

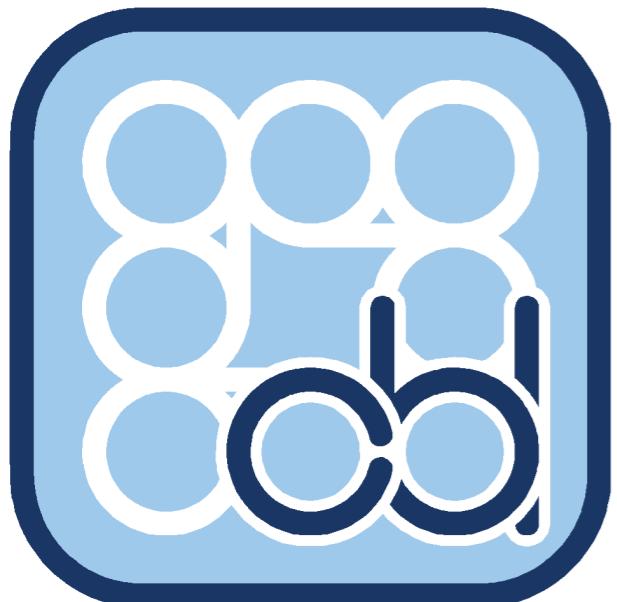
---

# Machine Learning with Objective Priors

---

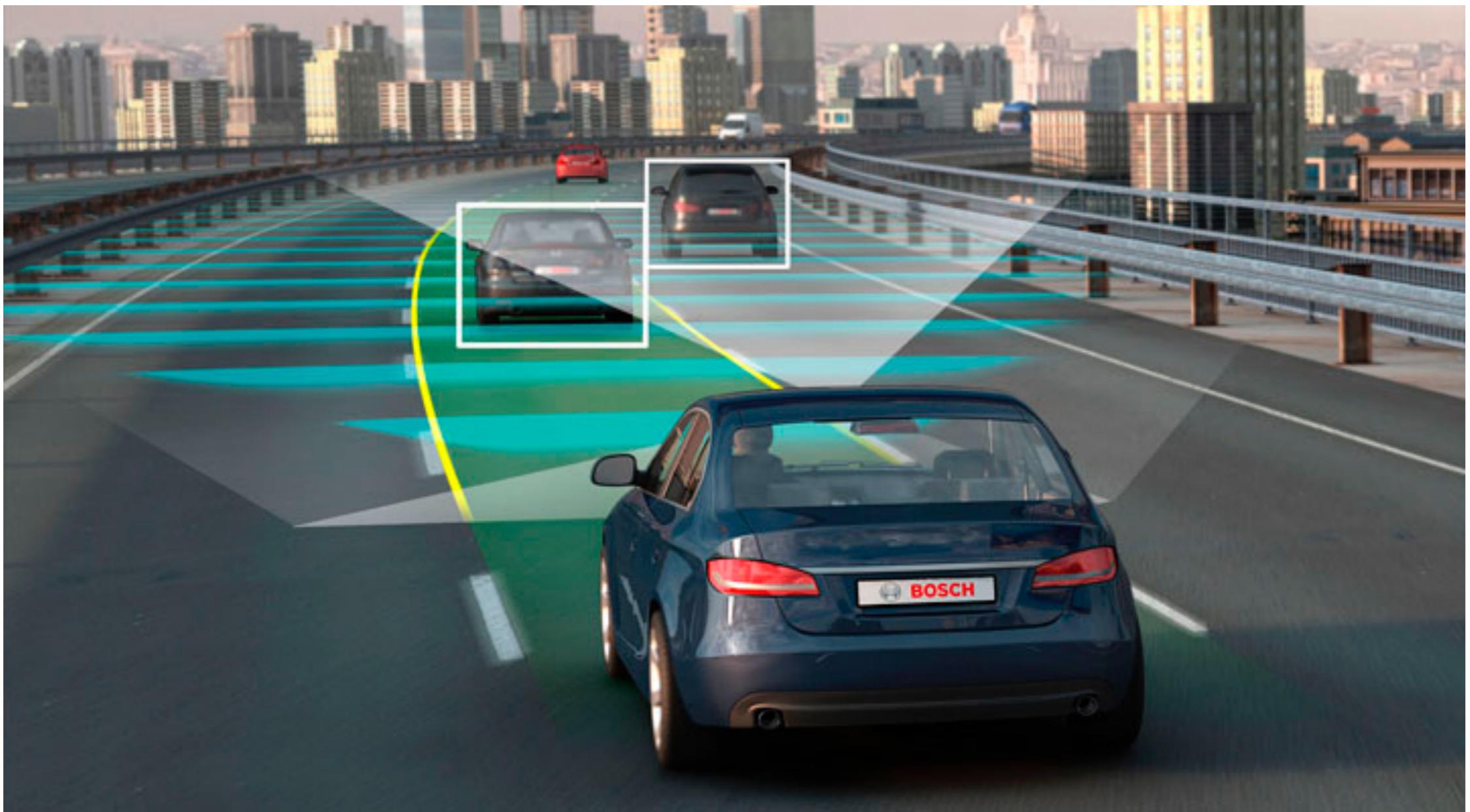
Eric Nalisnick

Div F Conference  
20.3.2019



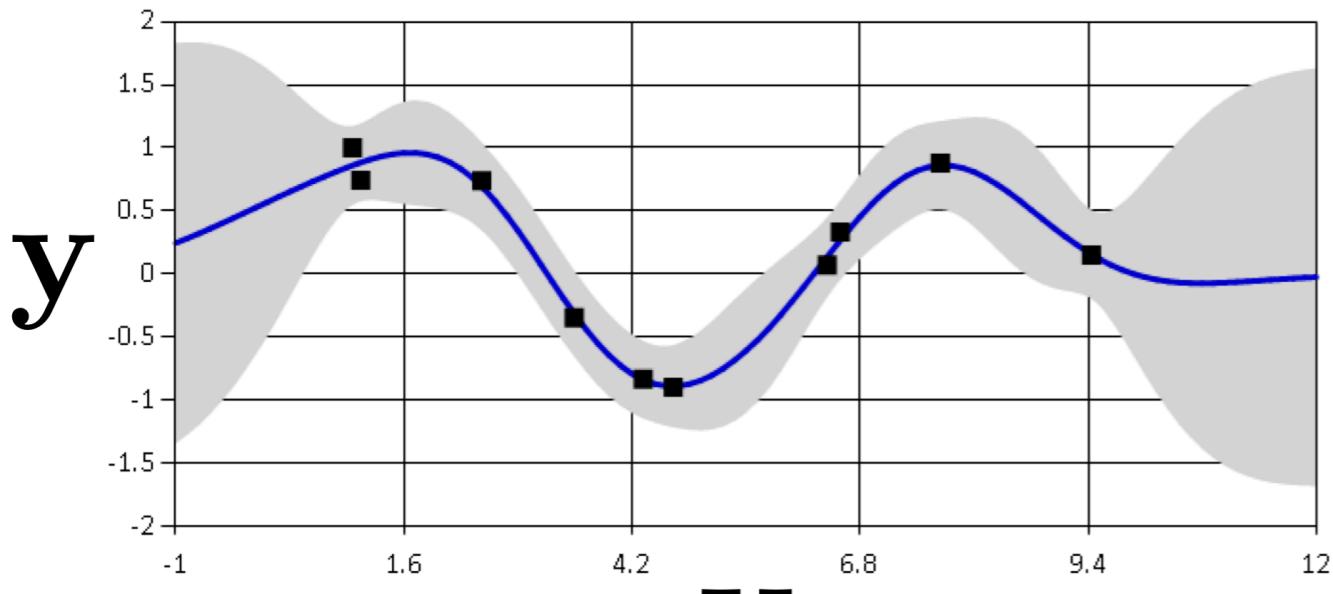
Computational and  
Biological Learning  
University of Cambridge

# Modern machine learning requires uncertainty estimation.

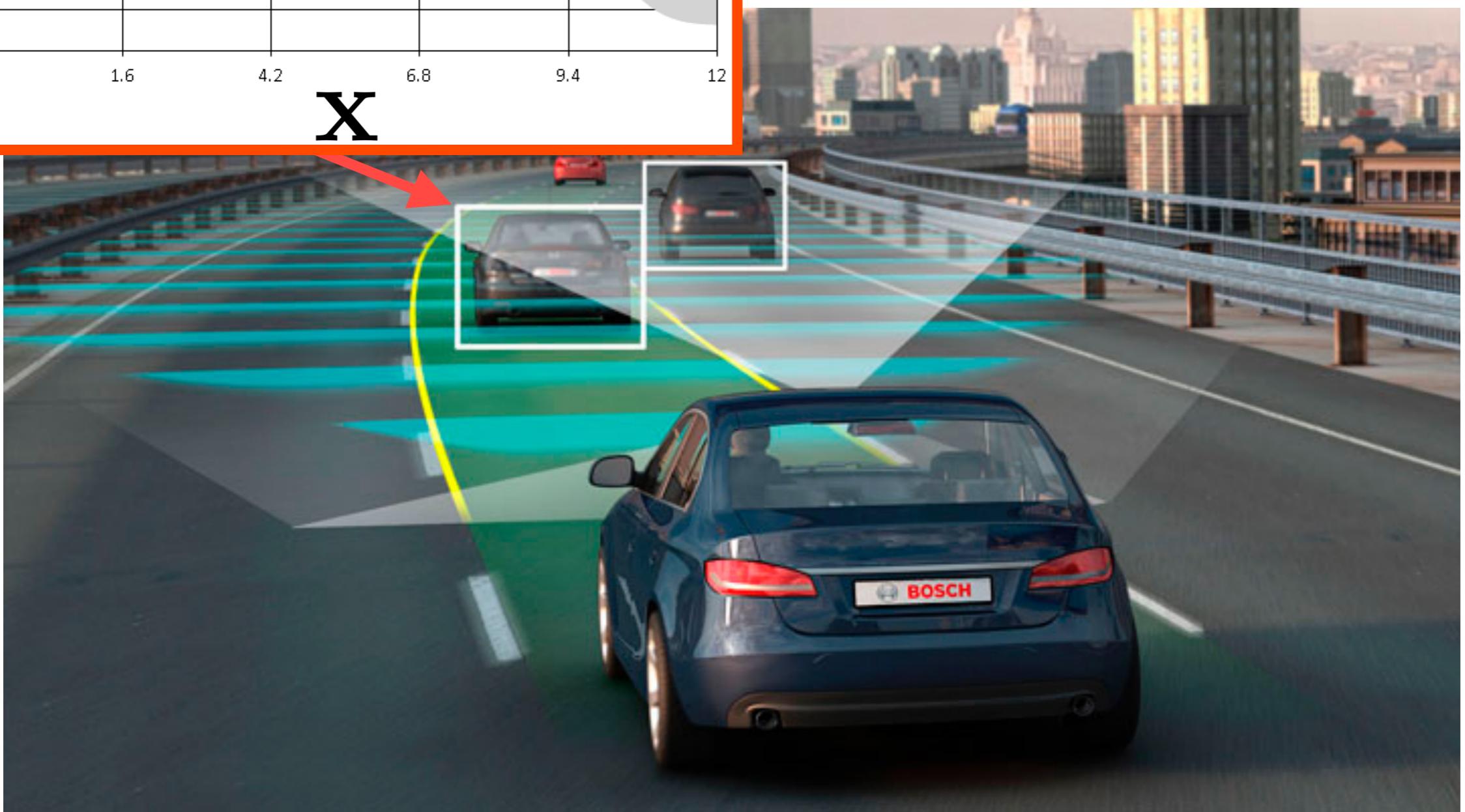


<https://www.rac.co.uk/drive/features/will-self-driving-cars-mean-we-wont-need-car-insurance-anymore/>

■ = data, — = prediction, □ = uncertainty



learning requires estimation.



