

Quiz :

$$p(x|\theta) = \frac{1}{\theta} \mathbb{1}(0 \leq x \leq \theta)$$

$$p(\theta) = \alpha \frac{\beta^\alpha}{\theta^{\alpha+1}} \mathbb{1}(\theta \geq \beta)$$

$$= \text{Pareto}(\alpha, \beta) \quad E[\theta] = \frac{\alpha \beta}{\alpha - 1}$$

Find the posterior

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

$$\prod_{i=1}^n p(x_i | \theta) \quad \text{prior}$$

$$\propto \alpha \left( \frac{\beta}{\theta} \right)^\alpha \mathbb{1}(\theta \geq \beta) \times \prod_{i=1}^n \left( \frac{1}{\theta} \right) \mathbb{1}(0 \leq x_i \leq \theta)$$

$$= \underbrace{\alpha \beta^\alpha}_{\text{does not depend on } \theta} \left( \frac{1}{\theta} \right)^{\alpha+1} \mathbb{1}(\theta \geq \beta) \left( \frac{1}{\theta} \right)^n \prod_{i=1}^n \mathbb{1}(0 \leq x_i \leq \theta)$$

does not depend on  $\theta$

$$\left( \frac{1}{\theta} \right)^{\alpha+n+1} \quad \alpha+1$$

new  $\alpha$  parameter

$$\alpha \left( \frac{1}{\theta} \right)^{\alpha+n+1} \mathbb{1}(\theta \geq \beta) \prod_{i=1}^n \mathbb{1}(0 \leq x_i \leq \theta)$$

$$\frac{1}{\theta^{\alpha+1}}$$

$$\mathbb{1}(\theta \geq \beta, \theta \geq x_1, \theta \geq x_2, \dots, \theta \geq x_n) \mathbb{1}(\theta \geq \beta_*)$$

$$\mathbb{1}(\theta \geq \beta)$$

$$\beta_* = \max(x_1, \dots, x_n, \beta)$$

$$P(\theta | x_{1:n}) = \text{Pareto}(\alpha+n, \max(x_1, \dots, x_n, \beta))$$

$\beta_*$

b.) MAP

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} P(\theta | x_{1:n})$$

$$= \beta_* = \max(x_1, \dots, x_n, \beta)$$

c.) MMSE

$$\hat{\theta}_{\text{MMSE}} = E[\theta | x_{1:n}] \Rightarrow$$

based on the  
parameters  
of the  
posterior.

$$= \frac{(\alpha+n) \max(x_1, \dots, x_n, \beta)}{(\alpha+n) - 1}$$

$$p(\theta) \propto \mathbb{1}(\theta \geq \beta) \times \theta^{-(\alpha+1)}$$

$$= \text{Pareto}(\alpha, \beta)$$

$$E[\theta] = \frac{\alpha \beta}{\alpha - 1}$$

$$p(\theta | X_{1:n}) = \underbrace{\left(\frac{1}{\theta}\right)^n}_{\text{likelihood}} \left(\frac{1}{\theta}\right)^{d+1}$$

$$= \left(\frac{1}{\theta}\right)^{\underbrace{n+d+1}_{\alpha}} \mathbb{1}(\theta \geq \beta_*)$$

- 
- Choice of Prior
  - Calibration of Prior Distributions
  - Theoretical Aspects
  - Posterior - Predictive

