

MFH3041

Bayesian Statistics, philosophy and Practice

\ Classical view / • Probability is relative frequency (thing to be measured).

• Data, y , generated "at random" from a probability model parameterised by a fixed unknown, θ . $P(y|\theta)$ completely specify the problem

• Data $y = (y_1, \dots, y_n)$ (i.i.d) $P(y|\theta)$ = likelihood

• inference by estimating true θ , maybe with an interval.

\ Bayesian view / • probability is subjective (it measures uncertainty or degree of belief)

• y is again "generated" at random from $P(y|\theta)$ ← expresses our uncertainty about y is we know θ . ($P(y|\theta)$ follows) (likelihood)

• we are also uncertain about θ . So we can specify $P(\theta)$

• we can form the joint $P(y, \theta)$

$$\text{Bayes thm } p(\theta|y) = \frac{P(\theta)p(y|\theta)}{P(y)}$$

$p(\theta)$ ← prior distribution ($\pi(\theta)$)

$p(\theta|y)$ ← posterior distri ($\pi(\theta|y)$)

$p(y|\theta)$ ← likelihood

$P(y)$ ← marginal likelihood (prior predictive).

\ Ex / ~ 17. 88 women [40, 50] have breast cancer.

~ 90% of time a mammogram is "correct".

r.g. is θ , where $\theta=1$ is 1 woman has breast cancer & $\theta=0$ otherwise

$$\pi(\theta=1) = 0.01 \quad \therefore \pi(\theta=0) = 0.99$$

Data y ← mammogram test $y=1$

$$\pi(\theta=1 | y=1) = \frac{\pi(\theta=1) p(y=1 | \theta=1)}{p(y=1)}$$

$\theta=1 \Rightarrow y=1$ "correct"

$\theta=0 \Rightarrow y=0$ "correct"

$$p(y=1 | \theta=1) = p(y=0 | \theta=0) = 0.9$$

$$p(y=1) = p(y=1 | \theta=1) p(\theta=1) + p(y=1 | \theta=0) p(\theta=0) \quad (\text{LOT P } \theta=1, 0 \text{ form a partition.})$$

$$\pi(\theta=1 | y=1) = \frac{0.1 \times 0.9}{\left(\frac{9}{10} \times \frac{1}{100} + \frac{1}{10} \times \frac{99}{100} \right)} = \frac{9}{9+99} = \frac{1}{12}$$

$$\pi(\theta | y) \propto \pi(\theta) p(y | \theta)$$

$$p(y) = \int_{-\infty}^{\infty} p(y | \theta) \pi(\theta) d\theta \quad p(y) = \int_{-\infty}^{\infty} p(y, \theta) d\theta$$

$p(y)$ is const in θ

we know that densities (pd's) integrate to 1.

$$y_i | \theta \text{ iid } \text{Ber}(\theta) \quad 0 \leq \theta \leq 1 \quad \text{let } \theta \sim \text{Unif}(0, 1)$$

Suppose $i=1, \dots, n$, find 2 posterior distri & write down its density.

By Bayes Theorem one mark

$$\pi(\theta | y) \propto \pi(\theta) p(y | \theta) \quad \boxed{1}$$

$$\propto \pi(\theta) \prod_{i=1}^n \theta^{y_i} (1-\theta)^{1-y_i}$$

$$= 1 \times \theta^{\sum y_i} (1-\theta)^{n - \sum y_i}$$

$$S = \sum_{i=1}^n y_i \quad [\text{Eng}]$$

$$\pi(\theta | y) \propto \theta^{\bar{y}} (1-\theta)^{n(1-\bar{y})}$$

$$x \sim \text{Beta}(a, b) \quad \text{pdg} \quad \frac{P(x) P(y)}{P(x) P(y)} \cancel{\frac{\Gamma(a+b)}{\Gamma(a) \Gamma(b)}} x^{a-1} (1-x)^{b-1}$$

matching powers with a Beta Distr

$$[a-1] = n\bar{y} \therefore a = n\bar{y} + 1, \quad b = 1 + n(1-\bar{y})$$

$$\theta \sim \text{Beta}(n\bar{y} + 1, 1 + n(1-\bar{y}))$$

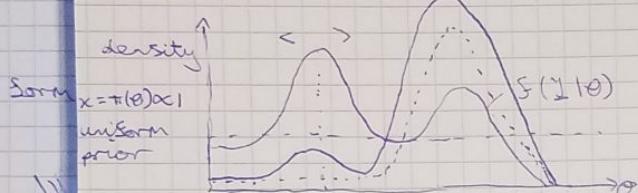
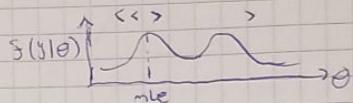
$$\text{density } \pi(\theta | y) = \frac{\Gamma(n+2)}{\Gamma(n\bar{y}+1) \Gamma(1+n(1-\bar{y}))} \theta^{n\bar{y}} (1-\theta)^{n(1-\bar{y})}$$

$$y_i | \theta \sim S(y_i; \theta)$$

y_1, \dots, y_n iid $S(y_i; \theta)$

likelihood $P(Y| \theta) = (S(Y| \theta))^n$

Bayes $\pi(\theta|Y) \propto \pi(\theta) S(Y| \theta)$



posterior inference (Bayesian inference) \equiv likelihood inference

Mervin Mulder is right

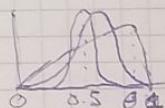
P is Paul is right

$$M_i \sim \text{Ber}(\theta_M) \quad P_i \sim \text{Ber}(\theta_P)$$

$$\text{is } \theta_M > \theta_P? \quad \theta_M = \hat{\theta}_M?$$

$$\text{likelihood } P(M| \theta_M) = \theta_M^m (1-\theta_M)^{n-m} \quad \hat{\theta}_M = 9/10$$

$$P(P| \theta_P) = \theta_P^p (1-\theta_P)^{n-p} \quad \hat{\theta}_P = 9/10$$



$$P(\theta_M), P(\theta_P)$$

Paul is an octopus and can't see stars and doesn't understand they're attached to any football games & \therefore is less likely to not be

superstitious

$$P(\theta_P > 0.5) \approx 0 \quad P(\theta_M > 0.5) \gg P(\theta_P > 0.5)$$

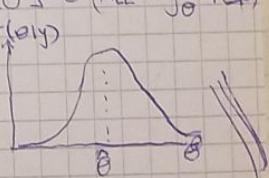
MAP (here) $<$ MLE

Suppose $P(y| \theta)$, $\pi(\theta)$, Bayes gives $\pi(\theta|Y)$

F1: compute an estimation, $\hat{\theta}$, for θ . Study bias etc. $E[\hat{\theta}] = \theta$ (MLE = $\arg\max_{\theta} P(Y|\theta)$)

$$\text{BI: } E[\theta|Y] = \int_{-\infty}^{\infty} \theta \pi(\theta|Y) d\theta$$

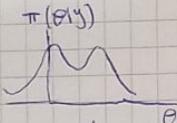
Maximum A posteriori estimate: (MAP) $\hat{\theta} = \arg\max_{\theta} \pi(\theta|Y)$



F2: An a.v. confidence interval for true θ : in an infinite sequence of repeated experiments with data y_1, y_2, \dots

$P(\theta \in [a_i, b_i])$ is α .

$\pi(\theta|y)$



B2: $\pi(\theta|y)$

$$\text{for any } a, b \quad \pi(a < \theta < b | y) = \int_a^b \pi(\theta|y) d\theta$$

if we choose a so that $\pi(a < \theta < b | y) = \alpha$

then $[a, b]$ is a $100\alpha\%$ credible interval for y

F3: consider a hypothesis test involving true θ . we might test $\theta \geq T$ or $\theta < T$ for some T . Data rejects one hypothesis in favour of another

B3: $\pi(\theta|y)$, $\pi(\theta \geq T|y)$, $\pi(\theta < T|y)$

$$\pi(\theta \geq T|y) = \int_T^\infty \pi(\theta|y) d\theta$$

H_0 (null) 'no effect' e.g. $X \sim \text{Ber}(\theta)$ $\theta = \frac{1}{2}$

$$y|x, \theta \sim N(\theta_0 + \theta_1 x, \sigma^2) \quad \theta_0 = 0$$

H_1 (Alternate) $\theta \sim \text{Beta}(a, b)$ $\theta_1 \sim N(\mu, \sigma^2)$

$$\text{Bayes Factor BF} \quad B = \frac{P(y|H_0)}{P(y|H_1)} = \frac{\int p(y|\theta, H_0) \pi(\theta, H_0) d\theta}{\int p(y|\theta, H_1) \pi(\theta, H_1) d\theta} =$$

$$\frac{\int y^n (\frac{1}{2})^n (\frac{1}{2})^{n-y} d\theta}{\int \binom{n}{y} \theta^y (1-\theta)^{n-y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1} d\theta} = \left(\frac{1}{2}\right)^n$$

$$\underline{y} = (y_1, \dots, y_n) \text{ iid } P(y|H_0)$$

prediction: estimate $\hat{\theta} = \theta$ plugin $\hat{\theta}$ to $P(y|\theta)$

predicting \hat{y} (y_{new}) involves assuming $P(\hat{y}|y) \approx P(\hat{y}|\hat{\theta})$

$$P(\hat{y}, \theta|y) = P(\hat{y}|\theta, y) P(\theta|y) =$$

$$P(\hat{y}|\theta) P(\theta|y)$$

$$P(\hat{y}|y) = \int P(\hat{y}, \theta|y) d\theta \text{ marginalisation} = \int_{-\infty}^{\infty} P(\hat{y}|\theta) \pi(\theta|y) d\theta$$

so desired, $P(\hat{y}|y)$ is called Z posterior predictive distribution

y_1, \dots, y_n coin tosses $y_j | \theta \sim \text{Ber}(\theta)$

prior for θ , $\pi(\theta) = 1 \quad \theta \in [0, 1]$

$$\sum y_i = n\bar{y} := s$$

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

$$\pi(\theta | y) = \frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} \theta^s (1-\theta)^{n-s}$$

Posterior predictive

$$p(\tilde{y} | y) = \int_0^1 p(\tilde{y} | \theta) \pi(\theta | y) d\theta = \frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} \int_0^1 \theta^{\tilde{y}} (1-\theta)^{s-\tilde{y}} \theta^s (1-\theta)^{n-s} d\theta =$$

$$\frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} \int_0^1 \theta^{s+\tilde{y}} (1-\theta)^{n-s-\tilde{y}+1} d\theta$$

is I a pds then $\int I d\theta = 1$

$$= \frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} \int_0^1$$

$$= \frac{\Gamma(n+2)\Gamma(s+\tilde{y}+1)\Gamma(n-s-\tilde{y}+2)}{\Gamma(s+1)\Gamma(n-s+1)\Gamma(n+3)} \int_0^1 \frac{\Gamma(n+3)}{\Gamma(s+\tilde{y}+1)\Gamma(n-s-\tilde{y}+2)} \theta^{s+\tilde{y}} (1-\theta)^{n-s-\tilde{y}+1} d\theta$$

$$P(\tilde{y} | y) = \frac{\Gamma(n+2)\Gamma(s+\tilde{y}+1)\Gamma(n-s-\tilde{y}+2)}{\Gamma(s+1)\Gamma(n-s+1)\Gamma(n+3)}$$

$$\Gamma(t+1) = t\Gamma(t) \quad \forall t, \quad \Gamma(1) = 1$$

y, n, s are all integers

$$P(\tilde{y} | y) = \frac{\tilde{y}!}{n+2} \times \frac{(s+\tilde{y})!}{s!} \times \frac{(n-s+1-\tilde{y})!}{(n-s)!}$$

Laplace's law of succession

Poss/
potentially infinite "repeatable" trials. probab is limiting
relative frequency of an event

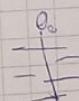
$H \sim \text{Ber}(\theta)$ toss n times, y number of heads

$$\theta(y) = y/n$$

$P(\hat{\theta}(y) > \frac{1}{2}) \times$ (question doesn't make sense)

$P(\hat{\theta}(y) - \varepsilon \leq \theta_c \leq \hat{\theta}(y) + \varepsilon) \times$ (question doesn't make sense since there's no randomness anywhere)

$$\hat{\theta}(Y), \quad P(\hat{\theta}(Y) - \varepsilon \leq \theta_c \leq \hat{\theta}(Y) + \varepsilon) \quad \checkmark$$



Poss/
A random quantity is a quantity whose value is unknown to you

Week 1

1a/ By Bayes thm: posterior: $\pi(\theta|y) \propto \pi(\theta) p(y|\theta)$

$$\propto \theta^{\alpha-1} (1-\theta)^{b-1} \prod_{i=1}^n \theta^{y_i-1} (1-\theta)^{n-y_i-1} \propto \theta^{\alpha+n-1} (1-\theta)^{b+n\bar{y}-n-1}$$
$$\left\{ \prod_{i=1}^n (1-\theta)^{y_i-1} = (1-\theta)^{\sum_{i=1}^n (y_i-1)} \right\}$$

$$\theta|y \sim \text{Beta}(\alpha+n, b+n\bar{y}-n)$$

∴ $\pi(\theta|y)$ is proportional to a Beta density

$$\pi(\theta|y) = \frac{\Gamma(\alpha+b+n\bar{y})}{\Gamma(\alpha+n)\Gamma(b+n\bar{y}-n)} \theta^{\alpha+n-1} (1-\theta)^{b+n\bar{y}-n-1} \quad \theta \in [0, 1]$$

$$1b/ E[\theta|y] = \int_0^1 \theta \pi(\theta|y) d\theta = \frac{n}{\alpha+n} = \frac{\alpha+n}{\alpha+b+n\bar{y}}$$

Three prisoners / A, B, C Alan thinks all equally likely to go

$$\text{Free} \therefore P(A) = P(B) = P(C) = \frac{1}{3}$$

No information about A is you tell me Z name \checkmark , some one to be executed

Jailer says truth - B will die

Note Alan thinks its Me or ~~One~~ $P(A) = \frac{1}{2}$

b be Event jailer says this happen $P(b|A) = P$

$$\text{not } x: P(A|b) \quad \left\{ x = \frac{P}{P+1} \right\} \quad x = \frac{P(b|A)P(A)}{\sum_{i=A,B,C} P(b|i)P(i)} =$$

$$\frac{P}{P+1}$$

$$\text{Alan arguing } P(A|b) = P(A) \therefore x = \frac{1}{3} \therefore \frac{1}{3} = \frac{P}{P+1} \therefore \frac{1}{3}P + \frac{1}{3} = P \therefore$$

$$\frac{2}{3}P = \frac{1}{3} \therefore 2P = 1 \therefore P = \frac{1}{2} \therefore P = \frac{1}{2}$$

$$\text{Alan thinks } P(A|b) \neq P(A) \therefore x = \frac{1}{2} \therefore \frac{P}{P+1} = \frac{1}{2} \therefore \frac{1}{2}P + \frac{1}{2} = P \therefore$$

$\frac{1}{2}P = \frac{1}{2} \therefore P = 1$ but this goes against commutative law

SHELF ② Elicit individual judgments

2) group discussion 3) group consensus

Tessary's prior Beta($\frac{1}{2}, \frac{1}{2}$)

$$\text{Beta: } \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\text{Beta}(1, 1) = \text{Unif}(0, 1)$$

$$\text{Beta}(a, b)$$

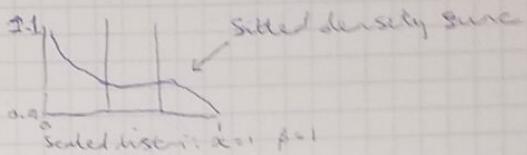
Median for quantity of interest (θ) \times

• Median most quantity is equally likely to be lower or higher
 $0.5 \rightarrow \theta$ is prob of tails

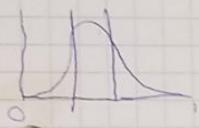
Quartiles for 2 quantity of interest (QoI) suits the Median \pm

2 median and 2 bounds

$0.25 \Delta 0.75$



options, see rottinham.ac.uk



go is median is 0.505 because on coin $\therefore [574, 563] = [\alpha, \beta]$

$\text{Beta}(a+y, b+n-y)$ n = tosses y = heads

Code on the week 2 slides 5

with 2 tails out of 5 you get 2 jerseys prior

its becoming dissident for Beta to have any influence at all because of your prior knowledge

$P(y|\theta)$ likelihood choose $\pi(\theta)$ be anything $\pi(\theta|y) \propto \pi(\theta) P(y|\theta)$
non-informative prior

$$\left(\text{scaled density } \int_{-\infty}^{\infty} P(y|\theta) \pi(\theta) d\theta \right)$$

Def 1.3 / is F is a class of sampling distris with params θ , $P(y|\theta)$ & P is a class of prior distris for θ then if $\forall \pi(\theta) \in P$, we have $\pi(\theta|y) \in F$ then F class P is conjugate for F

Ex / y is 1 is a coin toss is heads suppose n coin tosses

$y|\theta \sim \text{Bin}(n, \theta)$ $\pi(\theta|y) \propto \pi(\theta) P(y|\theta)$

$$\begin{aligned} \text{By Bayes theorem} \quad \pi(\theta|y) &\propto \pi(\theta) P(y|\theta) \\ &\propto \theta^{a-1} (1-\theta)^{b-1} \theta^y (1-\theta)^{n-y} \propto \theta^{a+y-1} (1-\theta)^{b+n-y-1} \end{aligned}$$

$\pi(\theta|y)$ is proportional to a Beta density $\textcircled{1}$

$$e^{lyn} \text{Beta}(ay, b+ny)$$

$\textcircled{2}$ Beta distri is conjugate for Binomial Models

$$E[\theta|y] = \frac{ay}{ay+b+ny} = \frac{ay}{ay+n} \quad \text{as } n \rightarrow \infty \quad E[\theta|y] \rightarrow \frac{y}{n} = \bar{y}$$

$$\theta = \arg\max \pi(\theta|y) = \frac{ay+1}{ay+n-2} \quad \text{tends to } \frac{y}{n} = \bar{y} \quad \text{as } n \text{ tends to } \infty \quad \textcircled{3}$$

measuring an object of length θ , instrument with known error

$$\text{var } \sigma^2 \quad \therefore \quad y_i|\theta, \sigma^2 \sim N(\theta, \sigma^2)$$

suppose have prior for $\theta \sim N(\theta_0, \sigma_0^2)$

given n measurements $y_1, \dots, y_n = \bar{y}$, find $\pi(\theta|y)$

By Bayes Theorem $\pi(\theta|y) \propto \pi(\theta) P(y|\theta)$

$$\propto \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left\{-\frac{1}{2\sigma_0^2}(\theta-\theta_0)^2\right\} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i-\theta)^2\right\}$$

$$\propto \exp\left\{-\frac{1}{2\sigma_0^2}(\theta-\theta_0)^2\right\} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i-\theta)^2\right\}$$

$$\propto \exp\left\{-\frac{1}{2}\left(\frac{1}{\sigma_0^2}(\theta-\theta_0)^2 + \frac{1}{\sigma^2} \sum_{i=1}^n (y_i-\theta)^2\right)\right\}$$

$$\propto \exp\left\{-\frac{1}{2}\left(\frac{1}{\sigma_0^2}(\theta^2 - 2\theta_0\theta + \theta_0^2) + \frac{1}{\sigma^2} \sum_{i=1}^n (y_i^2 - 2y_i\theta + \theta^2)\right)\right\}$$

$$\propto \exp\left\{-\frac{1}{2}\left(\frac{1}{\sigma_0^2}(\theta^2 - 2\theta_0\theta) + \frac{n}{\sigma^2}(\theta^2 - 2\bar{y}\theta)\right)\right\}$$

$$\propto \exp\left\{-\frac{1}{2}\left(\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)\theta^2 - 2\left(\frac{\theta_0}{\sigma_0^2} + \frac{n\bar{y}}{\sigma^2}\right)\theta\right)\right\}$$

$$\propto \exp\left\{-\frac{1}{2}\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)(\theta^2 - 2\left(\frac{\theta_0}{\sigma_0^2} + \frac{n\bar{y}}{\sigma^2}\right)\theta)\right\}$$

$$\propto \exp\left\{-\frac{1}{2}\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)(\theta - \frac{(\theta_0 + n\bar{y})}{(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2})})^2 - \left(\frac{(\theta_0 + n\bar{y})}{(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2})}\right)^2\right\}$$

$$\propto \exp\left\{-\frac{1}{2}\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)(\theta - \left(\frac{\theta_0 + n\bar{y}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}}\right))^2\right\} \quad \therefore$$

$$\theta|y \sim N(\theta_n, \sigma_n^2) \quad \text{if } x \sim N(\mu, \sigma^2) \text{ is } S(n) \propto \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \quad \therefore$$

$$\theta_n = \frac{\left(\frac{\theta_0}{\sigma_0^2} + \frac{n\bar{y}}{\sigma^2}\right)}{\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)}$$

$$\therefore \sigma_n^{-2} = \frac{1}{\left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)}$$