<https://www.rac.co.uk/drive/features/will-self-driving-cars-mean-we-wont-need-car-insurance-anymore/>

$$p(\theta | \mathbf{X}) = \frac{p(\mathbf{X} | \theta) \ p(\theta)}{p(\mathbf{X})}$$



THOMAS BAYES

Posterior

Likelihood

Prior

$$p(\theta|X) = \frac{p(X|\theta) p(\theta)}{p(X)}$$

Marginal Likelihood

(a.k.a. model evidence)  
(a.k.a. normalizing constant)  
(a.k.a. partition function)

$$p(X) = \int_{\theta} p(X|\theta) p(\theta) d\theta$$



THOMAS BAYES

$$p(\theta | \mathbf{X}) = \frac{p(\mathbf{X} | \theta) \ p(\theta)}{p(\mathbf{X})}$$



THOMAS BAYES

$$p(\theta | \mathbf{X}) = \frac{p(\mathbf{X} | \theta) p(\theta)}{p(\mathbf{X})}$$

**Garbage in:** arbitrary priors

**Garbage out:** uncontrollable error bars

~ Michael I. Jordan, *MLSS* (2017)



THOMAS BAYES

$$p(\theta | \mathbf{X}) = \frac{p(\mathbf{X} | \theta) p(\theta)}{p(\mathbf{X})}$$

**How should we think about priors  
for machine learning models?**



THOMAS BAYES

EXAMPLE

---

## Priors for Modeling Coin Tosses

---

---

# Priors for Modeling Coin Tosses

---

$$\mathbf{X} = \{X_n \in \{0, 1\}\}_{n=1}^N$$

$$p(\mathbf{X}|\theta) = \text{Binomial}(N, \theta)$$



$$p(\theta) = \text{Beta}(\alpha_0, \beta_0)$$



---

# Priors for Modeling Coin Tosses

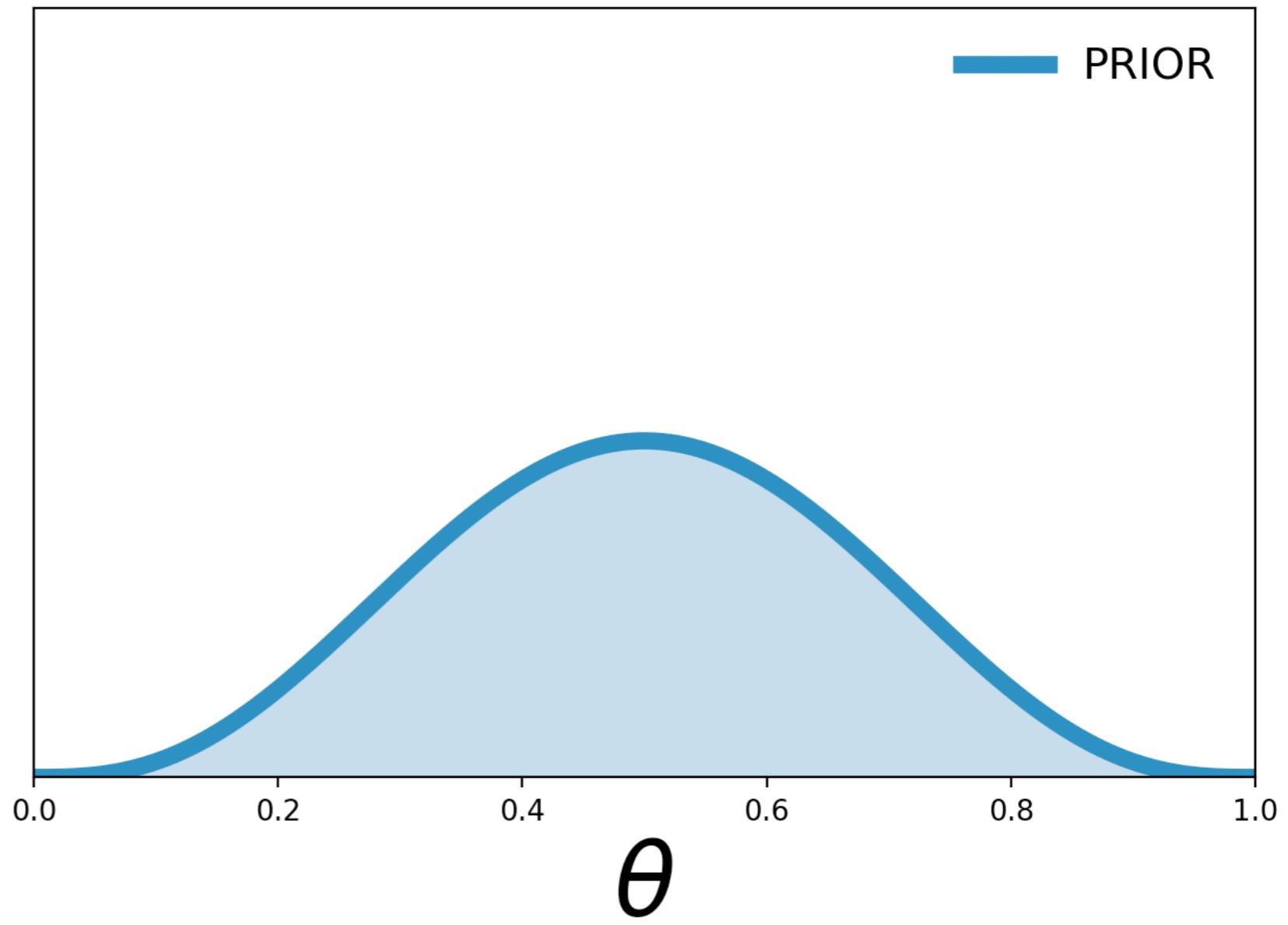
---

$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{ Beta}(\alpha_0, \beta_0)$$



