

# Single parameter inference

STA 427/527, Fall 2019, Xin Wang

## 2.1 Point Estimate

- In Bayesian inference, point estimator can be posterior mean, posterior median and posterior mode
- Loss function measures the “loss” generated by estimating  $\theta$  with the estimator  $\hat{\theta}$ .

*different loss functions lead to different  $\hat{\theta}$ s*

- Linear absolute loss:  $L_1(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$  *point estimator*
- Quadratic loss:  $L_2(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$
- Zero-one loss:  $L_3(\hat{\theta}, \theta) = \begin{cases} 0 & |\hat{\theta} - \theta| \leq \epsilon \\ 1 & |\hat{\theta} - \theta| > \epsilon \end{cases}$

*main one used in literature*

*(least common (often hard to calculate))*

- Expected loss:

*↳ this is a function of  $\hat{\theta}$*

$$E[L(\hat{\theta}, \theta) | \mathbf{y}] = \int L(\hat{\theta}, \theta) p(\theta | \mathbf{y}) d\theta$$

*↳ posterior dist. of  $\theta$*

- Bayesian estimators:

- Posterior mean:

*Goal is to find  $\hat{\theta}$  that minimizes expected loss*

$$\min_{\hat{\theta}} E[L_2(\hat{\theta}, \theta) | \mathbf{y}] = \min_{\hat{\theta}} \int (\hat{\theta} - \theta)^2 p(\theta | \mathbf{y}) d\theta$$

$$\rightarrow \int (\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2) p(\theta | \mathbf{y}) d\theta$$

$$= \int \hat{\theta}^2 p(\theta | \mathbf{y}) d\theta - \int 2\hat{\theta}\theta p(\theta | \mathbf{y}) d\theta + \int \theta^2 p(\theta | \mathbf{y}) d\theta$$

$$= \hat{\theta}^2 \int p(\theta | \mathbf{y}) d\theta - 2\hat{\theta} \int \theta p(\theta | \mathbf{y}) d\theta + C$$

*$C = \int \theta^2 p(\theta | \mathbf{y}) d\theta$  (constant w.r.t  $\hat{\theta}$ )*

$$= \hat{\theta}^2 - 2\hat{\theta} E(\theta | \mathbf{y})$$

$$\rightarrow \frac{d}{d\hat{\theta}} [\dots] = 2\hat{\theta} - 2E(\theta | \mathbf{y})$$

$$\rightarrow \hat{\theta} = E(\theta | \mathbf{y})$$

*↳ posterior mean*

- Posterior median:

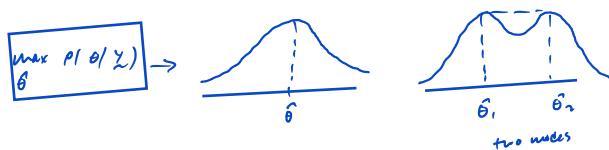
$$\min_{\hat{\theta}} E[L_1(\hat{\theta}, \theta) | \mathbf{y}] = \min_{\hat{\theta}} \int |\hat{\theta} - \theta| p(\theta | \mathbf{y}) d\theta$$

$$\Rightarrow \hat{\theta} = \int_{-\infty}^{\hat{\theta}} p(\theta | \mathbf{y}) d\theta = 0.5$$

*such that*

or written as  $p(\theta \leq \hat{\theta} | \mathbf{y}) = 0.5$

- Posterior mode:



→ Always good to check shape of posterior if we have derived it already

- The coin example

$$\rightarrow \text{Setup} \rightarrow Y \sim \text{Binomial}(100, \theta) \quad \theta \sim \text{Uniform}(0,1) \quad \{\text{Beta}(1,1)\}$$

$$\rightarrow \text{Data dist} \rightarrow p(Y|\theta) = \binom{100}{Y} \theta^Y (1-\theta)^{100-Y}$$

$$\rightarrow \text{prior dist} \rightarrow \pi(\theta) = 1$$

$$\rightarrow \text{posterior dist} \rightarrow p(\theta|y) \propto p(y|\theta) \pi(\theta)$$

$$\downarrow \propto \theta^Y (1-\theta)^{100-Y}$$

$$= \theta^{(y+1)-1} (1-\theta)^{101-y} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{shown previously}$$

$$\theta|y \sim \text{Beta}(y+1, 101-y)$$

$$\rightarrow \text{Posterior mean} \rightarrow E(\theta|y) = \frac{y+1}{(y+1) + (101-y)} = \frac{y+1}{102}$$