Non-informative prior

Weakly informative prior

Informative prior

posterior  $\propto$  likelihood

---


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

Requires that:  $p(\theta) \propto \text{constant}$

$$\int_{-\infty}^{\infty} p(\theta) d\theta \neq 1 \quad > \infty$$

For example:  $X_i \sim N(\mu, \sigma^2)$

- $\mu$  is unknown
- $\sigma^2$  is known

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

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

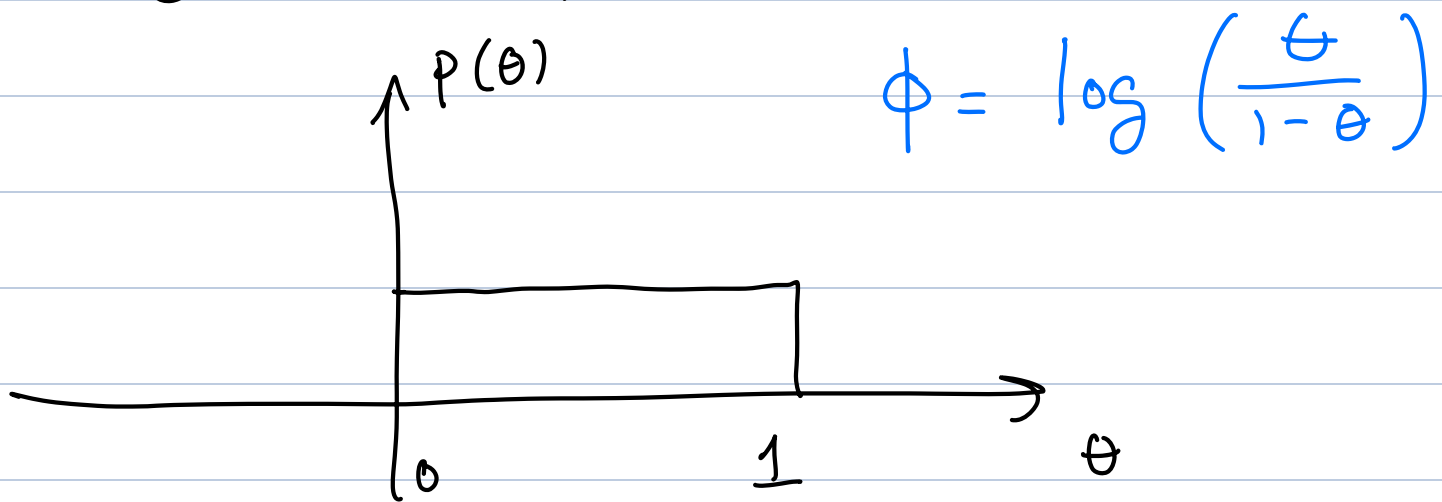
$$= N\left(\mu | \bar{X}_n, \frac{\sigma^2}{n}\right)$$