$\therefore \pi(\theta|y) \propto$ to a Normal density
 \therefore as $n \rightarrow \infty$: $\theta_n \rightarrow \bar{y}$

two data sets y_1, y_2

consider a prior, $\pi(\theta)$, that's 'uninformative'

$$\begin{aligned} \pi(\theta|y_1, y_2) &= \pi(\theta|y_1)\pi(y_1|\theta, y_2) = \pi(\theta|y_1)\pi(y_2|\theta) \\ &= \pi(\theta)\pi(y_1|\theta)\pi(y_2|\theta) \end{aligned}$$

'Objective' Bayes (O-Bayes)

$$P(y|\theta) \xrightarrow[\text{automatic choice}]{\quad} \pi(\theta)$$

• Not to move to posterior (uninformative)

• Have 'good' frequentist properties

$$P(\theta \in [a, b]) = 0.95$$

• Be transformation invariant

one approach is to 'uninformative prior' $\pi(\theta) \propto 1$

$\pi(\theta|y) \propto P(y|\theta)$ (Some invariance as for classical, different interpretations)

(Ex: A prior (or posterior) is proper if it is a valid probability distribution (eg its pdf integrates to 1) otherwise: it's improper)

Ex: $y \sim N(\theta, \sigma^2)$ $\pi(\theta) \propto 1$ $\theta \in (-\infty, \infty)$

$$\int_{-\infty}^{\infty} d\theta = \infty \therefore \pi \text{ is improper}$$

$$\text{but } \pi(\theta|y) \propto \exp\left\{-\frac{1}{2\sigma^2}(y-\theta)^2\right\} \sim N(y, \sigma^2)$$

\therefore posterior is proper

$y|\theta \sim \text{Ber}(\theta)$

uninformative prior $\pi(\theta) = 1$ $\theta \in [0, 1]$



(consider a reparameterisation $\phi = \sqrt{\theta}$ \therefore

$\pi(\phi) = 1$ $\phi \in [0, 1]$ \therefore

this approach is incoherent since $\pi(\theta) = 1 \Rightarrow \pi(\phi) \neq 1$

$$\text{as } \phi = g(\theta) \quad \pi_\phi(\phi) = \left| \det \left(\frac{\partial g^{-1}(\phi)}{\partial \phi} \right) \right| \pi_\theta(g^{-1}(\phi))$$

$$g(\theta) = \sqrt{\theta} \quad \therefore \quad g^{-1}(\phi) = \phi^2 \quad \left\{ g(\theta) = \sqrt{\theta} \therefore \forall \phi = \sqrt{\theta} \therefore \theta = \phi^2 \therefore g^{-1}(\phi) = \phi^2 \right\}$$

$$\therefore g^{-1}(\phi) = \phi^2 = \theta \quad \therefore \quad \frac{\partial g^{-1}(\phi)}{\partial \phi} = \frac{\partial}{\partial \phi} (\phi^2) = 2\phi \quad \therefore$$

$$\pi(\phi) = 2\phi \quad \begin{array}{c} \text{graph of } \pi(\phi) \\ \text{a triangle from } (0,0) \text{ to } (1,2) \end{array} \Rightarrow \pi(\phi) \uparrow$$

Des: For prob model $p(y|\theta)$ with scalar θ & expected info $J(\theta) = E\left[-\frac{\partial^2 L}{\partial \theta^2}\right]$, Jeffreys prior is a distri with density proportional to $J(\theta)^{1/2}$ $\pi(\theta) \propto J(\theta)^{1/2}$

Ex: $y|\theta \sim \text{Bin}(n, \theta)$ $\therefore p(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$

$$L(\theta) = \log(n) + y \log \theta + (n-y) \log(1-\theta) \quad \therefore$$

$$L'(\theta) = \frac{y}{\theta} - \frac{n-y}{1-\theta} \quad \therefore L''(\theta) = -\frac{y}{\theta^2} - \frac{n-y}{(1-\theta)^2}$$

$$\therefore J(\theta) = E\left[\frac{y}{\theta^2} + \frac{n-y}{(1-\theta)^2}\right] = \frac{E[y]}{\theta^2} + \frac{n-E[y]}{(1-\theta)^2} =$$

$$\frac{n\theta}{\theta^2} + \frac{n(1-\theta)}{(1-\theta)^2} = n\left(\frac{1}{\theta} + \frac{1}{1-\theta}\right) = n\left(\frac{1-\theta+\theta}{\theta(1-\theta)}\right) = \frac{n}{\theta(1-\theta)} \quad \therefore$$

$$\pi(\theta) \propto J(\theta)^{1/2} \propto \theta^{-1/2} (1-\theta)^{-1/2} \quad \therefore \theta \sim \text{Beta}(\frac{1}{2}, \frac{1}{2})$$

$$p(y|\theta) \quad \pi(\theta) \rightarrow \pi(\theta|y) \propto \pi(\theta) p(y|\theta)$$

$\pi(\theta)$ represent my or a scientist's beliefs. prior elicitation

$\pi(\theta)$ known psychological potentials
stick to observables

θ is abstract so try to elicit y .

be transparent do sensitivity analysis

Match uncertainty elicitation tool

die from lung cancer is more likely than breast cancer

book-11m person A most likely to be a librarian

X Heuristic / 1) Availability

2) Anchoring - is 400k in interval because he kept saying it?

3) Most people put librarian but there are more teachers & lawyers in the world ∴ More likely to be one of those

∴ 2) Representativeness: conditional probs are judged by symmetry (how similar A is to B) but $P(A|B) \neq P(B|A)$

e.g. prob of having this description giving they're a librarian
is quite high but prob they're a librarian given
this description is quite low

eg $P(A|\text{Teacher})$ low but $P(\text{Teacher}|A)$ high

$$\pi(\theta) \quad y|\theta \sim f(y|\theta) \text{ or } y|\theta \sim g(y|\theta) \quad \pi(\theta)$$

∴ try to infer $\pi(\theta) \rightarrow \pi(y)$

hard to think about $\pi(\theta) \quad p(y|\theta)$

$\pi(\theta)$ induces beliefs about y through prior predictive
distribs $p(y) = \int_{-\infty}^{\infty} p(y|\theta) \pi(\theta) d\theta$ $\pi(y) = \int_{-\infty}^{\infty} p(y,\theta) d\theta =$

$$p(y) = \int_{-\infty}^{\infty} p(y|\theta) \pi(\theta) d\theta$$

give a graph you think



∴ $p(y|\theta)$

$\pi(\theta)$ ∴ leads to:

$$\Sigma, E \Sigma \leftrightarrow P$$

↓ RV
R

Def/ a random quantity (r.q.) X is a number whose true val
is unknown (to you)

Def/ an event admits only two possible vals, True or False

we use 2 convention, if (A is a random event), that $A=1$ is A is True,

$A=0$ is A is False

\wedge 'and' - \vee 'or', \sim ('not' eg $\sim A$, \tilde{A})

$A \vee B$ (TRUE is A or B happen)

$A \wedge B$ (true is A and B happen)

$$\sim(A \vee B) = \tilde{A} \wedge \tilde{B} = \tilde{A} \tilde{B}$$

Des/

let x, y be real numbers, then $x \wedge y = \min(x, y)$ $x \vee y = \max(x, y)$

$$\hat{x} = 1 - x$$

this establishes a duality between r.g's and events

$$A \vee B = \sim(\tilde{A} \wedge \tilde{B}) = \sim(\tilde{A} \tilde{B}) = 1 - ((1-A)(1-B)) = 1 - (1-A-B+AB) =$$

$$A+B-AB$$

$$A \vee (B \vee C) = A \vee (\tilde{B} \wedge \tilde{C}) = \sim(\tilde{A} \wedge (\tilde{B} \wedge \tilde{C})) = 1 - ((1-A)(1-B)(1-C)) =$$

$$1 - ((1-A)(1-B-C+BC)) = 1 - (1-B-C+BC-A+AB+AC-ABC) =$$

$$A+B+C-AB-BC-AC+ABC$$

these questions will be in exam with probab infrom & then :-

change \vee 's into \wedge 's & \wedge 's into \vee 's

Des 2.4/ given a r.g., X your Expectation for X , $E[X]$ is that val, \bar{x} , you would choose on \bar{x} understanding that, you are committed (having chosen \bar{x}) to accepting any bet with gain $C(X-\bar{x})$, where C is both arbit and chosen by an "opponent"

Assumption: Coherence: we assume that you dont wish to lay down bets that with certainty result in a loss for you

Lemma 2.1/ Expectation is a linear func: $E\left[\sum_{i=1}^n a_i X_i\right] = \sum_{i=1}^n a_i E[X_i]$

for n r.g. X_i , constas a_i

$$\text{proof}/ E[aX] = aE[X] \quad E[X+Y] = E[X]+E[Y]$$

(1) let \bar{a} be $E[X]$, $\bar{a}_n = E[aX]$ we must accept all bets with gain $G = C_1(aX - \bar{a}) + C_2(X - \bar{x})$ \therefore

• Be transparent

• 1/0 sensitivity wrt \bar{x}

Set $C_2 = -\alpha C_1$ $C_T = C_1 \alpha X - C_1 \bar{x} - \alpha C_1 X + \alpha C_1 \bar{x} = C_1 (\alpha \bar{x} - \bar{x})$

by coherence (to avoid sure loss) $\alpha \bar{x} = \bar{x}$

$$(2) E[X] = \bar{x} \quad E[Y] = \bar{y} \quad Z = X+Y \quad \bar{z} = E[Z] \quad \bar{z} = \bar{x} + \bar{y}$$

$$\begin{aligned} G &= C_1(X+Y-\bar{z}) + C_2(X-\bar{x}) + C_3(Y-\bar{y}) = \\ &= (C_1+C_2)X + (C_1+C_3)Y - C_1\bar{z} - C_2\bar{x} - C_3\bar{y} \end{aligned}$$

$$C_1 = -C_2 = -C_3$$

$$G = -C_1\bar{z} + C_1\bar{x} + C_1\bar{y} = C_1(\bar{x} + \bar{y} - \bar{z})$$

To be coherent, $\bar{x} + \bar{y} = \bar{z}$ proves (2) \therefore

2 general proofs follows by induction \square

\Def 2.3 given r.g. X , your expectation for X , $E[X]$, is fixed
that \bar{x} , you would choose is you must suffer penalty $L = \frac{(X-\bar{x})^2}{k}$
where X is revealed. k is an arbit const

assumption: Coherence, you do not have \rightarrow preference for a given
penalty if you have 2 option so another that is certainly smaller

\Thm 2.1 Def 2.4 & 2.5 are equivalent / \proof let \bar{x} be $E[X]$

under 2.4 let \bar{x} be $E[X]$ under 2.5

under 2.5 gain $-(X-\bar{x})^2$ is preserved to $-(X-x)^2$ for $x \neq \bar{x}$

\therefore Gain $G_T = (X-x)^2 - (X-\bar{x})^2$ is preserved to 0 $\forall x \neq \bar{x}$

$$a = \bar{x}, b = x, c = \frac{1}{2}(a+b)$$

$$G_T = (X-b)^2 - (X-a)^2 = (X^2 - 2bX + b^2) - (X^2 - 2aX + a^2) = 2(a-b)X - (a^2 - b^2) =$$

$$2(a-b)(X - c) \quad \therefore E[G_T] > 0$$

$$\therefore \text{under 2.4 } E[G_T] = 2(a-b)(\bar{x} - c) > 0$$

$E[G_T]$ is positive i.e. $a > b \geq \bar{x} > c$ or $b > a \geq c > \bar{x}$

\bar{x} must be closer to \bar{x} (a) than arbit x (b) ($\forall x$)

letting $x \rightarrow \bar{x}$ $\bar{x} \rightarrow a$ $\bar{x} = \bar{x}$ is 2 optimal choice \square

Thm 2.2 (Boundedness of expectation) /

$$\inf x \leq E[x] \leq \sup x$$

Proof / first suppose $E[x] < \inf x$..

$$\inf x < -E[x] \Leftrightarrow x - \inf x < x - E[x] \Leftrightarrow$$

$(x - \inf x)^2 < (x - E[x])^2$ this implies we could reduce our penalty (\downarrow)

by changing $E[x]$ to $\inf x$ (indep of x) .. by coherence $E[x] \geq \inf x$

Suppose instead $E[x] > \sup x$ then $x - \sup x > x - E[x] \Leftrightarrow$

$(x - \sup x)^2 < (x - E[x])^2$ hence you can reduce your penalty independently

of x , by choosing $E[x] = \sup x$ instead ($\therefore \perp \text{X}$) ..

so by coherence $E[x] \leq \sup x$ \square

$$E[x] = \int x s(x) dx \quad \sum_x x p(x=x)$$

let A be an event

Def 2.6 (probab) / your probab for A ($P(A)$) is $E[A]$

Thm 2.3 / $0 \leq P(A) \leq 1$

Proof / this follows immediately from boundedness of expectation \square

Thm 2.4 / let A_1, \dots, A_n be incompatible (mutually exclusive) events

let A be \exists event that one of A happens $P(A) = P(A_1) + \dots + P(A_n)$

Proof / $A = A_1 \vee A_2 \vee \dots \vee A_n = A_1 + A_2 + \dots + A_n$ (incompatible)

$P(A) = P(A_1 + \dots + A_n) = P(A_1) + P(A_2) + \dots + P(A_n)$ by linearity of expectation

Thm 2.5 / for finite partition A_1, \dots, A_n then event $A = A_1 \vee \dots \vee A_n$

has probab 1

Proof / $P(A) = P(A_1) + \dots + P(A_n)$ {by thm 2.4} = 1 ..

$A = A_1 + \dots + A_n = 1$ so $P(A) = 1$ \forall in order to be coherent

{give this proof in exam}

\exists partition with $n=2$: A, \tilde{A} .. (thm 2.5) implies $P(\tilde{A}) = 1 - P(A)$

$$P(A \vee B) = P(\sim (\tilde{A} \wedge \tilde{B})) = P(1 - ((1-A)(1-B))) = P(A+B - AB) =$$

$$P(A) + P(B) - P(AB) \quad \text{by linearity of expectation} \quad \begin{matrix} \leftarrow \\ \text{mark in exam} \end{matrix}$$

\therefore be very careful

$$E[X], E[X|B] \quad P(A|B) = \frac{P(A \wedge B)}{P(B)}$$

\Des 2.7 / conditional Expectation // given $r=7 \times 8$ possible event H

Suppose you are subject to a penalty $L = H(\frac{X-\bar{x})^2}{k})$

where X is revealed. your conditional expectation, $E[X|H]$ is 2 and \bar{x}
that you would choose for this

\Des 2.8 / Conditional probability

if A is an event with H, then above: $P(A|H) = E[A|H]$

\Thm 2.6 (compound probab.) // for events H, A, 2 probabs $P(HA)$

$P(H) \otimes P(A|H)$ where \otimes is $P(HA) = P(A|H)P(H)$

\Proof // take $k=1$ 2 L we since having chosen $P(A|H)=x$,

$$P(H)=y \quad \text{and} \quad P(HA)=z \quad \text{is} \quad L = H(A-x)^2 + (H-y)^2 - (HA-z)^2$$

$$HA=1 \quad (H=A \equiv 1), \quad H=1, \quad A=0 \quad (HA=0), \quad H=0, \quad HA=0$$

$$\text{3 cases: } HA: L = u(x,y,z) = (-x)^2 + (1-y)^2 + (1-z)^2$$

$$H \tilde{A}: L = v(x,y,z) = x^2 + (1-y)^2 + z^2$$

$$\tilde{H}: L = w(x,y,z) = y^2 + z^2 \quad \text{here to show } z \rightarrow xy \text{ or there is a choice of } x, y, z \text{ that simultaneously reduces } u, v \in W$$

this only happens gradients $\nabla u, \nabla v, \nabla w$ are coplanar

$$\nabla u, \nabla v, \nabla w \text{ are coplanar} \Leftrightarrow J(x,y,z) = \begin{vmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} & \frac{\partial u}{\partial z} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} & \frac{\partial v}{\partial z} \\ \frac{\partial w}{\partial x} & \frac{\partial w}{\partial y} & \frac{\partial w}{\partial z} \end{vmatrix} = 0$$

$$\begin{vmatrix} -2(1-x) & -2(1-y) & -2(1-z) \\ 2x & -2(1-y) & 2z \\ 0 & 2y & 2z \end{vmatrix} = 8z(1-x) + 8xz(1-y) - 8xy(1-z) =$$

$$8(z-xy) = 0 \Leftrightarrow z = xy \quad \square$$

2 thm extends to expectation $E[X|H] = E[X|H](P(H)) \quad \{X \text{ an r.g.\}}$

Bayes thm // $P(A|H) = \frac{P(A)P(H|A)}{P(H)}$ follows immediately

$$Y_1, \dots, Y_n \text{ iid } Y_i | \theta \sim \text{Ber}(\theta) \quad \theta \sim \text{Beta}(a, b)$$

\Des 2.9 / a sequence Y_1, Y_2, \dots of r.g.'s is said to be exchangeable

is 2 joint probab distri of each sub-collection of n quantities

$(Y_{i_1}, \dots, Y_{i_n})$ is 2 same

\Ques 27 (Representation thm) / let y_1, y_2, \dots an infinite exchangeable sequence of r.v.'s with vals in S_n . Then \exists an measure F exist on 2 set of probab measures $\mathcal{Q}(S_n)$ (on S_n) st

for every n & subsets $E_1, \dots, E_n \in S_n$

$$P(Y_1 \in E_1, Y_2 \in E_2, \dots, Y_n \in E_n) = \int \mathcal{Q}(E_1) \dots \mathcal{Q}(E_n) dF$$

$$dF(\theta) = \int s(\theta) d\theta \quad \pi(\theta) \text{ exists}$$

$$P(\theta) = \int p(y|\theta) \pi(\theta) dy$$

Given n indep observations from some $s(y)$, true but unknown data generating process. Suppose we model y with $p(y|\theta)$

$\pi(\theta)$ Consider 2 Kullback-Leibler divergence (KL),

KL(θ) measures "distance" betw 2 probab distris

$$KL(\theta) = \|s(y) - p(y|\theta)\|_{KL} = E \left[\log \left(\frac{s(y)}{p(y|\theta)} \right) \right] = \int \log \left(\frac{s(y)}{p(y|\theta)} \right) s(y) dy$$

if likelihood is continuous w.r.t. θ not on 2 boundary, then \sim
 $n \rightarrow \infty : \hat{\theta}|y \rightarrow N(\theta_0, (nJ(\theta_0))^{-1})$ where $\theta_0 = \arg \min_{\theta} KL(\theta)$

$$\theta_0 = \arg \min_{\theta} KL(\theta) \quad \|s(y) - p(y|\theta)\|_{KL} \quad \hat{\theta} \text{ is } \underset{\theta}{\operatorname{Mfp}} \text{ (posterior mode)}$$

$\hat{\theta} \rightarrow \theta_0$ as $n \rightarrow \infty$ under 2 assumed condns is well known result

(consistency of MAP). $\theta|y$ as $n \rightarrow \infty : \theta|y \sim N(\theta_0, (nJ(\theta_0))^{-1})$

Let $L(\theta) = \log \pi(\theta|y)$, Taylor expansion around $\hat{\theta}$:

$$L(\theta) = L(\hat{\theta}) + \frac{1}{2} (\theta - \hat{\theta})^2 \frac{d^2}{d\theta^2} L(\theta) \Big|_{\theta=\hat{\theta}} + \frac{1}{6} (\theta - \hat{\theta})^3 \frac{d^3}{d\theta^3} L(\theta) \Big|_{\theta=\hat{\theta}} + \dots$$

$$\frac{1}{2} (\theta - \hat{\theta})^2 \left[\frac{d^2}{d\theta^2} \log \pi(\theta) + \sum_{i=1}^n \frac{d^2}{d\theta^2} \log p(y_i|\theta) \right]_{\theta=\hat{\theta}}$$

2nd term is \sim const + sum of n iid r.v.'s with negative mean

$\therefore s(y) = p(y|\theta_0)$ each of 2 n terms has mean $-J(\theta_0)$

else 2 mean is 2 negative \Rightarrow 2 2nd derivative of 2 1st at θ_0 ,

by law of large numbers 2 sum $\rightarrow -\infty$ as $n \rightarrow \infty$ i.e. as $n \rightarrow \infty : \hat{\theta} - \theta_0 \rightarrow 0$

away from $\hat{\theta}$, $L(\theta) \rightarrow -\infty \quad \pi(\theta|y) \rightarrow 0 \quad \therefore \pi(\theta) \propto e^{-L(\theta)}$

$$\pi(\theta|y) \propto \exp \left\{ -\frac{1}{2} n J(\theta_0) (\theta - \theta_0)^2 \right\} \quad \therefore \theta|y \sim N(\theta_0, (nJ(\theta_0))^{-1})$$