Posterior median  $\rightarrow$

posterior mode  $\rightarrow$  use R  $\rightarrow$  mode = qbeta(0.5, y+1, 101-y)

$\rightarrow$  mode = optimize(f = x\_beta, interval = c(0,1), maximum = TRUE)

$\hookrightarrow$  function(x) {dbeta(x, a, b)}

output  $\Rightarrow \$\text{maximum} = \text{mode}$



example  
Beta(5,1)

## 2.2 Interval Estimation

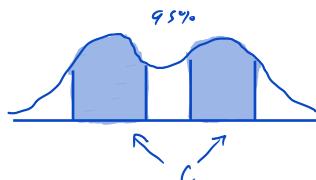
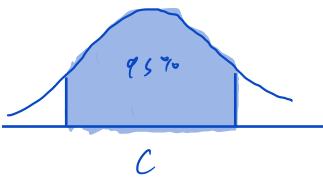
- Definition: A  $100(1-\alpha)\%$  credible set (interval)  $C$  is a subset of the parameter space  $\Theta$  such that

$$\alpha = 0.05 \Rightarrow 95\%$$

$$\int_C p(\theta|y) d\theta = 1 - \alpha$$

$$\int_{\Theta} p(\theta|y) d\theta = 1$$

If the parameter space is discrete, we replace the integral with sum.



(doesn't have to be a continuous range)

$\rightarrow$  Now the 95% represents an actual probability (which is different than 95% confidence in frequentist)

$\Rightarrow$  e.g. given the data we have, there is a 95% chance that our parameter is in the set

$\rightarrow$  said another way  $\rightarrow$   $P(L \leq \theta \leq U|y) = 1 - \alpha$

$\rightarrow$  In Bayesian, we no longer use confidence intervals b/c the interpretation of the confidence level  $100(1-\alpha)\%$  doesn't match what we want  
(i.e. capture rate of parameter w/ repeated sampling)  
 $\Rightarrow$  So instead we use a "credible set"

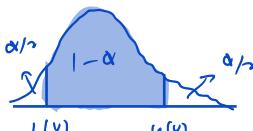
$\rightarrow$  If we simulate our data  $\Rightarrow$  then the interval is an approximation (Monte Carlo)

### • Equal tails interval

$\rightarrow$  Bounds  $L(y)$  &  $U(y)$  are just quantiles, so can easily find them if have posterior dist

$\rightarrow$  just use R or solve intervals

$$\int_{L(y)}^{L(y)} p(\theta|y) d\theta = \frac{\alpha}{2} \quad \left. \begin{array}{l} \\ \end{array} \right\} \Rightarrow \int_{U(y)}^{U(y)} p(\theta|y) d\theta = 1 - \alpha$$



Bounds can shift, just meet middle  $(1-\alpha)/2$

$\Rightarrow$  But there is an optimal credible set

- Based on quantiles of the posterior distribution.

- Not always the optimal, it could be wider if the posterior distribution is extremely skewed.

