

Lecture 1.pdf

Lecture 2.pdf

Lecture 3.pdf

Lecture 4.pdf

Lecture 5.pdf

Lecture 6.pdf

Lecture 7.pdf

Lecture 8.pdf

Lecture 9.pdf

Lecture 10.pdf

Lecture 11.pdf

Lecture 12.pdf

Lecture 13.pdf

Lecture 14.pdf

Lecture 15.pdf

Lecture 16.pdf

Lecture 17.pdf

Lecture 18.pdf

Lecture 19.pdf

Lecture 20.pdf

Lecture 21.pdf

Lecture 22.pdf

Conditional probabilities:

Given two events A and B in a probability space Ω we define:

$P(A \cap B)$ = probability that A and B happen simultaneously.
= joint probability of A and B

$P(A \cap B) = 0$ mutually exclusive events.

$P(A)$ = marginal prob. of A .

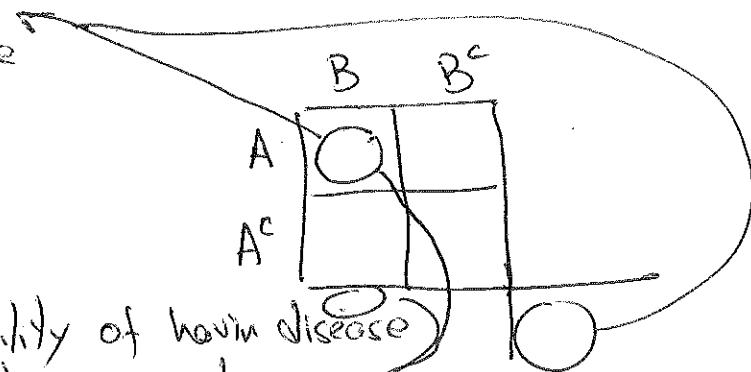
$P(B)$ = marginal prob. of B .

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

A = an individual randomly sampled has cancer

B = an individual randomly sample is positive in
a test for that cancer.

$P(A \cap B)$ = an individual has cancer and is positive
in the disease



$P(A|B)$ = probability of having disease
if test is positive

(2)

$$P(A \cap B) = P(B \cap A)$$

$$P(A|B) \neq P(B|A)$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad \text{Bayes theorem.}$$

Partitions: $\{H_1, H_2, \dots, H_k\}$ set of events such that

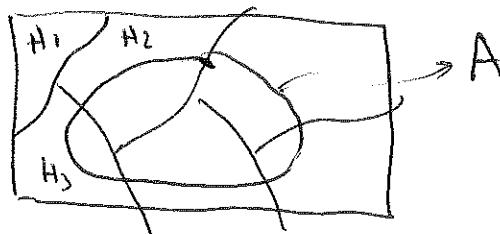
$$H_i \cap H_j = \emptyset$$

$$H_1 \cup H_2 \dots \cup H_k = \Omega$$

$$P(H_i | A) = \frac{P(A|H_i)P(H_i)}{\sum_{j=1}^k P(A|H_j)P(H_j)}$$

$P(A)$

Total probability law $P(A) = \sum_{j=1}^k P(A|H_j)P(H_j)$



Example: Data on incomes and educational level ③
for a sample of 30 year old males

$\{H_1, H_2, H_3, H_4\}$ events that correspond to a randomly sampled individual being in the 1st, 2nd, 3rd or 4th quartile of the income distribution

$$P(H_1) = \frac{1}{4}$$

$\{H_1, H_2, H_3, H_4\}$ represents a partition of the space

From the survey data we have

$$Pr(E|H_1) = 0.11$$

$$Pr(E|H_2) = 0.19$$

$$Pr(E|H_3) = 0.31$$

$$Pr(E|H_4) = 0.53$$

E: an individual has a college education

$$(1) \quad Pr(H_4|E) = \frac{\cancel{Pr(H_1|E)} P(H_4) Pr(E|H_4)}{\cancel{Pr(E)}}.$$

$$(2) \quad P(E) = Pr(E|H_1)P(H_1) + Pr(E|H_2)P(H_2) + \\ Pr(E|H_3)P(H_3) + Pr(E|H_4)P(H_4)$$

$$(1) + (2) \Rightarrow Pr(H_4|E) = 0.47.$$

(4)

Example: You have two operators for a machine that produces widgets. The number of widgets produced during a day follows a Poisson distribution, with a mean that depends on the operator. The mean for operator 1 is 10.5 and the mean for operator 2 is 7.8.

We collect one day's worth of widgets and want to determine which operator was working that day.

Call X = observation.

$$X \sim \text{Poisson}(\lambda)$$

$$p(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$p(X=x \mid \text{operator 1}) = \frac{e^{-10.5} (10.5)^x}{x!}$$

$$p(X=x \mid \text{operator 2}) = \frac{e^{-7.8} (7.8)^x}{x!}$$

$$Pr(\text{operator 1} | X=x) = \frac{Pr(X=x | \text{Op 1}) Pr(\text{Op 1})}{Pr(X=x | \text{Op 1}) Pr(\text{Op 1}) + Pr(X=x | \text{Op 2}) Pr(\text{Op 2})}$$

We need $Pr(\text{Op 1})$ ($Pr(\text{Op 2}) = 1 - Pr(\text{Op 1})$)

Option 1: I don't know so set $Pr(\text{Op 1}) = \frac{1}{2} = Pr(\text{Op 2})$

In that case

$$\Pr(\text{Op } 1 | \bar{X} = x) = \frac{\Pr(\bar{X} = x | \text{Op } 1)}{\Pr(\bar{X} = x | \text{Op } 1) + \Pr(\bar{X} = x | \text{Op } 2)} \quad (5)$$

\Rightarrow I will say that the operator that was in charge is the one that maximizes

$$\Pr(\bar{X} = x | \text{Op})$$

\hookrightarrow Maximum likelihood estimation (in its simplest form).

Option 2: If you know that one works part time and the other full time then:

$$\Pr(\text{Op } 1) = \frac{1}{3} \text{ (part time)}$$

$$\Pr(\text{Op } 2) = \frac{2}{3} \text{ (full time)}$$

$$\Pr(\text{Op } 1 | \bar{X} = x) = \frac{\Pr(\bar{X} = x | \text{Op } 1)}{\Pr(\bar{X} = x | \text{Op } 1) + 2\Pr(\bar{X} = x | \text{Op } 2)}$$

Independence

(b)

A and B are independent if

$$\Pr(A|B) = \Pr(A)$$

$$\hookrightarrow \Pr(A \cap B) = \Pr(A)\Pr(B)$$

Iterating conditionings

$$\begin{aligned}\Pr(A \cap B \cap C) &= \Pr(B | A \cap C) \Pr(A \cap C) \\ &= \Pr(B | A \cap C) \Pr(A | C) \Pr(C)\end{aligned}$$

Random variables:

Unknown numerical quantities about which we make probability statements.

Discrete random variables (countable outcome spaces)
Continuous random variable (\mathbb{N} countable outcome spaces)
 \mathbb{R}

Discrete random variables

(pmf) probability mass function

$$f(x) = \Pr(X=x) = F(x) - F(x-1).$$

(cdf) cumulative distribution function $\Pr(X \leq x) = F(x)$

(7)

- Continuous random variables
- (pdf) probability ~~density~~ function $f(x) = \frac{d}{dx} F(x)$
- (cdf) cumulative distribution function. $F(x) = \Pr(X \leq x)$

Summaries of random variables.

$$E(X) = \begin{cases} \sum_{x \in \Omega} x p(x) \\ \int_{\Omega} x p(x) dx \end{cases}$$

$$\text{Var}(X) = \begin{cases} \sum_{x \in \Omega} (x - E(X))^2 p(x) \\ \int_{\Omega} (x - E(X))^2 p(x) dx \end{cases}$$

Some random do not have finite means or variances !!

Bernoulli Trials:

D

Experiments:

- 1) Two possible outcomes (S/F, 0/1)
- 2) Repetcs are independent
- 3) Probability of outcomes is constant over time.

Bernoulli: $X_i \sim \text{Ber}(\theta) \quad X_i \in \{0,1\}$

$$P(X| \theta) = P(X) = \theta^x (1-\theta)^{1-x} \quad \theta = \text{success probability}$$

$$= \begin{cases} \theta & \text{if } X=1 \\ 1-\theta & \text{if } X=0 \end{cases}$$

Binomial: $X_1, X_2, \dots, X_n \sim \text{Ber}(\theta)$ and independent

$Z = X_1 + X_2 + \dots + X_n = \# \text{ of successes in } n \text{ trials}$

$$Z \sim \text{Bin}(n, \theta)$$

$$P(Z|n, \theta) = \binom{n}{z} \theta^z (1-\theta)^{n-z} \quad z = 0, 1, \dots, n$$

$$E(Z) = n\theta$$

$$\text{Var}(Z) = n\theta(1-\theta)$$

(2)

Geometric:

$Y = \# \text{ of trials before a success is observed in a Bernoulli experiment with success prob } \theta$

$$P(Y|θ) = θ(1-θ)^{Y-1} \quad Y = 1, 2, 3, \dots$$

Alternatively, define $W = \# \text{ of failures before the first success}$

$$Y = W + 1$$

$$P(W|θ) = θ(1-θ)^W \quad W = 0, 1, \dots$$

$$E(Y) = \frac{1}{θ}$$

$$\text{Var}(Y) = \frac{1-θ}{θ^2}$$

Negative Binomial

$$Y_1, Y_2, \dots, Y_K \sim \text{Geo}(θ)$$

$X = Y_1 + Y_2 + \dots + Y_K$ = number of trials to observe K successes

$$P(X|θ, K) = \binom{X-1}{K-1} θ^K (1-θ)^{X-K}$$

opposite to the
Binomial!!!

$$E(X) = \frac{K}{θ}$$

$$\text{Var}(X) = \frac{K(1-θ)}{θ^2}$$

Hyper-geometric (sampling uniformly without replacement) ③

$X = \#$ of blue balls in a sample of size n drawn from an urn that contains K blue balls and $N-K$ red balls.

$$P(X|n, K, N) = \frac{\binom{K}{X} \binom{N-K}{n-X}}{\binom{N}{n}}$$

$$E(X) = \cancel{\frac{Kn}{n}} \quad \frac{Kn}{N} \approx \begin{matrix} \# \text{ of balls in } x \text{ picks of blue ball} \\ \text{sample} \end{matrix} \quad \text{check!!!}$$

$$\text{Var}(X) = \cancel{\frac{Kn}{n^2}} \quad \cancel{\frac{N(N-1)}{N^2}} \quad \frac{Kn}{N} \quad \frac{(N-n)(N-K)}{N(N-1)}$$

Poisson: $X = \text{count of events}$

$$P(X|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad x=0, 1, \dots$$

Poisson process:

- Consider the number of events (phone calls) that you receive during a given hour
- Disjoint intervals are independent
- Cells arise independently
- The probability of two or more calls in a very short time period is very small
- The probability of getting exactly one call in a small interval is constant

If these 3 assumptions are satisfied then:

X = number of calls in one hour

$X \sim \text{Poisson}(\lambda)$ λ = mean number of calls in one hour.

$$E(X) = \lambda$$

$$\text{Var}(X) = \lambda$$

In practice, the fact that $E(X) = \text{Var}(X)$ is quite restrictive. It is often the case that

$\text{Var}(X) > E(X)$ overdispersion.

↳ If this is the case the Negative Binomial might be a better option.

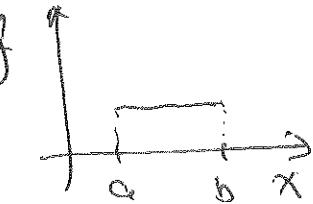
for modeling even if the data generation mechanism has nothing to do with Bernoulli trials!!!

(5)

Continuous random variables:

Uniform: $X \sim \text{Uni}[a, b]$

$$f(x|a,b) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$



$$E(X) = \frac{a+b}{2}$$

"All values have the same likelihood of occurring"

$$\text{Var}(X) = \frac{(b-a)^2}{12}$$

Exponential: interarrival times in a Poisson process

$$X \sim \text{Exp}(\mu)$$

$$f(x|\mu) = \begin{cases} \frac{1}{\mu} \exp\left\{-\frac{x}{\mu}\right\} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$E(X) = \mu$$

$$\text{Var}(X) = \mu^2$$

$$\begin{aligned} F(x|\mu) &= \int_0^x \frac{1}{\mu} \exp\left\{-\frac{t}{\mu}\right\} dt \\ &= 1 - \exp\left\{-\frac{x}{\mu}\right\} = \Pr(X \leq x) \end{aligned}$$

~~$P(X > t+s | X > s)$~~

$$P(X > t+s | X > s) = \frac{P(X > t+s, X > s)}{P(X > s)} = \frac{\Pr(X > t+s)}{\Pr(X > s)} = \exp\left\{-\frac{t+s}{\mu}\right\}$$

Gamma: $X_1, X_2, \dots, X_n \sim \text{Exp}(\mu)$

(6)

$Y = X_1 + X_2 + \dots + X_n \sim \text{Gamma}(n, \mu)$

$$f(y|n, \mu) = \begin{cases} \frac{1}{\Gamma(n)} \left(\frac{1}{\mu}\right)^n y^{n-1} \exp\left\{-\frac{y}{\mu}\right\}, & y \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$\Gamma(n)$ = gamma function.

$\Gamma(n) = (n-1)!$ if n is integer

$$\Gamma(n) = \int_0^\infty y^{n-1} \exp\{-y\} dy.$$

$$E(Y) = n\mu$$

$$\text{Var}(Y) = n\mu^2$$

(1)

Beta distribution: $X_1 \sim \text{Gamma}(a_1, b)$
 $X_2 \sim \text{Gamma}(a_2, b)$

$$Y = \frac{X_1}{X_1 + X_2} \sim \text{Beta}(a_1, a_2)$$

$$f(y) = \begin{cases} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1} (1-y)^{b-1} & \text{if } 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

If $a=b=1$ then you recover the $\text{Uni}[0,1]$ distribution!

$$\begin{aligned} E(Y) &= \int_{-\infty}^{\infty} y f(y) dy = \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^a (1-y)^{b-1} dy \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 y^a (1-y)^{b-1} dy. \end{aligned}$$

$$\begin{aligned} \boxed{\Gamma(a+1) = a\Gamma(a)} \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 y^{(a+1)-1} (1-y)^{b-1} dy \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} = \frac{a}{a+b} \end{aligned}$$

(2)

Normal distribution: $X \sim N(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

$F(x)$ is not available in closed form!

A common notation is:

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{t^2}{2}\right\} dt.$$

Note that $F(x)$ can be computed in terms of $\Phi(x)$ for any μ and σ .

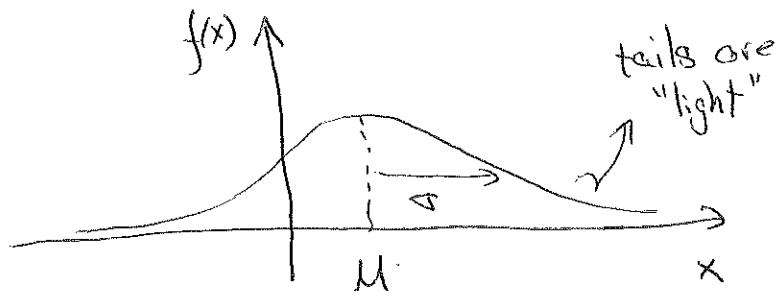
$$F(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{t-\mu}{\sigma}\right)^2\right\} dt$$

$$z = \frac{t-\mu}{\sigma} \quad dz = \frac{dt}{\sigma}$$

$$\Rightarrow F(x) = \int_{-\infty}^{\frac{x-\mu}{\sigma}} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z^2\right\} dz = \Phi\left(\frac{x-\mu}{\sigma}\right)$$

$$E(X) = \mu$$

$$\text{Var}(X) = \sigma^2$$

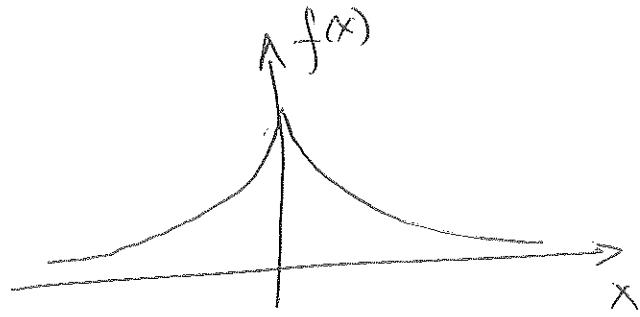


(3)

Double exponential distribution:

$$X \sim DExp(\mu, \lambda)$$

$$f(x) = \frac{1}{2\lambda} \exp\left\{-\frac{1}{\lambda}|x - \mu|\right\}$$



χ^2_n distribution: This is just a special case of the gamma.

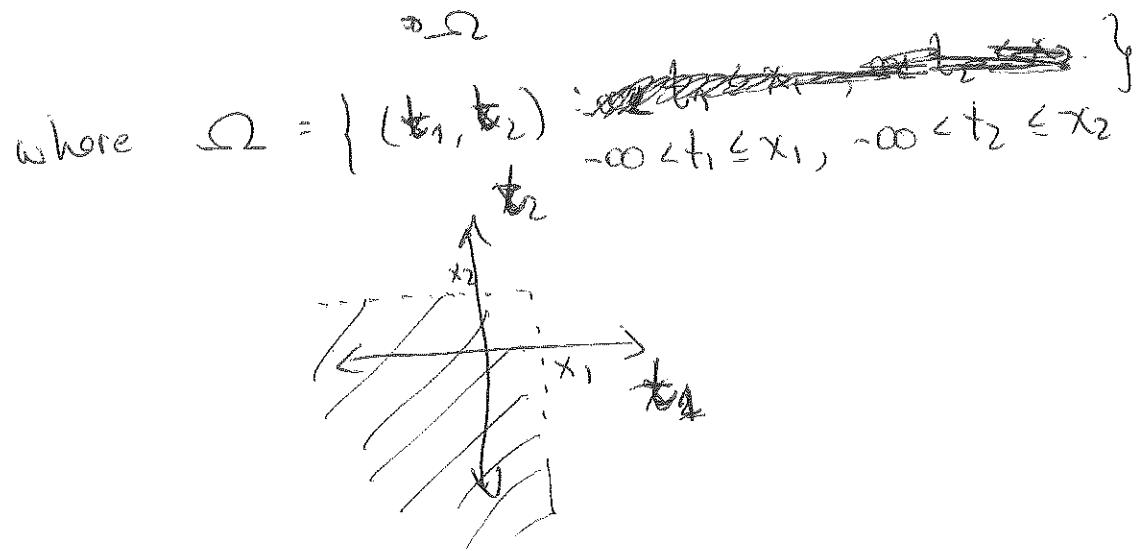
t distribution: we will discuss later when we deal with multiparameter models (inference for normal models).