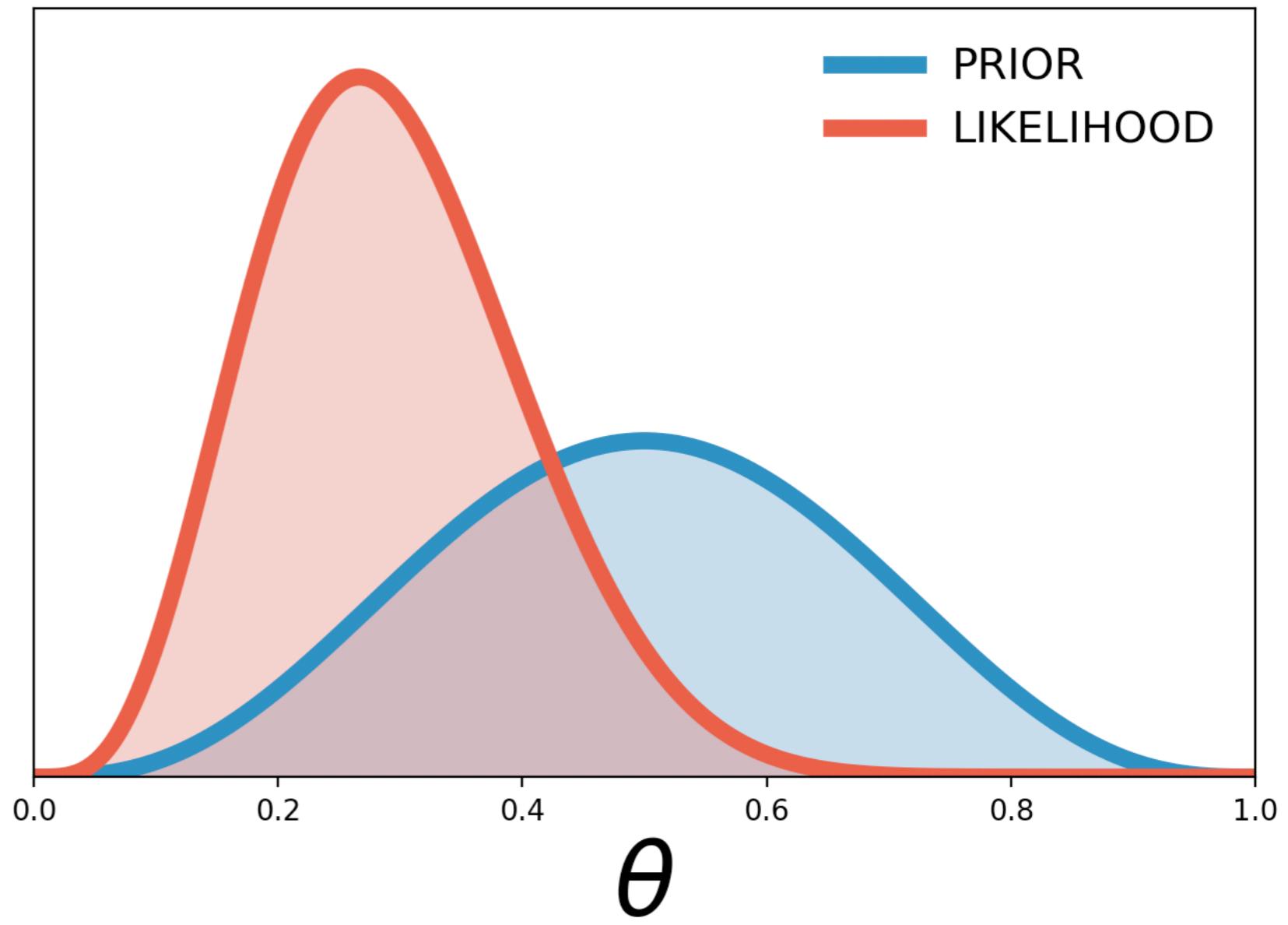
# Priors for Modeling Coin Tosses

$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{ Beta}(4, 4)$$



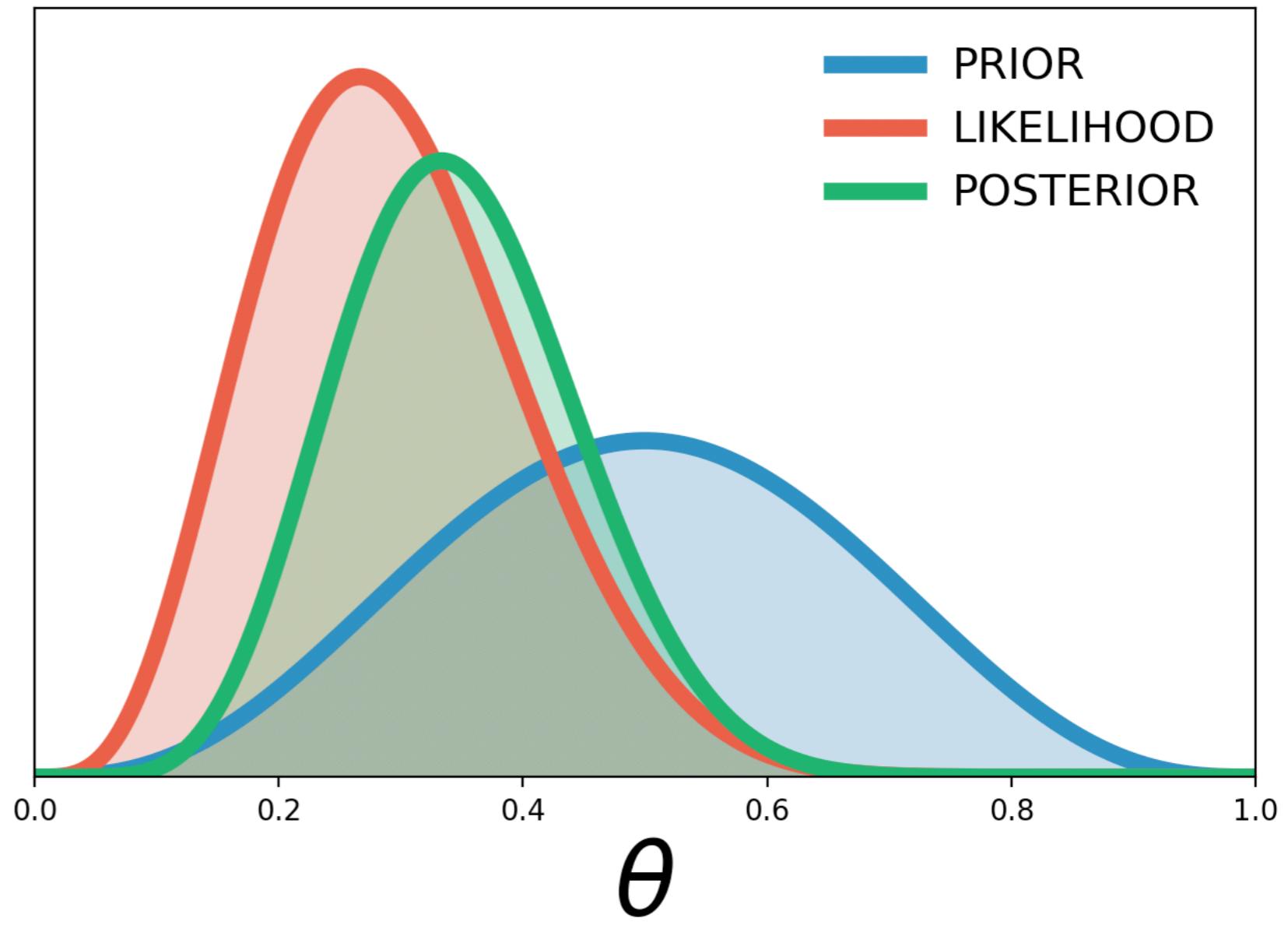
# Priors for Modeling Coin Tosses

$$p(\theta|X) \propto \text{Binomial}(15, \theta) \text{Beta}(4, 4)$$



# Priors for Modeling Coin Tosses

$$p(\theta|X) \propto \text{Binomial}(15, \theta) \text{Beta}(4, 4)$$



---

# Priors for Modeling Coin Tosses

---

$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{ Beta}(\alpha_0, \beta_0)$$

...But what if we have *even stronger* prior information?

# Priors for Modeling Coin Tosses

$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{ Beta}(\alpha_0, \beta_0)$$



**Harvey Dent**  
aka “Two-Face”  
*The Dark Knight, 2008*

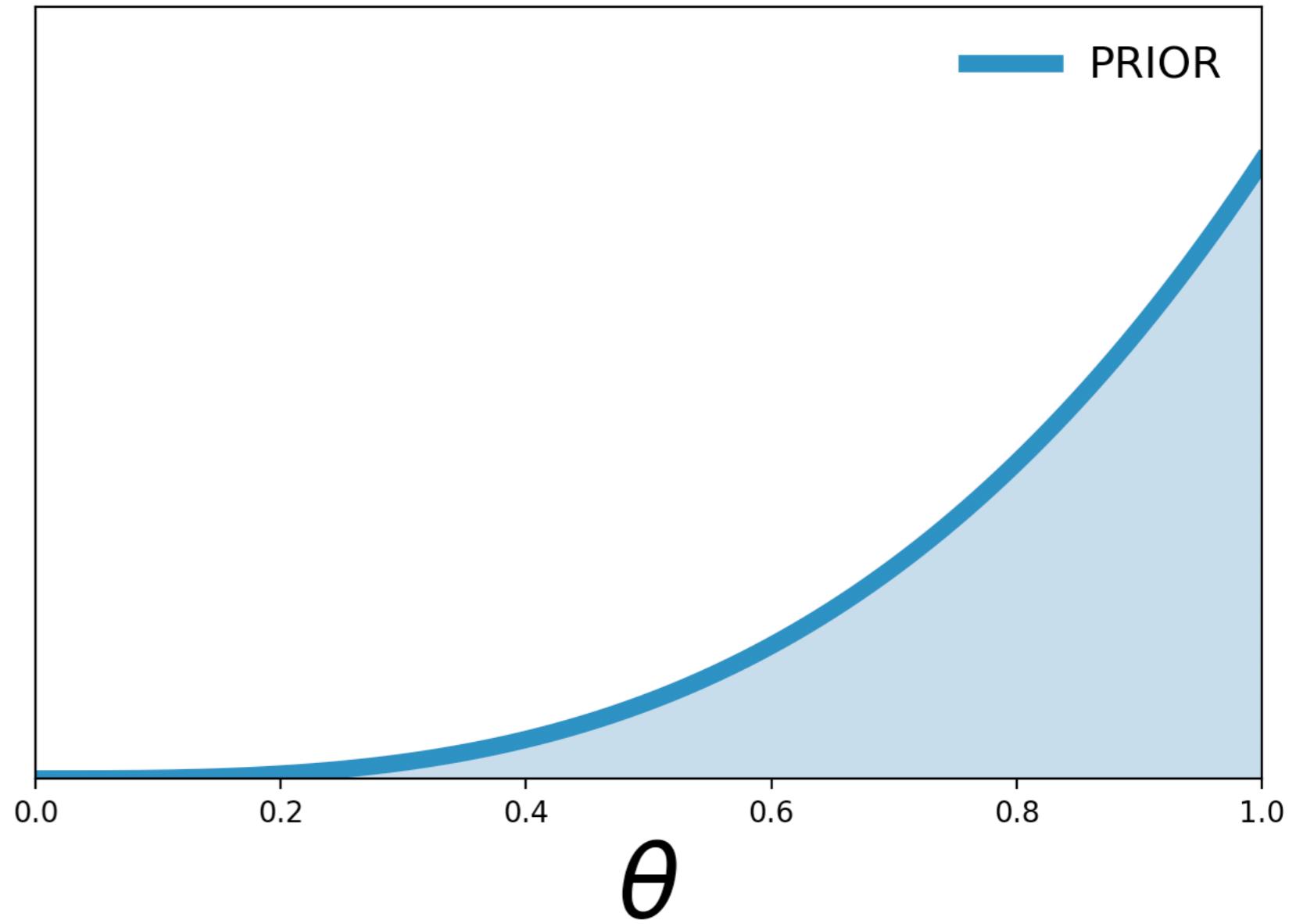
...But what if we have *even stronger* prior information?

# Priors for Modeling Coin Tosses

$$p(\theta | \mathbf{X}) \propto \text{Binomial}(N, \theta) \text{ Beta}(4, 1)$$



**Harvey Dent**  
aka “Two-Face”  
*The Dark Knight, 2008*



# Objective Priors

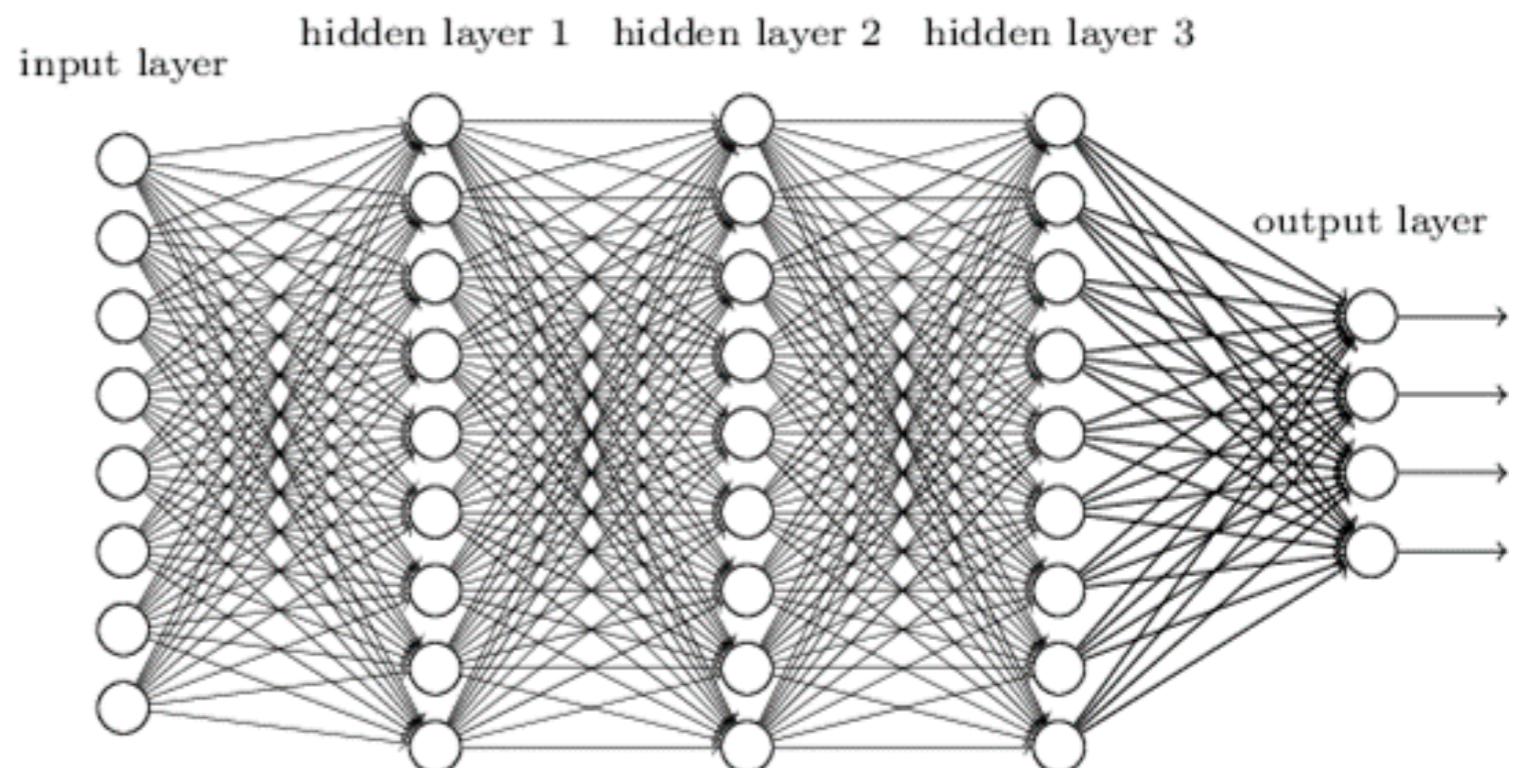
...But what if we have no information?

Say we are completely agnostic about  $\Theta$ ;  
What prior represents a *state of ignorance*?

# Objective Priors

...But what if we have no information?

Say we are completely agnostic about  $\Theta$ ;  
What prior represents a *state of ignorance*?

