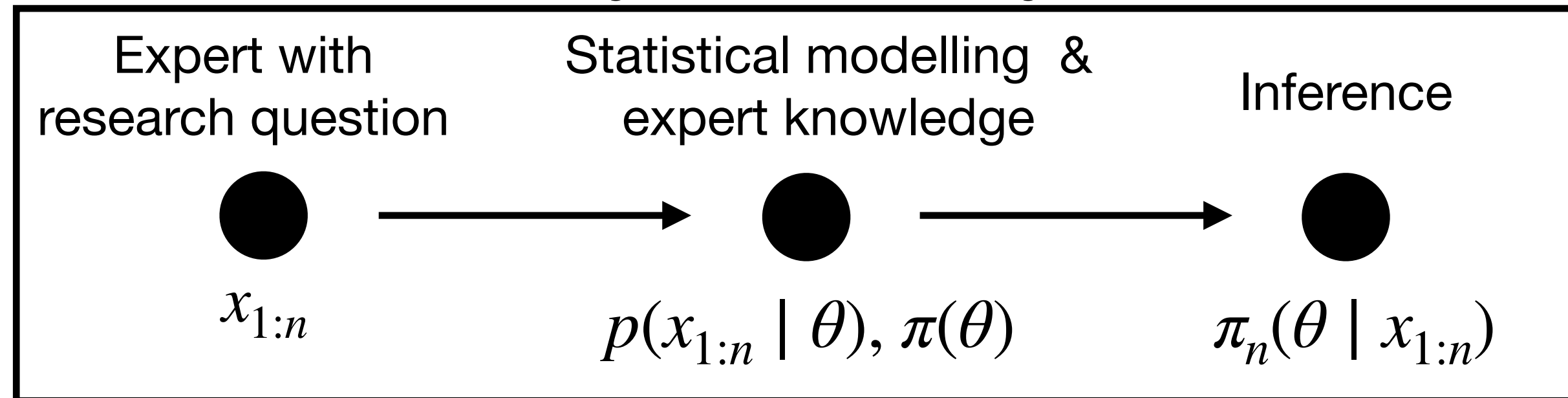
$\theta = (c_0, c_2, \dots, c_{J_1}, p_1, p_2, \dots, p_{J_2})^\top$
 $\pi(\theta) \sim$ hand-crafted by experts

(A2) ✓

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

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

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

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

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

$\pi_n(\theta | x_{1:n}) \longrightarrow$ computed exactly

(A2) ✓

(A3) ✓

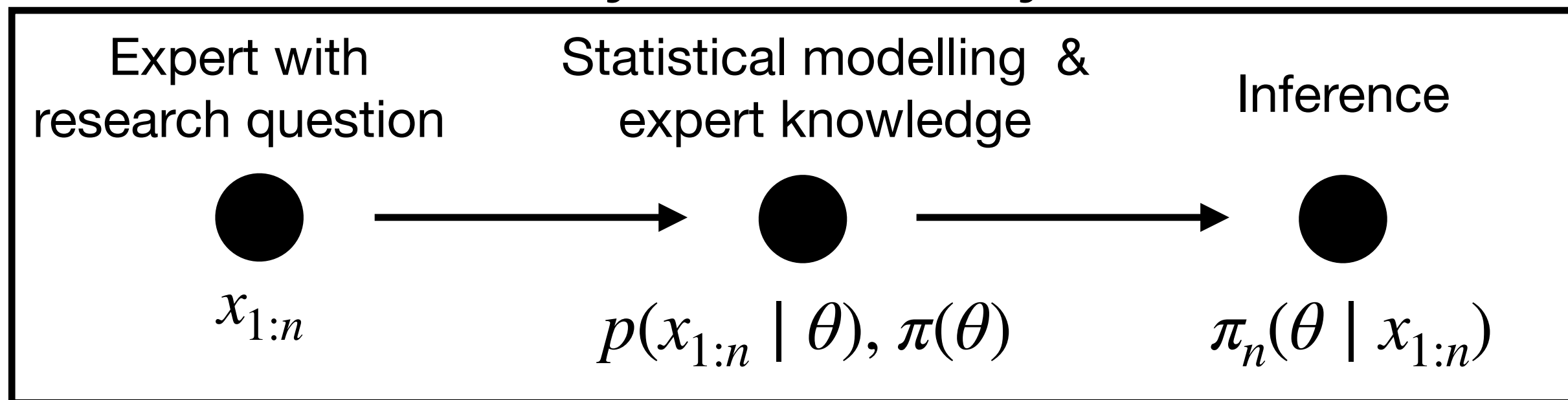
(A1) model well-specified

(A2) prior well-specified

(A3) computationally feasible

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Modern Bayesian ML



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

willingness to pay \uparrow p_j \uparrow pollutants \uparrow rooms, sqm, ... \uparrow measurement error \uparrow ε_i

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

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

$\pi_n(\theta | x_{1:n}) \rightarrow$ computed exactly

(A2) ✓

(A3) ✓

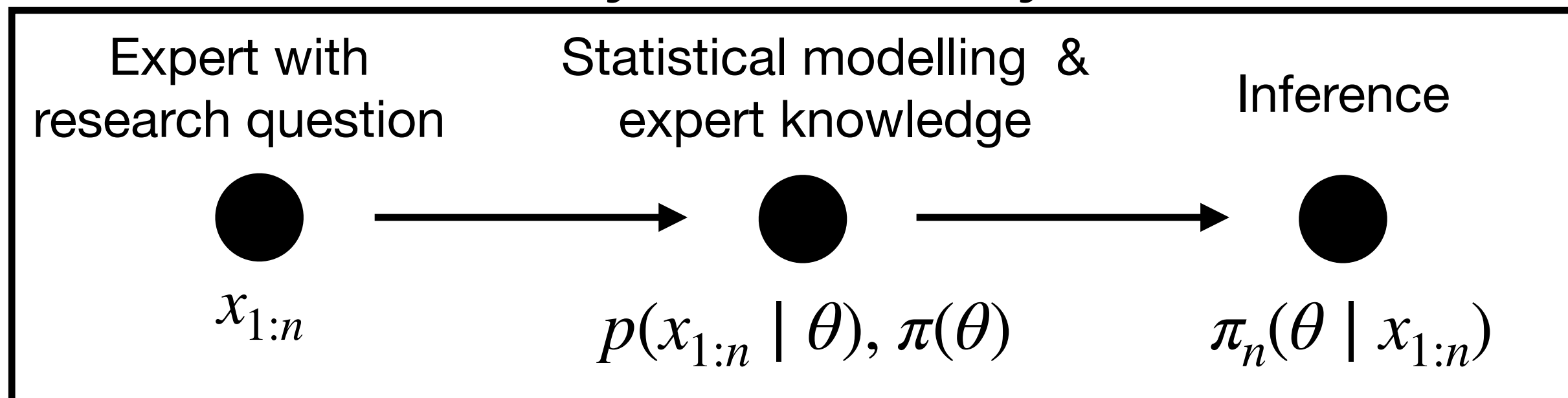
(A1) model well-specified

(A2) prior well-specified

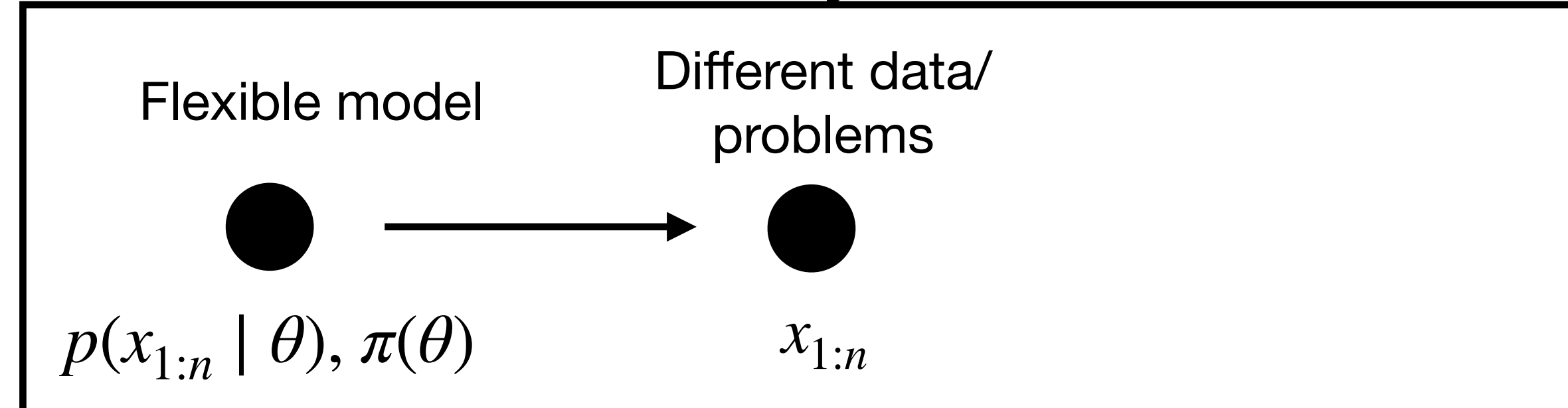
(A3) computationally feasible

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Modern Bayesian ML



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

willingness to pay \uparrow p_j \uparrow pollutants \uparrow rooms, sqm, ... \uparrow measurement error \uparrow ε_i

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

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

$\pi_n(\theta | x_{1:n}) \rightarrow$ computed exactly

(A2) ✓

(A3) ✓

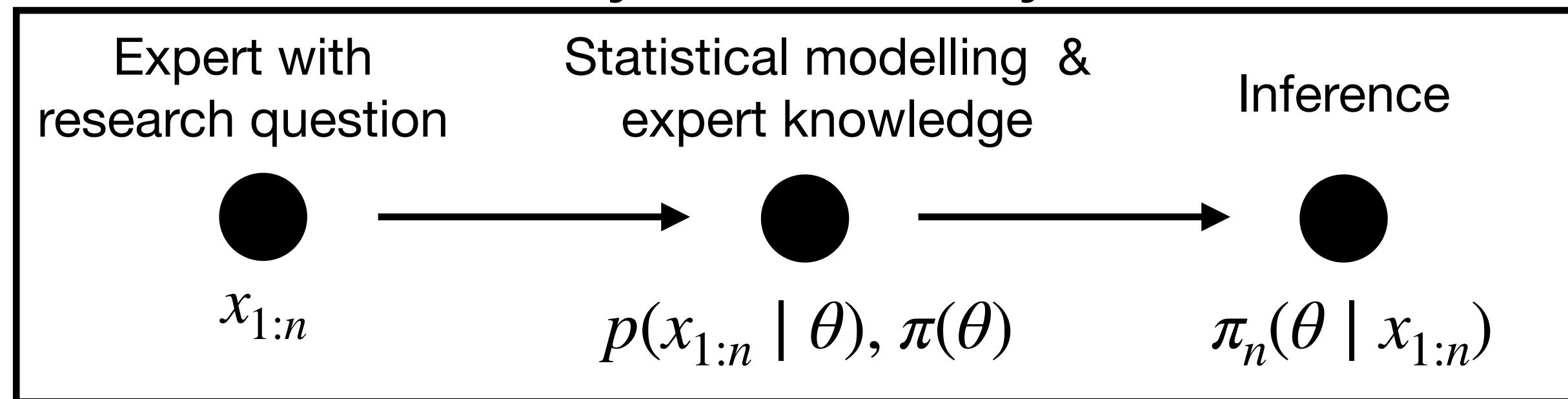
(A1) model well-specified

(A2) prior well-specified

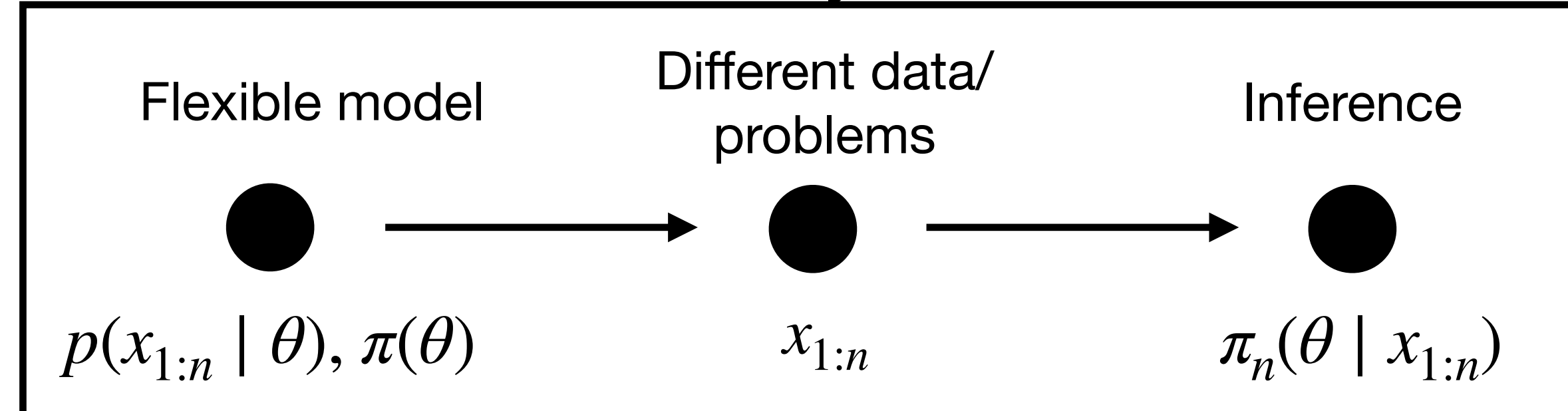
(A3) computationally feasible

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Modern Bayesian ML



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

willingness to pay \uparrow p_j \uparrow pollutants \uparrow rooms, sqm, ... \uparrow measurement error \uparrow ε_i

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

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

$\pi_n(\theta | x_{1:n}) \rightarrow$ computed exactly

(A2) ✓

(A3) ✓

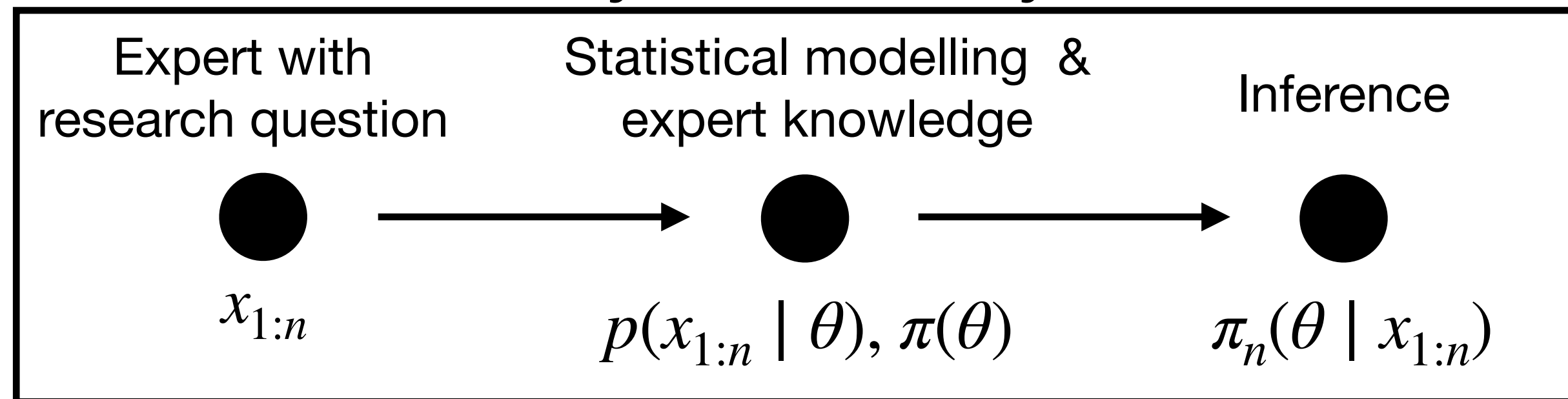
(A1) model well-specified

(A2) prior well-specified

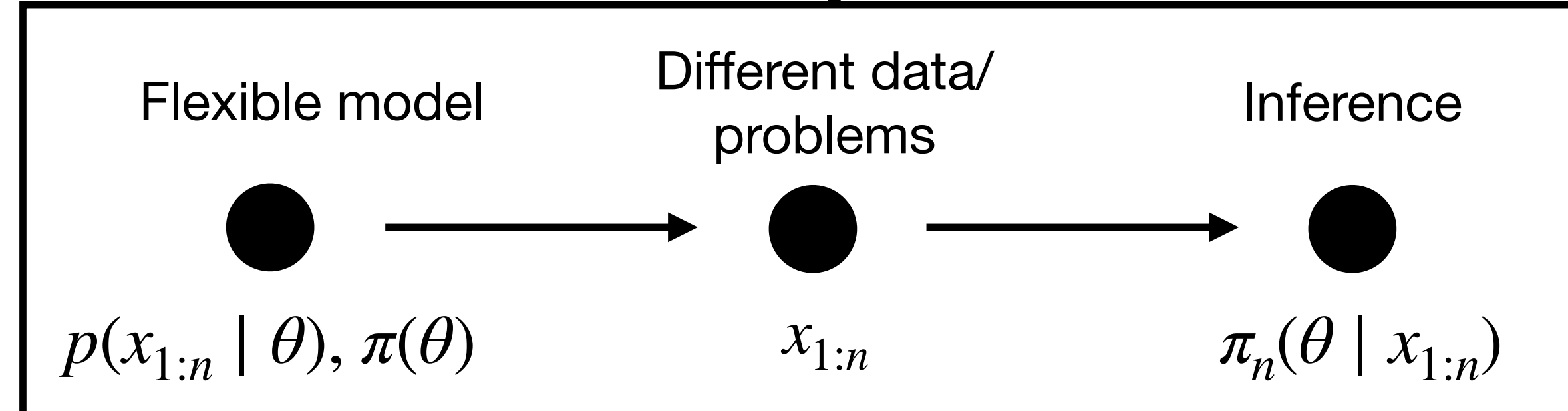
(A3) computationally feasible

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Modern Bayesian ML



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

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

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

parameters of interest incidental parameters

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

$\pi_n(\theta | x_{1:n}) \longrightarrow$ computed exactly

(A2) ✓

(A3) ✓

Pearce et al. (2020) [AISTATS]

Research Question: Does my algorithm improve prediction on regression tasks like Boston UCI data?

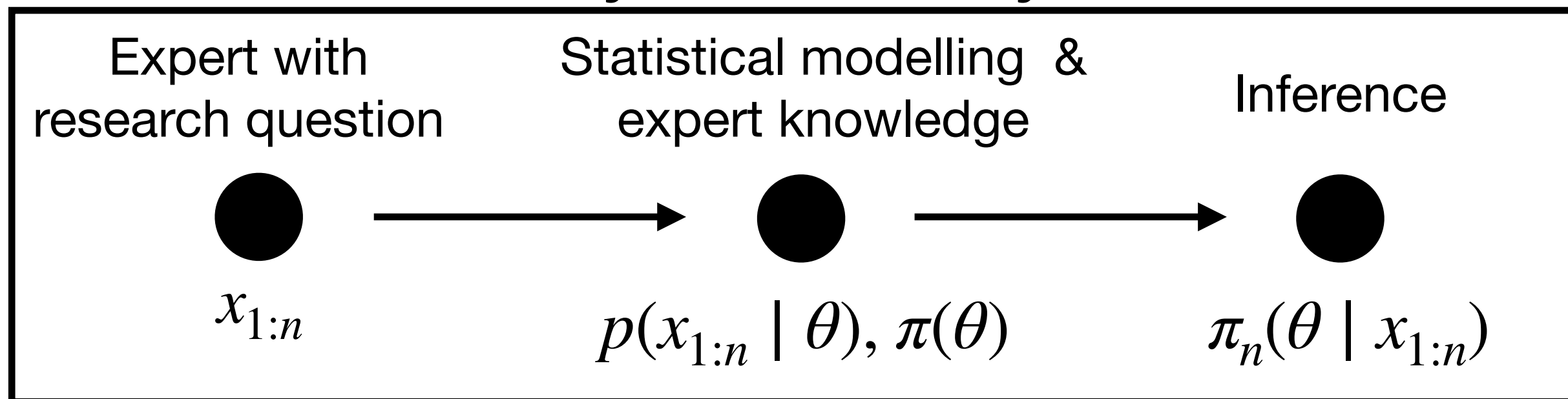
(A1) model well-specified

(A2) prior well-specified

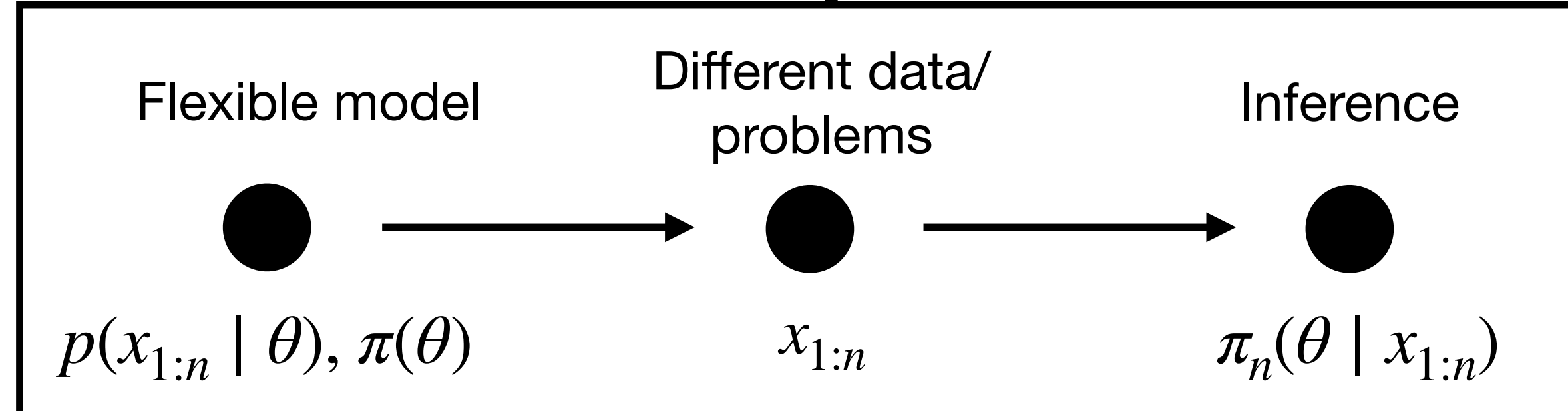
(A3) computationally feasible

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Modern Bayesian ML



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

\uparrow willingness to pay \uparrow pollutants \uparrow rooms, sqm, ... \uparrow measurement error

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

parameters of interest incidental parameters

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

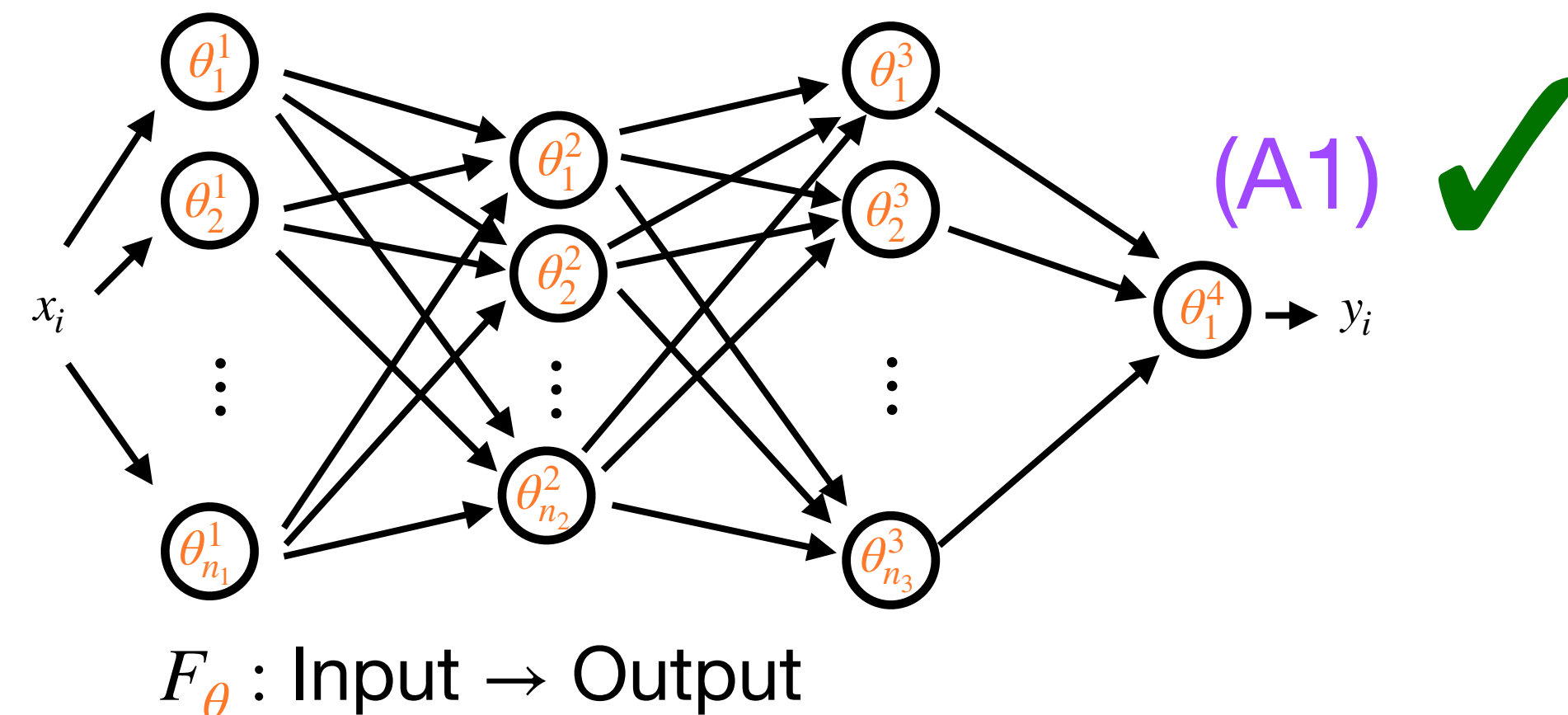
$\pi_n(\theta | x_{1:n}) \rightarrow$ computed exactly

(A2) ✓

(A3) ✓

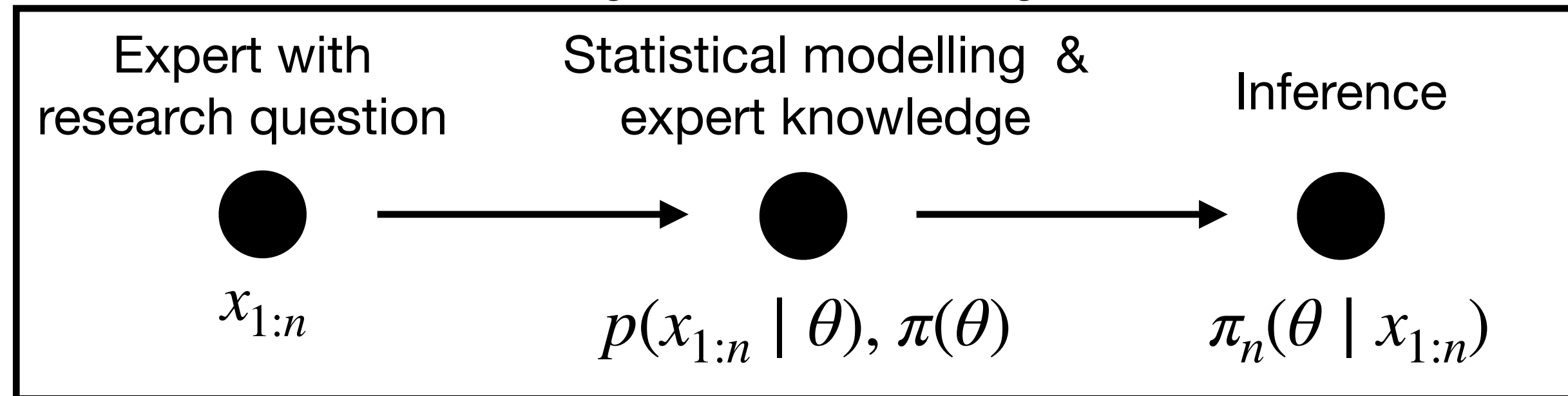
Pearce et al. (2020) [AISTATS]

Research Question: Does my algorithm improve prediction on regression tasks like Boston UCI data?

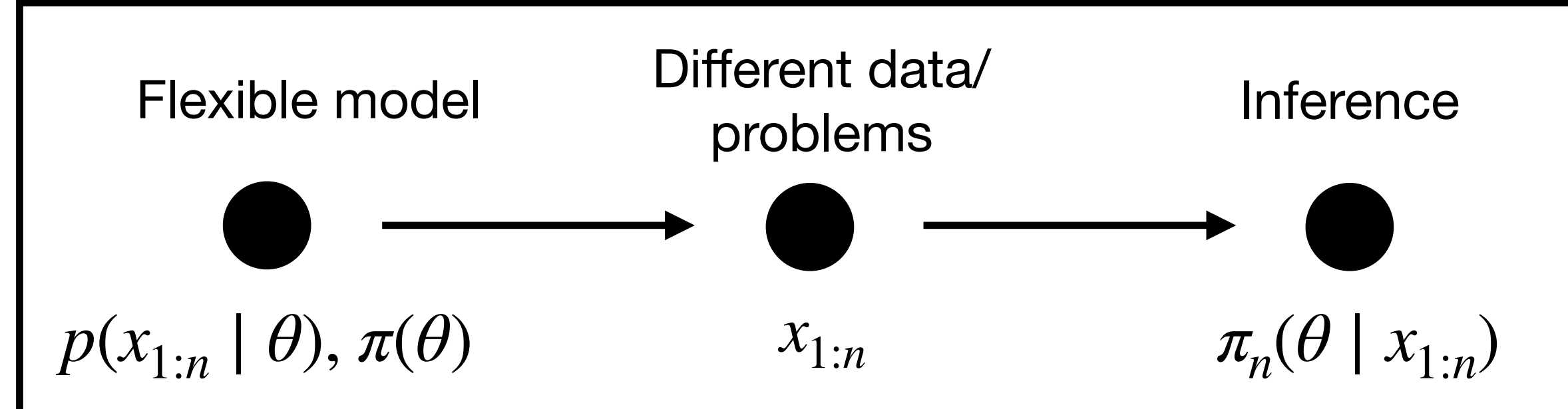


Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Modern Bayesian ML



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

\uparrow willingness to pay \uparrow pollutants \uparrow rooms, sqm, ... \uparrow measurement error

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

θ parameters of interest (blue) incidental parameters (orange)

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

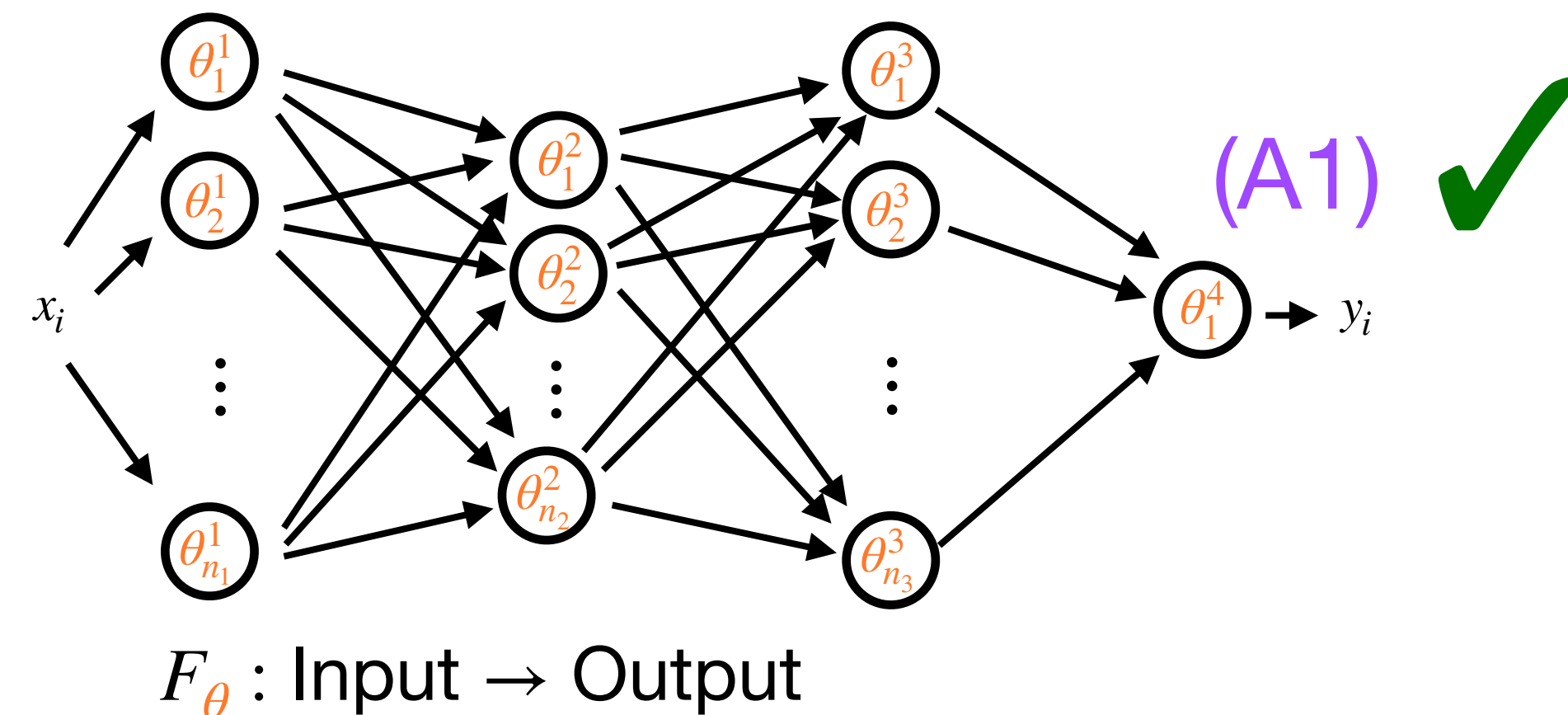
$\pi_n(\theta | x_{1:n}) \rightarrow$ computed exactly

(A2) ✓

(A3) ✓

Pearce et al. (2020) [AISTATS]

Research Question: Does my algorithm improve prediction on regression tasks like Boston UCI data?

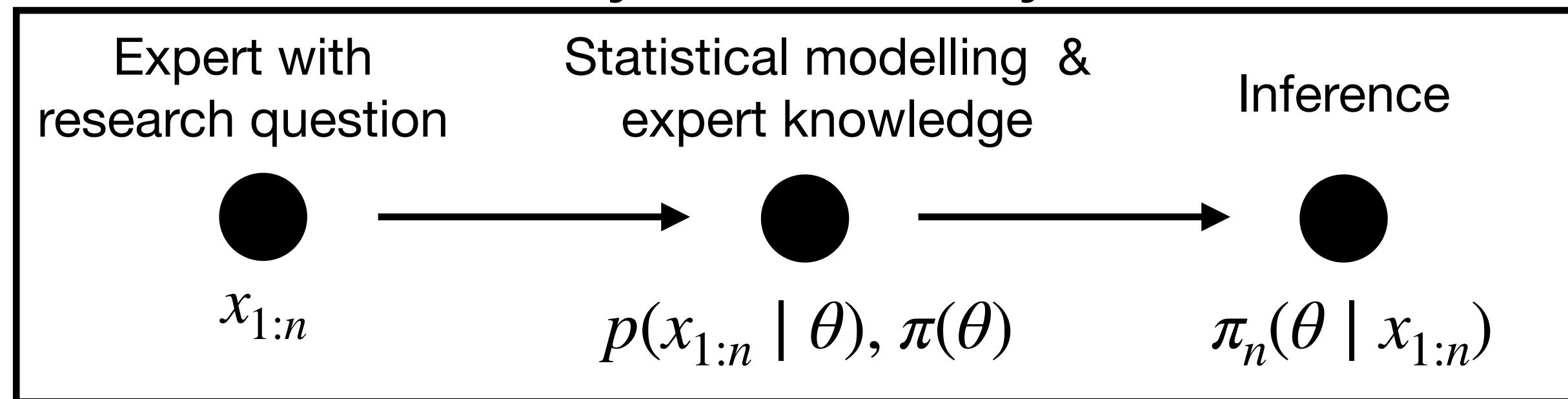


$\pi(\theta) \sim$ Normal

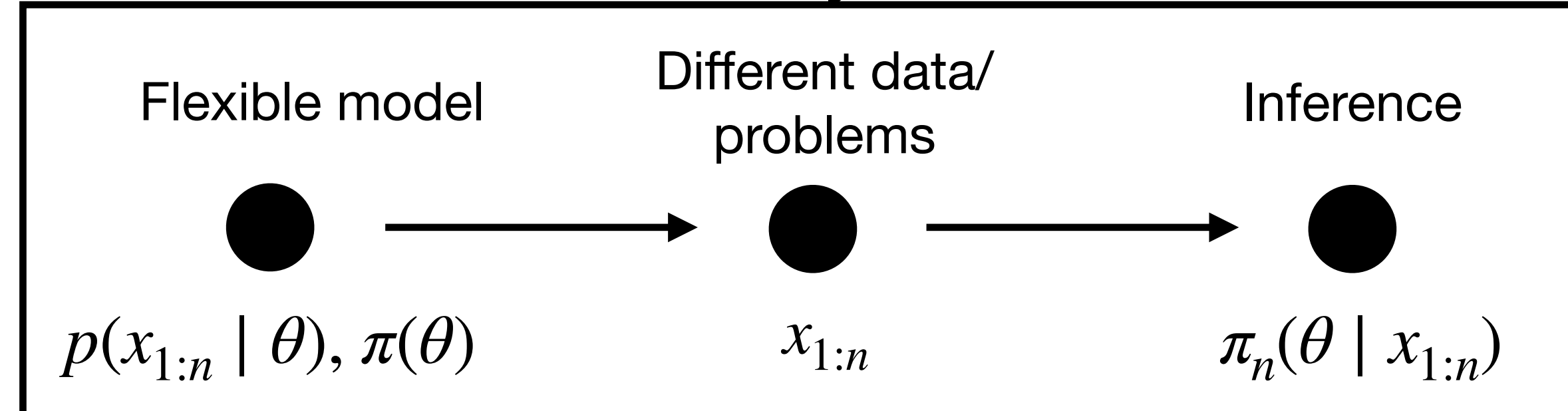
~~(A2)~~

Case Study: Regression with Boston Housing Data

Traditional Bayesian analysis in science



Modern Bayesian ML



Harrison & Rubinfeld (1978)

Research Question: influence of air pollution on house prices?

(A1) ✓

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

\uparrow willingness to pay \uparrow pollutants \uparrow rooms, sqm, ... \uparrow measurement error

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

θ components: p_j (parameters of interest), c_j (incidental parameters)

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

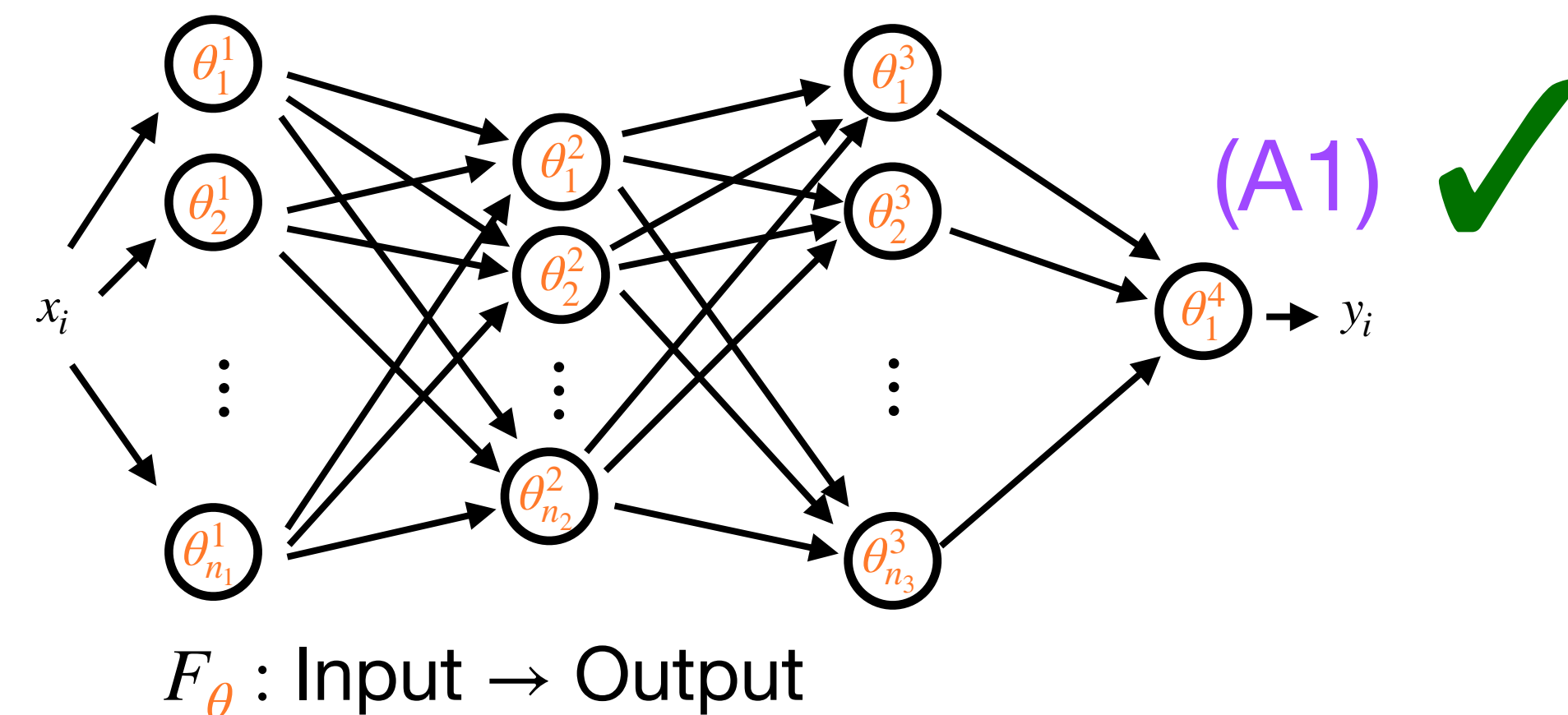
$\pi_n(\theta | x_{1:n}) \rightarrow$ computed exactly

(A2) ✓

(A3) ✓

Pearce et al. (2020) [AISTATS]

Research Question: Does my algorithm improve prediction on regression tasks like Boston UCI data?



$\pi(\theta) \sim$ Normal

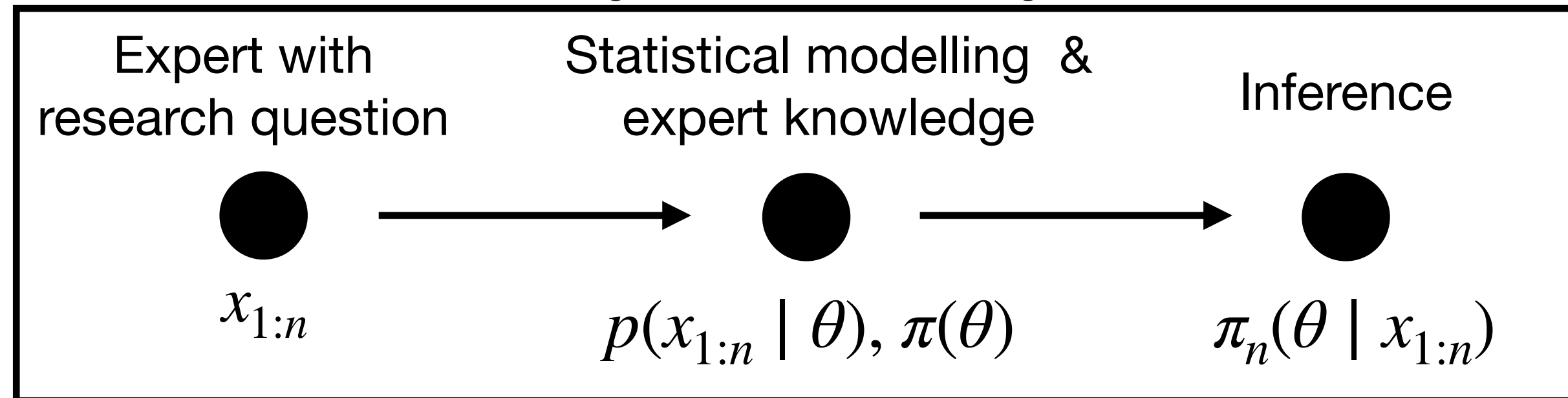
$\pi_n(\theta | x_{1:n}) \rightarrow$ coarse approximation

~~(A2)~~

~~(A3)~~

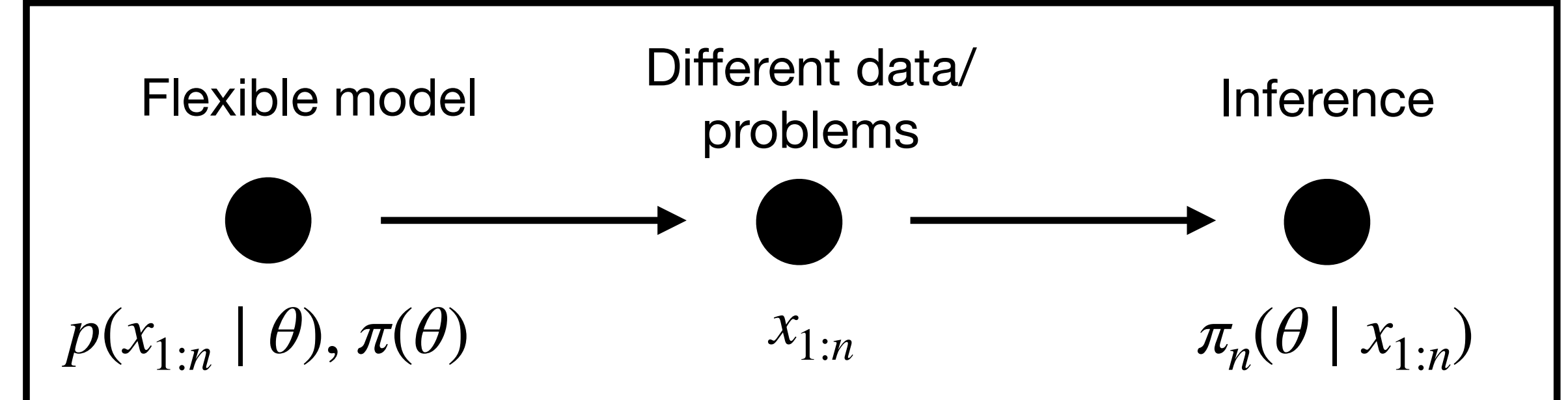
Assumptions & Foundations

Traditional Bayesian analysis in science



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

Modern Bayesian ML



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

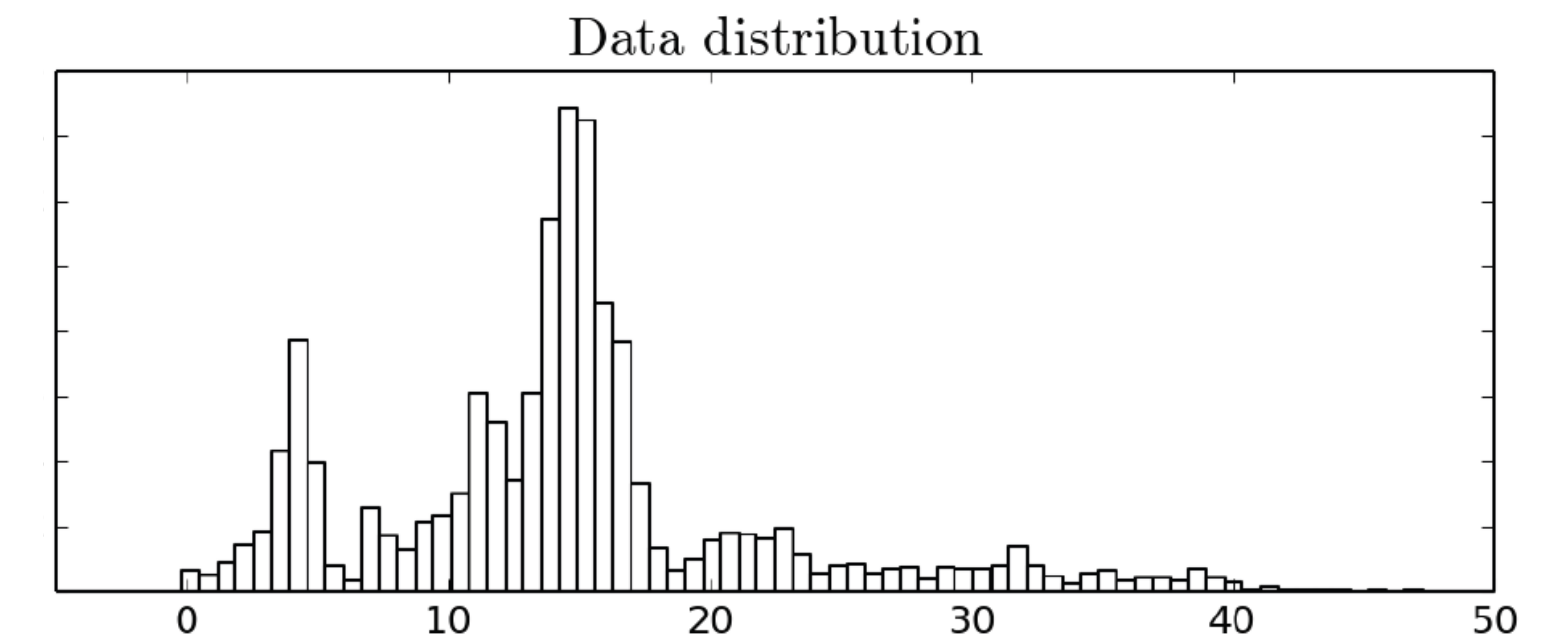


Ramifications

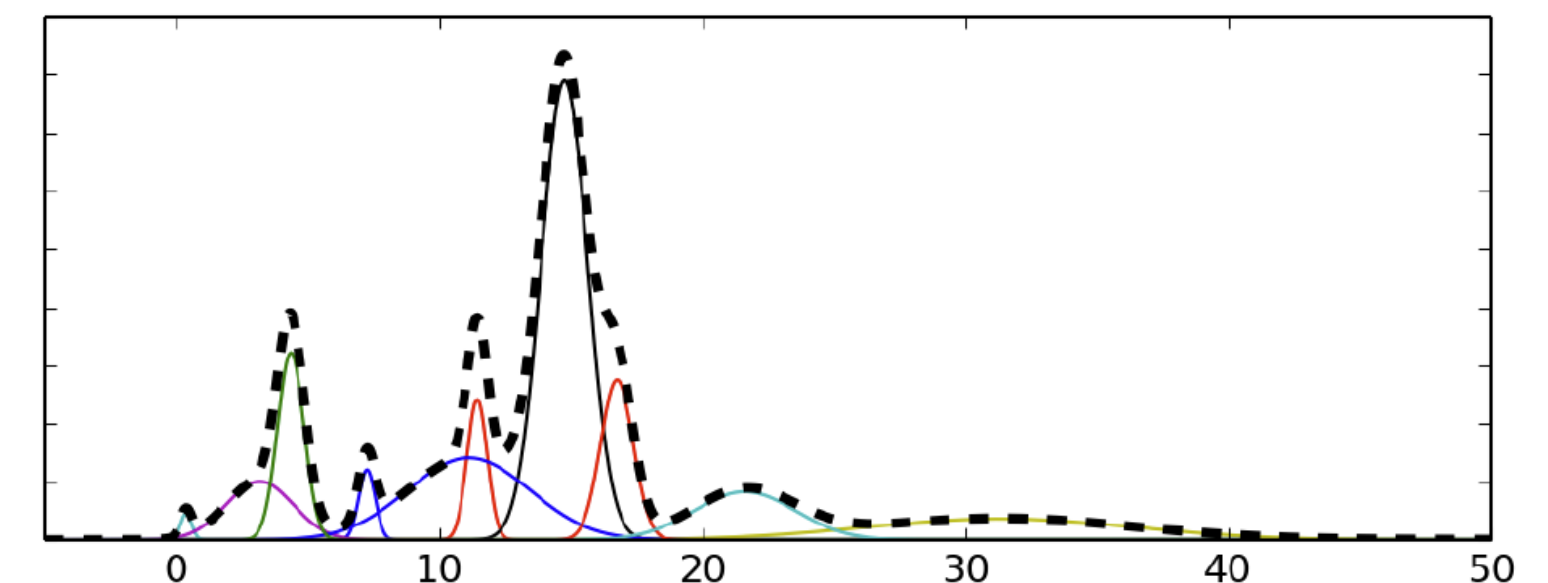
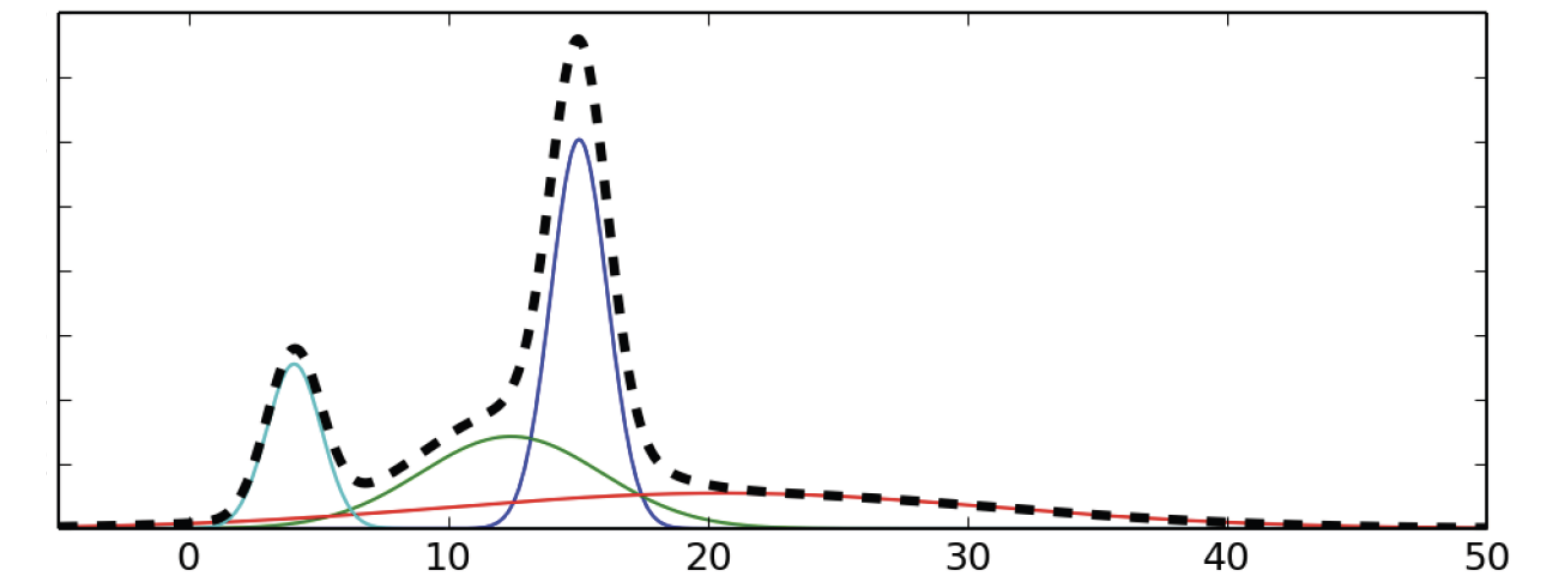
On the **Brittleness** of Bayesian Inference (Ohwadi et al., 2015)

Bayesian inference is generically ill-posed: [...]