(4)

Multivariate distributions:

Continuous case: Let X_1 and X_2 be two random variables with joint density $f(x_1, x_2)$ such that

$$\Pr(X_1 \leq x_1, X_2 \leq x_2) = \iint_{\Omega} f(t_1, t_2) dt_1 dt_2$$



You can define the marginal distributions and densities

$$f(x_1) = \int f(x_1, x_2) dx_2 \quad \text{and}$$

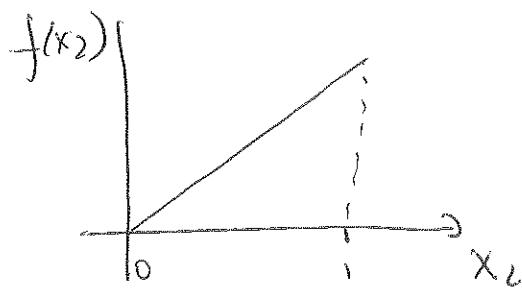
$$f(x_2) = \int f(x_1, x_2) dx_1$$

The conditional densities

$$f(x_1 | x_2) = \frac{f(x_1, x_2)}{f(x_2)} = \frac{f(x_2 | x_1) P(x_1)}{\int f(x_2 | x_1) P(x_1) dx_1}$$

(6)

$$f(x_2) = \int_{x_2}^{x_2} 1 dx_1 = 2x_2 \quad x_2 \in [0, 1]$$

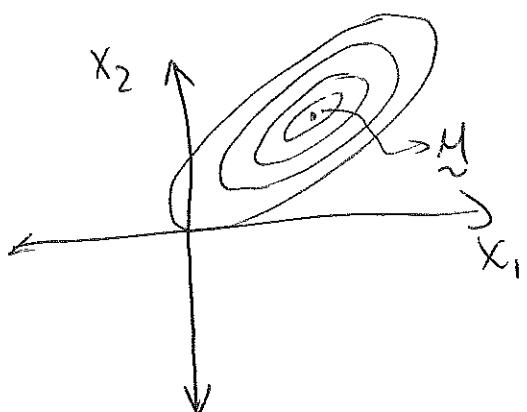


Multivariate normal distribution:

$$\underline{x} = (x_1, \dots, x_n) \sim N_n(\underline{\mu}, \Sigma) \quad \text{positive definite matrix}$$

$$f(\underline{x}) = \left(\frac{1}{\sqrt{2\pi}} \right)^n |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}) \right\}$$

In the bivariate case:

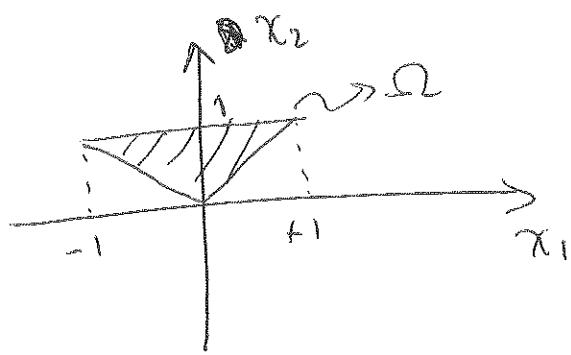


Σ gives you the "shape" \Rightarrow orientation and axes length.

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$$

ρ = correlation coefficient
 σ_i^2 = marginal variances
 μ_i = marginal means

(5)



$$f(x_2)$$

$$f(x_2 | x_1 = a)$$

$$f(x_1, x_2) = \begin{cases} 1 & \text{if } (x_1, x_2) \in \Omega \\ 0 & \text{otherwise} \end{cases}$$

$$f(x_2 | x_1) = \frac{f(x_1, x_2)}{\int f(x_1, x_2) dx_2}$$

$$f(x_1) = \int f(x_1, x_2) dx_2 = \int_{|x_1|}^1 1 dx_2 = 1 - |x_1|$$

$$f(x_2 | x_1) = \frac{1}{1 - |x_1|} \quad \begin{matrix} -1 \leq x_1 \leq 1 \\ |x_1| \leq x_2 \leq 1 \end{matrix}$$

(7)

Difference between independence and lack of correlation.

If $(X_1, X_2) \sim N(\mu, \Sigma)$ and Σ is diagonal.
 (meaning that the variables are uncorrelated) then
 they are also independent; and viceversa.

In general independence \Rightarrow lack of correlation.
 but not the other way around.

$$p = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1)\text{Var}(X_2)}} \quad \text{where } \text{Cov}(X_1, X_2) = E(X_1 X_2) - E(X_1)E(X_2)$$

$$E(X_1 X_2) = \iint x_1 x_2 f(x_1, x_2) dx_1 dx_2$$

$$\xrightarrow{\text{if independent}} = \iint x_1 x_2 f(x_1) f(x_2) dx_1 dx_2$$

$$= \int x_1 f(x_1) dx_1 \times \int x_2 f(x_2) dx_2$$

$$= E(X_1) E(X_2)$$

(8)

Conditional expectations and variances

$$E(X) = E[E(X|Y)]$$

$$\text{Var}(X) = E(\text{Var}(X|Y)) + \text{Var}(E(X|Y)).$$

$$\begin{aligned} E(X) &= \iint x f(x,y) dx dy = \iint x f(x|y) f(y) dx dy \\ &= \underbrace{\int \left[\int x f(x|y) dx \right]}_{g(y) = E(X|Y)} \circledast f(y) dy \end{aligned}$$

$$\begin{aligned} \cancel{\int g(y) dy} &= E(\\ &= \int g(y) f(y) dy = E(g(y)) \\ &= E[E(X|Y)] \end{aligned}$$

A machine makes a random number N of widgets in a day. Assuming that each is defective ~~of the~~ independently of the rest with probability θ , what is the expected number of defective widgets in a day?

(9)

If you know N , $Z = \#$ of defectives.

$$Z|N \sim \text{Bin}(N, \theta)$$

$$E(Z) = E[E(Z|N)] = E[N\theta] = \theta E(N).$$

Bernoulli Trials:

D

Experiments:

- 1) Two possible outcomes (S/F, 0/1)
- 2) Repetcs are independent
- 3) Probability of outcomes is constant over time.

Bernoulli: $X_i \sim \text{Ber}(\theta) \quad X_i \in \{0,1\}$

$$P(X| \theta) = P(X) = \theta^x (1-\theta)^{1-x} \quad \theta = \text{success probability}$$

$$= \begin{cases} \theta & \text{if } X=1 \\ 1-\theta & \text{if } X=0 \end{cases}$$

Binomial: $X_1, X_2, \dots, X_n \sim \text{Ber}(\theta)$ and independent

$Z = X_1 + X_2 + \dots + X_n = \# \text{ of successes in } n \text{ trials}$

$$Z \sim \text{Bin}(n, \theta)$$

$$P(Z|n, \theta) = \binom{n}{z} \theta^z (1-\theta)^{n-z} \quad z = 0, 1, \dots, n$$

$$E(Z) = n\theta$$

$$\text{Var}(Z) = n\theta(1-\theta)$$

(2)

Geometric:

$Y = \# \text{ of trials before a success is observed in a Bernoulli experiment with success prob } \theta$

$$P(Y|θ) = θ(1-θ)^{Y-1} \quad Y = 1, 2, 3, \dots$$

Alternatively, define $W = \# \text{ of failures before the first success.}$

$$Y = W + 1$$

$$P(W|θ) = θ(1-θ)^W \quad W = 0, 1, \dots$$

$$E(Y) = \frac{1}{θ}$$

$$\text{Var}(Y) = \frac{1-θ}{θ^2}$$

Negative Binomial

$$Y_1, Y_2, \dots, Y_K \sim \text{Geo}(θ)$$

$X = Y_1 + Y_2 + \dots + Y_K =$ number of trials to observe K successes

$$P(X|θ, K) = \binom{X-1}{K-1} θ^K (1-θ)^{X-K}$$

opposite to the
Binomial!!!

$$E(X) = \frac{K}{θ}$$

$$\text{Var}(X) = \frac{K(1-θ)}{θ^2}$$

Hyper-geometric (sampling uniformly without replacement) ③

$X = \#$ of blue balls in a sample of size n drawn from an urn that contains K blue balls and $N-K$ red balls.

$$P(X|n, K, N) = \frac{\binom{K}{X} \binom{N-K}{n-X}}{\binom{N}{n}}$$

$$E(X) = \cancel{\frac{Kn}{n}} \quad \frac{Kn}{N} \approx \begin{matrix} \# \text{ of balls in } x \text{ picks of blue ball} \\ \text{sample} \end{matrix} \quad \text{check!!!}$$

$$\text{Var}(X) = \cancel{\frac{Kn}{n^2}} \quad \cancel{\frac{N(N-1)}{N^2}} \quad \frac{Kn}{N} \quad \frac{(N-n)(N-K)}{N(N-1)}$$

Poisson: $X = \text{count of events}$

$$P(X|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad x = 0, 1, \dots$$

Poisson process:

- Consider the number of events (phone calls) that you receive during a given hour.
- Disjoint intervals are independent.
- Cells arise independently.
- The probability of two or more calls in a very short time period is very small.
- The probability of getting exactly one call in a small interval is constant.

If these 3 assumptions are satisfied then:

X = number of calls in one hour

$X \sim \text{Poisson}(\lambda)$ λ = mean number of calls in one hour.

$$E(X) = \lambda$$

$$\text{Var}(X) = \lambda$$

In practice, the fact that $E(X) = \text{Var}(X)$ is quite restrictive. It is often the case that

$\text{Var}(X) > E(X)$ overdispersion.

↳ If this is the case the Negative Binomial might be a better option.

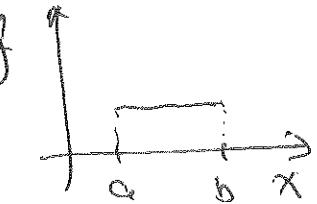
for modeling even if the data generation mechanism has nothing to do with Bernoulli trials!!!

(5)

Continuous random variables:

Uniform: $X \sim \text{Uni}[a, b]$

$$f(x|a,b) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$



$$E(X) = \frac{a+b}{2}$$

"All values have the same likelihood of occurring"

$$\text{Var}(X) = \frac{(b-a)^2}{12}$$

Exponential: interarrival times in a Poisson process

$$X \sim \text{Exp}(\mu)$$

$$f(x|\mu) = \begin{cases} \frac{1}{\mu} \exp\left\{-\frac{x}{\mu}\right\} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$E(X) = \mu$$

$$\text{Var}(X) = \mu^2$$

$$\begin{aligned} F(x|\mu) &= \int_0^x \frac{1}{\mu} \exp\left\{-\frac{t}{\mu}\right\} dt \\ &= 1 - \exp\left\{-\frac{x}{\mu}\right\} = \Pr(X \leq x) \end{aligned}$$

~~$P(X > t+s | X > s)$~~

$$P(X > t+s | X > s) = \frac{P(X > t+s, X > s)}{P(X > s)} = \frac{\Pr(X > t+s)}{\Pr(X > s)} = \exp\left\{-\frac{t+s}{\mu}\right\}$$

Gamma: $X_1, X_2, \dots, X_n \sim \text{Exp}(\mu)$

(6)

$Y = X_1 + X_2 + \dots + X_n \sim \text{Gamma}(n, \mu)$

$$f(y|n, \mu) = \begin{cases} \frac{1}{\Gamma(n)} \left(\frac{1}{\mu}\right)^n y^{n-1} \exp\left\{-\frac{y}{\mu}\right\}, & y \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$\Gamma(n)$ = gamma function.

$\Gamma(n) = (n-1)!$ if n is integer

$$\Gamma(n) = \int_0^\infty y^{n-1} \exp\{-y\} dy.$$

$$E(Y) = n\mu$$

$$\text{Var}(Y) = n\mu^2$$

(1)

Happiness data.

Each female of age 65 or over in the 1998 General Social Survey was asked whether or not they were generally happy. Let $Y_i = 1$ if the i -th woman said "happy" and $Y_i = 0$ if the i -th woman said "unhappy". The survey includes $n = 129$ women over 65 who provided answers.

We want to estimate the "prevalence" of happiness in women over 65 in the US.

Likelihood:

Assuming happiness is constant (and does not on covariates) and independence across observations.

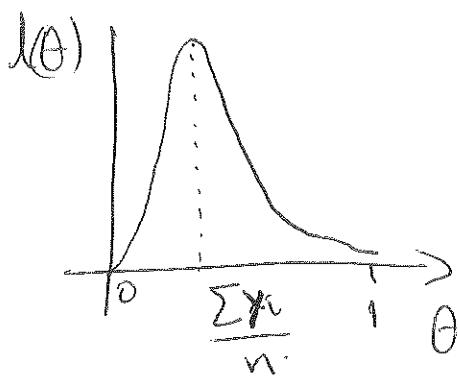
$$f(Y_1, \dots, Y_n | n, \theta) = \prod_{i=1}^n \theta^{Y_i} (1-\theta)^{1-Y_i} = \theta^{\sum Y_i} (1-\theta)^{n-\sum Y_i}$$

$$Z = \sum_{i=1}^n Y_i$$

$$f(Z | n, \theta) = \binom{n}{Z} \theta^Z (1-\theta)^{n-Z} = \binom{n}{\sum Y_i} \theta^{\sum Y_i} (1-\theta)^{n-\sum Y_i}$$

$$l(\theta; \mathbf{y}) = \theta^{\sum Y_i} (1-\theta)^{n-\sum Y_i} \quad 0 \leq \theta \leq 1$$

(2)



Classical statistics would use the MLE as a "good" estimate of the unknown θ

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta)$$

$$\Rightarrow \hat{\theta} = \frac{\sum x_i}{n} \text{ in this case.}$$

$$l(\theta) = P(X|\theta)$$

I want to combine that with a prior distribution $p(\theta)$ and use Bayes theorem.

What does it mean to elicit a prior? What is a prior?

We treat θ as random, but think of it as your knowledge about θ rather than θ being intrinsically random.

If you have no idea, you might want to say that all possible values have the same probability a priori

$$P(\theta) = \begin{cases} 1 & \text{if } 0 \leq \theta \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

$$P(\theta | \mathbf{y}) = \frac{P(\mathbf{y} | \theta) P(\theta)}{\int P(\mathbf{y} | \theta) P(\theta) d\theta} = \begin{cases} \frac{\theta^{\sum y_i} (1-\theta)^{n-\sum y_i}}{\int_0^1 \theta^{\sum y_i} (1-\theta)^{n-\sum y_i} d\theta}, & 0 \leq \theta \leq 1 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

Kernel of a Beta($\sum y_i + 1$,
 $n - \sum y_i + 1$)

$$= \frac{\theta^{\sum y_i} (1-\theta)^{n-\sum y_i}}{\Gamma(\sum y_i + 1) \Gamma(n - \sum y_i + 1)} = \frac{\Gamma(n+2)}{\Gamma(\sum y_i + 1) \Gamma(n - \sum y_i + 1)} \theta^{\sum y_i + n - \sum y_i} (1-\theta)^{n - \sum y_i + 1}$$

Assume that we have observed in the past values between 0.1 and 0.4 for θ (but never anything smaller or bigger). How do we incorporate this into the analysis.

For convenience, start with $\theta \sim \text{Beta}(a, b)$ and let's tweak a and b to reflect our prior knowledge.

$$E(\theta) = \frac{a}{a+b} = 0.25$$

Consider a few cases

a	b
0.1	0.3
1	3
10	30

(4)

My prior is going to be $a=10$ $b=30$

$$\begin{aligned}
 p(\theta|y) &:= \frac{\theta^{\sum y_i} (1-\theta)^{n-\sum y_i} \frac{\cancel{\Gamma(a+b)}}{\cancel{\Gamma(a)}\cancel{\Gamma(b)}} \theta^{a-1} (1-\theta)^{b-1}}{\int \theta^{\sum y_i} (1-\theta)^{n-\sum y_i} \frac{\cancel{\Gamma(a+b)}}{\cancel{\Gamma(a)}\cancel{\Gamma(b)}} \theta^{a-1} (1-\theta)^{b-1} d\theta} \\
 &= \frac{\theta^{\sum y_i + a-1} (1-\theta)^{n-\sum y_i + b-1}}{\int \theta^{\sum y_i + a-1} (1-\theta)^{n-\sum y_i + b-1} d\theta} \\
 &= \frac{\Gamma(n+a+b)}{\Gamma(\sum y_i + a) \Gamma(n-\sum y_i + b)} \theta^{\sum y_i + a-1} (1-\theta)^{n-\sum y_i + b-1}
 \end{aligned}$$

If you recognize the numerator as the kernel of a distribution you know, then you do not need to solve the integral explicitly because it has to be the reciprocal of the normalizing constant of that distribution

$$\begin{aligned}
 p(\theta|y) \propto p(y|\theta) p(\theta) &\propto \theta^{\sum y_i + a-1} (1-\theta)^{n-\sum y_i + b-1} \\
 \Rightarrow \theta|y &\sim \text{Beta}(a_n, b_n) \\
 a_n = \sum y_i + a &\quad b_n = n - \sum y_i + b
 \end{aligned}$$

To provide a point estimate we could use the (5)
 posterior mean:
 $a = b = 1$

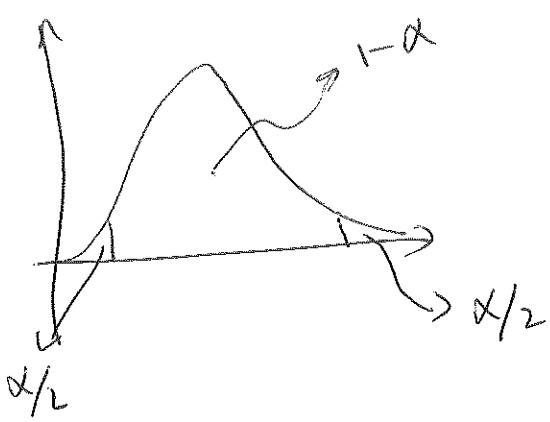
$$E(\theta|y) = \frac{a_n}{a_n + b_n} = \frac{\sum y_i + a}{n + a + b} = \text{point estimator}$$

$$= 0.1984$$

Interval estimator: Credible interval: an interval that receives high probability under the posterior.