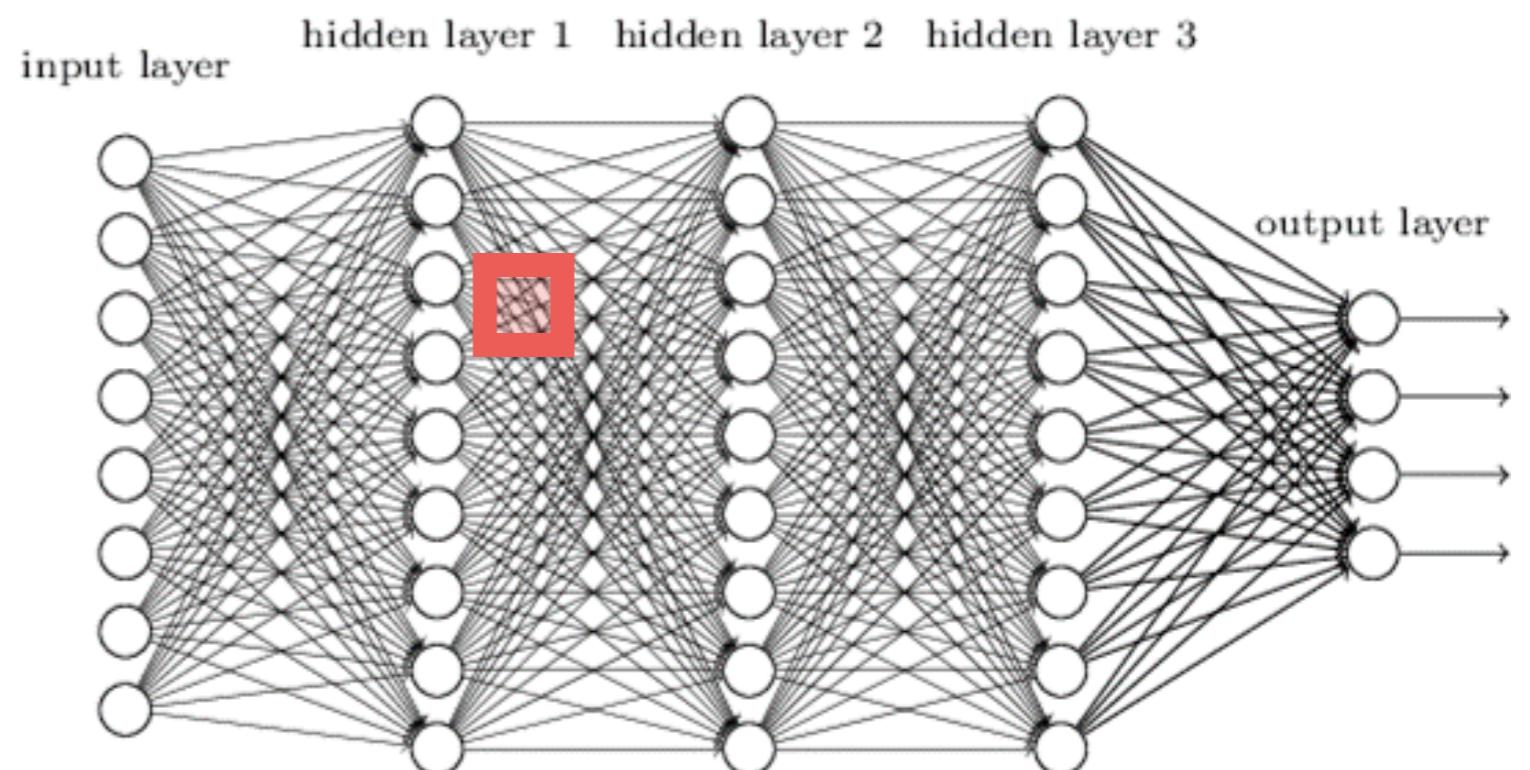


# Objective Priors

...But what if we have no information?

Say we are completely agnostic about  $\Theta$ ;  
What prior represents a *state of ignorance*?

What prior beliefs  
do I have about  
the 341st weight?

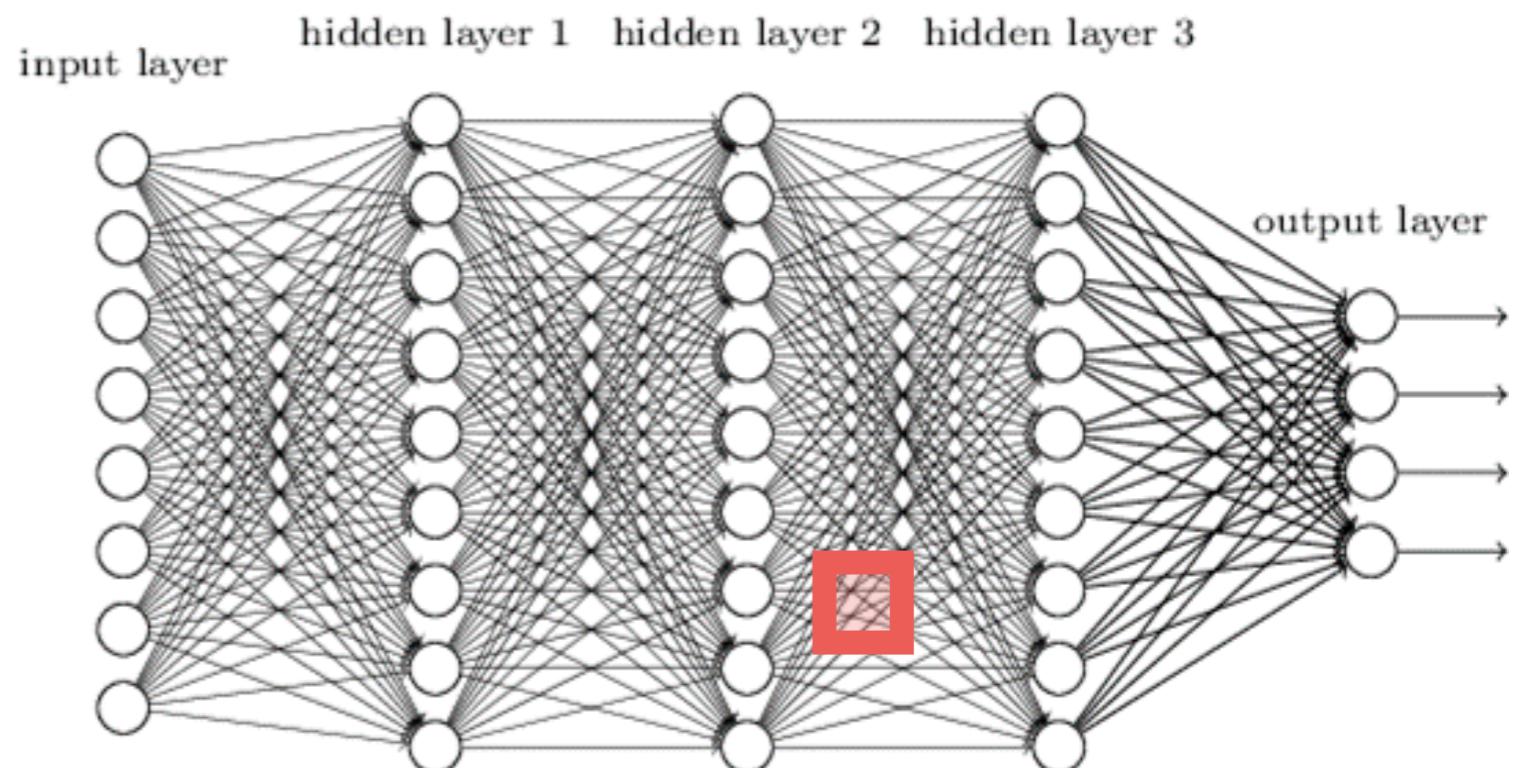


# Objective Priors

...But what if we have no information?

Say we are completely agnostic about  $\Theta$ ;  
What prior represents a *state of ignorance*?

What prior beliefs  
do I have about  
the 833rd weight?  
etc...



---

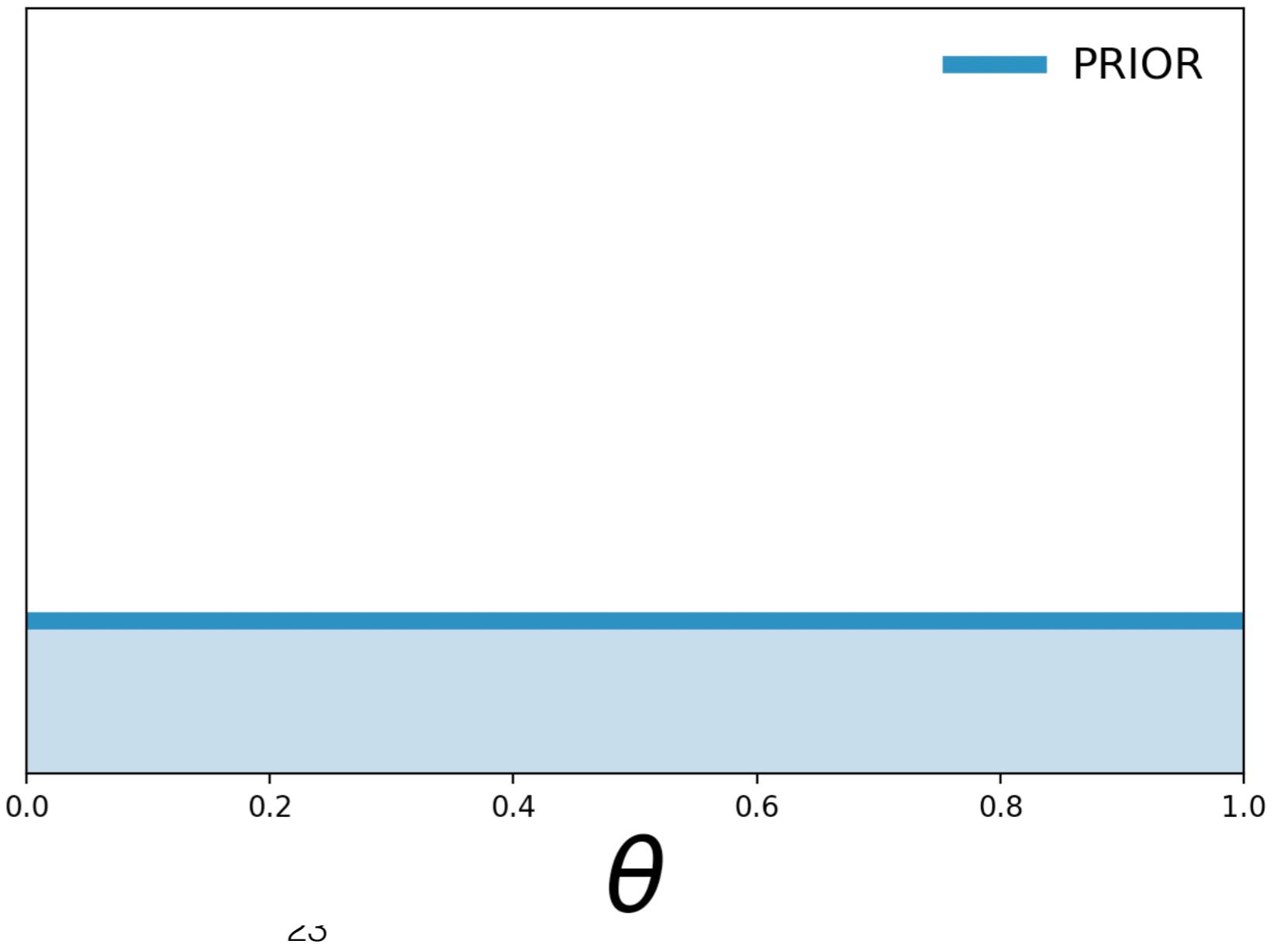
# Objective Priors for Coin Tosses

---

$$p(\theta | \mathbf{X}) \propto \text{Binomial}(N, \theta) \text{ Beta}(\alpha_0, \beta_0)$$