Improper prior can actually lead to a proper posterior distribution

$$X_i \sim \text{Bernoulli}(\theta)$$

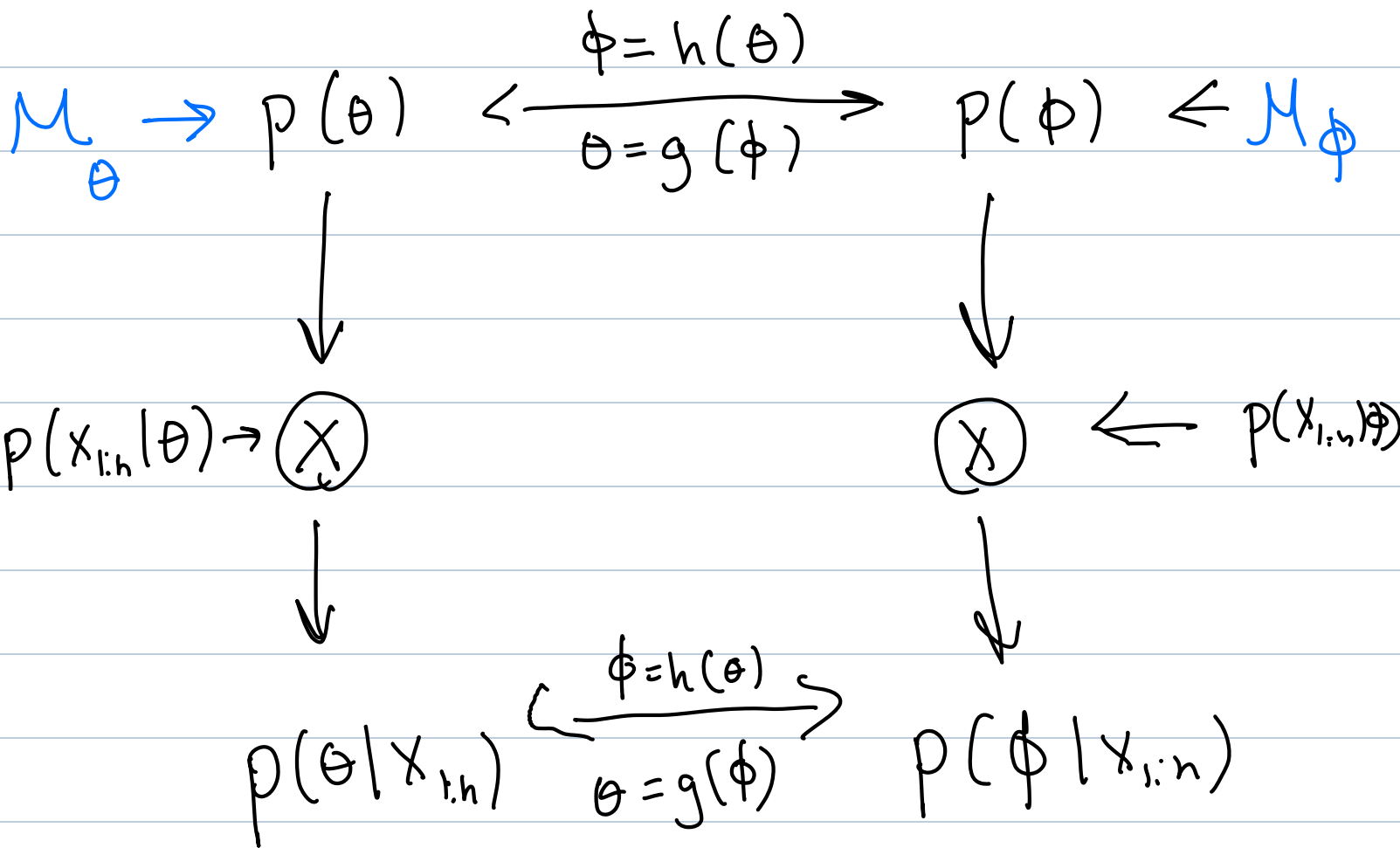
$\theta$ : probability of success

$$\theta \sim \mathcal{U}(0, 1)$$



- The problem of using a flat prior is that although it is "flat" in one parameterization, it may be very informative in another

Jeffreys' Principle



Jeffreys' Prior: Let  $\theta$  be a random variable that parameterizes another random variable  $X_1, \dots, X_n$ , where  $X_1, \dots, X_n$  are iid. The Jeffreys' prior is given by:

$$p(\theta) \propto \sqrt{I(\theta)}$$

$$I(\theta) = -\mathbb{E} \left[ \left( \frac{\partial^2 \log p(x|\theta)}{\partial \theta^2} \right) \middle| \theta \right]$$

Claim: The Jeffreys' prior satisfies Jeffreys' invariance principle

Change of variables Formulas

$$p(\phi) = p(\theta) \left| \frac{dg}{d\phi} \right|$$

where  $g(\phi) = h^{-1}(\phi)$

$$p(\phi) = p(\theta) \left| \frac{dh}{d\theta} \right|^{-1}$$

$$P_{\star}(\theta) = \sqrt{I(\theta)}$$

$$P_{\star}(\phi) = \sqrt{I(\phi)}$$

Inference in  $\theta$ -space:  $p(\theta | x_{1:n}) \propto p(x_{1:n} | \theta) P_{\star}(\theta)$

Inference in  $\phi$ -space:  $p(\phi | x_{1:n}) \propto p(x_{1:n} | \phi) P_{\star}(\phi)$

Start in  $\theta$ -space:

Jeffreys prior

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

$$= p(x_{1:n} | \theta) \left( -\mathbb{E} \left[ \left( \frac{\partial^2 \log p(x | \theta)}{\partial \theta^2} \right) \middle| \theta \right] \right)$$

Let's apply the change-of-variables formula:

posterior of  $\phi$

$$P(\phi | x_{1:n}) \propto P(\theta = g(\phi) | X_{1:n}) \sqrt{I(g(\phi))} \left| \frac{dh}{d\theta} \right|$$

what is this?

$$I(\theta) = -\mathbb{E} \left[ \left( \frac{\partial^2 \log p(x | \theta)}{\partial \theta^2} \right) \middle| \theta \right]$$

substitute

the transformation

$$\phi = h(\theta)$$

$$= -\mathbb{E} \left[ \left( \frac{\partial^2 \log p(x | \phi = h(\theta))}{\partial \phi^2} \right) \left| \frac{dh}{d\theta} \right|^2 \middle| \theta \right]$$

$$= I(\phi) \left| \frac{dh}{d\theta} \right|^2$$

$$P_{\star}(\theta) \propto \sqrt{I(\theta)} = \sqrt{I(\phi)} \left| \frac{dh}{d\theta} \right|$$