An easy interval is found by computing the $\alpha/2$ and $1 - \alpha/2$ quantile of the posterior for small α .

$$(0.1350, 0.2706)$$



(1)

Poisson model:

$$x_i \sim \text{Poisson}(\theta)$$

$$\theta \sim \text{Gamma}(a, b)$$

$$f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!} = \left(\prod_{i=1}^n \frac{1}{x_i!} \right) e^{-n\theta} \theta^{\sum_{i=1}^n x_i}$$

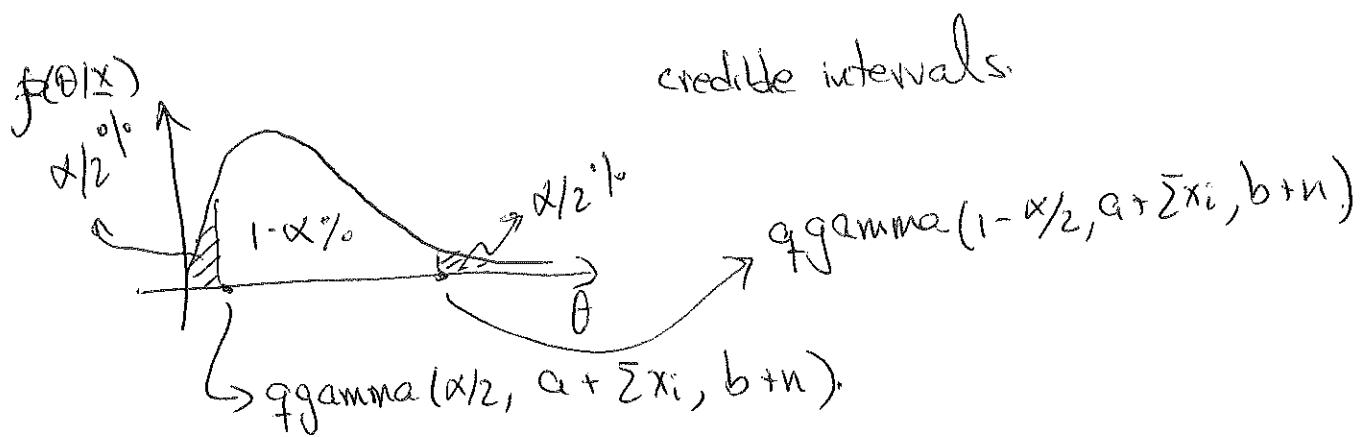
$$f(\theta) = \begin{cases} \frac{b^a}{\Gamma(a)} \theta^{a-1} \exp\{-b\theta\} & \theta > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$f(\theta | x_1, \dots, x_n) \propto \theta^{\sum x_i} e^{-n\theta} \theta^{a-1} \exp\{-b\theta\}$$

$$= \theta^{\sum x_i + a - 1} \exp\{-\theta(b+n)\}$$

$$\Rightarrow \theta | x_1, \dots, x_n \sim \text{Gamma}(a + \sum x_i, b + n).$$

$$E(\theta | x_1, \dots, x_n) = \frac{a + \sum x_i}{b + n}$$



Hypothesis testing:

(2)

$$H_0: \theta = \theta_0 \quad \text{vs} \quad H_a: \theta \neq \theta_0$$

$$\Pr(H_0) = 1 - \Pr(H_a) = 1/2$$

~~$$P(\bar{x}|\theta_0) p(x_i|\theta, H_0) = \left(\prod_{i=1}^n \frac{1}{x_i!} \right) \theta_0^{\sum x_i} \exp\{-n\theta_0\}$$~~

~~$$P(\bar{x}|\theta) p(x_i|\theta, H_a) = \left(\prod_{i=1}^n \frac{1}{x_i!} \right) \theta^{\sum x_i} \exp\{-n\theta\}$$~~

$$P(\theta|H_0) = S_{\theta_0}(\theta)$$

$$P(\theta|H_a) = \frac{b^a}{\Gamma(a)} \exp\{-b\theta\} \theta^{a-1}$$

$$\frac{a}{b} = \theta_0 \quad b = a \quad \Rightarrow \quad b = \frac{a}{\theta_0} \quad a = \gamma$$

$$= \left(\frac{\gamma}{\theta_0} \right) \frac{1}{\Gamma(\gamma)} \theta^{\gamma-1} \exp\left\{-\frac{\gamma}{\theta_0} \theta\right\}.$$

$$B_{0a} = \frac{\left(\prod_{i=1}^n \frac{1}{x_i!} \right) \theta_0^{\sum x_i} \exp\{-n\theta_0\}}{\int \left(\prod_{i=1}^n \frac{1}{x_i!} \right) \theta^{\sum x_i} \exp\{-n\theta\} \left(\frac{\gamma}{\theta_0} \right)^{\gamma} \frac{1}{\Gamma(\gamma)} \theta^{\gamma-1} \exp\left\{-\frac{\gamma}{\theta_0} \theta\right\} d\theta}$$

(3)

$$B_{\theta_0} = \frac{\int_{-\infty}^{\sum x_i} \theta_0^{\sum x_i} \exp\{-n\theta_0\} d\theta}{\int_{-\infty}^{\sum x_i + v - 1} \left(\frac{v}{\theta_0}\right)^v \frac{1}{\Gamma(v)} \int_{-\infty}^{\sum x_i + v - 1} \theta^{v-1} \exp\left\{-\left(n + \frac{v}{\theta_0}\right)\theta\right\} d\theta d\theta}$$

$$= \frac{\int_{-\infty}^{\sum x_i} \theta_0^{\sum x_i} \exp\{-n\theta_0\} d\theta}{\left(\frac{v}{\theta_0}\right)^v \frac{1}{\Gamma(v)} \Gamma(\sum x_i + v) \left(\frac{\sum x_i + v}{n + \frac{v}{\theta_0}}\right)^{\sum x_i + v}}$$

$$= \frac{\int_{-\infty}^{\sum x_i} \theta_0^{\sum x_i} \exp\{-n\theta_0\} d\theta}{\left(\frac{v}{\theta_0}\right)^v \left(\frac{n + \frac{v}{\theta_0}}{v + \sum x_i}\right)^{v + \sum x_i}}$$

How should Bayes Factor be interpreted?

Kass & Raftery 1995

Gaussian Beta:

(4)

$x_1, \dots, x_n \sim \text{IID}, \phi \sim N(\theta, \frac{1}{\phi})$ so that $\frac{\sigma^2}{\phi} = \frac{1}{\phi}$
"variance"

$\phi = \text{precision}$.

Assume ϕ is known.

$$f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \phi^{1/2} \exp \left\{ -\frac{\phi}{2} (x_i - \theta)^2 \right\}$$

$$= \left(\frac{1}{\sqrt{2\pi}} \right)^n \phi^{n/2} \exp \left\{ -\frac{\phi}{2} \sum_{i=1}^n (x_i - \theta)^2 \right\}$$

$$\sum_{i=1}^n (x_i - \theta)^2 = \sum_{i=1}^n (x_i^2 - 2x_i\theta + \theta^2) = n\theta^2 - 2\theta \bar{x} + \sum_{i=1}^n x_i^2$$
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$f(x_1, \dots, x_n | \theta) = \left(\frac{1}{\sqrt{2\pi}} \right)^n \phi^{n/2} \exp \left\{ -\frac{\phi}{2} (n\theta^2 - 2\theta n\bar{x} + \sum x_i^2) \right\}$$

$$\theta \sim N(\mu, \tau^2) \Rightarrow p(\theta) = \frac{1}{\sqrt{2\pi}} \frac{1}{\tau} \exp \left\{ -\frac{1}{2\tau^2} (\theta - \mu)^2 \right\}$$

$$P(\theta | x_1, \dots, x_n) \propto \exp \left\{ -\frac{\phi}{2} (n\theta^2 - 2\theta n\bar{x}) - \frac{1}{2\tau^2} (\theta^2 - 2\mu\theta) \right\} \quad (5)$$

$$= \exp \left\{ -\frac{1}{2} \left[\left(\phi n + \frac{1}{\tau^2} \right) \theta^2 - 2\theta \left(n\phi\bar{x} + \frac{\mu}{\tau^2} \right) \right] \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} \left(\phi n + \frac{1}{\tau^2} \right) \left[\theta - \left(\phi n + \frac{1}{\tau^2} \right)^{-1} \left(n\phi\bar{x} + \frac{\mu}{\tau^2} \right) \right]^2 \right\}$$

$$\theta | x_1, \dots, x_n \sim N \left(\frac{n\phi\bar{x} + \frac{\mu}{\tau^2}}{\left(\phi n + \frac{1}{\tau^2} \right)}, \left(\phi n + \frac{1}{\tau^2} \right)^{-1} \right)$$

①

Hypothesis testing: two proportions.

Date: x = number of suc in pop 1

n = number of trials in pop 1

y = number of suc in pop 2

m = number of trials in pop 2

Likelihood $x \sim \text{Bin}(n, \theta_1)$ $y \sim \text{Bin}(m, \theta_2)$ independent

Goal $H_0: \theta_1 = \theta_2$ vs $H_a: \theta_1 \neq \theta_2$

$$\Pr(H_0 | \tilde{x}, \tilde{y}) = \frac{\Pr(\tilde{x}, \tilde{y} | H_0)}{\Pr(\tilde{x}, \tilde{y})}$$

We are going to use Bayes theorem.

$$\Pr(H_0 | \tilde{x}, \tilde{y}) = \frac{\Pr(\tilde{x}, \tilde{y} | H_0) \Pr(H_0)}{\Pr(\tilde{x}, \tilde{y}) \Pr(H_0) + \Pr(\tilde{x}, \tilde{y} | H_a) \Pr(H_a)}$$

density density

$$= \frac{f(\tilde{x}, \tilde{y} | H_0) \Pr(H_0)}{f(\tilde{x}, \tilde{y} | H_0) \Pr(H_0) + f(\tilde{x}, \tilde{y} | H_a) \Pr(H_a)}$$

In the absence of other information, we might want to use

$$\Pr(H_0) = \Pr(H_a) = \frac{1}{2}$$

$$f(x, y | H_0) = \int f(x, y | \theta_1, \theta_2, H_0) P(\theta_1, \theta_2 | H_0) d\theta_1 d\theta_2 \quad (2)$$

$$f(x, y | \theta_1, \theta_2, H_0) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \binom{m}{y} \theta^y (1-\theta)^{m-y}$$

$$f(x, y | \theta, H_0) = \binom{n}{x} \binom{m}{y} \theta^{x+y} (1-\theta)^{n+m-x-y}$$

$$P(\theta_1, \theta_2 | H_0) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

$$P(\theta | H_0)$$

$$f(x, y | H_0) = \binom{n}{x} \binom{m}{y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \theta^{x+y+a-1} (1-\theta)^{n+m-x-y+b-1} d\theta$$

$$= \binom{n}{x} \binom{m}{y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(x+y+a)}{\Gamma(n+m+a+b)} \frac{\Gamma(n+m-x-y+b)}{\Gamma(n+m+a+b)}$$

→ Marginal likelihood (under model H_0)
 prior predictive distribution (under model H_0)

$$f(x, y | H_a) = \iint f(x, y | \theta_1, \theta_2 | H_a) p(\theta_1, \theta_2 | H_a) d\theta_1 d\theta_2 \quad (3)$$

$$f(x, y | \theta_1, \theta_2, H_a) = \binom{n}{x} \theta_1^x (1-\theta_1)^{n-x} \binom{m}{y} \theta_2^y (1-\theta_2)^{m-y}$$

$$P(\theta_1, \theta_2 | H_a) = P(\theta_1 | H_a) P(\theta_2 | H_a)$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta_1^{a-1} (1-\theta_1)^{b-1} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta_2^{a-1} (1-\theta_2)^{b-1}$$

→ Same a and b as under $H_0!!$

$$f(x, y | H_a) = \binom{n}{x} \binom{m}{y} \left[\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \right]^2 \times \int \theta_1^{x+a-1} (1-\theta_1)^{n-x+b-1} d\theta_1$$

$$\times \int \theta_2^{y+a-1} (1-\theta_2)^{m-y+b-1} d\theta_2$$

$$= \binom{n}{x} \binom{m}{y} \left[\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \right]^2 \frac{\Gamma(a+x) \Gamma(b+n-y)}{\Gamma(a+b+n)}$$

$$\frac{\Gamma(a+y) \Gamma(b+m-y)}{\Gamma(a+b+m)}.$$

$$\Pr(H_0 | \underline{x}, \underline{y}) = \frac{\frac{1}{2} \binom{n}{x} \binom{m}{y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(x+y+a) \Gamma(n+m-x-y+b)}{\Gamma(n+m+a+b)}}{\frac{1}{2} \binom{n}{x} \binom{m}{y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(x+y+a) \Gamma(n+m-x-y+b)}{\Gamma(n+m+a+b)} + \dots}$$

$\rightarrow \frac{\frac{1}{2} \binom{n}{x} \binom{m}{y} \left(\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \right)^2 \frac{\Gamma(a+x) \Gamma(b+n-x) \Gamma(a+y) \Gamma(b+m-y)}{\Gamma(a+b+n) \Gamma(a+b+m)}}{\frac{\Gamma(x+y+a) \Gamma(n+m-x-y+b)}{\Gamma(n+m+a+b)}}$

$$= \frac{\frac{\Gamma(x+y+a) \Gamma(n+m-x-y+b)}{\Gamma(n+m+a+b)}}{\frac{\Gamma(x+y+a) \Gamma(n+m-x-y+b)}{\Gamma(n+m+a+b)} + \frac{\Gamma(a+b) \Gamma(a+x) \Gamma(b+n-x) \Gamma(a+y) \Gamma(b+m-y)}{\Gamma(a)\Gamma(b) \Gamma(a+b+n) \Gamma(a+b+m)}}$$

Other ways to write $\Pr(H_0 | \underline{x}, \underline{y})$:

$$\Pr(H_0 | \underline{x}, \underline{y}) = \frac{1}{1 + \underbrace{\frac{\Pr(H_a)}{\Pr(H_0)} \frac{\Pr(\underline{x}, \underline{y} | H_a)}{\Pr(\underline{x}, \underline{y} | H_0)}}_{\text{Posterior odds}}}$$

$$\Rightarrow \frac{\Pr(H_0 | \underline{x}, \underline{y})}{\Pr(H_a | \underline{x}, \underline{y})} = \frac{\Pr(H_0)}{\Pr(H_a)} \frac{\Pr(\underline{x}, \underline{y} | H_0)}{\Pr(\underline{x}, \underline{y} | H_a)}$$

$\underbrace{\Pr(H_0)}_{\text{Posterior odds}} \quad \underbrace{\Pr(\underline{x}, \underline{y} | H_a)}_{\text{Prior odds}} \quad \underbrace{\frac{\Pr(H_0)}{\Pr(H_a)} \frac{\Pr(\underline{x}, \underline{y} | H_0)}{\Pr(\underline{x}, \underline{y} | H_a)}}_{\text{Bayes factor}}$

(5)

Count data and Poisson likelihoods:

Public health: estimating prevalence.

Data from different counties.

s_i = total population in the county.

x_i = # of recorded cases in the county.

You want to estimate the infection rate. θ

$$\hat{\theta} = \frac{\sum_{i=1}^n x_i}{\sum_{i=1}^n s_i} \text{ = maximum likelihood estimator under a Poisson model}$$

$x_i | \theta$ Poisson ($s_i \theta$) independent $i = 1, \dots, n$

$$P(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n \frac{(s_i \theta)^{x_i} e^{-s_i \theta}}{x_i!}$$

$$= \frac{\prod_{i=1}^n (s_i \theta)^{x_i}}{\prod_{i=1}^n x_i!} \theta^{\sum x_i} e^{-\theta \sum s_i}$$

Bayesian equivalent?

Prior? How about $\theta \sim \text{Gamma}(a, b)$ a priori

$$P(\theta) = \frac{b^a}{\Gamma(a)} \theta^{a-1} e^{-b\theta}$$

$$\begin{aligned}
 P(\theta | \mathbf{x}) &= \frac{\cancel{\prod (s_i)^{x_i}}}{\cancel{\prod x_i!}} \theta^{\sum x_i} \exp\{-\theta \sum s_i\} \frac{b^a}{\cancel{\Gamma(a)}} \theta^{a-1} \exp\{-b\theta\} \quad (6) \\
 &= \frac{\theta^{\sum x_i + a-1} \exp\{-(b + \sum s_i)\theta\}}{\int \theta^{\sum x_i + a-1} \exp\{-(b + \sum s_i)\theta\} d\theta} \\
 &\Rightarrow \theta | \mathbf{x} \sim \text{Gamma}(a + \sum x_i, b + \sum s_i)
 \end{aligned}$$

$$\phi \sim \text{Gamma}(a, b) \quad \theta | \phi \sim N(\mu, K\phi^{-1})$$

①

$p(\theta, \phi | \underline{x})$ is such that

$$\phi | \underline{x} \sim \text{Gamma}(a_n, b_n)$$

where $a_n = a + \frac{n}{2}$

$$b_n = b + \frac{1}{2} \sum_{i=1}^n x_i^2 + \frac{1}{2} \frac{\mu^2}{K} - \frac{1}{2} \left[\frac{n\bar{x} + \frac{\mu}{K}}{n + \frac{1}{K}} \right]^2$$

$$\theta | \phi, \underline{x} \sim N(\mu_n, K_n \phi^{-1})$$

$$\mu_n = \frac{n\bar{x} + \frac{\mu}{K}}{n + \frac{1}{K}}$$

$$K_n = n + \frac{1}{K}$$

$$\begin{aligned} p(\theta, \phi | \underline{x}) &= \frac{b_n^{a_n} \phi^{a_n-1} \exp\{-\phi b_n\}}{\Gamma(a_n)} \\ &\times \frac{1}{\sqrt{2\pi}} \frac{\phi^{1/2}}{K_n^{1/2}} \exp\left\{-\frac{\phi}{2K_n} (\theta - \mu_n)^2\right\} \\ &= \frac{1}{\sqrt{2\pi}} \frac{b_n^{a_n}}{\Gamma(a_n)} \frac{1}{K_n^{1/2}} \phi^{a_n + 1/2 - 1} \exp\left\{-\phi \left[b_n + \frac{(\theta - \mu_n)^2}{2K_n}\right]\right\} \end{aligned}$$

$$p(\theta | \underline{x}) = \int p(\theta, \phi | \underline{x}) d\phi = \int \frac{1}{\sqrt{2\pi}} \frac{b_n^{a_n}}{\Gamma(a_n) K_n^{1/2}} \phi^{a_n + 1/2 - 1} \exp\left\{-\phi \left[b_n + \frac{(\theta - \mu_n)^2}{2K_n}\right]\right\} d\phi$$

$$= \frac{1}{\sqrt{2\pi}} \frac{b_n^{a_n}}{\Gamma(a_n)} \frac{1}{K_n^{1/2}} \frac{\Gamma(a_n + 1/2)}{\left[b_n + \frac{(\theta - \mu_n)^2}{2K_n}\right]^{a_n + 1/2}}$$

This a (scaled and shifted t distribution.)

$$\begin{aligned}
 P(\theta | x) &= \frac{1}{\sqrt{2\pi}} \frac{\Gamma(a_n + \frac{1}{2})}{\Gamma(a_n)} \frac{b_n^{a_n}}{K_n^{\frac{1}{2}}} \frac{1}{b_n^{a_n + \frac{1}{2}}} \\
 &\quad \left[1 + \frac{(\theta - \mu_n)^2}{2b_n K_n} \right]^{-(a_n + \frac{1}{2})} \\
 &= \frac{1}{\sqrt{2\pi}} \frac{\Gamma(a_n + \frac{1}{2})}{\Gamma(a_n)} \left(\frac{1}{b_n K_n} \right)^{\frac{1}{2}} \left[1 + \frac{(\theta - \mu_n)^2}{2b_n K_n} \right]^{-(a_n + \frac{1}{2})}
 \end{aligned}$$

(2)

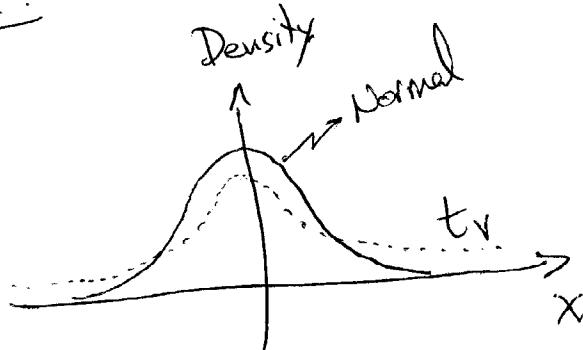
~~$\theta = \mu_n + \frac{1}{2} \ln \left(\frac{1}{b_n K_n} \right)$~~

$Z = (\theta - \mu_n) / \sqrt{\frac{b_n K_n}{a_n}}$

$$\Rightarrow P(Z | x) = \frac{1}{\sqrt{2\pi}} \frac{\Gamma(a_n + \frac{1}{2})}{\Gamma(a_n)} a_n^{-\frac{1}{2}} \left[1 + \frac{Z^2}{2a_n} \right]^{-(a_n + \frac{1}{2})}$$

$Z \sim t_{\underline{2a_n}}$ degrees of freedom

t distribution:



(3)

Point estimation.