Baby's weight $\tilde{y} \mid \mu, \sigma^2 \sim N(\mu, \sigma^2)$

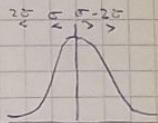
$\pi(\mu, \sigma^2) \quad \mu \mid \sigma^2 \sim N(\mu_0, \sigma^2/\tau)$

$$\pi(\sigma^2) = \frac{\beta^\alpha}{\Gamma(\alpha)} \left(\frac{1}{\sigma^2}\right)^{\alpha+1} \exp\left(-\frac{\beta}{\sigma^2}\right)$$

$\sigma^2 \sim IG(\alpha, \beta)$

$$E[\sigma^2] = \frac{\beta}{\alpha-1}$$

$\mu \pm 2\sigma$



$$\rho(\mu, \sigma^2) = \rho(\sigma^2) \rho(\mu \mid \sigma^2) = \frac{\beta^\alpha \sqrt{\tau}}{\Gamma(\alpha) \sqrt{2\pi}} (\sigma^2)^{-1/2} (\sigma^2)^{-(\alpha+1)} \cdot \exp\left\{-\frac{\beta}{\sigma^2}\right\} \exp\left\{-\frac{\tau}{2\sigma^2}(\mu - \mu_0)^2\right\}$$

$$= \frac{\beta^\alpha \sqrt{\tau}}{\Gamma(\alpha) \sqrt{2\pi}} (\sigma^2)^{-(\alpha+1+\frac{1}{2})} \exp\left\{-\frac{1}{2\sigma^2}(2\beta + \tau(\mu - \mu_0)^2)\right\}$$

y_1, \dots, y_n data: Bayes thm \Rightarrow

$$\pi(\mu, \sigma^2 \mid \mathbf{y}) \propto \pi(\mu, \sigma^2) \prod_{i=1}^n \frac{1}{\sigma} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mu)^2\right\}$$

$$\propto (\sigma^2)^{-(\alpha+1+\frac{1}{2}+\frac{n}{2})} \exp\left\{-\frac{1}{2\sigma^2}(2\beta + \tau(\mu - \mu_0)^2 + \sum_{i=1}^n (y_i - \mu)^2)\right\}$$

$$\propto (\sigma^2)^{-(\alpha+1+\frac{n+1}{2})} \exp\left\{-\frac{1}{2\sigma^2}(2\beta + \tau\mu^2 - 2\mu_0\tau\mu + \mu_0^2\tau + \sum_{i=1}^n y_i^2 - 2\bar{y}\mu + \mu^2)\right\}$$

$$\propto (\sigma^2)^{-(\alpha+1+\frac{n+1}{2})} \exp\left\{-\frac{1}{2\sigma^2}(2\beta + (\tau+n)\mu^2 - 2(\mu_0\tau + \bar{y}\mu) + \mu_0^2\tau + \sum_{i=1}^n y_i^2 + \tau\mu^2)\right\}$$

$$\therefore \alpha_n \rightarrow \alpha + \frac{n}{2} \quad \mu_n \rightarrow \frac{\mu_0\tau + \bar{y}\mu}{\tau + n} \quad \tau_n \rightarrow \tau + n \quad \beta_n = \beta + \frac{1}{2} \sum_{i=1}^n (y_i - \bar{y})^2 + \frac{n\tau}{2(n+\tau)} (\bar{y} - \mu_0)^2$$

\gg library(invgamma)

\gg BirthWeightPrior \leftarrow list(mu=3, alpha=6, beta=5, tau=16)

\gg posteriorUpdate \leftarrow function(Data, Prior){

\gg n \leftarrow length(Data) \gg Posterior \leftarrow list()

\gg Posterior\$alpha \leftarrow Prior\$alpha + n/2

\gg Posterior\$beta \leftarrow Prior\$beta + n

\gg Posterior\$mu \leftarrow (Prior\$mu * prior\$tau + sum(Data)) / (prior\$tau + n)

\gg Posterior\$beta + 0.5 * (n * var(Data) + (n * prior\$tau / (n * prior\$tau + sum(Data))) * (mean(Data) - prior\$mu)^2)

\gg return(Posterior) \gg }

\gg BirthWeightPrior \leftarrow list(mu=3, alpha=6, beta=5, tau=16)

\gg BirthWeightData \leftarrow rnorm(40, mean=3.5, sd=0.5)

\gg BirthWeightPosterior \leftarrow posteriorUpdate(Data = BirthWeightData, Prior = BirthWeightPrior)

$$\theta \in [-1.9, 1.9]$$

$$\# P(\tilde{y} > 4.54 | y) = \int_{-5.9}^{\infty} \int \int p(y|\mu, \sigma^2) \pi(\mu, \sigma^2 | \tilde{y}) d\mu d\sigma^2 d\tilde{y}$$

$$= \int \int \mathbb{1}(\tilde{y} > 4.54) p(y|\mu, \sigma^2) \pi(\mu, \sigma^2 | \tilde{y}) d\mu d\sigma^2$$

Draw μ, σ^2 from $\pi(\mu, \sigma^2 | \tilde{y})$ Draw y

Draw $y|\mu, \sigma^2$ repeat Count num > 4.54

$\pi(\sigma^2 | \tilde{y}) \sim \text{InvGamma}(a_n, b_n)$

$$\pi(\mu | \sigma^2, \tilde{y}) \sim N(\mu_n, \frac{\sigma^2}{n}) \quad \#$$

>> N=1000

>> samples <- rinvgauss(N, 18) # weight Prior \$alpha & beta weight Prior \$beta

$$\# \mu | \sigma^2 \sim N(\mu_0, \frac{\sigma^2}{n})$$

$P(Y > \text{Percy})$

>> phat <- sum(Y > 4.54) / N

$$E[g(\theta) | \tilde{y}] = \int_{-\infty}^{\infty} g(\theta) \pi(\theta | \tilde{y}) d\theta$$

$$\pi(\theta | \tilde{y}) \propto \pi(\theta) p(\tilde{y} | \theta) \int_{-\infty}^{\infty} p(\tilde{y} | \theta) \pi(\theta) d\theta$$

let $\theta \sim \text{Beta}(a, b)$, $y | \theta, f \sim \text{Gamma}(\theta + 1, \beta)$

$$y_1, \dots, y_n \Rightarrow \pi(\theta | \tilde{y}) \stackrel{\text{bycs}}{\propto} \pi(\theta) p(\tilde{y} | \theta) \propto \theta^{a-1} (1-\theta)^{b-1} \prod_{i=1}^n y_i \theta^{-\beta y_i} \\ \propto \theta^{a-1} (1-\theta)^{b-1} \exp\left\{\theta \sum_{i=1}^n \log y_i\right\} \quad \theta \in [0, 1]$$

Def 3.1 Suppose your probab distri is θ . θ is $P(\theta)$ let $\theta_1, \dots, \theta_N$ be N indep random draws from $P(\theta)$. \Rightarrow Monte Carlo

estimate for expectation of a real valued func of θ , $g(\theta)$

$$\text{where } E[g(\theta)] = \int g(\theta) dP(\theta) \quad \text{is } \bar{g}(\theta) = \frac{1}{N} \sum_{i=1}^N g(\theta_i)$$

Let $\theta \sim \text{Beta}(2, 2)$, $E[\theta] = \frac{1}{2}$ can estimate $E[\theta]$ by MC by sampling

$\theta_1, \dots, \theta_N$ from $\text{Beta}(2, 2)$ & averaging $g(\theta_i) = \theta_i$: $E[\theta]$

$$E[\theta] = \int_0^1 \theta \frac{\Gamma(4)}{\Gamma(2)^2} \theta(1-\theta) d\theta \approx \frac{1}{N} \sum_{i=1}^N \theta_i$$

>> theta = 0.5 >> N=10000 >> samples <- rbeta(N, 2, 2)

>> Mestimate <- mean(samples)

$\theta \in [-1.96, 1.96]$ and $\int_{\theta} e^{-\frac{1}{2}\theta^2} d\theta$

$$\therefore \int_{-1.96}^{1.96} e^{-\frac{1}{2}\theta^2} d\theta = \sqrt{2\pi} \int_{-1.96}^{1.96} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\theta^2} d\theta = \sqrt{2\pi} [\Phi(1.96) - \Phi(-1.96)] = 0.95\sqrt{2\pi}$$

$$\theta \sim \text{Unif}(-1.96, 1.96), P(\theta) = \frac{1}{2 \times 1.96} = \frac{1}{3.92} \therefore$$
$$\int_{-1.96}^{1.96} 3.92 e^{-\frac{1}{2}\theta^2} P(\theta) d\theta \quad g(\theta) = 3.92 e^{-\frac{1}{2}\theta^2} \text{ sample of } N(\theta)$$

average $g(\theta)$ as usual

$\gg \text{integral} \leftarrow 0.95 * \text{sqrt}(2 * \pi)$ $\gg N \leftarrow 10000$

$\gg \text{samples} \leftarrow \text{runif}(N, -1.96, 1.96)$

$\gg j \leftarrow \text{function(theta)}$ {

$\gg 3.92 * \exp(-0.5 * \text{theta}^2)$ $\gg }$

$\gg \text{MCestimate} \leftarrow \text{mean}(g(\text{samples}))$ $\gg \text{integral} \gg \text{MCestimate}$

$\gg \text{hist}(\text{samples}) \rightarrow \text{hist}(g(\text{samples}))$ \gg

$\gg \text{abs}(\text{integral} - \text{MCestimate})$

$$E[g(\theta)] \quad \bar{g}(\theta) = \frac{1}{N} \sum_{i=1}^N g(\theta_i) \quad \sigma_g^2 = \text{var}[g(\theta)]$$

$$\text{var}[\bar{g}(\theta)] = \text{var}\left[\frac{1}{N} \sum_{i=1}^N g(\theta_i)\right] = \frac{1}{N^2} \sum_{i=1}^N \text{var}[g(\theta_i)] = \frac{N \sigma_g^2}{N^2}$$

Sample $\text{var}[g(\theta)] = \frac{1}{n-1} \sum_{i=1}^n (g(\theta_i) - \bar{g})^2 = S^2 \therefore S^2$ is an unbiased estimator

for $\sigma_g^2 \therefore$ Monte Carlo standard Error is $\frac{S}{\sqrt{N}}$ $\gg \text{MCerror} = S/\sqrt{N}$

a) $\gg \text{abline}(v = \text{MCestimate}, col = 2)$

b) $\gg \text{abline}(v = \text{MCestimate} - 1.96 * \text{MCerror}, col = 4, lty = 2)$

c) $\gg \text{abline}(v = \text{MCestimate} + 1.96 * \text{MCerror}, col = 4, lty = 2)$

dy $\gg \text{abline}(v = E[\theta], col = 3, lwd = 1.4)$

$\gg \text{MCerror} \leftarrow \text{sd}(g(\text{samples}))/\sqrt{N}$

$\bar{g} = \frac{1}{N} \sum_{i=1}^N [g(\theta_i)] \quad \theta_i \sim P(\theta) \text{ is dim}(g(\theta)) = 1, \text{ for large } n, \text{ by CLT}$

$$\bar{g}(\theta) \sim N(E[g(\theta)], \text{var}[\frac{g(\theta)}{n}]) \quad \frac{\bar{g}(\theta) - E[g(\theta)]}{\sqrt{\text{var}[g(\theta)]/n}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

$\bar{g}(\theta) \pm z_{\alpha/2} \frac{S}{\sqrt{n}}$ is an α level C.I. for $E[g(\theta)]$ this idea

generalises for $y: \mathbb{R}^p \rightarrow \mathbb{R}^q$ $q < p$

$$\pi(\theta \in S | y) = \int_{\theta \in S} \pi(\theta | y) d\theta \quad g(\theta) = \begin{cases} 1 & \theta \in S \\ 0 & \text{otherwise} \end{cases}$$

$\pi(\theta \in S | y) = \int_S g(\theta) \pi(\theta | y) d\theta$ is we can sample $\theta_1, \dots, \theta_N$ from

$\pi(\theta | y)$ then our MC estimate for $\pi(\theta \in S | y)$ is $\hat{P} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}(\theta_i \in S)$

$\hat{P} \approx X \sim \text{Ber}(\hat{P})$ $\text{var}(X) = \hat{P}(1 - \hat{P})$ Monte Carlo Standard error is

$$\sqrt{\hat{P}(1 - \hat{P})/N}$$

Temps are $y_i | \theta \sim N(\theta, \sigma^2)$ $\theta \sim N(\theta_0, \sigma_0^2)$

$$\theta_0 = 20, \sigma_0^2 = 0.75^2 \quad \sigma^2 = 15.6 \quad Y_1, \dots, Y_n$$

$$P(\tilde{y} < 10.3 | y) = \int_{-\infty}^{10.3} P(\tilde{y} | \theta) d\tilde{y} = \int_{-\infty}^{10.3} \int_{-\infty}^{\infty} P(\tilde{y} | \theta) \pi(\theta | y) d\theta d\tilde{y} =$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mathbf{1}(\tilde{y} < 10.3) P(\tilde{y} | \theta) \pi(\theta | y) d\theta d\tilde{y} \quad \text{sample } \theta_i \text{ from } \pi(\theta | y)$$

$$\text{Sample } \tilde{y}_i \text{ from } P(\tilde{y} | \theta = \theta_i) = N(\theta_i, \sigma^2) \quad \{\gg \text{rnorm}\}$$

$$\theta | y \sim N(\theta_0, \sigma_n^2) \quad \frac{1}{\sigma_n^2} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}, \quad \theta_n = \frac{\theta_0 / \sigma_0^2 + n\bar{y} / \sigma^2}{1 / \sigma_0^2 + n / \sigma^2}$$

$\gg \text{Temps} \leftarrow \text{c}()$ $\gg \text{Temperature} \leftarrow \text{tibble}(\text{Temp} = \text{Temps})$

$\gg \text{ggplot}(\text{Temperature}) + \text{geom_hist}(\text{aes}(x = \text{Temp})) + \text{scale_x("Temp")}$

$\gg n \leftarrow 11 \quad \gg \text{Sigma0Sq} = 0.75^2 \quad \gg \text{SigmaSq} = \text{var}(\text{Temps}) \quad \gg \text{theta} = 20$

$\gg \text{Sigma} = \sqrt{\text{Sigma0Sq} + 1 / (1 / \text{Sigma0Sq} + n / \text{SigmaSq})}$

$\gg \text{stDev} = (\text{theta} / \text{Sigma0Sq} + n * \text{mean}(\text{Temps}) / \text{SigmaSq}) / \sqrt{1 / \text{Sigma0Sq} + n / \text{SigmaSq}}$

$\gg N \leftarrow 10000 \quad \gg \text{thetaSamples} \leftarrow \text{rnorm}(N, \text{mean} = \text{theta}, \text{sd} = \text{Sigma})$

$\gg \text{hist}(\text{thetaSamples}) \quad \gg \text{yhatSamples} \leftarrow \text{rnorm}(N, \text{mean} = \text{theta}, \text{sd} = \text{Sigma})$

$\gg \text{hist}(\text{yhatSamples}) \quad \gg \text{phat} \leftarrow \text{sum}(\text{yhatSamples} < 10.3) / N$

$\gg \text{MCerror} \leftarrow \sqrt{(\text{phat} * (1 - \text{phat}) / N)} \quad \gg \text{phat} \gg \text{MCerror} \gg \text{print}$

$\gg \text{paste}("[", \max(0, \text{phat} - 1.96 * \text{MCerror}), ", ", \text{phat} + 1.96 * \text{MCerror}, "]")$, sep = "")

$\gg \text{hist}(\text{thetaSamples}) \gg \text{quantile}(\text{thetaSamples}, \text{probs} = \text{c}(0.025, 0.975))$

$\gg \text{abline}(v = 18.54, \text{col} = 2, \text{lty} = 2) \gg \text{abline}(v = 21.03, \text{col} = 2, \text{lty} = 2)$

$\gg \text{quantile}(\text{yhatSamples}, \text{probs} = \text{c}(0.005, 0.995))$

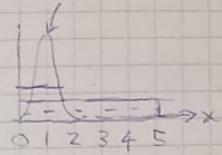
Suppose pdfs $g(x)$ hard (impossible) to sample from, $E[g(x)]$

$$x \in [0, 1] \quad \therefore \int_0^1 x g(x) dx \quad S(x) = 1$$

$$\therefore = \int_0^1 x g(x) S(x) dx \quad \text{set } g^*(x) = x g(x)$$

Samples $N \times$; $\sim \text{Unif}(0, 1)$ $\frac{1}{N} \sum g^*(x_i)$

$$\int_0^1 x g(x) dx = \int_0^1 x g(x) S(x) dx = \int_0^1 g^*(x) S(x) dx$$



$$S(x) \sim \text{Unif}(0, 1) \quad g^{**}(x) = 5xg(x)$$

$\gg g \leftarrow \text{function}(x) \{ \quad \gg x \leftarrow x * \text{dbeta}(x, 2, 2) \quad \gg \}$ $\rightarrow \text{library(tidyverse)}$

$\gg xs \leftarrow \text{seq}(\text{from} = 0, \text{to} = 5, \text{length} = 200)$

$\gg \text{gdata2} \leftarrow \text{tibble}(x = xs, g = \text{dbeta}(x, 2, 2), S = \text{rep}(1, \text{length}(xs)))$

$\gg \text{rplot2} \leftarrow \text{ggplot}(\text{gdata2}) + \text{geom_line}(\text{aes}(x = x, y = 1), \text{color} = "blue") +$
 $\text{geom_line}(\text{aes}(x = x, y = S), \text{color} = "red") \quad \gg \text{rplot2}$

$\gg N \leftarrow 100 \quad \gg \text{MCsamples} \leftarrow \text{gss}(\text{runif}(N, 0, 1))$

$\gg \text{MCestimate} \leftarrow \text{mean}(\text{MCsamples}) \quad \gg \text{MCerror} \leftarrow \text{sd}(\text{MCsamples}) / \text{sqrt}(N)$

$\gg \text{yplot2} \leftarrow \text{geom_vline}(\text{aes}(xintercept} = \text{MCestimate}), \text{color} = 3) + \text{geom_vline}(\text{aes}($

$xintercept = \text{MCestimate} - 1.96 * \text{MCerror}), \text{color} = 3, (\text{ty} = 2) + \text{geom_vline}(\text{aes}($

$xintercept = \text{MCestimate} + 1.96 * \text{MCerror}), \text{color} = 3, (\text{ty} = 2)$

Importance Sampling Theory / let X be on domain D & let $h(x)$ be a density for X with $h(x) \neq 0$ for $x \in D$ $X \sim S(x)$, $E[g(x)]$

$E[g(x)] = \int g(x) S(x) dx = \int \frac{g(x) S(x)}{h(x)} h(x) dx = E_h \left[\frac{g(x) S(x)}{h(x)} \right]$ \Rightarrow var of E_h is minimised when $h(x) \propto |g(x) S(x)|$ eg if $g(x) \geq 0 \therefore$

$E_h \left[\frac{g(x) S(x)}{h(x)} \right] = \text{const} \forall x \quad \therefore \text{var} \left[\frac{g(x)}{h(x)} \right] = 0$

need 1) $h(x) \geq 0$ when $g(x) \neq 0$

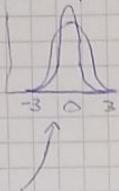
2) $h(x)$ nearly proportional to $|g(x) S(x)|$

3) $h(x)$ easy to sample from

4) $h(x)$ easy to evaluate

Importance Sampling Example / $[-50, 50]$ standard normal i.

$$\int_{-50}^{50} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$



$$\therefore g(x) = \frac{100}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$[-50, 50]$ standard normal i.

$$\int_{-50}^{50} g(x) s(x) dx \quad s(x) \sim \text{Unif}(-50, 50)$$

$$h(x) \sim t_1$$

$$h(x) \sim t_{30}$$

$$\int_{-\infty}^{\infty} \prod_{i=1}^{30} (-50 < x_i < 50) g(x_i) h(x_i) dx$$

t. fat tailed, t₃₀ similar to normal distri

Week 5 / Introduction to Random Number Generation

RNG's $x_{i+1} = (ax_i + c) \bmod M$ $\sim \text{Unif}(0, 1)$, M large

$\frac{x_i}{M}, \dots, \frac{x_N}{M}$ gets "Unif(0,1)"

\gg runs \Rightarrow Random Seed "Mersenne-Twister"

e.g. X cdfs $F(x) = P(X \leq x)$ $F(x) = W$ ($0, 1$) met $P(W \leq w)$

$w \in [0, 1]$ $P(W \leq w) = P(W > w) = 0$

for $w \in [0, 1]$ $P(W \leq w) = P(F^{-1}(W) \leq F^{-1}(w))$ is F^{-1} exists

$= P(F^{-1}(F(x)) \leq F^{-1}(w)) = P(X = F^{-1}(w)) = F(F^{-1}(w)) = w$

$S(w) = F^{-1}(w) = 1 - \log(1-w)$

is F^{-1} exists to obtain a sample with Z same distri as X .

