## Assignment Project Page Help Bayesian Methods Help

https://tutorcs.com

WeChat: CStutorcs
Lecture 9

### A Thought Experiment:

### Assignment Project Exam Help

- Event A := a randomly selected person is taller than 190cm. Then, <a href="https://tutofr@s-.com">https://tutofr@s-.com</a>
- Event B := a randomly selected person is female.
- Then, what are the odds that she is 190cm of tallered. Then, what are the odds that she is Pr(A|B) = ?

### Subjective Probability

## Assignment be wroject Exam Help

- When we learn that a female is selected, we may revise our belief about the person being taller than 190cm.

  The peally the person being taller than 190cm.
- Your belief before the additional information, Pr(A), is called the prior.
- Your (revised) belief after the information,  $\Pr(A|B)$ , is the posterior.

### **Bayes Theorem**



- $\begin{array}{c} \text{So, Pr}(A \cap B) / \text{Pr}(A) \text{ Pr}(B|A) \\ \text{Nttps.} / \text{Tutorcs.com} \\ \text{Pr}(A|B) = \frac{\text{Pr}(A) \text{ Pr}(B|A)}{\text{Pr}(B)} \end{array}$ 
  - ► TWI (Caversion) BaceSthedrem (CS
- Bayes theorem shows how the prior is updated to the posterior when there is additional information.

### Bayes Theorem for PDF

Reprint Bayes theorem:

### Assignment Project Exam Help

- Let  $\pi(x)$  be the marginal PDF of X.
- Let f(y|x) be the interginal PDF of Y given X = x.
- By Bayes theorem, we write

WeChat: 
$$\underset{p(y)}{\text{cst}} = \underbrace{\text{utors}}_{p(y)}$$

 $\blacktriangleright$   $\pi(x)$  is the prior on X and  $\pi(x|y)$  is the posterior on X given Y = v

### Bayes Theorem for PDF

# Assignment Project Exam Help $\pi(\theta|y) = \frac{\pi(\theta)f(y|\theta)}{p(y)}.$

- Interporior that to resist is the parties on  $\theta$  given Y = y.
- ▶ The parameter  $(\theta)$  may not be random. But, it is uncertain.
- Our subjective belief about  $\theta$  is represented as  $\pi(\theta)$  (prior belief, before data) and  $\pi(\theta|y)$  (posterior belief, after data).

- A coin is tossed and Y = 1 if the outcome is Head and otherwise Y = 0.
- Assignment,  $P_{\theta}$  (o) set (y Exam Help
  - Suppose you believe every  $\theta \in (0,1)$  is equally likely. Then, the prior is https: whiteleft https: whiteleft https: whiteleft https: whiteleft https:
  - Then, the posterior is

$$We^{\text{Coh}}\bar{a}t^{(\theta \in (0,1))}_{-\pi(\theta)}t^{\theta'}_{t}(1-\bar{e}S)^{1-\gamma/\rho(y)}_{t(y|\theta)}$$

 $\blacktriangleright \text{ If } y=1,$ 

$$\pi(\theta|y=1) = \mathbb{1}(\theta \in (0,1)) \cdot \theta/p(1)$$
$$= 2\theta \cdot \mathbb{1}(\theta \in (0,1))$$

• What is  $\pi(\theta|y=0)$ ?



▶ To be more flexible, we consider the beta prior, i.e.,

### Assignment roject Exam Help

Since  $\theta \sim Beta(\alpha, \beta)$ , we know that

https://tutorcs.com

WeChat: 
$$constant constant co$$

- Suppose you believe that  $\theta = 0.25$  on average, i.e.,  $E[\theta] = 1/4$ , and your belief is quite strong, e.g.,  $V[\theta] = 0.01$ . Then,  $(\alpha, \beta) \approx (4, 13)$ .
- ▶ As before, let's compute  $\pi(\theta|y=1)$ .