# Objective Priors for Coin Tosses

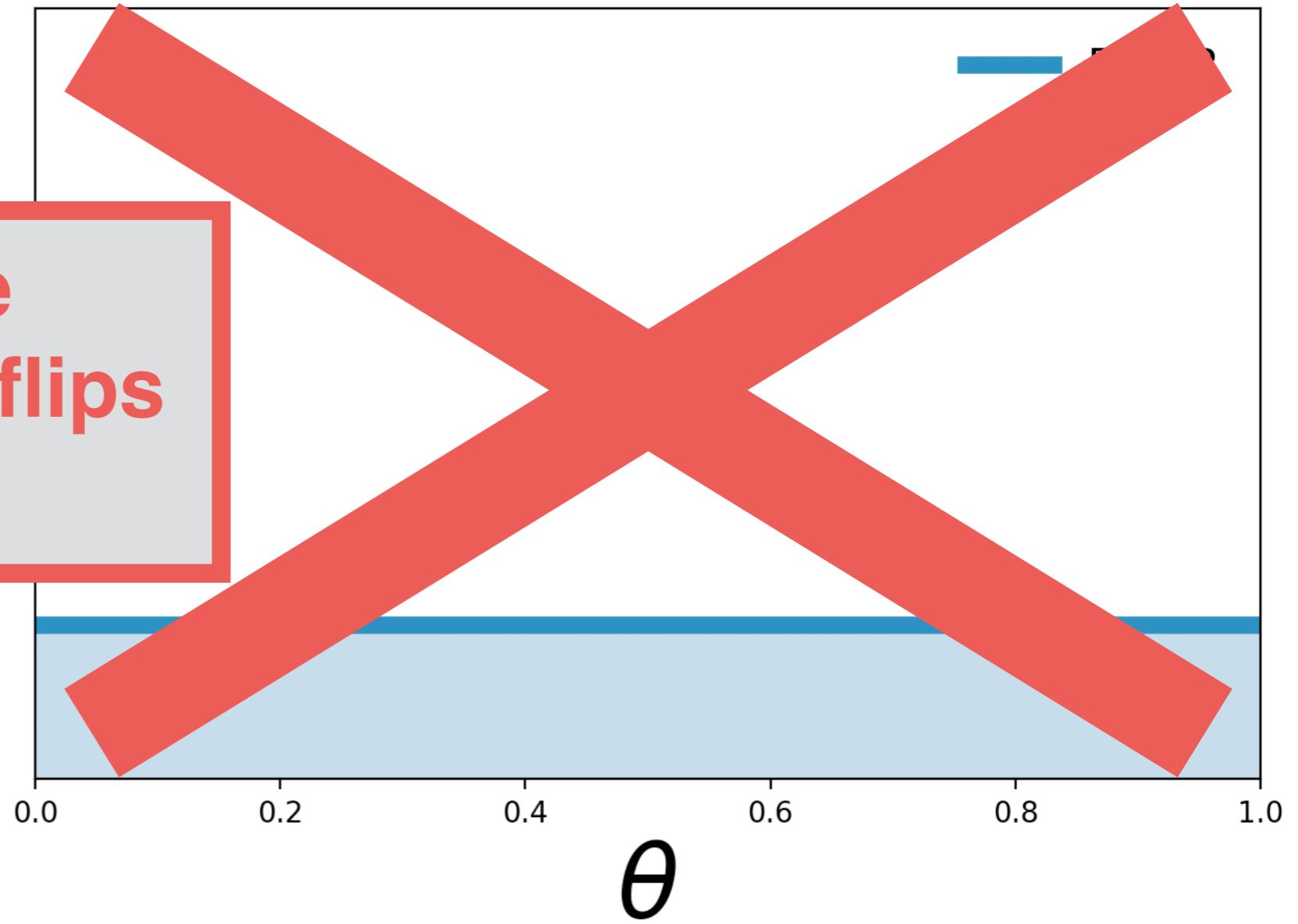
$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{ Beta}(1, 1)$$



# Objective Priors for Coin Tosses

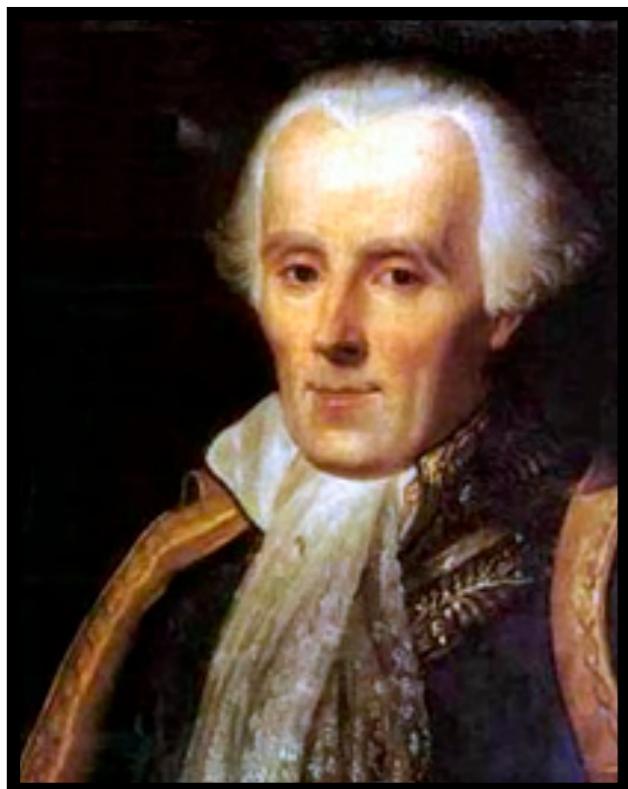
$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{ Beta}(1, 1)$$

Implies that we've  
already seen two flips  
one heads, one tails.

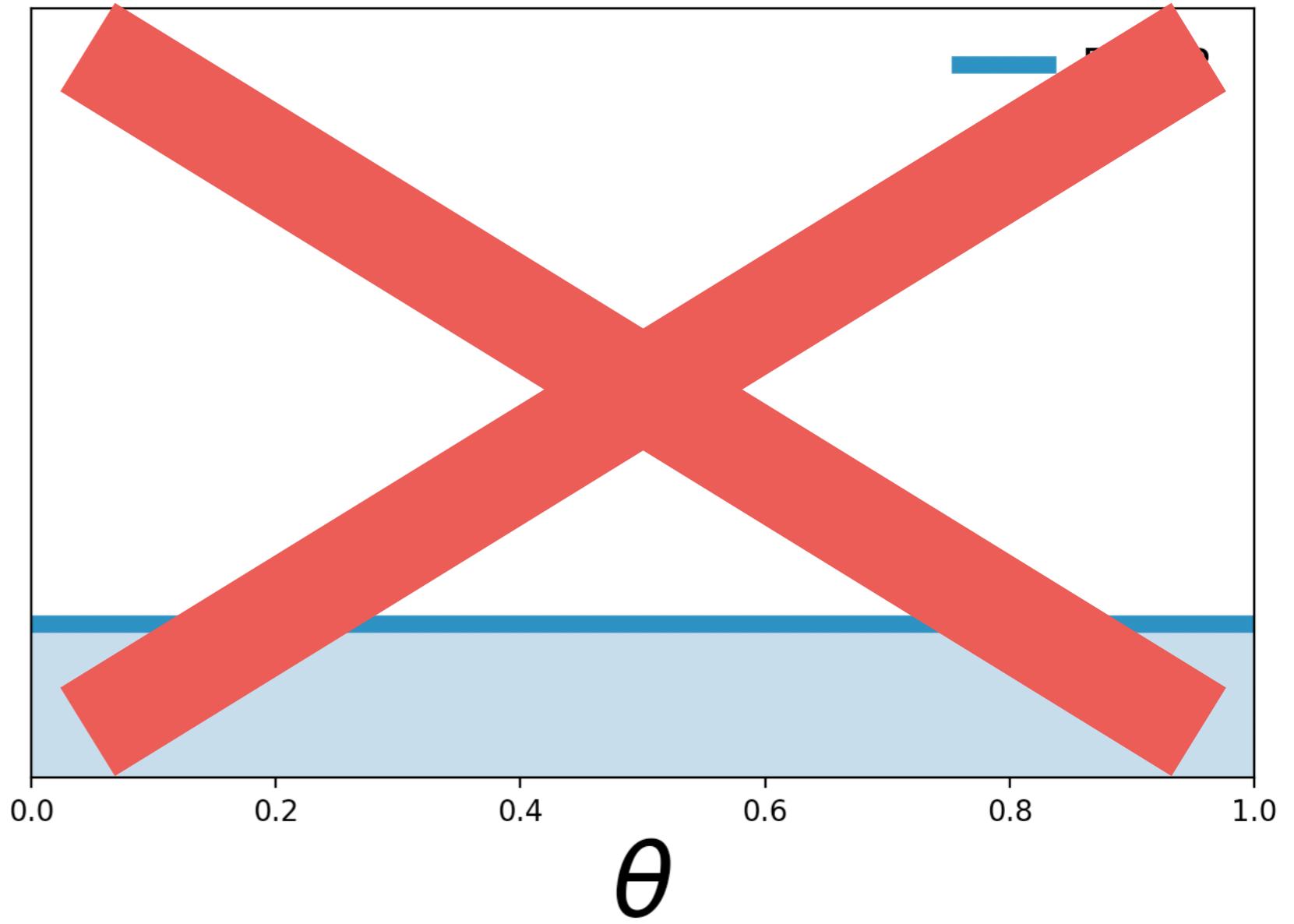


# Objective Priors for Coin Tosses

$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{ Beta}(1, 1)$$

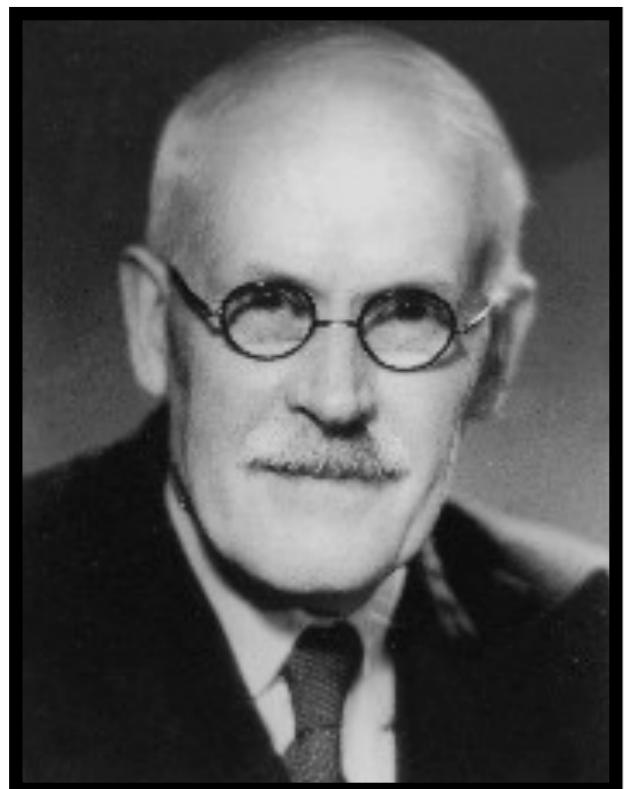


Pierre-Simon Laplace  
(1814)