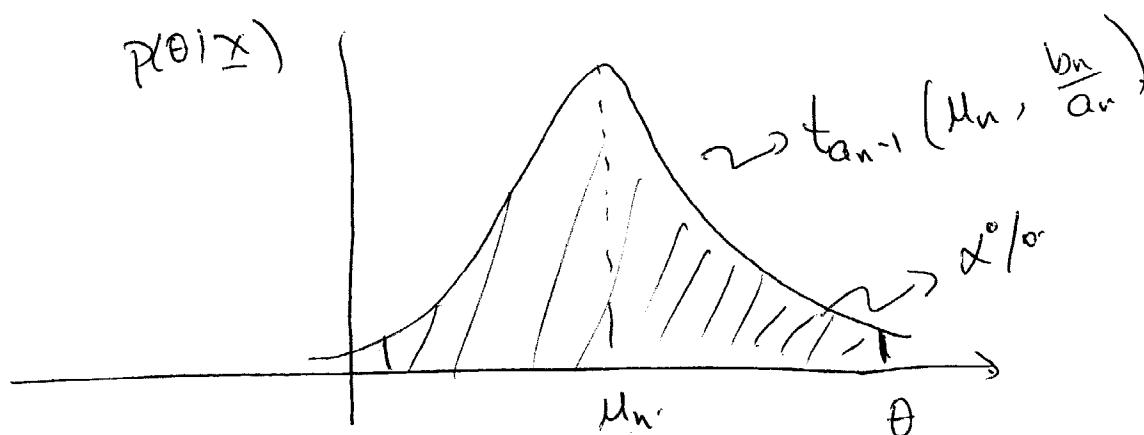
$$\tilde{\theta} = E(\theta | \underline{x}) \doteq \mu_n$$

$$E(\theta | \underline{x}) = E\left[E(\theta | \phi, \underline{x})\right] = E[\mu_n] = \mu_n.$$

$$\begin{aligned} \text{Var}(\theta | \underline{x}) &= \text{Var}(E(\theta | \phi, \underline{x})) + E(\text{Var}(\theta | \phi, \underline{x})) \\ &= \text{Var}(\mu_n) + E\left[k_n \phi^{-1}\right] \\ &= 0 + k_n \cdot \frac{b_n}{a_n - 1} = \frac{k_n b_n}{a_n - 1} \end{aligned}$$

$$E(\phi^{-1} | \underline{x}) = \frac{b_n}{a_n - 1}$$

Credible intervals



Because the distribution is symmetric, the HPD credible interval and the symmetric interval match!!

(4)

Hypothesis testing:

$$x_1, \dots, x_n \sim N(\mu_1, \phi_1^{-1}) \quad y_1, \dots, y_m \sim N(\mu_2, \phi_2^{-1})$$

with ϕ_1 and ϕ_2 known.

$$H_0: \mu_1 = \mu_2 \quad \text{vs} \quad H_a: \mu_1 \neq \mu_2$$

$$Pr(H_0) = 1 - Pr(H_a) = \gamma_2$$

$$\begin{aligned} P(x, y | H_a, \mu_1, \mu_2) &= \prod_{i=1}^n P(x_i | \mu_1) \prod_{j=1}^m P(y_j | \mu_2) \\ &= \left(\frac{1}{\sqrt{2\pi}} \right)^n \phi_1^{\gamma_2} \exp \left\{ -\frac{\phi_1}{2} \sum_{i=1}^n (x_i - \mu_1)^2 \right\} \times \\ &\quad \left(\frac{1}{\sqrt{2\pi}} \right)^m \phi_2^{\gamma_2} \exp \left\{ -\frac{\phi_2}{2} \sum_{j=1}^m (y_j - \mu_2)^2 \right\} \\ P(x, y | H_0, \mu) &= \left(\frac{1}{\sqrt{2\pi}} \right)^{n+m} (\phi_1 \phi_2)^{\gamma_2} \exp \left\{ -\frac{\phi_1}{2} \sum_{i=1}^n (x_i - \mu)^2 - \right. \\ &\quad \left. \frac{\phi_2}{2} \sum_{j=1}^m (y_j - \mu)^2 \right\} \end{aligned}$$

~~$\mu | H_0 \sim N(\gamma, \tau^2)$~~

$$\mu_1 | H_a \sim N(\gamma, \tau^2) \text{ independent}$$

$$\mu_2 | H_a \sim N(\gamma, \tau^2)$$

$$P(x, y | H_a) = \iint P(x, y | H_a, \mu_1, \mu_2) P(\mu_1, \mu_2 | H_a) d\mu_1 d\mu_2 \quad (5)$$

$$= \left(\frac{1}{\sqrt{2\pi}} \right)^n \phi_1^{y_1} \exp \left\{ -\frac{\phi_1}{2} \sum_{i=1}^n (x_i - \mu_1)^2 \right\} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\tau^2} (\mu_1 - \bar{\gamma})^2 \right\} d\mu_1$$

$$\times \left(\frac{1}{\sqrt{2\pi}} \right)^m \phi_2^{y_2} \exp \left\{ -\frac{\phi_2}{2} \sum_{j=1}^m (y_j - \mu_2)^2 \right\} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\tau^2} (\mu_2 - \bar{\gamma})^2 \right\} d\mu_2.$$

$$(1) = \left(\frac{1}{\sqrt{2\pi}} \right)^{n+1} (\phi_1 \phi_2)^{y_2} \int_{-\infty}^{\infty} \left\{ \frac{\phi_1}{2} \sum_{i=1}^n (x_i^2 - 2x_i \mu_1 + \mu_1^2) - \frac{1}{2\tau^2} (\mu_1^2 - 2\mu_1 \bar{\gamma} + \bar{\gamma}^2) \right\} d\mu_1$$

$$= \left(\frac{1}{\sqrt{2\pi}} \right)^{n+1} (\phi_1 \phi_2)^{y_2} \exp \left\{ -\frac{\phi_1}{2} \sum_{i=1}^n x_i^2 - \frac{\bar{\gamma}^2}{2\tau^2} \right\} \times$$

$$\exp \left\{ -\frac{1}{2} \left(n\mu_1^2 \phi_1 + \frac{\mu_1^2}{\tau^2} - 2\mu_1 \bar{x} \phi_1 - \frac{2\mu_1 \bar{\gamma}}{\tau^2} \right) \right\} d\mu_1$$

$$(3) = \int \exp \left\{ -\frac{1}{2} \left(n\phi_1 + \frac{1}{\tau^2} \right) \left(\mu_1^2 - 2\mu_1 \left(\frac{n\phi_1 \bar{x} + \frac{\bar{\gamma}}{\tau^2}}{n\phi_1 + \frac{1}{\tau^2}} \right) \right) \right\} d\mu_1$$

$$= \int \exp \left\{ -\frac{1}{2} \left(n\phi_1 + \frac{1}{\tau^2} \right) \left[\left(\mu_1 - \frac{n\phi_1 \bar{x} + \frac{\bar{\gamma}}{\tau^2}}{n\phi_1 + \frac{1}{\tau^2}} \right)^2 - \left(\frac{n\phi_1 \bar{x} + \frac{\bar{\gamma}}{\tau^2}}{n\phi_1 + \frac{1}{\tau^2}} \right)^2 \right] \right\} d\phi_1$$

$$= \exp \left\{ \frac{1}{2} \left[\frac{\left(n\phi_1 \bar{x} + \frac{\bar{\gamma}}{\tau^2} \right)^2}{n\phi_1 + \frac{1}{\tau^2}} \right] \right\} \sqrt{2\pi} \left(n\phi_1 + \frac{1}{\tau^2} \right)^{-\frac{1}{2}}$$

→ Complete!!!

(1)

Recap

$$x_1, \dots, x_n | \theta \stackrel{\text{iid}}{\sim} N(\theta, \phi^{-1}) \quad \phi \text{ known, } \theta \text{ unknown.}$$

$$\theta \sim N(\mu, \tau^2)$$

$$\Rightarrow \theta | x_1, \dots, x_n \sim N\left(\frac{n\phi\bar{x} + \frac{\mu}{\tau^2}}{n\phi + \frac{1}{\tau^2}}, \frac{1}{n\phi + \frac{1}{\tau^2}} \right)$$

$$E(\theta | x_1, \dots, x_n) = \frac{n\phi\bar{x} + \frac{\mu}{\tau^2}}{n\phi + \frac{1}{\tau^2}} = w\bar{x} + (1-w)\mu$$

$$w = \frac{n\phi}{n\phi + \frac{1}{\tau^2}}$$

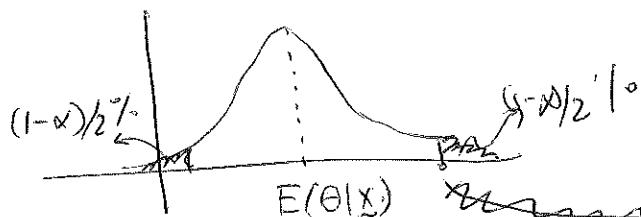
What is the posterior median in this case?

Same as before

What about posterior mode?

Same

What is the $\alpha\%$ posterior symmetric credible interval?



$$\frac{n\phi\bar{x} + \frac{\mu}{\tau^2}}{n\phi + \frac{1}{\tau^2}} \pm z_{\alpha/2} \sqrt{\frac{1}{n\phi + \frac{1}{\tau^2}}}$$

$$z_{\alpha/2} = \Phi^{-1}\left(1 - \frac{(1-\alpha)}{2}\right)$$

(2)

If $\alpha = 95\%$ then

$$\text{credible interval} \quad \frac{n\phi\bar{x} + \frac{1}{\tau^2}}{n\phi + \frac{1}{\tau^2}} \pm 1.96 \sqrt{\frac{1}{n\phi + \frac{1}{\tau^2}}}.$$

If $\tau^2 \rightarrow \infty$ this becomes

$$\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}} \quad (\text{remember } \phi = \frac{1}{\sigma^2})$$

Consider now the case θ is known and ϕ is unknown.

Remember the likelihood

$$P(x_1, \dots, x_n | \phi) = \left(\frac{1}{\sqrt{2\pi}} \right)^n \phi^{n/2} \exp \left\{ -\frac{\phi}{2} \sum_{i=1}^n (x_i - \theta)^2 \right\}$$

$\phi \sim \text{Gamma}(a, b)$

$$\Rightarrow P(\phi) = \frac{b^a}{\Gamma(a)} \phi^{a-1} \exp \{-b\phi\}$$

$$P(\phi | \mathbf{x}) \propto \phi^{n/2} \exp \left\{ -\phi \frac{\sum_{i=1}^n (x_i - \theta)^2}{2} \right\} \phi^{a-1} \exp \{-b\phi\}$$

$$= \phi^{n/2+a-1} \exp \left\{ -\phi \left[b + \frac{\sum_{i=1}^n (x_i - \theta)^2}{2} \right] \right\}.$$

$$\Rightarrow \phi | \mathbf{x} \sim \text{Gamma} \left(\frac{n}{2} + a, b + \frac{\sum_{i=1}^n (x_i - \theta)^2}{2} \right)$$

$$E(\phi | \mathbf{x}) = \frac{\frac{n}{2} + a}{b + \frac{\sum (x_i - \theta)^2}{2}}$$

$$f(x) = (x-1)^{\alpha} \Gamma(x-1)$$

(3)

~~Explain~~

$$E(\sigma^2 | x_1, \dots, x_n) = E\left(\frac{1}{\phi} | x_1, \dots, x_n\right)$$

$$= \int \frac{1}{\phi} \frac{\left(b + \frac{\sum(x_i - \theta)^2}{2}\right)^{a+\frac{n}{2}}}{\Gamma\left(\frac{n}{2} + a\right)} \phi^{a+\frac{n}{2}-1} \exp\left\{-\phi\left(b + \frac{\sum(x_i - \theta)^2}{2}\right)\right\} d\phi$$

$$= \frac{\left(b + \frac{\sum(x_i - \theta)^2}{2}\right)^{a+\frac{n}{2}}}{\Gamma\left(\frac{n}{2} + a\right)} \frac{\sum\left(\frac{n}{2} + a - 1\right)}{\left(b + \frac{\sum(x_i - \theta)^2}{2}\right)^{a+\frac{n}{2}-1}}$$

$$= \frac{b + \sum \frac{(x_i - \theta)^2}{2}}{a + \frac{n}{2} - 1}$$

Let's derive the distribution of τ^2 .

$$\eta = \frac{1}{\phi} \quad d\eta = \frac{d\phi}{\phi^2} \Rightarrow \phi = \frac{1}{\eta} \quad d\phi = \frac{d\eta}{\eta^2}$$

$$P(\eta) = \frac{b^a}{\Gamma(a)} \left(\frac{1}{\eta}\right)^{a-1} \exp\left\{-\frac{b}{\eta}\right\} \left\{ \frac{1}{\eta^2} \right\}$$

$$= \frac{b^a}{\Gamma(a)} \left(\frac{1}{\eta}\right)^{a+1} \exp\left\{-\frac{b}{\eta}\right\}$$

∴ The distribution of η is called the inverse-gamma distribution.

(4)

We are calling $\eta = \sigma^2$ so

$$P(\sigma^2) = \frac{b^a}{\Gamma(a)} \left(\frac{1}{\sigma^2} \right)^{a+1} \exp \left\{ -\frac{b}{\sigma^2} \right\} \quad \text{IG}(a, b)$$

Now with this prior:

$$\sigma^2 | x \sim \text{IG} \left(a + \frac{n}{2}, b + \frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2 \right)$$

Now, what if both θ and ϕ are unknown?

$$\begin{aligned} P(x_1, \dots, x_n | \theta, \phi) &= \left(\frac{1}{\sqrt{2\pi}} \right)^n \phi^{n/2} \exp \left\{ -\frac{\phi}{2} \sum (x_i - \theta)^2 \right\} \\ &= \left(\frac{1}{\sqrt{2\pi}} \right)^n \phi^{n/2} \exp \left\{ -\frac{\phi}{2} \left[\sum x_i^2 - 2n\bar{x}\theta + \theta^2 \right] \right\} \end{aligned}$$

$$\theta | \phi \sim N \left(\mu, \frac{\kappa}{\phi} \right).$$

$$\phi \sim \text{Gamma}(a, b).$$

\Rightarrow This is different from

$$\theta \sim N(\mu, \tau^2)$$

$$\phi \sim \text{Gamma}(a, b).$$

!!

$$P(\theta, \phi | x) \propto \phi^{n/2} \exp\left\{-\frac{\phi}{2} \left[\sum x_i^2 - 2n\bar{x}\theta + n\theta^2 \right] \right\} \left\{ \left(\frac{\phi}{K} \right)^{1/2} \exp\left\{-\frac{\phi}{2K} (\theta - \mu)^2\right\} \right\}$$

$$\phi^{a-1} \exp\{-b\phi\}$$

$$\propto \phi^{n/2+a+1/2-1} \exp\left\{-\frac{\phi}{2} \left[\sum x_i^2 - 2n\bar{x}\theta + n\theta^2 + \frac{\theta^2}{K} - \frac{2\theta\mu}{K} + \frac{\mu^2}{K} \right]\right\}$$

$$P(\phi | x) = \int P(\theta, \phi | x) d\theta$$

$$\propto \int \phi^{n/2+a+1/2-1} \exp\left\{-\frac{\phi}{2} \left[b + \sum x_i^2 - 2n\bar{x}\theta + n\theta^2 + \frac{\theta^2}{K} - \frac{2\theta\mu}{K} + \frac{\mu^2}{K} \right]\right\} d\theta$$

$$= \phi^{n/2+a+1/2-1} \exp\left\{-\frac{\phi}{2} \left[2b + \sum x_i^2 + \frac{\mu^2}{K} \right]\right\}$$

$$\left(\exp\left\{-\frac{\phi}{2} \left[\theta^2 \left(n + \frac{1}{K^2} \right) - 2\theta \left[n\bar{x} + \frac{\mu}{K} \right] \right]\right\} d\theta \right)$$

$$\boxed{\frac{1}{2\pi}} \left[\phi \left(1 + \frac{1}{K^2} \right) \right]^{-1/2}$$

$$\propto \phi^{n/2+a-1} \exp\left\{-\frac{\phi}{2} \left[2b + \sum x_i^2 + \frac{\mu^2}{K} \right]\right\}$$

$$\Rightarrow \phi | x \sim \text{Gamma} \left(\frac{n}{2} + a, b + \frac{\sum x_i^2}{2} + \frac{\mu^2}{2K} \right).$$

(50)

$$\left\{ \exp \left\{ -\frac{\phi}{2} \left[\theta^2 \left(n + \frac{1}{K} \right) - 2\theta \left[n\bar{x} + \frac{\mu}{K} \right] \right] \right\} d\theta \right.$$

$$= \left\{ \exp \left\{ -\frac{\phi}{2} \left(n + \frac{1}{K} \right) \left[\left(\theta - \frac{n\bar{x} + \frac{\mu}{K}}{n + \frac{1}{K}} \right)^2 - \left(\frac{n\bar{x} + \frac{\mu}{K}}{n + \frac{1}{K}} \right)^2 \right] \right\} d\theta \right.$$

$$= \exp \left\{ \frac{\phi}{2} \frac{\left(n\bar{x} + \frac{\mu}{K} \right)^2}{\left(n + \frac{1}{K} \right)} \right\} \underbrace{\left\{ \exp \left\{ -\frac{\phi}{2} \left(n + \frac{1}{K} \right) \left[\theta - \frac{n\bar{x} + \frac{\mu}{K}}{n + \frac{1}{K}} \right]^2 \right\} d\theta \right\}}$$

$$\sqrt{2\pi} \left[\phi \left(n + \frac{1}{K} \right) \right]^{-1/2}$$

$$\Rightarrow p(\phi | \mathbf{x}) \propto \phi^{\frac{n}{2} + \alpha - 1} \exp \left\{ -\frac{\phi}{2} \left[2b + \sum_{i=1}^n x_i^2 + \frac{\mu^2}{K} - \frac{\left(n\bar{x} + \frac{\mu}{K} \right)^2}{\left(n + \frac{1}{K} \right)} \right] \right\}$$

$$\Rightarrow \phi | \mathbf{x} \sim \text{Gamma} \left(\frac{n}{2} + \alpha, b + \frac{1}{2} \sum_{i=1}^n x_i^2 + \frac{1}{2} \frac{\mu^2}{K} - \frac{1}{2} \frac{\left(n\bar{x} + \frac{\mu}{K} \right)^2}{\left(n + \frac{1}{K} \right)} \right)$$

(6)

$$P(\theta | \phi, \mathcal{X}) = \frac{P(\theta, \phi | \mathcal{X})}{P(\phi | \mathcal{X})}$$

$$P(\theta | \phi, \mathcal{X}) \propto \frac{\phi^{\frac{n}{2} + \alpha + \frac{1}{2} - 1} \exp\left\{-\frac{\phi}{2} [zb + \sum x_i^2 - 2n\bar{x}\theta + \theta^2 + \frac{\theta^2}{K} - \frac{2\theta\mu}{K} + \frac{\mu^2}{K}]\right\}}{\phi^{\alpha/2 + \alpha - 1} \exp\left\{-\frac{\phi}{2} [zb + \sum x_i + \frac{\mu^2}{K}]\right\}}$$

$$= \phi^{1/2} \exp\left\{-\frac{\phi}{2} \left[\theta^2 \left(n + \frac{1}{K} \right) - 2\theta \left(n\bar{x} + \frac{\mu}{K} \right) \right]\right\}.$$

$$\propto \exp\left\{-\frac{\phi}{2} \left[\theta^2 \left(n + \frac{1}{K} \right) - 2\theta \left(n\bar{x} + \frac{\mu}{K} \right) \right]\right\}$$

$$\theta | \phi, \mathcal{X} \sim N\left(\frac{n\bar{x} + \frac{\mu}{K}}{\phi \left(n + \frac{1}{K} \right)}, \left[\phi \left(n + \frac{1}{K} \right) \right]^{-1}\right)$$

→ check!!!