▶ To solve this problem, we need to learn a useful trick:

$$\pi(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot \theta^{\alpha - 1} (1 - \theta)^{\beta - 1} \cdot \mathbb{1}(\theta \in (0, 1))$$

# $\pi(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot \theta^{\alpha - 1} (1 - \theta)^{\beta - 1} \cdot \mathbb{1}(\theta \in (0, 1))$ $Assignment^{\alpha} - \text{Project} \text{ Exam Help}$

where c is the normalising constant (i.e. it does not depend on  $\theta$ ).

https://tutorcs.com

$$\pi(\theta) = \mathbf{c} \cdot \theta^{\alpha - 1} (1 - \theta)^{\beta - 1} \cdot \mathbb{1}(\theta \in (0, 1)),$$

### We also know that $r = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)} = \frac{\Gamma(\alpha)}{\Gamma(\alpha)} = \frac{$

So, even if we simply write

$$\pi(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \cdot \mathbb{1}(\theta \in (0,1)),$$

we immediately know

$$\theta \sim Beta(\alpha, \beta)$$
.



Bayes theorem:

$$Assignment(1Project(E)xeam\theta)Help$$

$$\begin{array}{l} = \theta^{(\alpha+y)-1}(1-\theta)^{(\beta+1-y)-1} \cdot \mathbb{1}(\theta \in (0,1)) \\ \text{https://tutoffc.s.c.o.m} \end{array}$$

where  $\tilde{\alpha} = \alpha + y$  and  $\tilde{\beta} = \beta + 1 - y$ .

► We Ciohatity i cstutores

$$\pi(\theta|y) \propto \theta^{\tilde{\alpha}-1} (1-\theta)^{\tilde{\beta}-1} \cdot \mathbb{1}(\theta \in (0,1))$$

which means

$$heta| extbf{y} \sim extbf{Beta}( ilde{lpha}, ilde{eta})$$



Since the posterior is  $Beta(\tilde{\alpha}, \tilde{\beta})$ , the posterior mean and variance when y = 1, are

### Assignment Project Exam Help

$$V[\theta|y=1] = \frac{\tilde{\alpha}\tilde{\beta}}{\sqrt{\tilde{\rho}} + \tilde{\beta}/2(\tilde{\alpha} + \tilde{\beta} + 1)} = \frac{(\alpha + 1) \cdot \beta}{(\alpha + \beta + 1)^2(\alpha + \beta + 2)}$$
https://tubicrcs.com

Since  $(\alpha, \beta) = (4, 13)$ , we have

We Chat: 
$$_{1}^{[\theta]y=1}=0.29 > E[\theta]=0.25$$

- ▶ Since y = 1 observed, we believe  $\theta = \Pr(y = 1)$  should be larger.
- Since the data contradicts our prior, uncertainty around our belief on  $\theta$  increases. The variance increases.

### Example: Normal Model

▶ When  $X \sim \mathcal{N}(\mu, \sigma^2)$ , the PDF of X is given as

### Assignment Project, Exam Help

- Suppose  $y|\theta \sim \mathcal{N}(\theta, 1)$ . The PDF is https://tutorcs.com  $f(y|\theta) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y-\theta)^2\right)$
- ► The Confriggiven as SM(14100 TCS

$$\pi(\theta) = \frac{1}{\sqrt{2\pi 10^2}} \exp\left(-\frac{1}{2 \cdot 10^2} \theta^2\right)$$



### Example: Normal Model

Posterior density

Assignment 
$$\Pr_{x_{\theta},y_{\theta}}^{\pi(\theta|y)} = \exp\left(-\frac{1}{2\cdot 10^2}\theta^2\right) \exp\left(-\frac{1}{2}(y-\theta)^2\right)$$

$$https://tutorcs.com/{2(100/101)}^2 \left( \frac{-\frac{1}{2} \left( y^2 - 2y\theta + \theta^2 + \theta^2 / 100 \right)}{2(100/101)} \right)^2 \right)$$

This implies that the conditional distribution of  $\theta$  given yet from the with near (100/101) Cand variance 100/101, i.e.,

$$\theta | y \sim \mathcal{N}\left(\frac{100y}{101}, \frac{100}{101}\right)$$

► Hence,  $E[\theta|Y=y] = (\frac{100}{101}) y$  and  $V[\theta|Y=y] = \frac{100}{101}$ .

### Conjugate Priors

Previously, we considered an example where we observed Y given the unknown mean  $\theta$ , and  $\theta$  had a normal prior distribution. The posterior distribution turned out to be

### Assignment Project Exam Help

which has the property that the posterior distribution is of the same family/as the prior.