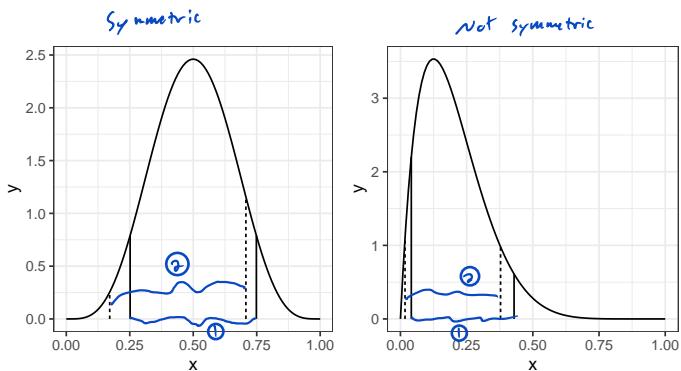
$\hookrightarrow$  Drawbacks of equal tails interval  
 $\rightarrow$  obviously want the narrowest interval

$\rightarrow$  then non-equal tails intervals

$\rightarrow$  Advantages

$\rightarrow$  If posterior dist is symmetric & unimodal, then equal tails interval is optimal  
(i.e. better than any other type of interval)

$\rightarrow$  Easy to find/calculat



① is equal tail

② not equal tail

equivalent to equal tails when posterior is unimodal & symmetric,  
but if don't have  $g^{-1}(P)$  function  
in R, hard to compute  $\Rightarrow$  Equal tails preferred

- **Highest posterior density (HPD):** A  $100(1 - \alpha)\%$  HPD region for  $\theta$  is a subset  $\mathcal{C} \in \Theta$  defined by  $\mathcal{C} = \{\theta : p(\theta|y) \geq k\}$ , where  $k$  is the largest number such that

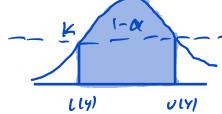
$$\int_{\theta : p(\theta|y) \geq k} p(\theta|y) d\theta = 1 - \alpha$$

D: advantages { Hard to compute if we don't have the inverse CDF for the posterior distribution  $\rightarrow$  would have to write code to search  $\rightarrow$  e.g. find  $\theta$ :  $p(\theta|y) \geq k$   
Not guaranteed to be an interval (could be a set)

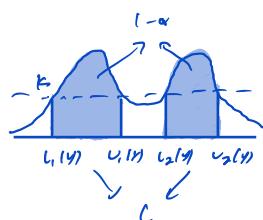
Advantages  $\rightarrow$  good + useful for skewed distributions

$\rightarrow$  probability statement  $\rightarrow p(L(y) \leq \theta \leq U(y)) = 1 - \alpha \Leftrightarrow \theta \in [L(y), U(y)] \Leftrightarrow p(\theta|y) \geq k$   
all density values outside of  $[L(y), U(y)]$  are smaller than  $k$

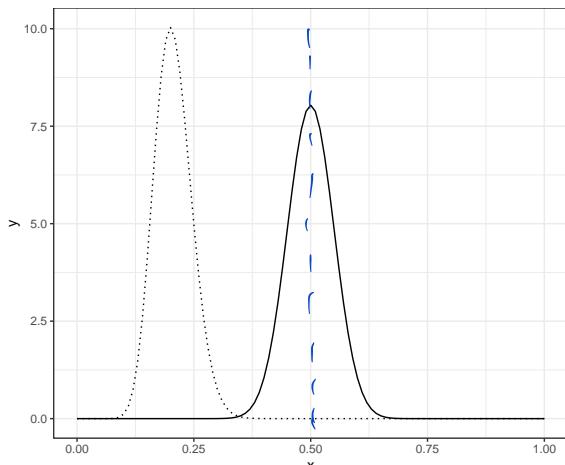
↙ simplest case (unimodal)



also can be for more complex sets



Example: In the coin example, consider  $y = 50$  and  $y = 30$



$\rightarrow \theta|y \sim \text{Beta}(y+1, 101-y) \quad \theta \in (0,1)$

$\rightarrow$  If  $y=50 \rightarrow \theta|y \sim \text{Beta}(51, 51)$   
 $\Rightarrow$  symmetric  $\Rightarrow$  equal tails interval  
is optimal

$\rightarrow$  If  $y=30 \rightarrow \theta|y \sim \text{Beta}(31, 71)$

$\Rightarrow$  not symmetric  $\Rightarrow$  need HPD interval  
to be optimal

$\rightarrow$  It is good practice to (roughly) graph the posterior density so that we know which type of interval will be narrowest