Sample U from $\text{Unif}(0, 1)$ - 2 compute $F^{-1}(U)$

$\forall x / X \sim \text{Exp}(\lambda) \quad S(x) = \lambda e^{-\lambda x} \quad F(x) = \int_0^x \lambda e^{-\lambda t} dt = 1 - e^{-\lambda x} \quad x \geq 0$

$U = F(x) = 1 - e^{-\lambda x} \quad \therefore x = -\frac{1}{\lambda} \log(1-U)$

\gg expsamples $\leftarrow \text{exp}(100, 1) \quad \gg$ u $\leftarrow \text{runif}(100)$

\gg expsamples2 $\leftarrow -1 * \log(1-u) \quad \gg$ par(mfrow=c(1, 1))

\gg plot(sort(expsamples), sort(expsamples2))

Consider discrete r.v. X with values x_1, \dots, x_m , probab. P_1, \dots, P_m

(w.l.o.g.) $x_1 < x_2 < \dots < x_m \quad \therefore F(x_j) = P_1 + P_2 + \dots + P_j : F_j$

$F(x) = F_j \quad x_j \leq x \leq x_{j+1} \quad \text{for } p \in (F_j, F_{j+1})$ F^{-1} doesn't exist

non-unique when $p = F_j$ inverse CDF for Discrete

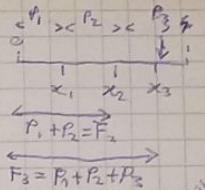
define $F_*^{-1}(p) = x_j \quad p \in (F_{j-1}, F_j] \quad \text{then } U \sim \text{Unif}(0, 1)$,

$X_* = F_*^{-1}(U) \quad \therefore$

(i) X_* can only take vals x_1, \dots, x_m

(ii) $P(X_* = x_j) = P(U \in (F_{j-1}, F_j]) = F_j - F_{j-1} = p_j$

• X_* has same distri as X .



Toss coin (prob $\frac{1}{2}$) 3 times $X = \# \text{num of heads} \sim \text{Bin}(3, \frac{1}{2})$

$$F_1 = P(X=x_1) = P(X=0) = \frac{1}{8} \quad F_2 = P(X=1) = \frac{3}{8} + \frac{1}{8} = \frac{4}{8}$$

$$F_3 = P(X=2) + F_2 = \frac{3}{8} + \frac{4}{8} = \frac{7}{8} \quad F_4 = 1 \quad x_1=0, x_2=1, \dots, x_4=3$$

Draw from $U \sim \text{Unif}(0,1)$ assign $X = F^{-1}(U)$ $F^{-1}(U) = x_j \quad U \in [F_{j-1}, F_j]$

» ThreeSamples ← function(N) { » Us ← runif(N)

 » ifelse(Us <= 1/8, 0, ifelse(Us <= 4/8, 1, ifelse(Us <= 7/8, 2, 3))) » }

» hist(ThreeSamples(1000), freq = FALSE) » hist(rbinom(1000, 3, prob = 0.5))

Box-Muller $X_1, X_2 \sim N(0,1) \quad P(x_1, x_2) = \frac{1}{2\pi} \exp\left\{-\frac{1}{2}(x_1^2 + x_2^2)\right\}$

$$X_1 = R \cos \theta, \quad X_2 = R \sin \theta, \quad |J| = R \quad \therefore$$

$$P(r, \theta) = \frac{1}{2\pi} r e^{-\frac{r^2}{2}} \Rightarrow r, \theta \text{ also indep} \quad P(\theta) = \frac{1}{2\pi} \quad P(r) = r e^{-\frac{r^2}{2}}$$

∴ $P(\theta) \sim \text{Unif}[0, 2\pi]$

$$\text{let } W = R^2, \quad \frac{\partial R}{\partial W} = \frac{1}{2\sqrt{W}} \quad \therefore P(W) = \frac{1}{2\sqrt{W}} e^{-\frac{W}{2}} = \frac{1}{2} e^{-\frac{W}{2}} \quad \therefore W \sim \text{Exp}\left(\frac{1}{2}\right)$$

To sample $X_1, X_2 \sim N(0,1)$ generate θ from $\text{Unif}(0, 2\pi)$

generate $W \sim \text{Unif}(0,1)$ $W = -2 \log(1-U) \quad X_1 = \sqrt{W} \cos \theta, \quad X_2 = \sqrt{W} \sin \theta$

» BoxMuller ← function(N) { » M ← matrix(0, nrow = N, ncol = 2)

 » M ← cbind(M, rnorm(N)) » M ← M * sqrt(2)

 » Us ← runif(N) » Ws ← -2 * log(1 - Us)

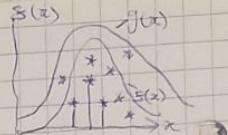
 » X1 ← sqrt(Ws) * cos(Us) » X2 ← sqrt(Ws) * sin(Us)

 » M[1:N, 1] ← X1 » M[1:N, 2] ← X2

 » return(M) » BoxMuller(100) » hist(BoxMuller(100))

» hist(» plot(sort(rnorm(1000)), sort(BoxMuller(1000))))

X density $\delta(x)$ can sample from $g(x)$
don't reject points not under $\delta(x)$



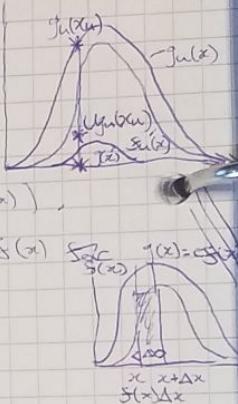
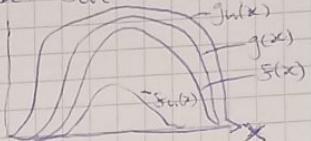
X p.d.f. $\delta(x)$ but $g(x) = c\delta(x)$ const. c

we can sample from

$$\text{since set } g_u(x) = k \cdot \delta_u(x) \geq g(x)$$

1) Draw X_u from $\delta_u(x)$, Uniform(0,1)

2) accept X_u (as a sample from $\delta(x)$) if $g_u(X_u) \leq g(X_u)$



claim /

x is a random pt from Z area under $\delta(x)$ ($X \sim \delta(x)$)

X_* has Z same distri as X . $P(X_* \in [x, x + \Delta x]) \approx \Delta x \delta(x)$ for $\delta(x) = c\delta_u(x)$

if Δx is small $\therefore \delta_{x*}(x) = F'_{x*}(x) =$

$$\lim_{\Delta x \rightarrow 0} \frac{F_{x*}(x + \Delta x) - F_{x*}(x)}{\Delta x} = \lim_{\Delta x \rightarrow 0} \frac{P(X_* \in [x, x + \Delta x])}{\Delta x} = \delta(x)$$

let $g(x) = c\delta(x)$ \therefore X_* is a random pt from Z area under $g(x)$

$$\text{because } P(X_* \in [x, x + \Delta x]) \approx \Delta x g(x) = \Delta x c\delta(x) = \frac{\Delta x c\delta(x)}{c} = \Delta x \delta(x)$$

Sampling rejection samples from a random pt from area under $\delta_u(x) = k \delta_u(x)$

accept if that pt is in Z area under $g(x)$

$$\text{acceptance Rate} = \frac{\text{area under } g}{\text{area under } \delta_u} = \frac{c}{k}$$

$X \sim \text{Beta}(a, b)$ sample X letting $A = a-1$, $B = b-1$

$$\delta(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^A (1-x)^B \quad g(x) = c \delta(x) \quad g(x) = x^A (1-x)^B$$

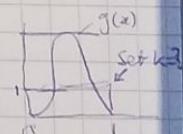
In a distri can sample from $\delta_u(x) = k \delta_u(x) \leq g(x) \forall x$

$\delta_u(x)$ Uniform on domain of $\delta(x)$. ~~so~~ $\delta_u(x) = 1$

choose k s.t. $k \delta_u(x) = k \leq g(x) \forall x$

$$\text{maximise } g(x) \quad h(x) = \log g(x) = A \log x + B \log(1-x) \quad ,$$

$$h'(x) = \frac{A}{x} - \frac{B}{1-x} = 0 \Leftrightarrow x = \frac{A}{A+B} \quad \therefore k = \max g(x) = g\left(\frac{A}{A+B}\right)$$



$\gg a=5 \gg b=1.2 \gg A=a-1 \gg B=b-1$
 $\gg g_{\text{sum}} \leftarrow \text{function}(x) \{ \gg (x^A)^*(1-x)^B \gg \}$
 $\gg t_{\text{bound}} \leftarrow g_{\text{sum}}(A/(A+B)) \# k \neq$
 $\gg xs \leftarrow \text{seq}(\text{from}=0, \text{by}=0.005, \text{to}=1) \gg x \leftarrow g_{\text{sum}}(xs)$
 $\gg \text{plot}(xs, g_{\text{sum}}, \text{type}='l', \text{xlab}="x", \text{ylab}="g(x)", \text{ylim}=(0, t_{\text{bound}}+0.5))$
 $\gg \text{points}(xs, dbeta(xs, a, b), \text{type}='l', \text{tly}=2)$
 $\gg \text{abline}(h=t_{\text{bound}}, \text{tly}=2, \text{col}=4) \gg \text{abline}(v=0, \text{tly}=1, \text{col}=1)$
 $\gg \text{abline}(v=1, \text{tly}=1, \text{col}=1) \gg \text{abline}(h=1)$
 $\gg Us \leftarrow \text{runif}(2) \gg is(bound * Us[2] \leq t_{\text{bound}} * Us[1]) \{$
 $\gg \text{points}(Us[1], bound * Us[2], \text{col}=3, \text{pch}=16) \gg \text{return}(Us[1]) \gg \}$
 $\gg \text{else} \{ \text{points}(Us[1], bound * Us[2], \text{col}=2, \text{pch}=4) \gg \text{return}(NA) \gg \}$
 $\gg \text{abline}() \gg \text{sum100} \leftarrow \text{sum}(1:100, \text{function}(k) \text{ adaptPoint}())$
 $\gg \text{sum}(!\text{is.na}(\text{sum100}))/100$

A density $S(x)$ is log concave if $\frac{d^2 \log S(x)}{dx^2} < 0 \forall x$

proposition / let $h(x) = \log j(x)$, suppose $c > H = \sup h''(x)$, &

$\hat{x} = \arg \max_x h(x)$. Then with $S_u(x) \geq \text{poly}_3 N(\hat{x}, -\frac{1}{H})$

$$g_u(x) = k S_u(x) = e^{h(\hat{x})} = e^{\frac{H}{2}(x-\hat{x})} \geq j(x) \quad \forall x$$

i.e. a rejection sampler is: i) Draw X_u from $N(\hat{x}, -\frac{1}{H})$

$$(\text{set } X_u = \hat{x} + \sqrt{-\frac{1}{H}} Z \quad Z \sim N(0, 1))$$

ii) accept X_u if $U_{\text{uni}}(X_u) \leq j(X_u) \quad U \sim \text{Unif}(0, 1)$

$$\begin{aligned} h(x) &\leq h(\hat{x}) + \int_{\hat{x}}^x h'(t) dt = h(\hat{x}) + \int_{\hat{x}}^x [h'(\hat{x}) + \int_{\hat{x}}^t h''(s) ds] dt \\ h(\hat{x}) &= \int_{\hat{x}}^x \int_{\hat{x}}^t h''(s) ds dt \leq h(\hat{x}) + \int_{\hat{x}}^x \int_{\hat{x}}^t H ds dt = h(\hat{x}) + \frac{H}{2} (x-\hat{x})^2 \end{aligned}$$

$$h_u(x) = h(\hat{x}) + \frac{H}{2} (x-\hat{x})^2 \therefore h_u(x) \geq h(x) \quad \forall x$$

$$\Rightarrow g_u(x) = \exp\{h_u(x)\} \geq j(x) \quad \forall x$$

$$g_u(x) = e^{h(\hat{x})} e^{\frac{H}{2}(x-\hat{x})^2} = e^{h(\hat{x})} \underbrace{\sqrt{\frac{2\pi}{-H}} \sqrt{\frac{-H}{2\pi}} e^{\frac{H}{2}(x-\hat{x})^2}}_{K} = K S_u(x)$$

$$X \sim \text{Beta}(a, b) \quad A = a - 1, \quad B = b - 1 \quad h'(x) = \frac{A}{x} - \frac{B}{1-x}, \quad \hat{x} = \frac{A}{A+B}$$

$$h''(x) = -\frac{A}{x^2} - \frac{B}{(1-x)^2} \quad \text{if } a, b > 1 \text{ then } h''(x) < 0$$

$$\text{if } S \text{ is log concave} \quad h''(x) = \frac{2A}{x^3} - \frac{2B}{(1-x)^3} = 0 \iff x = \frac{A^{1/3}}{A^{1/3} + B^{1/3}}$$

$$H = h''\left(\frac{A^{1/3}}{A^{1/3} + B^{1/3}}\right) = -(A^{1/3} + B^{1/3})^3$$

$$g(x) = e^{h(\hat{x})} e^{\frac{H}{2}(x-\hat{x})^2} = \hat{x}^A (1-\hat{x})^B e^{\frac{H}{2}(x-\hat{x})^2}$$

$$\text{Sample from } N\left(\frac{A}{A+B}, \frac{1}{(A^{1/3} + B^{1/3})^3}\right)$$

Rejection with log concave beta example

>> betaSams <- sumtn(n, a, b) { >> A = a - 1 >> B = b - 1 >> xhat <- A / (A + B)

$$>> H <- (-1)^3 (A^{1/3} + B^{1/3})^3$$

>> z <- rnorm(n) # Normal Sams(n) via Box-Muller ##

>> Xus <- xhat + z * (1 / sqrt(-H)) >> Us <- runif(n)

$$>> xhat <- A * log(xhat) + B * log(1 - xhat)$$

$$>> c <- gamma(a) * gamma(b) / gamma(a + b) >> k <- exp(xhat) * sqrt(z * pi / (-H))$$

>> xs <- seq(Srom = 0, by = 0.005, to = 1) >> gx <- dbeta(xs, a, b)

>> plot(xs, c * gx, type = 'l', xlab = "x", ylab = 'f(x)', ylim = c(0, k * dnorm(xhat, xhat, sqrt(-(1/H)))))

>> points(xs, k * dnorm(xs, xhat, sqrt(-(1/H))), type = 'l', col = 2)

>> print(paste("acceptance rate = ", c / k, sep = " "))

>> Xus[Us * exp(xhat) * exp((H/2) * (Xus - xhat)^2) <= (Xus^A) * ((1 - Xus)^B)]

>> # unding (betaSams) ##

>> betaSams(10000, 5, 1.2)

>> t_betaSams(10000, 5, 1.2) >> hist(t_betaSams, freq = F, breaks = 100, xlim = c(0, 1))

>> xs <- seq(Srom = 0, by = 0.005, to = 1) >> gx <- dbeta(xs, 5, 1.2)

>> plot(xs, gx, type = 'l', xlab = "x", ylab = 'g(x)', ylim = c(0, 5/12))

So $\int_0^1 x^5 (1-x)^3$ can be solved by Monte Carlo

Beta(5, 1.2)

$\{ \int g(x) S(x) dx \}$ $S(x)$ is density we sample from $E[g(x)] = \int g(x) S(x) dx$

$$X_1, \dots, X_n \sim S(x) \quad E[g(x)] \approx \frac{1}{N} \sum_{i=1}^N g(X_i)$$

\checkmark $X \sim \text{Beta}$

$$S(x) = \frac{\Gamma(x+\beta)}{\Gamma(x)\Gamma(\beta)} x^{x-1} (1-x)^{\beta-1}$$

$$\text{or } S(x)=1 \quad X \sim \text{Unif}(0,1)$$

2 Monte Carlo error shrinks at a rate proportional to $\frac{1}{\sqrt{N}}$

$$E[g(x)] \triangleq \int_{-\infty}^{\infty} g(x) S(x) dx \quad S(x)=0 \text{ is } \text{Prob } x \notin [0,1]$$

\checkmark (A+B) always use "exchangeable" instead of iid but is much weaker

$$\pi(\theta|y_j) \propto \pi(\theta) p(y_j|\theta)$$

menti

$$\checkmark 4 / \pi(\theta|y_j) \propto \pi(\theta) p(y_j|\theta) ?$$

$y = y_1, \dots, y_n$ exchangeable with $p(y_i|\theta) \triangleq \pi(\theta)$ \therefore to eval

$$p(y_{n+1} > 0.5 | y) \text{ must eval } \mathbb{I}(y_{n+1} > 0.5)$$

$$\begin{aligned} p(y_{n+1} > 0.5 | y) &= \int_{0.5}^{\infty} \int_{-\infty}^{\infty} p(y_{n+1}, \theta | y) d\theta dy = \int_{0.5}^{\infty} \int_{-\infty}^{\infty} p(y_{n+1} | \theta) \pi(\theta | y) dy d\theta \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mathbb{I}(y_{n+1} > 0.5) p(y_{n+1} | \theta) \pi(\theta | y) dy d\theta \end{aligned}$$

$$\checkmark \text{ (B) } \int g(x) S(x) dx \quad \mathbb{I}(y_{n+1} > 0.5) \text{ is } g(x) \quad p(y_{n+1}, \theta | y) \text{ is } S(x)$$

$$1) \theta_1, \dots, \theta_n \sim \pi(\theta | y)$$

$$2) y'_{n+1}, \dots, y'_{n+n} \sim p(y_j | \theta = \theta_j) \quad j = n+1, \dots, n+n$$

$$p(x, y) = p(x | y) p(y) = p(y | x) p(x)$$

\checkmark have $y = y_1, \dots, y_n$ exchangeable with $p(y_i | \theta) \triangleq \pi(\theta)$ to eval

$p(y_{n+1} > 0.5 | y)$ finally repeat the last 3 steps as many times and
and average the $\mathbb{I}(y'_{n+1} > 0.5)$

4/ $x \sim y | \theta \sim S(y|\theta) \rightarrow \pi(\theta) \xrightarrow{\text{Says}} \pi(\theta|y)$

$$E[g(\theta)|y] = \int g(\theta) \pi(\theta|y) d\theta$$

$$E[g(y)|y] = \int g(y) p(y|\theta) \pi(\theta|y) d\theta$$

Sample $\pi(\theta|y) \rightarrow p(y|\theta)$

$$\pi(\theta_1, \theta_2, \dots, \theta_d | y) = \pi(\theta_1 | \theta_2, \dots, \theta_d, y) \pi(\theta_2 | \theta_3, \dots, \theta_d, y) \dots \pi(\theta_d | y)$$

$$\pi(\theta_d | \theta_{d-1}, \dots, \theta_1, y) \pi(\theta_{d-1} | \theta_{d-2}, \dots, \theta_1, y) \dots \pi(\theta_2 | \theta_1, y) \pi(\theta_1 | y)$$

inverse CDF / Discrete/continuous

Unis $\rightarrow N$ (Ran Muster)

Unis \rightarrow Anything (Rejection Sampling $g(x) \propto S(x)$)

4/ See die inverse CDF method & Unis(0,1)

what will θX do? draw $U=0.345$ give? : 3

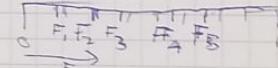
$$F(x) \quad F_1 = P(X=1) = \frac{1}{6}$$

$$F_2 = P(X \leq 2) = P(X=1) + P(X=2) = \frac{2}{6}$$

$$F_3 = \frac{3}{6} \quad F_4 = \frac{4}{6} \quad F_5 = \frac{5}{6} \quad F_6 = 1$$

$$P(X=1) \quad P(X=2)$$

< > < >



4/ inverse CDF See die

3/ $X \sim \text{Bin}(5, 0.25)$ convert $U=0.55 \sim \text{Unis}(0,1)$ into

3 same distri as X : 1

>> rbinom(0:5, 5, 0.25)

$\Rightarrow 0.237, 0.632, \dots$

$$\frac{1}{6} \quad \frac{1}{6}$$

$$F_1 = P(X=x_1) = P(X=0) = \binom{5}{0} 0.25^0 0.75^5 = 0.237$$

$$F_2 = P(X=x_2) + F_1 = 0.237 + P(X=1) = \binom{5}{1} 0.25^1 0.75^4 = 0.632$$

4/ $X \sim \text{Bin}(5, 0.25)$ convert $U=0.922 \sim \text{Unis}(0,1)$ into a sample from

2 same distri as X

$\therefore 0.237, 0.632, 0.876, 0.984, 0.991, 1 \quad \therefore 3$

5/ $U=0.102 \sim \text{Unis}(0,1) \quad \therefore 0$

16/2 rejection sampling algorithm involves 4 funcs
 S , S_u , g , g_u which otherwise are densities? S and S_u

17/ let $X \sim S(x)$ & $g(x) = cS(x)$ which or Z following is most precise?