Conjugacy is specific to prior and moder combination. For

- example,
   Normal prior is conjugate for Normal model for mean
- Beta prior is conjugate for Bernouilli model
- But, for instance, Beta prior is not conjugate for Normal model.
- ► Conjugate priors are useful: the posterior is a standard distribution with known properties, e.g, posterior mean and posterior variance are easy to compute. (as above)

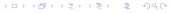
### Conjugate Priors: an example of non-conjugate prior

 Consider Beta prior for the normal model above. The posterior is

# Assignment Project Exam Help $\underset{\alpha \in Beta(\alpha,\beta)}{\text{Assignment Project Exam Help}}$

- Let  $\pi^*(\theta|y)$  denote the RHS, i.e.,  $\pi(\theta|y) = c \cdot \pi^*(\theta|y)$  for some c > 0.
- ▶ There is postandard distribution proportional to  $\pi^*(\theta|y)$ .
- In order to obtain the posterior, we need to find c such that

$$\int_0^1 \pi(\theta|y)d\theta = \int_0^1 c \cdot \pi^*(\theta|y)d\theta = 1$$



### Conjugate Priors: an example of non-conjugate prior

► Then, the posterior is

### Assignment Project Project Assig

Moreover, the posterior mean and posterior variance are

$$\underset{E[\theta|y]}{\text{https:}} / \underset{\pi(\theta|y)d\theta}{\text{tutor}} = \underset{0}{\text{second}} \underset{\int_{0}^{1} \pi^{*}(\tilde{\theta}|y)d\tilde{\theta}}{\text{second}}$$

$$\begin{array}{c} \mathbf{W}[\theta]\mathbf{y}] = \mathbf{E}[\theta - \mathbf{E}[\theta|\mathbf{y}])^{2}|\mathbf{y}] = \int_{0}^{1} (\theta - \mathbf{E}[\theta|\mathbf{y}])^{2} \left[ \frac{\pi^{*}(\theta|\mathbf{y})}{\int_{0}^{1} \pi^{*}(\tilde{\theta}|\mathbf{y})d\tilde{\theta}} \right] d\theta \\ \mathbf{E}[\theta]\mathbf{y}] = \mathbf{E}[\theta]\mathbf{y}]^{2}|\mathbf{y}] = \mathbf{E}[\theta]\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}] = \mathbf{E}[\theta]\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}] = \mathbf{E}[\theta]\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}] = \mathbf{E}[\theta]\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}]^{2}|\mathbf{y}^{2}|\mathbf{y}]^{2}|\mathbf{y}^{2}|\mathbf{y}^{2}|\mathbf{y}^{2}|\mathbf{y}|^{2}|\mathbf{y}^{2}|\mathbf{y}^{2}|\mathbf{y}^{2}|\mathbf{y}|^{2}|\mathbf{y}^{2}|\mathbf{y}^{2}|\mathbf{y}|^{2}|\mathbf{y}^{2}|\mathbf{$$

- ► The integrals need to be numerically evaluated. These days, computers are fast and algorithms are good. We will learn a few algorithms later.
- ► Still prior conjugacy is important: there are still limitations on what can be handled by computational algorithms.

### Prior vs Sample

Consider the model and prior as follows

## Assignment $\Pr{\mathbf{Project}^{1}}\mathbf{Exam}$ Help

where  $\tau$  and h are precisions (precision = inverse

variance). Assume  $\tau$  is known.

• Martia Provior Idaha Check S conjugate Turns out the posterior is

We Chather 
$$stuton$$
  $stuton$ 

- The posterior mean is a weighted average of the prior information and the sample information
- The posterior variance is the sum of prior and sample precisions.

### Prior vs Sample, and Improper prior

Reprint the posterior distribution;

Assignment 
$$Project Fxam-Help$$

where h is prior precision.

- If you are really suitable the value of the fore observing y. Then, your prior precision is high, i.e., h↑. The weight on the prior mean is large and the weight on the observation is small. So, E[f] y to T E[f]
   On the other hand, if you are uncertain about θ, h would be
- On the other hand, if you are uncertain about  $\theta$ , h would be small, and the weight on the observation would be relatively large.