## Comparison of Point estimates

- point estimates for frequentist & Bayesian will be very close if have a uniform prior
- Posterior mean → always is a compromise (i.e. weighted average) between the prior distribution + the likelihood function ( $E(\theta)$ )  
 ↳ e.g. algebraically can rewrite  

$$E(\theta|y) = \underbrace{k\bar{\theta}_{MLE}}_{\text{weight data}} + \underbrace{(1-k)\bar{\theta}}_{\text{weight prior}}$$

{ Applet in Examples later }

$$(1 \leq \theta) \quad \text{↳ for mean} \Rightarrow \bar{\theta}_{MLE} = \bar{y} \text{ sample mean}$$

$$\rightarrow k \text{ usually some fraction of sample size + hyperparameters}$$

$$\rightarrow \text{e.g. example } y|0 \sim \text{Binomial}(n, \theta) \quad \left. \begin{array}{l} \text{Post. mean} \\ \theta \sim \text{Beta}(\alpha, \beta) \end{array} \right\} \Rightarrow E(\theta|y) = \frac{y+\alpha}{n+\alpha+\beta} = \frac{y+1}{102} \quad \begin{matrix} \text{if } \alpha=\beta=1 \\ \text{known } n \text{ units} \end{matrix}$$

$$\downarrow = \frac{(\frac{n}{n+\alpha+\beta})y}{\text{weight data}} + \frac{(\frac{\alpha+\beta}{n+\alpha+\beta})\alpha}{\text{weight prior}}$$

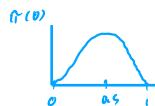
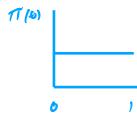
$$\rightarrow \text{fixing weight in data increases as } n \text{ increases}$$

$$\Rightarrow \text{more data} \Rightarrow \text{less influence of prior}$$

## Selection of prior

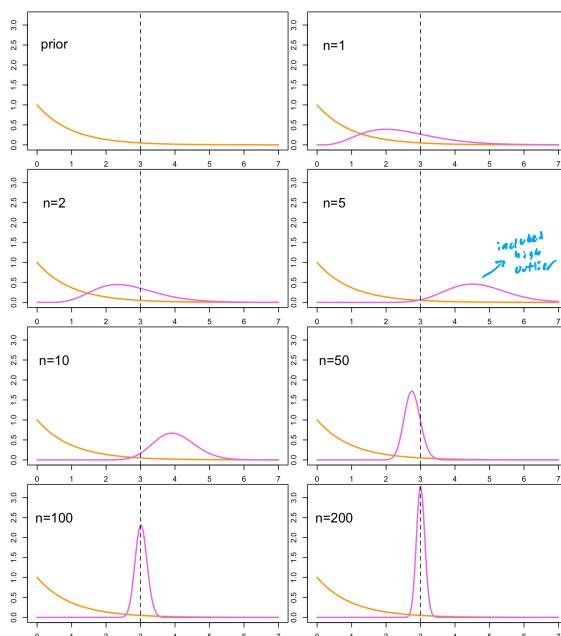
- Prior is a distribution that should match the parameter space
- If the parameter space is discrete, then the prior should be a discrete dist
- Also the ranges of  $\theta$  &  $\pi(\theta)$  should match

→ Example priors → uninformative uniform ( $0, 1$ ) vs  $\text{Beta}(\alpha=2, \beta=2)$   
 $\Rightarrow$  we suspect a higher prob of  $\approx 0.5$  based on previous knowledge



## Effect of choosing the wrong prior

- Choosing a really bad prior can impact the final result a lot
- with not a lot of data, a strong prior will dominate (i.e. heavy weight) the result  
 (↗ terms of the shape + location of the posterior)
- But if have lots of data, the data will dominate the resulting posterior instead of the prior