- (1) with two [...] arbitrarily close models and [...] data [one] may reach diametrically opposite conclusions; and*
- (2) any given prior and model can be slightly perturbed to achieve any desired posterior*



Posteriors based on two different priors



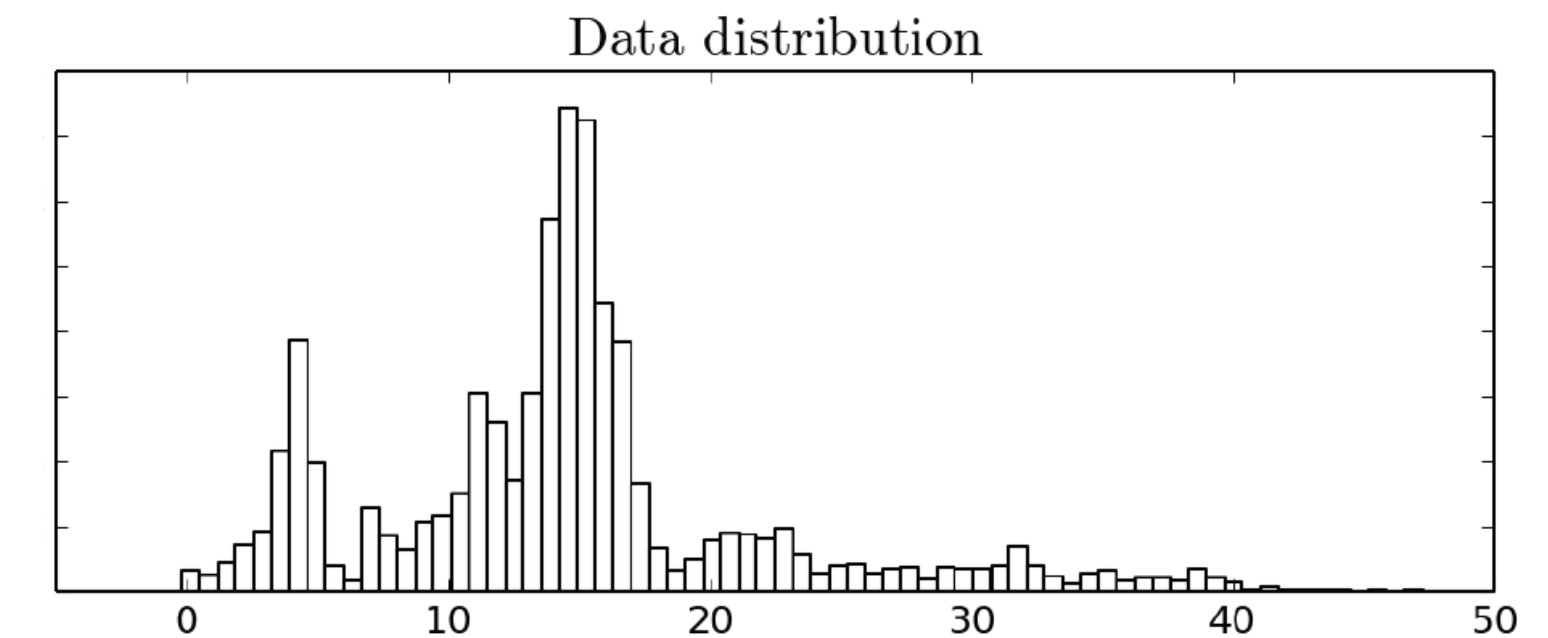
Ramifications

On the **Brittleness** of Bayesian Inference (Ohwadi et al., 2015)

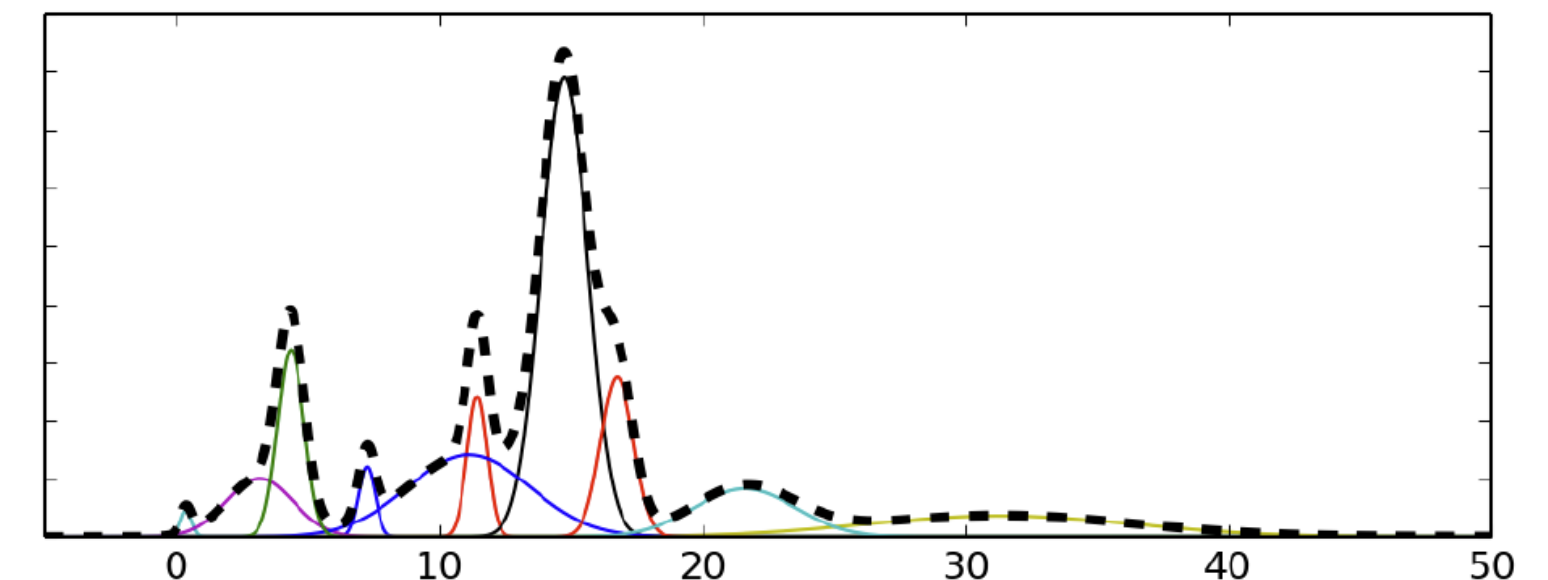
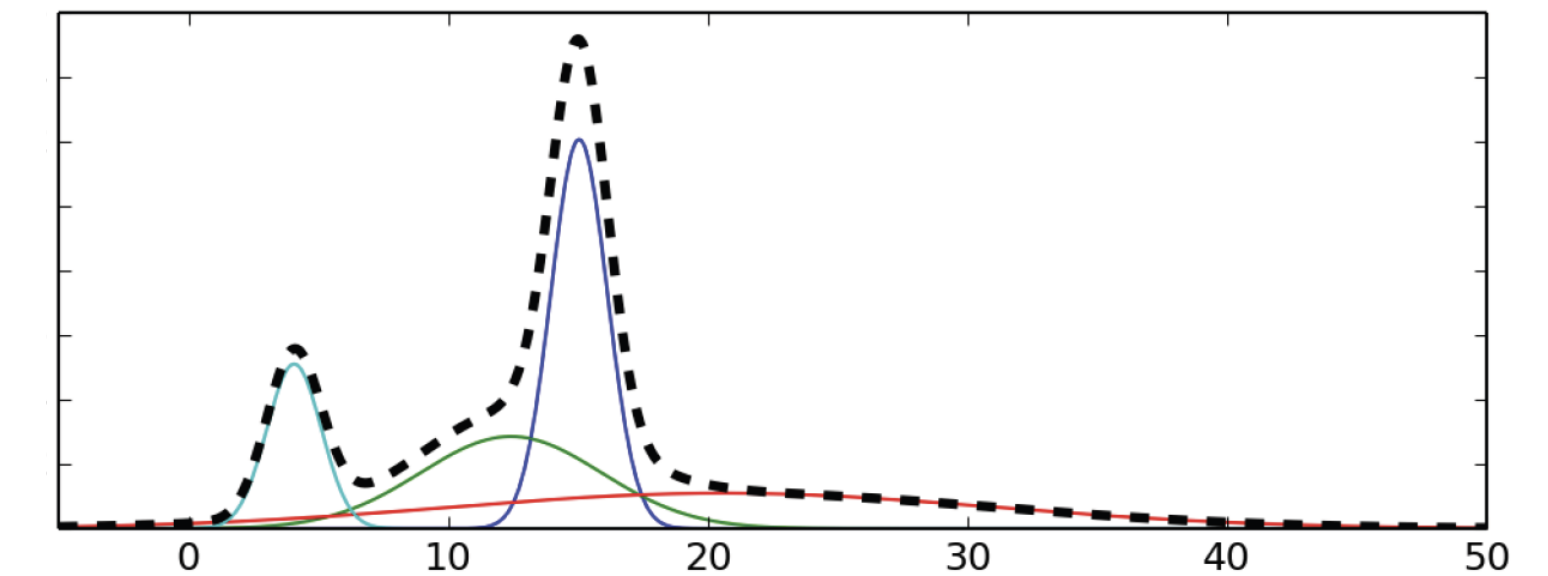
Bayesian inference is generically ill-posed: [...]

- (1) with two [...] arbitrarily close models and [...] data [one] may reach diametrically opposite conclusions; and*
- (2) any given prior and model can be slightly perturbed to achieve any desired posterior*

Plain English: Bayesian ML is often unreliable!



Posteriors based on two different priors



Ramifications: Post-Bayesian ML

Mathematical Foundations

1764—1786: **Bayes' Theorem**

1930: **DeFinetti's Representation Theorem**

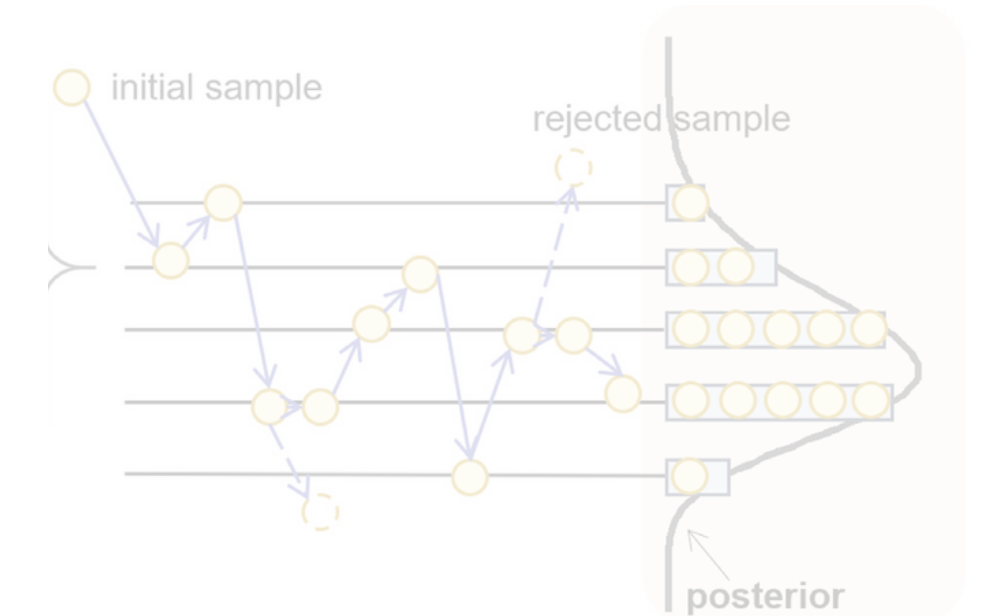
1950s/60s: **Savage Axioms, Birnbaum's Likelihood Principle, ...**

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

Computation

1960s/70s: **Computational Principles** (Markov chain Monte Carlo)

1980s onwards: computation becomes **feasible**



Bayesian ML

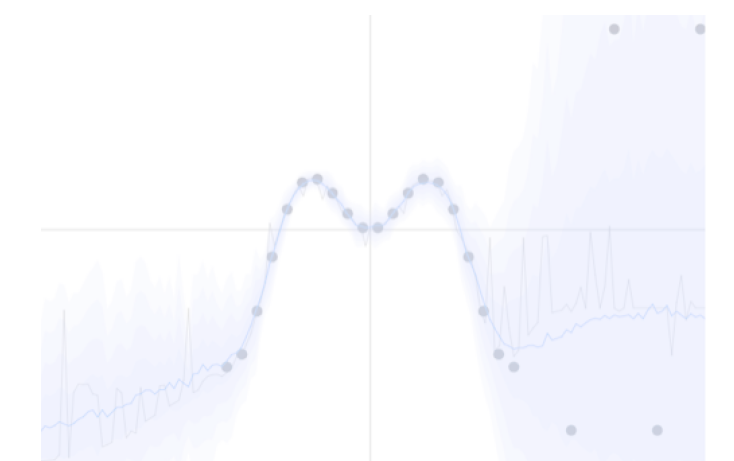
1990s: **ML** becomes data-driven

Classic perspective

Mid 1990s—mid 2000s: **Bayesian ML** emerges

on Bayesian ML

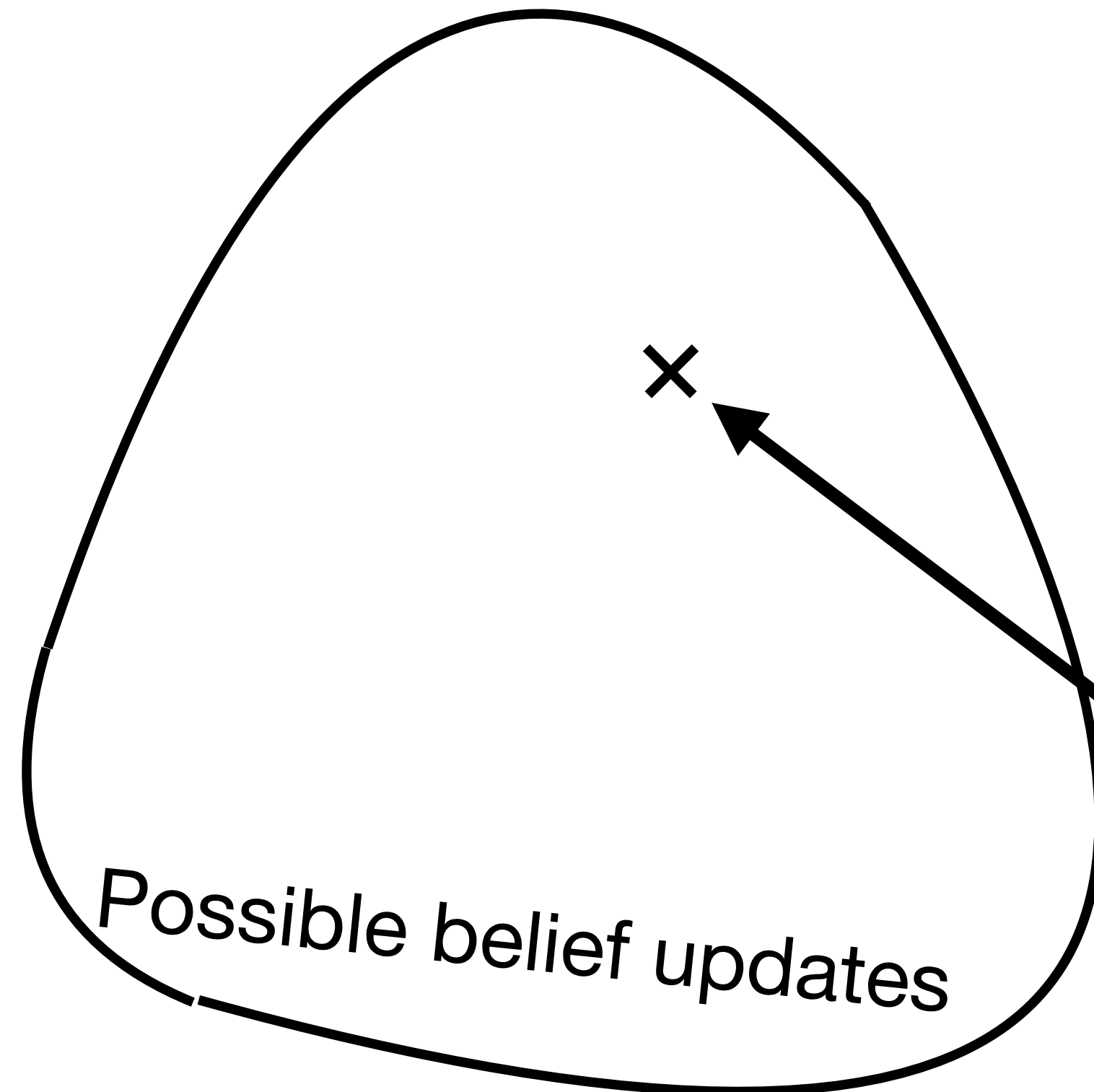
2000s—present: **Adaptations of Bayesian principles** for ML/DL



Collectively: Post-Bayesian ML

Post-Bayesian ML

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

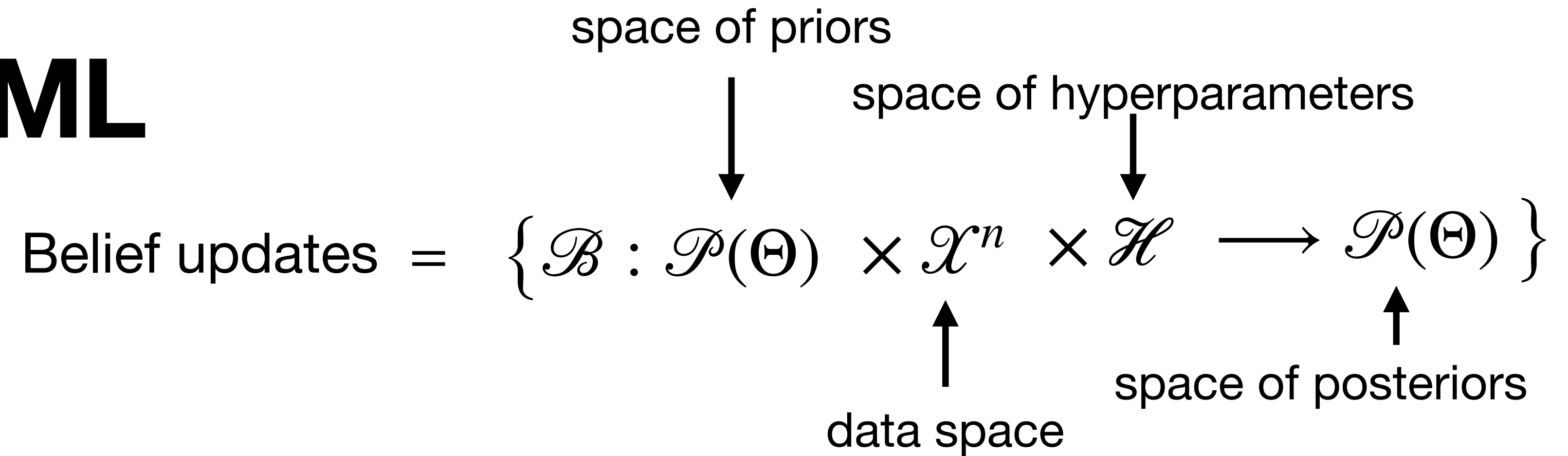


Bayes' Posterior

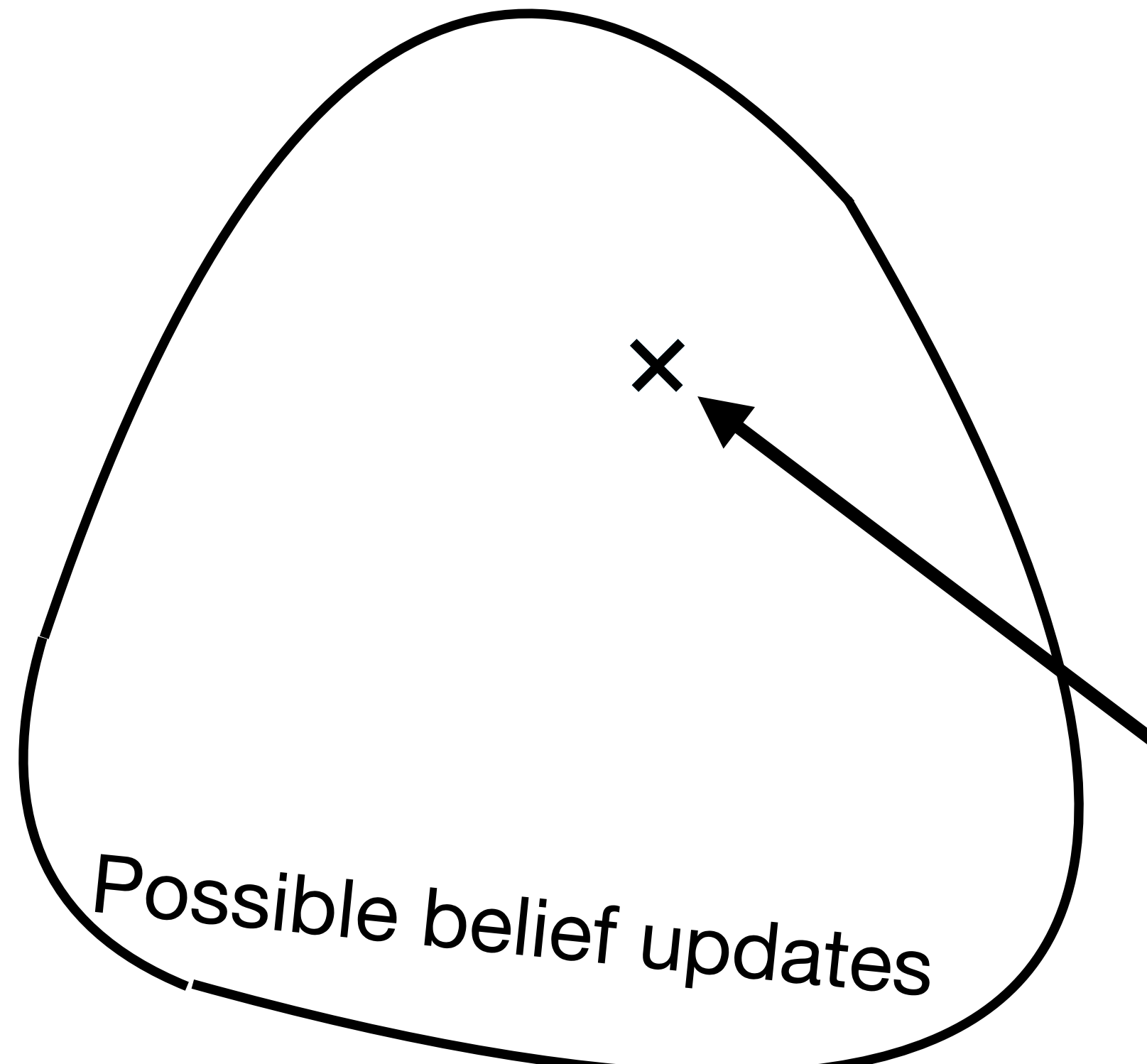
(A1), (A2), (A3)

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

Post-Bayesian ML



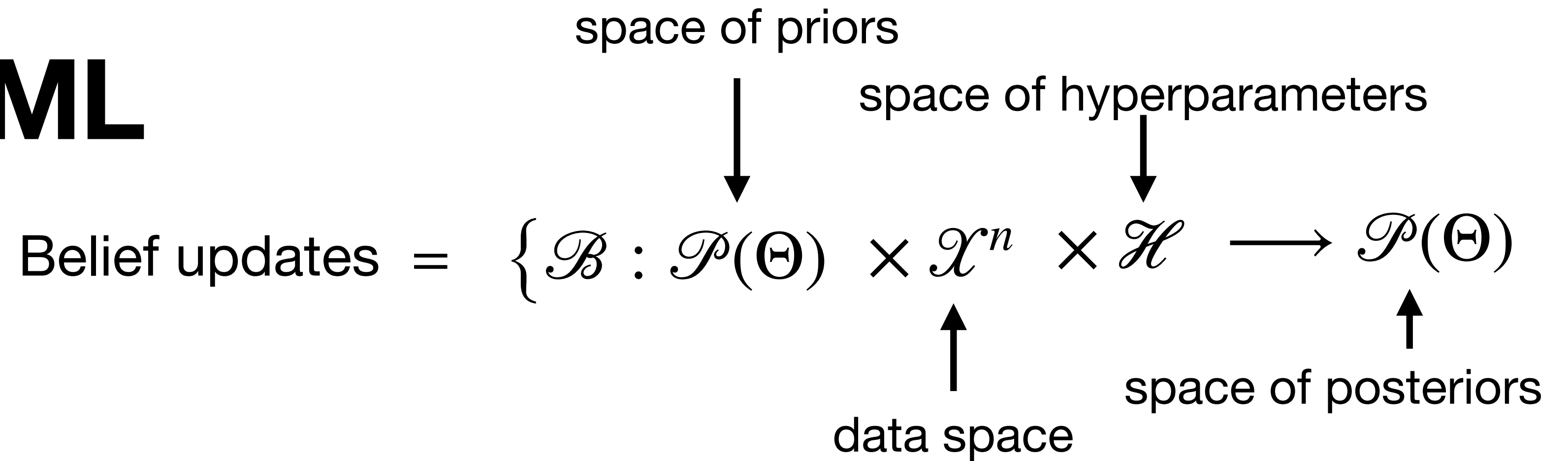
- (A1) model well-specified
- (A2) prior well-specified
- (A3) computationally feasible



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

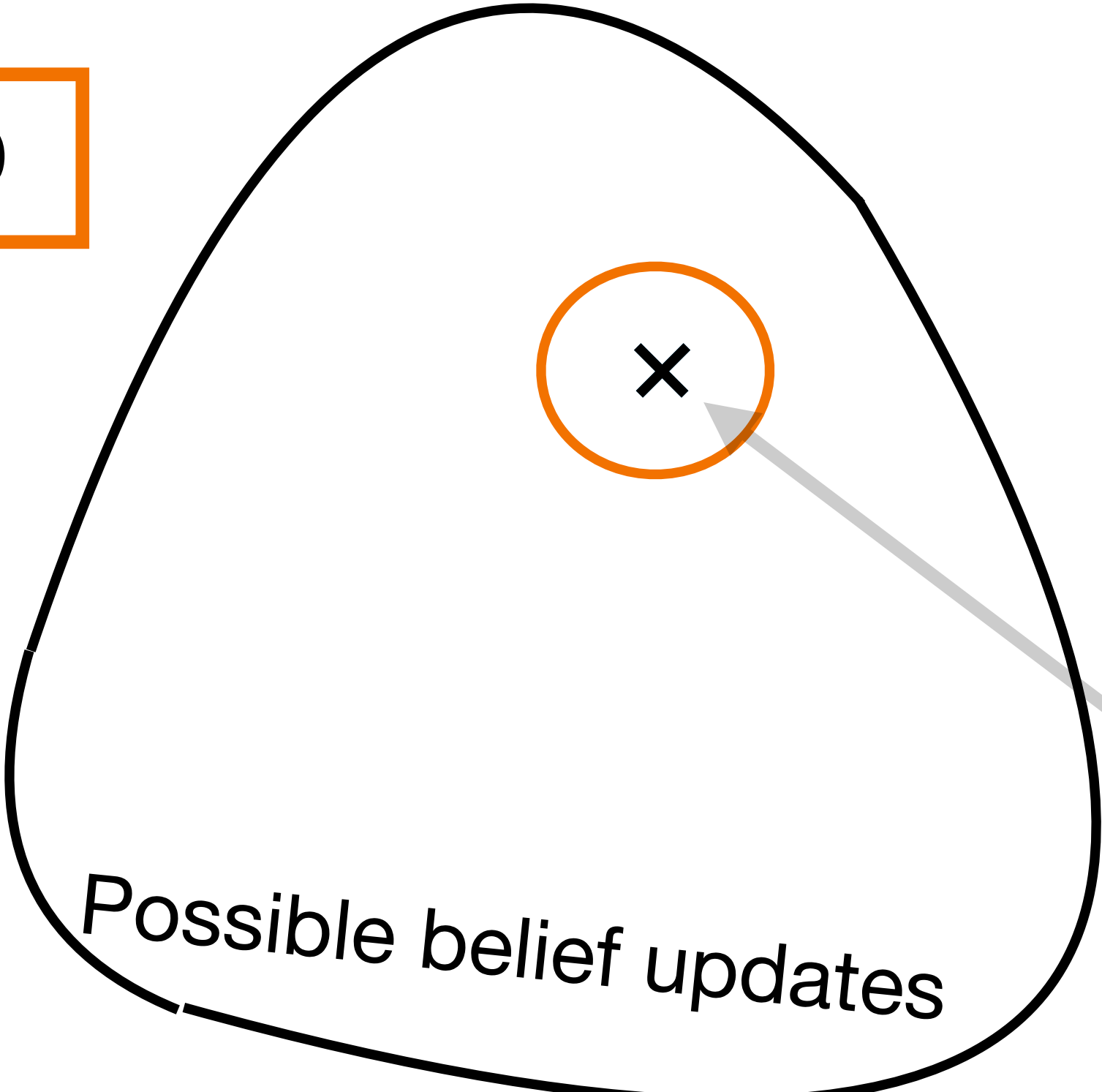
$$\pi_n(\theta | x_{1:n}) = \frac{p(x_{1:n} | \theta) \cdot \pi(\theta)}{\int p(x_{1:n} | \theta) \cdot \pi(\theta) d\theta}$$

Post-Bayesian ML



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

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

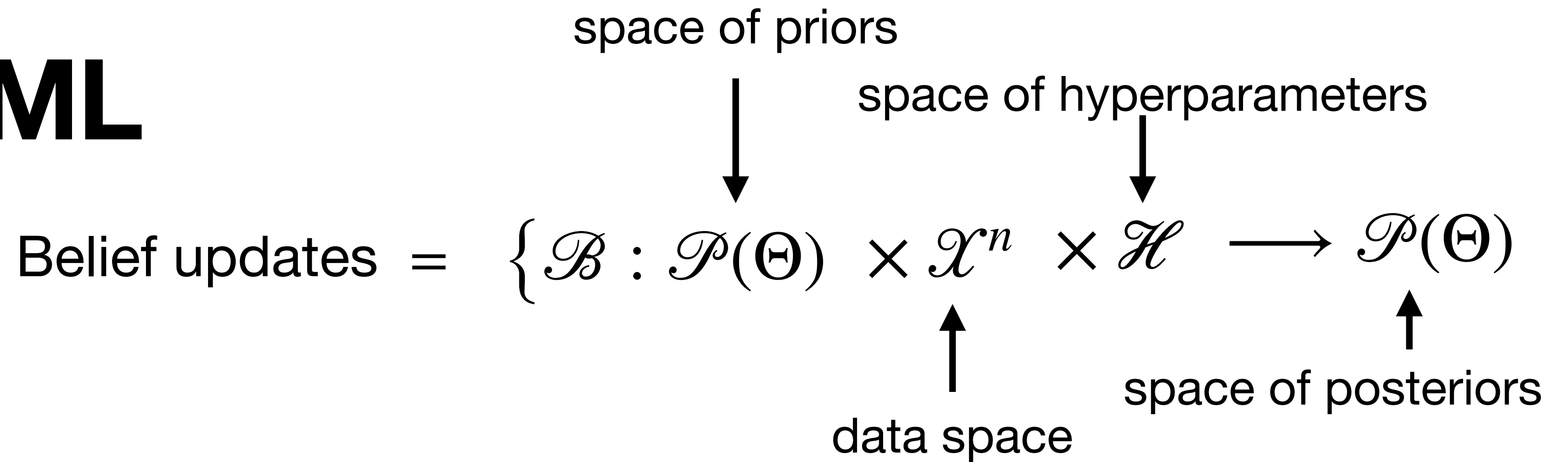


Bayes' Posterior (A1), (A2), (A3)

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

Post-Bayesian ML

Grünwald (2011); COLT
 Miller & Dunson (2015); JRSS-B
 Bhattacharya, Pati, & Yang (2019); Annals of Statistics
 Adlam et al. (2020); preprint
 Wenzel et al. (2020); ICML
 Aitchison (2021); ICLR
 ...
 McLatchie, Fong, Frazier, & Knoblauch. (2024); forthcoming



[Generally credited to Grünwald (2011)]

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

Power/Fractional/
Cold Posterior

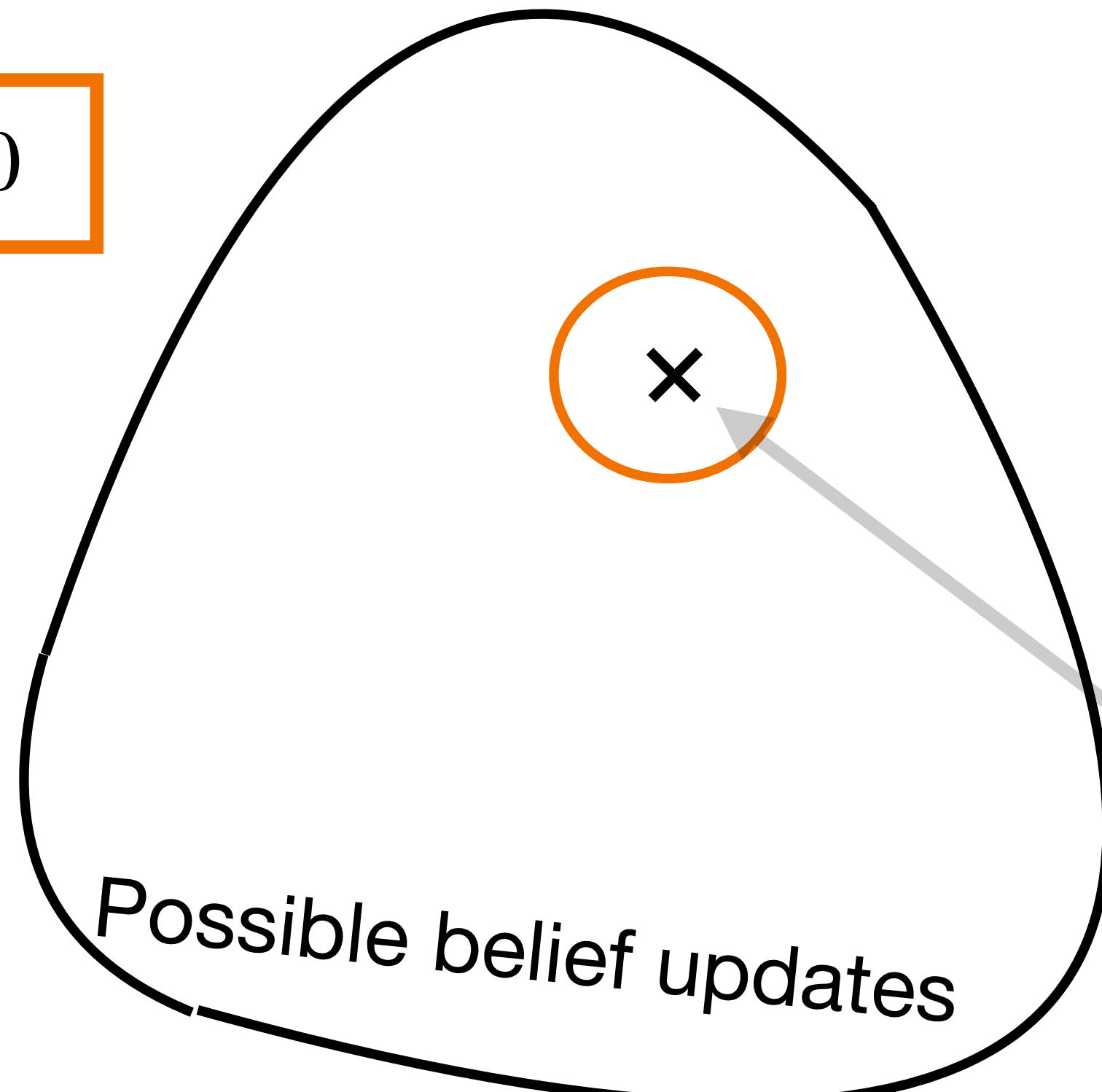
~~(A1)~~, (A2), (A3)

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

Bayes' Posterior

(A1), (A2), (A3)

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



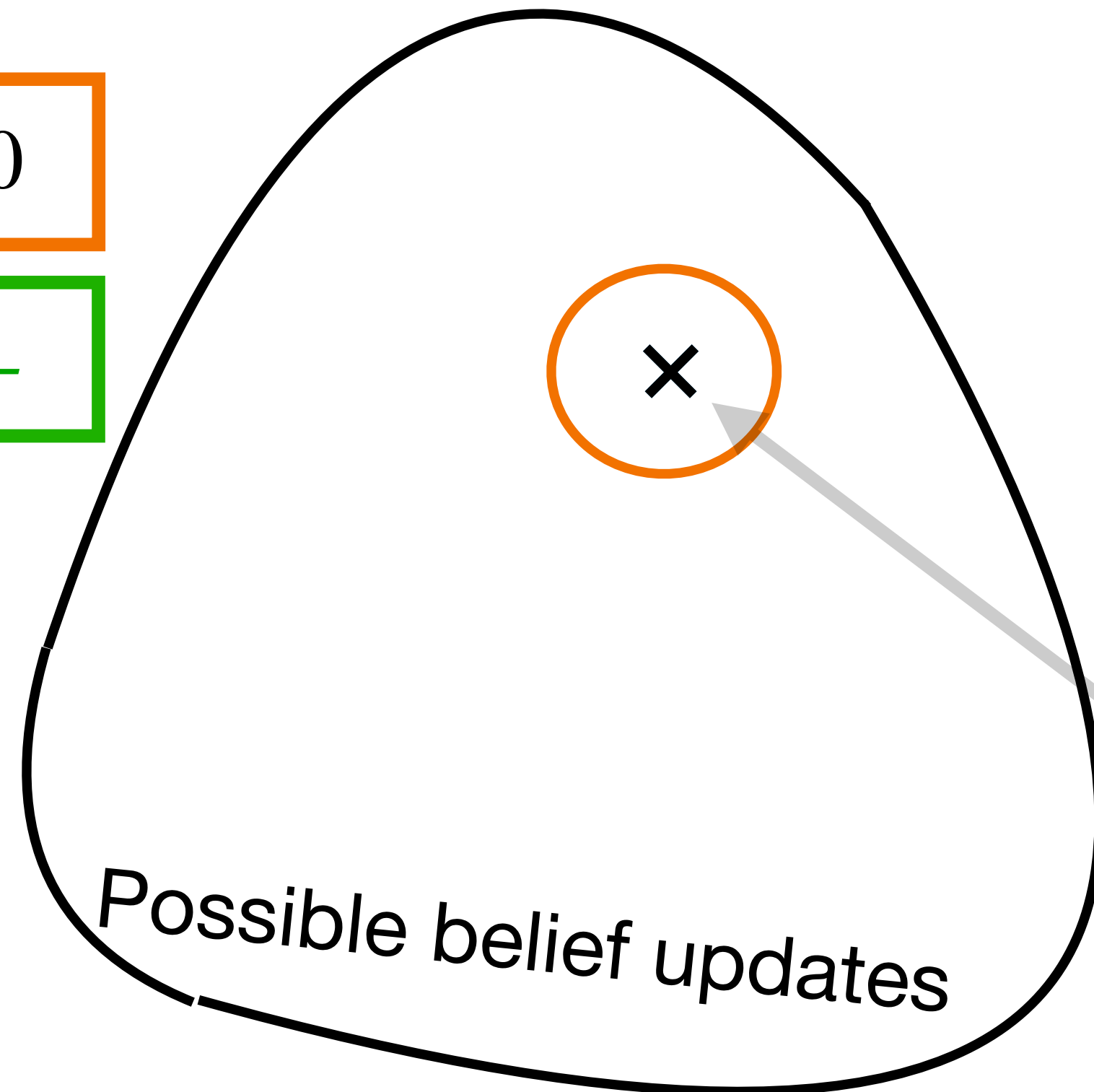
- (A1) model well-specified
- (A2) prior well-specified
- (A3) computationally feasible

Post-Bayesian ML

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

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

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



Power/Fractional/
Cold Posterior

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

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

Bayes' Posterior

(A1), (A2), (A3)

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

Post-Bayesian ML

Langford & Shawe-Taylor (2002); NeurIPS

Seeger (2002); ICML

Bissiri et al. (2016); JRSS-B

...