①

webcast.ucsc.edu

Back on our survey example:

How would we obtain point and interval estimates from a classical perspective for this same problem.

Point estimator: $\hat{\theta} = \frac{\sum y_i}{n}$ (Maximum likelihood est.).

Interval estimator: $\hat{\theta} \pm 2 \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$ \Rightarrow For a $\approx 95\%$ confidence interval.

The ^{classical} results are not too far away from the Bayesian. In fact, as n increases the differences diminish.

More generally, under some mild assumptions, point and interval estimates constructed using Bayesian methods behave like their frequentist counterparts.

(2)

Batch vs joint update.

Suppose you get an additional sample of $m=63$ women with $\sum x_i = 17$ positive answers.

Batch update.

$$\theta \sim \text{Beta}(a=10, b=30)$$

$$\begin{array}{c} n=129 \\ \downarrow \\ \sum y_i = 25 \end{array}$$

$$\theta \sim \text{Beta}\left(a = \frac{10+25}{35}, b = \frac{30+129-25}{134}\right)$$

Because of \Rightarrow
conjugacy
you stay in
the same
family.

$$\theta \sim \text{Beta}(a=52, b=180)$$

$$P_1(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)} = \frac{P(y|\theta)P(\theta)}{\int P(y|\theta)P(\theta)d\theta}$$

$$P_2(\theta|y, x) = \frac{P(x|\theta, y)P(\theta|y)}{\int P(x|\theta, y)P(\theta|y)d\theta}$$

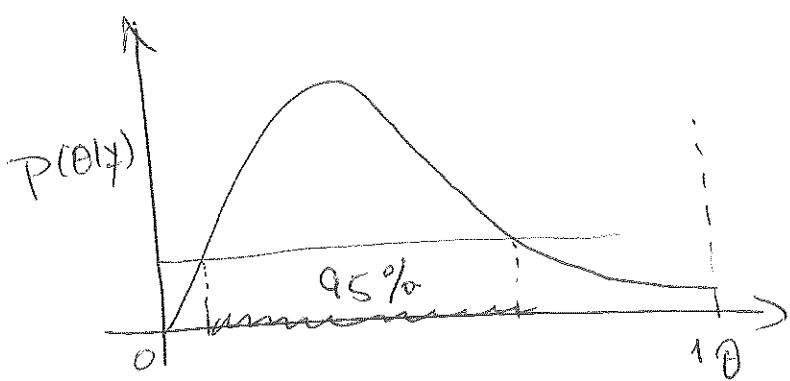
Note that $P(x|\theta, y) = P(y|\theta)$ in this case.

$$\frac{P(x|\theta)P(y|\theta)P(\theta)/P(y)}{\int P(x|\theta)P(y|\theta)P(\theta)/P(y)d\theta}$$

$$\frac{P(x|\theta)P(y|\theta)P(\theta)/P(y)}{\int P(x|\theta)P(y|\theta)P(\theta)/P(y)d\theta}$$

$$\frac{P(y|\theta) P(x|\theta) P(\theta)}{\int P(y|\theta) P(x|\theta) P(\theta) d\theta} = \frac{P(x, y|\theta) P(\theta)}{\int P(x, y|\theta) P(\theta) d\theta}. \quad (3)$$

HPD vs symmetric intervals:



HPD interval = shorter interval
that contains a pre-specified
amount of probability
(Highest Posterior Density)

They satisfy $P(\ell|y) = P(v|y)$
(in unimodal settings)

$$\int_a^v P(\theta|y) = 0.95$$

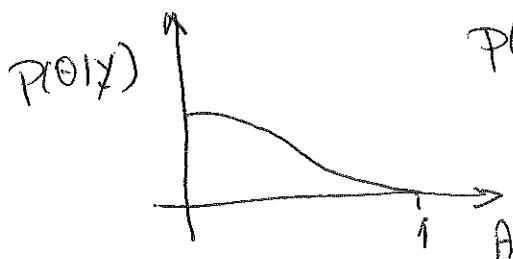
In symmetric posteriors, HPD and symmetric intervals
are the same.

$$\theta \sim \text{Beta}(1, 1) = \text{Uniform}[0, 1]$$

$$n = 129$$

$$\sum y_i = 0$$

$$\Rightarrow \theta|y \sim \text{Beta}(1, 130)$$



$$\begin{aligned} P(\theta|y) &= \frac{\Gamma(131)}{\Gamma(130)} \theta^{129} (1-\theta)^{129} \\ &= 130 (1-\theta)^{129} \end{aligned}$$

Mean / median / mode?

(4)

↳ utility functions / decision theory

Post Means are not invariant to transformation.

$$E(g(\theta)) \neq g(E(\theta))$$

However the median is invariant to monotone transformations.

Thinking about means of multiple parameters is easy;
defining a ^(joint) median for various parameters is tricky

In practice I will tend to use the posterior mean.

Hypothesis testing in two populations.

(5)

Why is ADHD diagnosis going up for younger kids.

Two populations (before and after cutoff date)

$$x = 15$$

$$y = 12$$

$$n = 87$$

$$m = 63$$

Before

After

$$x \sim \text{Bin}(n, \theta_1)$$

$$y \sim \text{Bin}(m, \theta_2)$$

Contrast $H_0: \theta_1 = \theta_2$

$H_a: \theta_1 \neq \theta_2$

From a Bayesian perspective we are going to treat the hypotheses as random quantities.

$$P(H_0) = 1 - P(H_a)$$

$$P(x, y | H_0, \theta_1, \theta_2)$$

$$P(\theta_1, \theta_2 | H_0)$$

$$P(x, y | H_a, \theta_1, \theta_2)$$

$$P(\theta_1, \theta_2 | H_a).$$

$$P(H_0 | x, y) = \frac{P(x, y | H_0) P(H_0)}{P(x, y | H_0) P(H_0) + P(x, y | H_a) P(H_a)}$$

$$P(x, y | H_0) = \int P(x, y | \theta_1, \theta_2, H_0) P(\theta_1, \theta_2 | H_0) d\theta_1 d\theta_2$$

①

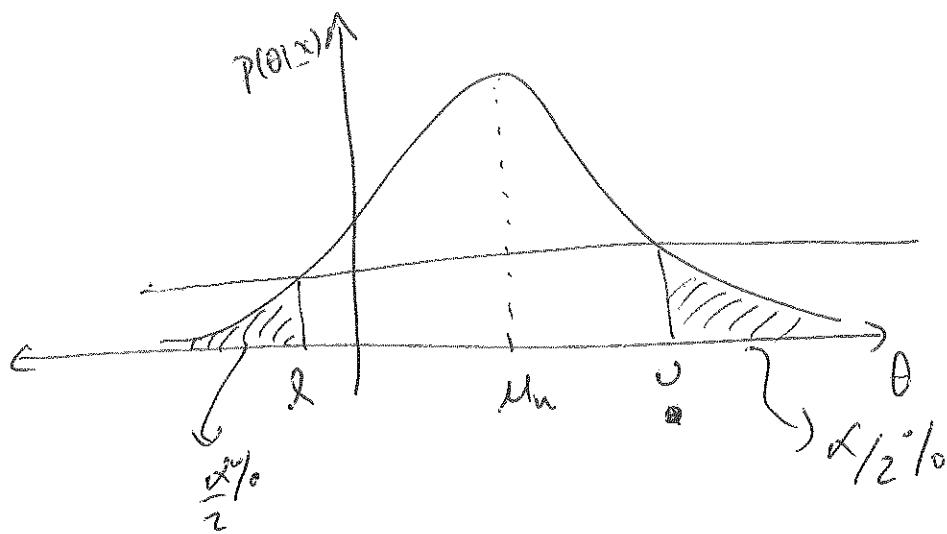
Credible interval

$$x_i - \theta \sim N(0, \sigma^2) \quad i=1, \dots, n$$

$$\theta \sim N(\mu, \tau^2)$$

$$\theta|x \sim N\left(\frac{\frac{n\bar{x}}{\sigma^2} + \frac{\mu}{\tau^2}}{\frac{n}{\sigma^2} + \frac{1}{\tau^2}}, \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\tau^2}}\right)$$

μ_n σ_n^2



For a $(1-\alpha)\%$
HPD credible interval

$$\frac{\theta - \mu_n}{\sigma_n} \mid x \sim N(0, 1)$$

$$\Pr\left(\frac{\theta - \mu_n}{\sigma_n} > z_{\alpha/2}\right) = \frac{\alpha}{2}$$

$$\Pr\left(\frac{\theta - \mu_n}{\sigma_n} < -z_{\alpha/2}\right) = 1 - \frac{\alpha}{2}$$

$$\Pr\left(\frac{\theta - \mu_n}{\sigma_n} < z_{\alpha/2}\right) = 1 - \frac{\alpha}{2}$$

$$\Rightarrow z_{\alpha/2} = q_{\text{norm}}\left(1 - \frac{\alpha}{2}\right)$$

(2)

For $\alpha = 5\%$ then

$$z_{0.975} = q_{\text{norm}}(0.975) \approx 1.964$$

$$\begin{cases} U = \bar{X}_n + 1.964 S_n \\ L = \bar{X}_n - 1.964 S_n \end{cases}$$

(3) Exponential families and conjugacy:

Let $f(x|\theta)$ be of the form:

$$f(x|\theta) = C(\theta) h(x) \exp\left\{ R^T(\theta) \cdot T(x)\right\}$$

where $x \in \mathcal{X}$ and $\theta \in \mathbb{R}^k$ such that the form of x is independent of θ

f is said to belong to the exponential family.

A particularly interesting case arises when $R(\theta) \in \mathbb{R}^k$ and $T(x) \in \mathbb{R}^k$ in which case we can ~~rearrange~~ rewrite.

$$f(x|\theta) = C(\eta) h(x) \exp\left\{ \eta^T \cdot T(x)\right\}$$

$$\text{where } \eta = R(\theta).$$

η the natural parameter of the distribution.

Examples:

$$x_i|\theta \sim \text{Ber}(\theta) \quad i=1, \dots, n$$

$$\Rightarrow f(x|\theta) = \theta^{x_1} (1-\theta)^{n-x_1} = (1-\theta)^n \left(\frac{\theta}{1-\theta}\right)^{\sum x_i}$$

$$= (1-\theta)^n \exp\left\{ \sum_{i=1}^n x_i \right\} \cdot \log \frac{\theta}{1-\theta} \}$$

$$\eta = \log \frac{\theta}{1-\theta} \quad T(x) = \sum_{i=1}^n x_i \quad C(\theta) = (1-\theta)^n \quad h(x) = 1$$

(4)

~~Max Likelihood Estimation~~

$x_i | \theta \sim N(\theta, \sigma^2)$ θ, σ^2 ~~Known~~ Unknown

$$f(x_1, \dots, x_n | \theta, \sigma^2) = \left(\frac{1}{\sqrt{2\pi}} \right)^n \left(\frac{1}{\sigma^2} \right)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2 \right\}$$

$$\begin{aligned} &= \left(\frac{1}{\sqrt{2\pi}} \right)^n \left(\frac{1}{\sigma^2} \right)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} (\sum x_i^2 - 2n\bar{x}\theta + \theta^2) \right\} \\ &= \left(\frac{1}{\sqrt{2\pi}} \right)^n \left(\frac{1}{\sigma^2} \right)^{n/2} \exp \left\{ -\frac{\theta^2}{2\sigma^2} \right. \\ &\quad \left. - \frac{\sum x_i^2}{2\sigma^2} + \frac{n\bar{x}\theta}{\sigma^2} \right\} \end{aligned}$$

$K=2$

$$T(\mathbf{x}) = \left(\sum_{i=1}^n x_i^2, \sum_{i=1}^n x_i \right) \quad R(\theta, \sigma^2) = \left(\frac{1}{2\sigma^2}, \frac{\theta}{\sigma^2} \right)$$

(5)

~~Another example~~

Note that if x_1, \dots, x_n are iid from $f(x_i|\theta)$ and $f(x_i|\theta)$ belong to the exponential family, then $f(\underline{x}|\theta)$ also belongs to the exponential family.

$$f(\underline{x}|\theta) = \prod_{i=1}^n c(\theta) h(x_i) \exp\left\{ \theta^\top T(x_i) R(\theta) \right\}$$

$$= [c(\theta)]^n \left[\prod_{i=1}^n h(x_i) \right] \exp\left\{ \left(\sum_{i=1}^n T(x_i) \right)^\top R(\theta) \right\}$$

Now consider a prior of the form:

~~$p(\theta) = [c(\theta)]^\lambda$~~

$$p(\theta|\lambda, \mu) = [c(\theta)]^\lambda h(\mu, \lambda) \exp\left\{ \frac{\mu^\top \theta}{\lambda} \right\}$$

Is this prior in the exponential family?

Consider the likelihood.

(6)

$$f(x|\eta) = C(\eta) h(x) \exp\{T(x)^T \eta\}$$

$$P(\eta) = [C^*(\eta)]^{-1} h^*(\mu, \lambda) \exp\{\mu^T \eta\}$$

$$P(\eta|x) \propto [C^*(\eta)]^{\lambda+1} \exp\{[T(x) + \mu]^T \eta\}$$

$$\Downarrow$$

$$P(\eta|x) = [C^*(\eta)]^{\lambda+1} h^*(\mu + T(x), \lambda+1) \exp\{[T(x) + \mu]^T \eta\}$$

↳ conjugate prior!!!

$$X_i | \lambda \sim \text{Poisson}(\lambda) \quad i=1, \dots n$$

⑦

$$\lambda \sim \text{Gamma}(a, b)$$

$$P(X|\lambda) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} = \left(\prod_{i=1}^n \frac{1}{x_i!} \right) e^{-n\lambda} \underbrace{\lambda^{\sum x_i}}_{\exp\{\sum x_i \log \lambda\}}$$

~~$\eta = \log \lambda$~~

$$T(X) = \sum X_i$$

$$c^*(\eta) = \exp\{-n \exp\{\eta\}\}$$

$$P(\lambda) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp\{-b\lambda\} \quad \lambda = \exp\{\eta\} \\ d\lambda = \exp\{\eta\} b d\eta.$$

$$P(\eta) = \frac{b^a}{\Gamma(a)} \exp\{(a-1)\eta\} \{ \exp\{-\exp\{\eta\} b\} \}$$

Complete the matching of terms.

①

Conjugate priors:

$$f(x|\eta) = C(\eta) h(x) \exp\{\eta^T T(x)\} \Rightarrow \text{Natural parameter exponential family.}$$

↳ Conjugate prior

$$P(\eta|\lambda, x_0) \propto \exp\{\eta^T x_0 + \lambda \log c(\eta)\}$$

⇒ Posterior

$$P(\eta|x) \propto \exp\{\eta^T (x_0 + T(x)) + (\lambda+1) \log c(\eta)\}$$

Updated parameters are

$$\lambda_n = \lambda + 1$$

$$x_n = x_0 + T(x)$$

For the Poisson likelihood:

$$P(x_1, \dots, x_n | \lambda) = \frac{e^{-\lambda} \lambda^{\sum x_i}}{\prod_{i=1}^n x_i!} \propto \lambda^{\sum x_i} \exp\{-\lambda n\} = \exp\left\{\sum x_i \log \lambda - n\lambda\right\}$$

$$\eta = \log \lambda$$

$$P(x_1, \dots, x_n | \eta) \propto \exp\{\eta \sum x_i - n \exp\{\eta\}\}$$

$$P(\eta) \propto \exp\{\eta x_0 - \log \exp\{\eta\}\}.$$

(2)

$$P(\gamma | \underline{x}) \propto \exp \left\{ \gamma \left(\frac{x_0}{\lambda} + \sum x_i \right) - (\gamma + n) \exp \{ \gamma \} \right\}$$

↓

$$\exp \left\{ \gamma \left(\lambda \left(\frac{x_0}{\lambda} \right) + \sum x_i \right) - (\gamma + n) \exp \{ \gamma \} \right\}$$

Back transform:

$$\gamma = \log \lambda \quad d\gamma = \frac{d\lambda}{\lambda}$$

$$\Rightarrow P(\lambda | \underline{x}) \propto \exp \left\{ \log \lambda \left[\underbrace{x_0}_{\lambda} + \sum x_i \right] - \underbrace{\lambda}_{(\gamma+n)} \right\} \frac{1}{\lambda}$$

$$\propto \lambda^{x_0 + \sum x_i - 1} \exp \{ -(\gamma+n) \lambda \}$$

$$\Rightarrow \lambda | \underline{x} \sim \text{Gamma}(x_0 + \sum x_i - 1, \gamma + n)$$

Let $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} f(x_i | \theta)$ and $E(x) = \xi(\theta)$

with f in the exponential family with conjugate

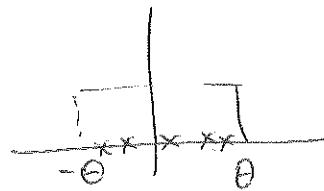
Prior $\pi(\theta) \propto \exp \{ \theta x_0 - \lambda \psi_0 \}$ then

$$E(\xi(\theta)) = \frac{x_0}{\lambda} \quad \text{and} \quad E(\xi(\theta) | \underline{x}) = \frac{x_0 + \cancel{\lambda \psi_0} + \cancel{n \bar{x}}}{\lambda + n}$$

Conjugacy outside the exponential family

(3)

$$X_1, \dots, X_n \stackrel{iid}{\sim} U[-\theta, \theta]$$



$$f(x_1, \dots, x_n | \theta) = \begin{cases} \left(\frac{1}{2\theta}\right)^n & \text{if } -\theta \leq x_i \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

$$= \begin{cases} \left(\frac{1}{2\theta}\right)^n & \text{if } \theta > \max_i \{ |x_i| \} \\ 0 & \text{otherwise} \end{cases}$$

~~not in the exponential family~~

∴ not in the exponential family.

$$f(x | \theta) = \begin{cases} \left(\frac{1}{2\theta}\right)^n & \theta > \max_i \{ |x_i| \} \\ 0 & \text{otherwise} \end{cases}$$

$$P(\theta) = \underbrace{\text{Beta}}_{\text{Beta}} \left(\frac{\beta x^\beta}{\theta^{\beta+1}} \right) \quad \theta > x \Rightarrow \text{Pareto}(\alpha, \beta)$$

$$P(\theta | x) \propto \left(\frac{1}{\theta}\right)^n \left(\frac{1}{\theta}\right)^{\beta+1} \quad \text{if } \theta > \max\{\alpha, \max_i \{ |x_i| \}\}$$

$$\alpha_n = \max\{\alpha, \max_i \{ |x_i| \}\} / \beta_n = \beta + n$$

$$\theta | \mathbf{x} \sim \text{Pareto}\left(\max\{\alpha, \max\{x_i\}\}, B+n\right) \quad (4)$$