⇒ If have several candidates for choice of prior, it is better to just use a uninformative prior rather than guessing/trying that our choice of a good prior is good

$$\begin{aligned} \text{In frequentist} \rightarrow & \quad Y_i = X_i^T \beta + \epsilon_i \Rightarrow \text{estimate } \hat{\beta} \\ (\text{regression}) \quad & \rightarrow \text{observe } \tilde{x}_i^T + \text{predict } \tilde{y}_i \rightarrow E[\tilde{y}_i] = \tilde{x}_i^T \hat{\beta} \\ & \quad \tilde{y}_i = E[\tilde{y}_i] + \tilde{\epsilon}_i \end{aligned}$$

## 2.3 Prediction

- Suppose  $\tilde{y}$  is an estimate of the future observation.

$$\begin{aligned} Y &= (Y_1, \dots, Y_n) \\ \text{Data distribution: } p(\mathbf{y}|\theta) & \quad \text{prior } \pi(\theta) \quad \left. \begin{array}{l} \text{posterior} \\ \text{Bayes' rule: based on the data we have} \end{array} \right\} \text{prob(y)} \\ \text{Prior: } \pi(\theta) & \\ \text{Future observation: } p(\tilde{y}|\mathbf{y}) & \quad \text{so far prediction, we are wanting to predict a new obs} \\ & \quad \text{given the data we currently have} \\ & \quad \downarrow \text{vector future value} \end{aligned}$$

- Definition:** The posterior predictive distribution of the future observation is

start

→ prior predictive distribution

$$p(\tilde{y}) = \int_{\Theta} p(\tilde{y}|\theta) \pi(\theta) d\theta$$

→ This is predicting a future  $y$  w/o data  
so only using knowledge from the prior

$$p(\tilde{y}|\mathbf{y}) = \int_{\Theta} p(\tilde{y}, \theta|\mathbf{y}) d\theta = \int_{\Theta} p(\tilde{y}|\theta) \pi(\theta) d\theta$$

future obs given dataset

data dist posterior

→ in Bayes, everything is written given  $\mathbf{y}$

→ this is the joint dist of future obs + our knowledge w/o  $\mathbf{y}$

Just sample more values if want a vector of future obs

integrate out  $\theta$  to get a marginal of  $\tilde{y}$  (the future obs), but still based on the data we have

need to know both of these for prediction

rewrite this as conditioned to use the info we have

→ predictive credible interval

$$\int_{U(Y)}^{U(Y)} p(\tilde{y}|\mathbf{y}) d\tilde{y} = 1 - \alpha$$

uses current data

- Example:** The coin example: Suppose a coin was tossed 100 times and 50 heads were obtained. What is the chance if another head is obtained for another toss?

$$\begin{aligned} \rightarrow \text{Posterior} &\rightarrow \theta|Y \sim \text{Beta}(50+1, 50-50) \\ &\rightarrow \text{if } Y=50 \rightarrow \sim \text{Beta}(51, 51) \rightarrow \text{semimatter fctn} = \frac{\text{Beta}}{\Gamma(51)\Gamma(51)} x^{50} (1-x)^{50-50} , 0 < x < 1 \end{aligned}$$

→ Posterior predictive distribution

$$\rightarrow p(\tilde{y}|Y) = ?$$

$$\begin{aligned} &= \int_0^1 p(\tilde{y}|\theta) \pi(\theta|Y) d\theta \\ &\quad \xrightarrow{\text{data dist from model (n=1 for this example)}} \\ &= \int_0^1 \theta^{\tilde{y}} (1-\theta)^{1-\tilde{y}} \frac{\Gamma(102)}{\Gamma(51)\Gamma(51)} \theta^{50} (1-\theta)^{50} d\theta \\ &= \frac{\Gamma(102)}{\Gamma(51)\Gamma(51)} \int_0^1 \theta^{(51+\tilde{y})-1} (1-\theta)^{(52-\tilde{y})-1} d\theta \\ &= \left[ \frac{\Gamma(51+\tilde{y}) \Gamma(52-\tilde{y})}{\Gamma(103)} \right] \end{aligned}$$

$$\rightarrow \tilde{y}|\theta \sim \text{Bernoulli} \rightarrow \text{if } Y=0 \rightarrow p(\tilde{y}=0|Y) = \frac{\Gamma(102)}{\Gamma(51)\Gamma(51)} \frac{\Gamma(51) \Gamma(52)}{\Gamma(103)} = \frac{51}{102} = 0.5$$

$$\Rightarrow \text{if } Y=1 \Rightarrow p(\tilde{y}=1|Y) = 1 - p(\tilde{y}=0|Y)$$