Knoblauch & Damoulas. (2018); ICML

Knoblauch, Jewson, & Damoulas et al. (2018); NeurIPS

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B

Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA

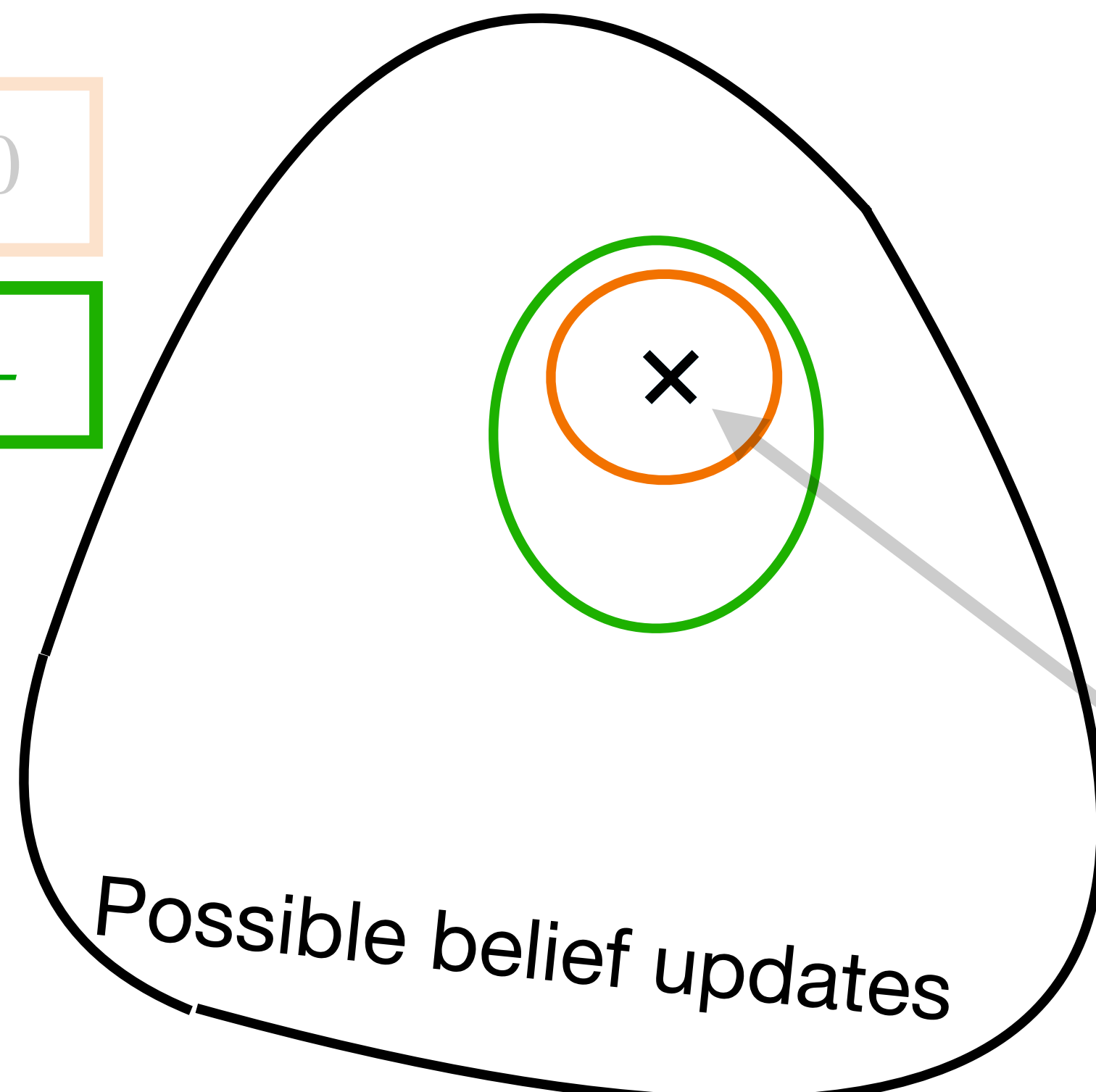
Altamirano, Briol, & **Knoblauch** (2023); ICML

Altamirano, Briol, & **Knoblauch** (2024); ICML spotlight

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

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

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



[Generally credited to Bissiri, Holmes & Walker (2016)]

Gibbs/Generalised/
Pseudo Posterior

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

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

Power/Fractional/
Cold Posterior

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

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

Bayes' Posterior

(A1), (A2), (A3)

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

Post-Bayesian ML

Optimisation-centric posteriors /
Generalised Variational Inference

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

$$q_n^*(\theta) = \arg \min_{q \in \mathcal{Q}} \underbrace{\mathcal{L}_{L, D}(q)}; \quad \mathcal{Q} \subseteq \mathcal{P}(\Theta)$$

$$= \int L(x_{1:n}, \theta) q(\theta) d\theta + D(q, \pi)$$

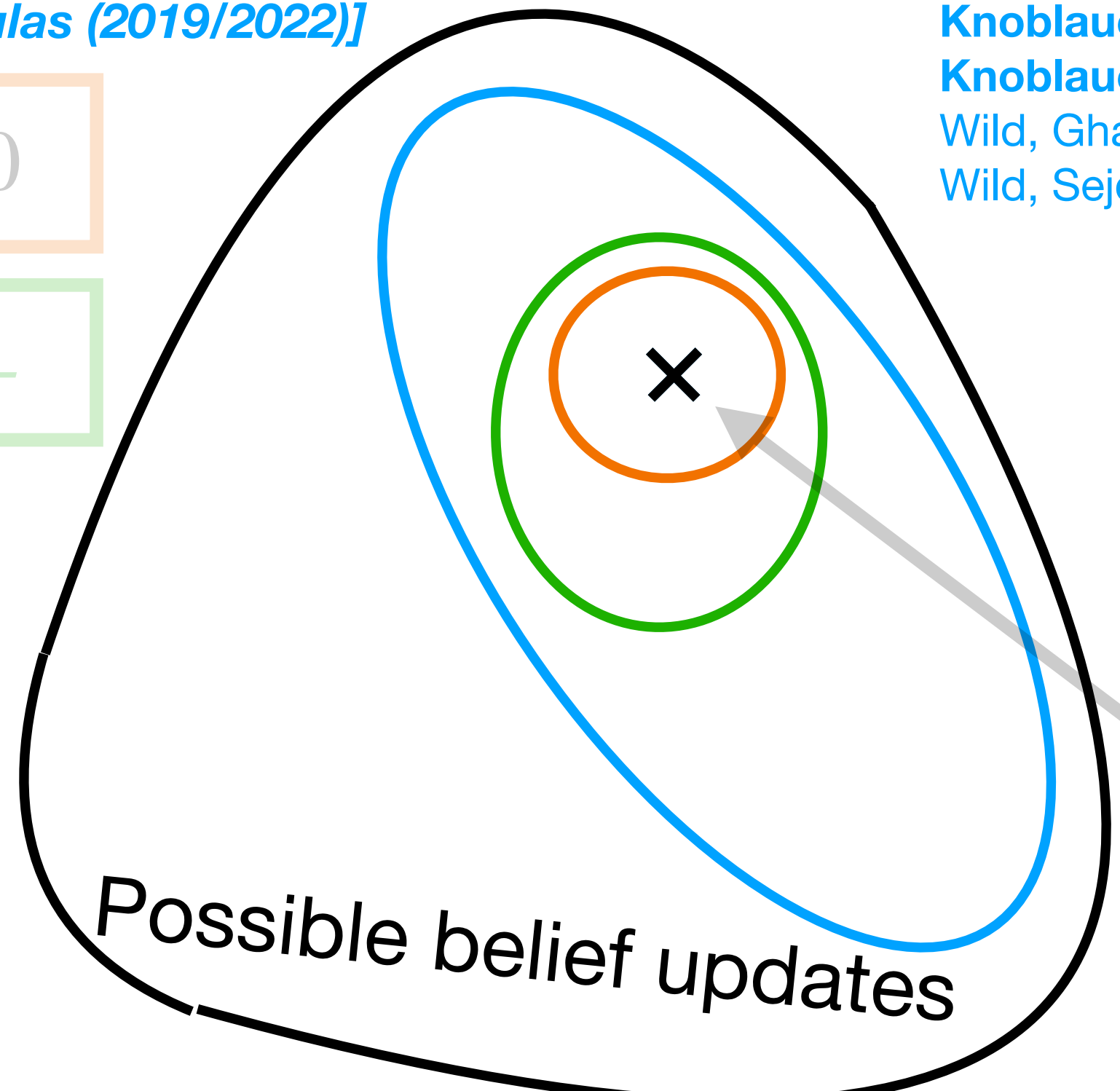
[Generally credited to **Knoblauch, Jewson, & Damoulas (2019/2022)**]

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

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

KL	→	D
$\mathcal{P}(\Theta)$	→	\mathcal{Q}

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



Gibbs/Generalised/
Pseudo Posterior

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

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

Xuan, Wu, Liu, & Lu (2024); UAI
Chi, Zhang, Yang, Ouyang, & Pei (2024); AAAI

...
Knoblauch, Jewson, & Damoulas (2019); NeurIPS
Knoblauch, Jewson, & Damoulas (2022); JMLR
Wild, Ghalebikesabi, Sejdinovic, & **Knoblauch (2023)**; NeurIPS (oral)
Wild, Sejdinovic, & **Knoblauch (2024)**; forthcoming

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

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

Bayes' Posterior

(A1), (A2), (A3)

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

Post-Bayesian ML

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

$$\pi_n^L(\theta \mid x_{1:n}) = \frac{\exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta)}{\int \exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta) d\theta} = \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \int L(x_{1:n}, \theta) q(\theta) d\theta + \text{KL}(q, \pi) \right\}$$

Husain & **Knoblauch** (2022); ALT

Knoblauch, Jewson, & Damoulas (2022); JMLR

Wild, Sejdinovic, & **Knoblauch** (2024); forthcoming

Wild, Ghalebikesabi, Sejdinovic, & **Knoblauch** (2024); NeurIPS (oral)

Post-Bayesian ML

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

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

$$= \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \int L(x_{1:n}, \theta) q(\theta) d\theta + \text{KL}(q, \pi) \right\}$$

~~(A1)~~ by using
robust loss L

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

Optimisation-centric posteriors /
Generalised Variational Inference (GVI) $\quad \quad \quad =: \mathcal{L}_{L, D}(q)$

Post-Bayesian ML

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

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

$$= \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \int L(x_{1:n}, \theta) q(\theta) d\theta + \text{KL}(q, \pi) \right\}$$

~~(A1)~~ by using
robust loss L

~~(A2)~~ by using
robust regulariser D

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

Optimisation-centric posteriors /
Generalised Variational Inference (GVI)

$$=: \mathcal{L}_{L, D}(q)$$

Post-Bayesian ML

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

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

$$= \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \int L(x_{1:n}, \theta) q(\theta) d\theta + \text{KL}(q, \pi) \right\}$$

~~(A3)~~ by optimising over a set $\mathcal{Q} \subseteq \mathcal{P}(\Theta)$

~~(A1)~~ by using robust loss L

~~(A2)~~ by using robust regulariser D

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

Optimisation-centric posteriors /
Generalised Variational Inference (GVI)

$$=: \mathcal{L}_{L, D}(q)$$

Post-Bayesian ML

Optimisation-centric posteriors /
Generalised Variational Inference

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

$$q_n^*(\theta) = \arg \min_{q \in \mathcal{Q}} \underbrace{\mathcal{L}_{L, D}(q)}; \quad \mathcal{Q} \subseteq \mathcal{P}(\Theta)$$

$$= \int L(x_{1:n}, \theta) q(\theta) d\theta + D(q, \pi)$$

Gibbs/Generalised/
Pseudo Posterior

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

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

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

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

$$\begin{array}{l} \text{KL} \longrightarrow D \\ \mathcal{P}(\Theta) \longrightarrow \mathcal{Q} \end{array}$$

Power/Fractional/
Cold Posterior

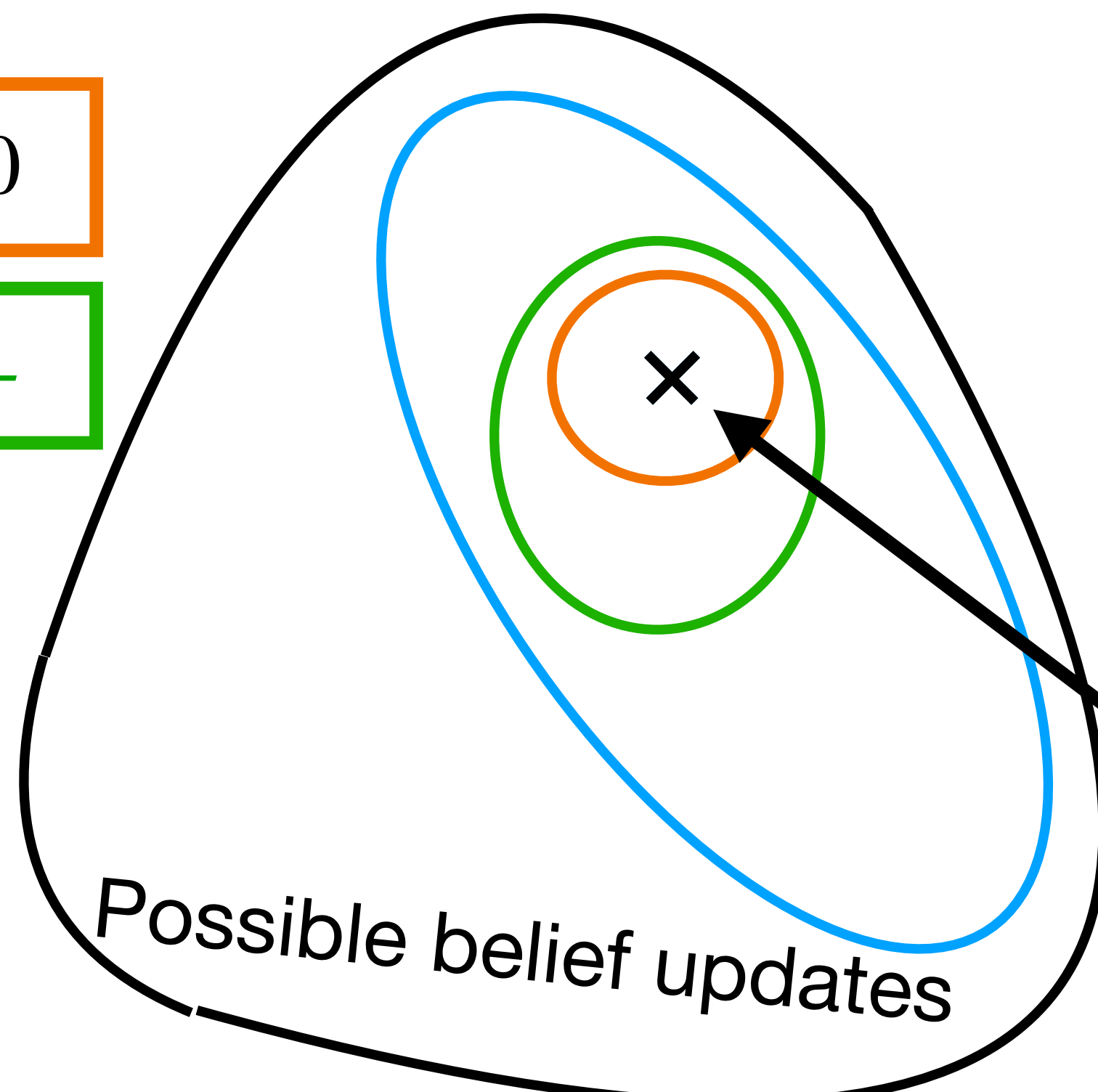
~~(A1)~~, (A2), (A3)

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

Bayes' Posterior

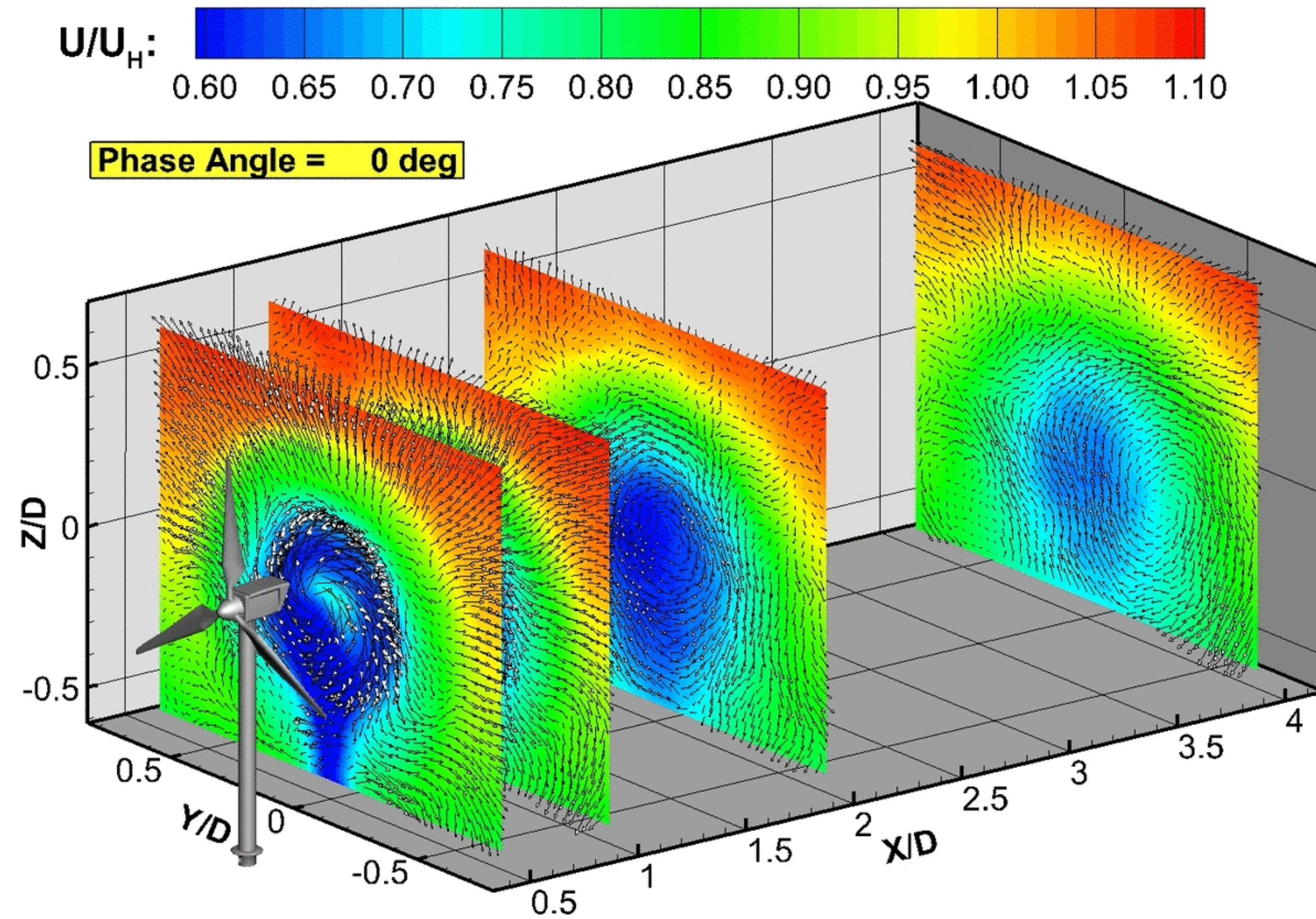
(A1), (A2), (A3)

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



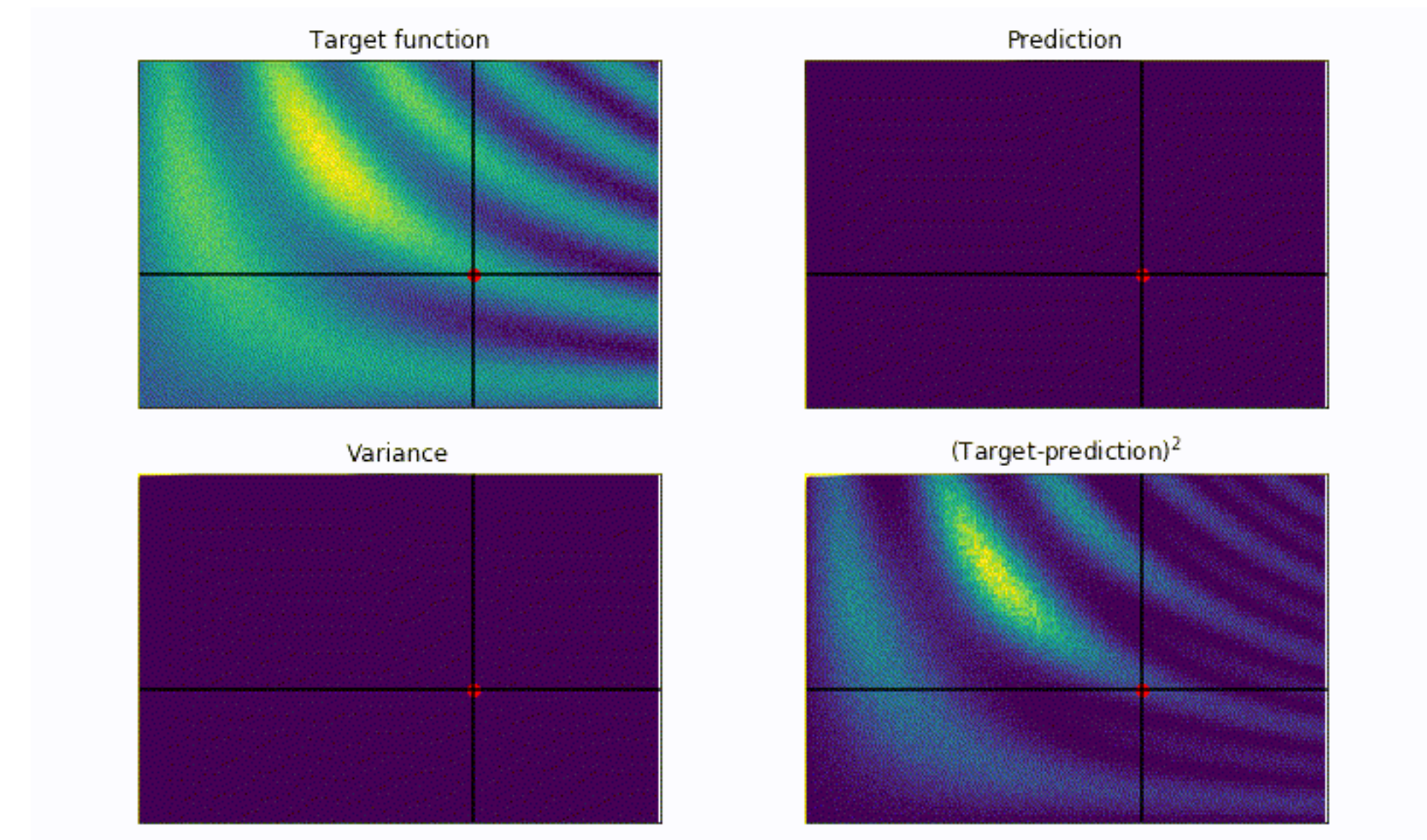
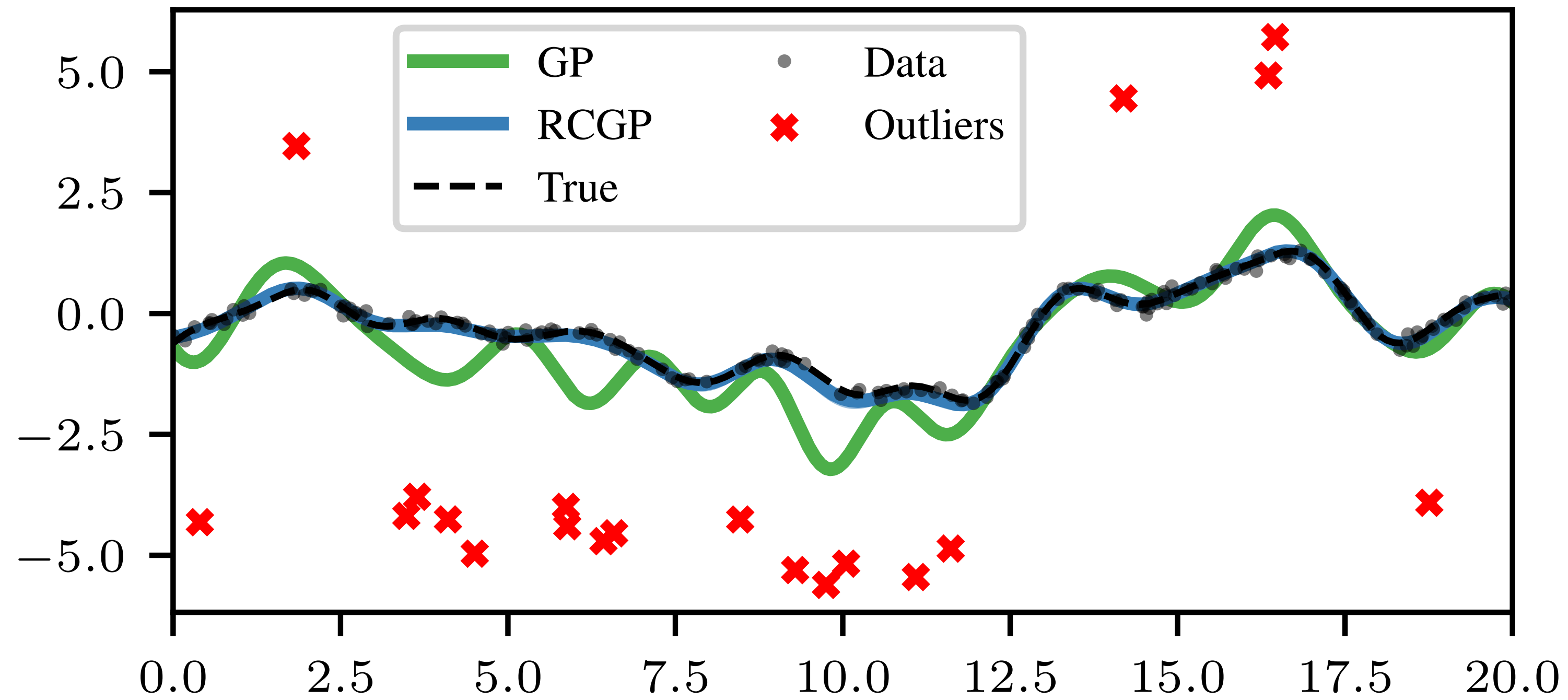
- (A1) model well-specified
- (A2) prior well-specified
- (A3) computationally feasible

Post-Bayesian ML: Recent Success Stories



Based on work of Dellaportas, **Knoblauch**, Damoulas, & Briol (2022); AISTATS (best paper award)

Post-Bayesian ML: Recent Success Stories



Based on work of Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

Post-Bayesian ML Research

① Foundations

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

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

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

Post-Bayesian ML Research

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

1 Foundations

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

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



Knoblauch & Damoulas (2018); ICML
Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Frazier*, Knoblauch*, & Drovandi (2024); preprint
McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

Post-Bayesian ML Research

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

① Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$



Knoblauch & Damoulas (2018); ICML
Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Frazier*, Knoblauch*, & Drovandi (2024); preprint
McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

② State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

Post-Bayesian ML Research

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

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

~~(A1)~~

(A2)

(A3)

Knoblauch & Damoulas (2018); ICML
Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Frazier*, Knoblauch*, & Drovandi (2024); preprint
McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification +
computation

~~(A1)~~

(A2)

~~(A3)~~

Schmon, Cannon, & Knoblauch (2020); AABI
Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
Dellaporta, Knoblauch, Damoulas, & Briol (2022); AISTATS (best paper award)
Altamirano, Briol, & Knoblauch (2023); ICML
Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch, Briol, & Murphy (2024); ICML

Post-Bayesian ML Research

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

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

**model misspecification +
computation**

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

Schmon, Cannon, & Knoblauch (2020); AABI
 Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
 Dellaporta, Knoblauch, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & Knoblauch (2023); ICML
 Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch,
 Briol, & Murphy (2024); ICML

3 The Future

$$q_n^*(\theta)$$

Post-Bayesian ML Research

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

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

**model misspecification +
computation**

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

Schmon, Cannon, & Knoblauch (2020); AABI
 Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
 Dellaporta, Knoblauch, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & Knoblauch (2023); ICML
 Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch, Briol, & Murphy (2024); ICML

3 The Future

$$q_n^*(\theta)$$

**model misspecification +
prior misspecification +
computation**

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

Husain & Knoblauch (2022); ALT
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Matsubara, Knoblauch, Briol, & Oates (2023); JASA
 Wild, Sejdinovic, & Knoblauch (2024); forthcoming
 Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2024); NeurIPS (oral)



Matias Altamirano (UCL)



Yann McLatchie (UCL)



Veit Wild (Oxford)



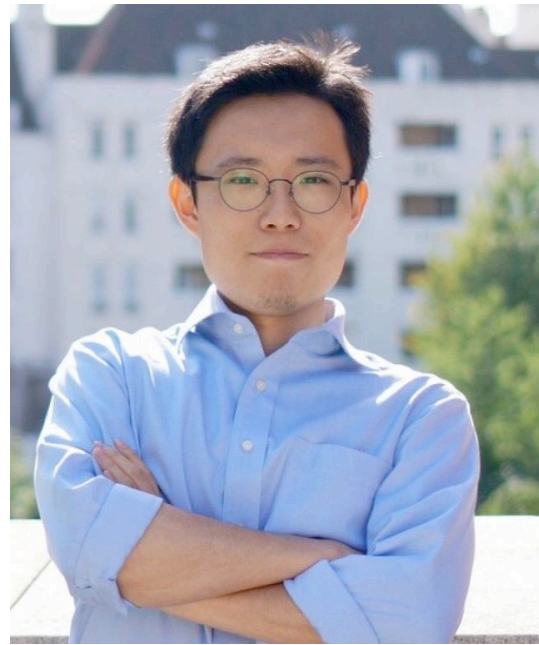
Dino Sejdinovic
(Oxford/Adelaide)



Kevin Murphy (DeepMind)



Chris Drovandi (QUT)



Takuo Matsubara
(Edinburgh)



Gerardo Duran-Martin
(QMU/Oxford)



Sahra Ghalebikesabi
(Oxford/DeepMind)



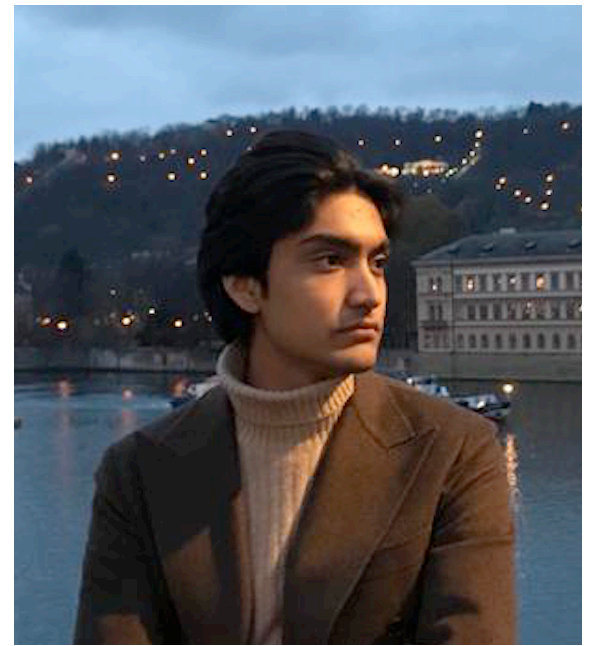
Edwin Fong (Hong Kong)



David Frazier (Monash)



Miheer Dewaskar
(Duke)



Hisham Husain
(Amazon)



Francois-Xavier Briol (UCL)



Chris Oates (Newcastle)



Jack Jewson
(UPF Barcelona/Monash)



Theo Damoulas
(Warwick)



Chris Tosh (Memorial
Sloan Kettering Institute)



Harita Dellaporta
(Warwick)



David Dunson (Duke)

Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification + computation

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

Schmon, Cannon, & Knoblauch (2020); AABI
 Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
 Dellaporta, Knoblauch, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & Knoblauch (2023); ICML
 Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch, Briol, & Murphy (2024); ICML

3 The Future

$$a^*(\theta)$$

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

model misspecification + prior misspecification + computation

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

Husain & Knoblauch (2022); ALT
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Matsubara, Knoblauch, Briol, & Oates (2023); JASA
 Wild, Sejdinovic, & Knoblauch (2024); forthcoming
 Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2024); NeurIPS (oral)

Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification + computation

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

Schmon, Cannon, & Knoblauch (2020); AABI
 Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
 Dellaporta, Knoblauch, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & Knoblauch (2023); ICML
 Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch, Briol, & Murphy (2024); ICML

3 The Future

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

$$a^*(\theta)$$

model misspecification + prior misspecification + computation

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

Husain & Knoblauch (2022); ALT
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Matsubara, Knoblauch, Briol, & Oates (2023); JASA
 Wild, Sejdinovic, & Knoblauch (2024); forthcoming
 Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2024); NeurIPS (oral)

Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

- Q1:** Can tuning λ improve robustness?
- Q2:** What L leads to robust posteriors π_n^L ?
- Q3:** How should we design/choose L ?

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

Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

Q1: Can tuning λ improve robustness?

Q2: What L leads to robust posteriors π_n^L ?

Q3: How should we design/choose L ?

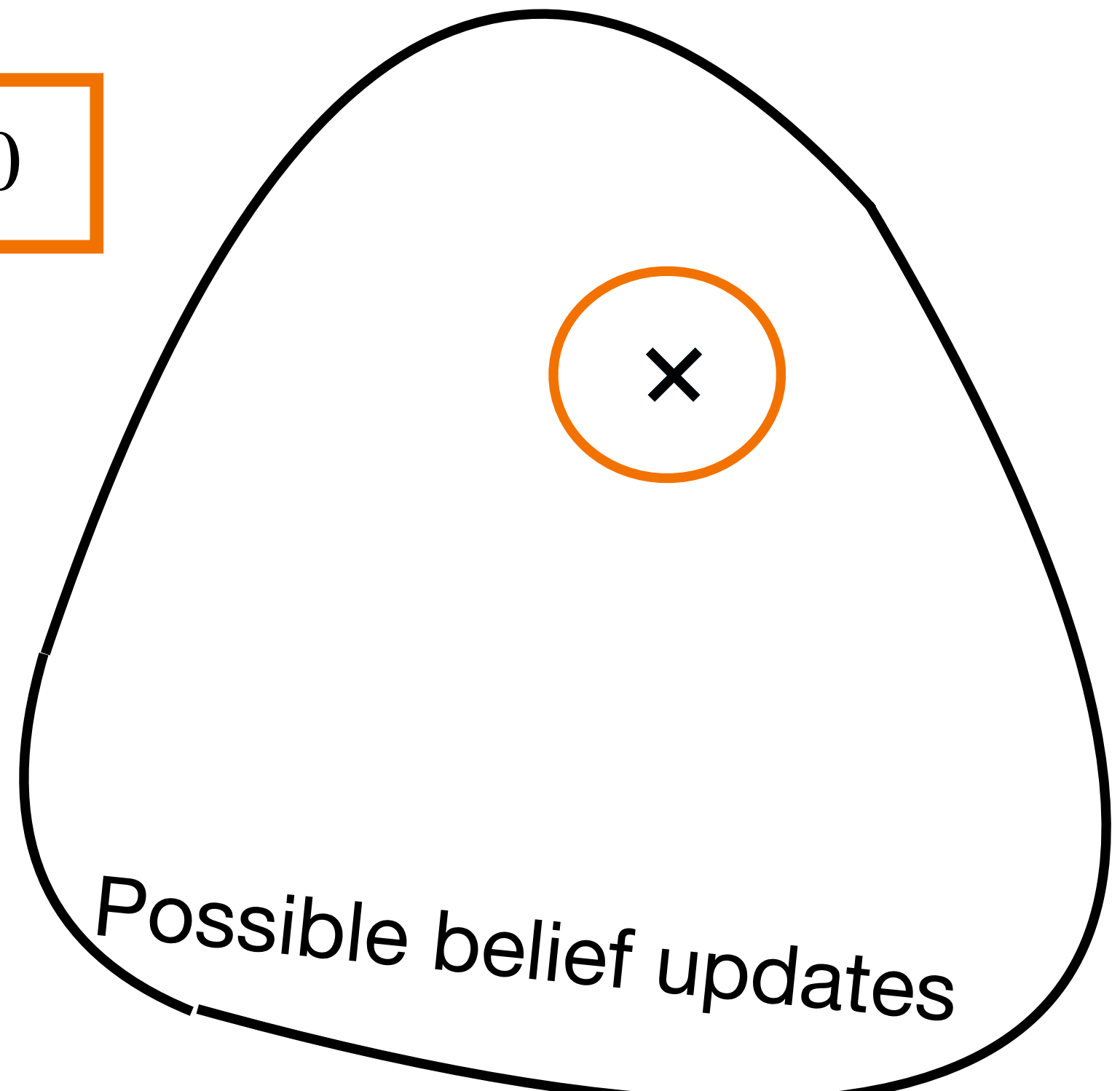
- (A1) model well-specified
- (A2) prior well-specified
- (A3) computationally feasible

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

Post-Bayesian ML

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

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



Power/Fractional/
Cold Posterior

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

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

Bayes' Posterior

(A1), (A2), (A3)

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



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

Q1: Can tuning λ improve Robustness?

(classical statistics)

Grünwald (2012); ALT
Holmes & Walker (2017); Biometrika
Miller & Dunson (2018); JRSS-B
Bhattacharya, Pati, & Yang (2019); Ann. Statist.
...



Frequent claim: λ can deliver robust predictions



(core ML)

Wenzel et al. (2020); ICML
Adlam et al. (2020); preprint
Noci et al. (2021); NeurIPS
Aitchison (2021); ICLR
...

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

Q1: Can tuning λ improve Robustness?

(classical statistics)

Grünwald (2012); ALT
Holmes & Walker (2017); Biometrika
Miller & Dunson (2018); JRSS-B
Bhattacharya, Pati, & Yang (2019); Ann. Statist.
...

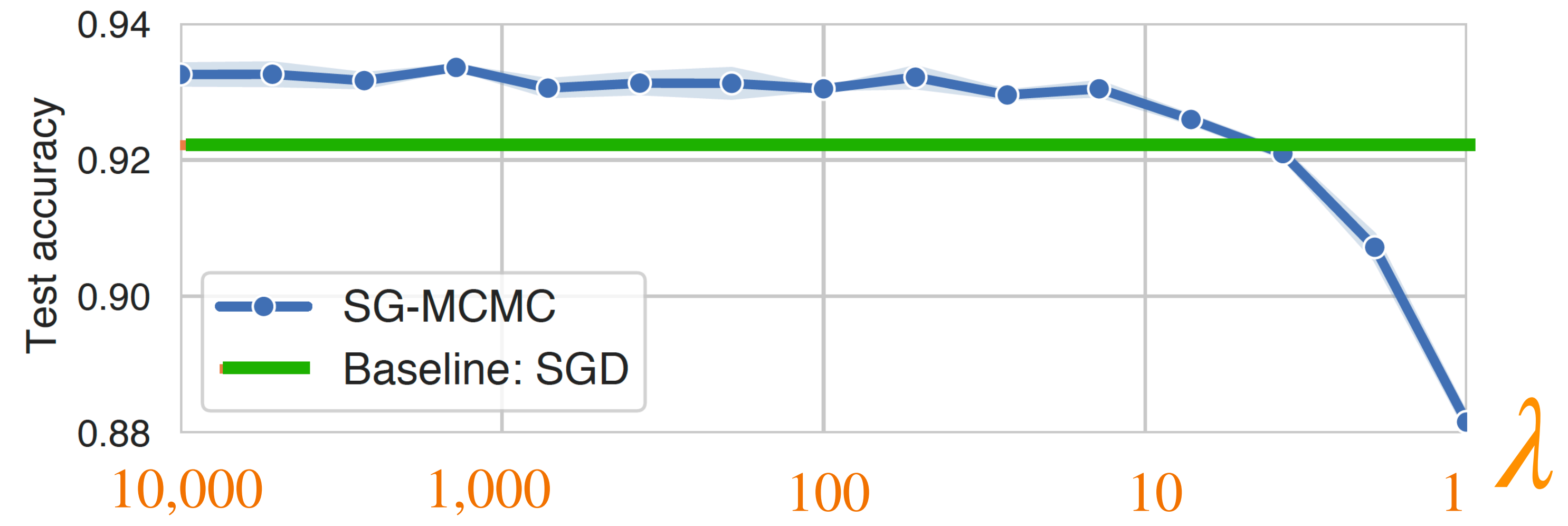
(core ML)

Wenzel et al. (2020); ICML
Adlam et al. (2020); preprint
Noci et al. (2021); NeurIPS
Aitchison (2021); ICLR
...

Frequent claim: λ can deliver robust predictions

How Good is the Bayes Posterior in Deep Neural Networks Really?

Florian Wenzel^{*1} Kevin Roth^{*+2} Bastiaan S. Veeling^{*+31} Jakub Świątkowski⁴⁺ Linh Tran⁵⁺
Stephan Mandt⁶⁺ Jasper Snoek¹ Tim Salimans¹ Rodolphe Jenatton¹ Sebastian Nowozin⁷⁺



Q1: Can tuning λ improve Robustness?

(classical statistics)

Grünwald (2012); ALT
Holmes & Walker (2017); Biometrika
Miller & Dunson (2018); JRSS-B
Bhattacharya, Pati, & Yang (2019); Ann. Statist.
...

(core ML)

Wenzel et al. (2020); ICML
Adlam et al. (2020); preprint
Noci et al. (2021); NeurIPS
Aitchison (2021); ICLR
...

Frequent claim: λ can deliver robust predictions

Unclear: Why should this be true?

Only regulates trade-off $\text{data} \leftrightarrow \text{prior}$

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

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

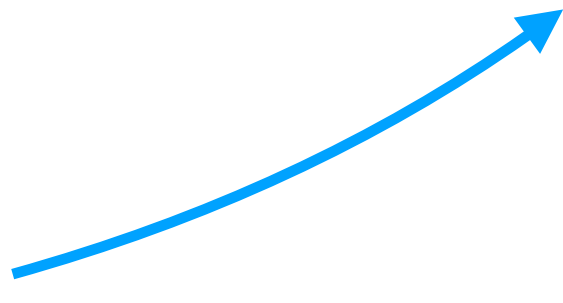
Q1: Can tuning λ improve Robustness?

Question: What is the predictively optimal λ ?

Posterior predictive = $p_n^\lambda(z) = \int p(z | \theta) \pi_n^{(\lambda)}(\theta | x_{1:n}) d\theta$

Predictively optimal λ : $\lambda^* = \operatorname{argmin}_{\lambda > 0} D_{\text{TV}}(q, p_n^\lambda)$

Data-generating density: $x_{1:n} \sim q(x_{1:n})$



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

Q1: Can tuning λ improve Robustness?

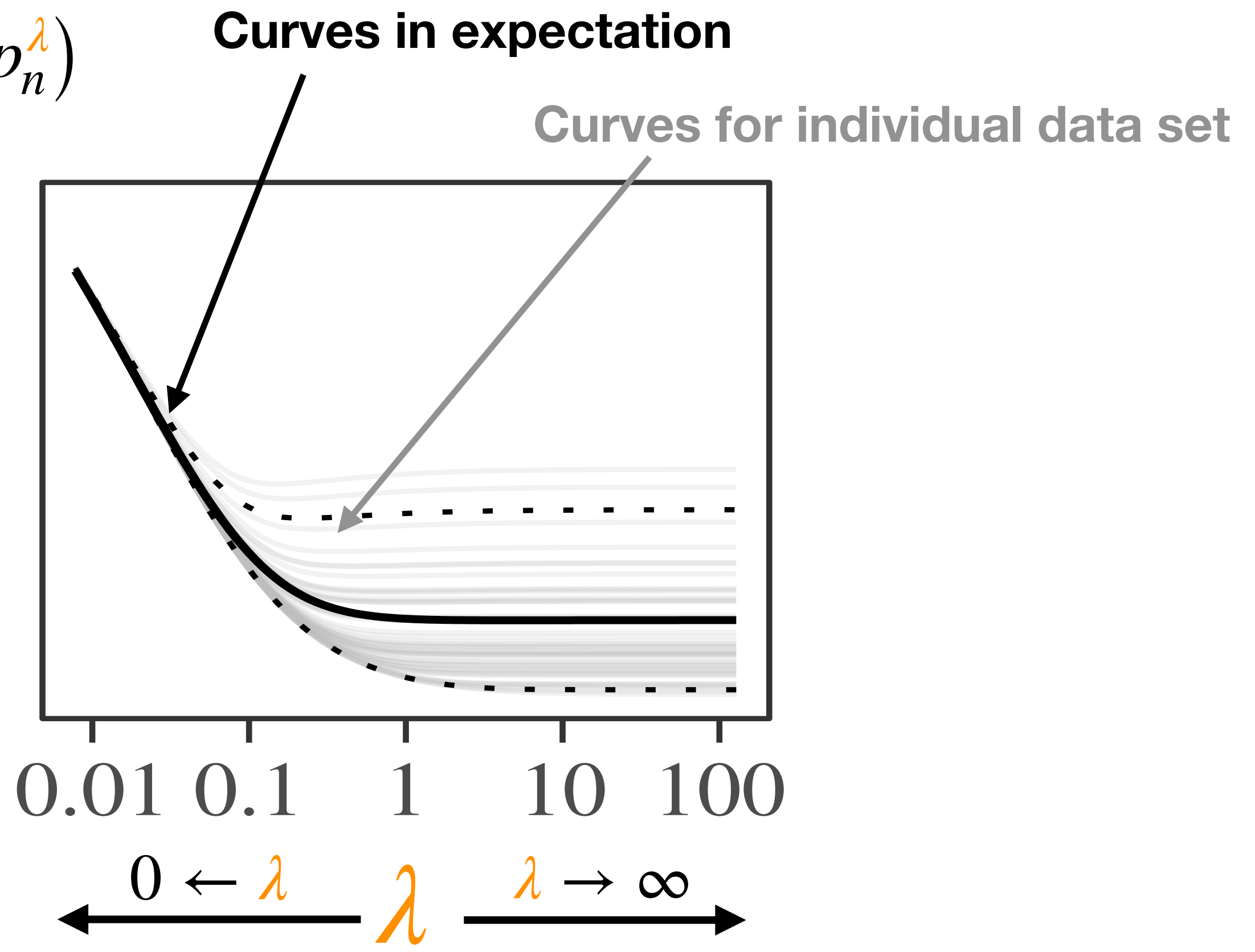
Question: What is the predictively optimal λ ?

Posterior predictive = $p_n^\lambda(z) = \int p(z | \theta) \pi_n^{(\lambda)}(\theta | x_{1:n}) d\theta$

Predictively optimal λ : $\lambda^* = \operatorname{argmin}_{\lambda > 0} D_{\text{TV}}(q, p_n^\lambda)$

Data-generating density: $x_{1:n} \sim q(x_{1:n})$

$D_{\text{TV}}(q, p_n^\lambda)$



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

Q1: Can tuning λ improve Robustness?

Question: What is the predictively optimal λ ?

Posterior predictive = $p_n^\lambda(z) = \int p(z | \theta) \pi_n^{(\lambda)}(\theta | x_{1:n}) d\theta$

Predictively optimal λ : $\lambda^* = \operatorname{argmin}_{\lambda > 0} D_{\text{TV}}(q, p_n^\lambda)$

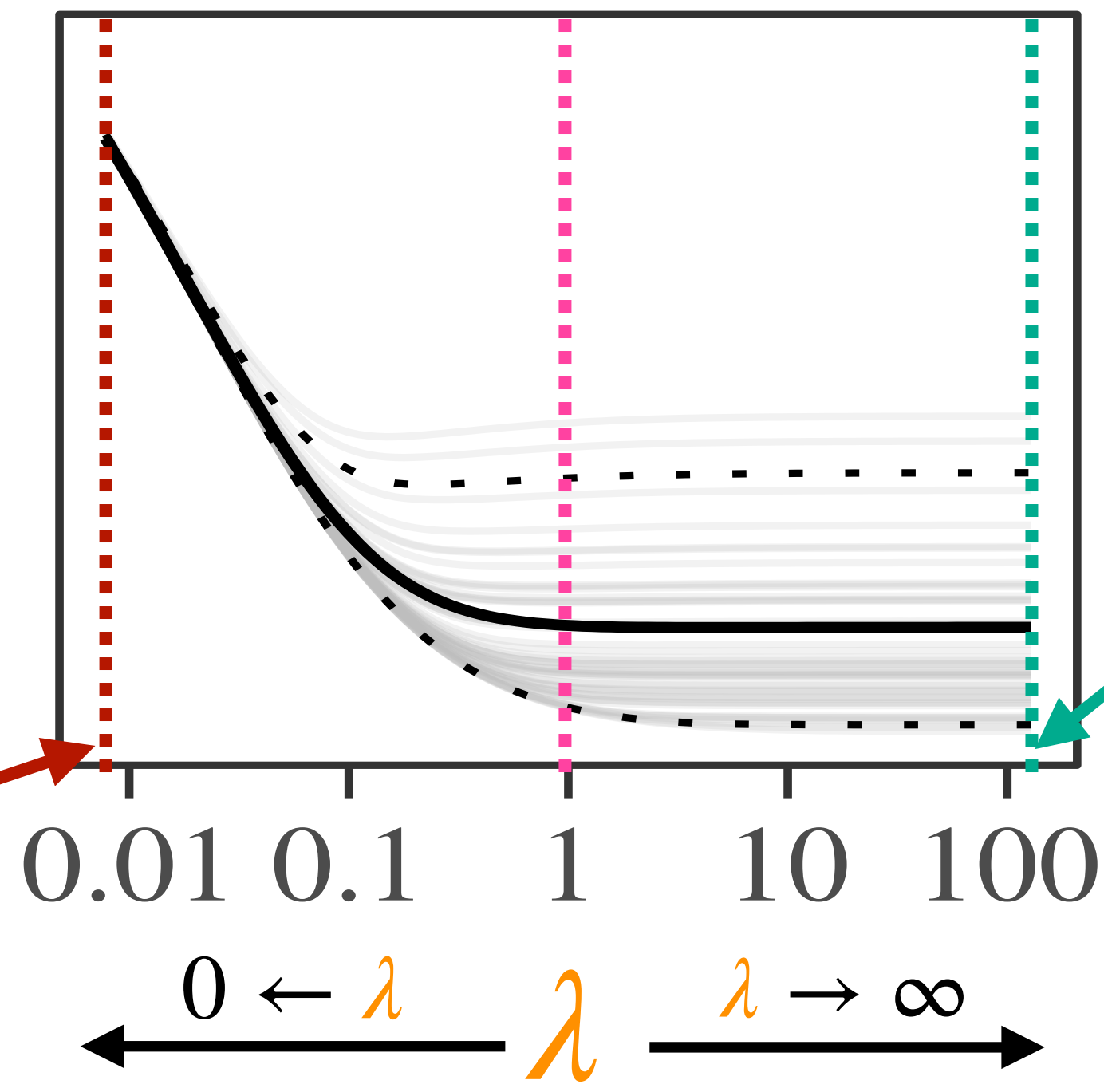
Data-generating density: $x_{1:n} \sim q(x_{1:n})$

$D_{\text{TV}}(q, p_n^\lambda)$

Prior predictive

Posterior predictive

\approx Plug-in predictive



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

Q1: Can tuning λ improve Robustness?

Question: What is the predictively optimal λ ?

Posterior predictive = $p_n^\lambda(z) = \int p(z | \theta) \pi_n^{(\lambda)}(\theta | x_{1:n}) d\theta$

Predictively optimal λ : $\lambda^* = \operatorname{argmin}_{\lambda > 0} D_{\text{TV}}(q, p_n^\lambda)$

Data-generating density: $x_{1:n} \sim q(x_{1:n})$

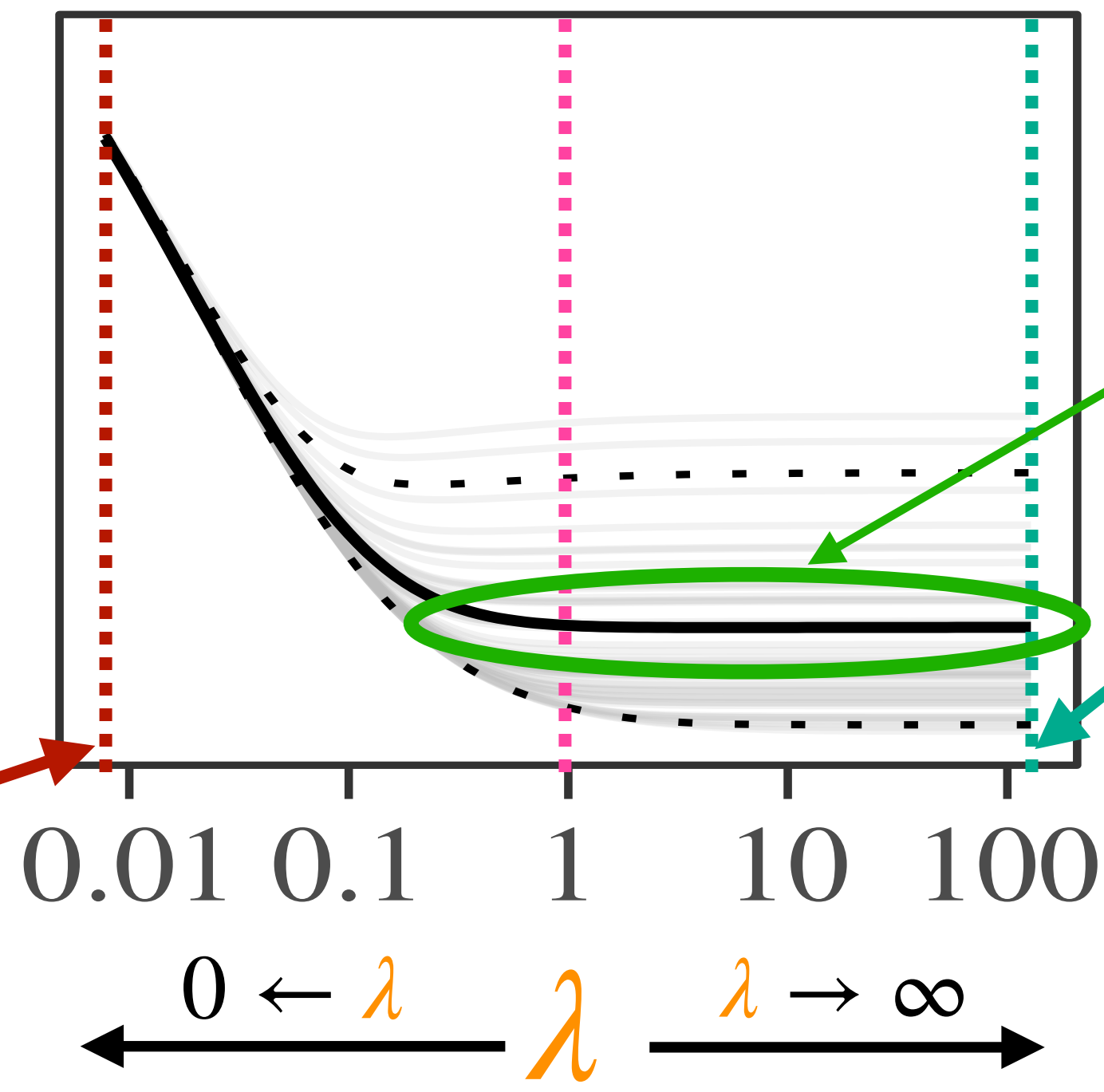
$D_{\text{TV}}(q, p_n^\lambda)$

Prior predictive

Posterior predictive

Flat region: all λ equally good for prediction

\approx Plug-in predictive



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

Q1: Can tuning λ improve Robustness?

Question: What is the predictively optimal λ ?