### Prior vs Sample, and Improper prior

Reprint the posterior distribution;

- where h is prior precision.

   The prior larts by the Sat Couring to weight on the prior mean.
- At the limit,  $\pi(\theta) \propto 1 \Rightarrow$  often called uniform/flat prior, or ever reginformative prostutores
- ► The flat prior is not a density because it does not integrate to 1. So, it is called an improper prior. But, the posterior is still a density.

#### More than one observation

- Let  $Y_1, Y_2, \ldots, Y_n$  be a random sample from  $f(\cdot | \theta)$  and we have  $\pi(\theta)$ .
- ▶ When only  $Y_1 = y_1$  is observed, the posterior is given as

### Assignment Project Exam Help

Nhen  $Y_2 = y_2$  is also observed,  $\pi(\theta|y_1)$  is the prior to be updated

https://tutores.com/
$$|\theta| f(y_2|\theta)$$

- ▶ When  $Y_3 = y_3$  is observed, similarly  $\pi W_1 + C_2 hat \theta y_1 C_2 f t t t + C_3 f (y_2 | \theta) f(y_2 | \theta) f(y_3 | \theta)$ new prior
- So, the posterior is proportional to the prior times the likelihood

$$\pi(\theta|y_1,\ldots,y_n)\propto \pi(\theta)\cdot\prod_{i=1}^n f(y_i|\theta)$$

#### More than one observation

► Suppose we have a random sample from the normal model

Assignment and roje the normal prior Help 
$$\theta \sim \mathcal{N}(d, h^{-1})$$
.

Then, applying the conjugacy recursively, we have the conjugacy recursively.

$$\mathbf{\overset{\theta|y_1,..,y_n \sim \mathcal{N}}{VeChat.}} \underbrace{\begin{bmatrix} \frac{h}{h+n_{\tau}} \end{bmatrix} d + \begin{bmatrix} \frac{n\tau}{h+n_{\tau}} \end{bmatrix} \begin{bmatrix} \frac{1}{n} \sum_{i=1}^{n} y_i \end{bmatrix}, (h+n\tau)^{-1}$$

- ▶ For all  $(d, h) \in \mathbb{R} \times \mathbb{R}_+$ , i.e., for all proper priors, as  $n \to \infty$ ,
  - the posterior mean gets close to the sample mean
  - $-\frac{1}{n}\sum_{i=1}^{n}y_{i}\stackrel{p}{\rightarrow}\theta_{0}$ , the true population mean (LLN)
  - the posterior precision explodes
  - So, the posterior will be degenerate at  $\theta_0$  (Consistency)

### Optimal Decision under Uncertainty

Let  $L(\theta, a)$  be the loss when you choose an action  $a \in A$ under  $\theta \in \Theta$ .

### Assignmentu-Projectin Exam Help

- $\triangleright$  But,  $\theta$  is unknown.
- You have prior  $\pi(\theta)$  and observe  $y_1, \ldots, y_n \stackrel{iid}{\sim} f(y|\theta)$ .

   Under axioms of Salvage (1954), it is obtained to choose

$$We Charmin \underset{a \in \mathcal{A}}{\text{long min}} \int_{\Theta} L(\theta, a) \pi(\theta | y_1, \dots, y_n) d\theta$$

which is called the Bayes action. (Subjective Expected Utility Theory)

#### Estimation as a Decision Problem

The econometrician would suffer the error square loss

## Assignment $\underset{\leftarrow}{\text{Project}}$

For this particular problem, it is optimal to choose

https://eutrores.com 
$$E[\theta|y_1,...,y_n]$$

- If  $L(\theta, a) = |\theta a|$ , it is optimal to report the posterior method: CStutorcs
- ➤ As seen here, the optimality depends on the loss function (the way you feel about the error you make)
- In the example of the normal model, the Bayesian estimate is the posterior mean under either of the loss functions above.

### Summary of Uncertainty

► The posterior variance or the posterior standard deviation can be reported as a summary of uncertainty around the estimate.