\hookrightarrow Pareto is not in the exponential family either
(if you consider both α and B as parameters)

Closed form expression outside conjugacy.

$$x_1, \dots, x_n | \theta \stackrel{\text{iid}}{\sim} N(\theta, \sigma^2)$$

$$P(\theta) = \frac{1}{2\lambda} \exp\left\{-\frac{1}{\lambda} |\theta - \mu|\right\}$$

$$P(\mathbf{x} | \theta) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \frac{1}{\lambda} \exp\left\{-\frac{1}{2\sigma^2} \left[\sum x_i^2 - 2\theta n \bar{x} + n\theta^2 \right]\right\}.$$

$$P(\theta | \mathbf{x}) \propto \exp\left\{-\frac{1}{2\sigma^2} \left[n\theta^2 - 2n\theta \bar{x} \right] - \frac{1}{\lambda} |\theta - \mu|\right\}$$

$$\text{For } \theta > \mu \quad |\theta - \mu| = \theta - \mu$$

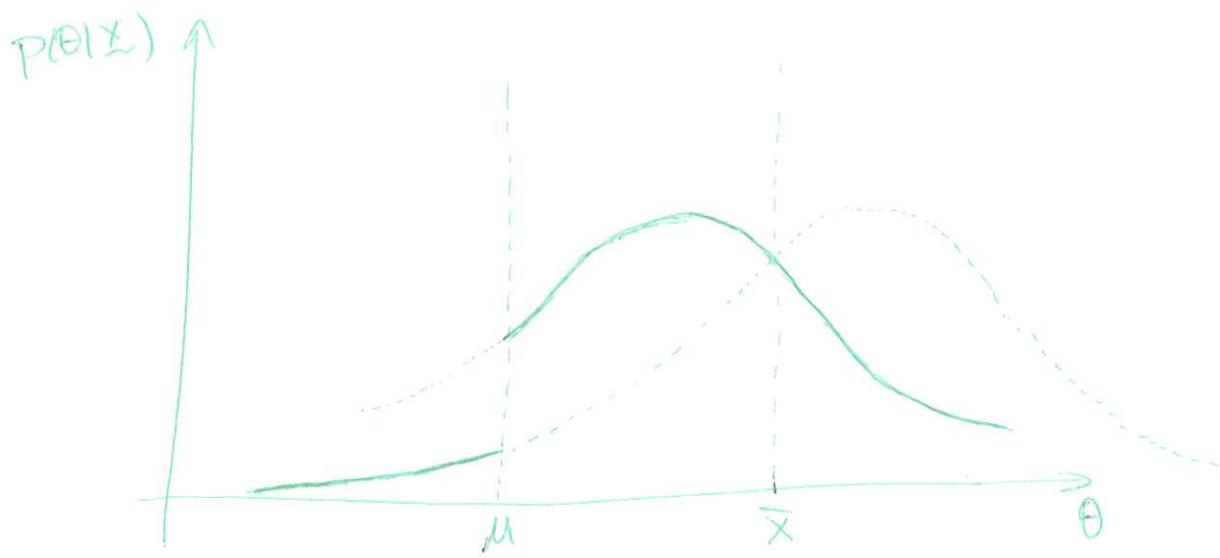
$$\Rightarrow P(\theta | \mathbf{x}) \propto \exp\left\{-\frac{1}{2} \left[\frac{n\theta^2}{\sigma^2} - \frac{2n\theta \bar{x}}{\sigma^2} - \frac{2}{\lambda} (\theta - \mu) \right]\right\}$$

$$\propto \exp\left\{-\frac{1}{2} \left[\frac{n\theta^2}{\sigma^2} - 2\theta \left[\frac{n\bar{x}}{\sigma^2} + \frac{1}{\lambda} \right] \right]\right\}.$$

$$\text{For } \theta < \mu:$$

$$P(\theta | \mathbf{x}) \propto \exp\left\{-\frac{1}{2} \left[\frac{n\theta^2}{\sigma^2} - 2\theta \left[\frac{n\bar{x}}{\sigma^2} - \frac{1}{\lambda} \right] \right]\right\}$$

5



$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{x^2}{\sigma^2}\right\} dx = \frac{1}{2} \quad ①$$

$$E[\mathbb{1}_{(-\infty, m)}(x)] \quad x \sim N(0, \sigma^2)$$

$$\mathbb{1}_{\Omega}(x) = \begin{cases} 1 & \text{if } x \in \Omega \\ 0 & \text{if } x \notin \Omega \end{cases}$$

How do we make it easier to compute expectations?
 Use the CLT (Central Limit Theorem) and the
 LLN (Law of Large Numbers)

LLN:

$$E(f(x)) \approx \int f(x) p(x) dx \stackrel{N \rightarrow \infty}{\approx} \sum_{i=1}^n f(x^{(i)}) \quad x^{(i)} \sim p$$

Computing integrals can be reduced to a problem of
 sampling from a probability distribution.

The power of simulation. Example:

②

		Test	
		Diseased	Healthy
Truth	Diseased	✓	Error
	Healthy	Error	✓

False Positive rate: $\Pr(+|\bar{D})$

Easy to get from data
but not what
you really care about

~~False~~ True positive rate: $\Pr(+|D)$

We care about $\Pr(D|+)$
 $\Pr(\bar{D}|-)$

$$\Pr(D|+) = \frac{\Pr(+|D)\Pr(D)}{\Pr(+|D)\Pr(D) + \Pr(+|\bar{D})\Pr(\bar{D})}$$

We will also need data on $\Pr(D)$ = incidence of the disease

Data:

x = # of positives in the test in a population of n_1 diseased individuals

y = # of positives in the test in a population of n_2 ~~healthy~~ individuals

z = # of diseased individuals in a sample of n subjects

(3)

independent

$$\begin{cases} X \sim \text{Bin}(n_1, \theta_1) \\ Y \sim \text{Bin}(n_2, \theta_2) \\ Z \sim \text{Bin}(m, \phi) \end{cases}$$

θ_1 : true positive rate = $\Pr(+|D)$
 θ_2 : false positive rate = $\Pr(+|\bar{D})$
 ϕ : incidence = $\Pr(D)$

η = True diagnostic probability = $\Pr(D|+)$
 η is my quantity of interest

$$\eta = \frac{\theta_1 \phi}{\theta_1 \phi + \theta_2 (1-\phi)}$$

Let's start with the interpretable parameters

$$\left. \begin{array}{l} \theta_1 \sim \text{Beta}(a_1, b_1) \\ \theta_2 \sim \text{Beta}(a_2, b_2) \\ \phi \sim \text{Beta}(a_3, b_3) \end{array} \right\} \text{independent}$$

A posteriori: they are independent

$$\left. \begin{array}{l} \theta_1|x \sim \text{Beta}(\overbrace{a_1+x}^{a_{1,n}}, \overbrace{b_1+n-x}^{b_{1,n}}) \\ \theta_2|y \sim \text{Beta}(a_2+y, b_2+n-y) \\ \phi|z \sim \text{Beta}(a_3+z, b_3+m-z) \end{array} \right\} \text{independent}$$

$$\begin{aligned} P(\theta_1, \theta_2, \phi | x, y, z) &\propto P(x|\theta_1)P(y|\theta_2)P(z|\phi)P(\theta_1)P(\theta_2)P(\phi) \\ &= [P(x|\theta_1)P(\theta_1)][P(y|\theta_2)P(\theta_2)][P(z|\phi)P(\phi)] \\ &= P(\theta_1|x)P(\theta_2|y)P(\phi|z) \end{aligned}$$

(4)

How would you proceed to obtain, for example, the posterior mean for η

$$\begin{aligned} \mathbb{E}(\eta | x, y, z) &= \int \int \int \frac{\theta_1 \phi}{\theta_1 \phi + \theta_2 (1-\phi)} \frac{\frac{\Gamma(a_{1n}+b_{1n})}{\Gamma(a_{1n}) \Gamma(b_{1n})}}{\sum (a_{1n})} \theta_1^{a_{1n}-1} (1-\theta_1)^{b_{1n}-1} \\ &\quad \frac{\frac{\Gamma(a_{2n}+b_{2n})}{\Gamma(a_{2n}) \Gamma(b_{2n})}}{\sum (a_{2n})} \theta_2^{a_{2n}-1} (1-\theta_2)^{b_{2n}-1} \\ &\quad \frac{\frac{\Gamma(a_{3n}+b_{3n})}{\Gamma(a_{3n}) \Gamma(b_{3n})}}{\sum (a_{3n})} \phi^{a_{3n}-1} (1-\phi)^{b_{3n}-1} d\theta_1 d\theta_2 d\phi \end{aligned}$$

$$\approx \frac{1}{N} \sum_{i=1}^N \frac{\theta_1^{(i)} \phi^{(i)}}{\theta_1^{(i)} \phi^{(i)} + \theta_2^{(i)} (1-\phi^{(i)})}$$

$$\theta_1^{(i)}, \theta_2^{(i)}, \phi^{(i)} \stackrel{iid}{\sim} P(\theta_1, \theta_2, \phi | x, y, z)$$

↓

$$\text{ind. } \left\{ \begin{array}{l} \theta_1^{(i)} \sim \text{Beta}(a_{1,n}, b_{1,n}) \\ \theta_2^{(i)} \sim \text{Beta}(a_{2,n}, b_{2,n}) \\ \phi^{(i)} \sim \text{Beta}(a_{3,n}, b_{3,n}) \end{array} \right.$$

The posterior of interest was "simple" in this case,
so we can directly simulate from the posterior. (5)

For most problems, however, the posterior might not have a simple form:

$$x_1, \dots, x_n \stackrel{\text{IID}}{\sim} N(\theta, \phi^{-1}) \quad \theta, \phi \text{ unknown.}$$

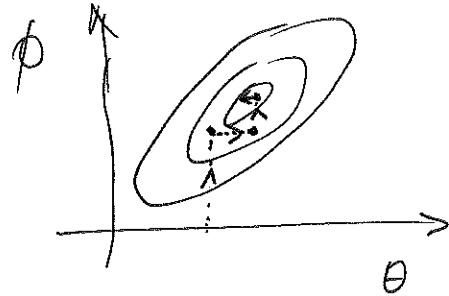
$$\begin{aligned} \theta &\sim N(\mu, \tau^2) \\ \phi &\sim \text{Gamma}(a, b) \end{aligned} \quad \left. \begin{array}{l} \Rightarrow \text{Not conjugate!!} \end{array} \right\}$$

Posterior:

$$P(\theta, \phi | \mathbf{x}) = \left(\frac{1}{2\pi} \right)^n \phi^{-n/2} \exp \left\{ -\frac{\phi}{2} \left(\sum x_i^2 - 2\theta n \bar{x} + n\theta^2 \right) \right\} \cdot \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\tau^2} (\theta^2 - 2\theta\mu + \mu^2) \right\} \frac{b^a}{\Gamma(a)} \phi^{a-1} \exp \left\{ -b\phi \right\}$$

Let's go for a simpler task: find posterior mode.

$$\arg \max_{(\phi, \theta)} P(\theta, \phi | \mathbf{x}) = \arg \max_{(\phi, \theta)} \log P(\theta, \phi | \mathbf{x})$$



ACM = Alternate conditional modes.

Diagnostics for MCMC algorithms

- Recall that, although Metropolis-Hastings algorithm work if we run them long enough, the fact that samples are dependent presents challenges:
 - Convergence.
 - Mixing.

Diagnostics for MCMC algorithms

- **Mixing:** As with every Monte Carlo algorithm an important question is how many samples should be generated in order to have estimates of the parameters with a reasonable level of accuracy.
 - However, because samples are (most usually positively) correlated, standard theory (standard central limit theorems, Chebyshev's inequality, etc.) does not apply here!
 - Mixing is also related to how fast the high probability areas of the posterior distribution are explored.

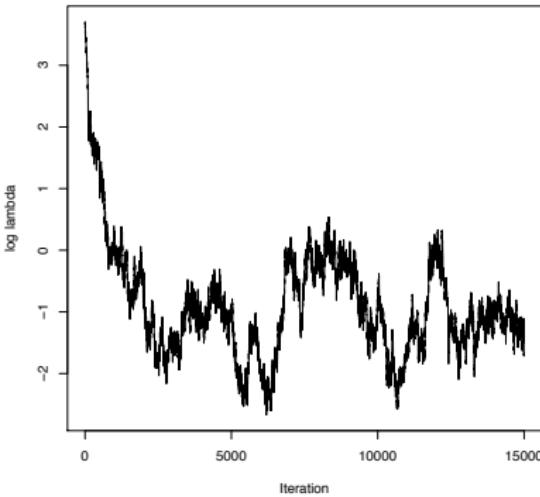
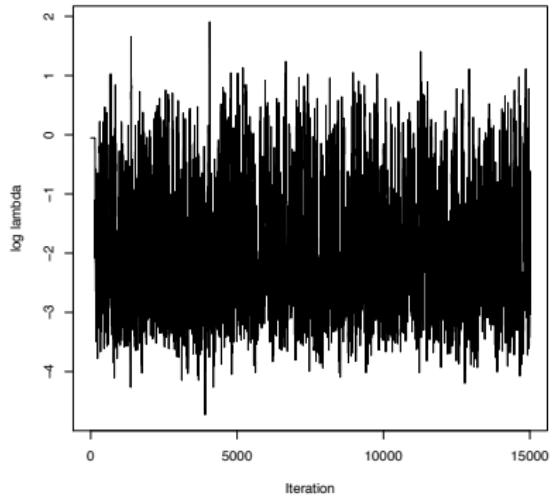
Diagnostics for MCMC algorithms

- **Convergence:** An additional question when using MCMC algorithms that is not a concern in standard Monte Carlo is how many of the initial observations need to be discarded before we can perform any inference.
 - Remember that, unless your initial state is sampled from the stationary distribution, the samples of your chain are only distributed as the stationary distribution when the number of iterations goes to infinity.
 - In practice, if your posterior is unimodal, you just want your initial state to be close to the high-probability areas.
 - When posterior are multimodal convergence can be a huge issue.

Diagnostics for MCMC algorithms

- These two questions are interrelated, but are not identical.
For example, you can accelerate convergence to the posterior by carefully selecting the initial state of your chain even if the chain mixes slowly. On the other hand, even if your initial state is really distributed from the stationary distribution (e.g., by using perfect sampling), you might need to run a very long chain to have an acceptable level of uncertainty in your estimates because the autocorrelation in the parameters of interest can be very high.
- Time series plots of the realized chains (the so-called *trace* plots) provide visual intuition about these two concepts, and can be used as an informal check (which should be supplemented with some of the formal techniques discussed later).

Diagnostics for MCMC algorithms



Checking for convergence

- Generally speaking, determining *a priori* how many iterations are needed for burn-in is not possible. Indeed, except for some special cases, general theory about the behavior of the algorithm (e.g., the value of the largest eigenvalue of the kernel function) is not easy to derive.
- For this reason, most practical algorithms to monitor convergence look at the time-series behavior of a small number of selected parameters (or functions of parameters).
 - Selecting which parameters need to be monitored is tricky and model dependent. In high-dimensional models I usually monitor either the likelihood or the posterior, a small (random) subset of the main first-stage parameters, and a couple of important hyperparameters.
 - All criteria here are implemented in the R package CODA.
- We typically cannot ensure convergence, just argue that there is no evidence of lack of convergence.

Checking for convergence (Multi-chain methods)

- As the name suggests multi-chain methods use $M > 1$ realizations of the algorithm (each of length B , typically started from “over dispersed” values) and compare the output of the chains. Once the statistical properties of the realizations appear to be similar, the chains are deemed to have converged.
- One example of this type of approach was introduced by Gelman & Rubin (1992) who propose to monitor a statistic that compares the average within-chain variability against the between-chains variability:

$$R = \left(\frac{B - 1}{B} + \frac{M + 1}{M} \frac{Z_B}{W_B} \right) \frac{\nu_B}{\nu_B - 2},$$

where Z_B is an estimate of between-chain variability, W_B is an estimate of within-chain variability and ν_B is the approximate number of degrees of freedom.

Checking for convergence (Multi-chain methods)

If $\zeta = h(\theta)$ is the quantity of interest being monitored, define

$$\bar{\xi}_m = \frac{1}{B} \sum_{b=1}^B \xi_m^{(b)} \quad \bar{\xi} = \frac{1}{M} \sum_{m=1}^M \bar{\xi}_m \quad s_m^2 = \frac{1}{B} \sum_{b=1}^B \left(\xi_m^{(b)} - \bar{\xi} \right)^2$$

Then

$$Z_B = \frac{1}{M} \sum_{m=1}^M \left(\bar{\xi}_m - \bar{\xi} \right)^2 \quad W_B = \frac{1}{M} \sum_{m=1}^M s_m^2$$

and

$$s_m^2 = \frac{1}{B} \sum_{b=1}^B \left(\xi_m^{(b)} - \bar{\xi} \right)^2 \quad \nu_B = \frac{2 \left(\frac{B-1}{B} W_B + \frac{M+1}{M} Z_B \right)^2}{W_B}$$

Checking for convergence (Multi-chain methods)

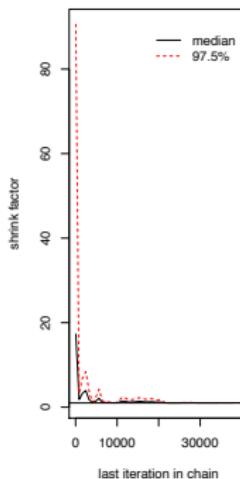
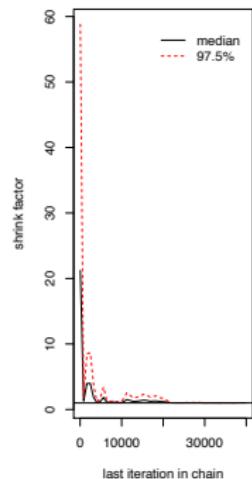
- We have $BZ_B/W_B \sim F_{M-1, 2W_B^2/\varpi_B}$ approximately, where

$$\varpi_B = \frac{1}{M^2} \left[\sum_{m=1}^M s_m^4 - \frac{1}{M} \left(\sum_{m=1}^M s_m^2 \right)^2 \right]$$

- Once the chains have converged we expect $R \rightarrow 1$, so the approximate distribution can be used to test convergence at different values of B .
- Samples from different chains are then combined together for computing any expectation of interest.

Checking for convergence (Multi-chain methods)

```
#... Run algorithm with the first set of initial  
values  
  
chain1.mcmc = as.mcmc(cbind(log(lambda.out),  
log(kappa2.out)))  
  
#... Run algorithm with the second set of  
initial values  
  
chain2.mcmc = as.mcmc(cbind(log(lambda.out),  
log(kappa2.out)))  
  
x = as.mcmc.list(list(chain1.mcmc, chain2.mcmc))  
  
x.gr = gelman.diag(x, autoburnin=F)  
  
gelman.plot(x)
```



Checking for convergence (Multi-chain methods)

- I tend to prefer multi-chain methods. One advantage is that they can be naively parallelized.
- However, finding over-dispersed initial points can be hard. For example, if the posterior is multimodal and you do not start the chain at least once on each mode you might never find out.
- For slow mixing chains, inferences based on a single long chain might be more accurate than the inferences based on many small chains put together. (However, this is an issue mostly if computational resources are limited.)