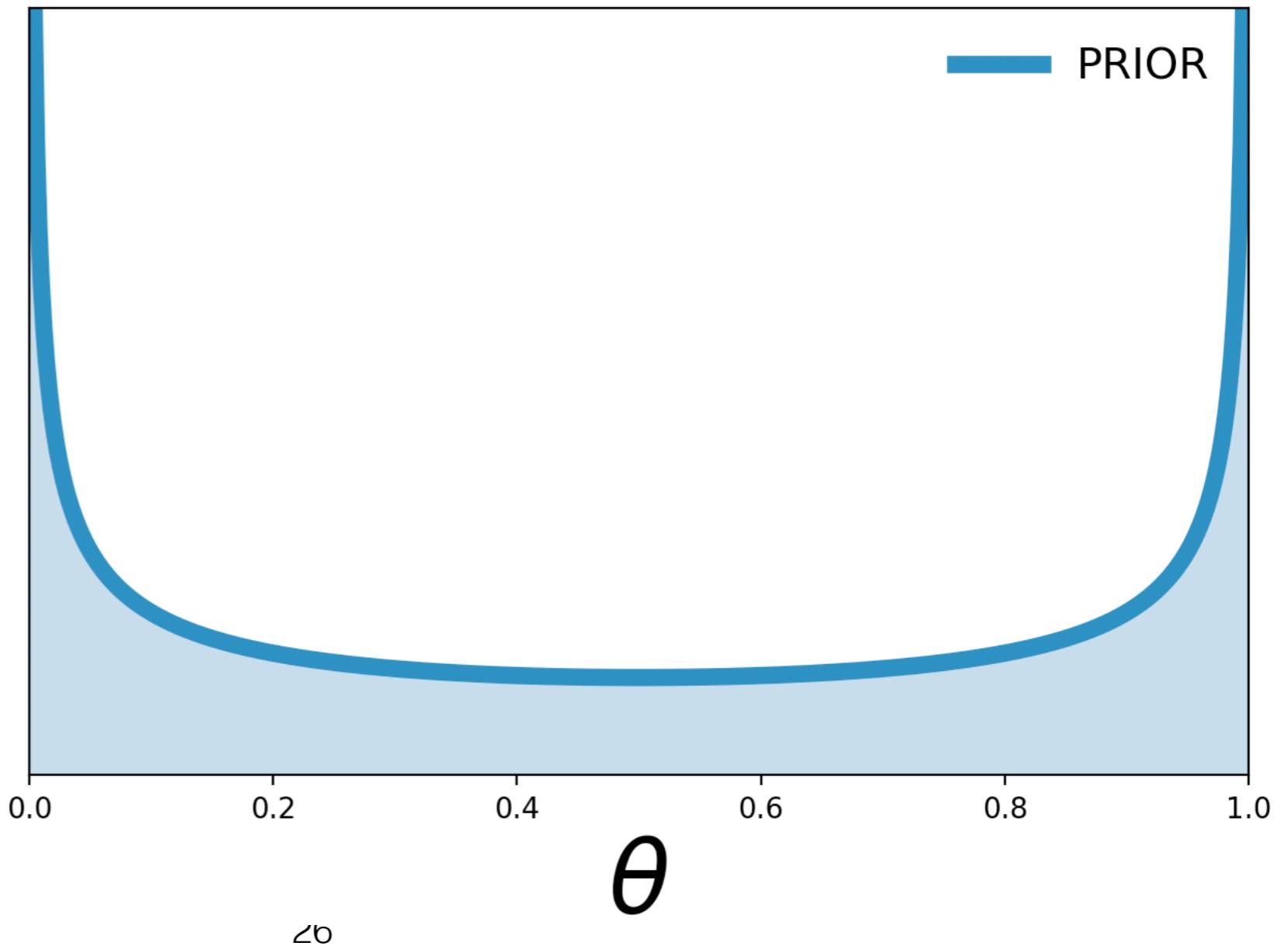


# Objective Priors for Coin Tosses

$$p(\theta|X) \propto \text{Binomial}(N, \theta) \text{Beta}(0.5, 0.5)$$



**Sir Harold Jeffreys**  
(1946)



---

# Objective Priors

---

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

---

# Objective Priors

---

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \operatorname{KLD}[p(\theta|\mathbf{X}) || p(\theta)] d\mathbf{X}$$

---

# Objective Priors

---

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \text{ KLD}[p(\theta|\mathbf{X}) || p(\theta)] d\mathbf{X}$$



***Maximize divergence between posterior and prior***

# Objective Priors

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \operatorname{KLD}[p(\theta|\mathbf{X}) || p(\theta)] d\mathbf{X}$$

Averaged over  
marginal likelihood

**Maximize divergence between  
posterior and prior**

---

# Objective Priors

---

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \operatorname{KLD}[p(\theta|\mathbf{X}) || p(\theta)] d\mathbf{X}$$

## More Examples

$$p(\theta) \propto 1$$

$$p(\theta) \propto \frac{1}{\sigma}$$

$$p(\theta) \propto \frac{1}{\sqrt{\lambda}}$$

Gaussian Mean

Gaussian Scale

Poisson Rate

---

# Objective Priors

---

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \operatorname{KLD}[p(\theta|\mathbf{X}) || p(\theta)] d\mathbf{X}$$

## Properties

- Invariant** to model reparametrization.
- Credible intervals** nearly match **confidence intervals**.\*

---

# Objective Priors

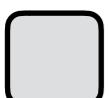
---

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \operatorname{KLD}[p(\theta|\mathbf{X}) || p(\theta)] d\mathbf{X}$$

## Problems / Obstacles



# Objective Priors

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \operatorname{KLD}[p(\theta|\mathbf{X}) || p(\theta)] d\mathbf{X}$$

## Problems / Obstacles

- Calculus of variations problem:** Solve for function



# Objective Priors

[Bernardo, 1979]

$$p^*(\theta) = \operatorname{argmax}_{p(\theta)} \mathbb{I}[\theta, \mathbf{X}]$$

$$= \operatorname{argmax}_{p(\theta)} \int_{\mathbf{X}} p(\mathbf{X}) \text{ KLD}[p(\theta|\mathbf{X}) \mid p(\theta)] d\mathbf{X}$$

## Problems / Obstacles

- Calculus of variations problem:** Solve for function
- Need posterior in closed-form**

## CONTRIBUTION

---

# **Approximating Objective Priors**

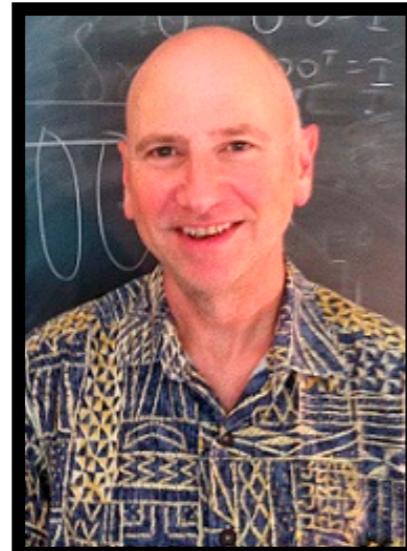
---

**Can we approximate objective priors so as to expand their use in machine learning?**

# Previous Work

**Markov Chain Monte Carlo** [Lafferty & Wasserman, 2001]:  
Use MCMC to draw samples from the objective prior.

$$p^*(\theta) \approx \frac{1}{S} \sum_{s=1}^S \delta[\hat{\theta}_s]$$



## Iterative Markov Chain Monte Carlo Computation of Reference Priors and Minimax Risk

**John Lafferty**  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213  
[lafferty@cs.cmu.edu](mailto:lafferty@cs.cmu.edu)

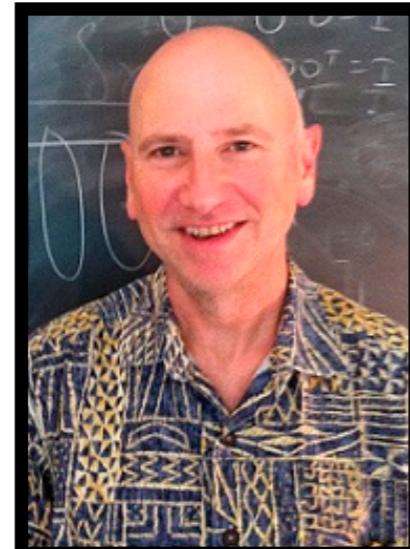
**Larry Wasserman**  
Department of Statistics  
Carnegie Mellon University  
Pittsburgh, PA 15213  
[larry@stat.cmu.edu](mailto:larry@stat.cmu.edu)

# Previous Work

**Markov Chain Monte Carlo** [Lafferty & Wasserman, 2001]:  
Use MCMC to draw samples from the objective prior.

$$p^*(\theta) \approx \frac{1}{S} \sum_{s=1}^S \delta[\hat{\theta}_s]$$

**Drawback:** Have to perform a second modeling step in order to evaluate the prior's density.



## Iterative Markov Chain Monte Carlo Computation of Reference Priors and Minimax Risk

John Lafferty  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213  
[lafferty@cs.cmu.edu](mailto:lafferty@cs.cmu.edu)

Larry Wasserman  
Department of Statistics  
Carnegie Mellon University  
Pittsburgh, PA 15213  
[larry@stat.cmu.edu](mailto:larry@stat.cmu.edu)

---

# A Parametric Approach

---

We propose optimizing a **parametric approximation**.

$$p^*(\theta) \approx p(\theta; \underline{\lambda}^*)$$

Optimize w.r.t.  $\lambda$  instead  
of function itself.

---

# A Parametric Approach

---

We propose optimizing a **parametric approximation**.

$$p^*(\theta) \approx p(\theta; \underline{\lambda}^*)$$

Optimize w.r.t.  $\lambda$  instead  
of function itself.

Objective is a **lower bound** on the mutual information.

$$\mathbb{I}[\theta, \mathbf{X}] \geq \mathbb{E}_{\theta_\lambda, \mathbf{X}} \left[ \log p(\mathbf{X}|\theta) - \max_{\theta_\lambda} \log p(\mathbf{X}|\theta) \right]$$

$$\approx \frac{1}{S} \sum_{s=1}^S \text{KLD} \left[ p(\mathbf{X}|\hat{\theta}_s) \parallel p(\mathbf{X}|\hat{\theta}_{\max}) \right]$$

Samples depend on  $\lambda$ .

---

# A Parametric Approach

---

BENEFITS

---

# A Parametric Approach

---

## BENEFITS

- Comparison to Lafferty & Wasserman's [2001] MCMC:** No need to perform a second modeling step to evaluate the prior.

---

# A Parametric Approach

---

## BENEFITS

- Comparison to Lafferty & Wasserman's [2001] MCMC:** No need to perform a second modeling step to evaluate the prior.
- Conjugate Objective Priors:** Can pick a prior family that is conjugate and then find the member that is nearest to a objective prior.

# A Parametric Approach

## BENEFITS

- Comparison to Lafferty & Wasserman's [2001] MCMC:** No need to perform a second modeling step to evaluate the prior.
- Conjugate Objective Priors:** Can pick a prior family that is conjugate and then find the member that is nearest to a objective prior.
- Proper Objective Priors:** If concerned b/c the true prior is improper, can use an approximation that is proper.

---

# A Parametric Approach

---

## BENEFITS

- Comparison to Lafferty & Wasserman's [2001] MCMC:** No need to perform a second modeling step to evaluate the prior.
  - Conjugate Objective Priors:** Can pick a prior family that is conjugate and then find the member that is nearest to a objective prior.
  - Proper Objective Priors:** If concerned b/c the true prior is improper, can use an approximation that is proper.
- 

**Gaussian mean:**

$$p^*(\mu) \propto 1$$

# A Parametric Approach

## BENEFITS

- Comparison to Lafferty & Wasserman's [2001] MCMC:** No need to perform a second modeling step to evaluate the prior.
- Conjugate Objective Priors:** Can pick a prior family that is conjugate and then find the member that is nearest to a objective prior.
- Proper Objective Priors:** If concerned b/c the true prior is improper, can use an approximation that is proper.

Gaussian mean:

$$p^*(\mu) \propto 1$$

Parametric  
Approximation:

$$N(\mu_\lambda, \sigma_\lambda)$$

# A Parametric Approach

## BENEFITS

- Comparison to Lafferty & Wasserman's [2001] MCMC:** No need to perform a second modeling step to evaluate the prior.
- Conjugate Objective Priors:** Can pick a prior family that is conjugate and then find the member that is nearest to a objective prior.
- Proper Objective Priors:** If concerned b/c the true prior is improper, can use an approximation that is proper.

**Gaussian mean:**

$$p^*(\mu) \propto 1$$

**Parametric  
Approximation:**

$$N(\mu_\lambda, \sigma_\lambda)$$

**Optimization  
Objective:**

$$\mathcal{J}(\lambda) = \|\sigma_\lambda\|_2^2$$

# A Parametric Approach

## BENEFITS

- Comparison to Lafferty & Wasserman's [2001] MCMC:** No need to perform a second modeling step to evaluate the prior.
- Conjugate Objective Priors:** Can pick a prior family that is conjugate and then find the member that is nearest to a objective prior.
- Proper Objective Priors:** If concerned b/c the true prior is improper, can use an approximation that is proper.