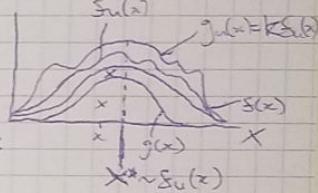
~~17/~~ x -coords of a random sample from Z area under $S(x)$ has Z same distri as X

$$X \sim S(x) \quad g(x) \quad S_u(x) \quad g_u(x)$$

$g(x) = cS(x)$ points in $g(x)$ which is under $S(x)$ has distri X

$S_u(x)$: Easy to sample

$$g_u(x) = kF_u(x)$$



18/ let $X \sim S(x)$ & $g(x) = cS(x)$. Z x -coord of a sample from under $g(x)$ has prob $S(x)$

19/ let S, S_u, g, g_u be 2 funcs used in rejection sampling. which of these is 2 target density of our sample?: S

$$g_u(x) \geq g(x) \quad \forall x$$

sample at random from area under $g_u(x)$ (x^*, y^*) is $y^* \leq g(x^*)$

(x^*, y^*) is a random pt from area under g

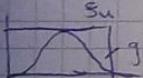
\therefore keep all samples under g . x -coords are samples from S

$U \sim \text{Unif}(0,1) \therefore U g_u(x^*) \sim \text{Unif}(0, g_u(x^*)) \therefore (x^*, U g_u(x^*))$ is random under $g_u(x)$ $\therefore U g_u(x^*) \leq g(x^*)$

19a/ let S, S_u, g, g_u be 2 funcs used in rejection sampling. which func do we sample pt randomly from 2 area underneath: g_u or g

11/ S, S_u, g, g_u which use to bound our target? g_u

12/ for optimal boundary box algo which is unifor density: S_u



13/ $g(x) = cS(x)$ & $g_u(x) = kS_u(x)$
 want c/k large

Sheet 5 / Week 5 / $S(x) \in \begin{cases} x^3 & x \in [0, 3] \\ 0 & \text{otherwise.} \end{cases}$

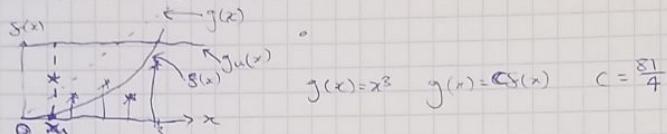
$$F(x) = \int_{-\infty}^x S(u) du \quad S(x) = kx^3 \quad \int_{-\infty}^{\infty} kx^3 dx = 1 \Rightarrow$$

$$\int_0^3 kx^3 dx = \left[\frac{kx^4}{4} \right]_0^3 = \frac{81}{4} k = 1 \quad \therefore k = \frac{4}{81}$$

$$F(x) = \int_0^x \frac{4}{81} t^3 dt = \left[\frac{1}{81} t^4 \right]_0^x = \frac{x^4}{81}$$

$$U = \frac{x^4}{81} \Rightarrow x = \sqrt[4]{81U} = 3U^{1/4} = F^{-1}(U)$$

>> Us ← runis(1000) >> Xs ← 3 * Us^(1/4)

>> hist(Xs) 

$$\backslash 3b / \quad g_u(x) = k S_u(x) \geq g(x) \quad \therefore S_u(x) = \frac{1}{3} \quad g(x) = x^3 \quad g(x) = Cg_u(x) \quad C = \frac{81}{4}$$

$$\therefore \frac{C}{k} = \frac{81}{4 \cdot 81} = \frac{1}{4} \quad \{ \text{Safety check: is it an acceptance rate? Is it between 0 \& 1?}$$

generate X_u from $S_u(x)$

$$(X_u, U_{gu}(x_u)) \sim \text{runis}(0, 1) \quad \therefore (X_u, U_{gu}(x_u)) = (X^*, Y^*)$$

>> accept X^* is $Y^* \leq g(x^*)$

>> reject ← function(N) { >> Xus ← runis(N, 0, 3) # 3 * runis(N) } ↗

>> Us ← runis(N) >> Xus [Us + 27 <= Xus^3] >> }

>> Nsum ← length(Sams)

>> while(Nsum < 1000) {

>> Sums ← c(Sams, reject(1000))

>> Nsums ← length(Sams) >> }

>> Sams ← Sums(1:1000)

>> plot(sort(Xs), sort(Sams))

→ produces a line $y=x$ to show they are sampled from the same model

$$\text{Sheet Week 5/ St8 } \backslash 6a / \quad g(x) \propto \begin{cases} (2+x)(2-x) & x \in [-2, 2] \\ 0 & \text{otherwise} \end{cases}$$

Let $g(x) = (2+x)(2-x)$ since $\subset St$ $g(x) \leq C g(x)$

$$\therefore g(x) = 4 - x^2 \quad \therefore \int_{-2}^2 g(x) dx = C = \left[4x - \frac{x^3}{3} \right]_{-2}^2 = 8 - \frac{8}{3} + 8 - \frac{8}{3} = \frac{32}{3} = C$$

\ 6d / Show $g(x)$ is log concave

\therefore Show $h''(x) \leq 0 \forall x$ $h(x) = \log(g(x)) \therefore$

$$h'(x) = \frac{-2x}{4-x^2} \quad \therefore h''(x) = \frac{-12x}{(4-x^2)^2} - \frac{4x^2}{(4-x^2)^3} < 0 \quad h''(x) = \frac{-2}{(4-x^2)} - \frac{4x^2}{(4-x^2)^2} < 0 \therefore$$

log concave

\ 6e / $\hat{x} = \arg \max g(x)$, $H = \sup h''(x) < 0$

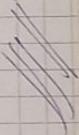
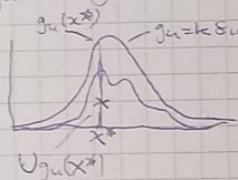
$$g_u(x) \text{ pds } N(\hat{x}, -\frac{1}{H}) \quad g_u(x) = k g_u(x) \quad g_u(x) \geq g(x) \quad \forall x$$

$$K = e^{h(\hat{x})} \sqrt{\frac{2\pi}{-H}}$$

$$\boxed{X^* = \sqrt{\frac{-H}{2}} Z + \hat{x} \quad U \sim \text{Unif}(0,1)}$$

accept X^* if $U g_u(x^*) \leq (2+x)(2-x)$

$$g_u(x) = \frac{\sqrt{-H}}{\sqrt{2\pi}} e^{-H(x-\hat{x})^2} = U g_u(x^*)$$



$$\theta = (\theta_1, \dots, \theta_d) \quad \pi(\theta | y) \quad \pi(\theta_1 | \theta_{-1}, y) \quad \pi(\theta_2 | \theta_{-2}, y)$$

$$\pi(\theta_d | \theta_{-d}, y)$$

$$y^{(i)} = (\theta_1^{(i)}, \theta_2^{(i)}, \dots, \theta_d^{(i)}) = (\theta_1^{(i)}, \theta_2^{(i)}, \theta_3^{(i)})$$

>> theta_j = p[3+j] <-- function(j, mu, sigma, tau2, years, njs) {

>> mom(1, mean(mu/tau2 + njs[j]*gbar[j]/years)/sigma2)/(1/tau2 + njs[j])

>> library(coda)

$$P(\theta_4 > \theta_1) = E[\mathbb{1}(\theta_4 > \theta_1)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mathbb{1}(\theta_4 > \theta_1) \pi(\theta_1, \theta_4 | y) d\theta_1 d\theta_4$$

Monte Carlo est. using 2 samples is ..

>> phat = sum(thetaSamples\$Theta_4 > thetaSamples\$Theta_1)/ntheta

$$\theta^{t+1} \sim \mathcal{N}(\theta^t, \Sigma) \quad \text{Propose } \theta^* \sim q(\theta^* | \theta^t)$$

>> accept with prior $f(\theta^*, \theta^{t+1})$

$$\text{metropolis ratio } \theta^* \sim \mathcal{N}(\theta^t, \Sigma) \quad r(\theta^*, \theta^{t+1}) = \frac{\pi(\theta^* | y) / \pi(\theta^t | y)}{\pi(\theta^{t+1} | y) / \pi(\theta^* | \theta^t)}$$

$$\log(\pi(\theta^* | y)) - \log(\pi(\theta^t | y))$$

$$\Sigma = \begin{pmatrix} S_L & 0 \\ 0 & S_x \end{pmatrix} \quad \frac{1}{\sqrt{N_{est}}} S_e(\rho) \quad \{ N_{est} \text{ is effective sample size} \}$$

Week 8 tutorial lecture / Data $y | \theta \sim \mathcal{S}(y; \theta)$

$$\text{prior } \pi(\theta) \xrightarrow{\text{Bayes}} \pi(\theta | y) \propto \pi(\theta) \prod_{i=1}^n \mathcal{S}(y_i; \theta)$$

$$p(\text{as } \theta \text{ explains } y) = \int_a^b \pi(\theta | y) d\theta \quad P(y_{\text{new}} | y) = \int p(y_{\text{new}} | \theta) \pi(\theta | y) d\theta$$

Linear Model: $y = \beta_0 + \beta_1 x \sim N(\beta_0 + \beta_1 x, \sigma^2)$ Can do 95% error bars: $\pm \frac{1}{2} \sigma$

disc ii: $\pi(\theta_1, \dots, \theta_n)$ is a prior \therefore when condition on a param: $\underbrace{\pi(\theta_1, \dots, \theta_n | y)}_{\pi(\theta_1, \dots, \theta_n | y)}$

$$\pi(\theta_1, \dots, \theta_n | y) \xrightarrow{\text{Bayes}} \pi(\theta_1 | y) \pi(\theta_2 | y) \dots \pi(\theta_n | y)$$

$$\pi(\theta_1 | y) = p(\theta_1 | \theta_2, \dots, \theta_n, y)$$

1/ Menti / 1/ what call property: $p(\theta^t | \theta^{t-1}, \theta^{t-2}, \dots, \theta^1) = p(\theta^t | \theta^{t-1})$

is Markov property (Markov chains)

2/ let X^1, X^2, \dots be a Markov chain. wants Z rule as

$$q(y | x) = p_{X^{(t+1)} | X^{(t)}}(y | x) ? \quad Z \text{ transition density}$$

3/ Markov chain X^1, X^2, \dots $P(x)$ st $\pi(y) = \int p(y | x) p(x) dx$

Z equilibrium distri

- \4/ what's 2 role of equilibrium distri for Bayesian Sampling
So 2 posterior is 2 equilibrium
- \5/ let X be a vec of params & want to sample from $P(X|y)$ by Gibbs sampling. Step 1? start at X^0 with $P(X^0|y) > 0$
- \6/ let X be a vec of params & want to sample from $P(X_i|y)$ by Gibbs sampling. Step 2? sample from $P(X_i^*|X_{-i}^0, y)$
since we've start val $X_{-i}^0 = X_1^0, \dots, X_d^0$
- \7/ let X be a vec of params & want to sample from $P(X|y)$ by Gibbs sampling. Step $j+1$? sample from $P(X_j^*|X_{-j}^0, y)$
since desired: $X_{-j}^0 = X_1^0, \dots, X_{j-1}^0, X_{j+1}^0, \dots, X_d^0$
- \8/ in 2 Metropolis Hastings algorithm, what is 2 role of 2 posterior distri? 2 equilibrium distri ss 2 marker chain
- \9/ in 2 metropolis hastings alg with $q(\cdot|\cdot)$
2 proposal distri
- \10/ $\pi(\theta^*)q(\theta^*|\theta^{t-1})/\pi(\theta^{t-1}q(\theta^{t-1}|\theta^*))$
is 2 Metropolis ratio
- \11/ what makes 2 Metropolis alg & a special case of Metropolis hastings? $q(y|x) = q(x|y) \therefore$

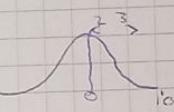
$$\frac{\pi(\theta^*|y)}{\pi(\theta^t|y)} \frac{q(\theta^t|\theta^{t-1})}{q(\theta^*|\theta^{t-1})} = \frac{\pi(\theta^*|y)}{\pi(\theta^{t-1}|y)}$$
 start at θ^{t-1} :: propose θ^* $q(\theta^*|\theta^{t-1})$
calc $r(\theta^*, \theta^{t-1})$
- accept: Set $\theta^t = \theta^*$ prob min(1, r) $\theta^t = \theta^{t-1}$ otherwise
- \12/ has this Markov chain converged? Yes
- \13/ is this Markov chain well mixed? Yes
- \14/ is this Markov chain well Mixed? Yes
- \15/ has this Markov chain converged? No
- \16/ has this Markov chain converged? No have to run for longer

\(\text{V7} \) / is this Markov chain well mixed? No

\(\text{V12} \) / has this Markov chain converged? No

```
>> library(brms) >> library(tidyverse)
>> options(mc.cores = parallel::detectCores())
>> library(bayesplot) >> library(reshape2)
>> load("~/Dropbox/Teaching/Bayes Courses/bayes.RData")
>> drinking <- as_tibble(drinking)
>> drinking <- drinking %>% rename(births = Births)
>> ? Melt
>> drinking %>% melt(id.vars = 'Cirrhosis') %>% ggplot() +
  geom_point(aes(x = value, y = Cirrhosis, colour = variable)) +
  scale_x_continuous(~ variables, scale = "free - x") +
  geom_smooth(aes(value, y = Cirrhosis, colour = variable), method = "(n)")
# y ~ N(a + b*x, sigma^2) #>> lm(Cirrhosis ~ Urban, data = drinking)
>> summary(lm(Cirrhosis ~ ...)) # isn't for bayes i.e. needs to bayes
>> mod1 <- brm(Cirrhosis ~ Urban, data = drinking)
>> summary(mod1) >> plot(mod1) >> mod1$prior # {P(ynew | y)} #
>> myge(drinking$Urban) #>> Urban ~ new & segf(expane = 27, te = 87, by = 1)
>> pred1 <- predict(mod1, newf(X))
>> urban <- tibble(Urban = segf(expane = 27, te = 87, by = 1))
>> pred1 <- predict(mod1, newdata = urban) >> pred1
# E[y] \int y P(y | f, \theta) \pi(\theta | f) d\theta dy P(Q_{2.5} < y_{new} < Q_{97.5}) = 0.95
>> plotData <- cbind(urban, pred1) >> plotData
>> ggplot(plotData) + geom_line(aes(x = Urban, y = Estimate), col = 2) +
  geom_line(aes(x = Urban, y = Q2.5), col = 2, lty = 2) +
  geom_line(aes(x = Urban, y = Q97.5), col = 2, lty = 2) +
  geom_point(data = drinking, aes(x = Urban, y = Cirrhosis), col = 4, pch = 16)
# only 3 pts outside 2 quartiles but 2 outliers are far (4 SD away)
```

```

>> mod2 <- brm(cirrhosis ~ urban + births + nribs + spirits, data = drinking)
>> summary(mod2) # check it's converged first: Rhat = 1
? traceplots are fine >> plot(mod2, N=3) # >> plot(mod2)#
>> predfull <- predict(mod2, newdata = drinking) >> predfull
>> subdata <- cbind(drinking, predfull) >> subdata <- subdata[,-7]
>> full_data <- melt(id.var = c("cirrhosis", "Estimate", "Q2.S", "Q17.S"))
? >? geom_hex() + geom_errorbar(...) + geom_point(y = Estimate) + geom_pointr(y = Cirrhosis) +
  facet_wrap(~ variables, scales = "free-x") # priors are currently
  default but need to be changed to reduce this max uncertainty#
>> urbanPrior <- set_prior("normal(0, 3)", class = "Intercept") # unless 100% certain set
  a prior mean 0 but can change sigma # 
>> birthPrior <- set_prior("normal(0, 3)", class = "births")
>> nribsPrior <- set_prior("normal(0, 3)", class = "nribs")
>> intercept <- set_prior("normal(0, 100)", class = "Intercept")
>> sdprior <- set_prior("normal(0, 10)", class = "sigma")
# can use normal, beta, gamma, t with 3 degrees of freedom #

```

mod2 <- brm(cirrhosis ~ ... , prior = c(urbanPrior, birthPrior, ..., nribsPrior, intercept, sdprior))

week 7

$$y = X\beta + \epsilon \quad \beta_1, \dots, \beta_p \quad \epsilon \sim N(0, \sigma^2 I) \quad y | \beta, \sigma^2 \sim N(X\beta, \sigma^2 I)$$

$$\pi(\beta, \sigma^2) \text{ objective prior } \pi(\beta, \sigma^2) \propto \frac{1}{\sigma^2} \quad (\beta \sim \text{Unif}(-\infty, \infty))$$

$$\pi(\beta, \sigma^2 | y) \propto \pi(\beta | \sigma^2, y) \pi(\sigma^2 | y) \quad \beta | \sigma^2, y \sim N(\underbrace{(X^T X)^{-1} X^T y}_{\hat{\beta}}, \underbrace{(X^T X)^{-1}}_{F})$$

$$\sigma^2 | y \sim \text{Inv}-\chi^2(n-k, S^2) \quad S^2 = \frac{1}{n-k} (y - \hat{X}\hat{\beta})^T (y - \hat{X}\hat{\beta})$$

me $\hat{\beta}$ is $\hat{\beta}$ me σ^2 is S^2 $E[\sigma^2 | y] = \frac{(n-k)S^2}{n-k-2}$ $\beta | \sigma^2 \sim N(\beta_0, V_0)$

$y_i \in \{0, 1\}$ $y_i | X_i, \theta \sim \text{Ber}(p | X_i, \theta)$

$$\pi(\beta | y) \propto \pi(\beta) p(y | \beta) \propto \pi(\beta) \prod_{i=1}^n \left(\frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}} \right)^{y_i} \left(\frac{1}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}} \right)^{1-y_i}$$

y) #

\ Bayesian Hierarchical Models / $p(y|\theta), \pi(\theta)$

Data Layer: $p(y|\theta, \sigma)$ (likelihood)

process Layer: $p(\theta|\sigma)$ Model parameters

prior layer $\pi(\sigma)$

Ex 3.10 / Normal Hierarchical model / J groups, measurements on individuals; in each group. (patients in J hospitals)

Data layer is $y_{ij} | \theta_j, \sigma^2 \sim N(\theta_j, \sigma^2)$

process layer $\theta_j | \mu, \tau^2 \sim N(\mu, \tau^2)$

prior layer $\pi(\mu, \tau^2, \sigma^2) = \pi(\sigma^2) \propto \frac{1}{\sigma^2} \quad \pi(\mu | \tau^2) \propto 1 \quad \pi(\tau^2) \propto \frac{1}{\tau^2}$

$\pi(\sigma^2) \pi(\mu | \tau^2) \pi(\tau^2)$

$$\theta_j | \sigma^2 = \theta_1, \dots, \theta_J \quad \pi(\theta | \mu, \sigma^2, \tau^2 | y) = \pi(\mu, \tau^2 | \theta, \sigma^2, y) \pi(\theta | \sigma^2, y) =$$

$$\pi(\mu | \tau^2, \theta) \pi(\tau^2 | \theta) \pi(\theta | \sigma^2, y) \pi(\sigma^2 | y)$$

$$\pi(\theta | \sigma^2, y) \propto \prod_{j=1}^J \pi(\theta_j | \sigma^2) \propto \prod_{j=1}^J \pi(\theta_j | \theta, \sigma^2) \propto$$

$$\pi(\theta | \sigma^2) \prod_{j=1}^J \int_{\theta_j} \pi(\theta_j | \theta, \sigma^2) d\theta_j \propto \int \pi(\theta | \mu, \tau^2) d\mu d\tau^2 \prod_{j=1}^J \int_{\theta_j} \pi(\theta_j | \theta, \sigma^2) d\theta_j$$

$$\propto \pi(\theta | \mu, \tau^2) \pi(\mu, \tau^2) d\mu d\tau^2$$

\ Des 3.2 / A Markov chain is a sequence of r.v.'s for which,

for any t , π distri $\theta^t | \theta^{t-1}, \theta^{t-2}, \dots, \theta^0$ only depends on θ^{t-1}

$$P(\theta^t | \theta^{t-1}, \dots, \theta^0) = P(\theta^t | \theta^{t-1})$$

let X be a r.v. (multidimensional) Markov chain $X^{(0)} \rightarrow X^{(1)} \rightarrow \dots$

with transition density $P(y|x) \equiv P_{X^{(t+1)}|X^{(t)}}(y|x)$

\ π equilibrium distri for a Markov Chain is π density.

i $\pi(x)$, for which $\pi(y) = \int y(x) \pi(x) dx$

State $X^{(t)}$ $P_{X^{(t+1)}}(x) = ?$, $P_{X^{(t+1)}}(y) = ?$

$$P_{X^{(t+1)}}(y) = \int P_{X^{(t+1)}|X^{(t)}}(y|x) dx = \int P_{X^{(t+1)}|X^{(t)}}(y|x) P_{X^{(t)}}(x) dx$$

$$= \int \pi(y|x) P_{X^{(t)}}(x) dx$$

$$\therefore P_{X^{(t+1)}}(\cdot) = P_{X^{(t)}}(\cdot) = \pi(\cdot)$$

is $\pi(x)$ is our equilibrium distri π is our target density

(we want π posterior) & if $X^{(t)}$ has marginal $\pi(x)$, then $X^{(1)}, \dots, X^{(N)}$ all have distri $\pi(\cdot)$

Need a $\psi(y|x)$ s.t π posterior is $\pi(x)$ (equilibrium distri).
 $\Delta X^{(t)}$ from $\pi(\cdot)$. Under weak conditions on $\psi(y|x)$:

$$P(X^{(t+1)}|X^{(t)}=x) \xrightarrow{t \rightarrow \infty} \pi(x) \therefore \text{MCMC is about-}$$

choosing $\psi(y|x)$ s.t π posterior is $\pi(x)$.