Posterior predictive = $p_n^\lambda(z) = \int p(z | \theta) \pi_n^{(\lambda)}(\theta | x_{1:n}) d\theta$

Predictively optimal λ : $\lambda^* = \operatorname{argmin}_{\lambda > 0} D_{\text{TV}}(q, p_n^\lambda)$

Data-generating density: $x_{1:n} \sim q(x_{1:n})$

$D_{\text{TV}}(q, p_n^\lambda)$

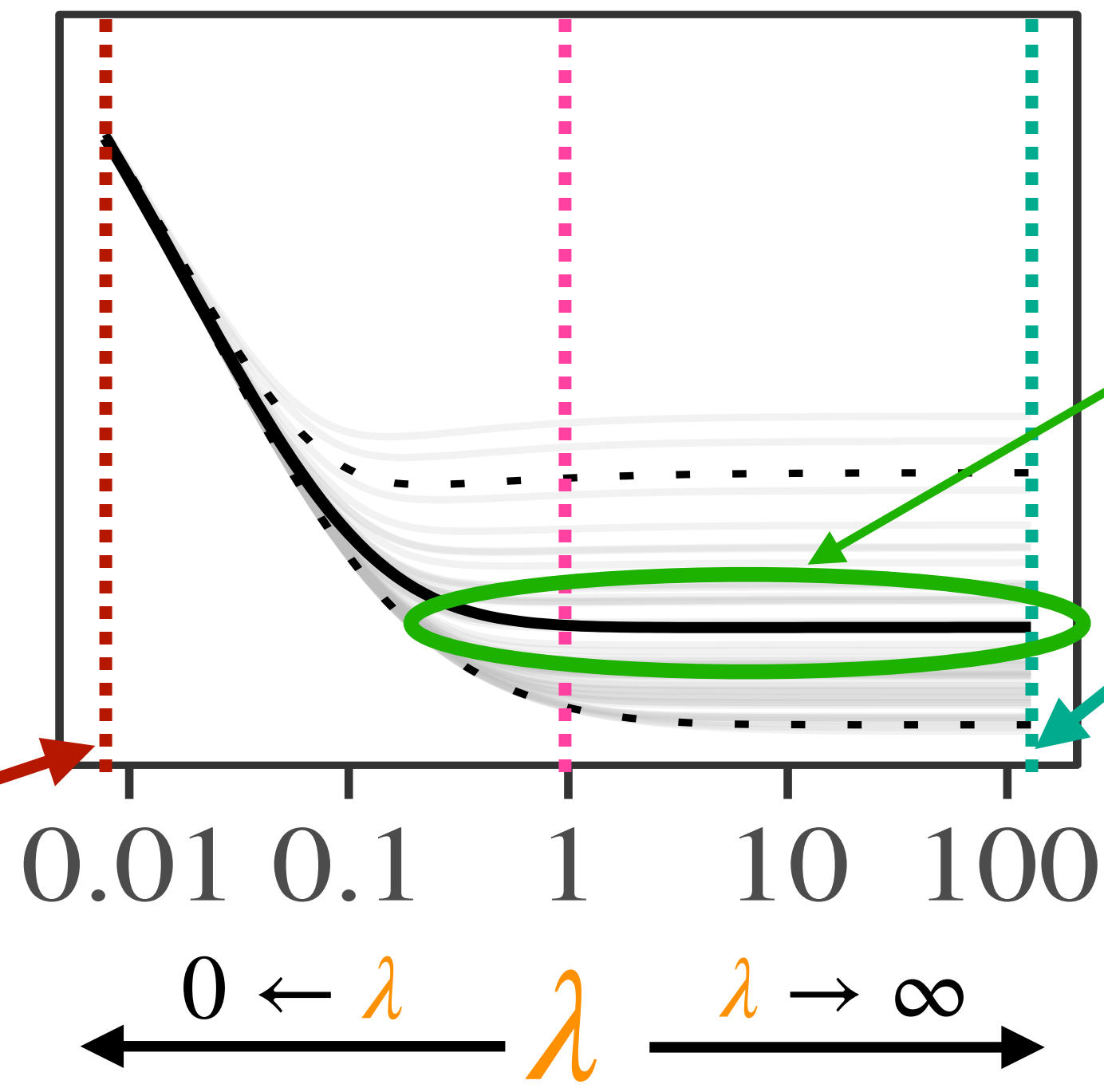
Theorem: these curves will always look that way.

Prior predictive

Posterior predictive

Flat region: all λ equally good for prediction

\approx Plug-in predictive



Q1: Can tuning λ improve Robustness?

Findings:

- (1) Ill-defined problem: minimiser λ^* doesn't exist
- (2) Flat region: infinitely λ yield (nearly) same predictive
- (3) Predictively no advantage over MLE / point estimators

Q1: Can tuning λ improve Robustness?

Findings:

- (1) Ill-defined problem: minimiser λ^* doesn't exist
- (2) Flat region: infinitely λ yield (nearly) same predictive
- (3) Predictively no advantage over MLE / point estimators

Conclusion: Normally, λ barely has an effect on robustness of predictions.

Q1: Can tuning λ improve Robustness?

Findings:

- (1) Ill-defined problem: minimiser λ^* doesn't exist
- (2) Flat region: infinitely λ yield (nearly) same predictive
- (3) Predictively no advantage over MLE / point estimators

Conclusion: Normally, λ barely has an effect on robustness of predictions.

Reason: As n grows, you almost predict from the plug-in predictive: $p_n^\infty \approx p_n^\lambda$

Q1: Can tuning λ improve Robustness?

Findings:

- (1) Ill-defined problem: minimiser λ^* doesn't exist
- (2) Flat region: infinitely λ yield (nearly) same predictive
- (3) Predictively no advantage over MLE / point estimators

Conclusion: Normally, λ barely has an effect on robustness of predictions.

Reason: As n grows, you almost predict from the plug-in predictive: $p_n^\infty \approx p_n^\lambda$

Cold Posterior Effect: not only deep learning; more general phenomenon

Q1: Can tuning λ improve Robustness?

Findings:

- (1) Ill-defined problem: minimiser λ^* doesn't exist
- (2) Flat region: infinitely λ yield (nearly) same predictive
- (3) Predictively no advantage over MLE / point estimators

Conclusion: Normally, λ barely has an effect on robustness of predictions.

Reason: As n grows, you almost predict from the plug-in predictive: $p_n^\infty \approx p_n^\lambda$

Cold Posterior Effect: not only deep learning; more general phenomenon

Possible Solution:

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

Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

Q1: Can tuning λ improve robustness?

Q2: What L leads to robust posteriors π_n^L ?

Q3: How should we design/choose L ?

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

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

Post-Bayesian ML

Gibbs/Generalised/
Pseudo Posterior

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

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

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

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

Power/Fractional/
Cold Posterior

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

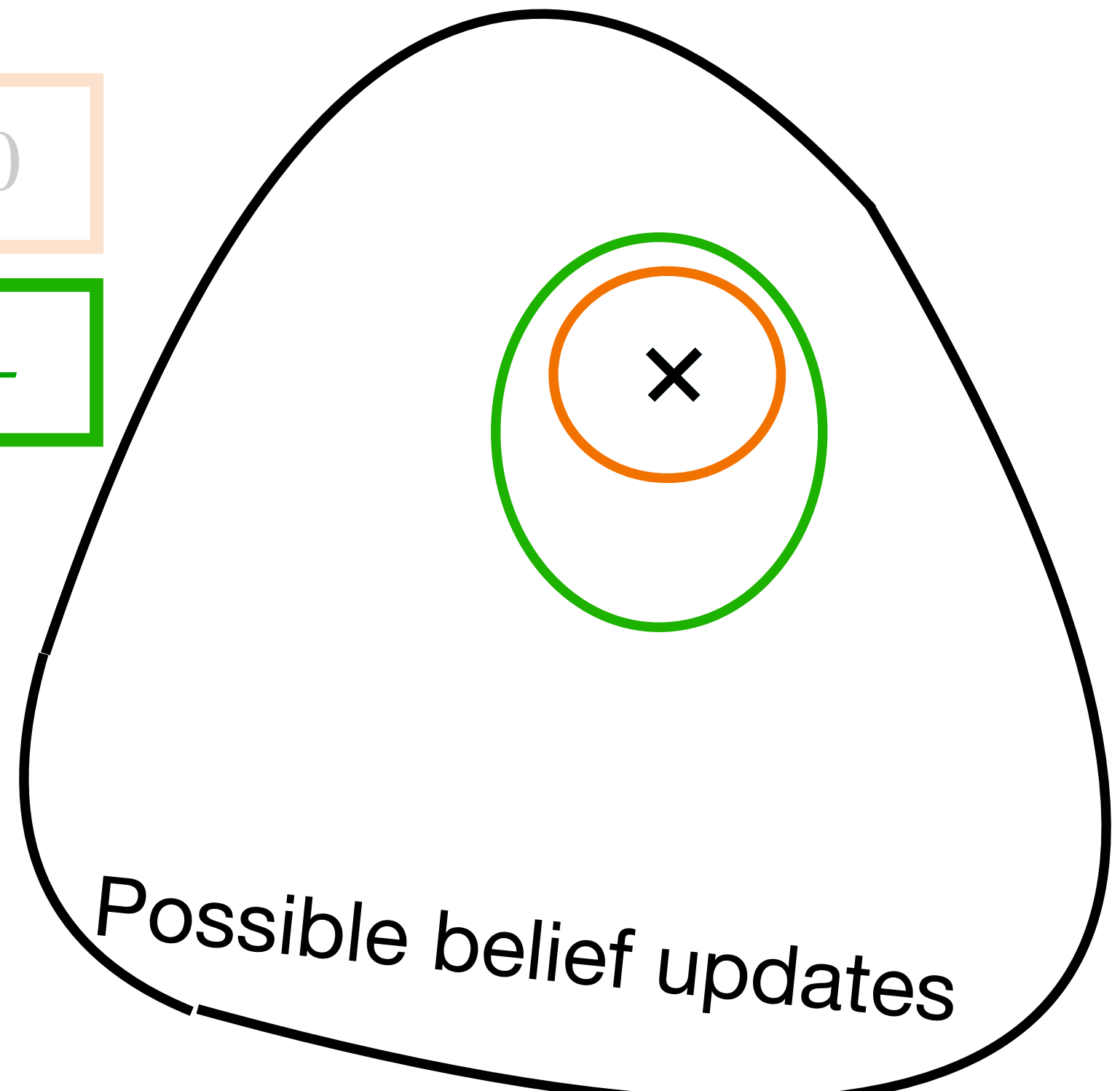
$$\pi_n^{(\lambda)}(\theta | x_{1:n}) = \frac{p(x_{1:n} | \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} | \theta)^\lambda \cdot \pi(\theta) d\theta}$$

Bayes' Posterior

(A1), (A2), (A3)

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

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



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

Q2: What **L** leads to robust posteriors?

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

Setting: for some small $\epsilon \geq 0$,

Data-generating probability distribution

ϵ -contamination distribution

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

Part of distribution our model captures

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

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

Q2: What **L** leads to robust posteriors?

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

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

Data-generating probability distribution

ε -contamination distribution

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

Part of distribution our model captures

What we want:

$$\left\{ \begin{array}{l} x_{1:n} \sim q_\varepsilon \longrightarrow \pi_n^L(\theta | x_{1:n}) \\ \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \approx \\ z_{1:n} \sim q_0 \longrightarrow \pi_n^L(\theta | z_{1:n}) \end{array} \right.$$

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

Q2: What L leads to robust posteriors?

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

Setting: for some small $\epsilon \geq 0$,

Data-generating probability distribution

ϵ -**contamination distribution**

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

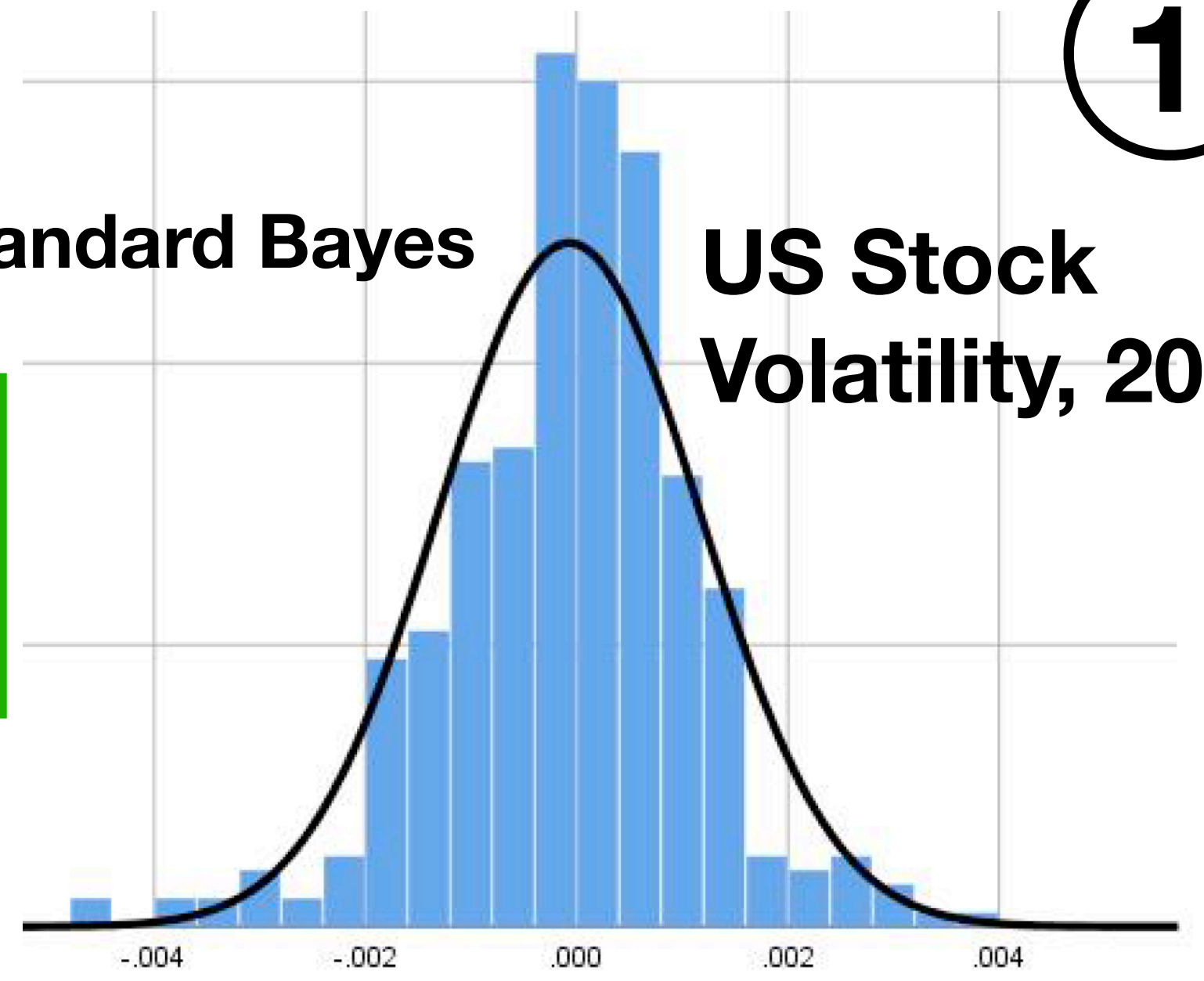
Part of distribution our model captures

What we want:

$$\left\{ \begin{array}{l} x_{1:n} \sim q_\epsilon \longrightarrow \pi_n^L(\theta | x_{1:n}) \\ z_{1:n} \sim q_0 \longrightarrow \pi_n^L(\theta | z_{1:n}) \end{array} \right. \approx$$

Standard Bayes

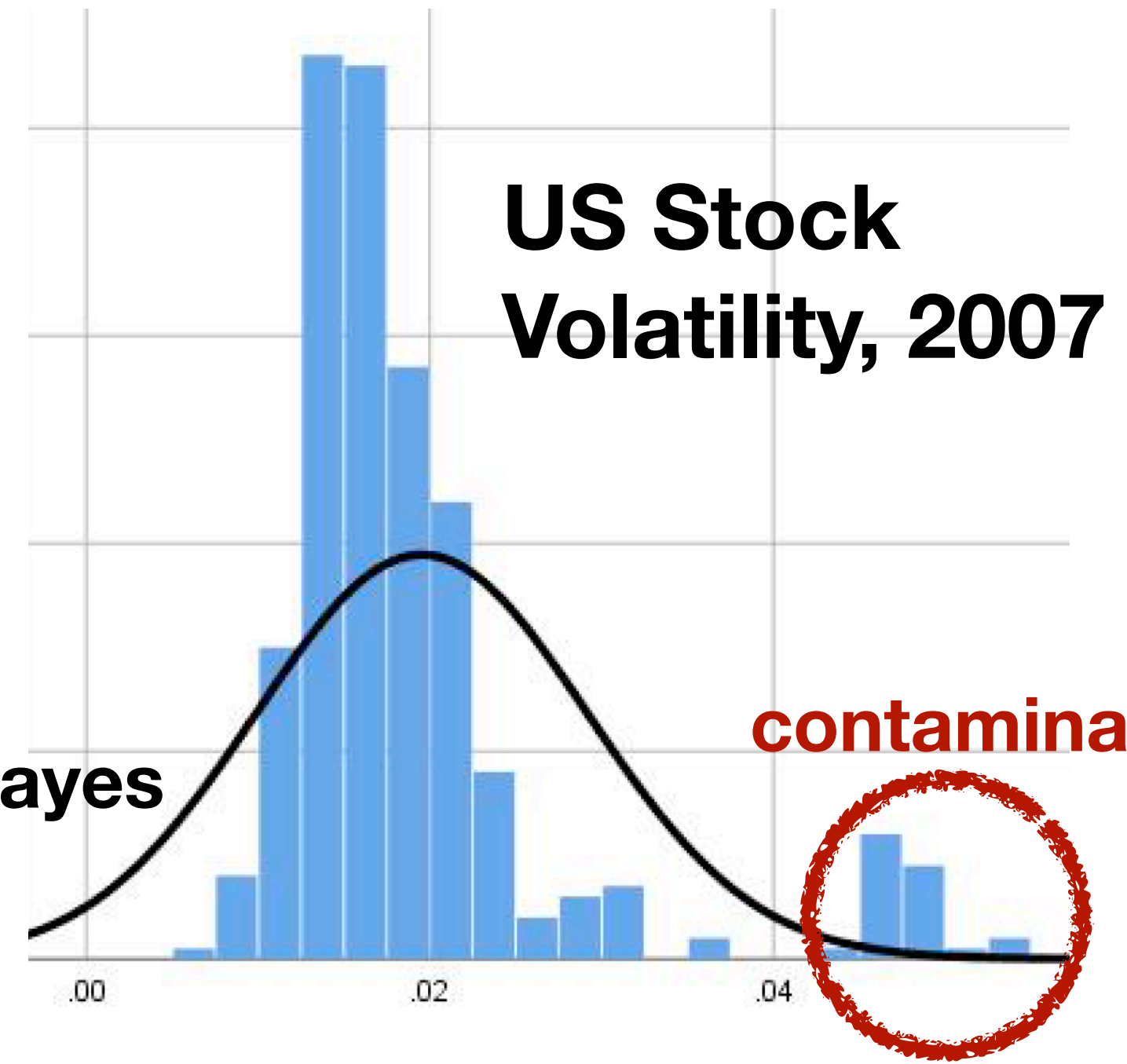
US Stock Volatility, 2004



US Stock Volatility, 2007

Standard Bayes

contamination



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

Q2: What L leads to robust posteriors?

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

Setting: for some small $\epsilon \geq 0$,

Data-generating probability distribution

ϵ -**contamination distribution**

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

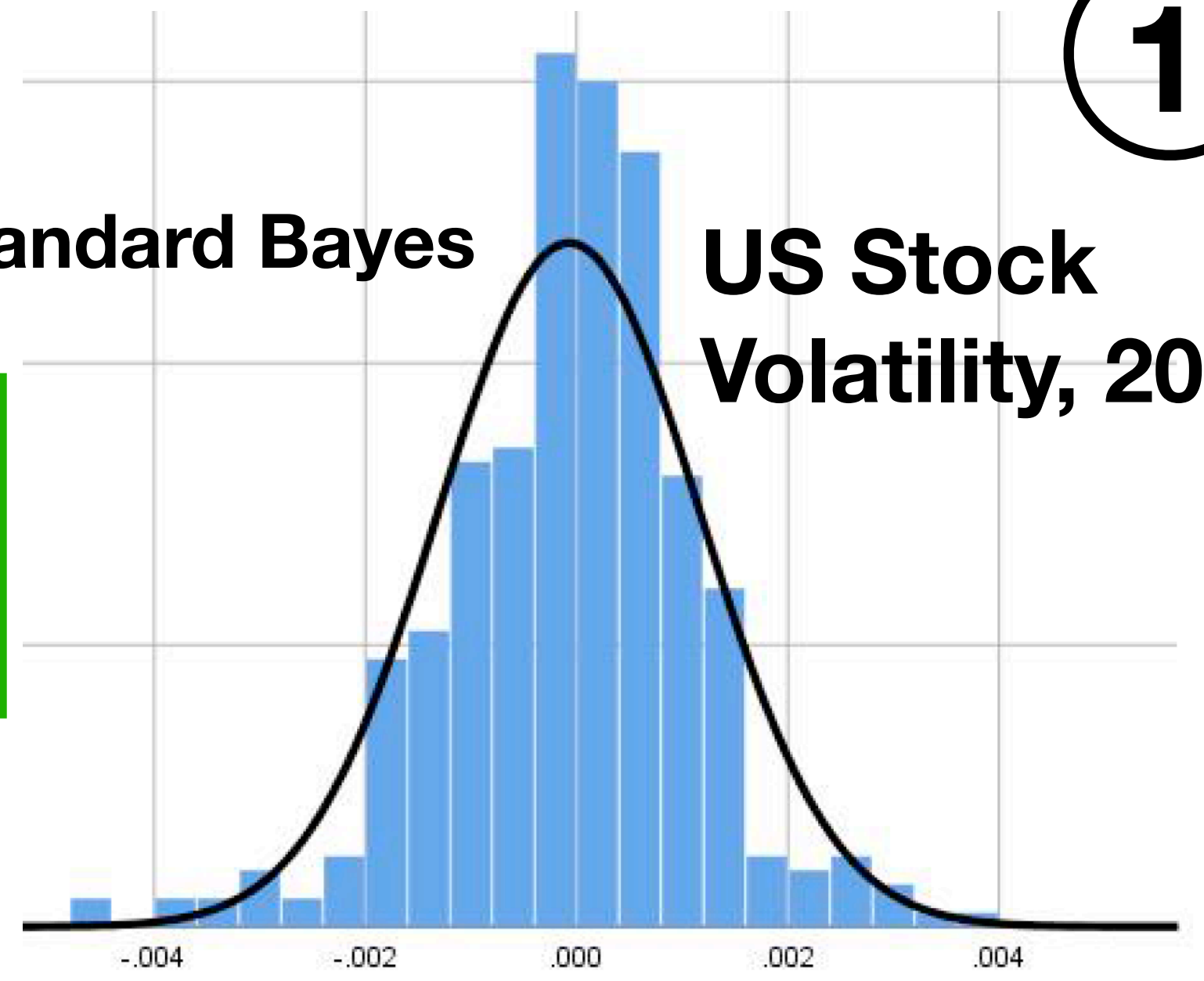
Part of distribution our model captures

What we want:

$$\begin{cases} x_{1:n} \sim q_\epsilon \longrightarrow \pi_n^L(\theta | x_{1:n}) \\ z_{1:n} \sim q_0 \longrightarrow \pi_n^L(\theta | z_{1:n}) \end{cases} \approx$$

Standard Bayes

US Stock Volatility, 2004

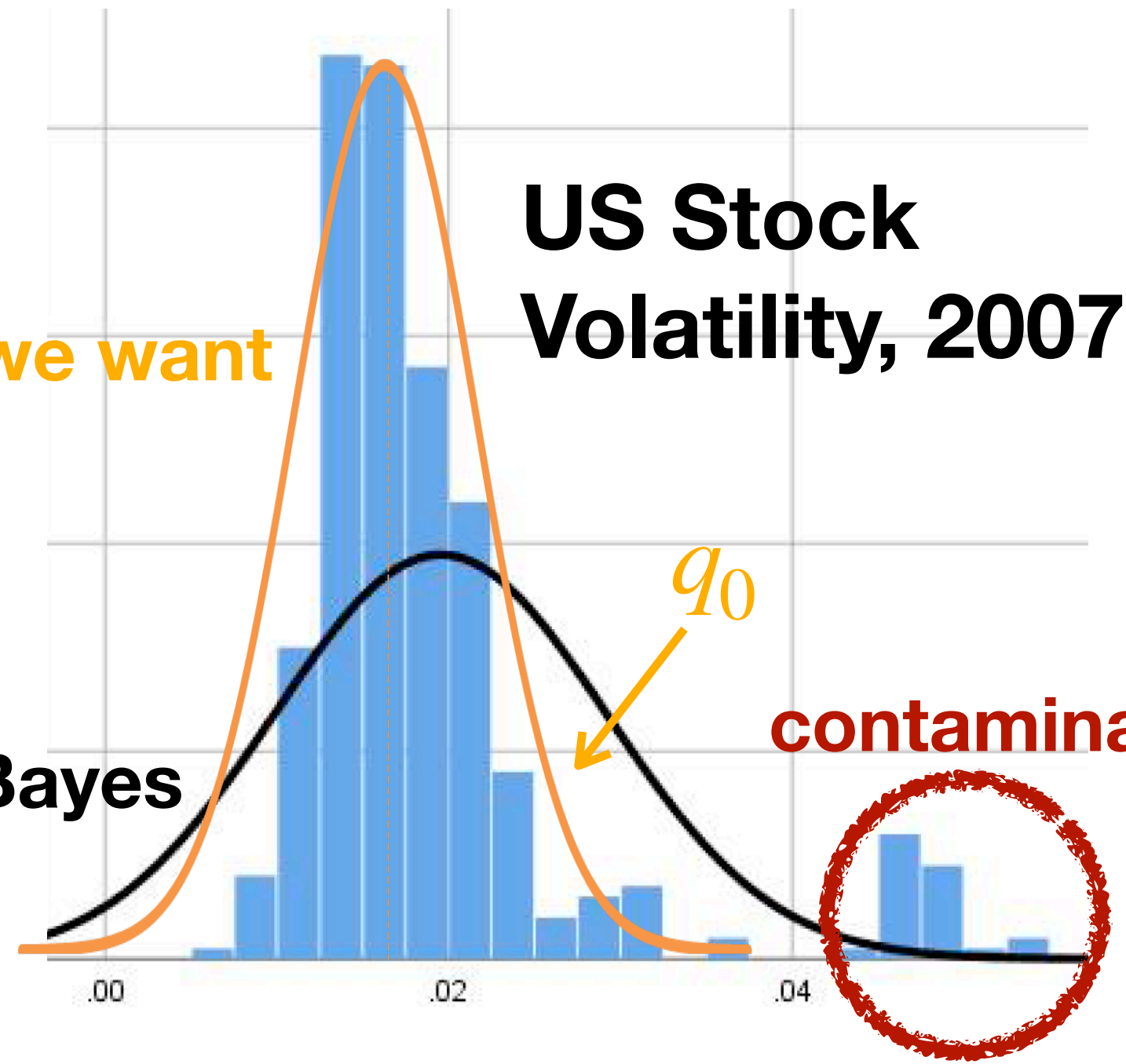


What we want

US Stock Volatility, 2007

Standard Bayes

contamination



Q2: What **L** leads to robust posteriors?

Setting:

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

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

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

Robustness:

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

Q2: What **L** leads to robust posteriors?

Setting:

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

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

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

Robustness:

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

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

Q2: What **L** leads to robust posteriors?

Setting:

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

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

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

Robustness:

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

Q2: What **L** leads to robust posteriors?

Setting:

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

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

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

Robustness:

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

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

Key quantity:

Loss robustness to **contamination**: $\frac{\partial}{\partial \varepsilon} L(\theta, x_{1:n}) \Big|_{\varepsilon=0}$

Q2: What **L** leads to robust posteriors?

Setting:

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

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

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

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

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

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

Key quantity:

Loss robustness to **contamination**: $\left. \frac{\partial}{\partial \varepsilon} L(\theta, x_{1:n}) \right|_{\varepsilon=0}$

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

Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

- Q1: Can tuning λ improve robustness?
- Q2: What L leads to robust posteriors π_n^L ?
- Q3: How should we design/choose L ?**

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

Q3: How should we choose **L**?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
Dewaskar, Tosh, **Knoblauch**, & Dunson (2023); preprint
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

Standard Bayes

↓

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

Q3: How should we choose **L**?

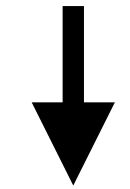
$$\pi_n^L(\theta | x_{1:n}) = \frac{\exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta)}{\int \exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta) d\theta}$$

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 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

$$n \cdot \text{KL}(q_\epsilon, p(\cdot | \theta))$$

$$x_i \sim q_\epsilon \approx$$

Standard Bayes



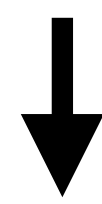
$$L(x_{1:n}, \theta) = \sum_{i=1}^n -\log p(x_i | \theta)$$

Q3: How should we choose L?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 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

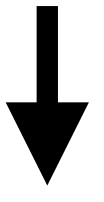
NOT robust to model misspecification ~~(A1)~~



$$n \cdot \text{KL}(q_\epsilon, p(\cdot | \theta))$$

$$x_i \sim q_\epsilon \approx$$

Standard Bayes



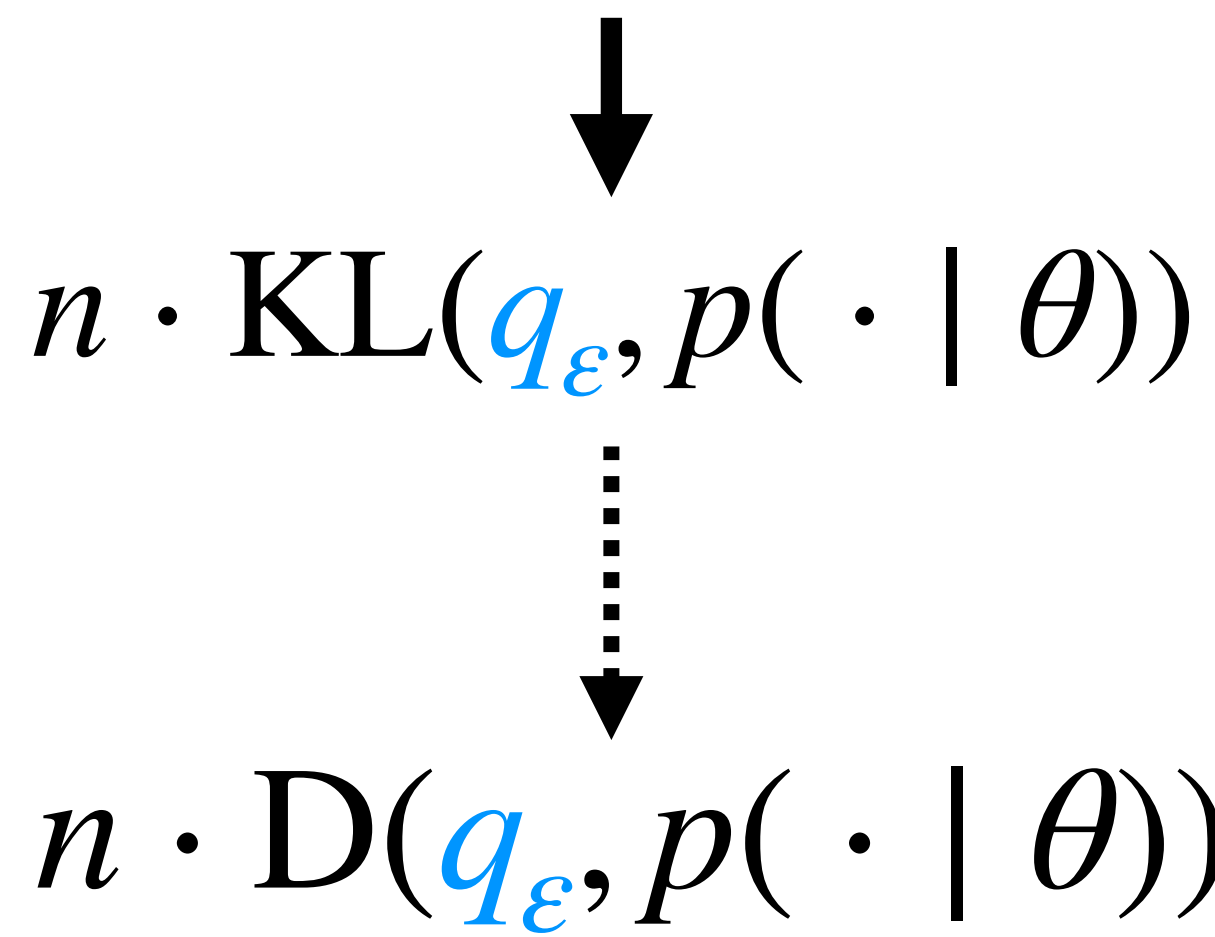
$$L(x_{1:n}, \theta) = \sum_{i=1}^n -\log p(x_i | \theta)$$

Q3: How should we choose L?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 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

NOT robust to model misspecification ~~(A1)~~



$$x_i \sim q_\epsilon \approx$$

Standard Bayes

↓

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

Robust discrepancy

$$D(q_\epsilon, p(\cdot | \theta)) \approx D(q_0, p(\cdot | \theta))$$

Q3: How should we choose L?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 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

NOT robust to model misspecification ~~(A1)~~

↓

$$n \cdot \text{KL}(q_\epsilon, p(\cdot | \theta))$$

$$x_i \sim q_\epsilon \approx$$

⋮

↓

$$n \cdot D(q_\epsilon, p(\cdot | \theta))$$

$$x_i \sim q_\epsilon \approx$$

Standard Bayes

↓

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

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

Robust discrepancy

$$D(q_\epsilon, p(\cdot | \theta)) \approx D(q_0, p(\cdot | \theta))$$

Q3: How should we choose **L**?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 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

NOT robust to model misspecification ~~(A1)~~

↓

$$n \cdot \text{KL}(q_\epsilon, p(\cdot | \theta))$$

$$x_i \sim q_\epsilon \approx$$

⋮

$$n \cdot D(q_\epsilon, p(\cdot | \theta))$$

$$x_i \sim q_\epsilon \approx$$

Standard Bayes

↓

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

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

Robust discrepancy

$$D(q_\epsilon, p(\cdot | \theta)) \approx D(q_0, p(\cdot | \theta))$$

Robust loss

⋮ → **L** is robust over all $c \in \mathcal{S}$

Q3: How should we choose **L**?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 Knoblauch*, Frazier*, & Drovandi (2024); preprint

Examples of this principle:



$$D(q_\epsilon, p(\cdot | \theta)) \approx D(q_0, p(\cdot | \theta)) \dashrightarrow L \text{ is robust over all } c \in \mathcal{S}$$

Q3: How should we choose **L**?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 Knoblauch*, Frazier*, & Drovandi (2024); preprint

Examples of this principle:

γ -Divergence	$x_i \sim q_\epsilon$	\approx	$L^\gamma(x_{1:n}, \theta)$	$\xrightarrow{\gamma \downarrow 0}$	$\sum_{i=1}^n -\log p(x_i \theta)$
β -Divergence	$x_i \sim q_\epsilon$	\approx	$L^\beta(x_{1:n}, \theta)$	$\xrightarrow{\beta \downarrow 0}$	$\sum_{i=1}^n -\log p(x_i \theta)$

Robust discrepancy

Robust loss

$$D(q_\epsilon, p(\cdot | \theta)) \approx D(q_0, p(\cdot | \theta)) \dashrightarrow L \text{ is robust over all } c \in \mathcal{S}$$

Q3: How should we choose L ?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 Knoblauch*, Frazier*, & Drovandi (2024); preprint

Examples of this principle:

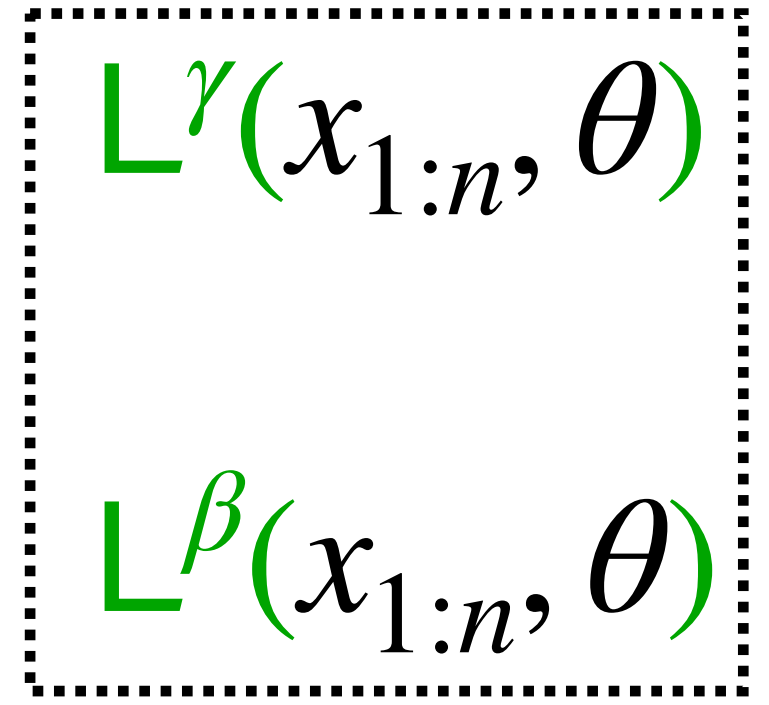
γ -Divergence

$$x_i \sim q_\epsilon \approx x_i \sim q_\epsilon$$

β -Divergence

$$x_i \sim q_\epsilon \approx x_i \sim q_\epsilon$$

For small γ / β ,
 $\pi_n^L(\theta | x_{1:n}) \approx$ Bayes w/o outliers



$$\xrightarrow{\gamma \downarrow 0} \sum_{i=1}^n -\log p(x_i | \theta)$$

$$\xrightarrow{\beta \downarrow 0} \sum_{i=1}^n -\log p(x_i | \theta)$$

Robust discrepancy

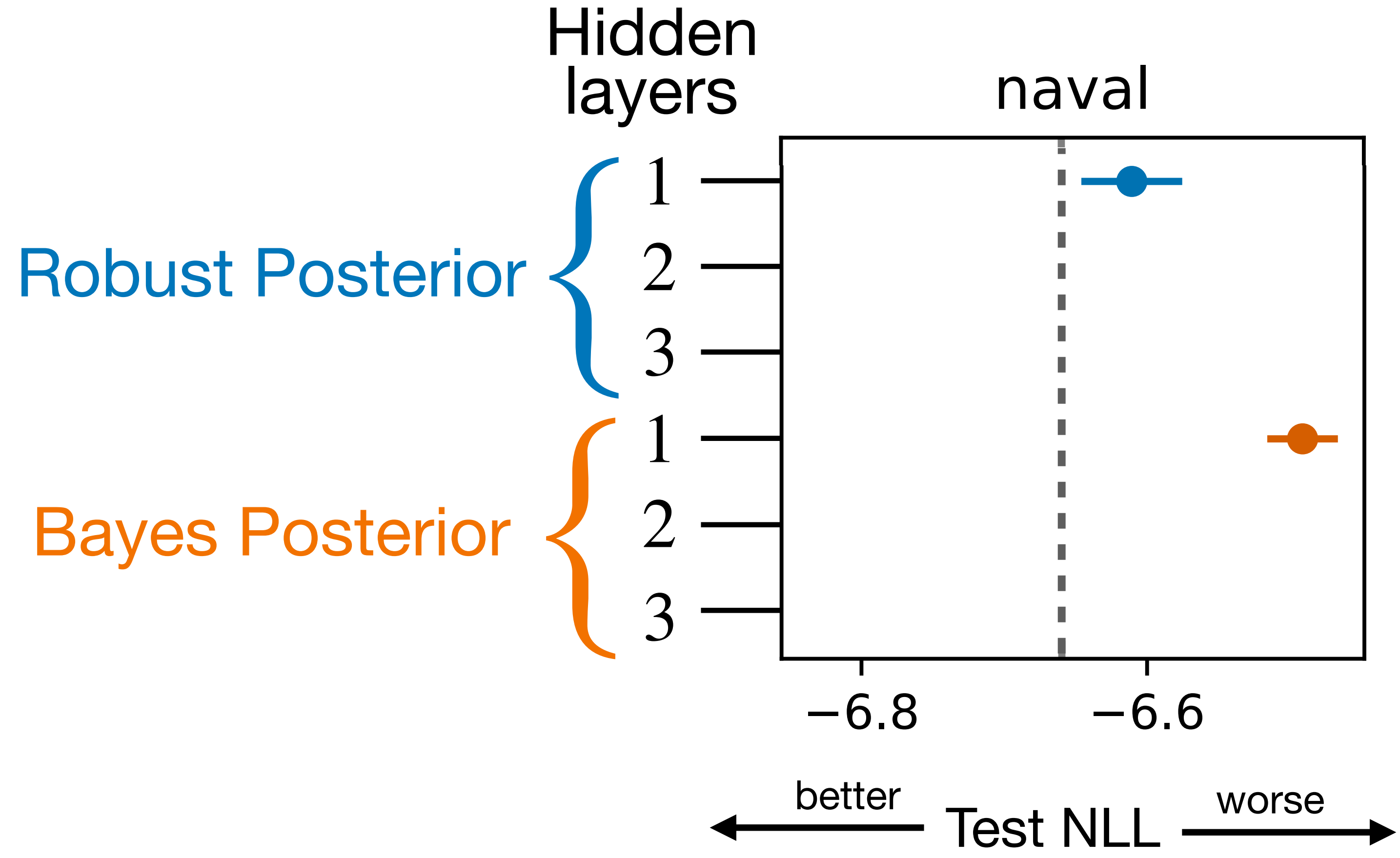
$$D(q_\epsilon, p(\cdot | \theta)) \approx D(q_0, p(\cdot | \theta))$$

Robust loss

L is robust over all $c \in \mathcal{S}$

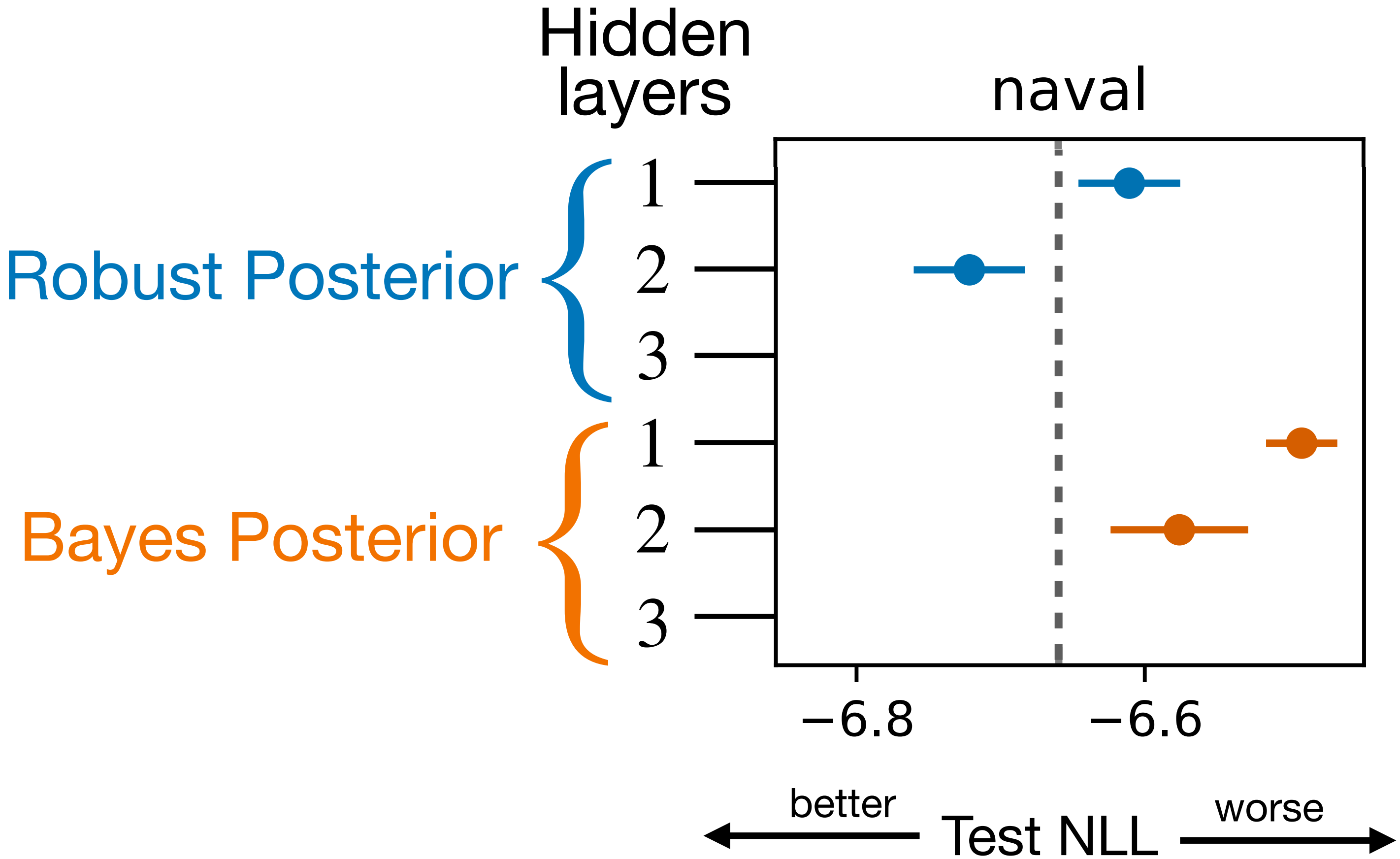
Q3: How should we choose L?

Example 1: Deep Gaussian Processes (γ -Divergence)



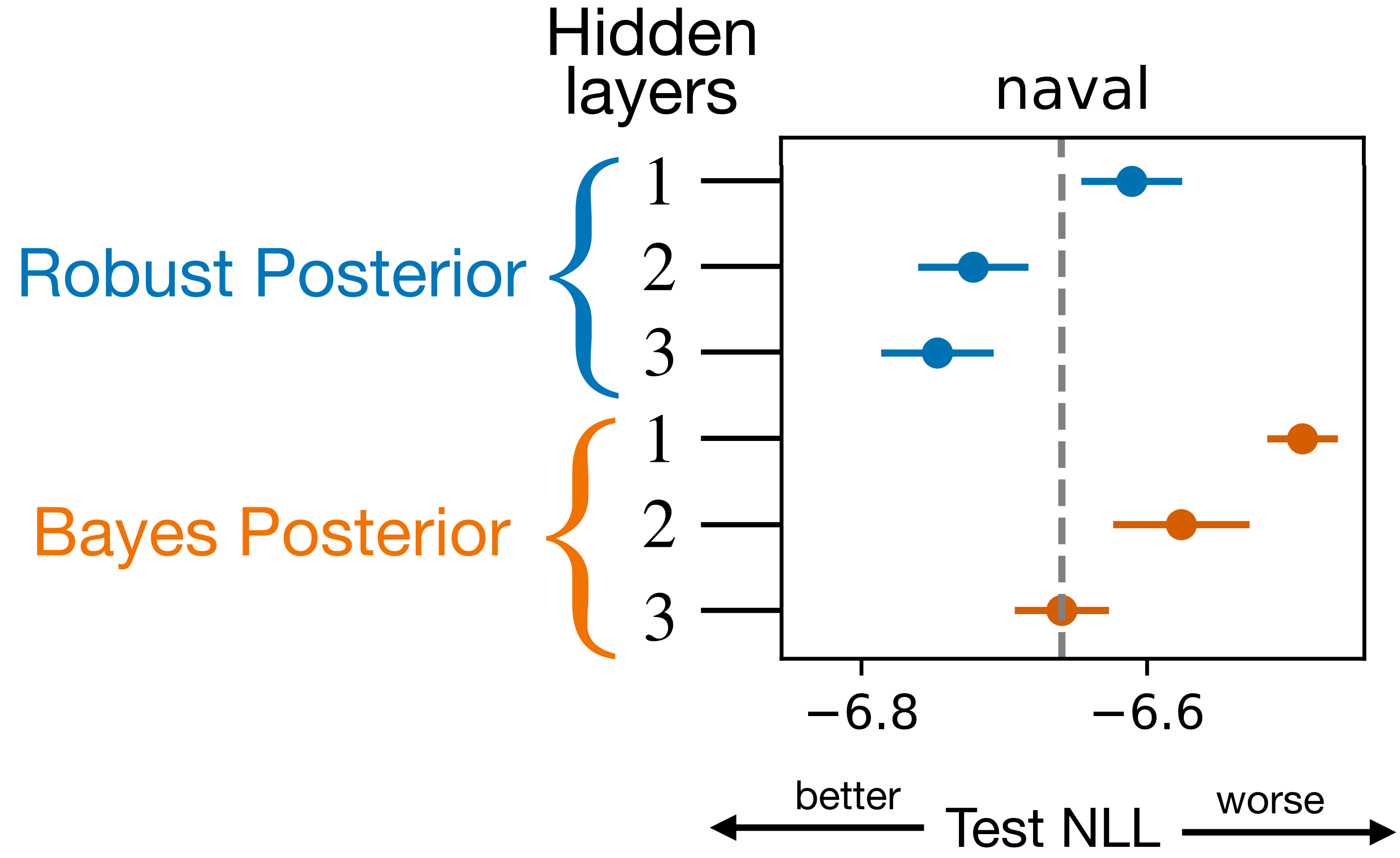
Q3: How should we choose L?

Example 1: Deep Gaussian Processes (γ -Divergence)



Q3: How should we choose L?

Example 1: Deep Gaussian Processes (γ -Divergence)



Q3: How should we choose **L**?