Assignment of crops of the last of the las

- ► There could be many such intervals (not unique).
- The narrowest one is generally preferred and we use here.
   Interpretation is natural. (compare with 95% confidence
- Interpretation is natural. (compare with 95% confidence interval!)
- ► For the normal model above, the 95% CI is

$$\hat{\theta}_B \pm 1.96 \sqrt{V(\theta|y_1,\ldots,y_n)}$$



Consider a regression model

where we assume that  $e_i|x_i, \beta \sim \mathcal{N}(0, 1/\tau)$ , and  $\beta$  and  $x_i$  are  $(K \times 1)$  vectors. Since  $e_i = y_i - x_i'\beta$ , we have https://tutorcs.com

 $nttps://tutores.eom, \\ y_i = x_i\beta | x_i, \beta \sim \mathcal{N}(0, 1/\tau),$ 

the likelihood is 
$$\prod_{i=1}^{N} \underbrace{\frac{\tau}{2\pi} \exp\left(-\frac{\tau}{2}e_i^2\right)}_{\text{PDF of } \mathcal{N}(0,1/\tau)} \propto \tau^{\frac{n}{2}} \exp\left(-\frac{\tau}{2}\sum_{i=1}^{n}(y_i - x_i'\beta)^2\right)$$

$$= au^{rac{n}{2}}\exp\left(-rac{ au}{2}(y-Xeta)'(y-Xeta)
ight)$$

## Assignment Project, Fx am Help

$$eta | au \sim \mathcal{N}\left(b_0, ( au \Sigma_0)^{-1}\right)$$

hart p.Snd /otal toll Csalas Quin K vector, and  $\Sigma_0$  is a  $K \times K$  positive definite matrix.

The prior is

$$\frac{1}{\pi(\tau,\beta)} e^{\frac{\tau}{2}} \underbrace{hat \frac{\lambda_0}{2} cstud}_{\propto \pi(\tau)} \underbrace{crcs}_{\frac{\tau}{2}(\beta-b_0)'\Sigma_0(\beta-b_0)} \underbrace{\frac{\tau}{2}(\beta-b_0)'\Sigma_0(\beta-b_0)}_{\propto \pi(\beta|\tau)}$$

Assignment  $\Pr_{\pi(\tau,\beta|y,X)} \sim \tau^{\frac{\alpha}{2}-1} \exp\left(-\frac{1}{2}\tau\right) \cdot \mathbb{I}(\tau>0)$  Help

$$https://tut_{\cdot \tau^{\frac{k}{2}} \exp\left(-\frac{\tau}{2}(\beta-b_0)'\Sigma_0(\beta-b_0)\right)}^{\frac{k}{2}\exp\left(-\frac{\tau}{2}(\beta-b_0)'\Sigma_0(\beta-b_0)\right)}$$

Toderive the posterior, it is useful to know  $(y - X\beta)'(y - X\beta) = (y - X\hat{\beta}_{LS})'(y - X\beta_{LS}) + (\beta - \hat{\beta}_{LS})'(X'X)(\beta - \hat{\beta}_{LS})$  where  $\hat{\beta}_{LS}$  is the OLS estimate, i.e.,  $\hat{\beta}_{LS} = (X'X)^{-1}X'y$ 

Some additional algebra shows

## 

where 
$$\sum_{n} = \sum_{n} \frac{1}{X^{n}} \frac{1}{X^{$$

- Note that the posterior mean of  $\beta$  is the weighted average of the OLS estimate  $\hat{\beta}_{LS}$  and the prior  $b_0$ .
- ▶ Moreover, conditional on  $\tau$ , the posterior precision of  $\beta$  is the sum of the OLS precision and the prior precision.

▶ If the econometrician has the mean squared error, i.e.,

$$L(\theta, a) = (\theta - a)'(\theta - a), \text{ for } a \in \Theta$$

### Assignment, Project be Exam by estimate is the posterior mean.

Each  $\beta_j$  for j = 1, ..., K is normally distributed under the descriptoristic form.

$$\beta_j | \tau, y, X \sim \mathcal{N}\left(b_{n,j}, V(\beta_j | \tau, y, X)\right)$$

- Where  $V(\beta)$  is the (j,j) element of  $(\tau \Sigma_n)^{-1}$ . • We often summarise the uncertainty about  $\beta_i$  by the
- We often summarise the uncertainty about  $\beta_j$  by the posterior standard deviation  $\sqrt{V(\beta_j|\tau,y,X)}$  or the (narrowest) 95% credible interval

$$b_{n,j} \pm 1.96 \sqrt{V(\beta_j | \tau, y, X)}$$
.