Checking for convergence (Single-chain methods)

- An alternative to multi-chain methods is to use a single long chain and see if the properties of the chain have “stabilized”.
- More specifically, we can divide the original chain of length B into M sections of approximately equal length and compare their spectral densities.
- In particular Geweke (1992) proposes to compare the first and second half of the chain (i.e., use $M = 2$). The algorithm is implemented in the R function `geweke.diag` and `geweke.plot`.

Mixing

- Again, I want to emphasize that mixing is a subtly different problem from convergence.
- Note that a chain that mixes slowly, if started on a low-probability region of the space, might take a long time to converge. However, if we start it in a high probability region convergence is much less of an issue.
- A simple way to decide how long the algorithm must be run *after* the burn-in period is to evaluate the variance of the estimators.

Mixing

- Let h be an integrable function. Recall that we aim at approximating $E_{\theta|y}\{h(\theta)\} \approx \frac{1}{B} \sum_{b=1}^B h(\boldsymbol{\theta}^{(b)})$
- If the samples $\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(B)}$ were independent, the variance of the estimator would be equal to $\frac{\text{Var}\{h(\theta)\}}{B}$. In practice a rough estimate of $\text{Var}\{h(\theta)\}$ can also be obtained by Monte Carlo.
- From Markov chain theory, it is easy to show that for an MCMC that has converged the variance of the same estimator is $\tau_h \frac{\text{Var}h(\theta)}{B}$ where τ_h is the integrated autocorrelation:

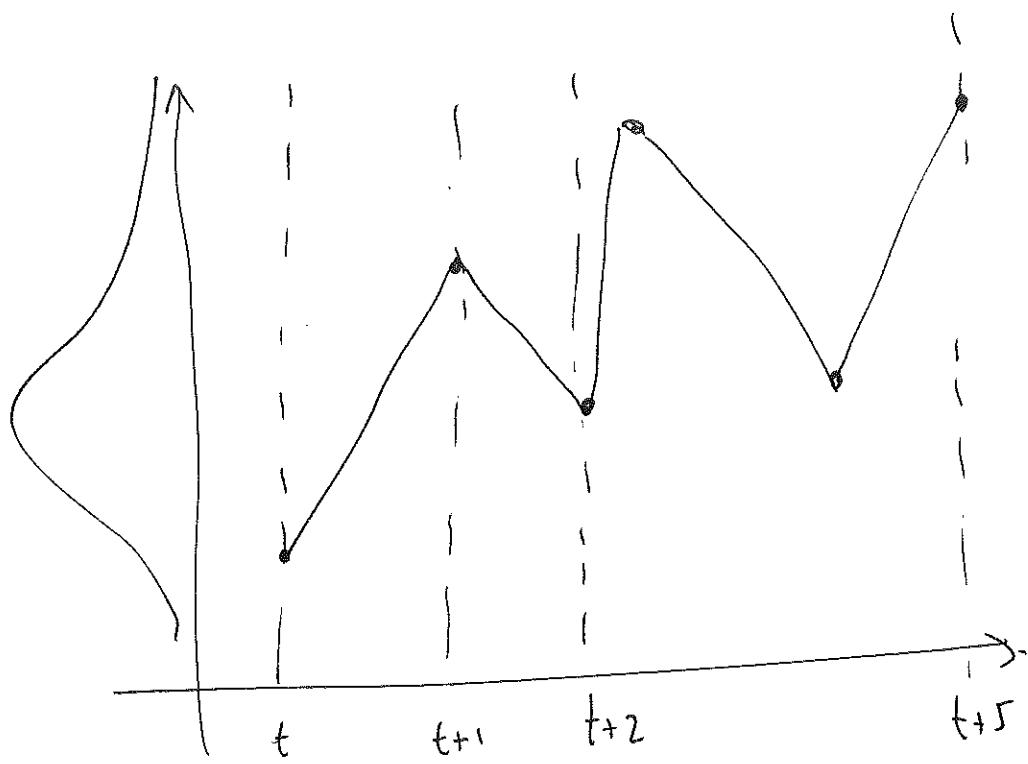
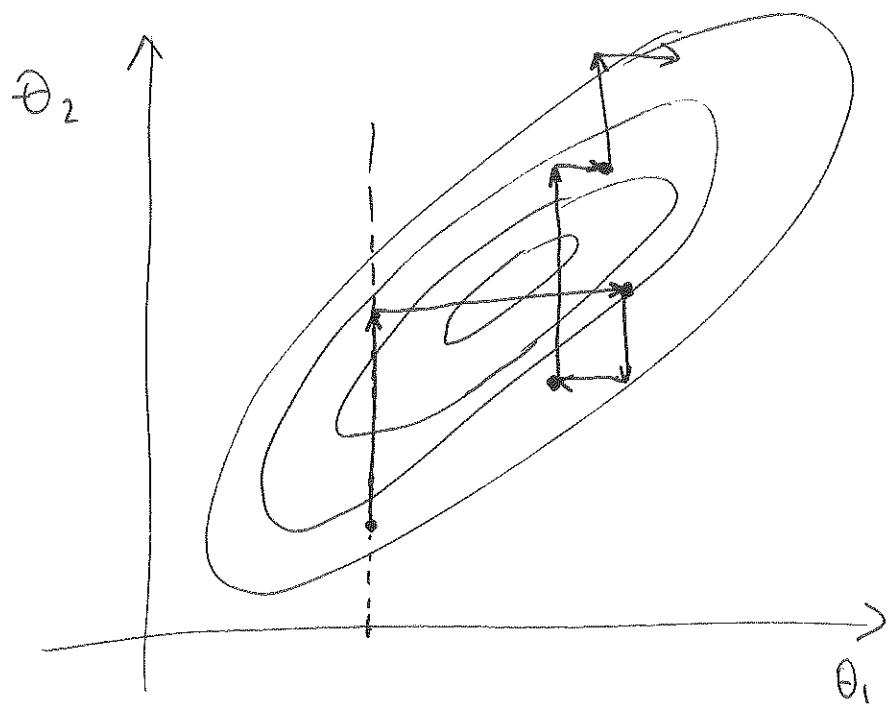
$$\tau_h = 1 + 2 \sum_{b=1}^{\infty} \text{Cor} \left\{ h(\boldsymbol{\theta}^{(1)}), h(\boldsymbol{\theta}^{(b)}) \right\}$$

- We call $\frac{B}{\tau_h}$ the equivalent sample size.

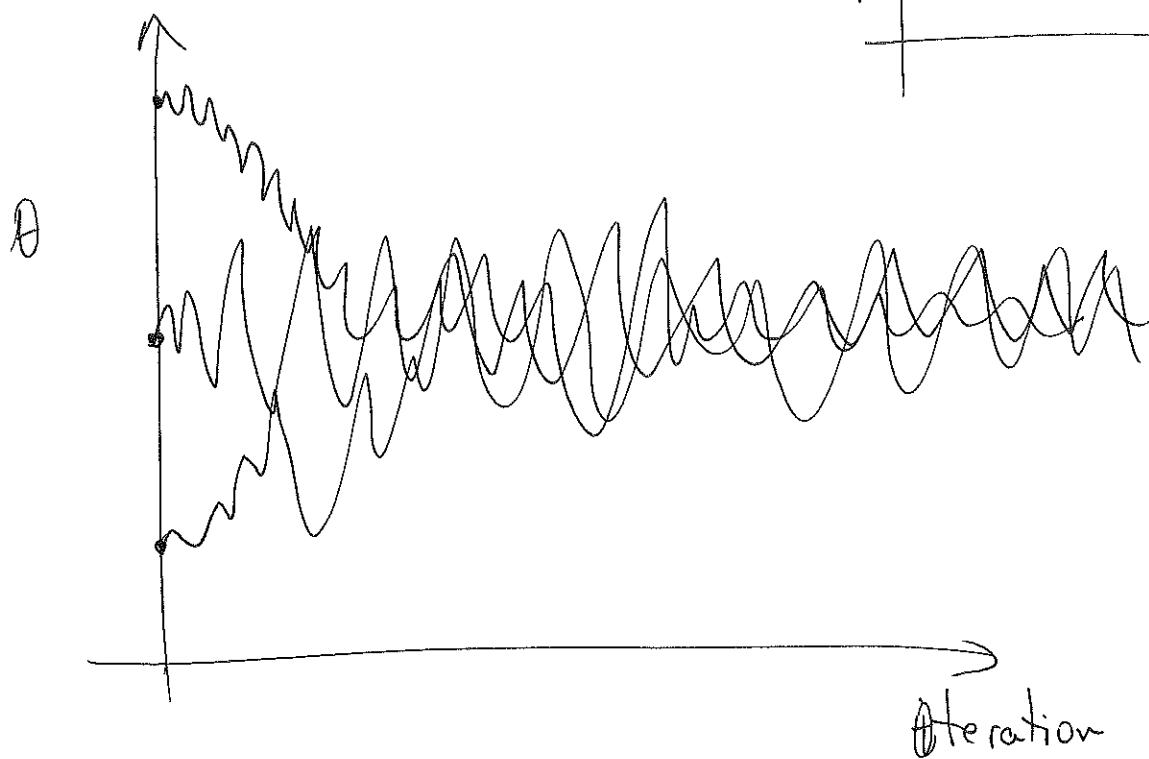
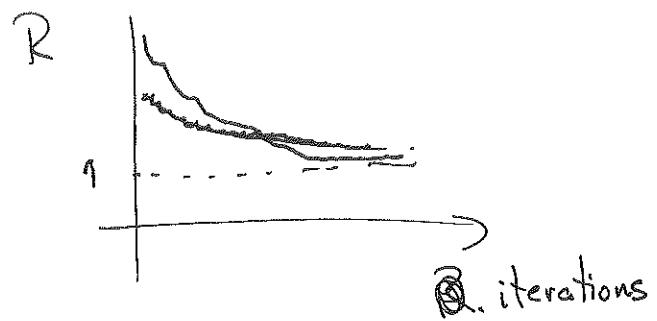
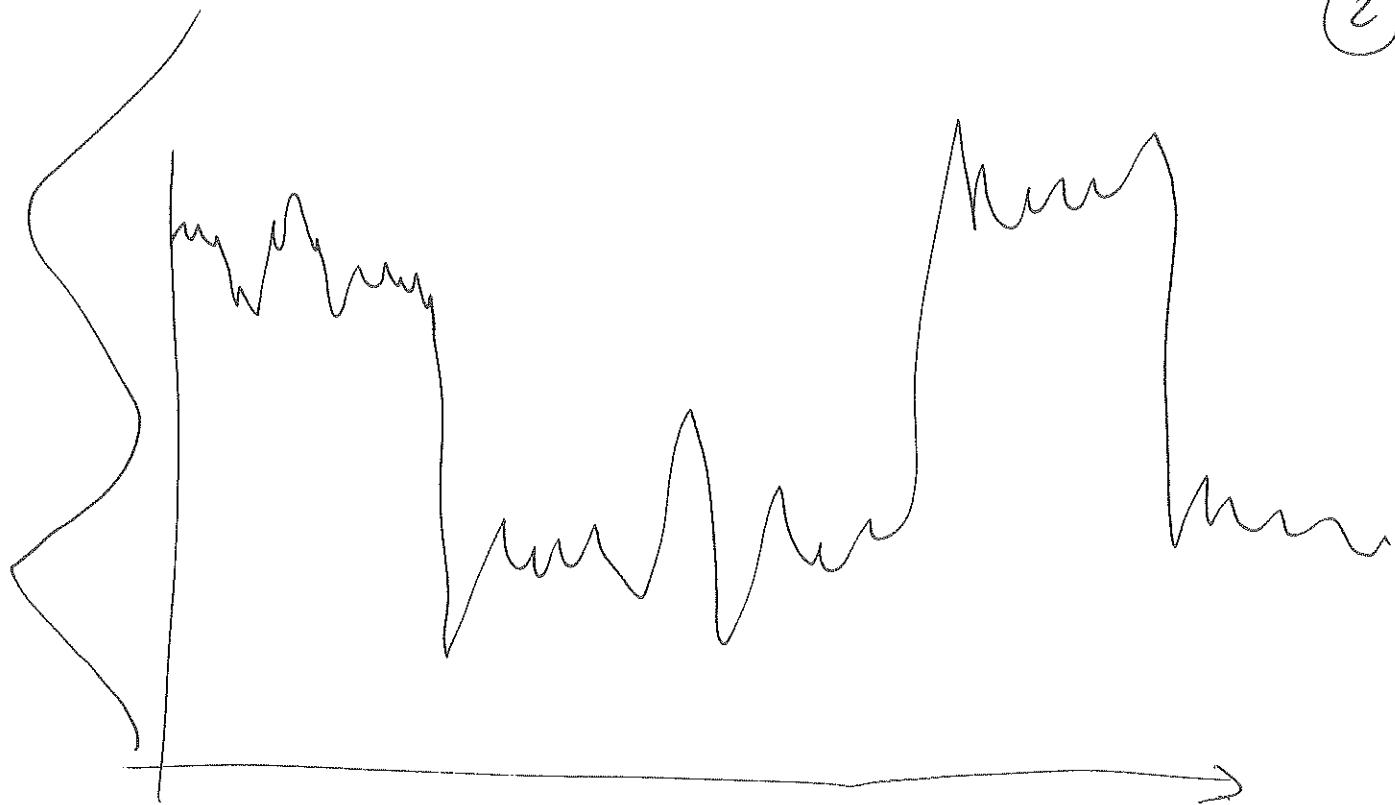
Mixing

- The equivalent sample size calculation is implemented in CODA through the function `effectiveSize`.
- Note that the equivalent sample size will be different for different parameters. That means that a given number of iterations might be enough for one parameter, but not for another.
- Some people like to “thin” their chain (retain every X number of iterations and drop the rest) as a way to decrease τ_g . However, this also reduces B !
 - In general, the tradeoff is not favorable.
 - In my opinion, unless storage is an issue, thinning is not a good idea.

①



(2)



(3)

X, Y be two random variables

$E(X+Y) = E(X) + E(Y)$ no matter whether X, Y are independent or not

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

↳ If $\underline{X_1, \dots, X_n}$ independent and identically distributed!

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \frac{n \text{Var}(X_i)}{n^2} = \frac{\text{Var}(X_i)}{n}.$$

(4)

Prior selection:

1) Subjective Bayes.

2) Objective Bayes

 ↳ { Invariance.
 "Minimal" information.

3) Robust Bayes.

 ↳ Bounded influence.

0.1

y_1, y_2, \dots, y_n
 \Downarrow
 data.

μ
 ρ
 σ^2

$$P(\mu, \rho, \sigma^2 | y_1, \dots, y_n) \leftarrow \text{ultimate goal}$$

$$\propto P(y_1, \dots, y_n | \mu, \rho, \sigma^2) P(\mu, \rho, \sigma^2)$$

\Downarrow

$$P(\mu) P(\rho) P(\sigma^2)$$

\Downarrow

$$P(y_t | \mu, \rho, \sigma^2, y_{t-1})$$

Gaussian

Unit [0, 1]

\Downarrow

Inverse Gamma

$$P(y_n, \dots, y_2 | y_1, \mu, \rho, \sigma^2) = P(y_n | y_{n-1}, \dots, y_1, \mu, \rho, \sigma^2) \times$$

$$P(y_{n-1} | y_{n-2}, \dots, y_1, \mu, \rho, \sigma^2) \times$$

... . . .

$$P(y_3 | y_2, \dots, y_1, \mu, \rho, \sigma^2) \times$$

$$P(y_2 | y_1, \mu, \rho, \sigma^2)$$

$$= \prod_{t=2}^n P(y_t | y_{t-1}, \mu, \rho, \sigma^2) = \prod_{t=2}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_t - \mu - \rho(y_{t-1} - \mu))^2 \right\}$$

$$y_t = \mu + \rho(y_{t-1} - \mu) + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma^2) \quad (0.2)$$

$$\hat{y}_t \sim N(\mu, \rho, y_{t-1}, \sigma^2) \sim N(\mu + \rho(y_{t-1} - \mu), \sigma^2)$$

(1)

Priors:

Subjective priors "Objective" priors Robust priors	Invariance Minimal information
--	-----------------------------------

Invariant priors: Jeffreys' priors

(Consistency across parameterizations)

Expected Fisher information:

Likelihood for your problem is $p(x|\theta)$ then

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

Example: $X \sim \text{Bin}(n, \theta)$

$$p(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

$$\Rightarrow l = \log p(x|\theta) = \log \binom{n}{x} + x \log \theta + (n-x) \log (1-\theta)$$

$$\frac{\partial l}{\partial \theta} = \frac{x}{\theta} - \frac{n-x}{1-\theta} \Rightarrow E_{x|\theta} \left[-\frac{x}{\theta^2} - \frac{(n-x)}{(1-\theta)^2} \right] =$$

$$\begin{aligned} \frac{\partial^2 l}{\partial \theta^2} &= -\frac{x}{\theta^2} - \frac{n-x}{(1-\theta)^2} \\ &= -\left[\frac{n\theta}{\theta^2} + \frac{n(1-\theta)}{(1-\theta)^2} \right] = \frac{n(1-\theta+\theta)}{\theta(1-\theta)} \\ &= \frac{n}{\theta(1-\theta)} \end{aligned}$$

The Jeffreys' prior is proportional to the squared root of the information number

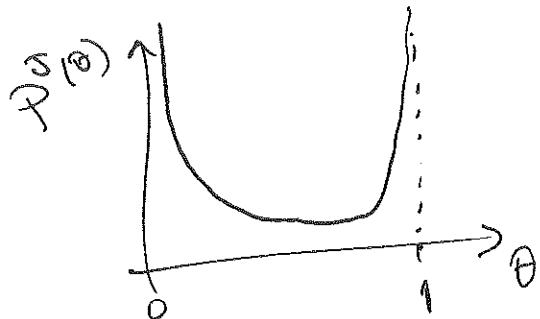
(2)

$$\tilde{P}^J(\theta) \propto \sqrt{I(\theta)}$$

For the binomial:

$$P^J(\theta) \propto \theta^{-1/2} (1-\theta)^{-1/2}$$

\Rightarrow The Jeffreys' prior for θ and a Binomial likelihood is a Beta($1/2, 1/2$)



$$\tilde{P}^J(\theta) = \frac{\Gamma(1)}{\Gamma(1/2)\Gamma(1/2)} \theta^{-1/2} (1-\theta)^{-1/2}$$

$$\gamma = \log \frac{\theta}{1-\theta}$$

What is the Jeffreys' prior for γ ?