**Gaussian mean:**

$$p^*(\mu) \propto 1$$

**Parametric  
Approximation:**

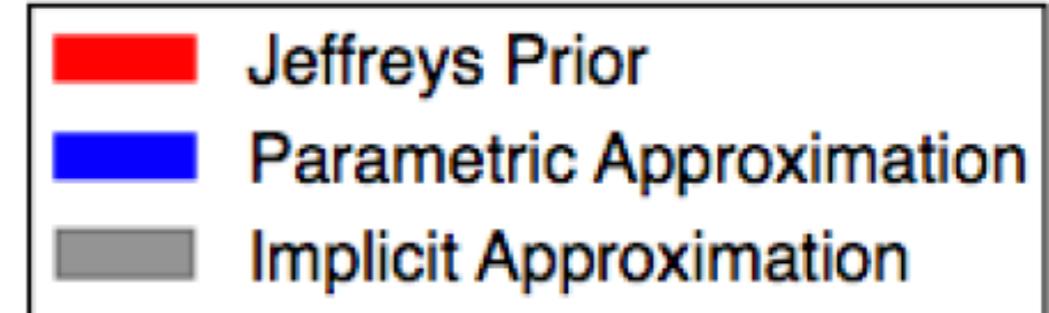
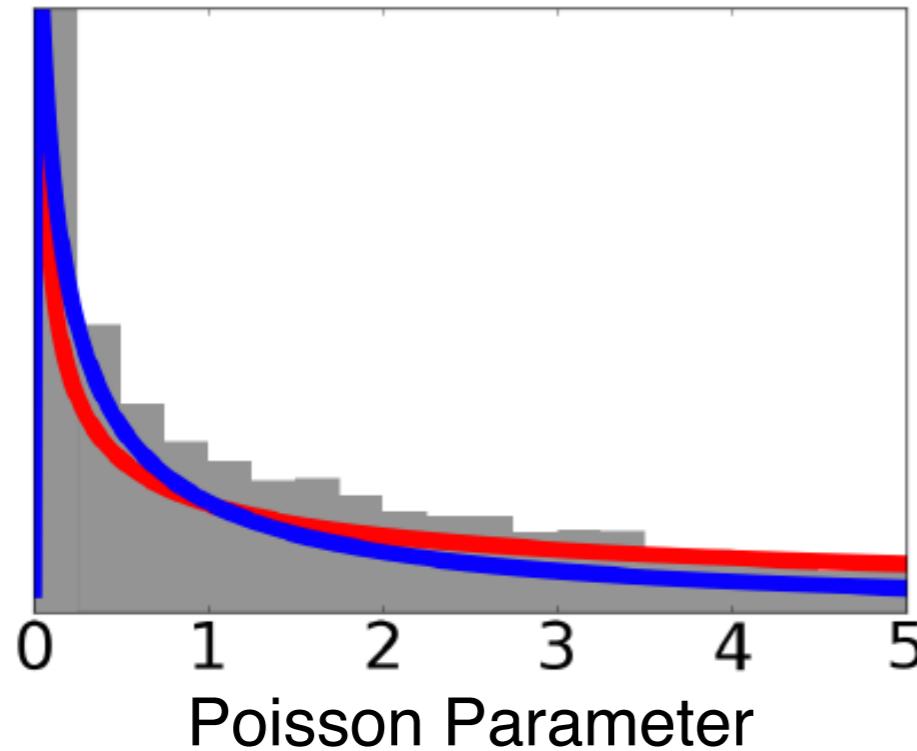
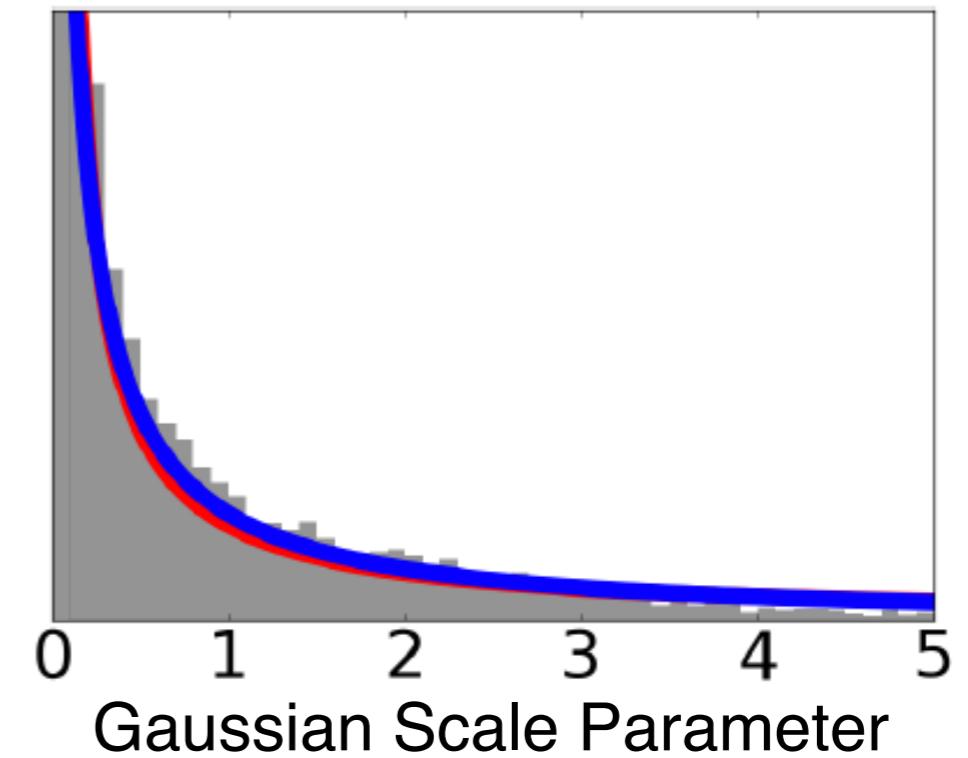
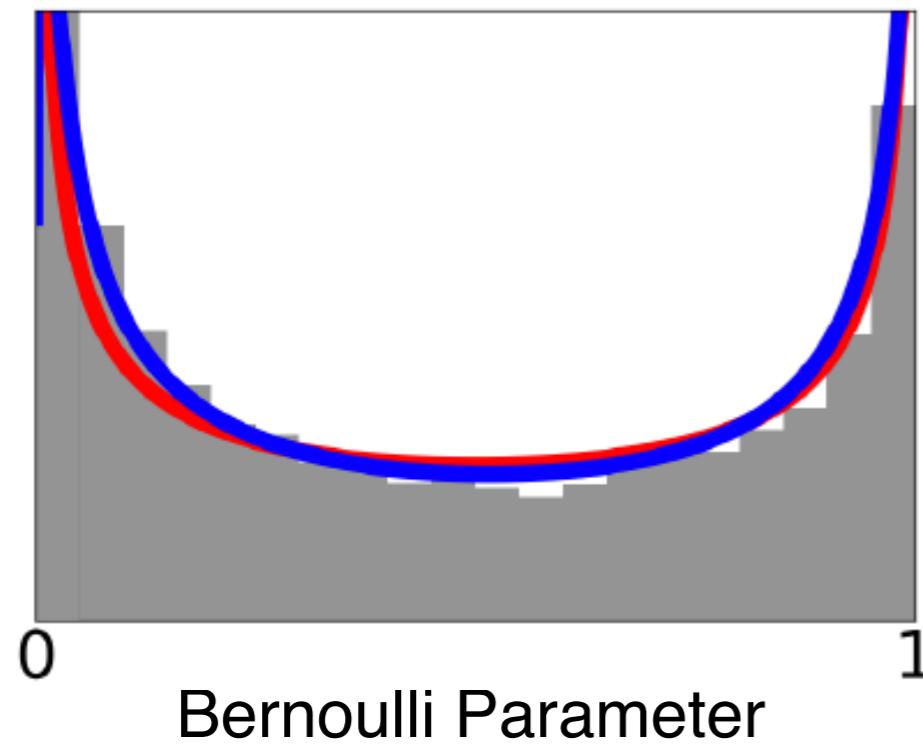
$$N(\mu_\lambda, \sigma_\lambda)$$

**Optimization  
Objective:**

$$\mathcal{J}(\lambda) = \|\sigma_\lambda\|_2^2$$

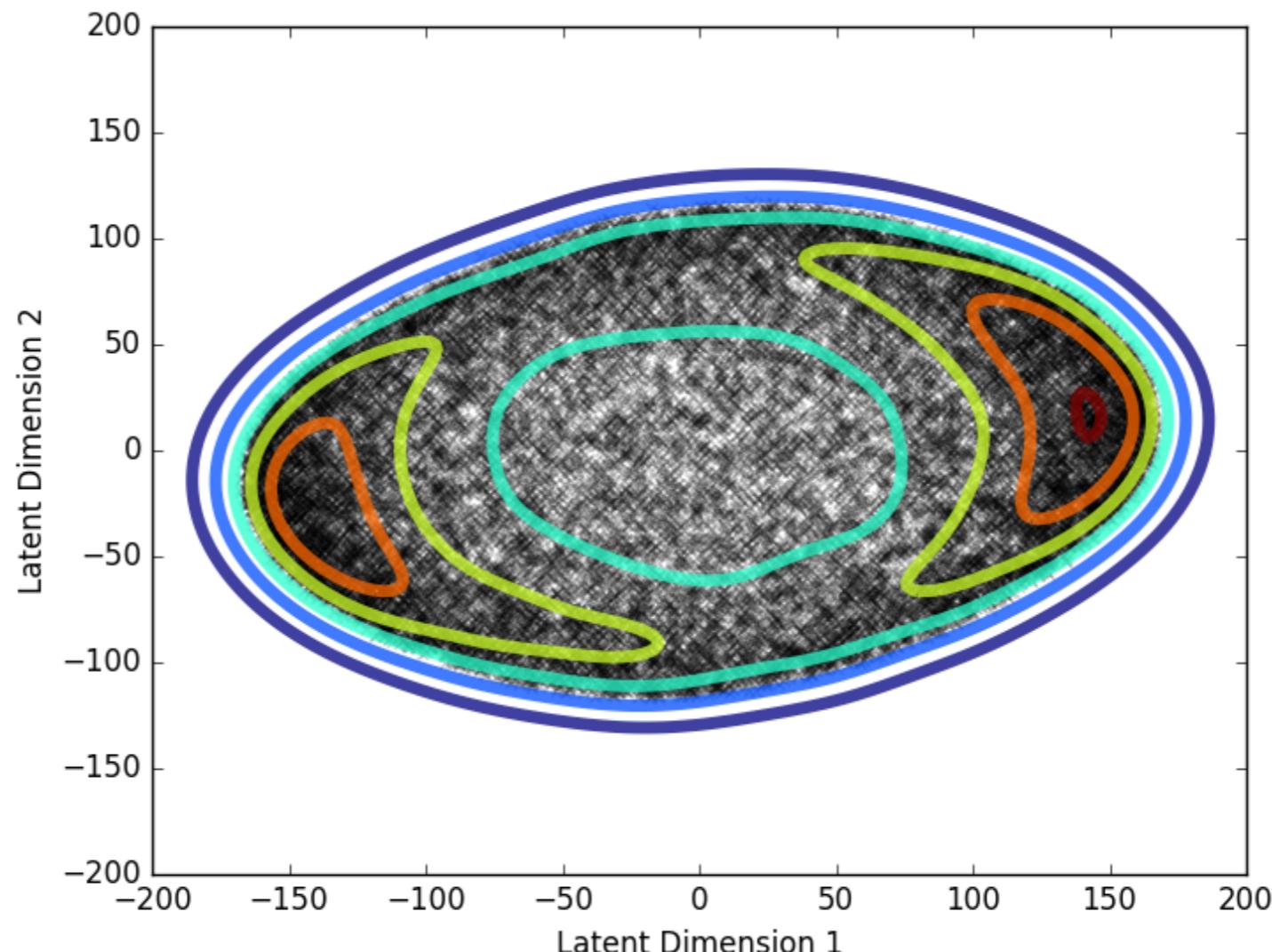
**RESULT:**  $p^* \propto 1 \approx N(\cdot, \infty)$

# Recovering Jeffreys Priors



See paper / dissertation for two-sample tests comparing approximation methods.