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

Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Dewaskar, Tosh, Knoblauch, & Dunson (2023); preprint
 Knoblauch*, Frazier*, & Drovandi (2024); preprint

Problems with these losses:

Generally intractable integrals
 analytic form for most exponential families
 otherwise: $\approx \frac{1}{S} \sum_{j=1}^S p(u_j | \theta)^\beta, u_j \sim p(u_j | \theta)$

γ -Divergence

$$x_i \sim q_\epsilon \approx$$

$$L^\gamma(x_{1:n}, \theta)$$

$$= n \cdot \int p(u | \theta)^{1+\beta} du + \dots$$

β -Divergence

$$x_i \sim q_\epsilon \approx$$

$$L^\beta(x_{1:n}, \theta)$$

$$= n \cdot \log \int p(u | \theta)^{1+\gamma} du + \dots$$

Robust discrepancy

$$(D(q_\epsilon, p(\cdot | \theta)) \approx D(q_0, p(\cdot | \theta)))$$

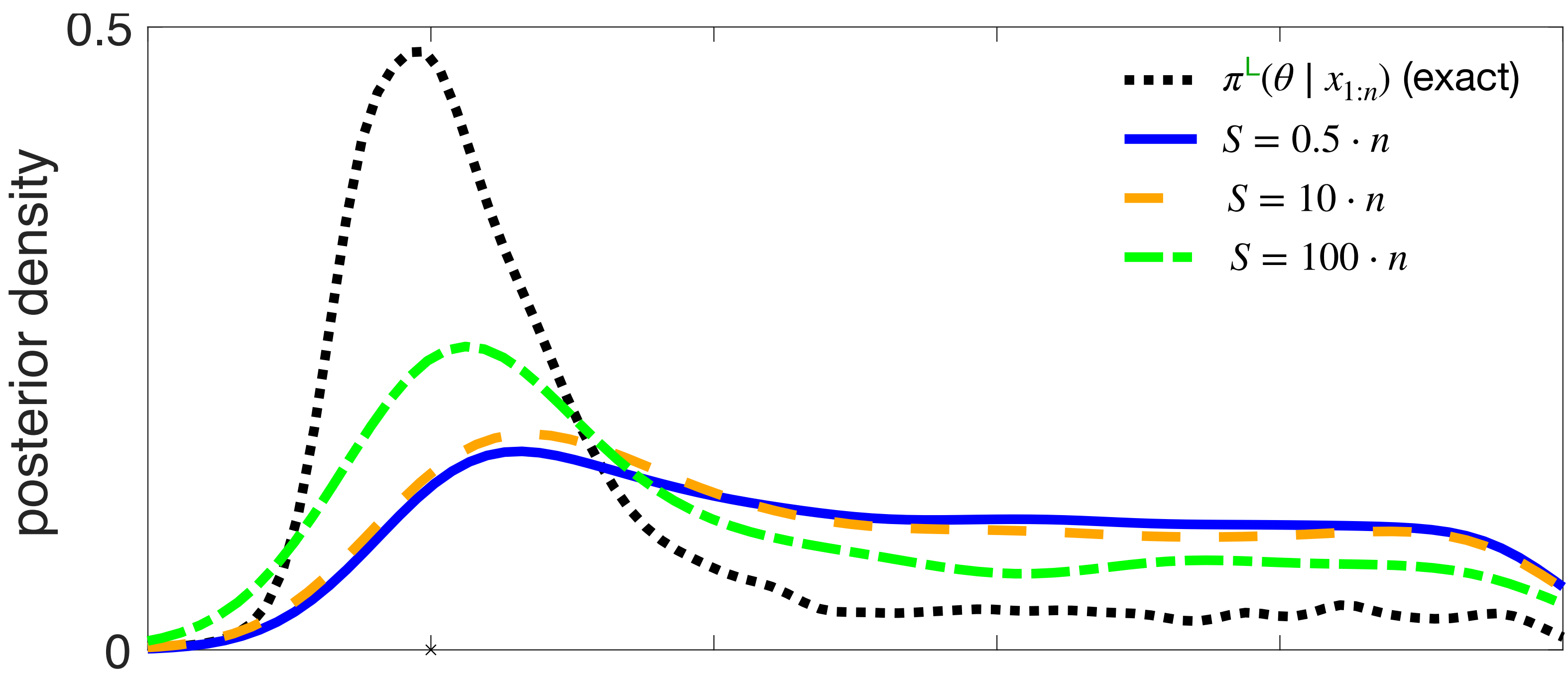
Robust loss

L is robust over all $c \in \mathcal{S}$

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

Q3: How should we choose **L**?

Question: How many simulations for good approximation quality of $\pi_n^L(\theta | x_{1:n})$?

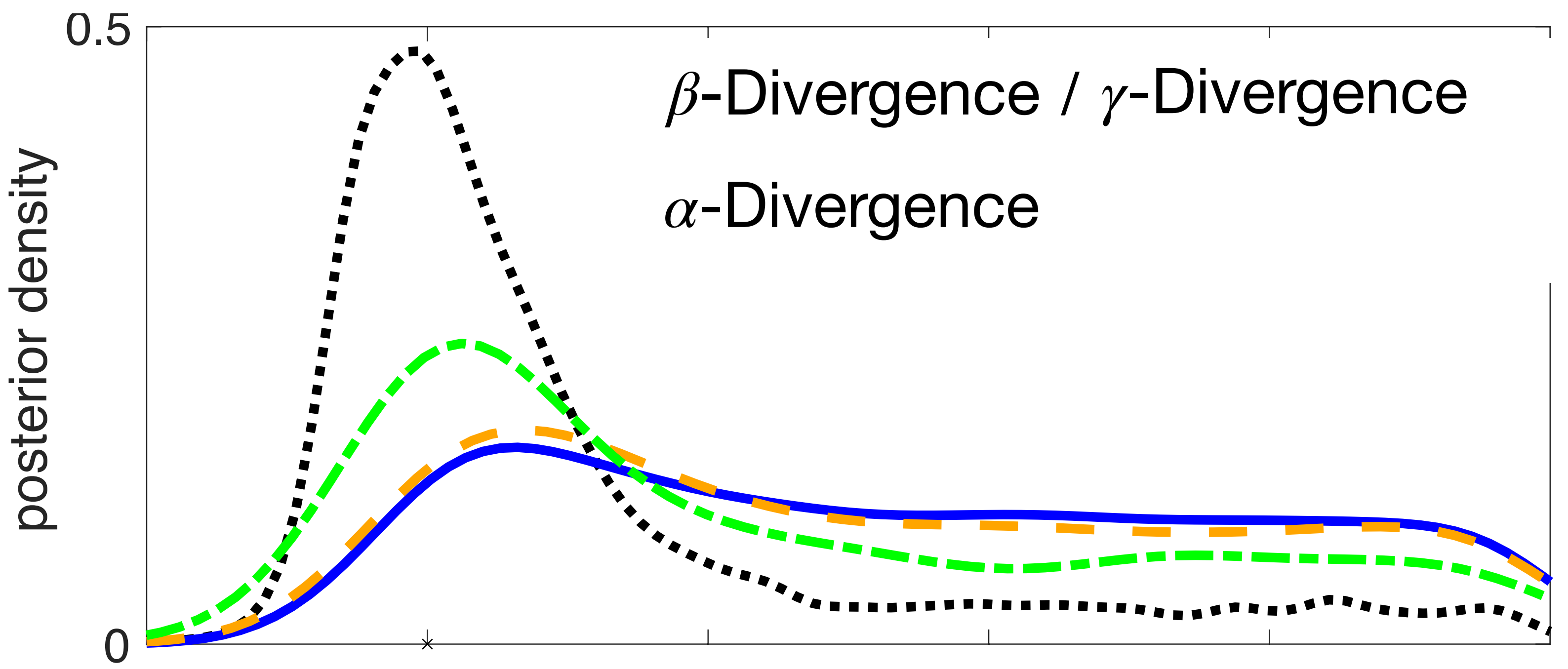


Terrible approximations, even for large S ! **X**

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

Q3: How should we choose L?

Question: How many simulations for good approximation quality of $\pi_n^L(\theta | x_{1:n})$?



$S = O(n^{2.5})$ ✘

$S = O(n^{2.25})$ ✘

Terrible approximations, even for large S ! ✘

Summary: Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$



Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

Q1: Can tuning λ improve robustness?	A: No. Work with L
Q2: What L leads to robust posteriors π_n^L ?	A: robust L \implies robust $\pi_n^L(\theta x_{1:n})$
Q3: How should we design/choose L ?	A: L = (robust) divergence

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

Summary: Foundations of Post-Bayesian ML Research

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$



Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

Q1: Can tuning λ improve robustness?	A: No. Work with L
Q2: What L leads to robust posteriors π_n^L ?	A: robust L \implies robust $\pi_n^L(\theta x_{1:n})$
Q3: How should we design/choose L ?	A: L = (robust) divergence

Robustness \iff **Tractability**

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

Post-Bayesian ML Research: State of the Art

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

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

Knoblauch & Damoulas (2018); CML
 Knoblauch, Jewson, & Damoulas (2022); NIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification + computation

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

Schmon, Cannon, & Knoblauch (2020); AABI
 Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
 Dellaporta, Knoblauch, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & Knoblauch (2023); ICML
 Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch, Briol, & Murphy (2024); ICML

3 The Future

$$q_n^*(\theta)$$

model misspecification + prior misspecification + computation

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

Husain & Knoblauch (2022); ALT
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Matsubara, Knoblauch, Briol, & Oates (2023); JASA
 Wild, Sejdinovic, & Knoblauch (2024); forthcoming
 Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2024); NeurIPS (oral)

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

Post-Bayesian ML Research: State of the Art

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

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification + computation

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

Schmon, Cannon, & **Knoblauch** (2020); AABI
 Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
 Dellaporta, **Knoblauch**, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & **Knoblauch** (2023); ICML
 Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, **Knoblauch**,
 Briol, & Murphy (2024); ICML

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

Post-Bayesian ML Research: State of the Art

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

2 State of the Art

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

**model misspecification +
computation**

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

Schmon, Cannon, & **Knoblauch** (2020); AABI
 Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
 Dellaporta, **Knoblauch**, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & **Knoblauch** (2023); ICML
 Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, **Knoblauch**,
 Briol, & Murphy (2024); ICML

Problem identified in 1 : Robustness ↔ Tractability

Q: Can we design losses L that are **both robust and tractable** ?

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

2

Post-Bayesian ML

Gibbs/Generalised/
Pseudo Posterior

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

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

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

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

Power/Fractional/
Cold Posterior

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

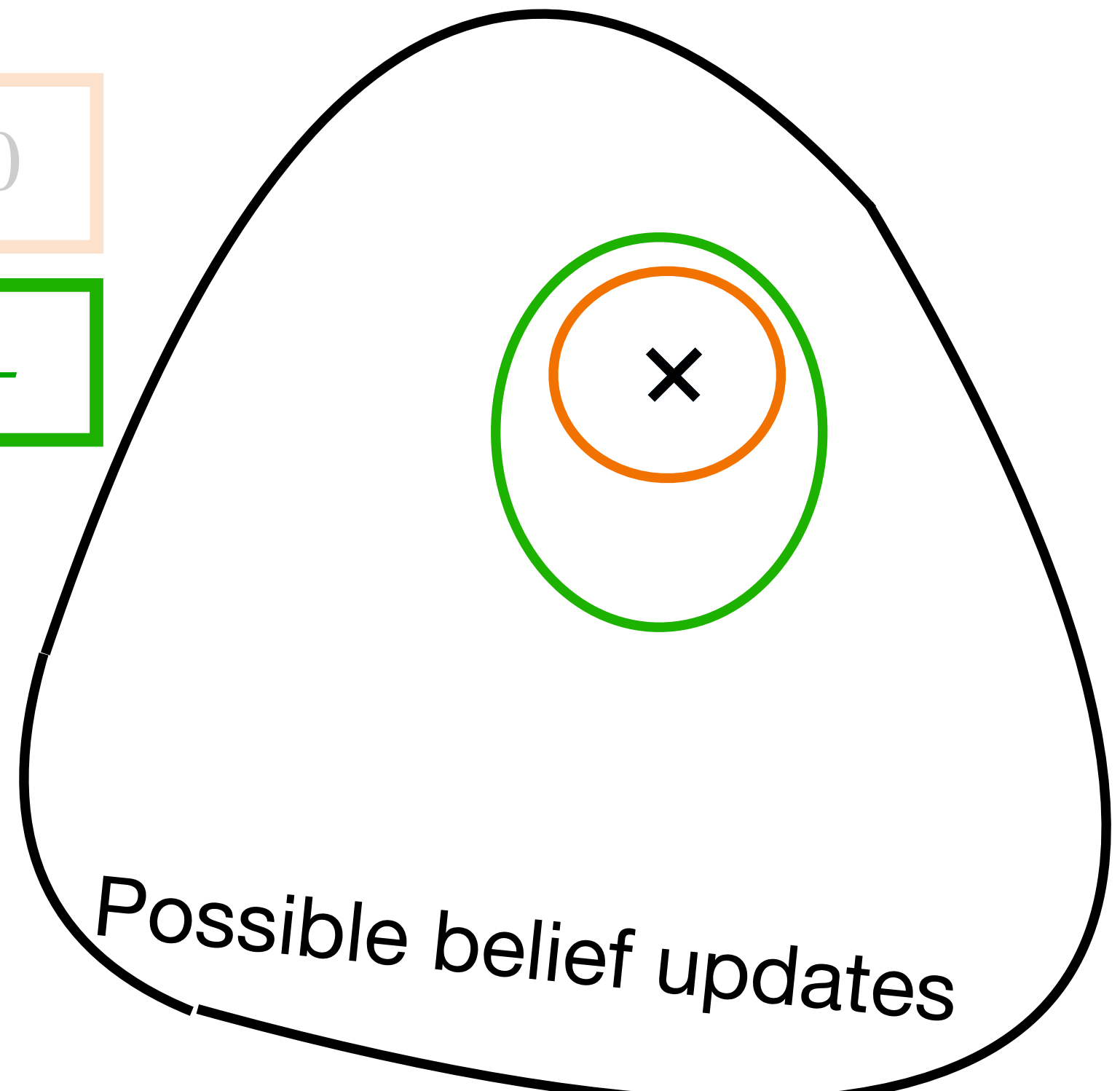
$$\pi_n^{(\lambda)}(\theta | x_{1:n}) = \frac{p(x_{1:n} | \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} | \theta)^\lambda \cdot \pi(\theta) d\theta}$$

Bayes' Posterior

(A1), (A2), (A3)

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

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



L for Robustness + Tractability

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

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

Inspiration: dealing with **unnormalised likelihoods**

Setting: $p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{can be evaluated}} / \underbrace{Z_\theta}_{= \int v(u | \theta) du \text{ (intractable integral)}}$

L for Robustness + Tractability

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

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
 Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
 Altamirano, Briol, & **Knoblauch** (2023); ICML
 Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

Inspiration: dealing with **unnormalised likelihoods**

Setting: $p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{can be evaluated}} / \underbrace{Z_\theta}_{= \int v(u | \theta) du \text{ (intractable integral)}}$

Brute Force: $u_j \sim v(u_j | \theta)$

$$\int v(u | \theta) du \approx \frac{1}{S} \sum_{j=1}^S v(u_j | \theta)$$

Tractability **X**

L for Robustness + Tractability

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

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
 Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
 Altamirano, Briol, & **Knoblauch** (2023); ICML
 Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

Inspiration: dealing with **unnormalised likelihoods**

Setting: $p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{can be evaluated}} / \underbrace{Z_\theta}_{= \int v(u | \theta) du \text{ (intractable integral)}}$

Next!

Brute Force: $u_j \sim v(u_j | \theta)$

$$\int v(u | \theta) du \approx \frac{1}{S} \sum_{j=1}^S v(u_j | \theta)$$

Tractability **X**

Smart strategy:

$$\implies L(x_{1:n}, \theta) = \text{Score Matching / Stein Discrepancies}$$

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

2

L for Robustness + Tractability

Suppose: $p(x_{1:n} | \theta) = v(x_{1:n} | \theta) / Z_\theta$

$\underbrace{\nabla_x \log v(x_{1:n} | \theta)}_{\text{can be evaluated}}$

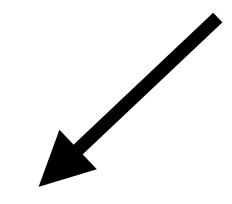
L for Robustness + Tractability



Suppose: $p(x_{1:n} | \theta) = v(x_{1:n} | \theta) / Z_\theta$

Stein / Hyvärinen score

$$\underbrace{\nabla_x \log v(x_{1:n} | \theta)}_{\text{can be evaluated}} = \frac{\nabla_x v(x_{1:n} | \theta)}{v(x_{1:n} | \theta)} = \frac{\nabla_x v(x_{1:n} | \theta) / Z_\theta}{v(x_{1:n} | \theta) / Z_\theta} = \underbrace{\nabla_x \log p(x_{1:n} | \theta)}_{\text{can be evaluated}}$$



L for Robustness + Tractability



Suppose: $p(x_{1:n} | \theta) = v(x_{1:n} | \theta) / Z_\theta$

Stein / Hyvärinen score

$$\underbrace{\nabla_x \log v(x_{1:n} | \theta)}_{\text{can be evaluated}} = \frac{\nabla_x v(x_{1:n} | \theta)}{v(x_{1:n} | \theta)} = \frac{\nabla_x v(x_{1:n} | \theta) / Z_\theta}{v(x_{1:n} | \theta) / Z_\theta} = \underbrace{\nabla_x \log p(x_{1:n} | \theta)}_{\text{can be evaluated}}$$

true data-generating density q_ϵ

$$\|\nabla_x \log p(x_{1:n} | \theta) - \nabla_x \log q_\epsilon(x_{1:n})\|_2^2$$

Score Matching

L for Robustness + Tractability



Suppose: $p(x_{1:n} | \theta) = v(x_{1:n} | \theta) / Z_\theta$

Stein / Hyvärinen score

$$\underbrace{\nabla_x \log v(x_{1:n} | \theta)}_{\text{can be evaluated}} = \frac{\nabla_x v(x_{1:n} | \theta)}{v(x_{1:n} | \theta)} = \frac{\nabla_x v(x_{1:n} | \theta) / Z_\theta}{v(x_{1:n} | \theta) / Z_\theta} = \underbrace{\nabla_x \log p(x_{1:n} | \theta)}_{\text{can be evaluated}}$$

true data-generating density q_ϵ

$$\underbrace{\|\nabla_x \log p(x_{1:n} | \theta) - \nabla_x \log q_\epsilon(x_{1:n})\|_2^2}_{\text{Score Matching}} \stackrel{+C}{=} \underbrace{\|\nabla_x \log p(x_{1:n} | \theta)\|_2^2 + 2\nabla \cdot \nabla_x \log p(x_{1:n} | \theta)}_{\text{can be evaluated}} = L(x_{1:n}, \theta)$$

Score Matching

can be evaluated

L for Robustness + Tractability

$$n \cdot D_F(q_\epsilon, p(\cdot | \theta))$$

Fisher Divergence

$$x_i \sim q_\epsilon \approx$$

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

Score Matching

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

L for Robustness + Tractability

$$n \cdot D_F(q_\epsilon, p(\cdot | \theta))$$

Fisher Divergence

X NOT a robust divergence

$$x_i \sim q_\epsilon \approx$$



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

Score Matching

L NOT robust

Barp, Briol, Duncan, Girolami, & Mackey (2019); NeurIPS
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

L for Robustness + Tractability

Can we generalise $D_F(q_\epsilon, p(\cdot | \theta))$
to make it robust?

$$n \cdot D_F(q_\epsilon, p(\cdot | \theta))$$

Fisher Divergence

X NOT a robust divergence

$$x_i \sim q_\epsilon \approx$$



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

Score Matching

L NOT robust

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

L for Robustness + Tractability

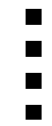
✓ CAN be made robust

Stein Discrepancy

$$n \cdot D_{SD}(q_\epsilon, p(\cdot | \theta))$$



Can we generalise $D_F(q_\epsilon, p(\cdot | \theta))$
to make it robust?



$$n \cdot D_F(q_\epsilon, p(\cdot | \theta))$$

Fisher Divergence

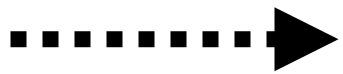
✗ NOT a robust divergence

$$x_i \sim q_\epsilon \approx$$

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

Score Matching

L NOT robust



Barp, Briol, Duncan, Girolami, & Mackey (2019); NeurIPS
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

L for Robustness + Tractability

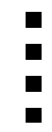
✓ CAN be made robust

Stein Discrepancy

$$n \cdot D_{SD}(q_\epsilon, p(\cdot | \theta))$$



Can we generalise $D_F(q_\epsilon, p(\cdot | \theta))$
to make it robust?



$$n \cdot D_F(q_\epsilon, p(\cdot | \theta))$$

Fisher Divergence

✗ NOT a robust divergence



$$x_i \sim q_\epsilon \approx$$

$$x_i \sim q_\epsilon \approx$$



L robust

Generalised Score Matching

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

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

Score Matching

L NOT robust

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

2

L for Robustness + Tractability

$$n \cdot D_{\text{SD}}(p(\cdot | \theta), q_\varepsilon) = \sup_{f \in \mathcal{F}_\theta} \left| \mathbb{E}_{X \sim q_\varepsilon} [f(X)] - \mathbb{E}_{X \sim p(X|\theta)} [f(X)] \right|$$

L for Robustness + Tractability

$$n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon) = \sup_{f \in \mathcal{F}_\theta} \left| \mathbb{E}_{X \sim q_\epsilon} [f(X)] - \mathbb{E}_{X \sim p(X|\theta)} [f(X)] \right|$$

~~$= 0$~~

Designed so that

$\mathbb{E}_{X \sim p(X|\theta)} [f(X)] = 0$ for all $f \in \mathcal{F}_\theta$

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

L for Robustness + Tractability

$$n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon) = \sup_{f \in \mathcal{F}_\theta} \left| \mathbb{E}_{X \sim q_\epsilon} [f(X)] - \mathbb{E}_{X \sim p(X|\theta)} [f(X)] \right| = \mathbb{E}_{X \sim q_\epsilon} [f^*(X)]$$

$= 0$

Designed so that

$$\mathbb{E}_{X \sim p(X|\theta)} [f(X)] = 0 \text{ for all } f \in \mathcal{F}_\theta$$

\mathcal{F}_θ ensures that supremum has a closed form solution

(weighted Langevin-Stein operator
+ Stein set as RKHS / $C^1 \cup L^2$)

L for Robustness + Tractability

$$n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon) = \sup_{f \in \mathcal{F}_\theta} \left| \mathbb{E}_{X \sim q_\epsilon} [f(X)] - \mathbb{E}_{X \sim p(X|\theta)} [f(X)] \right| = \mathbb{E}_{X \sim q_\epsilon} [f^*(X)] \quad x_i \sim q_\epsilon \approx L(x_{1:n}, \theta)$$

$= 0$

Designed so that

$$\mathbb{E}_{X \sim p(X|\theta)} [f(X)] = 0 \text{ for all } f \in \mathcal{F}_\theta$$

\mathcal{F}_θ ensures that supremum has a closed form solution

All $f(\cdot) \in \mathcal{F}_\theta$ depend on θ only via $w(\cdot) \cdot \nabla_x \log p(\cdot | \theta)$

(weighted Langevin-Stein operator
+ Stein set as RKHS / $C^1 \cup L^2$)

L for Robustness + Tractability

$$n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon) = \sup_{f \in \mathcal{F}_\theta} \left| \mathbb{E}_{X \sim q_\epsilon} [f(X)] - \mathbb{E}_{X \sim p(X|\theta)} [f(X)] \right| = \mathbb{E}_{X \sim q_\epsilon} [f^*(X)] \quad x_i \sim q_\epsilon \approx L(x_{1:n}, \theta)$$

$= 0$

Designed so that

$$\mathbb{E}_{X \sim p(X|\theta)} [f(X)] = 0 \text{ for all } f \in \mathcal{F}_\theta$$

\mathcal{F}_θ ensures that supremum has a closed form solution

All $f(\cdot) \in \mathcal{F}_\theta$ depend on θ only via $w(\cdot) \cdot \nabla_x \log p(\cdot | \theta)$

(weighted Langevin-Stein operator + Stein set as RKHS / $C^1 \cup L^2$)

Robustness can be evaluated

\implies L based on Stein Discrepancies = robust & tractable

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

L for Robustness + Tractability

Comparison: $\pi_n^L(\theta | x_{1:n}) \iff \pi_n(\theta | x_{1:n})$

All settings

$$\pi_n^L(\theta | x_{1:n})$$

robust via $w(\cdot)$

$\pi_n(\theta | x_{1:n})$ is not



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

L for Robustness + Tractability

Comparison: $\pi_n^L(\theta | x_{1:n}) \iff \pi_n(\theta | x_{1:n})$

All settings $\pi_n^L(\theta | x_{1:n})$ robust via $w(\cdot)$ $\pi_n(\theta | x_{1:n})$ is not ✓

$p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{tractable}} / \underbrace{Z_\theta}_{\text{intractable}}$ $\pi_n^L(\theta | x_{1:n})$ faster to compute than $\pi_n(\theta | x_{1:n})$ ✓

L for Robustness + Tractability

Comparison: $\pi_n^L(\theta | x_{1:n}) \iff \pi_n(\theta | x_{1:n})$

All settings $\pi_n^L(\theta | x_{1:n})$ robust via $w(\cdot)$ $\pi_n(\theta | x_{1:n})$ is not ✓

$p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{tractable}} / \underbrace{Z_\theta}_{\text{intractable}}$ $\pi_n^L(\theta | x_{1:n})$ faster to compute than $\pi_n(\theta | x_{1:n})$ ✓

$\underbrace{p(\cdot | \theta)}_{\text{tractable}}$ $\pi_n^L(\theta | x_{1:n})$ roughly as fast as $\pi_n(\theta | x_{1:n})$ ✓

L for Robustness + Tractability

Comparison: $\pi_n^L(\theta | x_{1:n}) \iff \pi_n(\theta | x_{1:n})$

All settings $\pi_n^L(\theta | x_{1:n})$ robust via $w(\cdot)$ $\pi_n(\theta | x_{1:n})$ is not ✓

$p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{tractable}} / \underbrace{Z_\theta}_{\text{intractable}}$ $\pi_n^L(\theta | x_{1:n})$ faster to compute than $\pi_n(\theta | x_{1:n})$ ✓

$\underbrace{p(\cdot | \theta)}_{\text{tractable}}$ $\pi_n^L(\theta | x_{1:n})$ roughly as fast as $\pi_n(\theta | x_{1:n})$ ✓

Next!

$p(\cdot | \theta)$ exponential family + $\pi(\theta)$ conjugate $\pi_n^L(\theta | x_{1:n})$??? Slower, surely ??? $\pi_n(\theta | x_{1:n})$

???

L for Robustness + Tractability

$$n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon) \stackrel{x_i \sim q_\epsilon}{\approx} L(x_{1:n}, \theta) = \sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))$$

↑
exponential family
(possibly with intractable normaliser)

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

L for Robustness + Tractability

$$\exp\{-n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon)\} \stackrel{x_i \sim q_\epsilon}{\approx} \exp\{-L(x_{1:n}, \theta)\} = \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\}$$

↑
exponential family
(possibly with intractable normaliser)

= unnormalised Squared Exponential / Gaussian in θ

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
 Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
 Altamirano, Briol, & **Knoblauch** (2023); ICML
 Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

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

L for Robustness + Tractability

$$\exp\{-n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon)\} \stackrel{x_i \sim q_\epsilon}{\approx} \exp\{-L(x_{1:n}, \theta)\} = \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\}$$

↑
 exponential family
 (possibly with intractable normaliser)

= unnormalised Squared Exponential / Gaussian in θ

conjugate prior!

$$\pi_n^L(\theta | x_{1:n}) \propto \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\} \underbrace{\exp\left\{(\theta - \mu_0)^\top \Lambda_0 (\theta - \mu_0)\right\}}_{\text{squared exponential prior}}$$

squared exponential prior

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

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
 Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
 Altamirano, Briol, & **Knoblauch** (2023); ICML
 Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

L for Robustness + Tractability

$$\exp\left\{-n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon)\right\} \stackrel{x_i \sim q_\epsilon}{\approx} \exp\left\{-L(x_{1:n}, \theta)\right\} = \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\}$$

↑
 exponential family
 (possibly with intractable normaliser)

= unnormalised Squared Exponential / Gaussian in θ

$$\pi_n^L(\theta | x_{1:n}) \propto \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\} \underbrace{\exp\left\{(\theta - \mu_0)^\top \Lambda_0 (\theta - \mu_0)\right\}}_{\text{squared exponential prior}}$$

conjugate prior!

squared exponential prior

$$= \mathcal{N}(\theta; \mu_L(x_{1:n}), \Sigma_L(x_{1:n}))$$

Closed form / conjugate Post-Bayesian posterior

L for Robustness + Tractability

$$\exp\{-n \cdot D_{SD}(p(\cdot | \theta), q_\epsilon)\} \stackrel{x_i \sim q_\epsilon}{\approx} \exp\{-L(x_{1:n}, \theta)\} = \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\}$$

↑
 exponential family
 (possibly with intractable normaliser)

= unnormalised Squared Exponential / Gaussian in θ

conjugate prior!

$$\pi_n^L(\theta | x_{1:n}) \propto \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\} \underbrace{\exp\left\{(\theta - \mu_0)^\top \Lambda_0 (\theta - \mu_0)\right\}}_{\text{squared exponential prior}}$$

$$= \mathcal{N}(\theta; \mu_L(x_{1:n}), \Sigma_L(x_{1:n}))$$

Closed form / conjugate Post-Bayesian posterior

even if $p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{tractable}} / \underbrace{Z_\theta}_{\text{intractable}} !!!$

Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
 Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
 Altamirano, Briol, & **Knoblauch** (2023); ICML
 Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)

L for Robustness + Tractability

Comparison: $\pi_n^L(\theta | x_{1:n}) \iff \pi_n(\theta | x_{1:n})$

All settings $\pi_n^L(\theta | x_{1:n})$ robust via $w(\cdot)$ $\pi_n(\theta | x_{1:n})$ is not ✓

$p(\cdot | \theta) = \underbrace{v(\cdot | \theta)}_{\text{tractable}} / \underbrace{Z_\theta}_{\text{intractable}}$ $\pi_n^L(\theta | x_{1:n})$ faster to compute than $\pi_n(\theta | x_{1:n})$ ✓

$\underbrace{p(\cdot | \theta)}_{\text{tractable}}$ $\pi_n^L(\theta | x_{1:n})$ roughly as fast as $\pi_n(\theta | x_{1:n})$ ✓

$p(\cdot | \theta)$ exponential family + $\pi(\theta)$ conjugate $\pi_n^L(\theta | x_{1:n})$ as fast as / faster to compute than $\pi_n(\theta | x_{1:n})$ ✓

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

2

L for Robustness + Tractability

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

Like sufficient statistics!
can be updated with simple algebra

L for Robustness + Tractability

$$\pi_n^L(\theta | x_{1:n}) = \frac{\exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta)}{\int \exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta) d\theta} = \mathcal{N}(\theta; \mu_L(x_{1:n}), \Sigma_L(x_{1:n}))$$

Like sufficient statistics!
can be updated with simple algebra

Graphical Modelling

Proposition 2. Consider $\mathcal{X} = \mathbb{R}^d$ and the Langevin Stein operator \mathcal{S}_{p_θ} in (5), where p_θ is the exponential family in (10), and a kernel $K \in C_b^{1,1}(\mathbb{R}^d \times \mathbb{R}^d; \mathbb{R}^{d \times d})$. Assuming the prior has a p.d.f. π , the LSD-Bayes generalised posterior has a p.d.f.

$$\pi_n^D(\theta) \propto \pi(\theta) \exp(-\beta n \{\eta(\theta) \cdot \Lambda_n \eta(\theta) + \eta(\theta) \cdot \nu_n\}),$$

where $\Lambda_n \in \mathbb{R}^{k \times k}$ and $\nu_n \in \mathbb{R}^k$ are defined as

$$\Lambda_n := \frac{1}{n^2} \sum_{i,j=1}^n \nabla t(x_i) \cdot K(x_i, x_j) \nabla t(x_j),$$

$$\nu_n := \frac{1}{n^2} \sum_{i,j=1}^n \nabla t(x_i) \cdot (\nabla_{x_j} \cdot K(x_i, x_j)) + \nabla t(x_j) \cdot (\nabla_{x_i} \cdot K(x_i, x_j)) + 2 \nabla t(x_i) \cdot K(x_i, x_j) \nabla b(x_j).$$

For a natural exponential family we have $\eta(\theta) = \theta$, and the prior $\pi(\theta) \propto \exp(-\frac{1}{2}(\theta - \mu) \cdot \Sigma^{-1}(\theta - \mu))$ leads to a generalised posterior

$$\pi_n^D(\theta) \propto \exp\left(-\frac{1}{2}(\theta - \mu_n) \cdot \Sigma_n^{-1}(\theta - \mu_n)\right),$$

where $\Sigma_n^{-1} := \Sigma^{-1} + 2\beta n \Lambda_n$ and $\mu_n := \Sigma_n^{-1}(\Sigma^{-1}\mu - \nu_n)$.

Changepoints

Proposition 3.1. If p_θ is given by (5), then

$$\pi_\omega^{\mathcal{D}^m}(\theta | x_{1:T}) \propto \pi(\theta) \exp(-\omega T [\eta(\theta)^\top \Lambda_T \eta(\theta) + \eta(\theta)^\top \nu_T]),$$

for $\Lambda_T = \frac{1}{T} \sum_{t=1}^T \Lambda(x_t)$, $\nu_T = \frac{2}{T} \sum_{t=1}^T \nu(x_t)$, and

$$\Lambda(x) = (\nabla r^\top m m^\top \nabla r)(x),$$

$$\nu(x) = (\nabla r^\top m m^\top \nabla b + \nabla \cdot (m m^\top \nabla r))(x).$$

Taking $\eta(\theta) = \theta$ and choosing a squared exponential prior $\pi(\theta) \propto \exp(-\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu))$, also makes $\pi_\omega^{\mathcal{D}^m}(\theta | x_{1:T})$ a (truncated) normal of the form

$$\pi_\omega^{\mathcal{D}^m}(\theta | x_{1:T}) \propto \exp\left(-\frac{1}{2}(\theta - \mu_T)^\top \Sigma_T^{-1}(\theta - \mu_T)\right),$$

for $\Sigma_T^{-1} = \Sigma^{-1} + 2\omega T \Lambda_T$ and $\mu_T = \Sigma_T^{-1}(\Sigma^{-1}\mu - \omega T \nu_T)$.

Gaussian Processes

Proposition 3.1. Suppose $f \sim \mathcal{GP}(m, k)$ and $\varepsilon \sim \mathcal{N}(0, I_n \sigma^2)$. Then, the RCGP posterior is

$$p^w(\mathbf{f} | \mathbf{y}, \mathbf{x}) = \mathcal{N}(\mathbf{f}; \mu^R, \Sigma^R),$$

$$\mu^R = \mathbf{m} + K(K + \sigma^2 J_w)^{-1}(\mathbf{y} - \mathbf{m}_w),$$

$$\Sigma^R = K(K + \sigma^2 J_w)^{-1} \sigma^2 J_w,$$

for $\mathbf{w} = (w(x_1, y_1), \dots, w(x_n, y_n))^\top$, $\mathbf{m}_w = \mathbf{m} + \sigma^2 \nabla_y \log(\mathbf{w}^2)$ and $J_w = \text{diag}(\frac{\sigma^2}{2} \mathbf{w}^{-2})$. The RCGP's posterior predictive over $f_\star = f(x_\star)$ at $x_\star \in \mathcal{X}$ is

$$p^w(f_\star | x_\star, \mathbf{x}, \mathbf{y}) = \int_{\mathbb{R}^n} p(f_\star | x_\star, \mathbf{f}, \mathbf{x}, \mathbf{y}) p^w(\mathbf{f} | \mathbf{y}, \mathbf{x}) d\mathbf{f}$$

$$= \mathcal{N}(f_\star; \mu_\star^R, \Sigma_\star^R),$$

$$\mu_\star^R = m_\star + \mathbf{k}_\star^\top (K + \sigma^2 J_w)^{-1}(\mathbf{y} - \mathbf{m}_w),$$

$$\Sigma_\star^R = k_{\star\star} - \mathbf{k}_\star^\top (K + \sigma^2 J_w)^{-1} \mathbf{k}_\star.$$

L for Robustness + Tractability

$$\pi_n^L(\theta | x_{1:n}) = \frac{\exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta)}{\int \exp\{-L(x_{1:n}, \theta)\} \cdot \pi(\theta) d\theta} = \mathcal{N}(\theta; \mu_L(x_{1:n}), \Sigma_L(x_{1:n}))$$

Next!

Like sufficient statistics!
can be updated with simple algebra

⇒ Feasible for on-line problems!

Graphical Modelling

Proposition 2. Consider $\mathcal{X} = \mathbb{R}^d$ and the Langevin Stein operator \mathcal{S}_{p_θ} in (5), where p_θ is the exponential family in (10), and a kernel $K \in C_b^{1,1}(\mathbb{R}^d \times \mathbb{R}^d; \mathbb{R}^{d \times d})$. Assuming the prior has a p.d.f. π , the LSD-Bayes generalised posterior has a p.d.f.

$$\pi_n^D(\theta) \propto \pi(\theta) \exp(-\beta n \{\eta(\theta) \cdot \Lambda_n \eta(\theta) + \eta(\theta) \cdot \nu_n\}),$$

where $\Lambda_n \in \mathbb{R}^{k \times k}$ and $\nu_n \in \mathbb{R}^k$ are defined as

$$\Lambda_n := \frac{1}{n^2} \sum_{i,j=1}^n \nabla t(x_i) \cdot K(x_i, x_j) \nabla t(x_j),$$

$$\nu_n := \frac{1}{n^2} \sum_{i,j=1}^n \nabla t(x_i) \cdot (\nabla_{x_j} \cdot K(x_i, x_j)) + \nabla t(x_j) \cdot (\nabla_{x_i} \cdot K(x_i, x_j)) + 2 \nabla t(x_i) \cdot K(x_i, x_j) \nabla b(x_j).$$

For a natural exponential family we have $\eta(\theta) = \theta$, and the prior $\pi(\theta) \propto \exp(-\frac{1}{2}(\theta - \mu) \cdot \Sigma^{-1}(\theta - \mu))$ leads to a generalised posterior

$$\pi_n^D(\theta) \propto \exp\left(-\frac{1}{2}(\theta - \mu_n) \cdot \Sigma_n^{-1}(\theta - \mu_n)\right),$$

where $\Sigma_n^{-1} := \Sigma^{-1} + 2\beta n \Lambda_n$ and $\mu_n := \Sigma_n^{-1}(\Sigma^{-1}\mu - \nu_n)$.

Changepoints

Proposition 3.1. If p_θ is given by (5), then

$$\pi_\omega^{\mathcal{D}^m}(\theta | x_{1:T}) \propto \pi(\theta) \exp(-\omega T [\eta(\theta)^\top \Lambda_T \eta(\theta) + \eta(\theta)^\top \nu_T]),$$

for $\Lambda_T = \frac{1}{T} \sum_{t=1}^T \Lambda(x_t)$, $\nu_T = \frac{2}{T} \sum_{t=1}^T \nu(x_t)$, and

$$\Lambda(x) = (\nabla r^\top m m^\top \nabla r)(x),$$

$$\nu(x) = (\nabla r^\top m m^\top \nabla b + \nabla \cdot (m m^\top \nabla r))(x).$$

Taking $\eta(\theta) = \theta$ and choosing a squared exponential prior $\pi(\theta) \propto \exp(-\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu))$, also makes $\pi_\omega^{\mathcal{D}^m}(\theta | x_{1:T})$ a (truncated) normal of the form

$$\pi_\omega^{\mathcal{D}^m}(\theta | x_{1:T}) \propto \exp\left(-\frac{1}{2}(\theta - \mu_T)^\top \Sigma_T^{-1}(\theta - \mu_T)\right),$$

for $\Sigma_T^{-1} = \Sigma^{-1} + 2\omega T \Lambda_T$ and $\mu_T = \Sigma_T^{-1}(\Sigma^{-1}\mu - \omega T \nu_T)$.

Gaussian Processes

Proposition 3.1. Suppose $f \sim \mathcal{GP}(m, k)$ and $\varepsilon \sim \mathcal{N}(0, I_n \sigma^2)$. Then, the RCGP posterior is

$$p^w(\mathbf{f} | \mathbf{y}, \mathbf{x}) = \mathcal{N}(\mathbf{f}; \mu^R, \Sigma^R),$$

$$\mu^R = \mathbf{m} + K(K + \sigma^2 J_w)^{-1}(\mathbf{y} - \mathbf{m}_w),$$

$$\Sigma^R = K(K + \sigma^2 J_w)^{-1} \sigma^2 J_w,$$

for $\mathbf{w} = (w(x_1, y_1), \dots, w(x_n, y_n))^\top$, $\mathbf{m}_w = \mathbf{m} + \sigma^2 \nabla_y \log(\mathbf{w}^2)$ and $J_w = \text{diag}(\frac{\sigma^2}{2} \mathbf{w}^{-2})$. The RCGP's posterior predictive over $f_\star = f(x_\star)$ at $x_\star \in \mathcal{X}$ is

$$p^w(f_\star | x_\star, \mathbf{x}, \mathbf{y}) = \int_{\mathbb{R}^n} p(f_\star | x_\star, \mathbf{f}, \mathbf{x}, \mathbf{y}) p^w(\mathbf{f} | \mathbf{y}, \mathbf{x}) d\mathbf{f}$$

$$= \mathcal{N}(f_\star; \mu_\star^R, \Sigma_\star^R),$$

$$\mu_\star^R = m_\star + \mathbf{k}_\star^\top (K + \sigma^2 J_w)^{-1}(\mathbf{y} - \mathbf{m}_w),$$

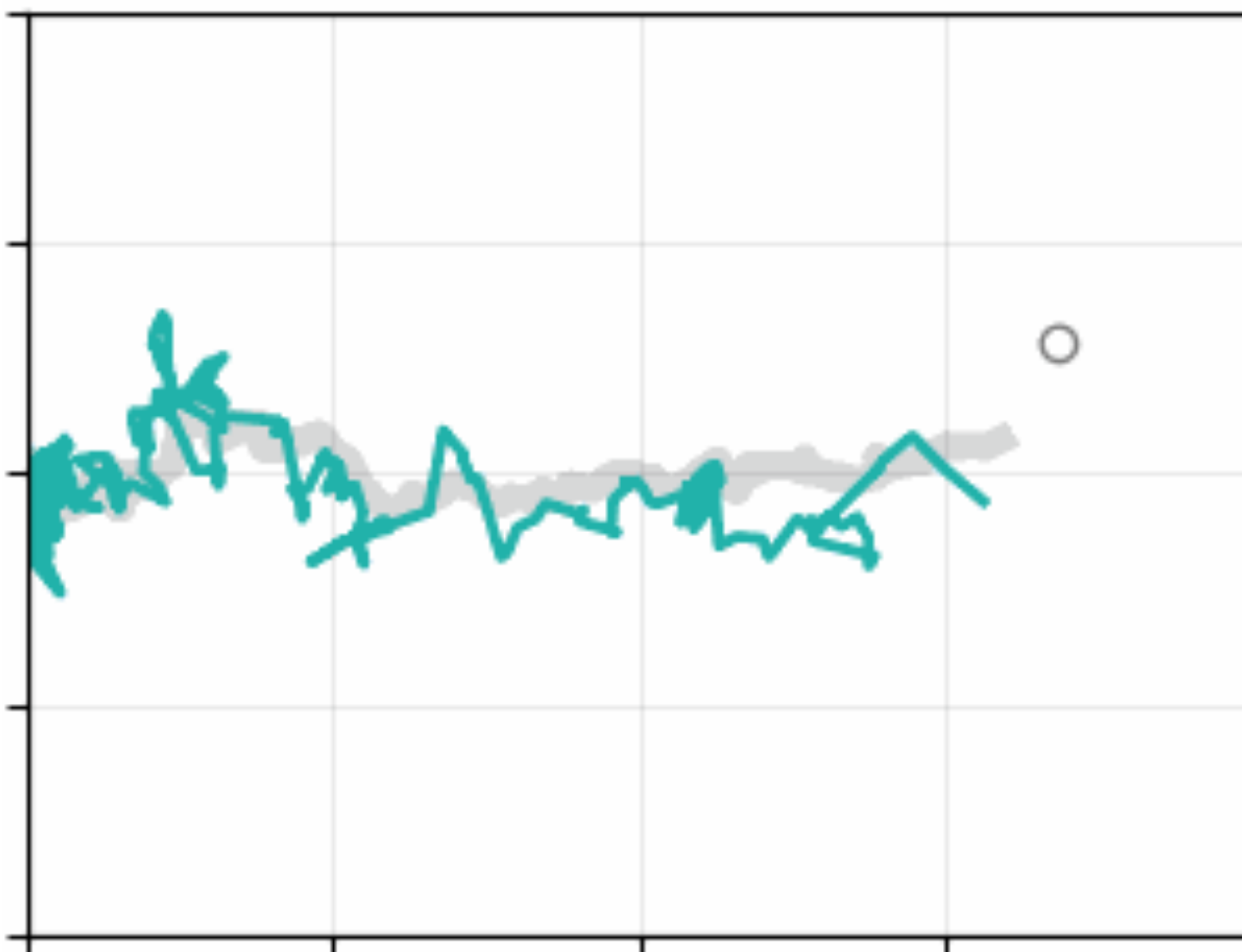
$$\Sigma_\star^R = k_{\star\star} - \mathbf{k}_\star^\top (K + \sigma^2 J_w)^{-1} \mathbf{k}_\star.$$

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

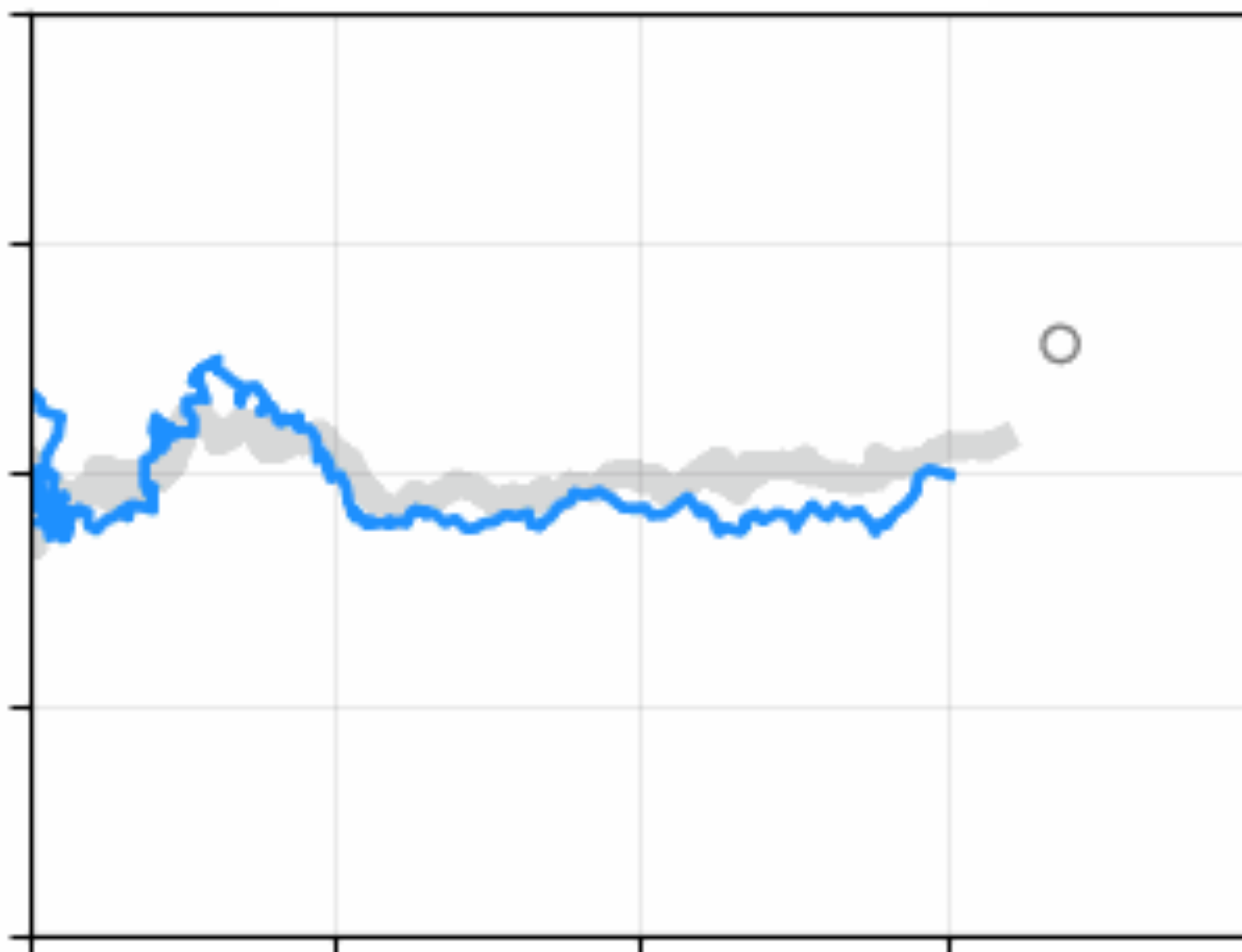
L for Robustness + Tractability: Kalman Filter

⇒ Closed form Post-Bayesian posteriors: feasible for on-line problems!

Standard Kalman Filter



Robust version



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

Summary: State of the Art in Post-Bayesian ML

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

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification + computation

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

Schmon, Cannon, & **Knoblauch** (2020); AABI
Matsubara, **Knoblauch**, Briol, & Oates (2022); JRSS-B
Dellaporta, **Knoblauch**, Damoulas, & Briol (2022); AISTATS (best paper award)
Altamirano, Briol, & **Knoblauch** (2023); ICML
Altamirano, Briol, & **Knoblauch** (2024); ICML (spotlight)
Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, **Knoblauch**,
Briol, & Murphy (2024); ICML

Q: Can we design losses L that are **both robust and tractable** ?

A: Yes! Stein Discrepancies. (And weighted likelihoods.)