Choosing $\psi(y|x)$ s.t $X^{(t)}$ $P_{X^{(t)}|X^{(t-1)}=x}$ converges to $\pi(x)$ (quickly)
Checking convergence has happened.

Def 3.3 / let $\phi(y|x) = P_{X^{(t+1)}|X^{(t)}=x}(y|x)$ be = transition density for
a Markov chain $\psi(\cdot|\cdot)$ is said to satisfy "detailed balance"
for a distri $P(\cdot)$ is $\psi(y|x)P(x) = \psi(x|y)P(y)$

Theorem 3.1 / if $\psi(\cdot|\cdot)$ satisfies detailed balance for
 $P(\cdot)$, then $P(\cdot)$ is an equilibrium distri for Z Markov chain
process / need to show $\int \psi(y|x)P(x)dx = P(y)$. Given detailed balance
 $\int \psi(y|x)P(x)dx = \int \psi(x|y)\phi(y)dx = P(y) \int \psi(x|y)dx = P(y) \cdot 1 = P(y)$

If $\psi(\cdot|\cdot)$ is an MCMC algorithm ($\Delta \psi(\cdot|\cdot)$) it satisfies detailed balance
for Z target (posterior distri), then Z Markov chain is a
sequence of correlated samples from Z target.

Gibbs Sampler / $\theta = (\theta_1, \dots, \theta_d)$, data y , $\pi(\theta|y)$

$$\theta^{t+1} = (\theta_1^{t+1}, \dots, \theta_d^{t+1}) \quad (\psi(y|x))$$

Algorithm / For $j=1, \dots, d$ Sample θ_j^t from $\pi(\theta_j | \theta_{-j}^{t-1}, y)$

$$\theta_{-j}^{t-1} := (\theta_1^t, \dots, \theta_{j-1}^t, \theta_{j+1}^{t-1}, \dots, \theta_d^{t-1})$$

θ_1^t from $\pi(\theta_1 | \theta_2^t, \dots, \theta_d^t, y)$

θ_2^t from $\pi(\theta_2 | \theta_1^t, \theta_3^t, \dots, \theta_d^t, y)$

$\bullet \theta_3^t$ from $\pi(\theta_3 | \theta_1^t, \theta_2^t, \theta_4^t, \dots, \theta_d^t)$

$(\theta^t | \theta^{t-1}) \quad \pi(\theta | y)$ all d 'full conditionals'

$$\text{Ex: } \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \sim N\left(\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \cdot \Sigma \right)$$

$$\pi(\theta) \propto 1$$

$$\theta | y \propto \exp\left\{-\frac{1}{2}\left(\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} - \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}\right)^T \sum^{-1} \left(\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} - \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}\right)\right\}$$

$$\theta | y \sim N\left(\begin{pmatrix} y_1 \\ y_2 \end{pmatrix}, \begin{pmatrix} 1-\rho^2 \\ \rho \end{pmatrix}\right)$$

$$\pi(\theta_1 | \theta_2, y) = \pi(\theta_2 | \theta_1, y) = \theta_1 | \theta_2, y \sim N(y_1 + \rho(\theta_2 - y_2), 1 - \rho^2)$$

$$\theta_2 | \theta_1, y \sim N(y_2 + \rho(\theta_1 - y_1), 1 - \rho^2)$$

Gibbs: Start $(\theta_1^0, \theta_2^0) = (6, 0)$.

$\theta_1^{t+1}, \theta_2^t, \theta_3^t, \dots, \theta_n^t$ via Z Gibbs Sampler

$(\theta_1^0 | \theta_2^0, y \rightarrow \theta_1^1 | \theta_2^0, y \rightarrow \theta_2^1 | \theta_1^1, y \dots)$

>> # Gibbs Sampler Code #

Ex 3.2 / Each step of Z Gibbs Sampler satisfies detailed balance (for Z posterior).

Propose / Use relation $\propto Y(y|x) \propto e^y \propto e^{t-y}$

Updating a single elem (single Gibbs step) has transition probab $Y(y|x_i) = p_j(y_j|x_{-j}) \prod_{i \neq j} \delta(y_i - x_i)$

$\delta(\cdot)$ dirac delta, p_j is jth full conditional

$$(\pi(\theta_j^t | \theta_{-j}^{t-1}, y)) \quad Y(y|x_i) p(x) = p_j(y_j|x_{-j}) \prod_{i \neq j} \delta(y_i - x_i) p(x) =$$

$$(p(x_1, \dots, x_{j-1}, y_j, x_{j+1}, \dots, x_d) / \int p(x) dx_j) \prod_{i \neq j} \delta(y_i - x_i) p(x)$$

$$= (p(y_1, \dots, y_d) / \int p(y) dy_j) \prod_{i \neq j} (\delta(y_i - x_i) p(y_1, \dots, y_{j-1}, y_j, y_{j+1}, \dots, y_d))$$

$$= p(y) \psi(x|y) \quad \square$$

NHM (normal hierarchical model) $y_{ij} | \theta_j, \sigma^2, \mu, \tau^2 \sim N(\mu, \tau^2)$ $j=1, \dots, n$

$\theta_j | \mu, \tau^2 \sim N(\mu, \tau^2)$ $\pi(\mu | \tau^2) \propto 1$, $\pi(\sigma^2) \propto \frac{1}{\sigma^2}$ $\pi(\tau^2) \propto \frac{1}{\tau^2}$

$\theta_j \sim N(\bar{\theta}_j, \hat{\sigma}_{\theta_j}^2)$

$$\hat{\sigma}_{\theta_j}^2 = \left(\frac{\mu}{\tau^2} + \frac{n_j}{\sigma^2} \bar{y}_{j*} \right) / \left(\frac{1}{\tau^2} + \frac{n_j}{\sigma^2} \right) \quad \hat{\sigma}_{\theta_j}^2 = 1 / \left(\frac{1}{\tau^2} + \frac{n_j}{\sigma^2} \right)$$

$$\bar{y}_{j*} = \frac{1}{n_j} \sum_{i=1}^{n_j} y_{ij} \quad \mu | \tau^2, \sigma^2, \bar{\theta}, y \sim N(\bar{\theta}, \tau^2/n_j)$$

$$\pi(\mu | \tau^2, \sigma^2, \bar{\theta}, y) \propto \pi(\mu | \tau^2, \sigma^2, \bar{\theta}) p(y | \mu, \tau^2, \sigma^2, \bar{\theta})$$

$$\propto \pi(\mu | \tau^2, \sigma^2) \prod_{j=1}^J \pi(\theta_j | \mu, \tau^2, \sigma^2) \quad \therefore Y \text{ has scaled inverse chi-}$$

squared distri with params $\nu \in \mathbb{S}^2$ $\nu \sim \text{Inv-}\chi^2(\nu, \sigma^2)$ is:

$$p(y | \nu, \sigma^2) = \frac{(2\pi)^{\nu/2}}{\Gamma(\nu/2)} \sum_{j=1}^J \psi^{-\nu/2+1} e^{-y_j^2/\nu}$$

$$\sigma^2 | \mu, \tau^2, \theta, y \sim \text{Inv-}\chi^2(n, \hat{\sigma}^2)$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n \sum_{i=1}^J (y_{ij} - \theta_i)^2 \quad \tau^2 | \mu, \sigma^2, \theta, y \sim \text{Inv-}\chi^2(J-1, \hat{\tau}^2)$$

$$\hat{\tau}^2 = \frac{1}{J-1} \sum_{j=1}^J (\theta_j - \mu)^2$$

Metropolis-Hastings / (i) Down-stating pt θ^0 st $\pi(\theta^0 | y) > 0$,
from (any) distri $p_0(\theta)$

(ii) A t=1, 2, ... Suppose θ^{t-1} is known: sample θ^t from a proposal
distri $q(\theta^t | \theta^{t-1})$

(iii) Compute 2 Metropolis ratio $r(\theta^t, \theta^*) = \frac{\pi(\theta^* | y) / q(\theta^* | \theta^{t-1})}{\pi(\theta^{t-1} | y) / q(\theta^{t-1} | \theta^*)}$

(iv) set $\theta^t = \begin{cases} \theta^* & \text{prob } r \\ \theta^{t-1} & \text{otherwise} \end{cases}$

generate $U \sim \text{Unif}(0, 1)$ set θ^t is $U \leq r$

$$\pi(\theta | y) \propto \pi(\theta) p(y | \theta)$$

$$= \frac{\pi(\theta) p(y | \theta)}{\int p(y | \theta) \pi(\theta) d\theta} = p(y) \quad \left\{ \int p(y | \theta) \pi(\theta) d\theta = p(y) \right\}$$

$$\frac{\pi(\theta^* | y)}{\pi(\theta^{t-1} | y)} = \frac{\pi(\theta^*) p(y | \theta^*) / p(y)}{\pi(\theta^{t-1}) p(y | \theta^{t-1}) / p(y)} = \frac{\pi(\theta^*) p(y | \theta^*)}{\pi(\theta^{t-1}) p(y | \theta^{t-1})} \quad \left\{ \text{only need } \pi(\theta) p(y | \theta) \right\}$$

comment $q(\cdot | \cdot)$ is symmetric $q(\theta^* | \theta^{t-1}) = q(\theta^{t-1} | \theta^*)$

$$r = \frac{\pi(\theta^* | y)}{\pi(\theta^{t-1} | y)} \quad (\text{Metropolis Algorithm})$$

$$\theta^* = \theta^{t-1} + D \quad D \sim \text{MVN}(0, \Sigma) \quad \underset{\text{choose}}{\Sigma}$$

Random Walk proposal

Random Walk Metropolis

$\pi(\theta | y) \rightarrow \theta^{(1)} \rightarrow \theta^{(2)} \rightarrow \dots \theta^{(N)}$ has my MCMC converged? \hat{R}
 $\psi = \frac{1}{N} \sum_{i=1}^N \psi_i$ $E(\psi)$, $\hat{R}(\psi)$ near 1 demonstrates convergence
 $\psi = g(\theta | y)$ \therefore Do I have enough samples?
 $E[\psi] = \frac{1}{N} \sum_{i=1}^N \psi_i$ Monte Carlo ψ_i from $\theta^{(i)}$ (Discarding burn-in samples)
 $\frac{1}{\sqrt{N}} \text{Se}(\psi)$ is wrong! Effective Sample Size N_{eff} is st
 $\text{MC error} = \frac{1}{\sqrt{N_{\text{eff}}}} \text{SD}(\psi)$

N_{eff} estimated by our software from your $MCMC < N$ ($<< N$)
 {sometimes much less than N : MCMC $<< N$ sometimes}
 good MCMC has N_{eff} close to N (but not N)

{use hamiltonian Monte Carlo for CW?}

Generate θ^* from $q(\theta^* | \theta^{t-1})$ acceptance rate (based on average r) is too low, chain gets stuck.

$$q(\theta^* | \theta^{t-1}) \sim N(\theta^{t-1}, \Sigma)$$

idea is 'tune' Σ (automatically) so acceptance rates are 'good' (in 1 dimensional param θ want ≈ 0.44 , in high dimensions ≈ 0.23)

Adaptive MCMC: (1) set $\Sigma^* = \Sigma_0, k = 4$

(2) run Metropolis Hastings for m steps with proposal

$$q(\theta^* | \theta^{t-1}) = N(\theta^{t-1}, k \Sigma^*)$$

(3) set $\Sigma^* = \text{Var}[\theta^1, \dots, \theta^m]$ & choose k st 2 acceptance rate is close to 'optimal'

(4) Repeat 2 & 3 a fixed number of times until 2 acceptance rate is optimal. Then Σ^* has converged.

(5) Run standard N step MH with $q(\theta^* | \theta^{t-1}) = N(\theta^{t-1}, k \Sigma^*)$

HMC (Hamiltonian Monte Carlo) is 'Auxiliary variable MCMC'

Desire 'Auxiliary variables', so st $p(\theta, \phi | y) = p(\theta | y)p(\phi)$
target becomes $\pi(\theta, \phi | y)$ keeping only θ

Samples gives Samples from $\pi(\theta | y)$

Require ϕ_j for each θ , $p(\phi)$ is called 'Momentum'

we specify a 'Mass' matrix M ($d \times d$)

- Scalar ϵ
- Integer number of leapfrog steps, L

(1) Draw ϕ^{t-1} from $\phi \sim N(0, M)$ Set $\phi^* = \phi^{t-1}$, $\theta^* = \theta^{t-1}$

(2) Repeat L leapfrog steps: (a) $\phi^* = \phi^* + \frac{1}{2} \epsilon \frac{d \log \pi(\theta^* | y)}{d \theta^*}$

(b) Update θ position vec using momentum $\theta^* = \theta^* + \epsilon M^{-1} \phi^*$

(c) $\phi^* = \phi^* + \frac{1}{2} \epsilon \frac{d \log \pi(\theta^* | y)}{d \theta^*}$

(3) Metropolis ratio $r = \frac{\pi(\theta^* | y)p(\theta^{t-1})}{\pi(\theta^{t-1} | y)p(\theta^*)}$

(4) $\theta^t = \begin{cases} \theta^* & \text{prob min}(1, r) \\ \theta^{t-1} & \text{otherwise} \end{cases}$

Intuition: θ^* is in flat area of $\pi(\theta | y)$,

$\frac{d \pi(\theta | y)}{d \theta} = 0$ leapfrog steps 'skate' along θ

space with const velocity.

Moving in direction of decreasing posterior, $\frac{d \log \pi(\theta^* | y)}{d \theta^*} < 0$.

These take smaller & smaller steps until change direction

Move in direction of increasing posterior $\pi(\theta^* | y)' > 0$

Momentum increases, take bigger steps

>> options(mc.cores = parallel::detectCores())

>> models >> library(brms) library(bayesplot) library(reshape)

>> ggplot

>> fit1 <- brm(Murder ~ low-income + unemployed + pop + I(pop^2), data = muroles)

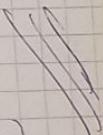
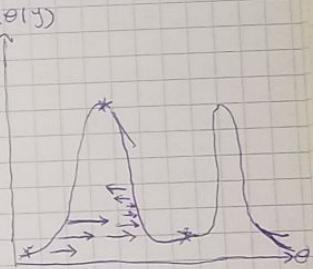
fit1
 $y_i | \beta, \sigma^2 \sim N(\beta_0 + \beta_1 x_1 + \dots, \sigma^2)$ $\pi(\beta, \sigma^2)$ $p(y | \text{data})$

$P(y | \beta = \hat{\beta}, \sigma = \hat{\sigma})$

>> summary(fit1) # that is not 1 but biggest # look at vs section

Sample size (ESS) for each variable

>> plot(fit1) ?brm >> brm(..., nits = "0")



can run for longer

$$y_i | x_i, \beta_i, \sigma^2 \sim N(\beta_0 + \beta_{pop} \text{pop} + \beta_{low} \text{low} + \beta_{unem} \text{Unem}, \sigma^2)$$

$$y \in [0, 40] \quad \text{low} = O(10) \quad \therefore \beta_{low} \sim O(1)$$

$$\text{Unem} = O(1), \beta_{unem} \sim O(10)$$

$$\beta_{pop} \in [-2, 2], \beta_{unem} \in [-2, 2]$$

$$\text{pop} = C(10^6) \quad \therefore \beta_{pop} \sim O(10^{-5}) \quad N(0, 2 \times 10^{-5})$$

`pop_prior <- set_prior("normal(0, 0.00002)", class = "b", coefs = "pop")`

`fit >> sitr2 <- brm(..., prior = pop_prior) >> summary(sitr2) >> plot(sitr2)`

$$\gg \text{Murbes\$pop} \leftarrow \text{Murbes\$pop}/10^6$$

`>> pop_prior <- setprior("normal(0, 10)", class = "b", coefs = "pop")`

charts are 4 & trace plots look good

`preds <- predict(sitr2, restdata = murbes) >> preds`

`>> posterior = samples(sitr2)`

`>> samples <- posterior.samples(sitr2)`

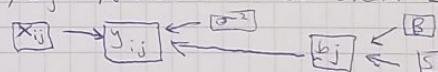
$$\gg \text{sum(samples\$b[, 1]} / 4000 \Rightarrow 0.127 \# \text{probab} \#$$

\ menti / 1/ $y_{ij} | b_j, \sigma \sim N(b_j^T X, \sigma^2)$, $b_j \sim N(B, S^2)$ & a valid

prior comprises what type of Model? Bayesian hierarchical regression

Model

2/ $y_{ij} | b_j, \sigma \sim N(b_j^T X, \sigma^2)$, $b_j \sim N(B, S^2)$ what about σ prior
distribution? $\pi(\sigma, B, S)$



3/ $y_{ij} | b_j, \sigma \sim N(b_j^T X, \sigma^2)$, $b_j \sim N(B, S^2)$, $\pi(\sigma, B, S)$ what is σ Model? or
 b_j called? process layer



4/ ... $\pi(\sigma, B, S)$ what can σ variables in σ data identified by j ?

Grouping variables

5/ what need to be aware of setting $b_j \sim MVN(B, S^2)$, $\pi(\sigma, B, S)$ with
wrms? B must be distri $N(0, var)$ process layer wrt σ

distri with mean 0 $\hat{B}_j \sim N(0, \Sigma) \quad N(B X + f_j X) \quad B_j = B + \hat{B}_j$

6/ $y_{ij} | b_j, \sigma \sim D(g^{-1}(y_{ij}), \sigma)$, $y_{ij} = b_j^T X_{ij}$, $b_j \sim N(B, S^2)$, $\pi(\sigma, B, S)$ what

is η ? \rightarrow linear predictor

\checkmark what is D ? \rightarrow family

\checkmark what is $g^{-1}(\eta_{ij})$? \rightarrow mean of D

\checkmark why do we need a link? η is link \therefore to make sure η means D

D is within η support of η distri

make sure η restrictions on D mean of D are met

\gg View(verbagg) \gg dim(verbagg)

(8.1.2) \gg verbset1 \leftarrow brm(... ~ Anger + Gender + btype + situ + (1 | id) + (1 | item), data = verbagg, tr, family = bernoulli())

(-1-) \rightarrow (predictors to group | Corr group variables)

$$\beta_1 \mid \beta_2 \dots \sim \text{Ber}((\phi(\beta, x))) \quad p_i = g^{-1} \quad g(z) = \logit(z) = \frac{z}{1 - e^{-z}}$$

$$\logit(\beta, x) = \beta_0 + \beta_1 x_1 + \dots + \beta_4 x_4 + \beta_{0j} + \beta_{1j} x_{1j} + \dots + \beta_{pj} x_{pj}$$

$$\beta_j \sim N(0, \Sigma)$$

$$\beta_{0,id} \sim N(0, \Sigma_{id}) \quad \beta_{item} \sim N(0, \Sigma_{item})$$

$$x_1 = \text{Anger has Coeffs of } b_1 \quad x_2 = \text{Gender Coeffs of } b_2$$

id x_3 btype Coeffs of b_3 x_4 situ Coeffs of b_4

\gg verbset1 \leftarrow brm(... + (Anger | item), ...)

for # \gg brm(... + (0 + Anger | item), ...)