▶ What if the prior is not conjugate? Then there is no analytic expression for the posterior...

Assignment of the parameter,  $f(\theta)$  is the parameter,  $f(\theta)$  the posterior is

$$\pi(\theta|\mathbf{y}) \propto \pi(\theta) f(\mathbf{y}|\theta)$$

**Metaps**(y//tutstarbleam( $n\theta$ ) is the density of y given  $\theta$ .

We are often interested in the posterior moment;

WeChat: 
$$cstutorcs$$

where  $h(\theta)$  is a measurable function of  $\theta$ , e.g.,  $h(\theta) = \theta$  gives the posterior mean,  $h(\theta) = (\theta - E[\theta|y])^2$  the posterior variance, and so on.

Suppose we can draw  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(S)}$  from  $\pi(\theta|y)$ . As suppose we can draw  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(S)}$  from  $\pi(\theta|y)$ .

## https://tsspaces.com

- So, we may evaluate the integrals using a simulated sample from the posterior (Monte Carlo).
- Fine Hayes an tomput stiff, leso dially a Markov chain Monte Carlo approach is widely used.
- We outline the Metropolis-Hastings algorithm.

Let  $q(\tilde{\theta}|\theta)$  be a conditional density, called the proposal density.

### Assignment Project-Human Help

$$\underset{\text{https://tutorcs.com}}{\text{https://tutorcs.com}} \delta \sim \underset{\text{https://tutorcs.com}}{\text{diag}} s.com$$

and let  $\theta^{(s)} := \theta^{(s-1)}$  with probability 1 - Q.

► The normalising constants are cancelled out and the posterior ratio is computable.

- Note that the sequence  $\theta^{(1)},\dots,\theta^{(S)}$  is a Markov chain, Assigned for ly extheoretical property of the previous pr
  - Theoretically, the chain  $\theta^{(1)},\ldots,\theta^{(S)}$  has a property good enough to approximate the posterior, i.e., regardless of the initial post  $\theta^{(r)}$  tutores.com

Wechat: 
$$Cstutores$$

► The Metropolis Hastings algorithm is an example of a Markov chain Monte Carlo (MCMC) method.

In practice, the quality of an MCMC outcome can heavily depend on the initial point  $\theta^{(0)}$ , the proposal density  $\theta^{(1)}$ . The proposal density  $\theta^{(1)}$  is a function of the model  $\theta^{(1)}$ .

- For example, if  $\theta^{(0)}$  is far away from the posterior distribution, it may take a long time for the chain to reach to the suppost of the best of the community of the comm
- So, we exclude some early draws that seem to be a 'searching' phase rather than draws from the posterior, i.e., we discard the burn-in draws
- If the proposal density is similar to the posterior density, the algorithm works well. But, in practice, it is hard to know the shape of the posterior.

The normal distribution is often used. In particular,  $q(\tilde{\theta}|\theta)$  is the density of  $\mathcal{N}(\theta,\Omega)$ . Then,

# Assignment Project Exam Help $q(\tilde{\theta}|\theta) \propto \exp\left[-\frac{1}{2}\left(\frac{\theta-\theta}{\sigma_{\theta}}\right)\right].$

► the since  $q(\theta|\theta)$  =  $q(\theta|\theta)$ , the acceptance rate simplifies to

### WeChat: "cstutores

- ► The algorithm is called the Gaussian Metropolis-Hastings algorithm.
- Often t distribution is used.

- Metropolis-Hastings algorithms require the researcher to ASS1 to propose the backstand B
  - If  $\sigma_{\theta}$  is too big, Q would be small (as  $\tilde{\theta}$  would usually be far from  $\theta^{(s-1)}$ ) and the algorithm may not often update. That is, it is too be a long time.

  - ► Either case, the chain does not effectively explore the posterior.

It is important to choose  $\sigma_{\theta}$  that enables the chain to ASS1 explore the quarter of energy but it is not always easily proceed.