The Future of Post-Bayesian ML

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification

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

Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification +
computation

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

Schmon, Cannon, & Knoblauch (2020); AABI
 Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
 Dellaporta, Knoblauch, Damoulas, & Briol (2022); AISTATS (best paper award)
 Altamirano, Briol, & Knoblauch (2023); ICML
 Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt, Knoblauch, Briol, & Murphy (2024); ICML

3 The Future

$$q_n^*(\theta)$$

model misspecification +
prior misspecification +
computation

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

Husain & Knoblauch (2022); ALT
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Matsubara, Knoblauch, Briol, & Oates (2023); JASA
 Wild, Sejdinovic, & Knoblauch (2024); forthcoming
 Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2024); NeurIPS (oral)

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

3

The Future of Post-Bayesian ML

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

3 The Future

$$q_n^*(\theta)$$

model misspecification +
prior misspecification +
computation

~~(A1)~~

~~(A2)~~

~~(A3)~~

Husain & **Knoblauch** (2022); ALT
Knoblauch, Jewson, & Damoulas (2022); JMLR
Matsubara, **Knoblauch**, Briol, & Oates (2023); JASA
Wild, Sejdinovic, & **Knoblauch** (2024); forthcoming
Wild, Ghalebikesabi, Sejdinovic, & **Knoblauch** (2024); NeurIPS (oral)

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

The Future of Post-Bayesian ML

Optimisation-centric posteriors / Generalised Variational Inference

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

$$q_n^*(\theta) = \arg \min_{q \in \mathcal{Q}} \underbrace{\mathcal{L}_{L, D}(q)}; \quad \mathcal{Q} \subseteq \mathcal{P}(\Theta)$$

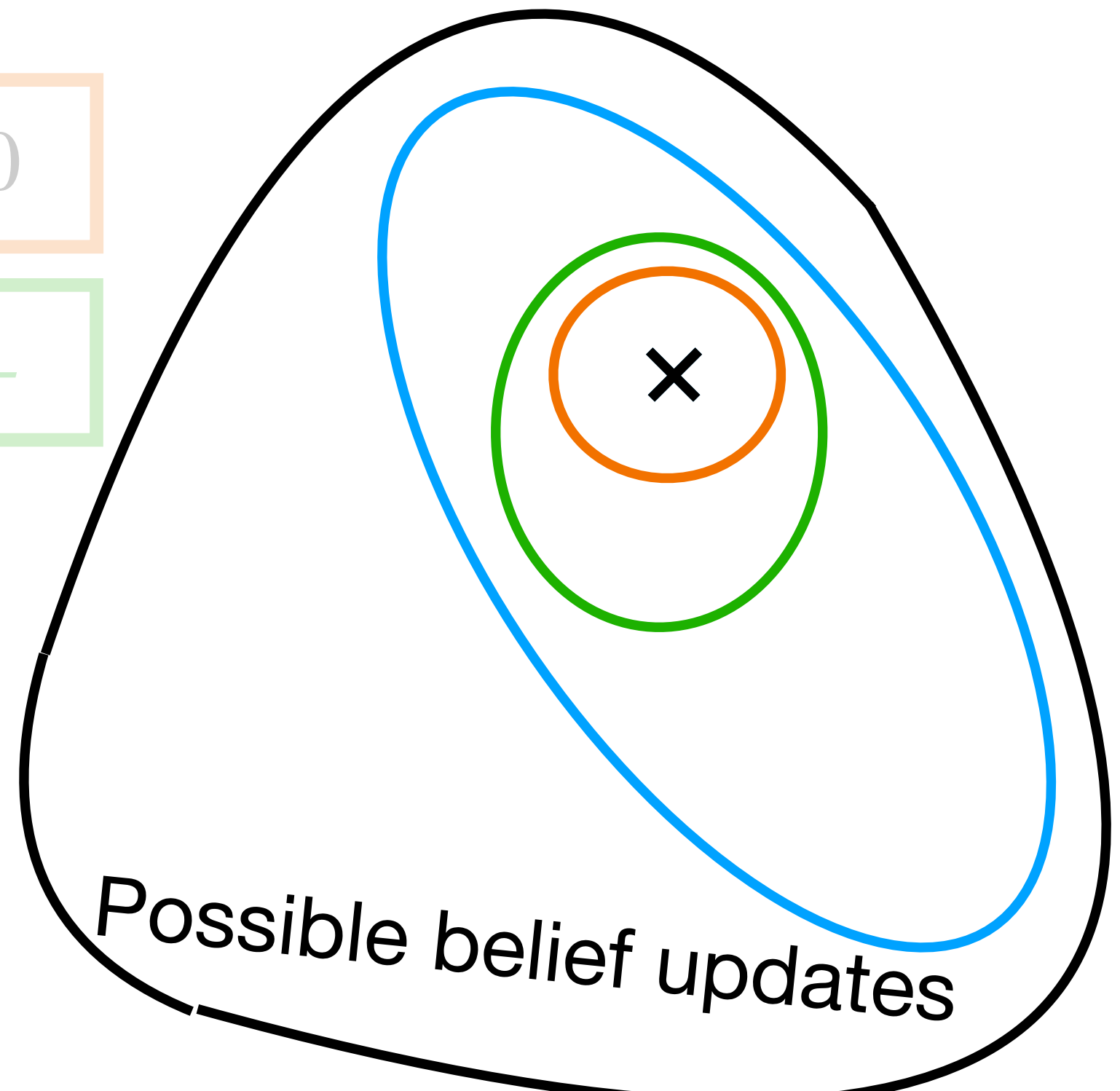
$$= \int L(x_{1:n}, \theta) q(\theta) d\theta + D(q, \pi)$$

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

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

$$\begin{matrix} \text{KL} & \longrightarrow & D \\ \mathcal{P}(\Theta) & \longrightarrow & \mathcal{Q} \end{matrix}$$

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



Gibbs/Generalised/Pseudo Posterior

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

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

Power/Fractional/Cold Posterior

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

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

Bayes' Posterior

(A1), (A2), (A3)

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

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

The Future of Post-Bayesian ML

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

$$= \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \int L(x_{1:n}, \theta) q(\theta) d\theta + \text{KL}(q, \pi) \right\}$$

~~(A3)~~ by optimising over a set $\mathcal{Q} \subseteq \mathcal{P}(\Theta)$

~~(A1)~~ by using robust loss L

~~(A2)~~ by using robust regulariser D

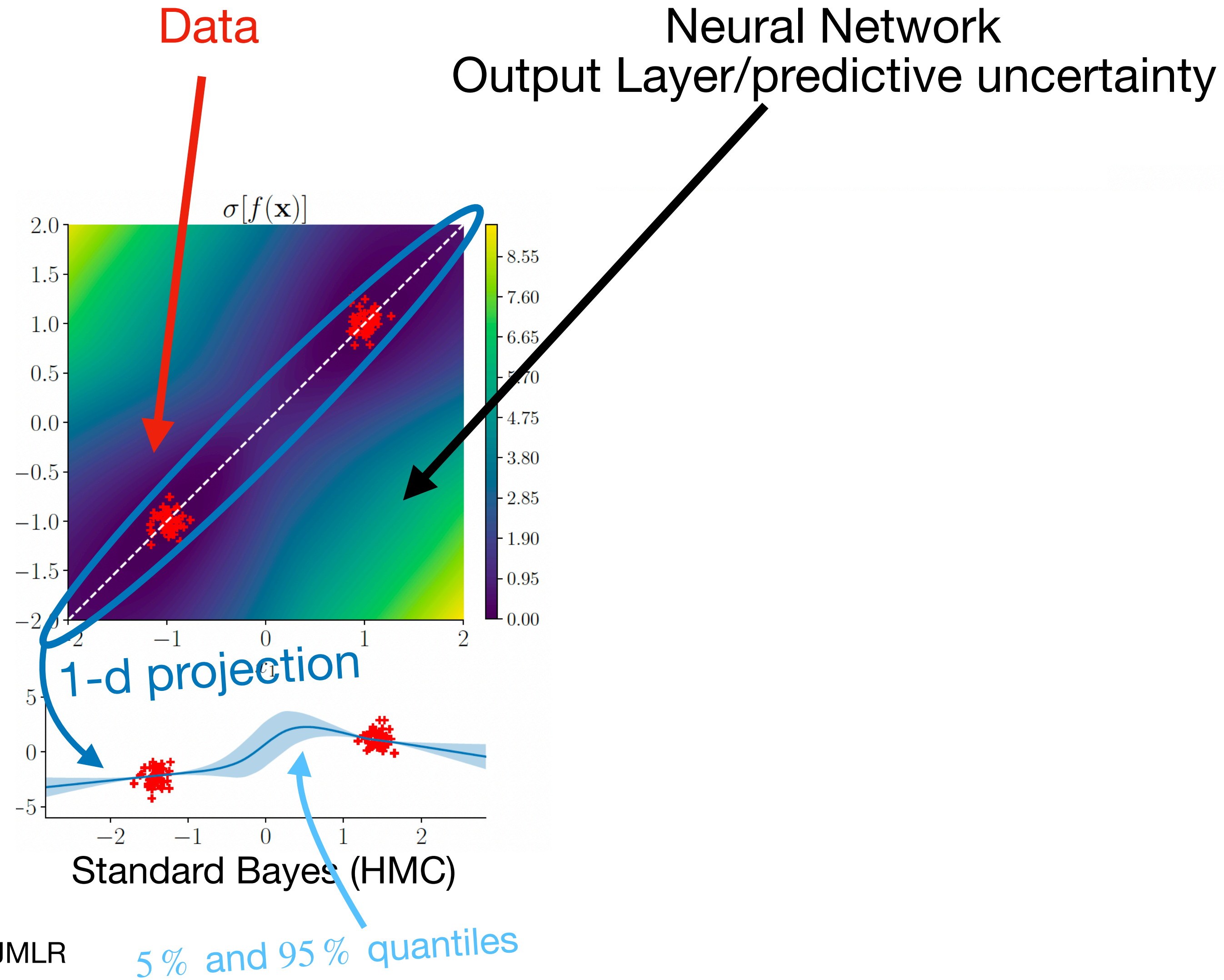
- ~~(A1)~~ model well-specified
- ~~(A2)~~ prior well-specified
- ~~(A3)~~ computationally feasible

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

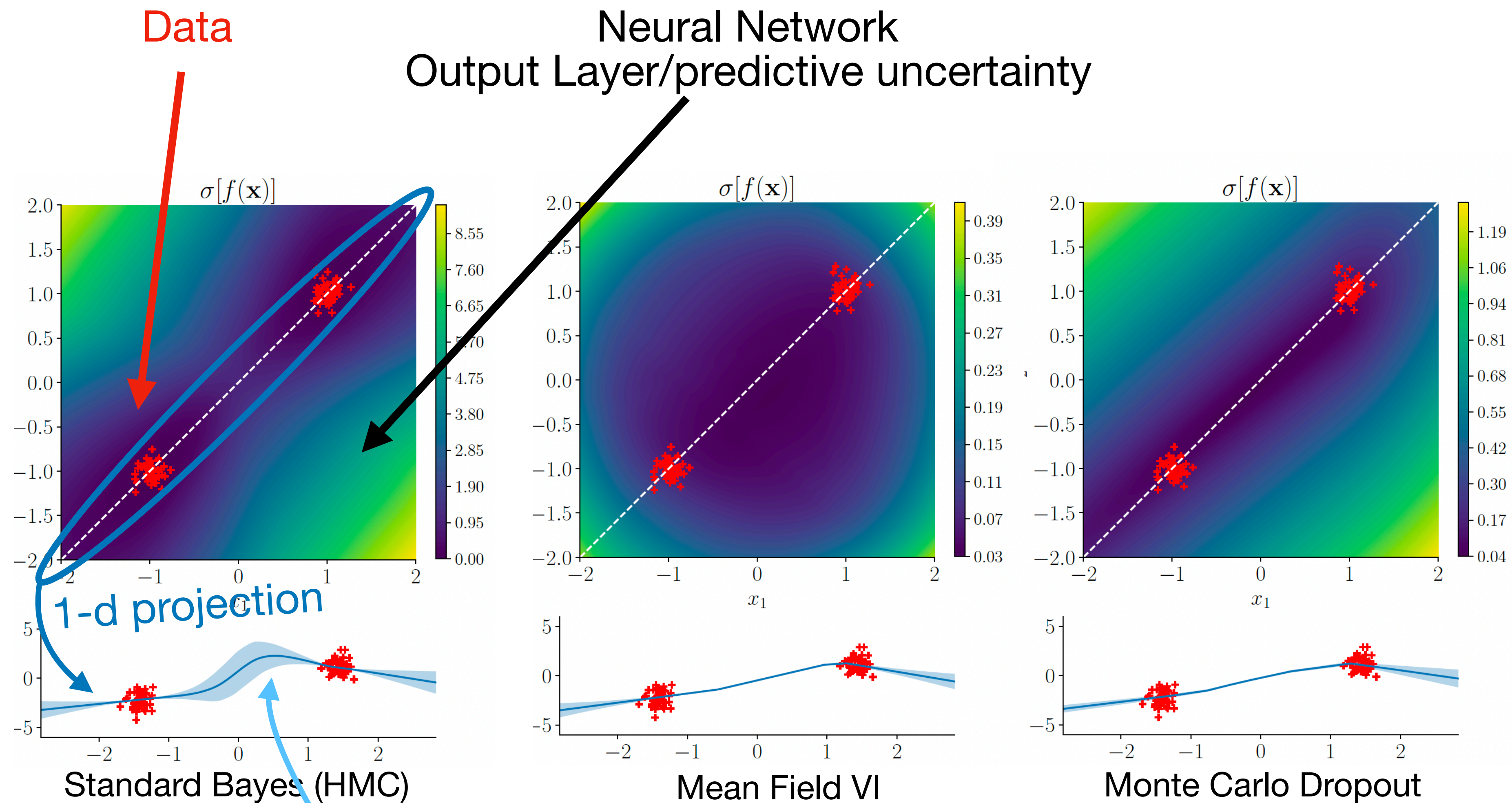
Optimisation-centric posteriors / Generalised Variational Inference (GVI) $\quad \quad \quad =: \mathcal{L}_{L, D}(q)$

Husain & **Knoblauch** (2022); ALT
Knoblauch, Jewson, & Damoulas (2022); JMLR
 Wild, Sejdinovic, & **Knoblauch** (2024); forthcoming
 Wild, Ghalebikesabi, Sejdinovic, & **Knoblauch** (2024); NeurIPS (oral)

The Future of Post-Bayesian ML

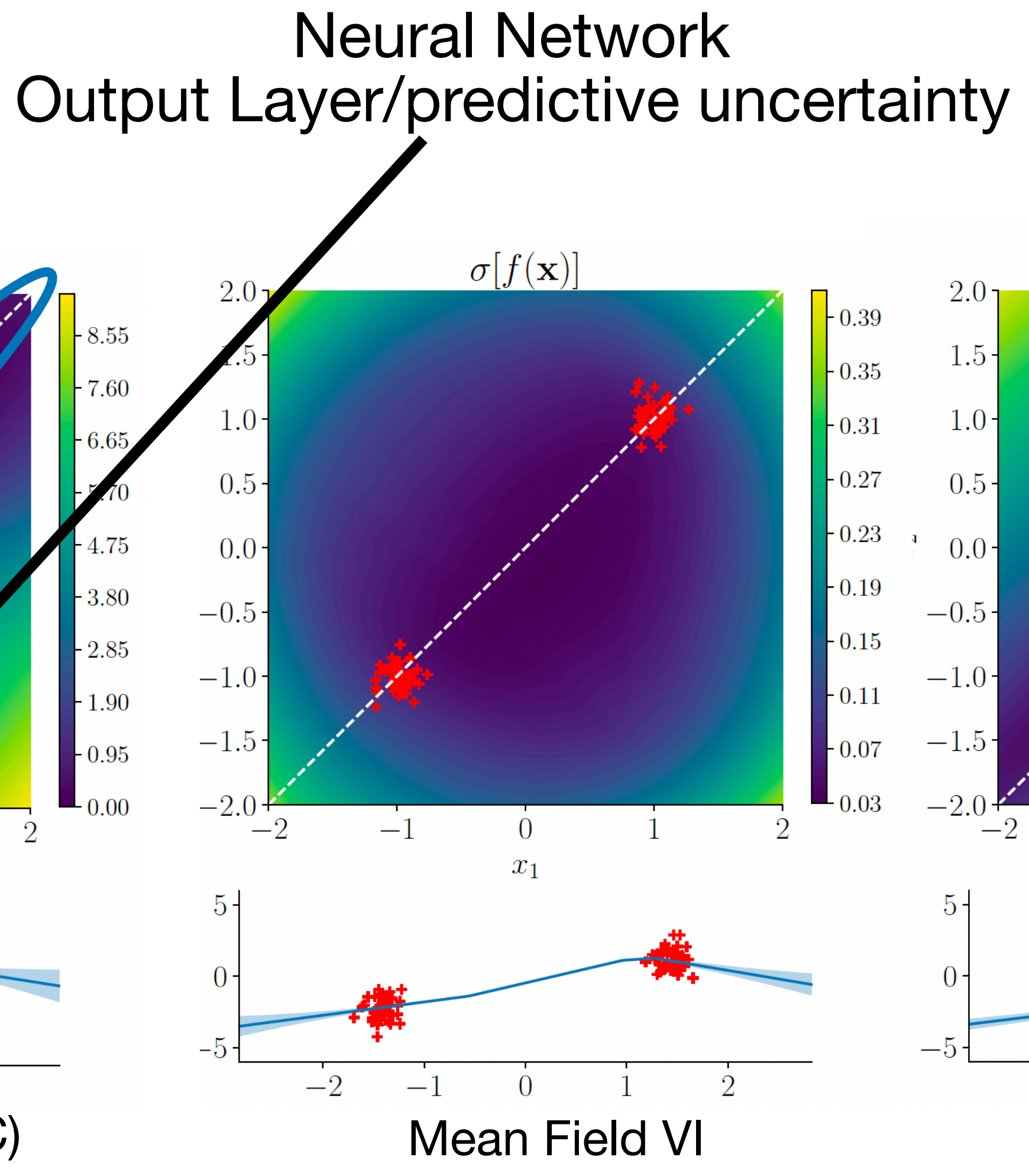
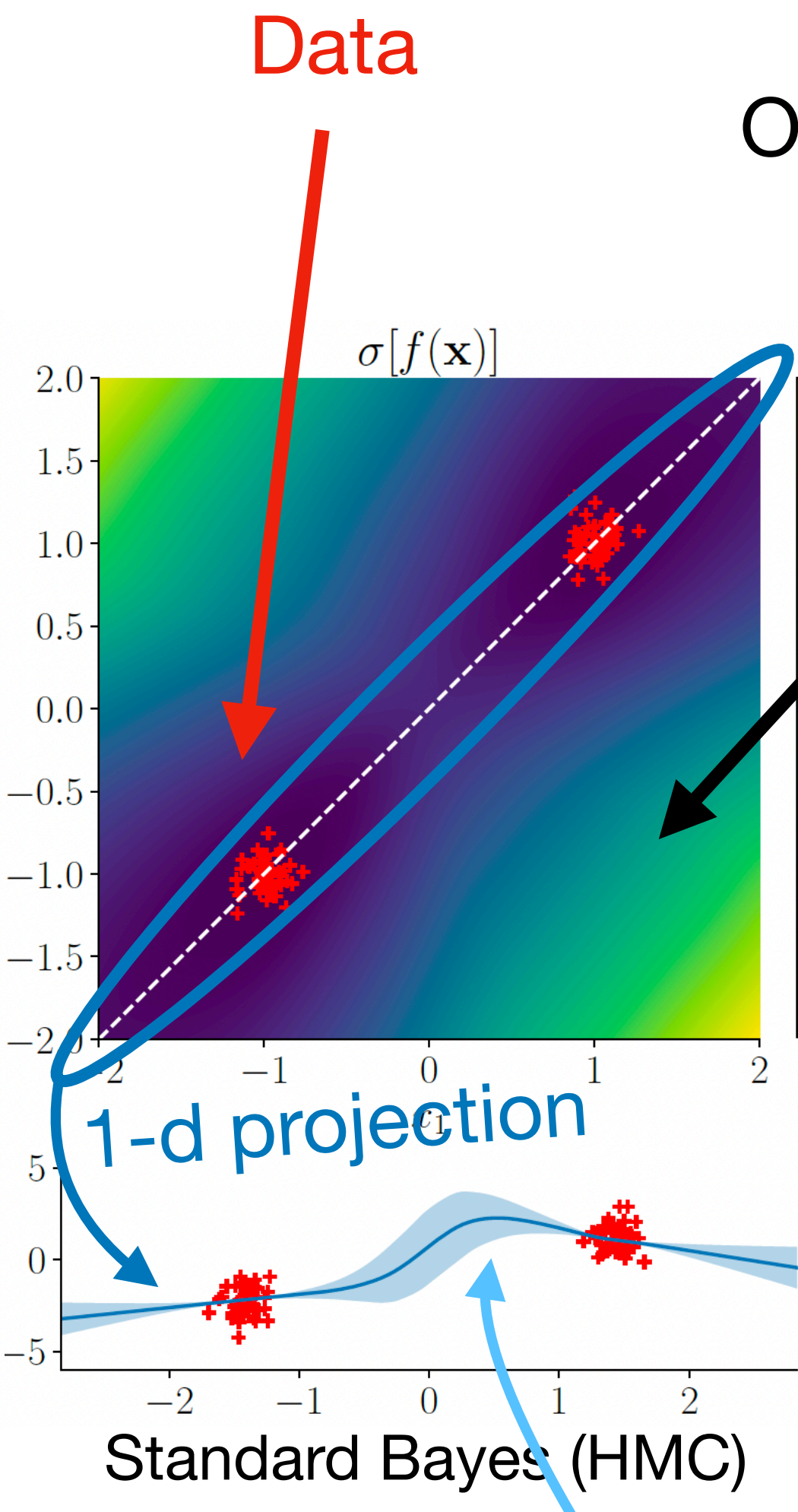
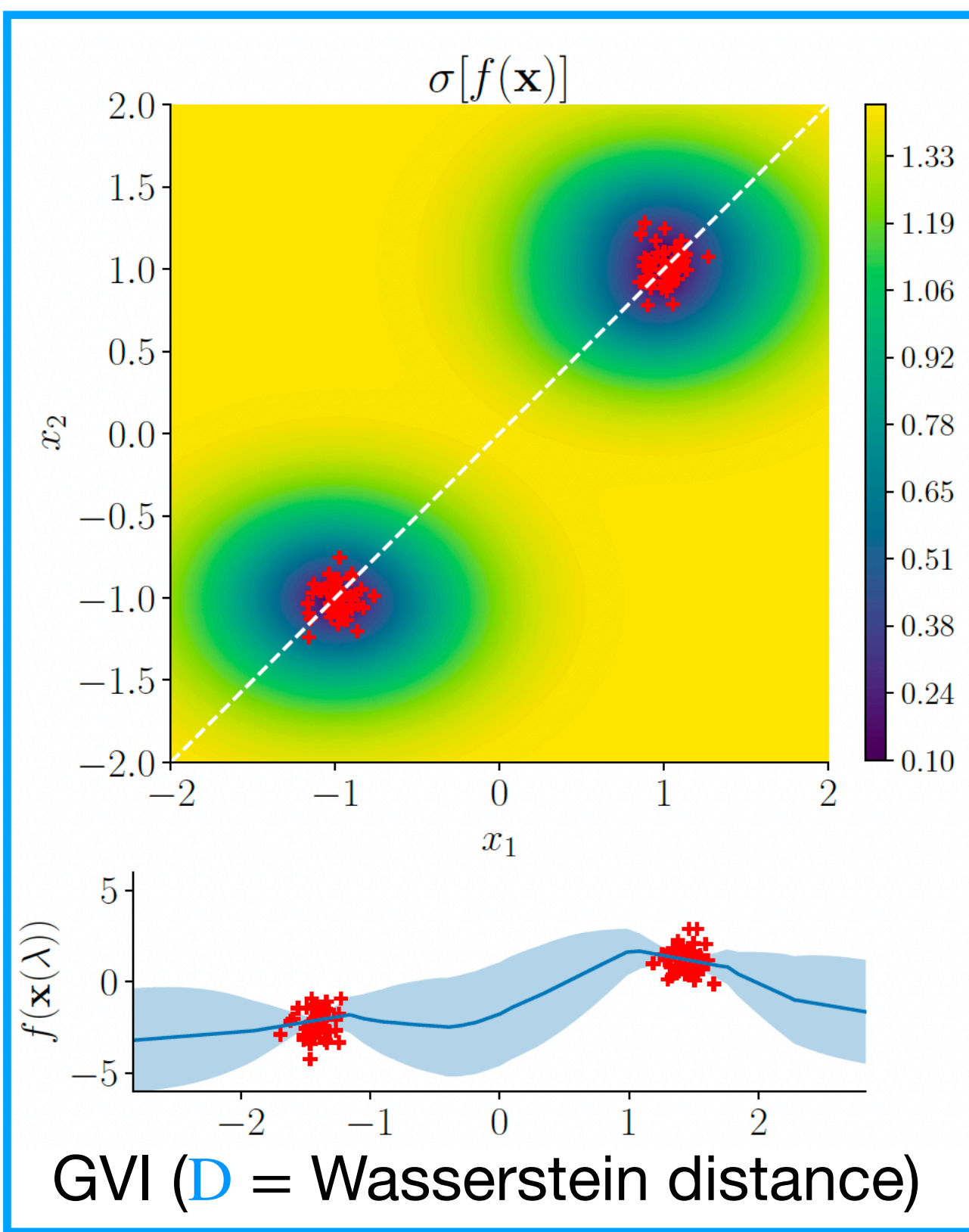


The Future of Post-Bayesian ML



The Future of Post-Bayesian ML

D = Wasserstein Distance
Improves Uncertainty



Knoblauch, Jewson, & Damoulas (2022); JMLR
Wild, Hu, & Sejdinovic, NeurIPS (2022)

5% and 95% quantiles

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

3

The Future of Post-Bayesian ML

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

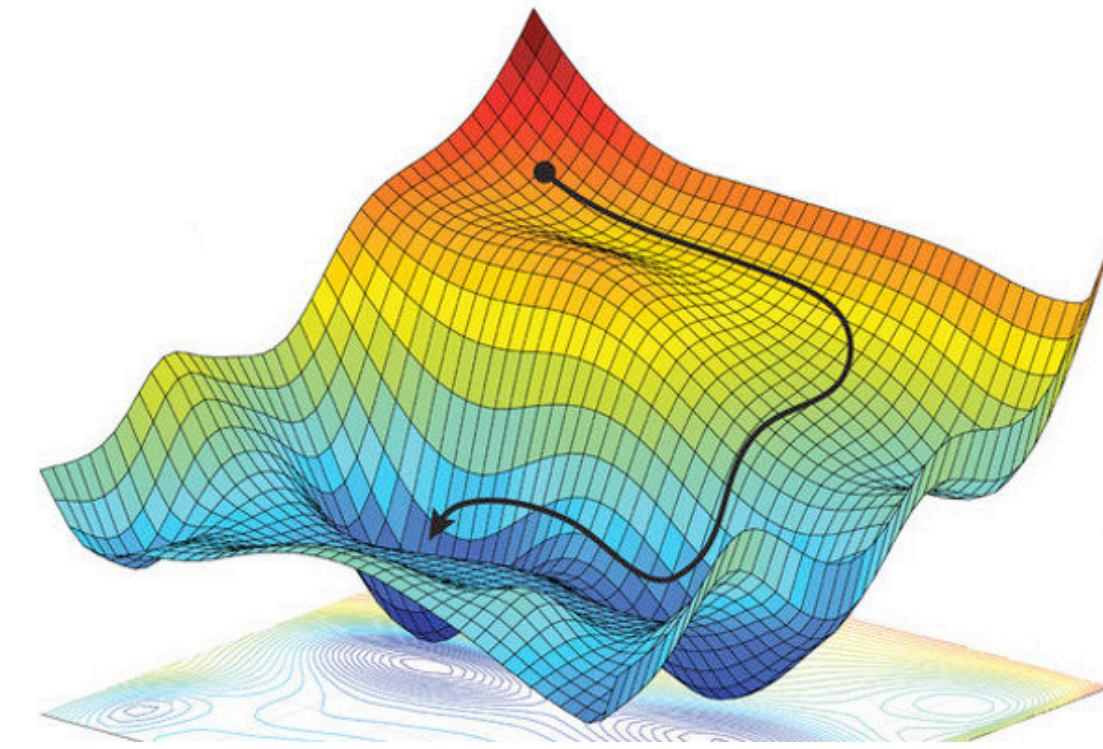
New Problem: Unclear how to compute $q_n^*(\theta)$ unless \mathcal{Q} is a parameterised set of distributions

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

3

The Future of Post-Bayesian ML

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

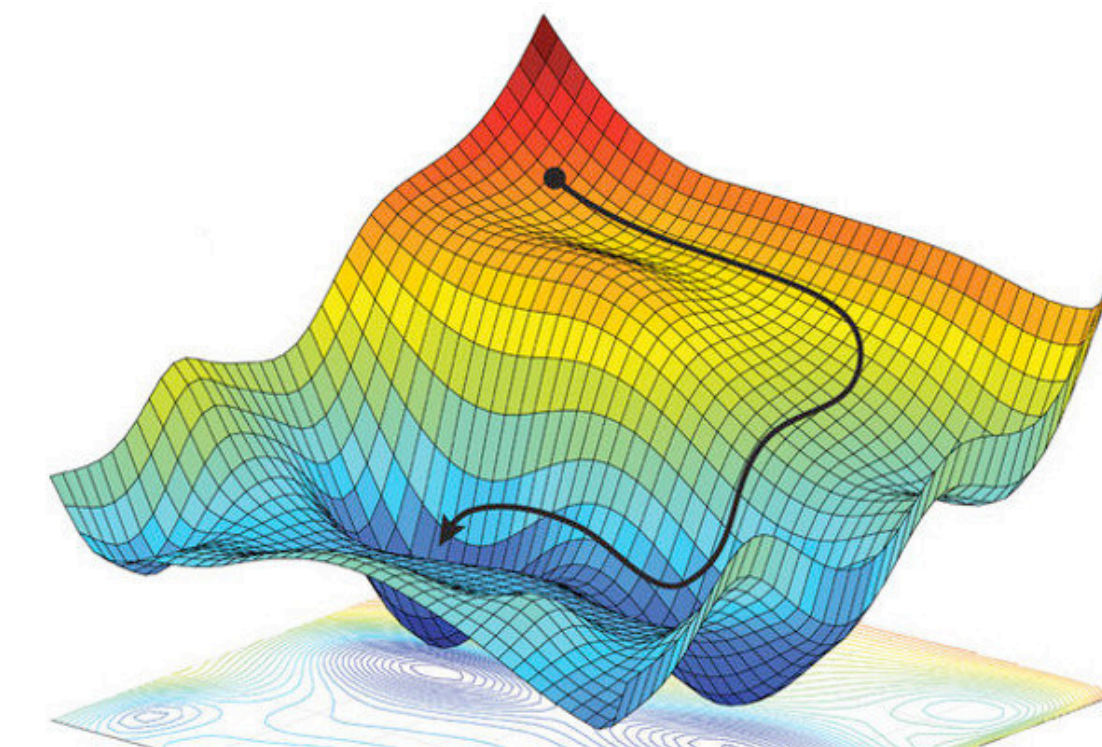


New Problem: Unclear how to compute $q_n^*(\theta)$ unless \mathcal{Q} is a parameterised set of distributions

New Algorithmic Solution: Infinite-dimensional Gradient Descent on $\mathcal{P}_2(\Theta)$

The Future of Post-Bayesian ML

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



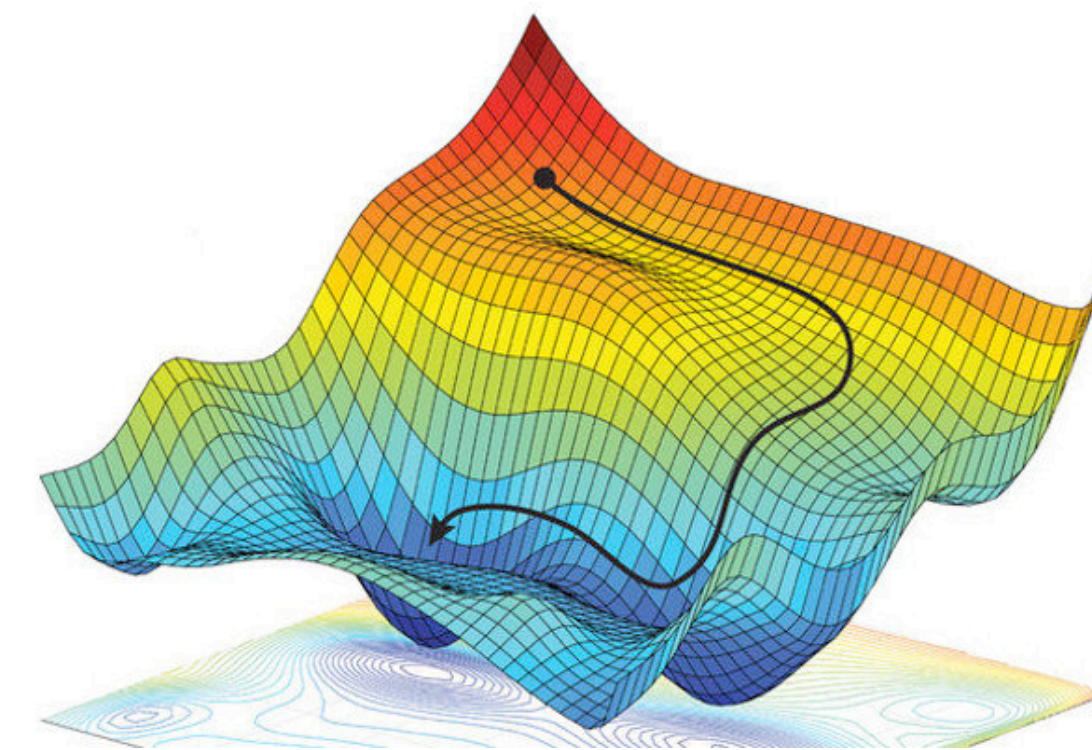
New Problem: Unclear how to compute $q_n^*(\theta)$ unless \mathcal{Q} is a parameterised set of distributions

New Algorithmic Solution: Infinite-dimensional Gradient Descent on $\mathcal{P}_2(\Theta)$

Implementation: Wasserstein Gradient Flow on $q \mapsto \int \mathbf{L}(x_{1:n}, \theta) q(\theta) d\theta + \mathbf{D}(q, \pi)$

The Future of Post-Bayesian ML

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



New Problem: Unclear how to compute $q_n^*(\theta)$ unless \mathcal{Q} is a parameterised set of distributions

New Algorithmic Solution: Infinite-dimensional Gradient Descent on $\mathcal{P}_2(\Theta)$

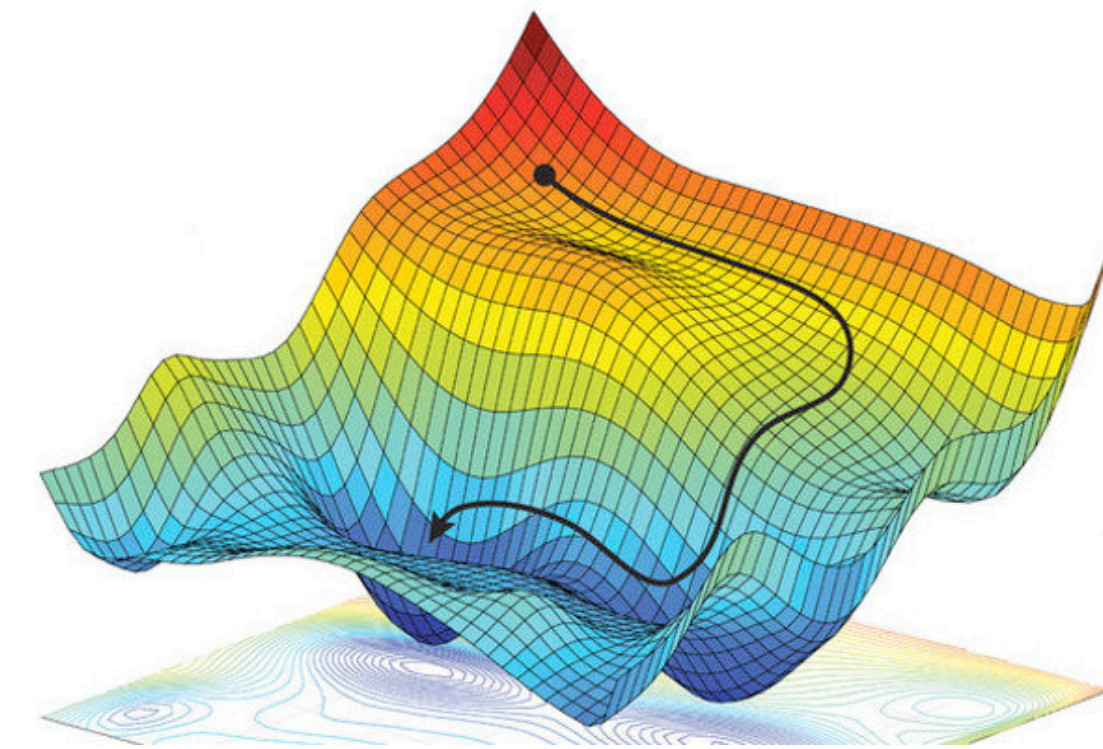
Implementation: Wasserstein Gradient Flow on $q \mapsto \int \mathbf{L}(x_{1:n}, \theta) q(\theta) d\theta + \mathbf{D}(q, \pi)$

Result: Sampling algorithm that doesn't need access to its target $q_n^*(\theta)$

First of its kind.

The Future of Post-Bayesian ML

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



New Problem: Unclear how to compute $q_n^*(\theta)$ unless \mathcal{Q} is a parameterised set of distributions

New Algorithmic Solution: Infinite-dimensional Gradient Descent on $\mathcal{P}_2(\Theta)$

Implementation: Wasserstein Gradient Flow on $q \mapsto \int \mathbf{L}(x_{1:n}, \theta) q(\theta) d\theta + \mathbf{D}(q, \pi)$

Result: Sampling algorithm that doesn't need access to its target $q_n^*(\theta)$

Bonus: recovers various deep learning algorithms + emblematic for why Post-Bayesian ML matters

First of its kind.

Will revisit this in a second.

The Future of Post-Bayesian ML

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

A Selection of unresolved challenges:

- How should we choose **D**?

The Future of Post-Bayesian ML

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

A Selection of unresolved challenges:

- How should we choose **D**?
- Can we characterise all divergence-based **L** that provide conjugate posteriors?

The Future of Post-Bayesian ML

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

A Selection of unresolved challenges:

- How should we choose **D**?
- Can we characterise all divergence-based **L** that provide conjugate posteriors?
- How should we **formalise** robustness to prior misspecification?

The Future of Post-Bayesian ML

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

A Selection of unresolved challenges:

- How should we choose **D**?
- Can we characterise all divergence-based **L** that provide conjugate posteriors?
- How should we **formalise** robustness to prior misspecification?
- How should we **choose between** different Post-Bayesian methods?

The Future of Post-Bayesian ML

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

A Selection of unresolved challenges:

- How should we choose **D**?
- Can we characterise all divergence-based **L** that provide conjugate posteriors?
- How should we **formalise** robustness to prior misspecification?
- How should we **choose between** different Post-Bayesian methods?
- (basic) asymptotics

The Future of Post-Bayesian ML

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

A Selection of unresolved challenges:

- How should we choose **D**?
- Can we characterise all divergence-based **L** that provide conjugate posteriors?
- How should we **formalise** robustness to prior misspecification?
- How should we **choose between** different Post-Bayesian methods?
- (basic) asymptotics
- What is the overlap / difference with Martingale Posteriors?
- ...

... But rather than expanding on technical challenges, let me end with a morality tale

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

3

Morality Tale: Why Post-Bayesian ML matters

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

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

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

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

3

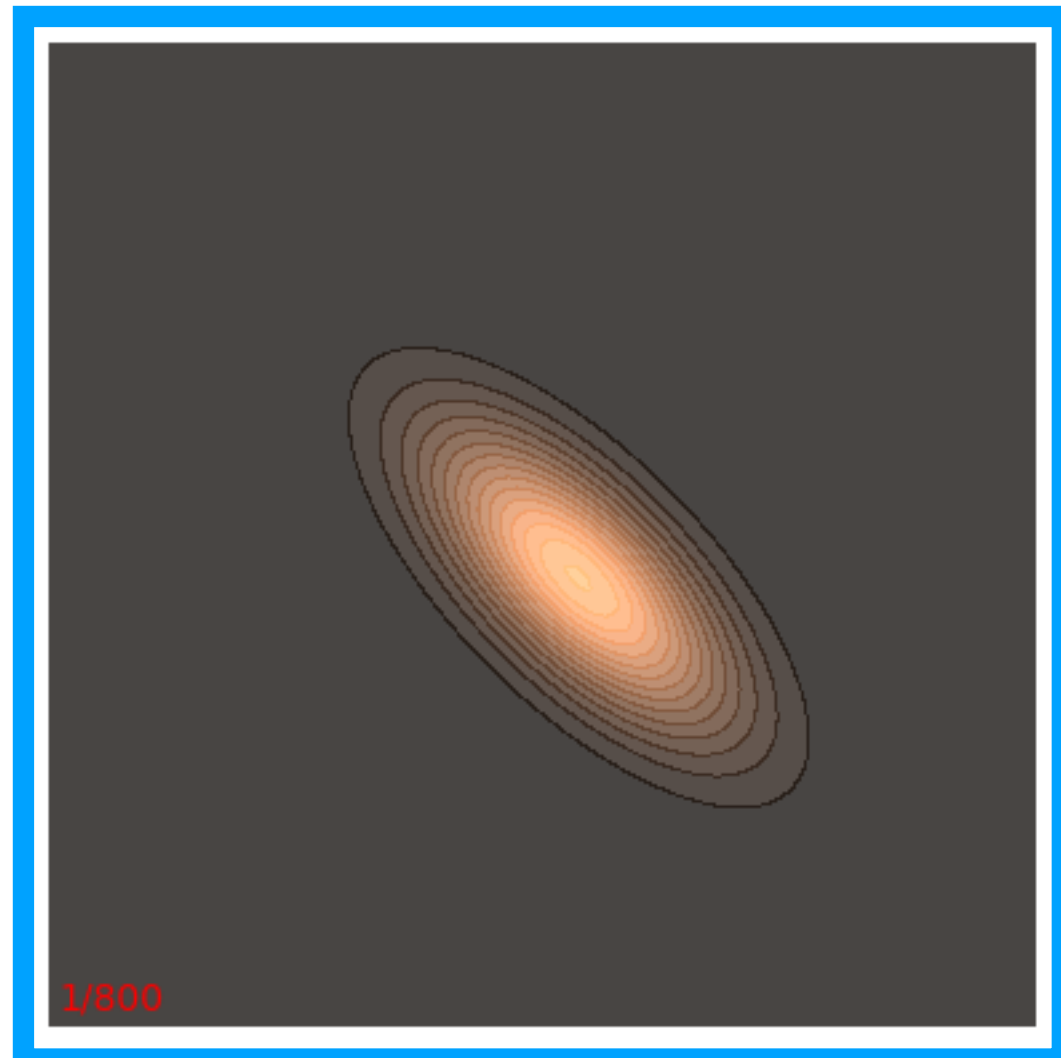
Morality Tale: Why Post-Bayesian ML matters

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

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

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

Wasserstein Gradient Flow = Langevin Diffusion



Converges to
well-defined density

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

Morality Tale: Why Post-Bayesian ML matters

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

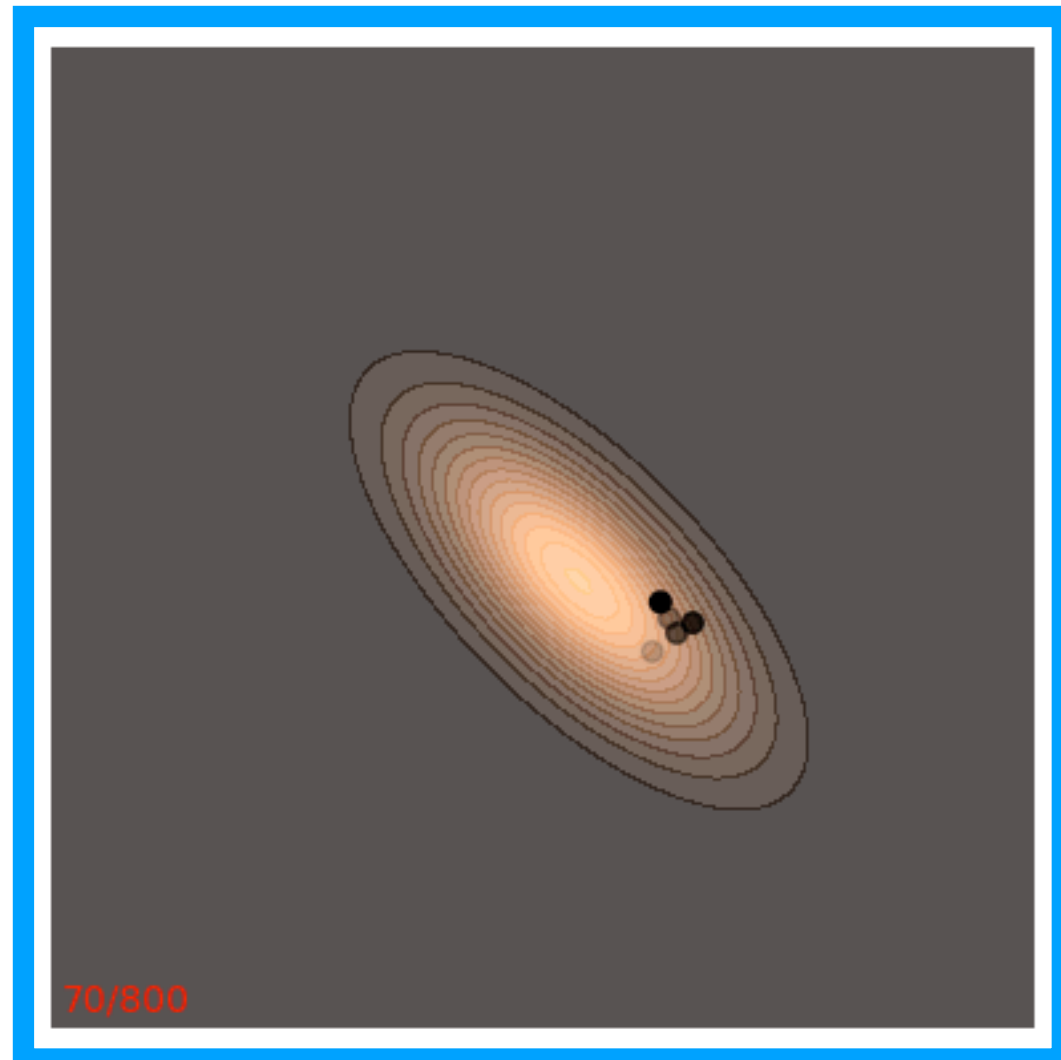
Objective: $q \mapsto \mathbb{E}_{\theta \sim q} \left[-\log p(x_{1:n} | \theta)^\lambda \right] + \mathbf{KL}(q, \pi) \xrightarrow{\lambda \rightarrow \infty} q \mapsto \mathbb{E}_{\theta \sim q} \left[-\log p(x_{1:n} | \theta) \right]$

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

Deep Ensemble (DE) ($\lambda \rightarrow \infty$)

Wasserstein Gradient Flow = Langevin Diffusion

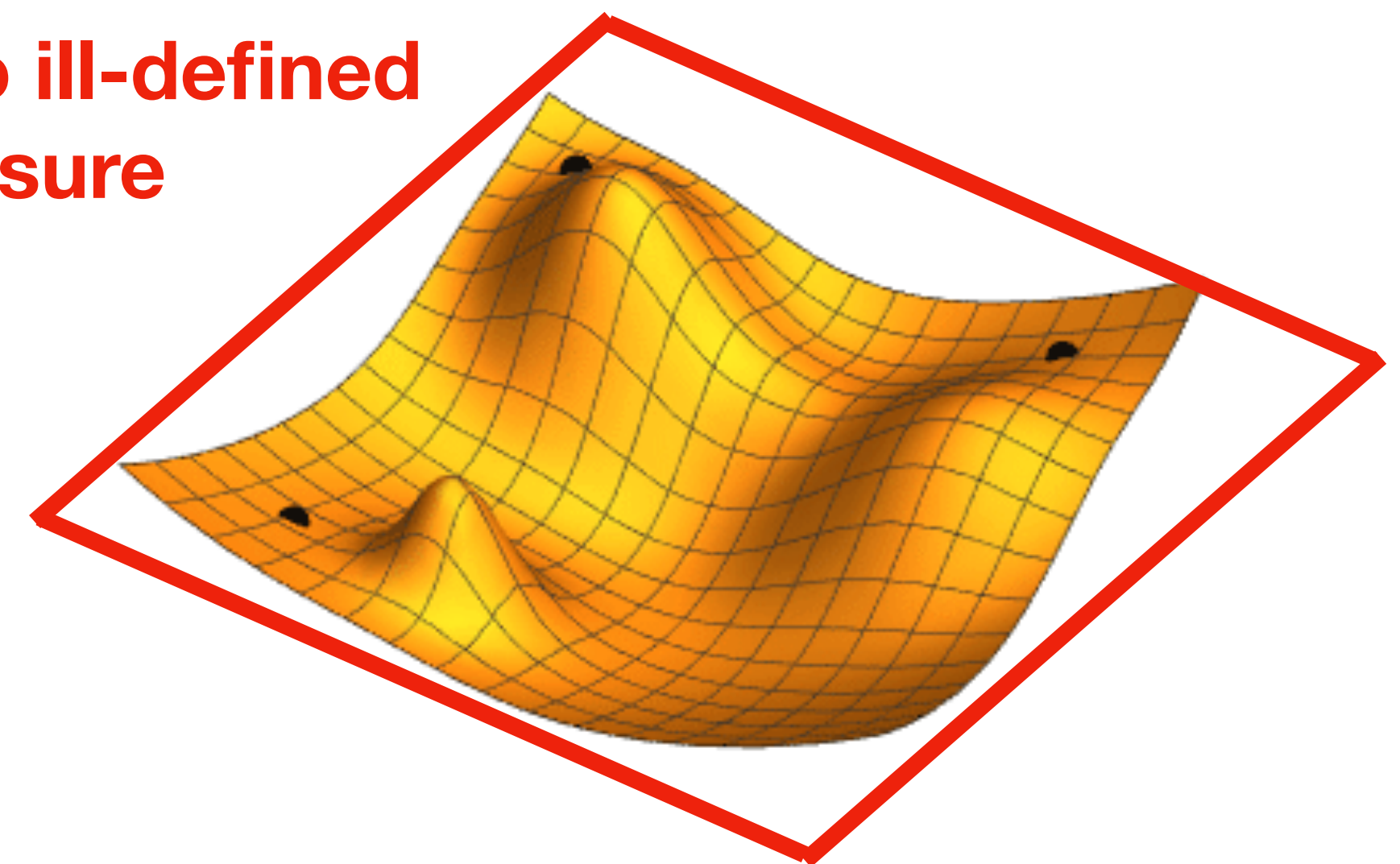
Wasserstein Gradient Flow = DE



Converges to well-defined density

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

Converges to ill-defined discrete measure



Morality Tale: Why Post-Bayesian ML matters

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 by a group specialised in Bayesian ML (2020) @ NeurIPS
(cited > 650 times according to Google scholar)**

Morality Tale: Why Post-Bayesian ML matters

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 by a group specialised in Bayesian ML (2020) @ NeurIPS
(cited > 650 times according to Google scholar)**

Unfortunately, as we just saw this is not correct.



Conclusion: Why Post-Bayesian ML matters

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

Conclusion: Why Post-Bayesian ML matters

- I. In practice, orthodox Bayesian ML has already been abandoned
(Bayes posterior: **prior regulariser, densities** ;
Deep Ensembles: **no prior regulariser, discrete measures**)
- II. But as a field, we often don't pay enough attention to the ramifications,
and lag behind in developing coherent Post-Bayesian formalisms
and a mathematical grammar to talk about them.