\gg plot(verbset1, nslk=FALSE) \gg summary(verbset1)

\gg preds \leftarrow predict(verbset1, newdata = verbagg, pred)

\gg g_class \gg a_classifier \leftarrow preds[, "Estimate"] > 0.5

\gg confmat \leftarrow confusionMatrix(..., ...)

\gg Confmat \gg sum(diag(confmat))/sum(confmat)

\gg Tonsit \leftarrow brm(TC ~ Co.Star.Age - Difference - (1 | TC ~ age), data = Tonsit)

\gg summary(Tonsit) \gg plot(Tonsit)

\gg imprior \leftarrow set_prior("normal(0, 10)", class = "Intercept")

\gg Tonsit1 \leftarrow brm(..., prior = imprior) \gg summary(Tonsit1) \gg plot(Tonsit1)

\gg tau_prior \leftarrow

We tried to improve the first model by factoring in recent form, that is weighting the average goals in the previous 5 games more heavily than less recent results. We calculated recent form as a team's average goals from the last 5 games, ignoring whether they were played at home or away. We tried weighting form at 25%, 50% and 75%. However all 3 attempts had worse results than the original model (where $w=0$). We believe this is due to the fact that ignoring whether the team was playing home or away when calculating recent form was adding more noise than accuracy to the predictions.

We next tried calculating recent form separately for home and away games, using a similar formula as before, but considering only home or away games. We again tested for 25%, 50% and 75% weighting of recent form and found that 25% gave the best predictions, and was an improvement on the first model.

Calculating form separately for home and away removes the noise induced when we calculated it as a single variable.

We tried again to improve our model further by including conceded goals by the opponent to our 3rd model.

We attempted this by adding goals conceded by the away team into the home team's expected goal number, and vice versa. This would then mean our model took into account the opponent the teams were facing. We added in two factors for conceded goals, conceded away goals and conceded form, both relative to home and away matches.

As previously stated, we added the conceded goals factor to our previous model, leading to our expected goal numbers calculated like so. We weighted all 4 factors equally, and also 12.5%. Weighting on form is as that worked well with model 3. The results showed that although the weightings were ~~too~~ too high, adding this factor did not improve the results.

$$y_i | \beta, \sigma^2 \sim N(\beta x, \sigma^2) \quad y_i | \beta, \phi \sim D(g^{-1}(\phi), \phi) \quad J = \beta^T x \quad \text{Tr}(\phi, \beta)$$

$$\pi(\beta, \sigma^2) \quad \pi(\phi, \beta) \quad \text{brm}(y \sim \dots, \text{data})$$

• brm(b5(y ~ ... , sigma ~ ...), data)

• formula sheet: only need normals table

must remember normal formulae

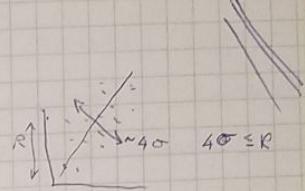
$$N(\beta x, \sigma^2) \Rightarrow y = \beta x + \epsilon \quad \epsilon \sim N(0, \sigma^2) \quad \therefore y - \beta x$$

$$\text{var}(y) = \sigma^2 \quad R \text{ is range of data} \quad \therefore \sigma \leq \frac{R}{4}$$

$$\text{Set } \sigma \sim N(0, (\frac{R}{12})^2)$$

$$\text{rank year prior} \sim N(1100, 500)$$

$$\text{degree year prior} \sim N(750, 400)$$



$$\eta \sim N(\beta_0 + \beta_1 \text{rank} + \beta_2 \cdot \text{rank.years} + \beta_3 \text{Degree} + \beta_4 \text{Degree.years}^2 +$$

$$\cdots b_{ij} + b_{ij} \text{Rank}_j + \cdots b_{ij} \text{Degree.years}_j, \sigma^2) \quad b_{ij} \sim N(0, \Sigma)$$

① give Monte Carlo est + err from a param

② give MC est + err of "some prediction"

③ give MC... as Pr(parameter being prediction thing < 77.5)

Est. Error Q2.5

$$\mathbb{E}_{\text{samples}}[\text{param}] \quad \text{sd(samples)} \quad P(\text{Param} \leq 0.025) \leftarrow 0.975$$

$$\mathbb{E}[g(\theta) | y] = \int g(\theta) \pi(\theta | y) d\theta \rightarrow \theta, \dots, \text{on from } \text{Tr}(\theta | y)$$

$$\approx \frac{1}{N} \sum_{i=1}^N g(\theta_i) \quad \mathbb{E}[g(y) | y] = \iint g(y) p(y | \theta) \pi(\theta | y) d\theta = \mathbb{E}[g(y(x)) | y] =$$

$$\iint g(y(x)) p(y(x) | \theta) \pi(\theta | y) d\theta \quad \text{new data} = (x_1, \dots, x_n) \quad \theta_1, \dots, \theta_n \sim \pi(\theta | y)$$

$$y_1(x_1) | \theta, \dots, y_n(x_n) | \theta$$

$$y_1(x_1) | \theta, \dots, y_n(x_n) | \theta \quad \text{is wanted } E(y(x))$$

$$\mathbb{E}(y(x_{\text{men, pros}}) - y(x_{\text{men, pros}}))$$

$$\mathbb{P}(\theta > 0 | y) \text{ or } \mathbb{P}(y_1 > y_2 | y) \quad \text{if: } g(\theta) = \mathbb{1}(\theta > 0)$$

$$\mathbb{P}(\theta > 0 | y) = \mathbb{E}[\mathbb{1}(\theta > 0)] = \int \mathbb{1}(\theta > 0) \pi(\theta | y) d\theta \approx \frac{1}{N} \sum_{i=1}^N \mathbb{1}(\theta_i > 0)$$

$$\text{if: } P(y_1 > y_2 | y) \quad \therefore P(y_1 > y_2 | y) = \iint \mathbb{1}(y_1 > y_2) p(y | \theta) \pi(\theta | y) d\theta dy$$

$$\frac{1}{N} \sum \mathbb{1}(y_1 > y_2)$$

$$\text{sd}(\text{MC esti}) = \frac{\text{sd}(g(\theta))}{\sqrt{N_{\text{ess}}}}$$

BulkESS >> TailESS

for $E[\text{param}]$ use BulkESS

is want $P(\text{param} > 7.5)$ use TailESS

is want $P(\theta > 0)$ use bulk vs \sum_i tail vs \int

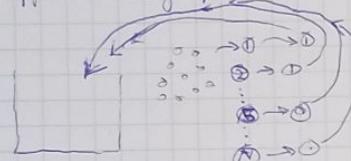
use smallest bulk ESS

$$\text{let } P(\theta > 0 | y) = \hat{P} \quad \therefore \text{MC} = \sqrt{\hat{P}(1-\hat{P})/N_{\text{ess}}} = \text{MC error}$$

Week 3 proof of Z representation theorem / Y_1, Y_2, \dots exchangeable

binary RG's $\therefore Y_i = 1$ is heads, 0 is tails

Suppose assign $P(Y_i=1)=\theta$ exchangeability $\Rightarrow P(Y_i=1)=\theta \quad i=1, \dots, N$



let Z be val of a random counter in Z

bucket $P(Z=1) = \theta$ (by exchangeability)

say θ_k is Z proportion of 1's in Z bucket $P(Z=1|\theta_k) = \theta_k$ (don't

know θ_k) possible vals are $\theta_i = i/N \quad i=0, 1, \dots, N$ Suppose set

$$P(\theta_k = \theta_i) = p_i \quad \text{LOT} \Rightarrow P(Z=1) = \sum_{i=0}^N P(Z=1|\theta_k) P(\theta_k = \theta_i) = \sum_{i=0}^N \theta_i p_i =$$

$$\int_0^1 dF_N(\theta_k) \quad F_N \text{ is cdf giving probab } p_i \text{ to } \theta_k = i/N$$

Suppose n tosses, let X_n be number of 1's in first n tosses.

Exchangeability $\Rightarrow P(X_n=k)$ is Z same for observing k heads

in any sample of n tosses $P(X_n=k)$ must be Z same as our

probab for drawing k counters labelled 1 from our bucket

when we take n out w/ random still have $F_n(\theta)$ giving probab

P to i/N 1's in Z bucket. is we knew θ : $P(X_n=k|\theta) =$

$$\binom{N\theta}{k} \binom{N(1-\theta)}{n-k} / \binom{N}{n} \quad \text{LOT} \Rightarrow P(X_n=k) = \int_0^1 \frac{\binom{N\theta}{k} \binom{N(1-\theta)}{n-k}}{\binom{N}{n}} dF_N(\theta)$$

when N gets large Sampling with/without replacement become

$$Z \text{ same as } N \rightarrow \infty \quad P(X_n=k) = \int_0^1 \binom{n}{k} \theta^k (1-\theta)^{n-k} dF_n(\theta)$$

Helly's thm \Rightarrow any sequence F_n such as this (i/N)

$$\text{Converges to a unique limit } F \quad P(X_n=k) = \int_0^1 \binom{n}{k} \theta^k (1-\theta)^{n-k} dF(\theta)$$

$$= \int_0^1 \binom{n}{k} \theta^k (1-\theta)^{n-k} \delta(\theta) d\theta$$

Week 7 / Full Conditional for μ / $\pi(\mu | \tau^2, \sigma^2, \theta, y) \propto$

$$\pi(\mu | \tau^2, \sigma^2, \theta) P(y | \mu, \tau^2, \sigma^2, \theta) \propto \pi(\mu | \tau^2, \sigma^2, \theta) \propto$$

$$\propto \pi(\mu | \tau^2, \sigma^2) \pi(\theta | \mu, \tau^2, \sigma^2) \propto \prod_{j=1}^J \exp\left\{-\frac{1}{2\tau^2} (\theta_j - \mu)^2\right\} \propto$$

$$\exp\left\{-\frac{1}{2\tau^2} \sum_{j=1}^J (\theta_j - \mu)^2\right\} \propto \exp\left\{-\frac{1}{2\tau^2} \sum_{j=1}^J (\mu^2 - 2\theta_j \mu)\right\} \propto$$

$$\exp\left\{-\frac{1}{2\tau^2} (J\mu^2 - 2y \bar{\theta})\right\} \propto \exp\left\{-\frac{J}{2\tau^2} (\mu - \bar{\theta})^2\right\}$$

$$\mu | \tau^2, \sigma^2, \theta, y \sim N(\bar{\theta}, \frac{\tau^2}{J})$$

Full conditional for σ^2 /

Des 4-3 / A random quantity Ψ , has scaled inverse chi-squared distri with params ν, s^2 ; $\Psi | \nu, s^2 \sim \text{Inv-}\chi^2(\nu, s^2)$ is

$$p(\Psi | \nu, s^2) = \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} s^{2\nu} \Psi^{-(\nu/2+1)} e^{-\frac{s^2}{2\Psi}}$$

$$\pi(\sigma^2 | \theta, \mu, \tau^2, y) \propto \pi(\sigma^2 | \theta, \mu, \tau^2) P(y | \theta, \mu, \tau^2, \sigma^2) \propto$$

$$\frac{1}{\sigma^2} \prod_{j=1}^J \exp\left\{-\frac{1}{2\sigma^2} (y_{ij} - \theta_j)^2\right\} \propto \frac{1}{\sigma^2} \prod_{j=1}^J \frac{1}{\sigma^2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n_j} (y_{ij} - \theta_j)^2\right\} \propto$$

$$\frac{1}{\sigma^2} \frac{1}{\sigma^2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{j=1}^J \sum_{i=1}^{n_j} (y_{ij} - \theta_j)^2\right\} \quad \sqrt{s} = \sqrt{\frac{1}{\sigma^2} \sum_{j=1}^J \sum_{i=1}^{n_j} (y_{ij} - \theta_j)^2}$$

$$\pi(\sigma^2 | \theta, \mu, \tau^2, \theta, y) \propto (\sigma^2)^{(-1 - \frac{J}{2})} e^{-\frac{n\sigma^2}{2}} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^J \sum_{i=1}^{n_j} (y_{ij} - \theta_j)^2$$

Set

$$\sigma^2 | \theta, \mu, \tau^2, \theta, y \sim \text{Inv-}\chi^2(n, \hat{\sigma}^2)$$

Full conditional for τ^2 / $\pi(\tau^2 | \mu, \sigma^2, \theta, y) \propto \{ \text{by Bayes} \}$

$$\pi(\tau^2 | \mu, \sigma^2, \theta) P(y | \sigma^2, \mu, \theta) \quad \therefore \{ y \sim N(\mu, \sigma^2) \text{ indep of } \tau^2 \}$$

$$\propto \pi(\tau^2 | \mu, \sigma^2) \pi(\theta | \tau^2, \mu, \sigma^2) \propto \frac{1}{\tau^2} \prod_{j=1}^J \frac{1}{\tau^2} \exp\left\{-\frac{1}{2\tau^2} (\theta_j - \mu)^2\right\}$$

$$\propto \frac{1}{\tau^{J+1}} \exp\left\{-\frac{1}{2\tau^2} \sum_{j=1}^J (\theta_j - \mu)^2\right\} \propto (\tau^2)^{-(\frac{J}{2} + \frac{1}{2})} \exp\left\{-\frac{1}{2\tau^2} \sum_{j=1}^J (\theta_j - \mu)^2\right\}$$

$$\tau^2 | \theta, \mu, \sigma^2, y \sim \text{Inv-}\chi^2(J-1, \hat{\tau}^2) \quad \hat{\tau}^2 = \frac{1}{J-1} \sum_{j=1}^J (\theta_j - \mu)^2$$

Metropolis-Hastings Satissg. detailed balance / Transition

density $\psi(y|x) = \psi(\theta^* | \theta^{*-1}) \propto \text{Norm } q(\theta^* | \theta^{*-1}) = \psi(y|\theta)$

$$\psi(y|x) \pi(x) \stackrel{\text{D.B. detailed balance}}{=} \psi(x|y) \pi(y), \quad \{ \text{except probab } p^* = \min(1, r) \}$$

$$\pi(y|x) = p_{\theta=\theta^{*-1}}(y|x) = \int p(y, \theta^* | x) d\theta^* = \int p(y | \theta^*, x) \underbrace{q(\theta^* | x)}_r d\theta^*$$

$$= p^* \delta(y - \theta^*) + (1 - p^*) \delta(y - x) \quad p^* := \begin{cases} \theta^* | x, \theta^* > 1 \\ 1 - \theta^* | x, \theta^* \leq 1 \end{cases} \quad \begin{cases} \theta^* \in D_1(x) \\ r(x, \theta^*) \theta^* \in D_0(x) \end{cases}$$

$$\text{with } D_1(x) = \{ \theta^* : r(x, \theta^*) > 1 \} \quad D_0(x) = \{ \theta^* : r(x, \theta^*) \leq 1 \} \quad \xrightarrow{\text{Posterior}}$$

$$\text{to satisfy detailed balance } \psi(y|x) \pi(x) = \psi(x|y) \pi(y)$$

$$\psi(y|x) \pi(x) = \int (p^* \delta(y - \theta^*) + (1 - p^*) \delta(y - x)) q(\theta^* | x) \pi(x) d\theta^* =$$

$$\int_{D_0(x)}^{} \{r(x, \theta^*)\delta(y - \theta^*) + (1 - r(x, \theta^*))\delta(y - x)\} q(\theta^* | x) \pi(x) d\theta^* +$$

$$\int_{D_1(x)}^{} \delta(y - \theta^*) q(\theta^* | x) \pi(x) d\theta^* =$$

$$\int_{D_0(x)}^{} r(x, \theta^*) \delta(y - \theta^*) q(\theta^* | x) \pi(x) d\theta^* + \int_{D_0(x)}^{} (1 - r(x, \theta^*)) \delta(y - x) q(\theta^* | x) \pi(x) d\theta^*$$

$$\begin{aligned} & q(y|x) \pi(x) \mathbb{I}(y \in D_1(x)) = \\ & r(x, y) q(y|x) \pi(x) \mathbb{I}(y \in D_0(x)) + \delta(y - x) \left[\int_{D_0(x)}^{} q(\theta^* | x) \pi(x) d\theta^* - \right. \\ & \left. \int_{D_0(x)}^{} r(x, \theta^*) q(\theta^* | x) \pi(x) d\theta^* \right] + q(y|x) \pi(x) \mathbb{I}(y \in D_1(x)) \end{aligned}$$

$$\left\{ r(x, y) = \frac{q(x|y) \pi(y)}{q(y|x) \pi(x)} \right\} \text{ by }$$

$$= q(y|x) \pi(x) = q(x|y) \pi(y) \mathbb{I}(y \in D_0(x)) +$$

$$\delta(y - x) \left[\int_{D_0(x)}^{} q(\theta^* | x) \pi(x) d\theta^* - \frac{\int_{D_0(x)}^{} q(x|\theta^*) \pi(\theta^*) q(\theta^* | x) \pi(x) d\theta^*}{\int_{D_0(x)}^{} q(\theta^* | x) \pi(x) d\theta^*} \right] +$$

$$q(y|x) \pi(x) \mathbb{I}(y \in D_1(x)) \quad \left\{ \begin{array}{l} \vdots = q(x|y) \pi(y) \\ \text{2nd term} \end{array} \right.$$

clearly symmetric in x, y .

$$q(x|y) \pi(y) : \quad \mathbb{I}(y \in D_0(x))$$

$$\mathbb{I}(y \in D_1(x))$$

$$y \in D_1(x) \Leftrightarrow r(x, y) > 1 \Leftrightarrow \frac{1}{r(y, x)} > 1$$

$$\Leftrightarrow r(y, x) \leq 1 \quad x \in D_0(y) \quad y \in D_0(x) \Leftrightarrow x \in D_1(y)$$

$\therefore q(x|y) \pi(y) = q(y|x) \pi(x)$ \square Metropolis Hastings has Z equilibrium
distributions Z posterior distri

\Week 8/ Stan for Z NHM/ $y_{ij} | \theta_j, \sigma^2 \sim N(\theta_j, \sigma^2) \quad j=1, \dots, J$

$$\theta_j | \mu, \tau^2 \sim N(\mu, \tau^2) \quad \pi(\mu, \tau^2, \sigma^2)$$

$$(x_i, y_{ij}) \quad \underbrace{y_{ij}}_{x} \quad y_{ij} \sim N(a + b x_i, \sigma^2) \quad y_{ij} | x_i, a, b, \sigma^2 \sim N(a + b x_i, \sigma^2)$$

default priors uniform priors $\pi(a, b) \propto 1 \quad \pi(\sigma^2) \propto \frac{1}{\sigma^2}$ but

with Stan: $a \sim ? \quad b \sim ? \quad \sigma^2 \sim ?$ what ever we want

For instead: $y_{ij} | x_i, \beta_i, \sigma^2 \sim N(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}, \sigma^2)$ is Bayesian

Likelihood $\therefore a \sim \beta_0 \quad b \sim \beta_1 \quad \therefore p \sim ? \quad \& \sigma^2 \sim ?$

$\Rightarrow \text{rstan}; \text{extract}(\dots) \quad \% .do \text{ rstan}; \text{extract} \quad \text{to tell R which package}$

extract lives in since it doesn't know which to choose %

$\text{DATA} y_{ij} | \theta_j, \sigma^2, \mu, \tau^2 \sim N(\theta_j, \sigma^2)$
anything anything what ever params need in this distri
process layer in brms must always be $\theta_j \sim N(0, \tau^2)$, $\pi(\tau^2, \sigma^2, \text{prior})$
 $\pi(\theta_j \sim N(0, \tau^2))$ with prior $\pi(\tau^2, \sigma^2, \text{prior})$

$\therefore \text{model: } y_{ij} \sim N(\mu + \theta_j, \sigma^2) \quad \theta_j \sim N(0, \tau^2) \quad \pi(\mu, \tau^2, \sigma^2)$
any thing after μ is one of those population params, so b are
params in Z data layer mean only, sigma are Z params in Z data
layer that are not in Z mean, Sd's are standard deviation of Z
group effects which have mean zero

$\gg \text{plot(rats_z)}$ $\gg \text{plot(rats_z, pars="diet")}$ \therefore specify Z params by Z
group, get plot of Z different group effects
 $\gg \text{DietEffects} \leftarrow \text{ranes(rats_z, summary=TRUE)}$ \therefore gives mean estis of
each diet
 $\gg \text{mains} \leftarrow \text{summary(rats_z)}$ $\gg \text{mains\$fixed}$
 $\gg \text{mains\$fixed[, "Estimate"]} + \text{DietEffects\$diet[, "Estimate", 1]}$ \therefore gives
Z means had before and can do similar with samples.

can set for Z params of Z data layer not in Z mean in this case &
(so $\mu \in Z$ are in Z mean?) \therefore set a prior

choosing class SigmaFit knows im fitting a gaussian Model, it
auto limit it so not allowed less than zero eg choose $N(0, 100)$ it's
only half normal since truncated at 0

for covariabel group.variable to Z right, put variable where have
groups. So variable diet, different level Δ like to fit our Model,
using, differently for each of those groups - invokes Z
hierarchical Model. Z params in x-variable, like intercept
+ params each get Z normal distri for each group with Mean Zero &
Some Common variance, common var has a prior can set via 1)
Z "Sd" class & there becomes a global mean term for x-variable,
that's estied & takes role of our μ in Z normal hierarchical model

$$y_i | \beta, \sigma^2 \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$y_i | \beta, \sigma^2 \sim N(\beta_0 + \beta_1 x_i, \sigma^2) \leftarrow$ Data layer

• prior layer $\pi(\beta, \sigma^2)$

① first brm model: intercept 24 s, $\sigma^2 = 51$

each day increase by 11

reaction time longer, more sleep deprived, all by different amounts

Δ starts from different place, good sign since don't want single

regression ∵ don't have a bunch of pts random about a line, seems like we have a bunch of lines

$$j=1, \dots, J \quad \text{Data layer: } y_{ij} | \beta_j, x_{ij}, \sigma^2, \beta_j \sim N(\beta_{0j} + \beta_{1j} x_{ij}, \sigma^2)$$

β 's intercept & slope are from each group ∵

Process Layer: $\beta | B, \Sigma_B \sim MVN(B, \Sigma_B)$ $B = (\beta_0, \beta_1)$
↑ b_{0j} has mean zero in process layer

Prior layer: $\pi(B, \Sigma_B, \sigma^2)$

brm version: $y_{ij} | \beta_j, x_{ij}, \sigma^2, B, \Sigma_B \sim N(\beta_{0j} + \beta_{1j} x_{ij}, \sigma^2)$

∴ B is zero mean where B is, b need to come in mean ∵ b

group level deviation from it, β_{0j} has mean zero in b_{0j} layer

global b regression term for x_{ij} ∵ b (general days for a subject)

$\beta_{1j} x_{ij}$ is (group level) process level deviation from this line

$\beta_{1j} x_{ij}$ prior: $\pi(\Sigma_B, \sigma^2, B) \quad B = \begin{pmatrix} b_0 \\ b_1 \end{pmatrix}$

∴ $\beta | \Sigma_B \sim MVN(0, \Sigma_B)$ prior: $\pi(\Sigma_B, \sigma^2, B) \quad B = \begin{pmatrix} b_0 \\ b_1 \end{pmatrix}$

it's, say all regression params ∵ just regressed on days, all have a group level

are halving sd says 2 std deviation around 2 thing fitted is actually

at smaller Δ is explainable by something else, e.g. like Subject

variability population level effects haven't change

now have group level effects, giving Sigma (Σ_B), getting sd of

intercept, sd of days is bottom right diagonal. Correlation is there too

• all tested well

>> plot(sit-sleep, N=3) % plot limited plots!