- plugging  $I(\theta)$  back to our change of variables, we have

$$p(\phi | x_{1:n}) \propto p(x_{1:n} | \phi) \sqrt{I(\phi)}$$

which is the desired result

## Calibrate a Prior using Moment Matching

- The idea of moment matching is to calibrate a prior distributions theoretical moments to the moments provided by the domain expert
- Solving a system of equations
- Prior has its own parameters and we want to "estimate" them:

$$p(\theta; \alpha)$$

Bayesian model parameter

hyperparameters of the prior

- Moment matching (for only two moments):

$$\mathbb{E}[\theta; \alpha] \approx \mu_{\text{expert}}$$

$$\mathbb{V}[\theta; \alpha] \approx \sigma_{\text{expert}}^2$$

- Let's do an example:

Suppose we have a Bayesian model with an unknown probability parameter  $\theta$  (where  $0 \leq \theta \leq 1$ ). We assume  $\theta$  a priori follows a Beta distribution

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

$$\Leftrightarrow p(\theta) \propto \theta^{\alpha_0 - 1} (1 - \theta)^{\beta_0 - 1}$$

Domain expert tells you that they believe with 95% confidence that

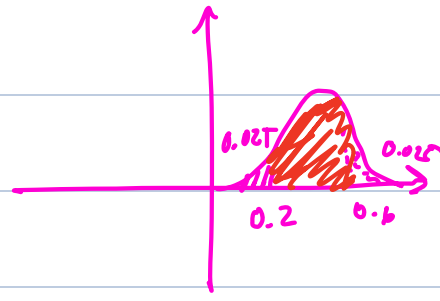
$$\theta \in [0.2, 0.6]$$

$$0.2 = \mu_{\text{expert}} - 2\sigma_{\text{expert}} \approx F(z \leq 0.025)$$

$$0.6 = \mu_{\text{expert}} + 2\sigma_{\text{expert}} \approx F(z \leq 0.75)$$

$z \sim N(\mu_e, \sigma_e^2)$

$$\mu_{\text{expert}} = 0.4$$

$$\sigma_{\text{expert}}^2 = 0.01$$


For a Beta distribution, the mean and variance are given by:

$$0.4 = E[\theta; \alpha_0, \beta_0] = \frac{\alpha_0}{\alpha_0 + \beta_0}$$

$$0.01 = V[\theta; \alpha_0, \beta_0] = \frac{\alpha_0 \beta_0}{(\alpha_0 + \beta_0)^2 (\alpha_0 + \beta_0 + 1)}$$

↑  
set equal  
to domain  
expert  
moments

We then solve this system of equations for  $\alpha_0$  and  $\beta_0$  to obtain the hyperparameters:

$$\alpha_0 = \left( \frac{\mu_{\text{expert}} (1 - \mu_{\text{expert}})}{\sigma_{\text{expert}}^2} - 1 \right) \mu_{\text{expert}}$$

$$\beta_0 = \left( \downarrow \right) (1 - \mu_{\text{expert}})$$

## Theoretical Aspects

Point estimation:

assume there is  
some ground truth  
 $\theta_0$

$$\hat{\theta} = g(X_1, \dots, X_n)$$

function of the  
random sample

$$\hat{\theta} \xrightarrow{\text{a.s.}} \theta_0$$

$$\xrightarrow{P} \theta_0$$

$$\xrightarrow{d} \theta_0$$

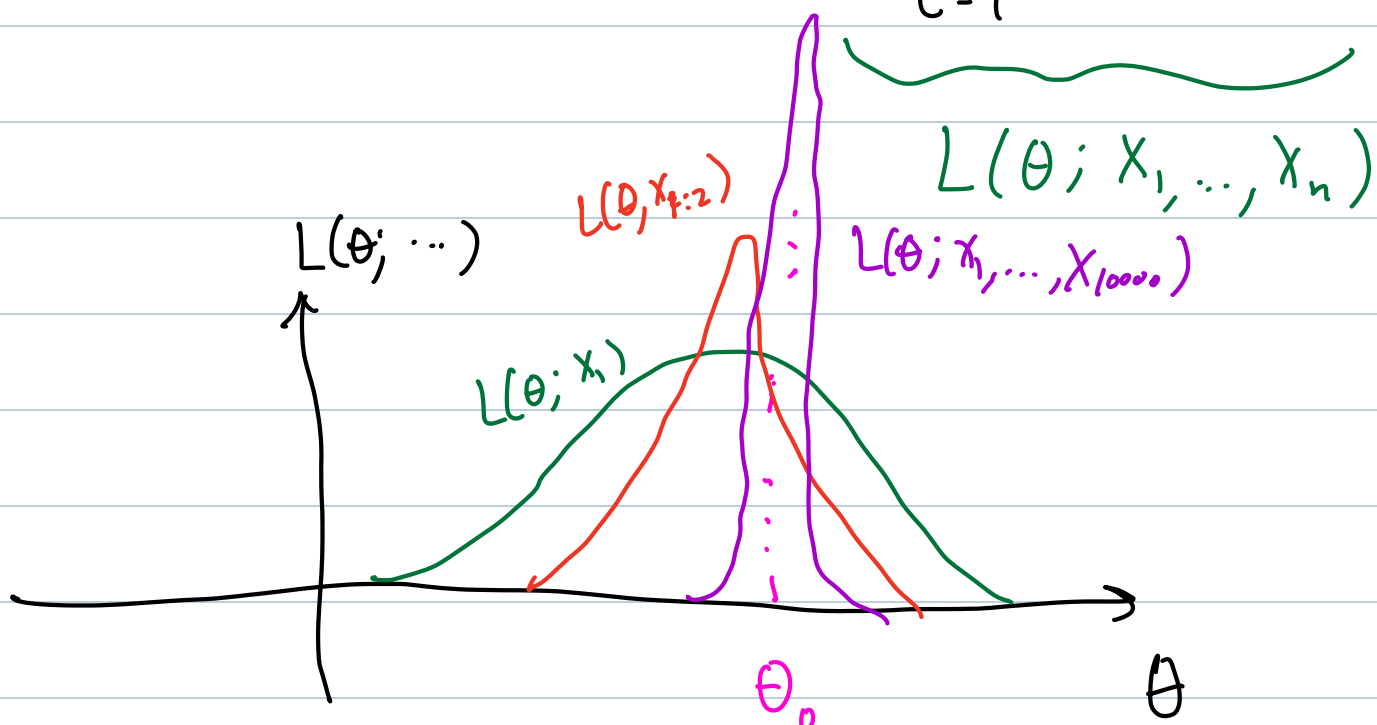
Limiting distribution of MLE

$$\hat{\theta}_{MLE} \sim \mathcal{N}(\theta_0, I^{-1}(\theta))$$

↑  
approximately  
distributed  
according to...

Since in Bayesian estimation, the posterior is the product of the prior and the likelihood:

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



Intuition: As  $n \rightarrow \infty$ , the likelihood dominates and approaches a  $\delta(\theta - \theta_0)$

delta function  
centered at  
 $\theta_0$

Question: What happens to  $p(\theta | X_{1:n})$  for large values of  $n$ ?

• Start with a Taylor expansion of  $\log p(\theta | X_{1:n})$  at  $\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} \log p(\theta | X_{1:n})$ :

$$\log p(\theta | X_{1:n}) = \log p(\hat{\theta}_{\text{MAP}} | X_{1:n}) + (\theta - \hat{\theta}_{\text{MAP}}) \left[ \frac{\partial \log p(\theta | X_{1:n})}{\partial \theta} \right]_{\theta = \hat{\theta}_{\text{MAP}}} + \frac{1}{2} (\theta - \hat{\theta}_{\text{MAP}})^2 \left[ \frac{\partial^2 \log p(\theta | X_{1:n})}{\partial \theta^2} \right]_{\theta = \hat{\theta}_{\text{MAP}}} + \text{higher order terms}$$

$$\text{Higher-order terms} \propto O\left(|\theta - \hat{\theta}_{\text{MAP}}|^3\right)$$

Observations:

- The term ① is constant  $\rightarrow$  does not affect the resulting posterior
- The term ② is 0 because  $\hat{\theta}_{\text{MAP}}$  is the maximizer of  $\log p(\theta | X_{1:n})$

- The higher-order terms (3.) will diminish in comparison to the second-order term as  $n \rightarrow \infty$  because  $\hat{\theta}_{MAP} \rightarrow \hat{\theta}_{MLE} \xrightarrow{p} \theta_0$ .