Conclusion: Why Post-Bayesian ML matters

- I. In practice, orthodox Bayesian ML has already been abandoned
(Bayes posterior: **prior regulariser, densities** ;
Deep Ensembles: **no prior regulariser, discrete measures**)
- II. But as a field, we often don't pay enough attention to the ramifications,
and lag behind in developing coherent Post-Bayesian formalisms
and a mathematical grammar to talk about them.
- III. This in turns leads to incorrect claims and conclusions.
(**'Deep Ensembles are Bayesian'**)

References: Post-Bayesian ML

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

1 Foundations

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

$$\pi_n^L(\theta | x_{1:n})$$

model missp
~~(A1)~~ (A2)

@ ISBA 2024



Knoblauch & Damoulas (2018); ICML
 Knoblauch, Jewson, & Damoulas (2018); NeurIPS
 Frazier*, Knoblauch*, & Drovandi (2024); preprint
 McLatchie, Fong, Frazier, & Knoblauch (2024); forthcoming

2 State of the Art

$$\pi_n^L(\theta | x_{1:n})$$

model misspecification +
 computational
~~(A1)~~ (A2)

@ ICML 2024



Schmon, Cannon, & Knoblauch (2020); AABI
 Matsubara, Knoblauch, Briol, & Oates (2022); JRSS-B
 Dellaporta, Knoblauch, Damoulas, & Briol (2022);
 AISTATS (best paper award)
 Altamirano, Briol, & Knoblauch (2023); ICML
 Altamirano, Briol, & Knoblauch (2024); ICML (spotlight)
 Duran-Martin, Altamirano, Shestopaloff, Sanchez-Betancourt,
 Knoblauch, Briol, & Murphy (2024); ICML

3 The Future

$$q_n^*(\theta)$$

model missp
 prior misspe
 compu
~~(A1)~~ ~~(A2)~~ ~~(A3)~~

Husain & Knoblauch (2022); ALT
 Knoblauch, Jewson, & Damoulas (2022); JMLR
 Matsubara, Knoblauch, Briol, & Oates (2023); JASA
 Wild, Sejdinovic, & Knoblauch (2024); forthcoming
 Wild, Ghalebikesabi, Sejdinovic, & Knoblauch (2023);
 NeurIPS (oral)

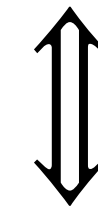
Foundations



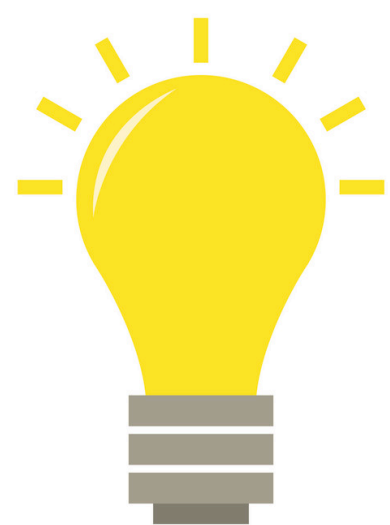
Mathematical
Foundations

1930: **DeFinetti's Representation Theorem** [cf. Hewitt & Savage (1955), Diaconis & Freedman (1984, 1987)]

For all $k \leq n$ and all permutations σ , $p(x_1, x_2, \dots, x_k) \stackrel{\forall \sigma}{=} p(x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(k)})$



There is a parameter space Θ and $\pi(\theta)$ s.t. $p(x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(k)}) \stackrel{\forall \sigma}{=} \int \prod_{i=1}^k p(x_i | \theta) \pi(\theta) d\theta$



\implies If data exchangeable, there MUST be

- 1) model $p(\cdot | \theta)$
- 2) prior $\pi(\theta)$

that represent the data through a Bayesian approach.

Justification for
Bayesianism

Foundations



Mathematical Foundations

1930: **DeFinetti's Representation Theorem** [cf. Hewitt & Savage (1955), Diaconis & Freedman (1984, 1987)]

For all $k \leq n$ and all permutations σ , $p(x_1, x_2, \dots, x_k) \stackrel{\forall \sigma}{=} p(x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(k)})$



There is a parameter space Θ and $\pi(\theta)$ s.t. $p(x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(k)}) \stackrel{\forall \sigma}{=} \int \prod_{i=1}^k p(x_i | \theta) \pi(\theta) d\theta$

\implies **Problem:** Does NOT tell you what $\pi(\theta)$ and $p(\cdot | \theta)$ are!



1) π generally depends on $p(\cdot | \theta)$ and n

2) θ generally ∞ -dimensional

[e.g. $p(x_i | \theta) = \theta(x_i)$ for θ a probability density on x_i]

Not how we practice Bayesianism

Foundations

Mathematical Foundations



1954: **Savage Axioms** connects Bayes' Theorem to **Decision Theory**

Actions $\mathcal{A} = \{a : \text{States} \longrightarrow \text{Consequences}\}$ $a_2 \lesssim a_1 \iff a_1$ preferred to a_2

\lesssim satisfies Savage Axioms on \mathcal{A}



\exists utility function $u : \text{Consequences} \rightarrow \mathbb{R}$ and π on States s.t.

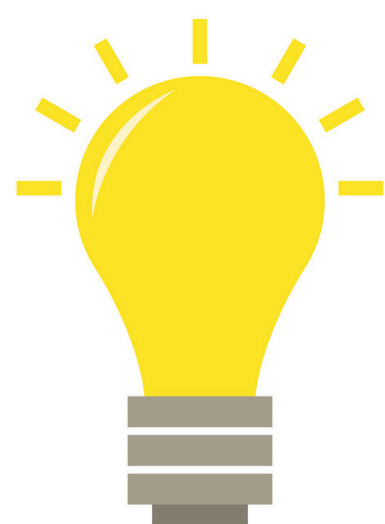
$$\forall a_1, a_2 \in \mathcal{A} : a_1 \lesssim a_2 \iff \int u[a_1(s)] \pi(s) ds \leq \int u[a_2(s)] \pi(s) ds$$

$$\forall a_1, a_2 \in \mathcal{A} : a_1 \lesssim a_2 \text{ given } x_{1:n} \iff \int u[a_1(s)] \pi_n(s | x_{1:n}) ds \leq \int u[a_2(s)] \pi_n(s | x_{1:n}) ds$$

Foundations



Mathematical Foundations



- ⇒ Prior $\pi(s)$ = beliefs IMPLIED by rational agent's preferences
- ⇒ Bayes' Posterior $\pi_n(s | x_{1:n})$ = rational agent's belief update given data

\lesssim satisfies Savage Axioms on \mathcal{A}



∃ utility function u : Consequences $\rightarrow \mathbb{R}$ and π on States s.t.

$$\forall a_1, a_2 \in \mathcal{A} : a_1 \lesssim a_2 \iff \int u[a_1(s)] \pi(s) ds \leq \int u[a_2(s)] \pi(s) ds$$

$$\forall a_1, a_2 \in \mathcal{A} : a_1 \lesssim a_2 \text{ given } x_{1:n} \iff \int u[a_1(s)] \pi_n(s | x_{1:n}) ds \leq \int u[a_2(s)] \pi_n(s | x_{1:n}) ds$$

Foundations



Mathematical Foundations

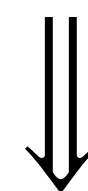


Problem: Prior $\pi(s)$ is defined on **States** (not parameters θ !)

\Rightarrow parameter space $\Theta =$ relevant **State** $\iff x_{1:n} \sim p_{\theta^*}(x_{1:n})$ for some $\theta^* \in \Theta$

\Rightarrow Does NOT tell you what $\pi(\theta)$ and $p(\cdot | \theta)$ are!

\lesssim satisfies Savage Axioms on \mathcal{A}



\exists utility function $u : \text{Consequences} \rightarrow \mathbb{R}$ and π on **States** s.t.

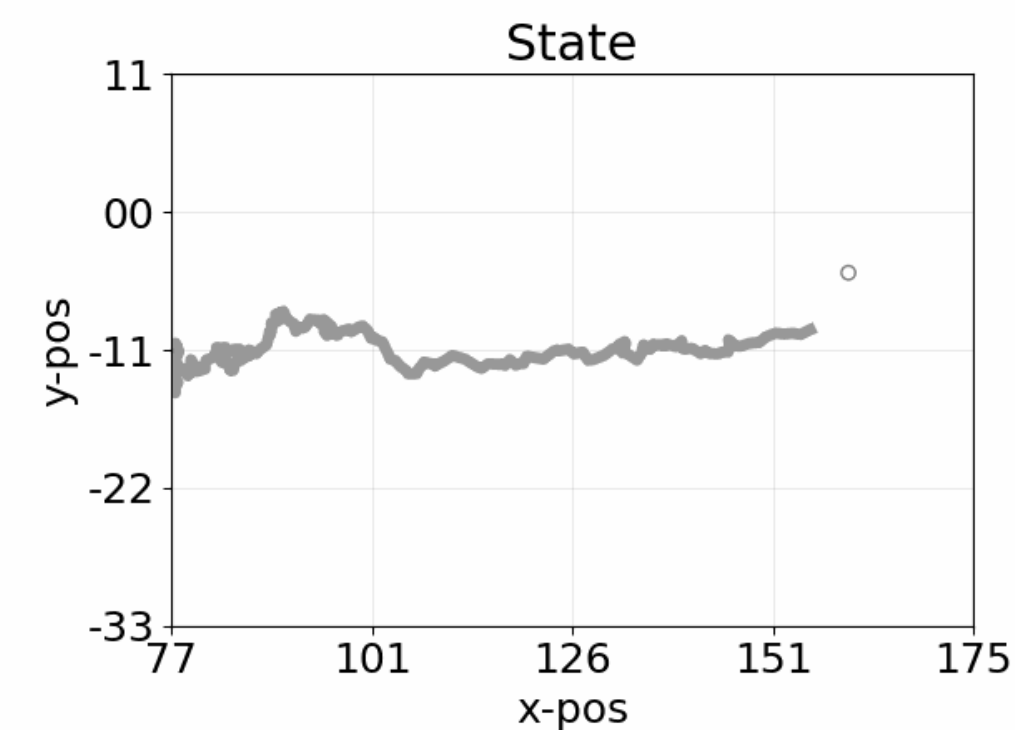
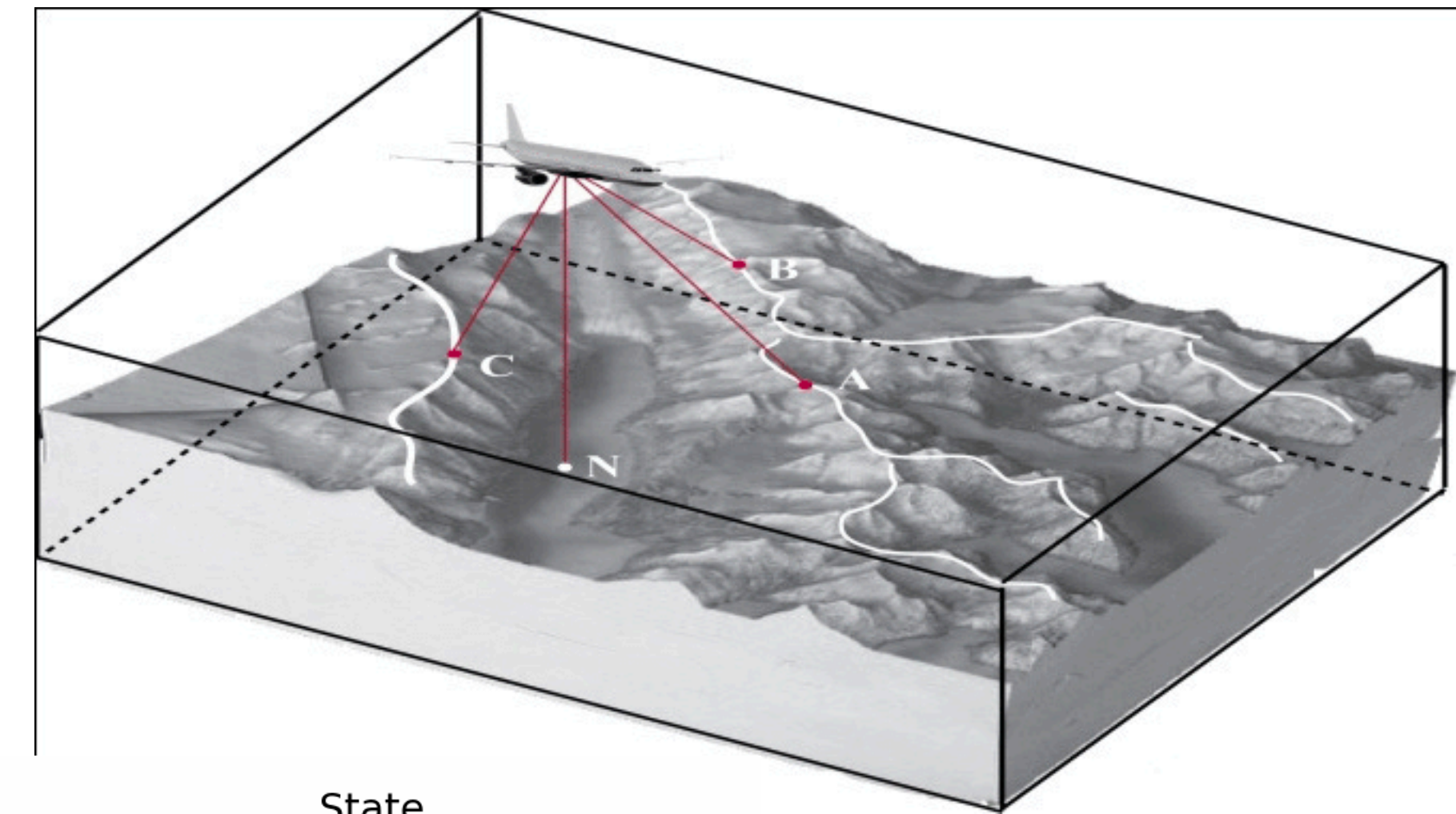
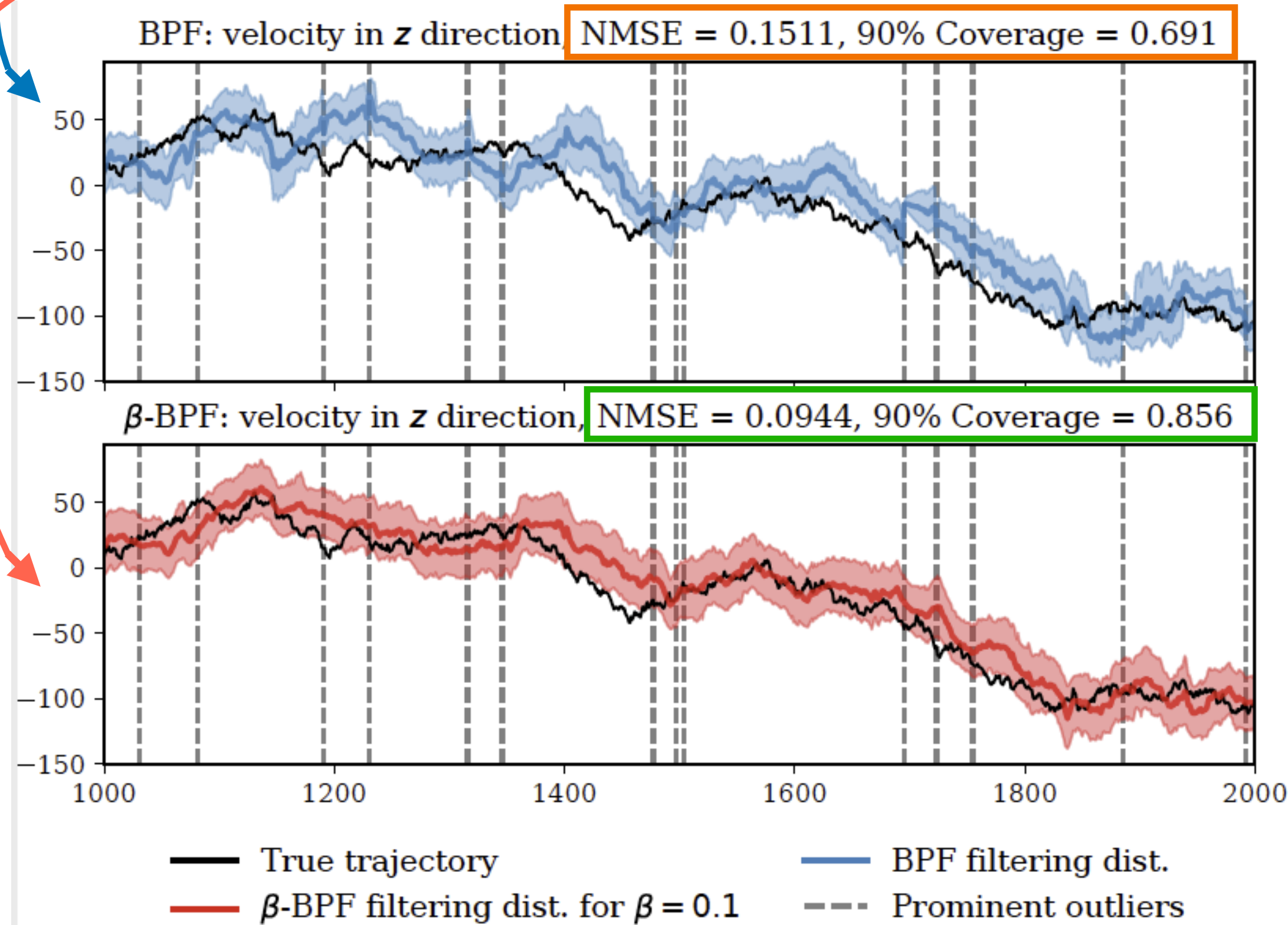
$$\forall a_1, a_2 \in \mathcal{A} : a_1 \lesssim a_2 \iff \int u[a_1(s)] \pi(s) ds \leq \int u[a_2(s)] \pi(s) ds$$

$$\forall a_1, a_2 \in \mathcal{A} : a_1 \lesssim a_2 \text{ given } x_{1:n} \iff \int u[a_1(s)] \pi_n(s | x_{1:n}) ds \leq \int u[a_2(s)] \pi_n(s | x_{1:n}) ds$$

Post-Bayesian ML: Success Stories

Standard Kalman Filter for Terrain Aided Navigation (Drone over Terrain Map)

Robustified version



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

Q1: Can tuning λ improve Robustness?

Question: What is the predictively optimal λ ?

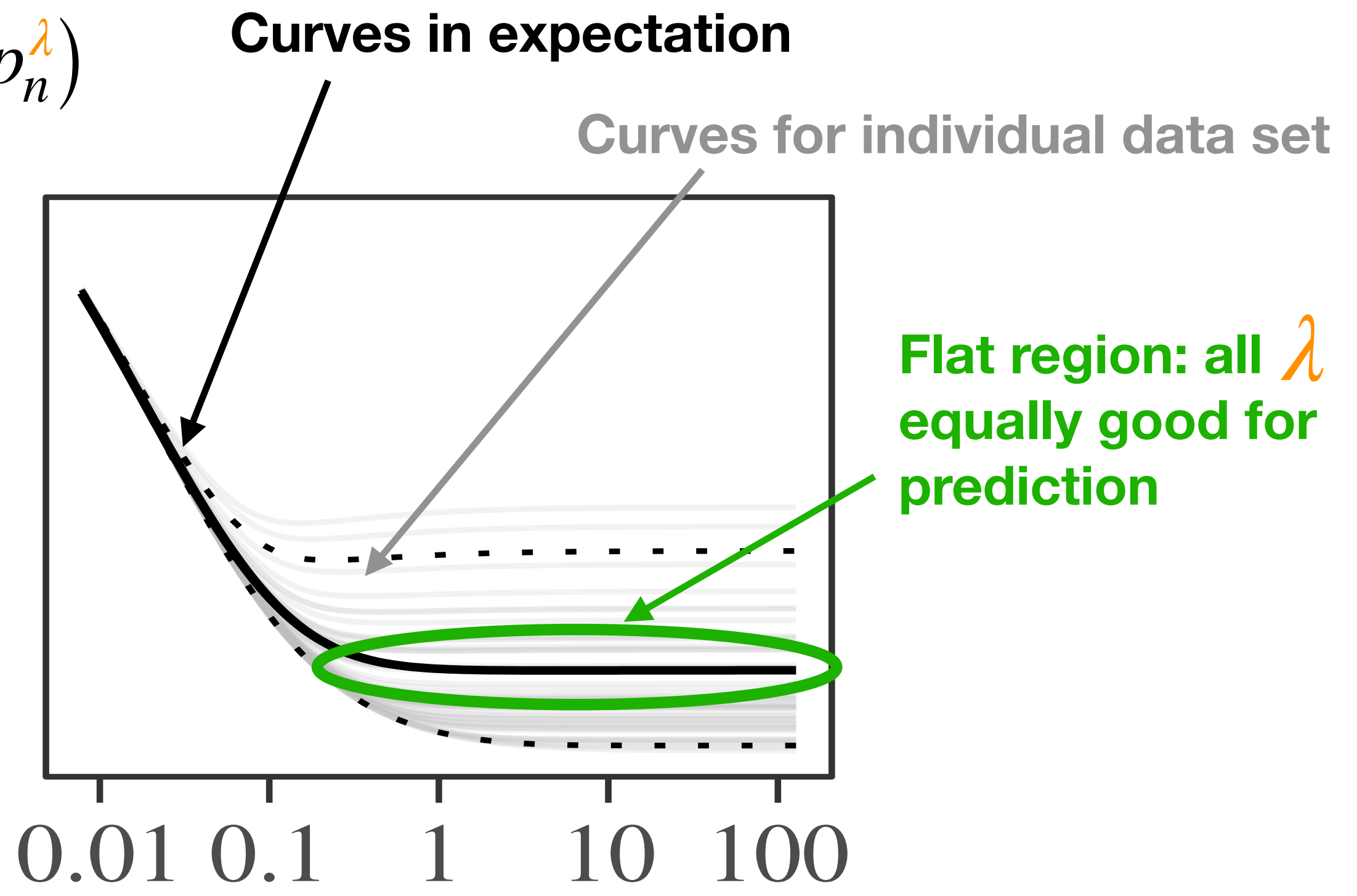
Posterior predictive = $p_n^\lambda(z) = \int p(z | \theta) \pi_n^{(\lambda)}(\theta | x_{1:n}) d\theta$

Predictively optimal λ : $\lambda^* = \operatorname{argmin}_{\lambda > 0} D_{\text{TV}}(q, p_n^\lambda)$

Data-generating density: $x_{1:n} \sim q(x_{1:n})$

$D_{\text{TV}}(q, p_n^\lambda)$

Theorem: these curves will always look that way.



$\pi_n^{(0)} = \pi \leftarrow 0 \leftarrow \lambda \quad \lambda \rightarrow \infty \rightarrow \pi_n^{(\infty)} = \text{Dirac at MLE}$

What λ leads to Robustness?

Findings:

- (1) Ill-defined problem: minimiser λ^* doesn't exist
- (2) Flat region: infinitely many choices yield (nearly) same predictive
- (3) Predictively, there is no advantage over MLE / point estimators

Q: Why does this happen?

p^* = oracle predictive induced by θ^*

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} \lim_{n \rightarrow \infty} n^{-1} - \log p(x_{1:n} | \theta)$$

p_n^∞ = MLE predictive induced by $\hat{\theta}_n$

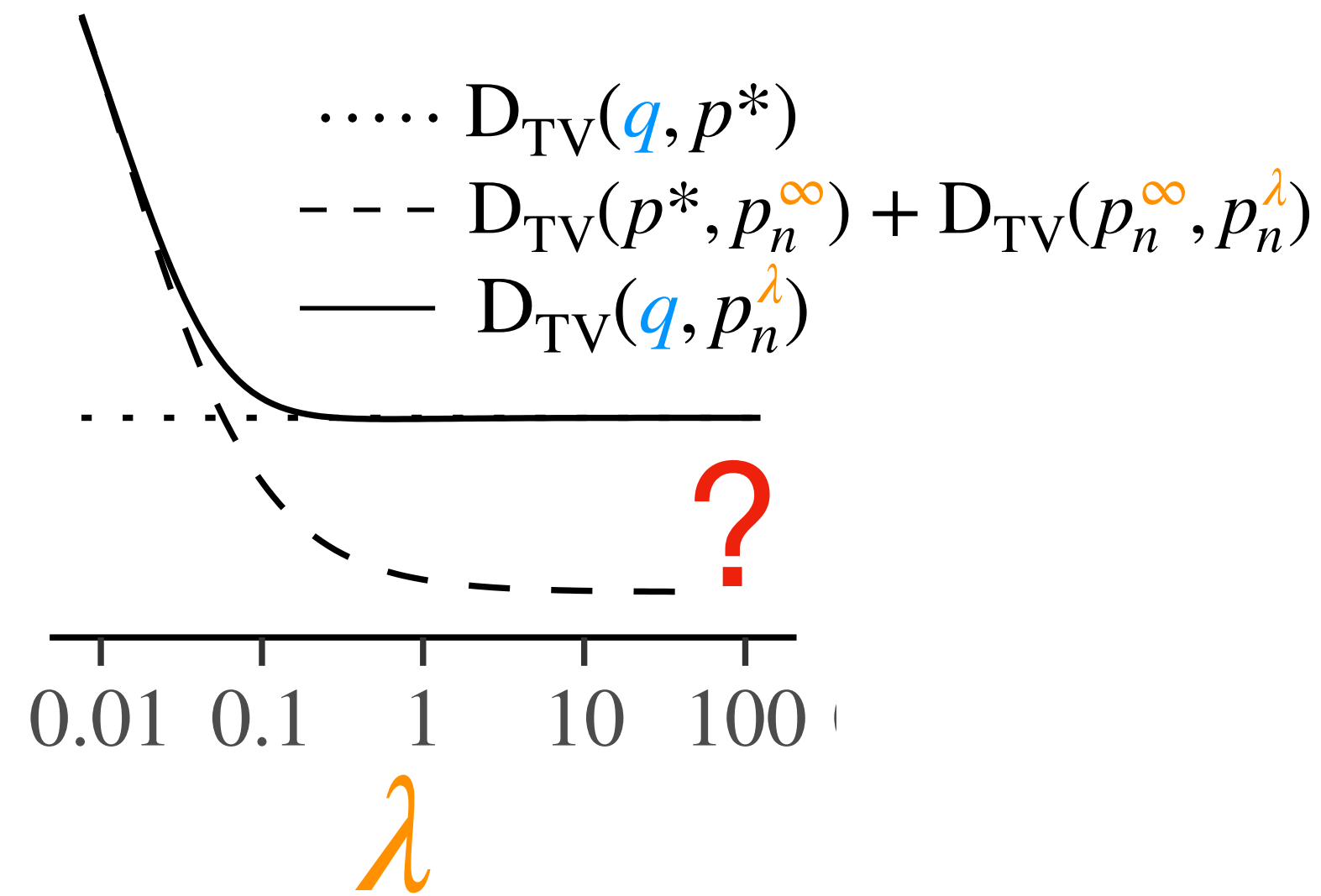
$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} - \log p(x_{1:n} | \theta)$$

$$\begin{aligned}
 D_{\text{TV}}(q, p_n^\lambda) &\leq \underbrace{D_{\text{TV}}(q, p^*)}_{\text{(irreducible error)} = \text{constant}} + \underbrace{D_{\text{TV}}(p^*, p_n^\infty)}_{\text{(error of MLE/point estimator)} \lesssim \nu_n \text{ if MLE converges at rate } \nu_n} + \underbrace{D_{\text{TV}}(p_n^\infty, p_n^\lambda)}_{\text{(difference MLE vs } p_n^\lambda) \lesssim \varepsilon_n \text{ if } \pi_n^{(\lambda)} \text{ concentrates at rate } \varepsilon_n} \\
 &\text{w.h.p}
 \end{aligned}$$

In practice / experimentally:

- goes to 0 MUCH faster than $\nu_n + \varepsilon_n$
- even works WITHOUT concentration and consistency

What λ leads to Robustness?



Similar for CV and $D = \text{KL}$ ✓

Exponentially fast ?

$$\boxed{D_{\text{TV}}(q, p_n^\lambda)} \stackrel{\text{w.h.p}}{\leq} \underbrace{D_{\text{TV}}(q, p^*)}_{\text{(irreducible error) = constant}} + \underbrace{\left| D_{\text{TV}}(p^*, p_n^\infty) + D_{\text{TV}}(p_n^\infty, p_n^\lambda) \right|}_{\substack{\lesssim \nu_n \text{ if MLE converges at rate } \nu_n \\ \lesssim \varepsilon_n \text{ if } \pi_n^{(\lambda)} \text{ concentrates at rate } \varepsilon_n}}$$

- ? In practice / experimentally:
- goes to 0 MUCH faster than $\nu_n + \varepsilon_n$
 - even works WITHOUT concentration and consistency

What λ leads to Robustness?

$$\pi_n^{(\lambda)}(\theta | x_{1:n}) = \frac{p(x_{1:n} | \theta)^\lambda \cdot \pi(\theta)}{\int p(x_{1:n} | \theta)^\lambda \cdot \pi(\theta) d\theta} = \operatorname{argmin}_{q \in \mathcal{P}(\Theta)} \left\{ \lambda \cdot \int -\log p(x_{1:n}, \theta) q(\theta) d\theta + \operatorname{KL}(q, \pi) \right\}$$

$0 < \lambda \ll 1$



$\lambda \gg 1$

How to square with **our finding** that λ is largely irrelevant?

Grünwald (2012); ALT
 Holmes & Walker (2017); Biometrika
 Miller & Dunson (2018); JRSS-B
 Bhattacharya, Pati, & Yang (2019); Ann. Statist.

Wenzel et al. (2020); ICML
 Adlam et al. (2020); preprint
 Noci et al. (2021); NeurIPS
 Aitchison (2021); ICLR

essential	Parameter uncertainty	incidental
incidental	Predictive uncertainty	essential
$n > d$	Model complexity	$d \ll n$
good	Prior quality	bad
theory	Evidence	Experimental

Prior quality should be VERY important

What L leads to Robustness?

Ghosh & Basu (2016); AISM
 K., Jewson, & Damoulas (2018); NeurIPS
 Boustati, Akyildiz, Damoulas, & Johansen (2020); NeurIPS
 K., Jewson, & Damoulas (2022); JMLR
 Frazier*, K.*, & Drovandi (2024); preprint

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

Q2: Which discrepancies $D(q, p(\cdot | \theta))$ should we construct $L(x_{1:n}, \theta)$ from?

Initial Work:

$$D^\beta(q, p(\cdot | \theta)) = \underbrace{\int p(z | \theta)^{1+\beta} dz}_{\text{Generally intractable...}} - \frac{1+\beta}{\beta} \underbrace{\int p(z | \theta)^\beta q(z) dz}_{\approx \frac{1}{n} \sum_{i=1}^n p(x_i | \theta)^\beta} + \frac{1}{\beta} \underbrace{\int q^{1+\beta}(z) dz}_{\text{Independent of } \theta}$$

Generally intractable...
 (and hard to approximate)

Sometimes: tractable as $A(\theta, \beta)$

$$\approx \frac{1}{n} \sum_{i=1}^n p(x_i | \theta)^\beta \quad \text{Independent of } \theta$$

Issues:



- Intractability
- Numerical instability
- Scale dependence

$$L(x_{1:n}, \theta) = n \cdot A(\theta, \beta) + \frac{\beta + 1}{\beta} \sum_{i=1}^n \underbrace{p(x_i | \theta)^\beta}_{\text{scale}} \xrightarrow{\beta \downarrow 0} \sum_{i=1}^n -\log p(x_i | \theta)$$

What **L** leads to Robustness?

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

Futami, Sato & Sugiyama (2018); AISTATS
 Cherrief-Abdellatif & Alquier (2020); AABI
 Pacchiardi, Khoo, & Dutta (2021); preprint
 Frazier*, K.*, & Drovandi (2024); preprint

Q2: Which discrepancies $D(q, p(\cdot \mid \theta))$ should we construct $\mathbf{L}(x_{1:n}, \theta)$ from?

Other Proposals: $D^\gamma(q, p(\cdot \mid \theta)) = \underbrace{\log \int p(z \mid \theta)^{1+\gamma} dz}_{\text{Generally intractable... (and hard to approximate)}} - \frac{1+\gamma}{\gamma} \underbrace{\log \int p(z \mid \theta)^\gamma q(z) dz}_{\approx \frac{1}{n} \sum_{i=1}^n p(x_i \mid \theta)^\gamma} + \frac{1}{\gamma} \log \int q^{1+\gamma}(z) dz$

(γ -Divergence)

Like β -divergence
Include Hellinger Divergence as well!

Bias!

Independent of θ

$D_k^2(q, p(\cdot \mid \theta)) = \underbrace{\mathbb{E}_{x \sim q(x), x' \sim q(x')} [k(x, x')]}_{\text{Independent of } \theta} - 2 \underbrace{\mathbb{E}_{x \sim q(x), x' \sim p(x|\theta)} [k(x, x')]}_{\text{Intractable (but easy to approximate!)}} + \underbrace{\mathbb{E}_{x \sim p(x|\theta), x' \sim p(x|\theta)} [k(x, x')]}_{\text{Intractable (but easy to approximate!)}}$

(MMD²)

Intractable Integrals
Not defined for conditional models
 (But useful for intractable/simulation models)

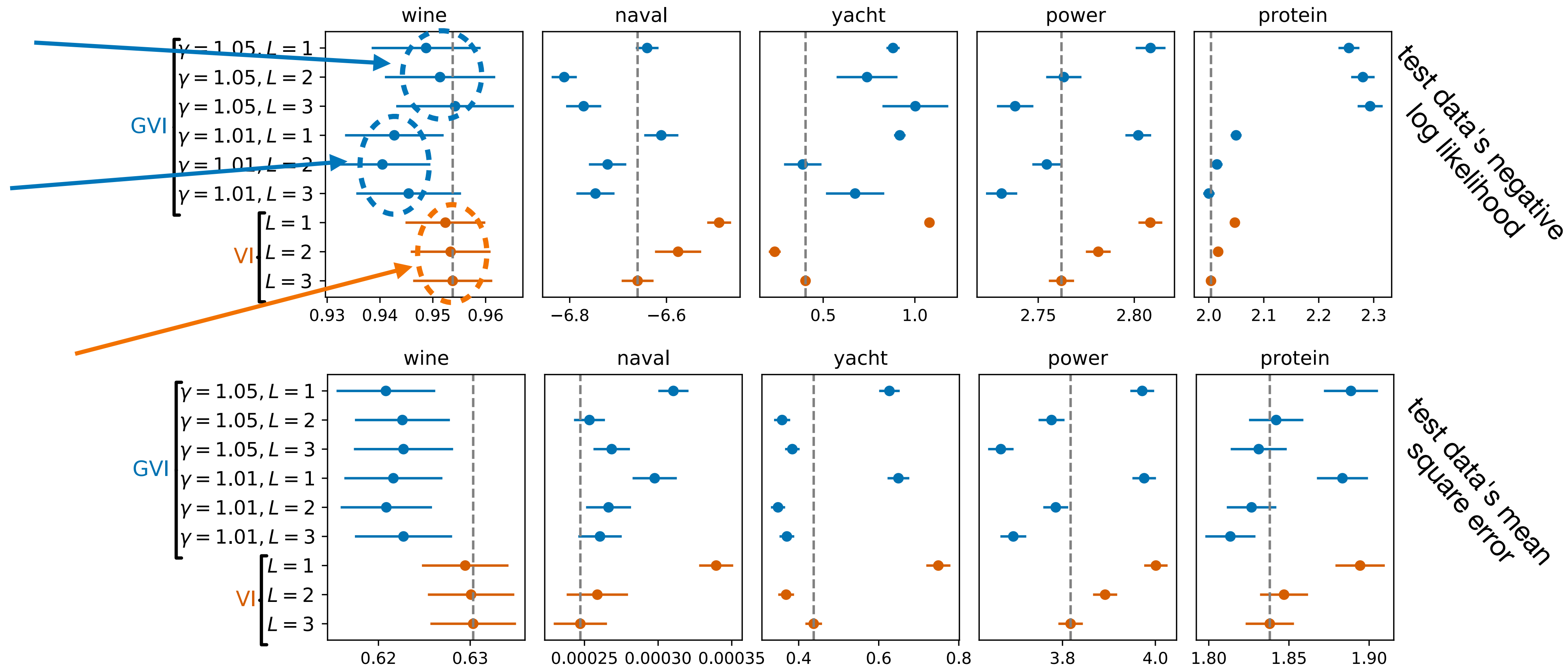
What **L** leads to Robustness?

Example 2: Deep Gaussian Processes (γ -Divergence)

$L(x_{1:n} | \theta) \approx D_\gamma(p_0 || p_\theta)$
 (γ -divergence [$\gamma = 0.05$])

$L(x_{1:n} | \theta) \approx D_\gamma(p_0 || p_\theta)$
 (γ -divergence [$\gamma = 0.01$])

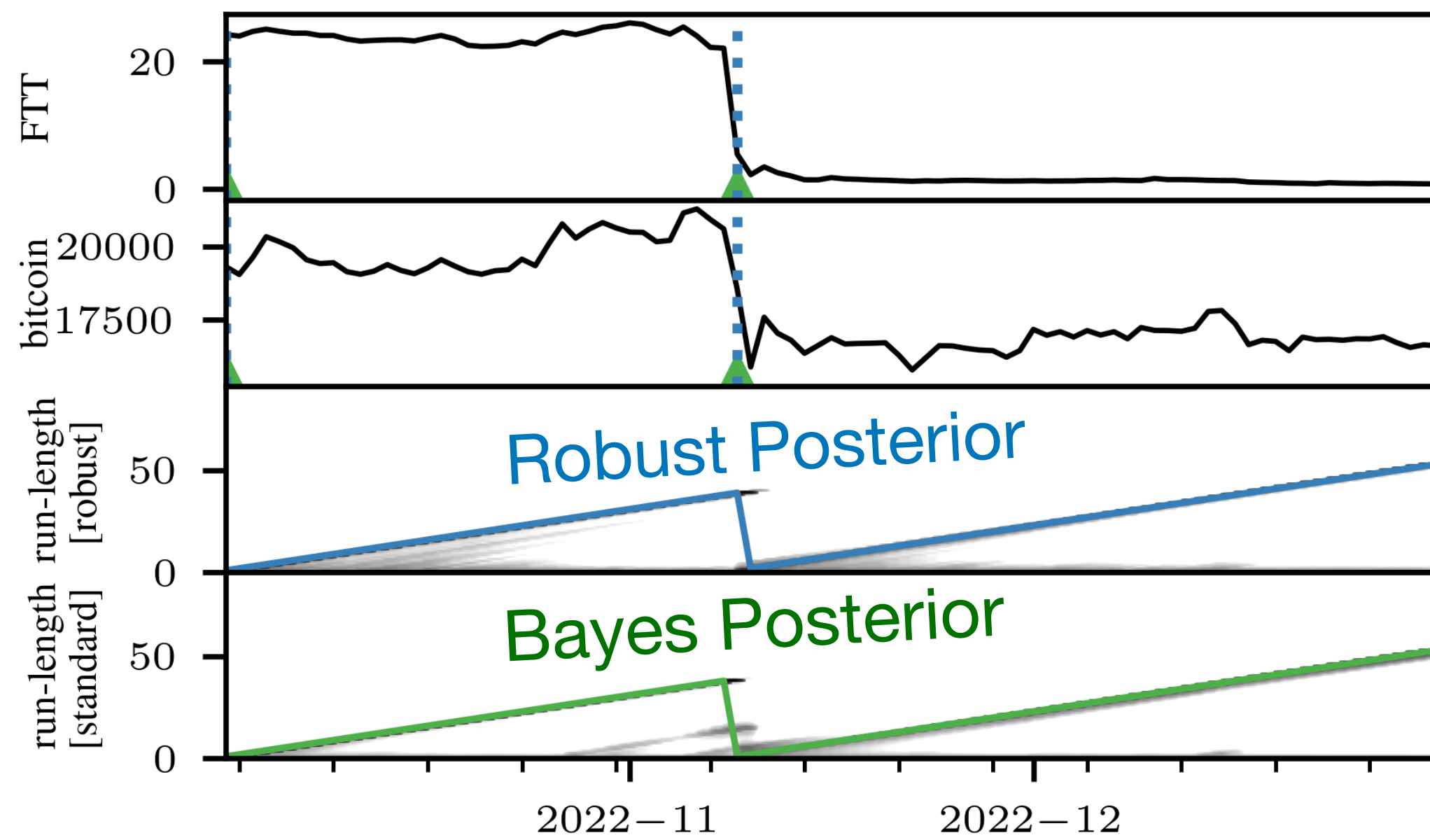
$-\log p(x_{1:n} | \theta) \approx \text{KLD}(p_0 || p_\theta)$



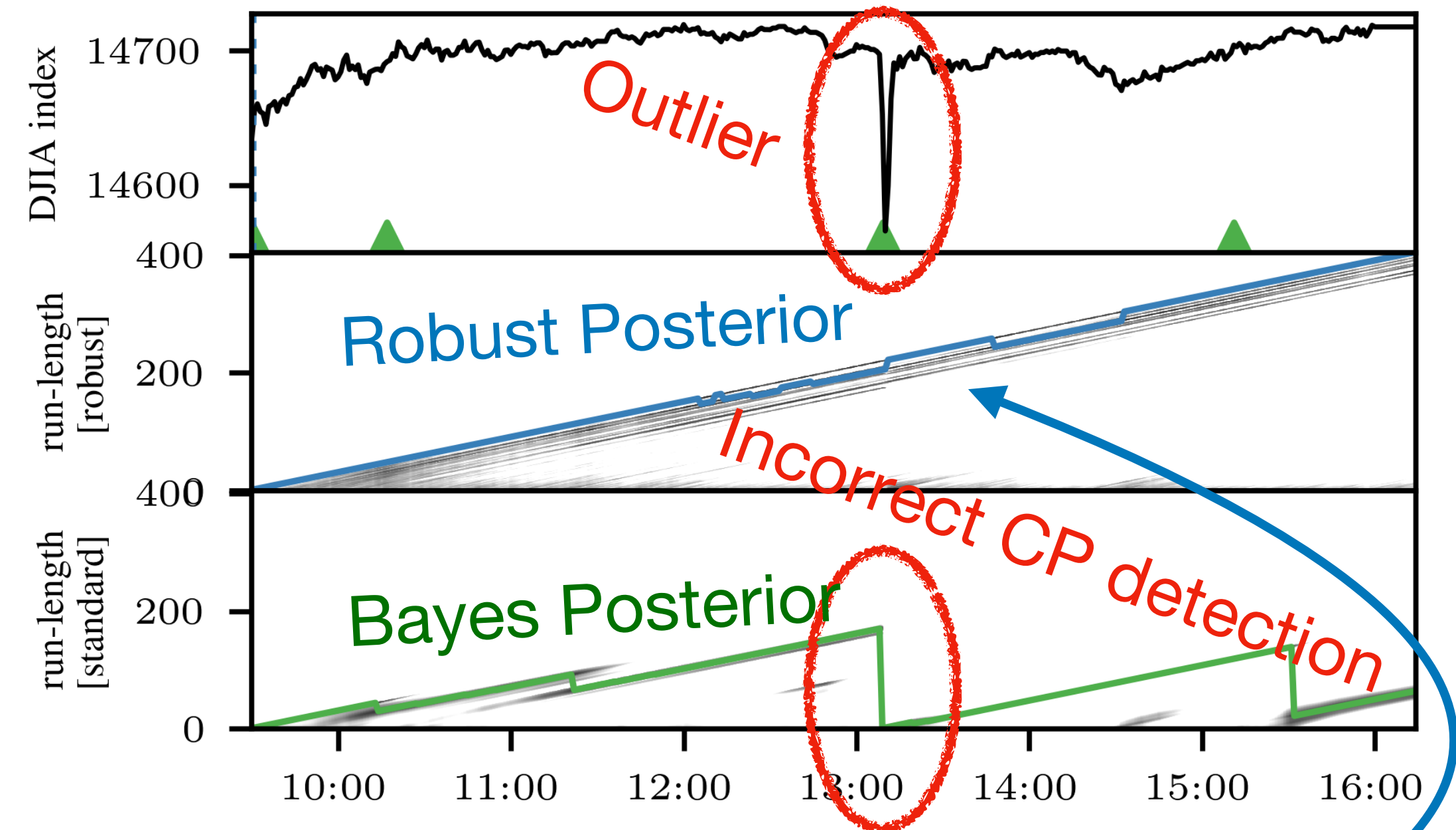
Q2: What **L** leads to Robustness?

Example 1: Bayesian On-line Changepoint Detection (β -Divergence)

FTX / cryptocurrency crash



Twitter flash crash



K. & Damoulas (2018); ICML
K., Jewson, & Damoulas (2018); NeurIPS

Theorem: the robust algorithm cannot declare a change point after a single outlier.

Approximating L

Q3: How does approximating $L(x_{1:n}, \theta)$ affect $\pi_n^L(\theta | x_{1:n})$?

Setting: $L(x_{1:n}, \theta)$ involves **intractable components**

(like integrals $I(\theta) = \int f(u) p(u | \theta) du$)

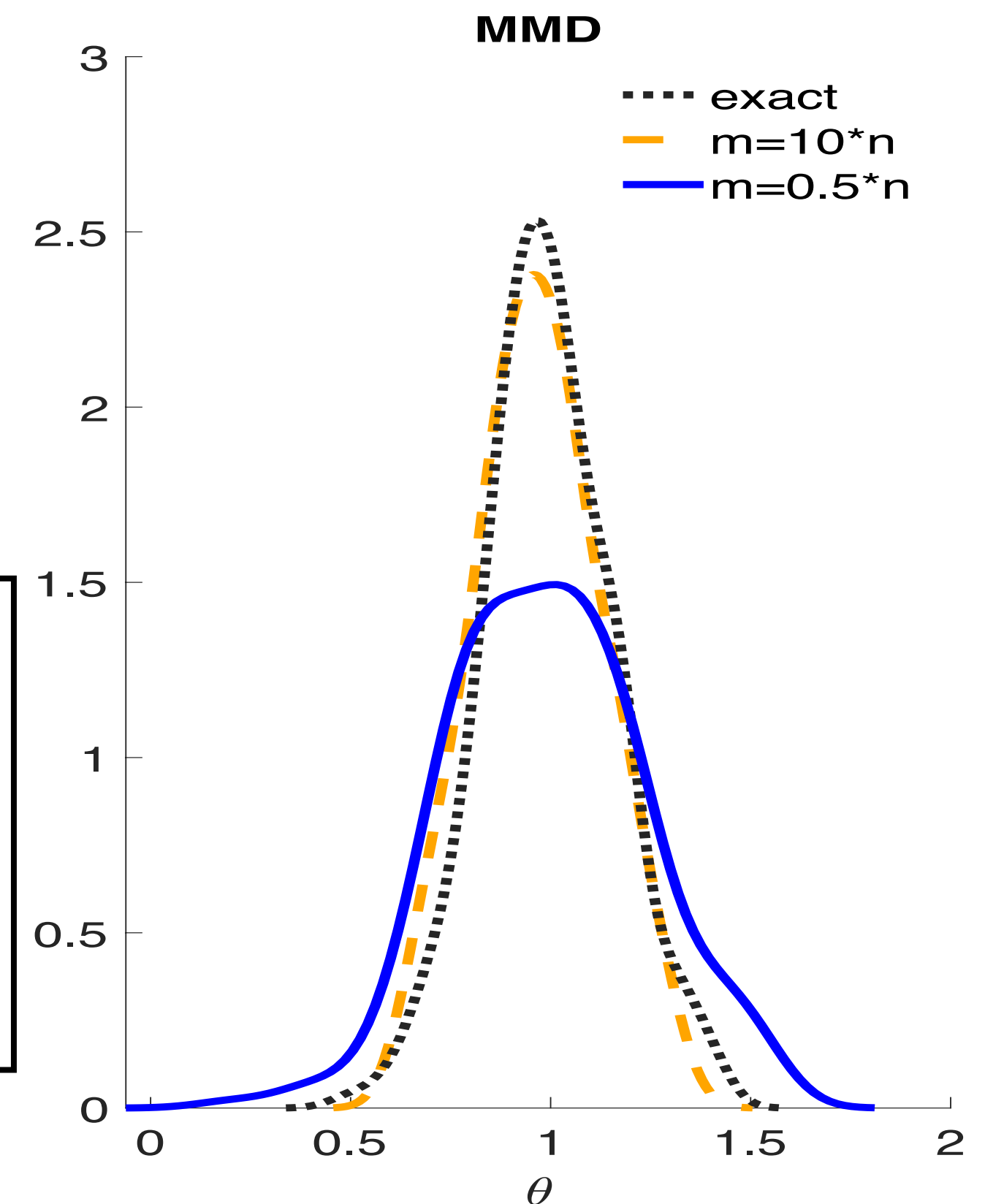
$$L_m(u_{1:m}, x_{1:n}, \theta) \approx L(x_{1:n}, \theta)$$

(e.g. $I(\theta) \approx \frac{1}{m} \sum_{j=1}^m f(u_j)$ and $u_j \sim p(u_j | \theta)$)

Example:

$$L^k(x_{1:n}, \theta) = n \cdot \iint k(u, u') p(u | \theta) p(u' | \theta) dud u' - 2 \sum_{i=1}^n \int k(x_i, u) p(u | \theta) du + C$$

$$L_m^k(u_{1:m}, x_{1:n}, \theta) = \frac{1}{n} \sum_{j=1}^n \sum_{l=1}^n k(u_j, u_l) - 2 \sum_{i=1}^n \frac{1}{m} \sum_{j=1}^m k(x_i, u_j) + C$$



Approximating L

Q3: How does approximating $L(x_{1:n}, \theta)$ affect $\pi_n^L(\theta | x_{1:n})$?

Assumption 1: There are $\kappa_1 > 0$, $\kappa_2 > 0$ and $\{\sigma_n^2(\theta)\}: \Theta \rightarrow \mathbb{R}_+\}_{n=1}^\infty$ so that

$$\begin{aligned} & \underbrace{\sqrt{\mathbb{E}_{u_{1:m} \sim p(u_{1:m}|\theta)} [L_m(u_{1:m}, x_{1:n}, \theta)] - L(x_{1:n}, \theta)}}_{= \text{bias}_m(\theta)} \lesssim \sigma_n^2(\theta) \cdot m^{-\kappa_1} \\ & \underbrace{\sqrt{\mathbb{E}_{u_{1:m} \sim p(u_{1:m}|\theta)} \left[\left\{ L_m(u_{1:m}, x_{1:n}, \theta) - \mathbb{E}_{u'_{1:m} \sim p(u_{1:m}|\theta)} [L_m(u'_{1:m}, x_{1:n}, \theta)] \right\}^2 \right]}}_{= \text{variance}_m(\theta)} \lesssim \sigma_n^2(\theta) \cdot m^{-\kappa_2} \end{aligned}$$

a.s. bounded by RHS for m large enough

Rates of convergence

Assumption 2: $\{\sigma_n^2(\theta)\}: \Theta \rightarrow \mathbb{R}_+\}_{n=1}^\infty$ satisfies some prior + posterior integrability conditions

$$\left(\implies \int \text{bias}_m(\theta) \pi_n^L(\theta | x_{1:n}) d\theta \lesssim m^{-\kappa_1}, \quad \int \text{variance}_m(\theta) \pi_n^L(\theta | x_{1:n}) d\theta \lesssim m^{-\kappa_2} \right)$$

Approximating L

Q3: How does approximating $L(x_{1:n}, \theta)$ affect $\pi_n^L(\theta | x_{1:n})$?

Actual target posterior: $\pi_{n,m}^L(\theta | x_{1:n}) \propto \exp \left\{ -\mathbb{E}_{u_{1:m} \sim p(u_{1:m}|\theta)} [L_m(u_{1:m}, x_{1:n}, \theta)] \right\} \cdot \pi(\theta)$

Theorem 1: Under **Assumptions 1 + 2**, for all $\xi \in [0, 2]$ and any fixed $x_{1:n}$,

a.s. bounded by RHS for
 m large enough

$$\int \|\theta\|^\xi \left| \pi_{n,m}^L(\theta | x_{1:n}) - \pi_n^L(\theta | x_{1:n}) \right| d\theta \lesssim m^{-\min\{\kappa_1, \kappa_2\}}$$

(also implies convergence in TV and TV of all ξ th moments)

Rate of convergence:

slower between bias and variance decay

\implies trading off bias vs variance improves approximation

Approximating L

Q3: How does approximating $L(x_{1:n}, \theta)$ affect $\pi_n^L(\theta | x_{1:n})$?