Sd params are group level params Cor-Subject (Correlation in that matrix)

$\gg \text{plot}(\text{sit-sleep}, N=7, \text{par5} = \text{"Subject"})$ same as before & get Γ 's
 one for each subject & get them for slopes of 7 days is 2
 Slope variable: see coefs for each subject
 $\gg \text{range}$ > random effects for us isn't random: randomness is in
 2 params).

for intercept, seeing different subjects

$\gg \text{b-prior}$ > use default $N(0, 10)$ default for regression params.
 b-mu has population params are class B apart from intercept (intercept
 does get this prior too), class Sd for group level params

final prior on Sigma

Week 12 tutorial

Ch 2 / intro to Bayesian machinery

By Bayes $\pi(\theta|y) \propto \pi(\theta) p(y|\theta)$
 avoiding 2 marginal likelihood (not predictive) $P(y) = \int_{-\infty}^{\infty} p(y|\theta) \pi(\theta) d\theta$

\therefore Some is proportional to something we know (Normal, Gamma, Beta, Binomial)

Conjugate: Form $\pi(\theta) \propto \pi(\theta|y)$

Bayesian inference / Estimation: $E[g(\theta)|y] = \int_{-\infty}^{\infty} g(\theta) \pi(\theta|y) d\theta$

use MAP when sampling from $\pi(\theta|y)$ hard: arg $\max_y \pi(\theta|y)$

Bayesian Credible intervals: $[a, b] = \int_a^b \pi(\theta|y) d\theta = P(a \leq \theta \leq b) = \alpha \quad (100\%)$

credible interval (not confidence interval)

prediction $P(Y_{\text{new}}|y) = \int_0^{\infty} P(Y_{\text{new}}|\theta) \pi(\theta|y) d\theta$
 \propto likelihood

subjective priors: Goals improper prior. $\int \pi(\theta) d\theta \neq 1$

uniformative prior $\pi(\theta) \propto 1 \quad \theta \in [c, \infty]$ improper

Jessoy's prior $\pi(\theta) \propto f(\theta)^{-1}$

subjective probability theory: Event random quantity duality

Event is TRUE/FALSE. Can be treated as 1/0 is random because we

don't know if 1/0 $A \vee B$ 'V' or $A \wedge B$ 'W' and

$\tilde{A} \sim \text{not } A \quad x \vee y = \max(x, y) \quad x \wedge y = \min(x, y) \quad \sim x = 1 - x$

$A \wedge B = AB \quad (min(A, B)) \quad A \vee B = \sim(\tilde{A} \wedge \tilde{B}) = 1 - (1 - A)(1 - B)$

prove $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

- Week 9 videos / Take 2 hierarchical model checking
- temperature took less 6 vals \therefore literature looks at how you treat it
- that as an ordered factor, like recipe A, B, C are factors but in no natural ordering, but $175 < 185 \therefore$ Should come into your model somewhere
- 15 replicates \therefore 15 different experiments
- ms/ don't use $\gg (\text{temp} | \text{recipe} * \text{replicate})$ \because since in R its
 $\text{recipe} * \text{replicate} + (\text{recipe})(\text{replicate})$
So use $\gg (\text{temp} | \text{recipe}: \text{replicate})$ \therefore is correct in R its
 $(\text{recipe})(\text{replicate})$
- 45 combinations of replicate & recipe
(x.variable | group) is x variable changes with Z group
45 levels $\therefore \gg \text{recipe}: \text{replicate}$ has 45 possible vals of this
- a) grouping variable
want to take out group level effects A1S, B1S, C1S then predict
replicate 15 from Z fit of Z old data
can plot 95% credible intervals so sine to have 1 pt out of 15
b) outside Z error bars \therefore expect 5%

Find good model by write down & justifying subjectively, write down priors can justify subjectively. Show through these plots, what you've done makes sense for data

$$y_i | \beta, \sigma^2 \sim N\left(\sum_{j=0}^p \beta_j X_{ij}, \sigma^2\right) \quad j_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} \sim N(\gamma_i, \sigma^2)$$

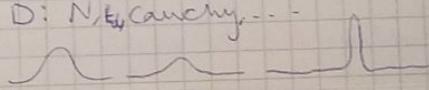
$$y_i | \beta, \theta \sim D(g^{-1}(\gamma_i), \theta) \quad D \text{ has mean } g^{-1}(\gamma_i)$$

$$(G(\alpha, \beta) \quad E[G] = \frac{\alpha}{\beta}, G(\frac{\alpha}{\beta}, \theta(\alpha, \beta)) \}$$

$g(\cdot)$ is called 'link function' $\gamma \in (-\infty, \infty) \therefore$ range of g lists map $\gamma \rightarrow$ support of D via g^{-1} $y(x) = \log(x)$ is a popular link

when $E[y]$ must be > 0

y_i : continuous support: $(-\infty, \infty)$ $D: N, t_k, \text{Cauchy}, \dots$



y : Continuous Support: $(0, \infty)$ D: G, lognormal
 Support: $[a, b]$ D: Beta, scaled beta

y : Discrete support $\{0, 1\}$ D: Bernoulli
 $\{0, 1, \dots, N\}$ D: Binomial
 $\{0, 1, \dots\}$ D: Poisson, Negative Binomial

y : Rain $\{0, R\}$ D: Hurdle model, Zero-inflated Models
 (R is continuous > 0)

choice of D is called Z family, default is family = gaussian
 \Rightarrow ? brms family

$$\begin{aligned} j_{ij} | J_{ij}, \theta &\sim D(g^{-1}(J_{ij}), \theta) \quad j = 1, \dots, J \\ \text{Data layer: } J_{ij} &= b_0 + b_1 x_{1j} + \dots + b_p x_{pj} + \beta_{0j} + \beta_{1j} x_{1j} + \dots + \beta_{pj} x_{pj} \end{aligned}$$

$g(\cdot)$ known

process layer: $\beta_j | \Sigma \sim MVN(0, \Sigma)$

prior: $\pi(b, \Sigma, \theta)$

\Rightarrow ? brms-package

\Rightarrow brm(formula, data, family=D, prior=?) // if set prior=NULL then uses default priors that depend on formula & family //.

Should always be specifying own priors //.

Normal family is gaussian, link "identity" is only one that makes sense

once got data find what type it is (eg count, binary, positive, normal, continuous) then make choice of D (distribution wanted?)

Set or this tailed, poisson, negative binomial for counts) then link function choices default or somewhere else, then what we

Z other params then which predictors to put in, which

variables to group on so were doing $\beta_j | I \sim MVN(0, \Sigma)$, what's prior on group variables, predictions, other params.)

done MCMC, give samples from posterior, gives Monte Carlo estis of predictive distns or of params for inference

prove $P(A \vee B) = P(A) + P(B) - P(A \cap B)$

$y_i \in \{0, 1\}$

$g: [0, 1]$

logistic R

Week 2 v

Model Str

$y_j | f_{yj}(\theta_j)$

$\theta_j \sim MVN$

prior dis

$\beta_j \sim MVN(0, \Sigma)$

Model

Data (e

predicti

poste

doing n

discre

o In Ser

then s

Model

(X, Y)

x, y

(Days)

j = {1

2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31}

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31

$y_i \in \{0, 1\}$ ~ Bernoulli $\sim \text{Bern}(p = \frac{\text{exp}(\beta_0 + \beta_1 x_i)}{1 + \text{exp}(\beta_0 + \beta_1 x_i)})$

$$y_i \in [0, 1] \text{ or } (0, 1) \quad g = \text{logit}(p) = \frac{E[y_i]}{1 - E[y_i]} = \beta_0 + \beta_1 x_i$$

Logistic Regression

Watch 2 videos / good Bayesian analysis: write down & justify for

model: Model structure & priors (30s-40s)

$$\text{prior: } \beta_0, \beta_1, \sigma^2 \sim N(\beta_0, \beta_1, \sigma^2)$$

$$\beta_j \sim \text{MVN}(\mathbf{B}, \mathbf{I})$$

$$\text{prior distri: } \pi(\mathbf{B}, \mathbf{I}, \sigma^2)$$

$$\mathbf{B} \sim \text{IW}(100, 0.01)$$

• Convergence (10-20%)

• Model Checking: posterior predictive plots for (hopefully withheld) data (out of sample - is in-sample justified since not much data). Are prediction intervals sensible etc

then • posterior densities param plots, are they sensible, is Z param doing anything, is so can we take it out & test. Is Z posterior different to Z prior (at param level)



• Inference: $E[\hat{y}_j | \mathbf{y}] \in [\hat{y}] \in [\hat{y}] \in [\hat{y}] \in [\hat{y}]$

• Show posterior densities & comment

then sensitivity analysis

• Modelling: $y_i = \text{Total}_i \quad i=1, \dots, 30 \quad y_i | \mathbf{x}_i \sim \text{Poisson}(\lambda(\mathbf{x}_i))$

$$(\lambda \text{ inverse link}) \quad \log \lambda(\mathbf{x}_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

x_1 high-temp x_2 low-temp, x_3 precip x_4 snowday

(Days Grouping?)

$j = \{1, 2\}$ 1 - weekday, 2 - weekend $n_1 = 9, n_2 = 21$

$$y_{i,j} | \beta_j, \mathbf{x}_{ij} \sim \text{Poisson}(\lambda(\beta_j, \mathbf{x}_{ij})) \quad \log \lambda(\beta_j, \mathbf{x}_{ij}) = \beta_{0j} + \beta_{1j} x_{1,i,j} + \beta_{2j} x_{2,i,j} + \beta_{3j} x_{3,i,j} + \beta_{4j} x_{4,i,j}$$

$$\beta_{0j} + \beta_{1j} x_{HT,j} + \beta_{2j} x_{LT,j} + \beta_{3j} x_{P,j} + \beta_{4j} x_{S,j}$$

$$\beta_j \sim N(\mathbf{B}, \mathbf{I}) \quad \text{prior: } \pi(\mathbf{B}, \Sigma)$$

r-group[... Intercept]

Estimate sd q2.5

r-group[... , someX]

rane\\$[c(4,56), 'Estimate']

{> strsplit}

Same priors & grouping is collusion & utterly impossible without it
(& same chains)

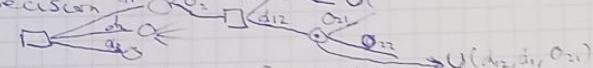
Solve week 10 cyclists data for high marks when doing coursework

Menti/11/ which is not a necessary ingredient of a decision problem?

Data

Q2/ on a decision tree, what shape are the decision nodes? Squares

Decision node: d_1, d_2, d_3 at $O_1: U(d_1, A_1, O_1, O_2)$ $E[U(O_2 | d_1, A_1)]$



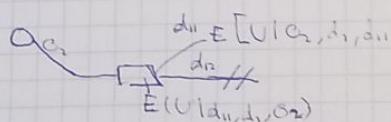
Q3/ chance nodes are? Circles

Q4/ starting at the RHS of the tree, when is a chance node must:

Mark with the conditional expectation of the leaves given the branch you are on

Q4/ starting at the RHS of the tree, when is a decision node must:

Mark with the maximum utility to the right



Q5/ a utility func maps rewards (of any kind) to the real line & has

2 properties. What are they? 2 magnitude of utility respects

2 presence ordering 2 utility of a gamble is equal to the

expected utility

Q6/ a utility func always exists? TRUE

Q7/ how unique is your utility func? it is unique up to positive

linear transformation (least $U=0$ $U(best)=1$)

Let's do some for many $U(x)$'s. What's happening
 $E[\dots]$, if we have n coins & generate ID_1 's heads/tails $E[\dots]$?
 Because of $E[U(\text{ID}_1)] = 50.5$ $E[U(\text{ID}_2)] = 50$
 Let's do it for $n=2$ events x : $(x_1, x_2) \sim \text{Uniform}(0, 1)$
 $C_1 = 0.1(3 + 0.75)$ with $r_1 = 3$, $r_2 = 10$, $r_3 = 15$ C_1 or C_2 , C_3 ?

Week 10 video / Modelling: $\ln(\lambda) = \ln(b) + b_0 + b_1 x_{HT, i} + b_2 x_{HT, j} + \dots$

$F_i \sim N(0, \Sigma)$

need prior for $\pi(\lambda, \Sigma)$

$$E[\lambda] = \lambda = e^{f_0 + b_1 x_{HT, i} + \dots + b_n x_{HT, n}} = e^{\log \lambda + b_1 x_{HT, i} + \dots + b_n x_{HT, n}} \dots$$

$$\begin{aligned} e^{\log \lambda} &\in [5000, 20000] \\ N(0, 5) &\quad (\text{majorise } N(4, 3) \text{ for sensitivity analysis but } \Sigma \text{ and } \lambda) \\ b_i &\sim N(0, 0.05) \quad \{ \gg \exp(0.17C(35, 81)) \Rightarrow 33, 32, 24 \dots \} \\ N(0, 0.05) &\quad \{ \end{aligned}$$

$$\begin{aligned} \Sigma &= \text{vec}(\Sigma) \times \Sigma_{kk} \times \text{vec}(\Sigma)^\top \\ &\Rightarrow \text{setpred}(\text{class} = "SD") \quad t_3(0, 5) \\ N(0, 0.05) &\quad \text{non-intercept SDs} \end{aligned}$$

Agreement / Convergence: precipitation amount is two orders of magnitude lower : expect to start at 2 orders of magnitude higher
 $\therefore N(0, 5)$

- Run up to 30 seconds per chain { very bad Rhat #1: don't use }
- Always marks for trace plots to check can interpret { trace plots are starting in different places & getting stuck } - hasn't converged
- 4 chains not converging mean also longer to have to warm up :
- allow scan to run for longer $\gg \text{iter} = 10000, \text{warmup} = 8000$
- { 3 runs per chain } 10000 samples, 2000 runs per chain .
- Started in different places & not converged
- maybe starting in bad place? { unless present Rhat's, they're bad, }

non-existent ESS}

Missing everywhere :: poisson has same Mean & Var :: eq at Mean
25000, 2 variances 25000

$y|\beta, x \sim D(g^{-1}(1), \phi)$ {never seen a poisson regression work}
Mean priors

$$y|\beta, x \sim \text{NegBin}(\lambda(\beta, x), \phi) \quad p(y=k|\lambda, \phi) \propto \left(\frac{\lambda}{\phi+\lambda}\right)^k \left(\frac{\phi}{\phi+\lambda}\right)^\phi$$

\rightarrow Poisson as $\phi \rightarrow \infty$

$$E[y] = \lambda \geq 0 \therefore \text{can still use } \log(\lambda|x, \beta|) = \sum_{i=0}^p \beta_i x_i$$

$$\text{var}[y] \neq \lambda \therefore \text{var}[y] = \lambda(1 + \frac{1}{\phi}) \text{ for } \phi \geq 0$$

$$\pi(\beta, \Sigma) = \pi(\beta, \sigma) \therefore \text{nonneed: } \pi(\beta, \sigma, \phi)$$

$\gg \text{sqrt}(\text{varys}(10)) \Rightarrow 1600, 7000 \therefore$ is 60 error is really bad

$\gg \text{sqrt}(\text{varys}(100)) \therefore$ gives 20% error not too bad %

$\gg \text{sqrt}(\text{varys}(600)) \therefore$ error bars getting tight is 5000 ± 400 %

\therefore Set mean=200, var=300 eg $200+3*300=800 \therefore$ set as prior

\therefore trace plots look better, signs things have converged, Rhads good

\therefore very good argument for convergence, trace plots excellent,

no low ESS's \therefore can do MonteCarlo & inference, is Z model is ok

\gg preds \therefore error bars are pretty big, for lower vals Model is fairly

skillful, everythings hitting Z error bars \therefore model is ok, only one

missing pt %.

Model checking: error bars too big :: decrease var of shape prior

{low temp is mean zero & unsure of sign \therefore chance low temp is -

effectively not a thing \therefore back into taking it out potentially

get rid of snowday: clearly not doing anything Σ hasn't moved from Z

prior {in regression & Bayesian regression, shouldn't take out two

variables at one time, could be they're perfectly correlated, Σ is perfect

was spread b/w them} {could put strong prior on snow eg mean 1)

negative, say people don't cycle in it?

P is \gg MCMC-pairs shows high correlation, its saying high-temp &

Low-temp are doing Z same sorts of thing & would be a good argument to kick one of them out. First kick out snow day then look again at this distri for low-temp & see if there's any correlation.

Now it's converged don't need such a hard warming, decrease it.
∴ $\text{ssitter} = 6000$, $\text{warmup} = 5000$

∴ only one pt outside. Model believable, still bit big uncertainty gives probab got sign wrong from Z expectation, st it has a negative on cyclists is 0.188 ∴ fairly sure of it

but for low-temp its 0.387 probab its Z other sign

∴ low-temp does seem to have an effect just a high chance of removing it

{Focus on Z biggest change (just make one) & demonstrate Z effect it made (staying sensible)}

>> mcmc (third_model) ∴ precip effect looks tiny

∴ only 1 pt missing ∴ convertible with it

>> summary(Sens_Model) ∵ Rhats are fine, were converged, a much smaller model

conclusions insensitive to model choices {document decisions made}
Sens Model then change something to check sensitivity

∴ Sens Model is Sorth Model

Inference: ∴ prob that increase in High-Temp leads to more cyclists is 1 - 0.169 ∴ 83.1%.

definitely reduce cyclists by precip increase

{learn how to plot random effects

to predict probab of an event ∴ need to do a posterior prediction,

comes from Z samples, & comes directly from Z pred func
eg >> predict(...) ∴ give it Z new data you want, it gives Expedtion as Z probab, Z standard deviation of Z samples (Est.Error) ∴ need