- So, a number of adaptive methods have been proposed. For example, Haario, Saksman, and Tamminen (2001).
- hard to tune the proposal density. For the Gaussian MH, a K dimensional covariance matrix has to be chosen.
- For a high-dimensional problem, Gibbs sampler is widely used. CStutorcs
- ldea is to recursively update one component or a small dimensional sub-vector of  $\theta$ .

Gibbs sampler: at each  $s^{th}$  iteration,

$$\begin{array}{c} Assign{ \begin{tabular}{l} Assign{ \begin{tabular}{l} Assign{ \begin{tabular}{l} Assign{ \begin{tabular}{l} Assign{ \begin{tabular}{l} Assign{tabular}{l} Assign{ \begin{tabular}{l} Assign{tabular}{l} Assign{tabular}$$

- At each-substep, the Gibbs sampler updates each  $\theta_k$  using either the (conditional) conjugacy or the MH algorithm (called H with Gibbs) STUTOTCS
- Each sub-step may update more than one component, i.e., it could update a block.

#### **Model Selection**

### Assiwe model Project Mexicanity Melp

- ▶ We may assume  $y \sim f(y|\theta)$  with unknown  $\theta \in \Theta$ ; Model 1.
- Or, we could assume  $\psi \sim q(\mathbf{y}|\psi)$  with unknown  $\psi \in \Psi$ ; Model DS.//TUTOTCS.COM
- ▶ We have two competing models (theories, assumptions),  $M := \{1, 2\}$
- Shows we have conditional drops  $\pi(\psi|m=2)$ .

#### **Model Selection**

▶ Conditional on m = 1, recall that the posterior of  $\theta$  is

$$Assignment for the marginal likelihood of model 
$$\begin{array}{l} \pi(\theta|y,m=1) = \frac{\pi(\theta|m=1)f(y|\theta)}{\sqrt{\pi(\theta|m=1)f(y|\theta)}} = \frac{\pi(\theta|m=1)f(y|\theta)}{h(y|m=1)} \\ \text{Assignment} for the marginal likelihood of model } \\ m=1. \end{array}$$$$

Similarly, conditional on m = 2,  $\frac{\text{https://tutorcs}}{\pi(\psi|y, m = 2)} = \frac{\pi(\theta|m = 2)g(y|\psi)}{\int \pi(\tilde{\psi}|m = 2)g(y|\tilde{\psi})d\tilde{\psi}} = \frac{\pi(\theta|m = 2)g(y|\psi)}{h(y|m = 2)}$ 

the marginal likelihood of model my 18 CStull Orcs

► Then, the (marginal) likelihood ratio of model 1 relative to model 2 is

$$B_{1,2} := \frac{h(y|m=1)}{h(y|m=2)}$$

which is called the **Bayes factor**.



#### Model Selection

Suppose we believe model  $m \in \mathcal{M}$  is true with prior probability  $\pi(m)$ . Then, the posterior probability of model m=1 is given as

## Assignment Project, he wan Help $\frac{1}{\pi(m=1)h(y|m=1) + \pi(m=2)h(y|m=2)}$

- The posterior odd ratio of model 1 to model 2 is  $\frac{\pi(m=1|y)}{\pi(m=2|y)} = \frac{\pi(m=1)h(y|m=1)}{\pi(m=2)h(y|m=2)}$
- That is the rosterior odd Salid Hurror od Sratio × Bayes factor
- ► If the posterior odd ratio > 1, we believe model 1 is more reasonable.
- ▶ BIC and AIC are rough approximations of this formal Bayesian model selection.

### **Bayesian Asymptotics**

Under mild conditions, posterior is consistent, i.e., the posterior asymptotically degenerates at the true

ASSIGNMENT (1941) OF COLUMN TO THE P

- The posterior and the sampling distribution of MLE are asymptotically equivalent under some regularity conditions on the model and the prior (Bernstein von Mises theorem)
- For large sample, therefore, the Bayesian analysis is repust to the choice of prior. Moreover, since MLE is evident, Bayes a biso efficient LOTCS
- Asymptotic robustness can be used as a reaction to the criticism that Bayesian analysis is essentially subjective (depends on the prior).