Assumption 3: Standard (mild) regularity conditions for posterior concentration of $\pi_n^L(\theta | x_{1:n})$

Also, we have $m = m(n)$ and there are $\eta_1 > 0, \eta_2 > 0$ so that

$$m(n)^{-\kappa_1} \underbrace{\int \sigma_n^2(\theta) \pi_n^L(\theta | x_{1:n}) d\theta}_{\text{a.s. of same order for } n \text{ large enough}} \asymp m(n)^{-\eta_1} \quad m(n)^{-\kappa_2} \int \sigma_n^2(\theta) \pi_n^L(\theta | x_{1:n}) d\theta \asymp m(n)^{-\eta_2}$$

a.s. of same order for n large enough

e.g., if $\sigma_n^2(\theta) = \text{constant}$, then $\kappa_1 = \eta_1$

$$\left(\implies \int \text{bias}_m(\theta) \pi_n^L(\theta | x_{1:n}) d\theta \text{ doesn't diverge as } n \rightarrow \infty \right)$$

Approximating L

Q3: How does approximating $L(x_{1:n}, \theta)$ affect $\pi_n^L(\theta | x_{1:n})$?

Theorem 2: Under **Assumptions 1–3**, for $m = m(n)$, and any $M_n \rightarrow \infty$, as $n \rightarrow \infty$

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\int \left| \min_{\theta' \in \Theta} L^\infty(\theta') - L^\infty(\theta) \right| \pi_{m,n}^L(\theta | x_{1:n}) d\theta > M_n / \min \left\{ \sqrt{n}, m(n)^{\min\{\eta_1, \eta_2\}} \right\} \right) = 0$$

$$L^\infty(\theta) = \lim_{n \rightarrow \infty} \mathbb{E}_{x_{1:n}} \left[\frac{1}{n} L(x_{1:n}, \theta) \right]$$

Standard concentration rate of un-approximated posterior $\pi_n^L(\theta | x_{1:n})$

Need to choose $m(n) \asymp (\sqrt{n})^{1/\min\{\eta_1, \eta_2\}}$ to avoid slower concentration due to $L_m(u_{1:m}, x_{1:n}, \theta) \approx L(x_{1:n}, \theta)$

Tells us how good loss approximation needs to be
Helps evaluate which losses are computationally infeasible

Stein Discrepancies

$$D_{\text{Stein}}(p(\cdot | \theta), q_\epsilon) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{X \sim q_\epsilon} [f(X)] - \mathbb{E}_{X \sim p(X|\theta)} [f(X)] \right| = \mathbb{E}_{X \sim q_\epsilon} \left[f_\theta^*(X) \right] \approx \frac{1}{n} \sum_{i=1}^n f_\theta^*(x_i) = \frac{1}{n} \mathbf{L}(x_{1:n}, \theta)$$

\uparrow
 $\mathcal{F} = \left\{ \mathcal{A}_{p(\cdot|\theta)}(f) : f \in \mathcal{F}_0 \right\}$

\uparrow
 Closed form depends on $p(\cdot | \theta)$ only via $\nabla_x \log p(\cdot | \theta)$!
 (For all $\mathcal{A}_{p(\cdot|\theta)}$ and \mathcal{F}_0 below)

Stein Operator $\mathcal{A}_{p(\cdot|\theta)}$

Langevin-Stein operator

$$\mathcal{A}_{p(\cdot|\theta)}(f)(x) = f(x) \cdot \nabla_x \log p(x | \theta) + \nabla \cdot f(x)$$

Diffusion Stein operator

$$\mathcal{A}_{p(\cdot|\theta)}(f)(x) = f(x) \cdot m(x)^T \nabla_x \log p(x | \theta) + \nabla \cdot f(x)$$

can be tuned for robustness

Stein Set \mathcal{F}_0

$$\mathcal{F}_0 = \left\{ f \in C^1(\mathcal{X}, \mathbb{R}) \cap L^2(\mathcal{X}; p(\cdot | \theta)) : \|f\|_{L^2(\mathcal{X}; p(\cdot|\theta))} \leq 1 \right\}$$

$$\mathcal{F}_0 = \left\{ f \in C^1(\mathcal{X}, \mathbb{R}) \cap L^2(\mathcal{X}; p(\cdot | \theta)) : \|f\|_{L^2(\mathcal{X}; p(\cdot|\theta))} \leq 1 \right\}$$

$$\mathcal{F}_0 = \left\{ f \in \text{RKHS}(K) : \|f\|_K \leq 1 \right\}$$

$D_{\text{Stein}}(p(\cdot | \theta), q_\epsilon)$

Fisher Divergence (FD)/
Score Matching

Weighted Fisher Divergence/
Diffusion Score Matching (DSM)
(robust for right choice of m)

Kernel Stein Discrepancy (KSD)
(robust for right choice of m)

Computationally Efficient Beliefs: Methodology

$$\pi_n^{\mathcal{L}}(\theta | x_{1:n}) = \frac{\exp\{-\mathcal{L}(x_{1:n}, \theta)\} \cdot \pi(\theta)}{\int \exp\{-\mathcal{L}(x_{1:n}, \theta)\} \cdot \pi(\theta) d\theta} \propto \exp\left\{-\sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))\right\} \underbrace{\exp\left\{(\theta - \mu_0)^\top \Lambda_0 (\theta - \mu_0)\right\}}_{\text{squared exponential prior}}$$

exponential family

$$\stackrel{!}{=} \mathcal{N}(\theta; \mu_{\mathcal{L}}(x_{1:n}), \Sigma_{\mathcal{L}}(x_{1:n}))$$

$$p(x | \theta) = h(x) \cdot \exp\{T(x)^\top \theta - A(\theta)\}$$

$f^*(x)$ is a 2nd order polynomial

$$\mathcal{L}(p_\theta(x_{1:n})) = \sum_{i=1}^n f_\theta^*(x_i) = \sum_{i=1}^n f^*(\nabla_x \log p(x_i | \theta)) \stackrel{!}{=} \sum_{i=1}^n a_i + b_i \cdot (\nabla_x \log p(x_i | \theta)) + c_i \cdot (\nabla_x \log p(x_i | \theta))^2$$

$$= \sum_{i=1}^n \left[a_i \tilde{h}_i + b_i \tilde{h}_i^2 \right] + \left[b_i \tilde{T}_i + 2c_i \tilde{h}_i \tilde{T}_i \right]^\top \theta + \theta^\top \left[c_i \tilde{T}_i \tilde{T}_i^\top \right] \theta$$

$$= \nabla_x \log h(x_i) + \nabla_x T(x_i)^\top \theta$$

$$= \tilde{h}_i$$

$$= \tilde{T}_i$$

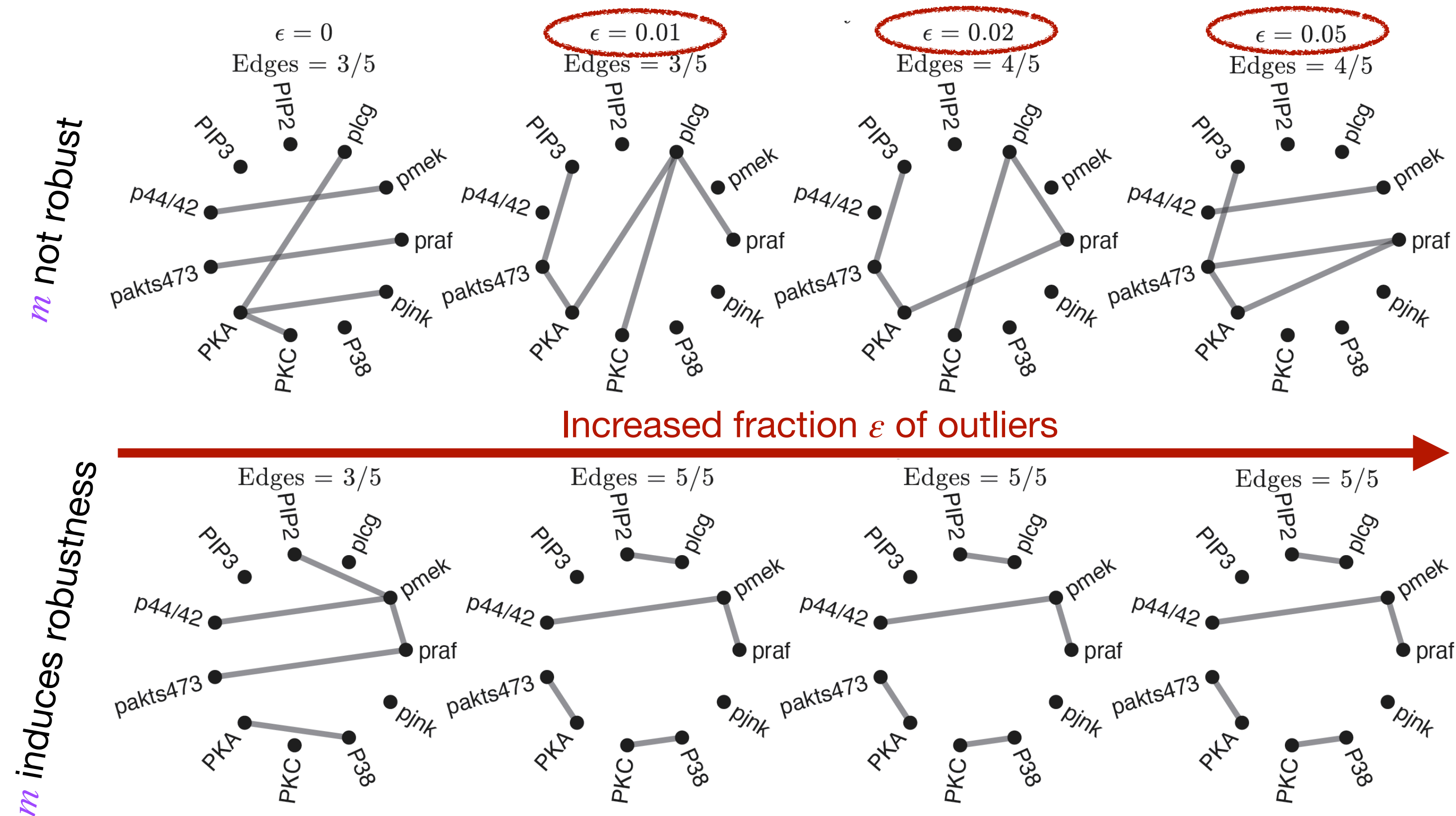
$$\stackrel{+C}{=} \sum_{i=1}^n (\theta - \mu(x_i))^\top \Lambda(x_i) (\theta - \mu(x_i))$$

$\exp\{-\dots\}$ of this looks like a Gaussian in θ
 \Rightarrow Put squared exponential prior on θ !

L for Robustness & Computation: Graphical Models

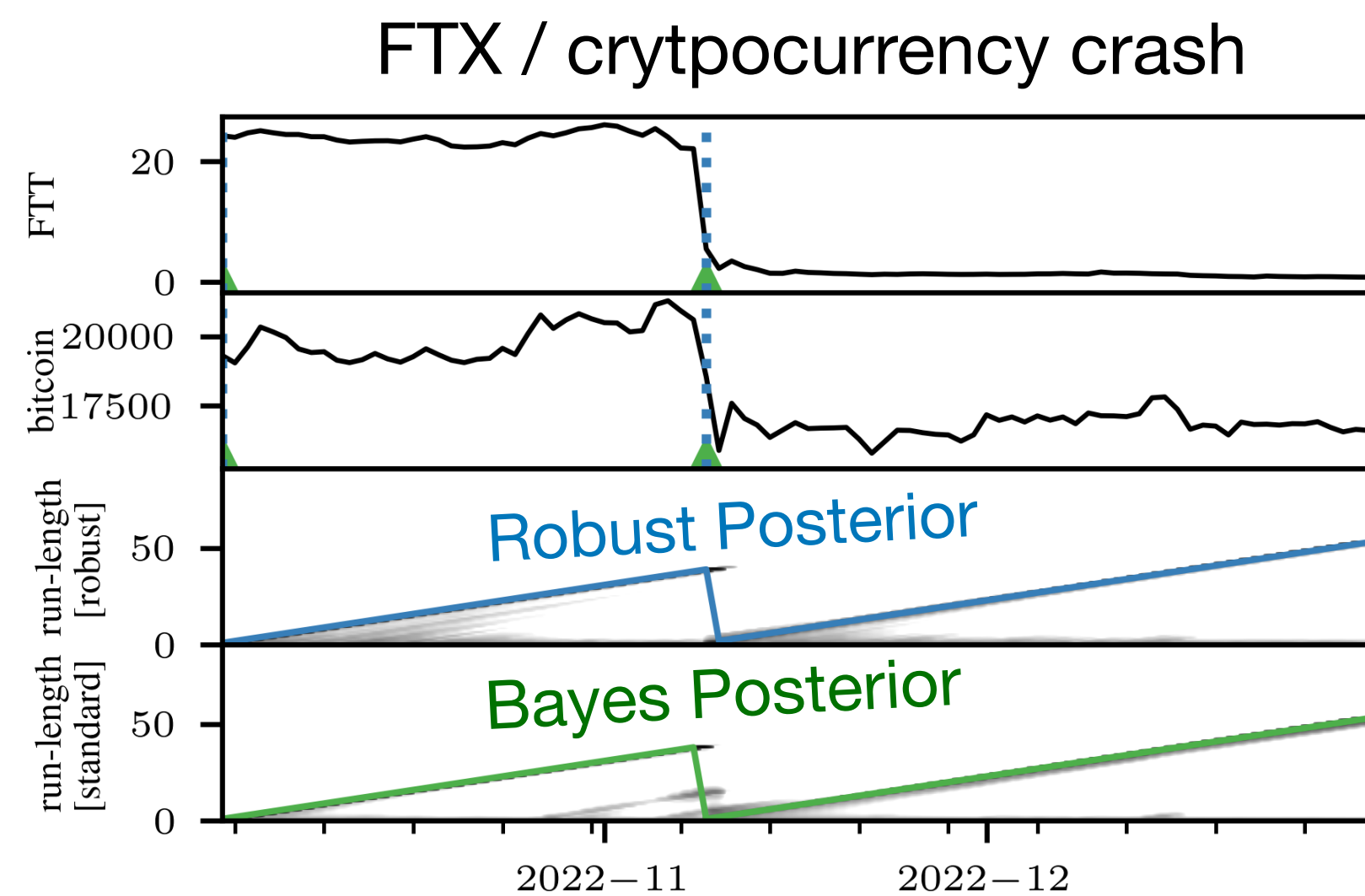
$$p(x_i | \theta) \propto \exp \left\{ - \sum_{d=1}^D \theta_d x_{i,(d)} - \sum_{d < d'}^D \theta_{d,d'} x_{i,(d)} x_{i,(d')} \right\}; \quad \theta_d > 0, \theta_{d,d'} \geq 0. D = 11 \text{ (proteins)}, \quad n = 7466 \text{ (number of cells measured)}$$

L = **Kernel Stein Discrepancy** (robust / non-robust loss depending only on $\nabla_x \log p(x_{1:n} | \theta)$)

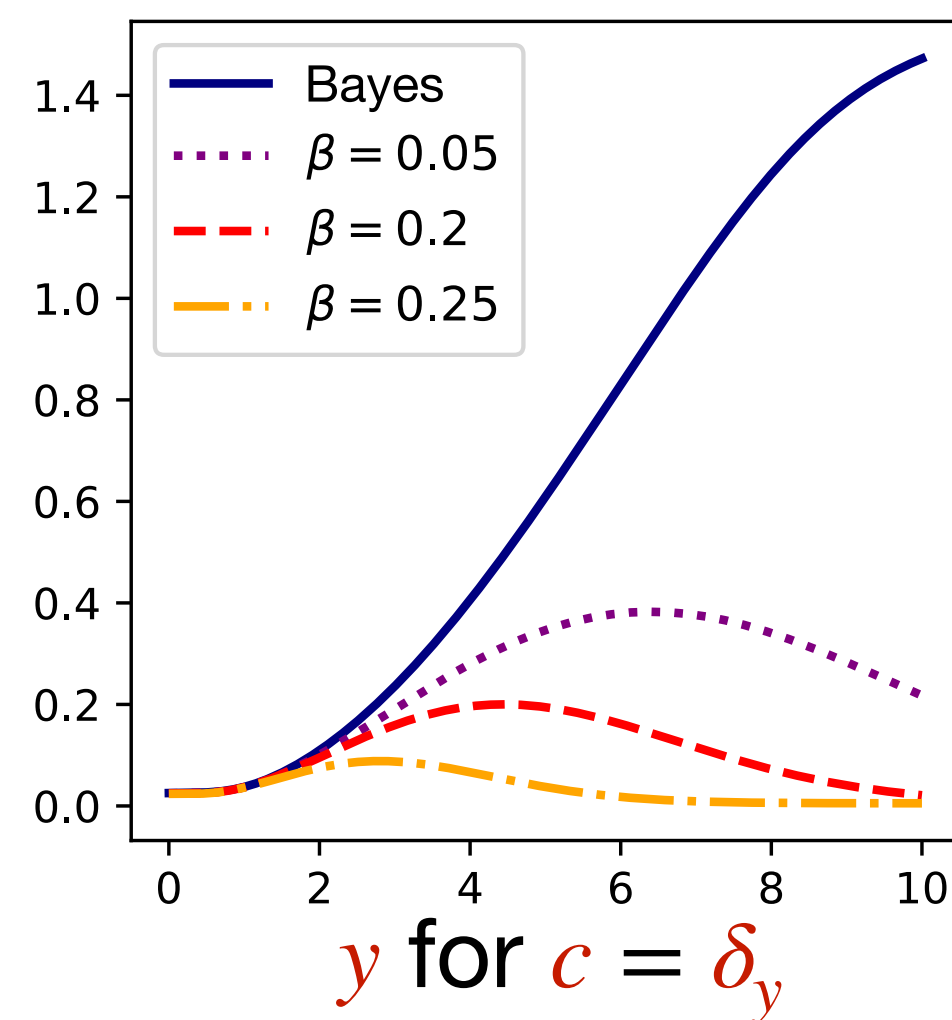


Choices of L

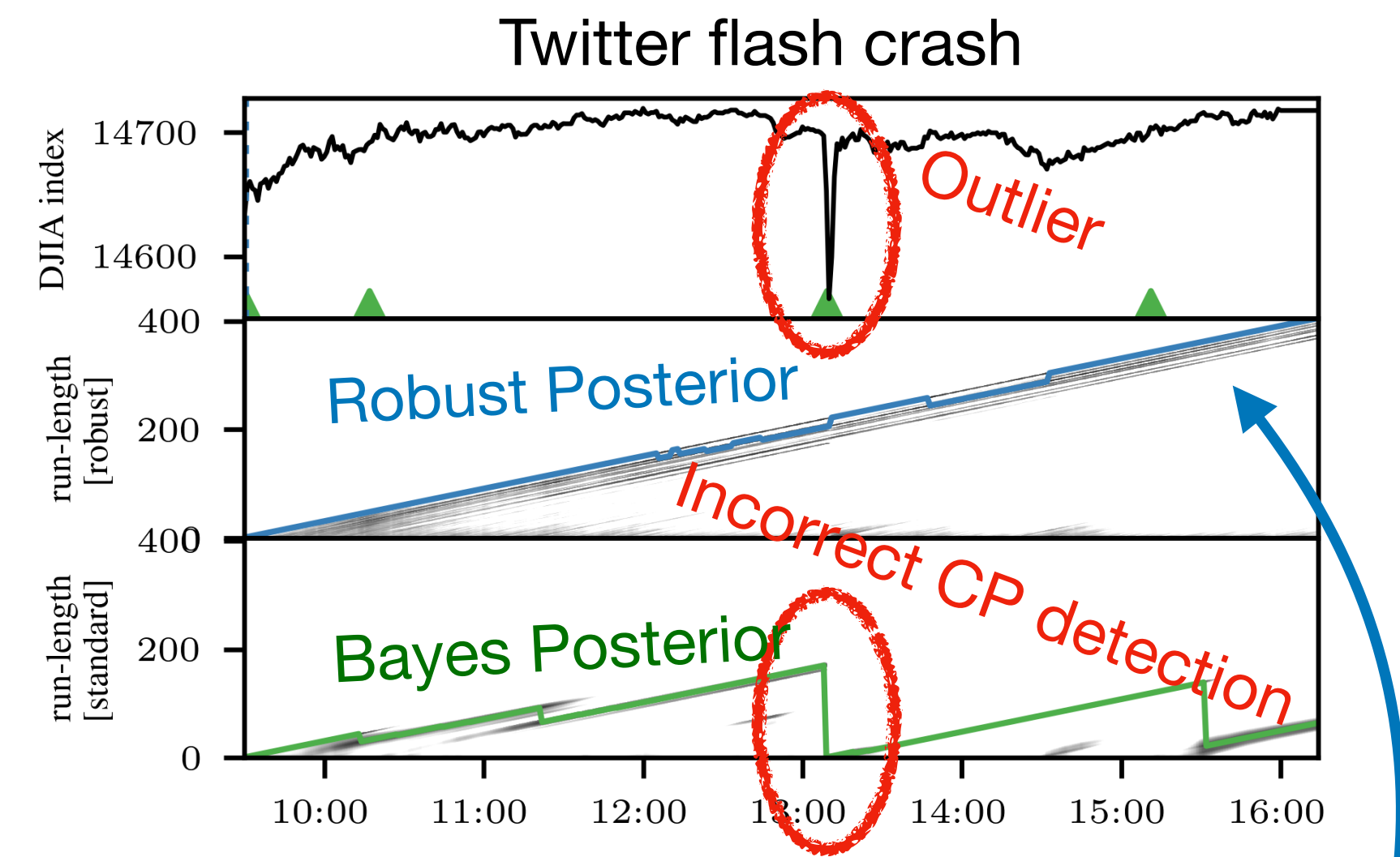
Example 1: Bayesian On-line Changepoint Detection (β -Divergence)



$$\lim_{\varepsilon \rightarrow 0} D_{FR}(\pi_n^L(\theta | q_\varepsilon), \pi_n^L(\theta | q_0))$$



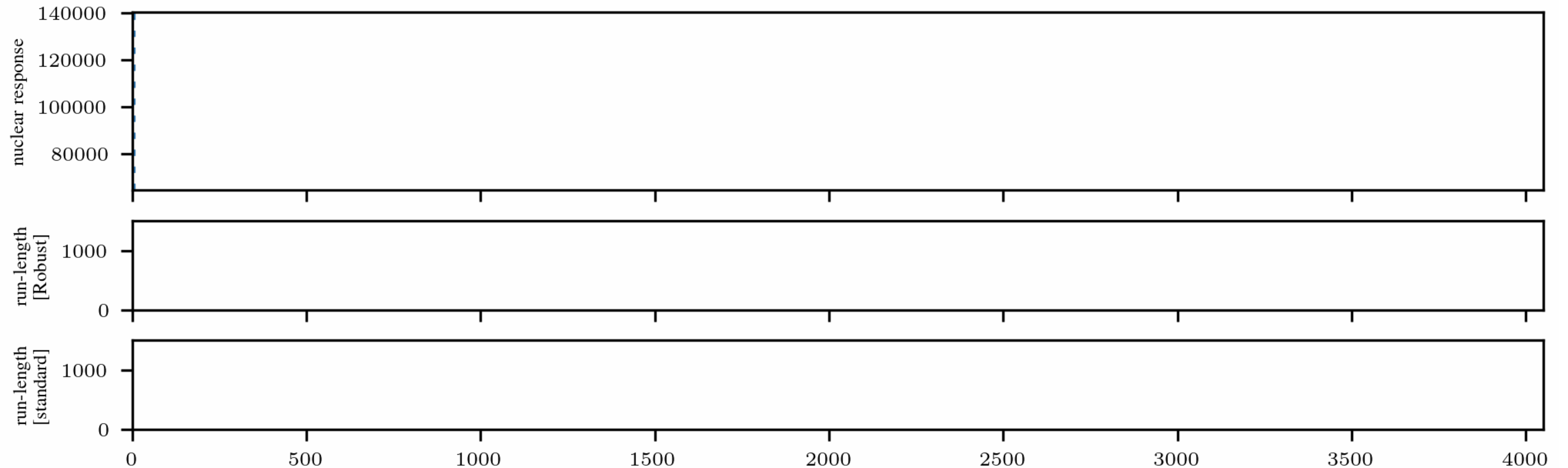
$$q_\varepsilon = (1 - \varepsilon) \cdot q_0 + \varepsilon \cdot c$$



Theorem: the robust algorithm cannot declare a change point after a single outlier.

L for Robustness & Computation: Changepoints

L = **Weighted Fisher Divergence** (robust / non-robust loss depending only on $\nabla_x \log p(x_{1:n} | \theta)$)



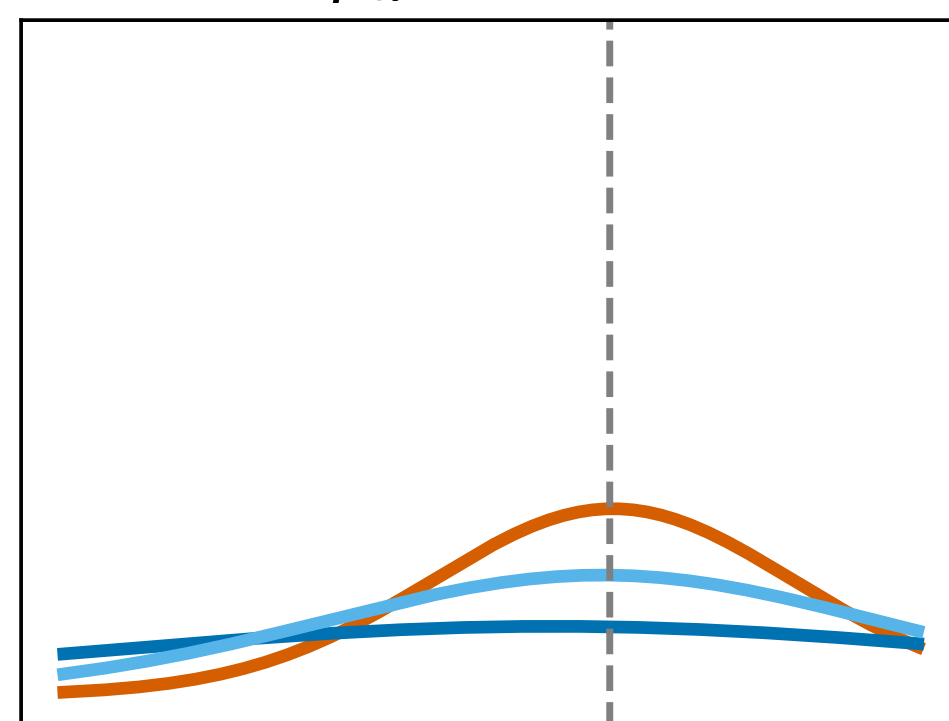
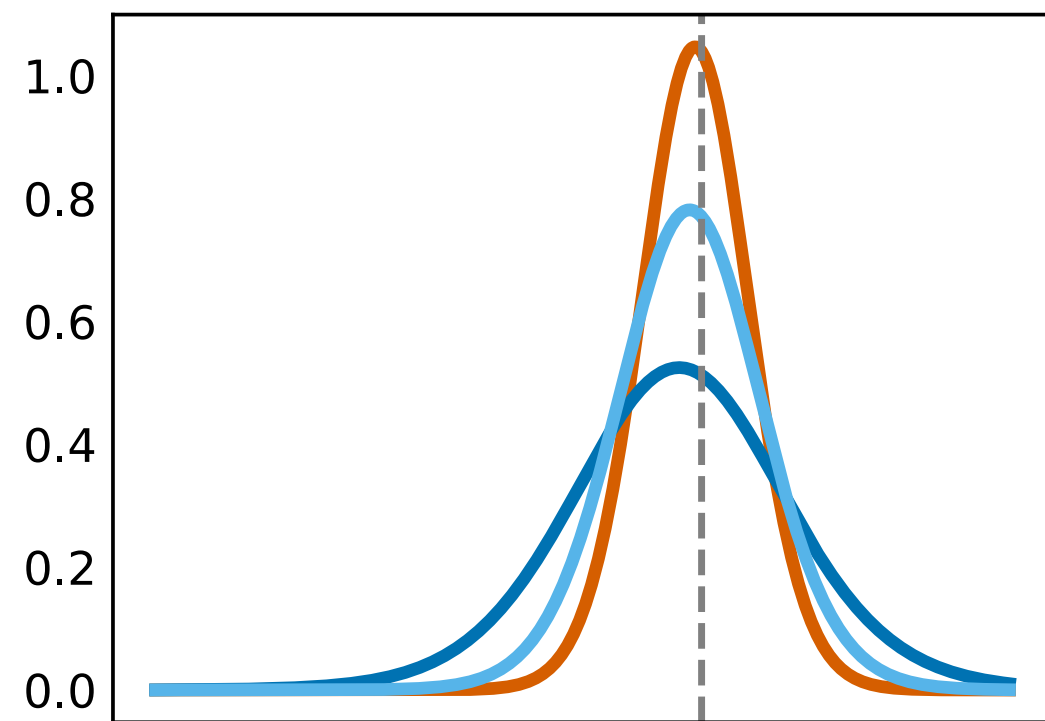
Post-Bayesian ML: Optimisation-centric posteriors

Assumptions

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

(Prior well-specified)

(Prior misspecified)

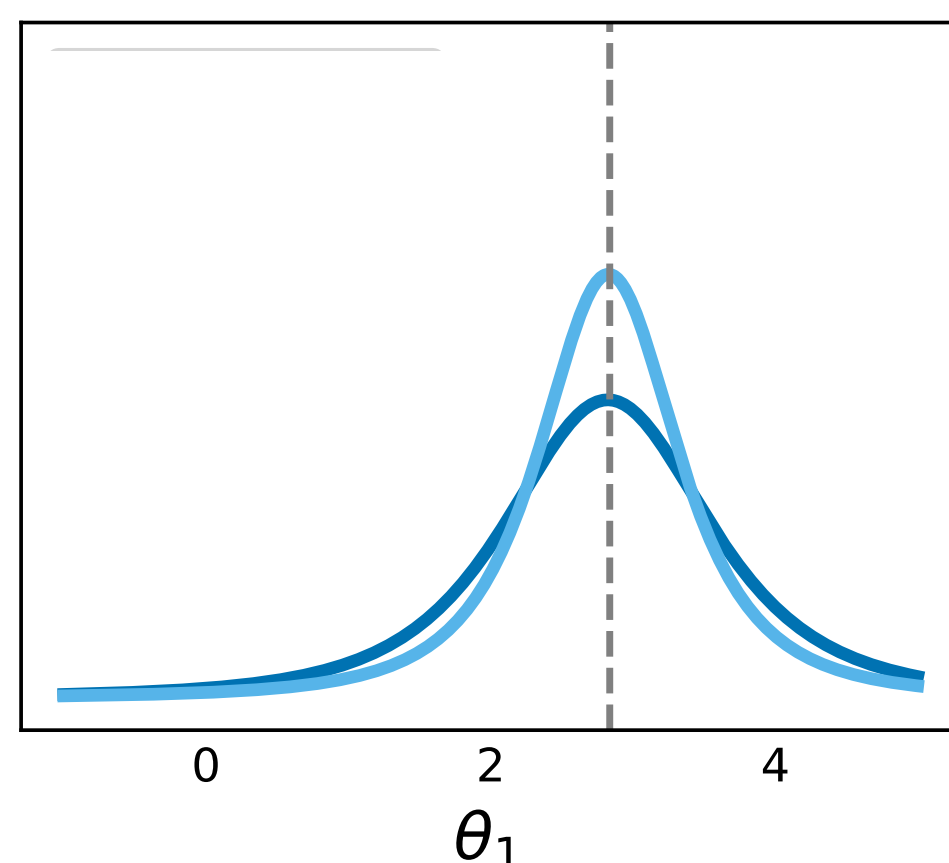
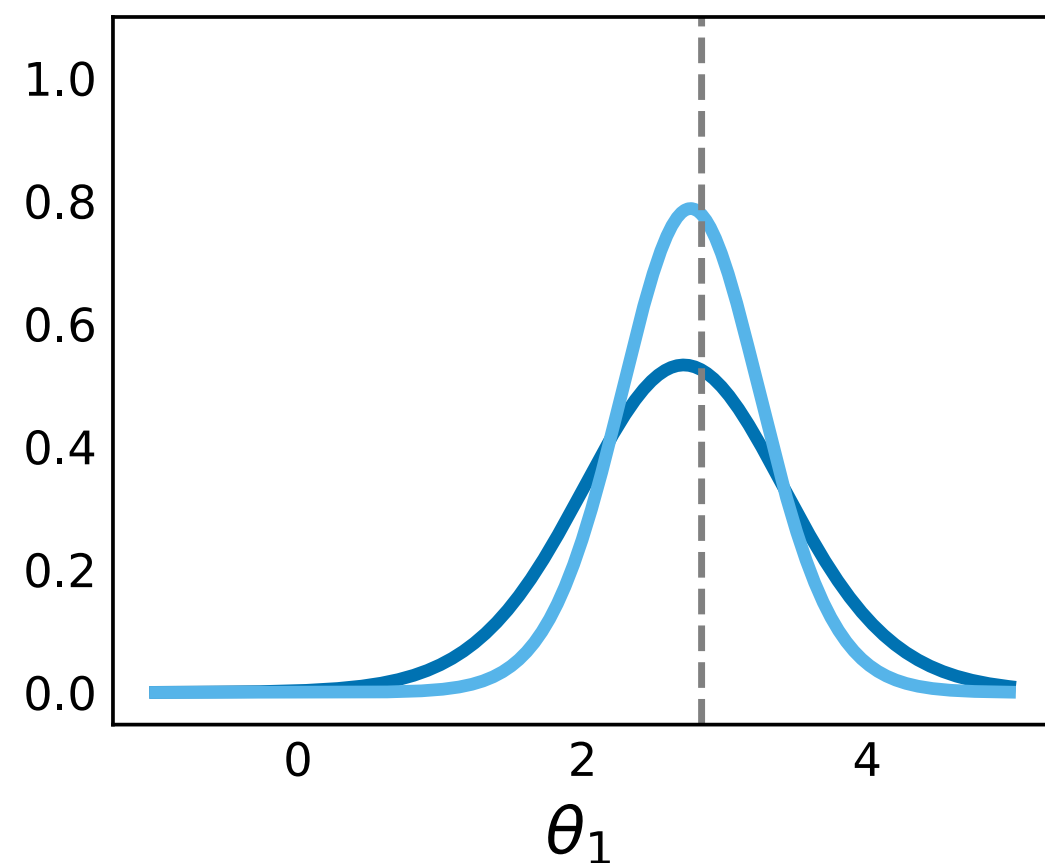


Mean Field VI for:

Standard Bayes

Power posterior; $\lambda = 0.6$

Power posterior; $\lambda = 0.3$



GVI/Optimisation-centric:

$D = \text{Renyi divergence}; \alpha = 0.6$

$D = \text{Renyi divergence}; \alpha = 0.3$

$$q_n^*(\theta) = \operatorname{argmin}_{q \in \text{Normals}} \left\{ \mathcal{L}_{L, D, \pi}^{(\lambda)}(q) \right\}$$

$$- \int \log p(x_{1:n} | \theta) q(\theta) d\theta + D^{\parallel}(q, \pi)$$

Case Study: Why Post-Bayesian ML matters

Objective: $q \mapsto \mathbb{E}_{\theta \sim q} [-\log p(x_{1:n} | \theta)] + \frac{1}{\lambda} \cdot \text{KL}(q, \pi) \xrightarrow{\lambda \rightarrow \infty} q \mapsto \mathbb{E}_{\theta \sim q} [-\log p(x_{1:n} | \theta)]$

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

Deep Ensemble ($\lambda \rightarrow \infty$)

Wasserstein Gradient Flow:

Step 1: Sample $\theta_k(0) \sim \pi, k = 1, 2, \dots, K$ Brownian motions

Step 2: For $t \in [0, T]$, evolve as $\{B_n(t)\}_{t \geq 0}$

$$d\theta_k(t) = - \left(\lambda \cdot \nabla_{\theta} \mathbf{L}(x_{1:n}, \theta_k(t)) - \nabla \log \pi(\theta_k(t)) \right) dt + \sqrt{2dB_k(t)}$$

$$q_n^*(\theta) = \pi_n^{(\lambda)}(\theta | x_{1:n}) \approx \frac{1}{K} \sum_{n=1}^K \theta_k(T) \text{ as } T \rightarrow \infty, K \rightarrow \infty.$$

Converges to well-defined density

Wasserstein Gradient Flow =
Deep Ensemble Algorithm

Step 1: Sample $\theta_k(0) \sim \pi, k = 1, 2, \dots, K$

Step 2: For $t \in [0, T]$, evolve as

$$d\theta_k(t) = -\lambda \cdot \nabla_{\theta} \mathbf{L}(x_{1:n}, \theta_k(t))$$

$$\text{Deep Ensemble} = \frac{1}{K} \sum_{n=1}^K \theta_k(T)$$

Converges to ill-defined atomic measure

The Future of Post-Bayesian ML: Success Stories

- $q \mapsto \mathbb{E}_{\theta \sim q} [\mathbf{L}(x_{1:n}, \theta)] \longrightarrow$ **Deep Ensemble**
- $q \mapsto \mathbb{E}_{\theta \sim q} [\mathbf{L}(x_{1:n}, \theta)] + w_2 \cdot \mathbf{KL}(q, \pi) \longrightarrow$ **Cold Posterior** ($w_2 \ll 1$) / **Bayes Posterior** ($w_2 = 1$)
- $q \mapsto \mathbb{E}_{\theta \sim q} [\mathbf{L}(x_{1:n}, \theta)] + w_1 \cdot \mathbf{MMD}^2(q, \pi) + w_2 \cdot \mathbf{KL}(q, \pi) \longrightarrow$ **BDL Ensemble with repulsive particles**

Infinite-dimensional gradient descent / Wasserstein Gradient Flow:

Step 1: Sample $\theta_k(0) \sim \pi, k = 1, 2, \dots, K$

Step 2: Evolve via SDE for $t \in [0, T]$ as

$$d\theta_k(t) = - \left(\nabla_{\theta} \mathbf{L}(x_{1:n}, \theta_k(t)) - w_1 \nabla \mu_{\pi}(\theta_k(t)) - w_2 \nabla \log \pi(\theta_k(t)) + \frac{w_1}{K} \sum_{j=1}^K w_1 \kappa(\theta_k(t), \theta_j(t)) \right) dt + \sqrt{2w_2} dB_k(t)$$

independent Brownian motions $\{B_n(t)\}_{t \geq 0}$

$$q_n^*(\theta) \approx \frac{1}{K} \sum_{n=1}^K \theta_k(T) \text{ as } T \rightarrow \infty, K \rightarrow \infty.$$

Adversarial Robustness: How to think about general **D**?

Theorem: $L(x_{1:n}, \theta)$ bounded + measurable, $\lambda > 0$, \mathcal{Q} is closed, bounded, convex:

$$\inf_{q \in \mathcal{Q}} \left\{ \mathbb{E}_{\theta \sim q} [\lambda \cdot L(x_{1:n}, \theta)] + D(q, \pi) \right\} = \sup_{h \in F_b(\Theta)} \left\{ \underbrace{\inf_{q \in \mathcal{Q}} \mathbb{E}_{\theta \sim q} [\lambda \cdot L(x_{1:n}, \theta) + \lambda \cdot h(\theta)]}_{\text{Perturbed / worsened loss}} - \underbrace{\Delta_{D, \pi}(\lambda \cdot h)}_{\text{Penalty for perturbation}} \right\}$$

inner minimisation over (perturbed) loss

outer maximisation (of adversary) over perturbations

$h \in F_b(\Theta) = \{ \text{bounded + measurable functions } f: \Theta \rightarrow \mathbb{R} \}$

Legendre-Fenchel dual of $D(\cdot, \pi) = \sup_{\rho \in \mathcal{P}(\Theta)} \left\{ \int_{\Theta} h(\theta) \rho(\theta) d\theta - D(\rho, \pi) \right\} = \Delta_{D, \pi}(h)$
 (penalty for increasing $h(\theta)$ highest where $\pi(\theta)$ largest)

Adversarial Robustness: How to think about general **D**?

Theorem: $L(x_{1:n}, \theta)$ bounded + measurable, $\lambda > 0$, \mathcal{Q} is closed, bounded, convex:

$$\inf_{q \in \mathcal{Q}} \left\{ \mathbb{E}_{\theta \sim q} [\lambda \cdot L(x_{1:n}, \theta)] + D(q, \pi) \right\} = \sup_{h \in F_b(\Theta)} \left\{ \underbrace{\inf_{q \in \mathcal{Q}} \mathbb{E}_{\theta \sim q} [\lambda \cdot L(x_{1:n}, \theta) + \lambda \cdot h(\theta)]}_{\text{Perturbed / worsened loss}} - \underbrace{\Delta_{D, \pi}(\lambda \cdot h)}_{\text{Penalty for perturbation}} \right\}$$

Example 1: $D = \text{KL}$,

$$\Delta_{D, \pi}(\lambda \cdot h) = \log \int \exp \{ \lambda \cdot h(\theta) \} \pi(\theta) d\theta$$

Example 2: $D = \chi^2$,

$$\Delta_{D, \pi}(\lambda \cdot h) = \lambda \cdot \int h(\theta) \pi(\theta) d\theta + \frac{\lambda^2}{4} \text{Var}_{\theta \sim \pi} (h(\theta))$$

Example 3: $D = \text{IPM}_{\mathcal{W}}$,

$$\Delta_{D, \pi}(\lambda \cdot h) = \begin{cases} \lambda \cdot \int h(\theta) \pi(\theta) d\theta & \text{if } h \in \mathcal{W} \\ \infty & \text{if } h \notin \mathcal{W} \end{cases}$$

π = **preferences**. More density/mass where we want perturbations to be more expensive for adversary.

D = **cost/penalty function**. Determines actual cost of perturbations to adversary

How to approach computation with general D ?

More Elegant Strategy: analytical solutions

↳ Advantage: further insight on effect of D

↳ Disadvantage:

How could we fix this?

◦ only applicable for a small selection of D

◦ computationally intractable

$$q_n^*(\theta) = \nabla f^* \left(Z - \lambda \cdot L(x_{1:n}, \theta) \right) \pi(\theta) = \arg \min_{q \in \mathcal{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q} \left[\lambda \cdot L(x_{1:n}, \theta) \right] + D_f(q, \pi) \right\}$$

$$D_f(q, \pi) = \mathbb{E}_{\theta \sim \pi} \left[f \left(\frac{q(\theta)}{\pi(\theta)} \right) \right]$$

$$f(x) = x^2 - 1$$

Z = normaliser defined via $\int \nabla f^* \left(Z - \lambda \cdot L(x_{1:n}, \theta) \right) \pi(\theta) d\theta = 1$

f^* is the Fenchel conjugate of f , $f^*(x) = \sup_{x' \in \mathbb{R}} \{ \langle x, x' \rangle - f(x') \}$

Example: $D_f = \chi^2$, $L \geq 0$

$$q_n^*(\theta) = \frac{1}{2} \max \{ 0, Z - \lambda \cdot L(x_{1:n}, \theta) \} \pi(\theta)$$

Post-Bayesian ML

Gibbs/Generalised/Quasi/Pseudo Posterior

Optimisation-centric posteriors / GVI

$$q_n^*(\theta) = \operatorname{argmin}_{q \in \mathcal{Q}} \left\{ \mathcal{L}_{L, D, \pi}^{(\lambda)}(q) \right\}; \quad \mathcal{Q} \subseteq \mathcal{P}(\Theta)$$

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

Martingale Posterior

For $i = 1, 2, \dots$

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

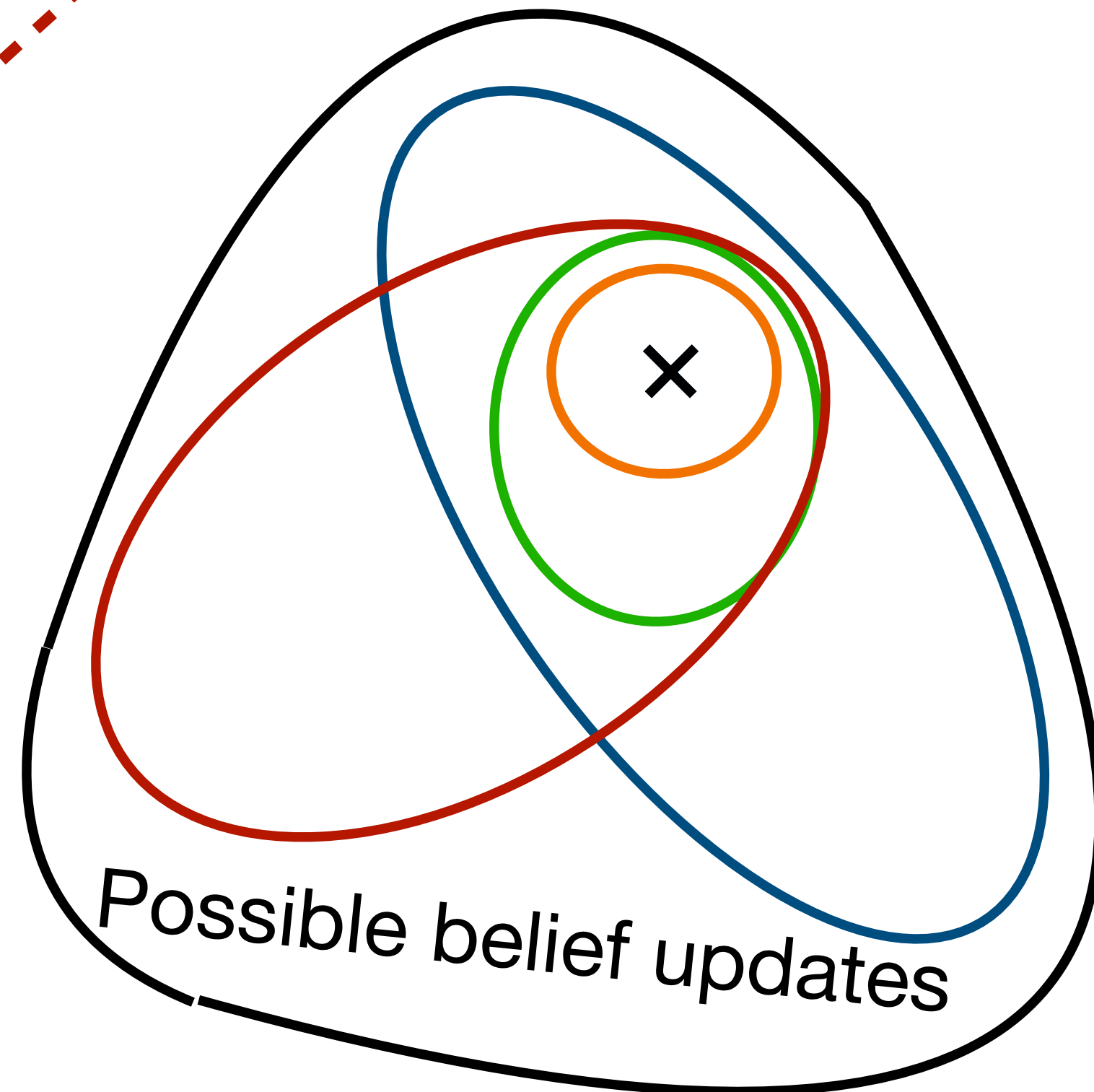
$$\theta^\infty = \operatorname{argmin}_{\theta \in \Theta} L([x_{1:n}, X_{n+1:\infty}], \theta)$$

Power/Fractional/Cold Posterior

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

Bayes' Posterior

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



Post-Bayesian ML

For $i = 1, 2, \dots$

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

$$\theta^\infty = \operatorname{argmin}_{\theta \in \Theta} -\log p(\{x_{1:n} \cup X_{n+1:\infty}\} \mid \theta)$$

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

Bayes' posterior predictive:

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

Bayes' Posterior

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

(Doob's Consistency
Theorem)

=

Post-Bayesian ML

For $i = 1, 2, \dots$

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

$$\theta^\infty = \operatorname{argmin}_{\theta \in \Theta} L(\{x_{1:n} \cup X_{n+1:\infty}\}, \theta)$$

Generalised Bayes' posterior predictive:

$$p(x_{n+1} \mid x_{1:n}) = \int p(x_{n+1} \mid \theta) \pi_n^{(\lambda, L)}(\theta \mid x_{1:n}) d\theta$$

(If parameter θ indexes
a model $p(\cdot \mid \theta)$)

Gibbs/Generalised/Quasi/Pseudo Posterior

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

(Doob's
Consistency)

=

$$\theta^\infty \sim \pi_n^{(\lambda, L)}(\theta \mid x_{1:n})$$

Post-Bayesian ML

For $i = 1, 2, \dots$

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

$$\theta^\infty = \operatorname{argmin}_{\theta \in \Theta} L(\{x_{1:n} \cup X_{n+1:\infty}\}, \theta)$$

Martingale Posterior

$$\theta^\infty \sim \pi_n^{(L,p)}(\theta \mid x_{1:n})$$

Choose predictive & parameter/loss:

$$p(x_{n+1} \mid x_{1:n}) = \text{Your choice}$$

$$L(x_{1:n}, \theta) = \text{Loss determining } \theta$$

Martingale condition:

$$\mathbb{E}_{X_{n+1} \sim p(X_{n+1} \mid x_{1:n})} [p(z \mid \{x_{1:n} \cup X_{n+1}\}) \mid x_{1:n}] = p(z \mid x_{1:n})$$

Post-Bayesian ML

For $i = 1, 2, \dots$

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

$$\theta^\infty = \operatorname{argmin}_{\theta \in \Theta} L(\{x_{1:n} \cup X_{n+1:\infty}\}, \theta)$$

Predictive replaces model & prior specification

Good specification

=

predictive capturing uncertainty about observables

Depends on choices for predictive & loss

Martingale Posterior

$$\theta^\infty \sim \pi_n^{(L,p)}(\theta \mid x_{1:n})$$

Assumptions of Bayesian inference

- (A1) model well-specified N/A?
- (A2) prior well-specified N/A?
- (A3) inversion computationally feasible ✓

Post-Bayesian ML

Optimisation-centric posteriors / GVI

$$q_n^*(\theta) = \operatorname{argmin}_{q \in \mathcal{Q}} \left\{ \mathcal{L}_{L, D, \pi}^{(\lambda)}(q) \right\}; \quad \mathcal{Q} \subseteq \mathcal{P}(\Theta)$$

Gibbs/Generalised/Quasi/Pseudo Posterior

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

Martingale Posterior

For $i = 1, 2, \dots$

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

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

Power/Fractional/Cold Posterior

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

Bayes' Posterior

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