Conclusion: The second-order term becomes the dominant term in the Taylor series expansion and  $\gg$  by analysis

$$\theta \sim \mathcal{N} \left( \hat{\theta}_{MAP}, \frac{1}{\underbrace{I(\theta)}_{\text{Fisher information}} + \underbrace{J(\theta)}_{\text{Fisher information of the prior}}} \right)$$

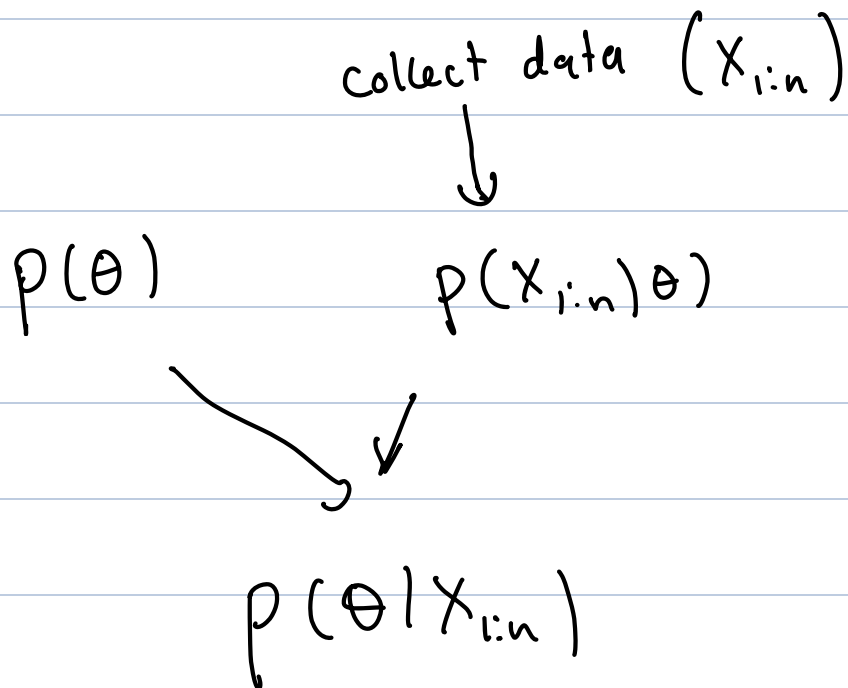
$-\mathbb{E} \left[ \frac{\partial^2 \log p(\theta)}{\partial \theta^2} \right]$

$$-\mathbb{E} \left[ \frac{\partial^2 \log p(x_{1:n} | \theta)}{\partial \theta^2} \right]$$

$$\begin{aligned} \log p(\theta | x_{1:n}) &= \log p(x_{1:n} | \theta) \\ &+ \log p(\theta) \\ &- \log p(x_{1:n}) \end{aligned}$$

Theorem: Bernstein von-Mises theorem

# Posterior - Predictive Checks



- We want to know if our Bayesian is well-calibrated to the observed data after fitting the model

Predictive distribution

$$p(\tilde{x} | X_{1:n}) = \int p(\tilde{x}, \theta | X_{1:n}) d\theta$$

apply chain rule of probability

$$= \int \underbrace{p(\tilde{x} | \theta, X_{1:n})}_{\text{}} p(\theta | X_{1:n}) d\theta$$

$\tilde{x} \perp\!\!\!\perp X_{1:n}$  when  $\theta$  is given

$$p(\tilde{x} | \theta, X_{1:n}) \equiv p(\tilde{x} | \theta)$$

$$p(\tilde{x} | X_{1:n}) = \int p(\tilde{x} | \theta) p(\theta | X_{1:n}) d\theta$$

$$= \mathbb{E} [p(\tilde{x} | \theta) | X_{1:n}]$$

conditional expectation

$S$ : # of samples

$$\theta^{(1)}, \dots, \theta^{(S)} \sim p(\theta | X_{1:n})$$

$$\approx \frac{1}{S} \sum_{s=1}^S p(\tilde{x} | \theta^{(s)})$$

Posterior - predictive check  $\Rightarrow$  test to see if my historical data is contained in the support of the predictive distribution