# Finding the Variational Autoencoder's Objective Prior



$$p(\mathbf{z})$$

---

# **Conclusions**

---

---

# Priors for Machine Learning

---

## A SPECTRUM OF PRIORS

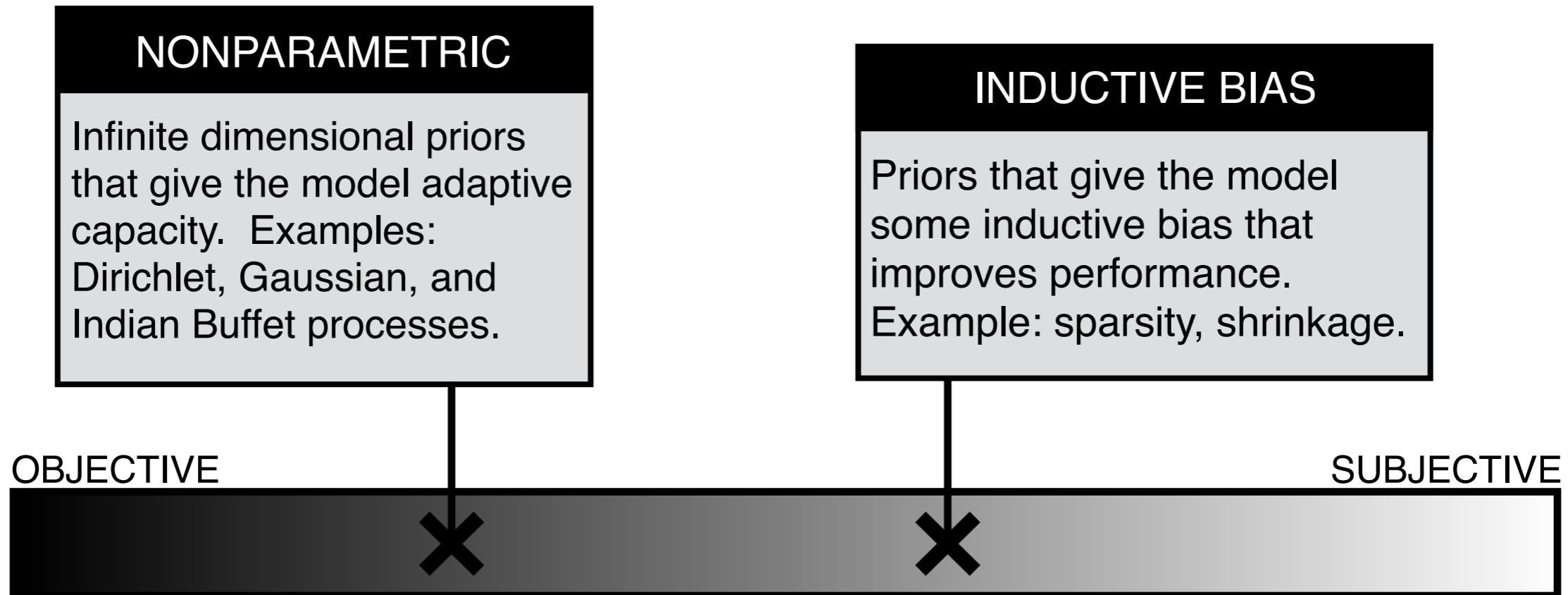
OBJECTIVE

SUBJECTIVE



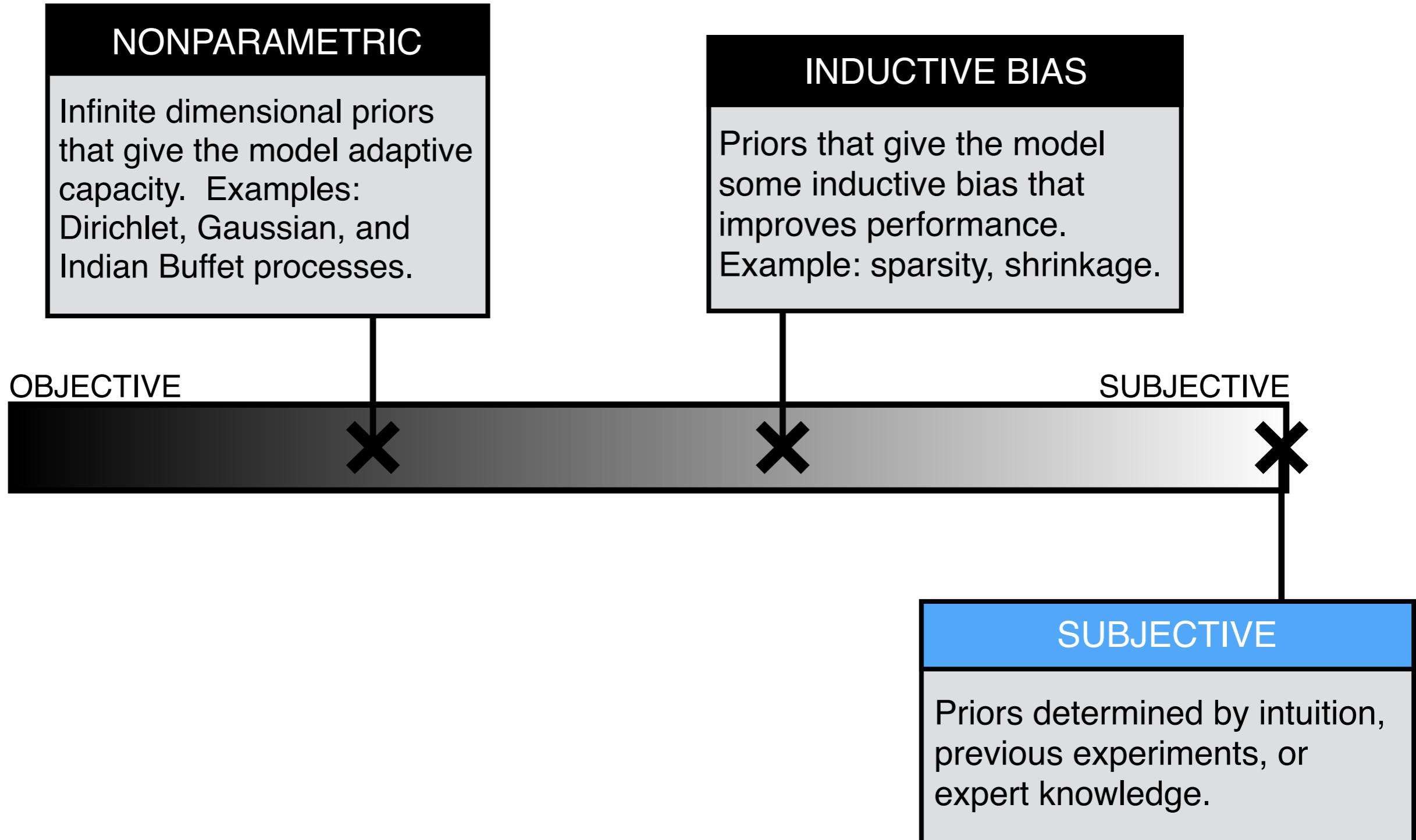
# Priors for Machine Learning

## A SPECTRUM OF PRIORS



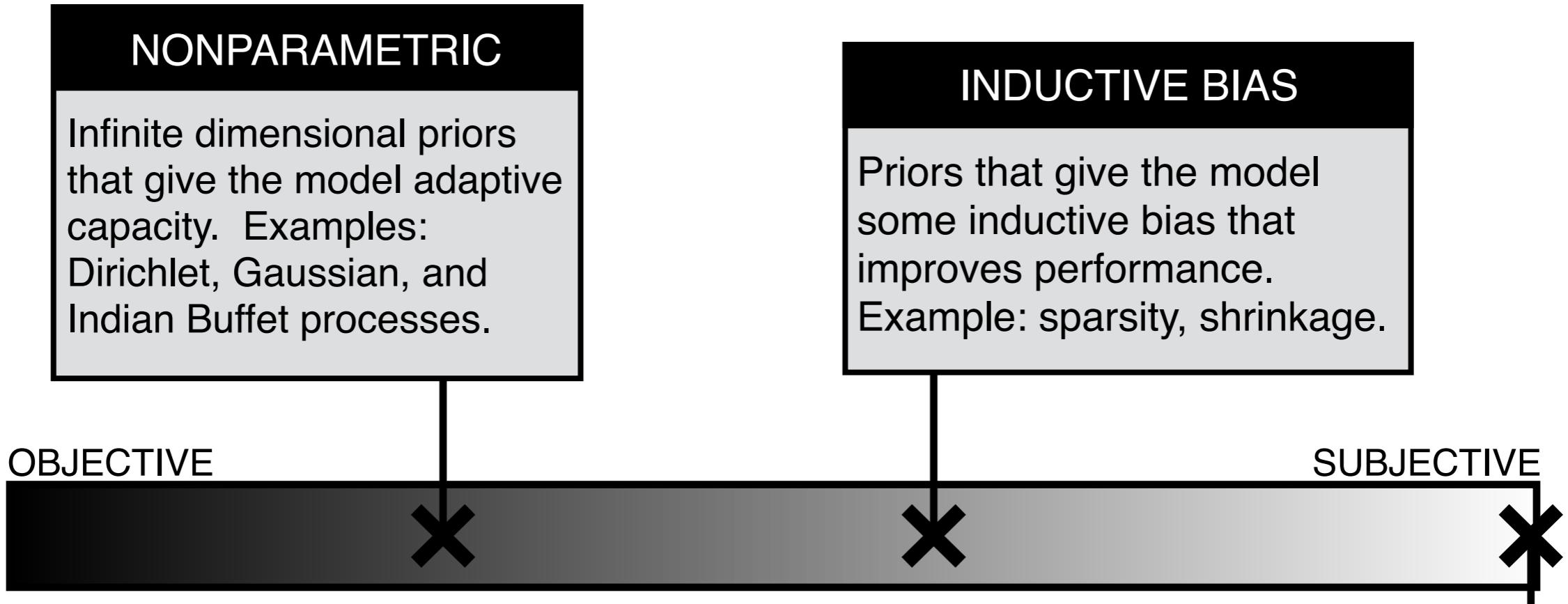
# Priors for Machine Learning

## A SPECTRUM OF PRIORS



# Priors for Machine Learning

## A SPECTRUM OF PRIORS



Published as a conference paper at ICLR 2019

## FUNCTIONAL VARIATIONAL BAYESIAN NEURAL NETWORKS

Shengyang Sun<sup>\*†</sup>, Guodong Zhang<sup>\*†</sup>, Jiaxin Shi<sup>\*‡</sup>, Roger Grosse<sup>†</sup>

<sup>\*</sup>University of Toronto, <sup>†</sup>Vector Institute, <sup>‡</sup>Tsinghua University

{ssy, gdzhang, rgrosse}@cs.toronto.edu, shijx15@mails.tsinghua.edu.cn

## SUBJECTIVE

Priors determined by intuition, previous experiments, or expert knowledge.

# Priors for Machine Learning

## A SPECTRUM OF PRIORS

### NONPARAMETRIC

Infinite dimensional priors that give the model adaptive capacity. Examples: Dirichlet, Gaussian, and Indian Buffet processes.

### INDUCTIVE BIAS

Priors that give the model some inductive bias that improves performance. Example: sparsity, shrinkage.

OBJECTIVE

SUBJECTIVE

THIS TALK

OBJECTIVE

Priors that ‘follow the data,’ perhaps having frequentist properties. Examples: diffuse, Jeffreys, reference, empirical Bayes.

SUBJECTIVE

Priors determined by intuition, previous experiments, or expert knowledge.

---

# Thank you. Questions?

---

## Learning Approximately Objective Priors

In collaboration with...



**Eric Nalisnick**  
Department of Computer Science  
University of California, Irvine  
[enalisni@uci.edu](mailto:enalisni@uci.edu)

**Padhraic Smyth**  
Department of Computer Science  
University of California, Irvine  
[smyth@ics.uci.edu](mailto:smyth@ics.uci.edu)

Padhraic Smyth