$\downarrow = 0.5$

$\Rightarrow$  If  $Y=50 \Rightarrow$  our posterior for  $\theta$  suggests equal prob for heads & tails  
 $\rightarrow$  so it makes sense that prob of tails for  $\tilde{y}=0.5$

## Posterior predictive distribution

→ In the previous example, the posterior dist had a nice, closed form solution  
 $\Rightarrow$  could easily do prediction

→ If it wasn't nice  $\Rightarrow$  can use simulation to obtain the posterior predictive distribution

$\left. \begin{array}{l} \\ \end{array} \right\}$  using simulation

→ Suppose  $p(\theta|y)$  is the posterior & we have/can get samples from  $p(\theta|y)$

$$\rightarrow p(y|\theta) = \int_0^\infty p(y|\theta) p(\theta|y) d\theta$$

→ 2 steps to generate sample of  $\tilde{Y}$

for  $m=1, 2, \dots, M \rightarrow$  th with sample point generated

$$1) \text{ Simulate } \theta^{(m)} \sim p(\theta|y)$$

2) Simulate  $\tilde{y}^{(m)} \sim p(y|\theta^{(m)}) \rightarrow$  data dist, just getting a new  $y$  value,  $\tilde{y}$   
 $\rightarrow$  now  $y$  conditioned on previous draw at  $\theta$

$\leftarrow$  repeat

$\rightarrow$  then we have  $\tilde{y}^{(1)}, \tilde{y}^{(2)}, \dots, \tilde{y}^{(M)}$

(dependent on previous draw)  
 $\Rightarrow$  Separately updating the posterior

→ Gain example  $\rightarrow N=10,000$  ( $y=50$ )

$$\rightarrow \text{step 1) } \theta^{(1)} \sim \text{Beta}(51, 51)$$

$$\rightarrow \text{step 2) } \tilde{y}^{(1)} \sim \text{Bernoulli}(\theta^{(1)})$$

$$\Rightarrow \tilde{y}^{(1)}, \dots, \tilde{y}^{(10000)} \rightarrow$$
 this is called a "Monte Carlo" sample  $\Rightarrow$  simulating independent samples from the posterior (regular posterior or posterior predictive)

$\left. \begin{array}{l} \\ \end{array} \right\}$  → then we can approximate probabilities of interest using our simulated/generated sample

$\rightarrow$  e.g. what's the probability of having 2 heads if toss the coin 5 times??

$$\rightarrow \text{simulate } \tilde{y}^{(1)} \sim \text{Bernoulli}(\theta^{(1)}) \text{ 5 times}$$

$\leftarrow$  this is just getting a sample rather than an individual  $\tilde{y}$  at a time

A Markov Chain (with learning) is when have dependent samples  
 $\text{est. } \theta^{(m)} \sim p(\theta^{(m-1)}|y)$

$\left. \begin{array}{l} \\ \end{array} \right\} (\theta^{(1)} \parallel \theta^{(0)})$   
 $\rightarrow$  Monte Carlo just means applying discrete formulas to approximate continuous integrals

$$0.9) E[\theta|y] = \frac{1}{M} \sum \theta^{(m)} | y$$

$$P(\theta \leq \theta_0 | y) = \frac{1}{M} \sum I(\theta^{(m)} \leq \theta_0 | y)$$

## form of posterior predictive dist

→ A conjugate prior has nothing to do w/ the form of the posterior predictive distribution

→ Usually the form of the posterior predictive is not nice

$\Rightarrow$  which is why we often use simulation to obtain future samples