(3)

Procedure 1:

change of variables on $\pi^J(\theta)$.

$$\eta = \log \frac{\theta}{1-\theta} \Rightarrow \exp\{\eta\} = \frac{\theta}{1-\theta} \Rightarrow \exp\{\eta\} - \theta \exp\{\eta\} = \theta$$

$$\Rightarrow \theta = \frac{\exp\{\eta\}}{1 + \exp\{\eta\}}$$

$$\Rightarrow \theta = \frac{1}{1 + \exp\{-\eta\}}$$

$$= (1 + \exp\{-\eta\})^{-1}$$

$$d\theta = (+1)(1 + \exp\{-\eta\})^{-2} \exp\{-\eta\}$$

$$= \frac{\exp\{-\eta\}}{(1 + \exp\{-\eta\})^2}$$

$$\hat{P}^J(\eta) = \left[\Gamma(1/2) \right]^2 \left[\frac{1}{1 + \exp\{-\eta\}} \right]^{-1/2} \left[\frac{\exp\{-\eta\}}{1 + \exp\{-\eta\}} \right]^{-1/2} \left[\frac{\exp\{-\eta\}}{1 + \exp\{-\eta\}} \right]^2$$

$$= \left[\Gamma(1/2) \right]^2 \frac{\exp\{-\frac{1}{2}\eta\}}{1 + \exp\{-\eta\}}$$

Homework: check that procedure 2 gives the same prior

$$p(x|\eta) = \binom{n}{x} \left(\frac{1}{1 + \exp\{-\eta\}} \right)^x \left(\frac{\exp\{-\eta\}}{1 + \exp\{-\eta\}} \right)^{n-x}$$

$$\pi^J(\eta) \propto \left[-\mathbb{E}_{x|M} \left[\frac{\partial^2}{\partial \eta^2} \log p(x|\eta) \right] \right]^{1/2}$$

Consider now $X \sim \text{Neg Bin}(K, \theta)$ (4)
 ↗ number of successes.
 ↗ # of trials

$$P(X|\theta) = \binom{X-1}{K-1} \theta^K (1-\theta)^{X-K}$$

$$l = \log P(X|\theta) = \log \binom{X-1}{K-1} + K \log \theta + (X-K) \log(1-\theta)$$

$$\frac{\partial l}{\partial \theta} = \frac{K}{\theta} - \frac{X-K}{1-\theta}$$

$$\frac{\partial^2 l}{\partial \theta^2} = -\frac{K}{\theta^2} - \frac{(X-K)}{(1-\theta)^2}$$

$$\mathbb{E}_{X|\theta} \left(\frac{\partial^2 l}{\partial \theta^2} \right) = -\mathbb{E}_{X|\theta} \left[\frac{K}{\theta^2} + \frac{(X-K)}{(1-\theta)^2} \right] = -\left[\frac{K}{\theta^2} + \frac{E(X)-K}{(1-\theta)^2} \right]$$

$$E(X) = \frac{K}{\theta} \quad \Rightarrow \quad I(\theta) = \frac{K}{\theta^2} + \frac{K(\frac{1}{\theta}-1)}{(1-\theta)^2}$$

$$= \frac{K((1-\theta)^2 + \theta^2(\frac{1}{\theta}-1))}{\theta^2(1-\theta)^2} = K \frac{1-2\theta+\theta^2+\theta^{-2}}{\theta^2(1-\theta)^2}$$

$$= K \frac{(1-\theta)}{\theta^2(1-\theta)^2} = \frac{K}{\theta^2(1-\theta)}$$

$$P^J(\theta) \propto \frac{1}{\theta(1-\theta)^{1/2}} = \theta^{-1} (1-\theta)^{-1/2}$$

↳ Improper prior!! (Does not integrate to anything finite!!!
 Not a density!!!)

(5)

Even though the prior is improper;

$$P(\theta|x) \propto \theta^k (1-\theta)^{x-k} \theta^{-1} (1-\theta)^{-1/2}$$

$$\Rightarrow P(\theta|x) \propto \theta^{k-1} (1-\theta)^{x-k+1/2}$$

↳ Kernel of a Beta($k, x-k+1/2$)

General result:

- * For location parameters the invariant prior is uniform.
- * For scale parameters the invariant prior is uniform on the log scale.

Example: $x_1, \dots, x_n \sim N(\mu, \sigma^2)$

- * If μ unknown but σ^2 known.

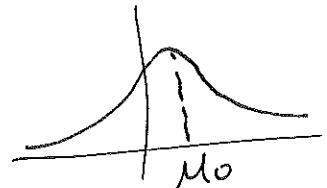
$$P^J(\mu) \propto 1$$

- * If σ^2 is unknown but μ is known

$$P^J(\sigma) \propto \frac{1}{\sigma}$$

For many models in the exponential family you can get the Jeffreys' prior as a limit of conjugate priors

$$P(\mu | \mu_0, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(\mu - \mu_0)^2\right\}$$



Objective Bayes: Jeffreys prior

- For location families where $p(y | \theta) = f(y - \theta)$ for some density f , the Jeffreys prior is simply

$$\pi^{J,L}(\theta) \propto 1.$$

- For scale families where $p(y | \theta) = \frac{1}{\theta} f\left(\frac{y}{\theta}\right)$ for some density f , the Jeffreys prior is

$$\pi^{J,S}(\theta) \propto \frac{1}{\theta}.$$

This is equivalent to placing a uniform prior on the log scale.

- In both cases the Jeffreys prior is improper, and we need to carefully check that the corresponding posterior is proper!

Objective Bayes: Jeffreys prior

- The Jeffreys priors for location and scale families can be obtained as the limit of proper priors:
 - For the location family, consider a Gaussian prior

$$\pi^L(\theta) = \frac{1}{\sqrt{2\pi}\tau} \exp \left\{ -\frac{1}{2\tau^2}(\theta - \mu)^2 \right\}$$

Note that $\lim_{\tau \rightarrow \infty} \pi^L(\theta) = \pi^{J,L}(\theta)$.

- For the scale family, consider a (slightly reparameterized) inverse Gamma family

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

where $E\{\theta\} = \beta$. Here $\lim_{\alpha \rightarrow 0} \pi^S(\theta) = \pi^{J,S}(\theta)$.

Objective Bayes: Jeffreys prior

- **Example (Normal distribution with unknown mean and variance):** Let $y | \mu, \phi \sim N_p(\mu, \phi)$, so that $\phi = \sigma^2$ is the variance. The Hessian in this case is

$$\begin{vmatrix} -\frac{1}{\phi} & -\frac{1}{\phi^2} \sum_{i=1}^n (y_i - \mu) \\ \frac{1}{\phi^2} \sum_{i=1}^n (y_i - \mu) & \frac{n}{2\phi^2} - \frac{1}{\phi^3} \sum_{i=1}^n (y_i - \mu)^2 \end{vmatrix}.$$

Since $E\{\sum_{i=1}^n (y_i - \mu)\} = 0$ and $E\{\sum_{i=1}^n (y_i - \mu)^2\} = n\phi$,

$$|\mathcal{I}(\mu, \phi)| = \frac{1}{\phi^3}$$

and

$$\pi^J(\mu, \phi) \propto \phi^{-3/2}.$$

Objective Bayes: Jeffreys prior

- Multivariate Jeffreys priors sometimes have undesirable behaviors.
- An alternative is the so-called independence Jeffreys prior (originally proposed by Jeffreys himself) is very common.
- The independence Jeffreys $\pi^{IJ}(\mu, \phi)$ assumes prior independence and uses the corresponding univariate Jeffreys priors, $\pi^J(\mu, \phi) = \pi^J(\mu)\pi^J(\phi)$.
- In the case of the normal model with unknown mean and variance this reduces to

$$\pi^{IJ}(\mu, \phi) \propto \frac{1}{\phi}.$$

- Independence Jeffreys priors are not invariant to transformations that involve more than one parameter.

Objective Bayes: Reference priors

- Another (similar) solution is to construct the prior sequentially.
- Let $y | \theta \sim p(y | \theta_1, \theta_2)$ where θ_1 is the parameter of interest and θ_2 is the nuisance parameter.
 - Use $p(y | \theta_1, \theta_2)$ to construct the Jeffreys prior for θ_2 assuming a fixed value of θ_1 to obtain $\pi^J(\theta_2 | \theta_1)$.
 - Compute $\tilde{p}(y | \theta_1) = \int p(y | \theta_1, \theta_2) \pi^J(\theta_2 | \theta_1) d\theta_2$.
 - Construct the Jeffreys prior for θ_1 using $\tilde{p}(y | \theta_1)$.
- The order of the parameters matters, and you might get a different prior depending on what is the main parameter of interest.
- Can be generalized to more than two groups of parameters.
- Can be justified as maximizing the posterior information.

Bayesian hypothesis testing

- Improper priors such as those that often arise from non-informative approaches are often appropriate for both point and interval estimation.
- However, their use in problems that involve models of different dimensions (as is the case in hypothesis testing) is extremely problematic.
- To illustrate this, let's consider Barlett's paradox.

Barlett's paradox

- **Example:** Let y_1, \dots, y_n be such that $y_i | \theta \sim N(\theta, 1)$ and consider testing the hypotheses $H_0 : \theta = 0$ against $H_1 : \theta \neq 0$ under the prior $\theta | H_1 \sim N(0, \tau^2)$.

The corresponding Bayes factor (which we already calculated in our introductory example) is

$$B_{10}(\tau^2, \bar{y}) = (1 + n\tau^2)^{-1/2} \exp \left\{ \frac{1}{2} \frac{n^2 \bar{y}^2}{n + \frac{1}{\tau^2}} \right\},$$

and it is easy to verify that $\lim_{\tau^2 \rightarrow \infty} B_{10}(\tau^2, \bar{y}) = 0$ and therefore $\lim_{\tau^2 \rightarrow \infty} p(H_0 | y_1, \dots, y_n) = 1$ for any values of n and \bar{y} .

Barlett's paradox

- This appears to be paradoxical because we would expect the Bayes factor to prefer H_1 when \bar{y} is large (and indeed $\lim_{\bar{y} \rightarrow \infty} B_{10}(\tau^2, \bar{y}) = \infty$). Also, we had shown that Bayesian model selection is consistent as $n \rightarrow \infty$.
- In conclusion, no matter how strong the information we have in our sample, the Bayes factor constructed using a flat prior will always favor the null model.
 - The same phenomenon appears in most comparisons that involve testing point vs. composite hypotheses using improper priors on the parameter being tested.
- More generally, this result highlights that priors can have a big impact for model selection, even for large sample sizes.

Barlett's paradox

- The key to understand Barlett's paradox is to notice that the order in which you take the limits matter! If you first take $\tau^2 \rightarrow \infty$ first, the Bayes factor is ill behaved.
- Another way to think about Barlett's paradox is that it arises because improper priors are defined up to an arbitrary normalizing constant, the Bayes factor depends on such constant and their value is therefore undefined.
 - This does not happen in estimation problems because the arbitrary constant appears both in the numerator and the denominator of the posterior distribution, canceling out.
- Another way to think about it is that the improper prior essentially puts probability 0 on H_1 a priori.
- Barlett's paradox does not arise when testing composite vs. composite hypotheses where the same improper prior is used under both hypotheses

Barlett's paradox

- **Example:** Again, let y_1, \dots, y_n be such that $y_i | \theta \sim N(\theta, 1)$ and consider testing the hypotheses $H_0 : \theta \leq 0$ against $H_1 : \theta > 0$ under the prior $\theta \sim N(0, \tau^2)$. The posterior odds in this case reduce to

$$\frac{p(H_0 | y_1, \dots, y_n)}{p(H_1 | y_1, \dots, y_n)} = \frac{\int_{-\infty}^0 \exp \left\{ -\frac{1}{2} \left(n + \frac{1}{\tau^2} \right) \left(\theta - \frac{n\bar{y} + \frac{\mu}{\tau^2}}{n + \frac{1}{\tau^2}} \right)^2 \right\}}{\int_0^\infty \exp \left\{ -\frac{1}{2} \left(n + \frac{1}{\tau^2} \right) \left(\theta - \frac{n\bar{y} + \frac{\mu}{\tau^2}}{n + \frac{1}{\tau^2}} \right)^2 \right\}}.$$

so that

$$0 < \lim_{\tau^2 \rightarrow \infty} \frac{p(H_0 | y_1, \dots, y_n)}{p(H_1 | y_1, \dots, y_n)} < \infty$$

for any finite values of \bar{y} and n . **Taking a limit is not a problem here!**

Barlett's paradox

- **Example (cont):** Intuitively, the improper prior is not a problem here because the (undefined!) normalizing constant of the prior appears in both the numerator and the denominator, so it cancels out!

More generally, traditional improper priors used in estimation will work well for testing composite vs. composite hypotheses, but NOT for testing point vs. composite hypotheses. Different (“well calibrated”) default priors are required in that setting.

Beyond conjugate models

- When you do not have conjugacy the process is the same, but you often cannot make much progress.
- There are a few exceptions. A well known one is the use of double exponential priors for the location parameter of a Gaussian likelihood,

$$p(y_1, \dots, y_n | \theta) = \left(\frac{1}{2\pi} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (y_i - \theta)^2 \right\},$$
$$\pi(\theta) = \frac{1}{2\tau} \exp \left\{ -\frac{1}{\tau} |\theta - \mu| \right\}.$$

Beyond conjugate models

- The posterior is in this case:

$$\begin{aligned} p(\theta | y_1, \dots, y_n) &= \frac{a(n, \bar{y})}{a(n, \bar{y}) + b(n, \bar{y})} \text{TN} \left(\lambda_1(n, \bar{y}), \frac{1}{n}, -\infty, \mu \right) \mathbb{I}(\theta < \mu) \\ &\quad + \frac{b(n, \bar{y})}{a(n, \bar{y}) + b(n, \bar{y})} \text{TN} \left(\lambda_2(n, \bar{y}), \frac{1}{n}, \mu, \infty \right) \mathbb{I}(\theta \geq \mu) \end{aligned}$$

where

$$\lambda_1(n, \bar{y}) = \bar{y} + \frac{1}{n\tau} \quad \lambda_2(n, \bar{y}) = \bar{y} - \frac{1}{n\tau}$$

and

$$\begin{aligned} a(n, \bar{y}) &= \exp \left\{ -\frac{\mu}{\tau} + \frac{n}{2} \lambda_1^2(n, \bar{y}) \right\} \Phi \left(\{\mu - \lambda_1(n, \bar{y})\} \sqrt{n} \right) \\ b(n, \bar{y}) &= \exp \left\{ \frac{\mu}{\tau} + \frac{n}{2} \lambda_2^2(n, \bar{y}) \right\} \{1 - \Phi \left(\{\mu - \lambda_2(n, \bar{y})\} \sqrt{n} \right)\} \end{aligned}$$

Beyond conjugate models

- The posterior mean is

$$\frac{a(n, \bar{y})}{a(n, \bar{y}) + b(n, \bar{y})} \left[\lambda_1(n, \bar{y}) - \frac{1}{\sqrt{n}} \frac{\phi(\sqrt{n}\{\mu - \lambda_1(n, \bar{y})\})}{\Phi(\sqrt{n}\{\mu - \lambda_1(n, \bar{y})\})} \right] + \\ \frac{b(n, \bar{y})}{a(n, \bar{y}) + b(n, \bar{y})} \left[\lambda_2(n, \bar{y}) + \frac{1}{\sqrt{n}} \frac{\phi(\sqrt{n}\{\mu - \lambda_2(n, \bar{y})\})}{1 - \Phi(\sqrt{n}\{\mu - \lambda_2(n, \bar{y})\})} \right],$$

where ϕ and Φ are the density and cdf of the standard normal distribution. Note that this is highly non-linear!

- One key feature is that the prior has bounded influence! \Rightarrow
Posterior distribution is robust to misspecification of the prior.

$$\lim_{\mu \rightarrow \infty} = \bar{y} + \frac{1}{n\tau} \qquad \qquad \lim_{\mu \rightarrow -\infty} = \bar{y} - \frac{1}{n\tau}$$

(Note the difference with a Gaussian prior.)

Bayesian Assymptotics

- If θ is continuous and the observations are iid, then under **mild regularity conditions** the posterior distribution is approximately normal

$$\theta | y \sim N\left(\hat{\theta}(y), H^{-1}(\hat{\theta}(y))\right),$$

where $\hat{\theta}(y)$ represents the maximum likelihood estimator of θ and $H(\hat{\theta}(y))$ is the observed information matrix.

- Sometimes called the “Bayesian CLT”.
- Other variants are possible, all of them are based on **Laplace approximations** to the posterior

Bayesian Assymptotics - Laplace expansions

- Consider an iid univariate sequence y_1, \dots, y_n where $y_i | \theta \sim p(y | \theta)$ and do a Taylor expansion of the log-likelihood:

$$\log \{\ell(\theta)\} \approx \log \ell\left(\hat{\theta}\right) + \frac{1}{2} \left. \frac{\partial^2}{\partial \theta^2} \log \ell\left(\hat{\theta}\right) \right|_{\theta=\hat{\theta}} \left(\theta - \hat{\theta}\right)^2$$

where $\left. \frac{\partial^2}{\partial \theta^2} \log \ell\left(\hat{\theta}\right) \right|_{\theta=\hat{\theta}}$ is the **observed information matrix**.

- For the prior, you have two options:
 - You can either do a similar Taylor expansion of the prior.
 - Under regularity conditions, the likelihood concentrates as n grows, so the prior can be discarded.

Bayesian Assymptotics - Laplace expansions

- This result has a few important consequences:
 - For large sample sizes, Bayesian and frequentist point and interval estimation procedures yield essentially the same results! (This is not true for hypothesis testing!)
 - For large sample sizes, the form of the prior does not matter much! (Again, this is not true for hypothesis testing!)
 - Most “reasonable” Bayesian point estimators are consistent (in a classical, frequentist sense).
 - For large sample sizes, the approximation can be used to easily construct point and interval estimates.
- The quality of the approximation might depend on the parameterization.

The Metropolis-Hastings algorithm

- We move on now to discuss specific MCMC algorithms used in Bayesian computation. We start with the Metropolis-Hastings (MH) algorithm, which encompasses many others as special cases.
- The idea behind the MH algorithm is similar to that behind the rejection algorithm:
 - Pick a proposal distribution $q(\vartheta | \theta)$ (think of it as *almost* the transition kernel of your Markov chain) that will be used to generate potentially new states.
 - The chain either stays on the old value or moves to this new proposed values according to a certain probability, that is chosen to ensure that the chain is reversible!

The Metropolis-Hastings algorithm

Given $\theta^{(t)}$

- ① Generate $\vartheta^{(t+1)} \sim q(\vartheta | \theta^{(t)})$.
- ② Take

$$\theta^{(t+1)} = \begin{cases} \vartheta^{(t+1)} & \text{with probability } \rho(\vartheta^{(t+1)}, \theta^{(t)}) \\ \theta^{(t)} & \text{with probability } 1 - \rho(\vartheta^{(t+1)}, \theta^{(t)}) \end{cases},$$

where

$$\rho(\vartheta^{(t+1)}, \theta^{(t)}) = \min \left\{ 1, \frac{p(\vartheta^{(t+1)} | \mathbf{y})}{p(\theta^{(t)} | \mathbf{y})} \frac{q(\theta^{(t)} | \vartheta^{(t+1)})}{q(\vartheta^{(t+1)} | \theta^{(t)})} \right\}.$$

The Metropolis-Hastings algorithm

- $\rho(\vartheta^{(t+1)}, \theta^{(t)})$ is called the Metropolis-Hastings acceptance probability.
- Note that the algorithm depends on the ratio $\frac{p(\vartheta^{(t+1)} | \mathbf{y})}{p(\theta^{(t)} | \mathbf{y})}$, so it works even if $p(\theta | \mathbf{y})$ is known only up to a normalizing constant.
- There are a number of variants of the algorithm!

The Metropolis-Hastings algorithm

In the random-walk Metropolis-Hastings, given $\theta^{(t)}$:

- ① Generate $\vartheta^{(t+1)} \sim g(|\vartheta - \theta^{(t)}|)$, where g is a density.
- ② Take

$$\theta^{(t+1)} = \begin{cases} \vartheta^{(t+1)} & \text{with probability } \rho(\vartheta^{(t+1)}, \theta^{(t)}) \\ \theta^{(t)} & \text{with probability } 1 - \rho(\vartheta^{(t+1)}, \theta^{(t)}) \end{cases},$$

where

$$\rho(\vartheta^{(t+1)}, \theta^{(t)}) = \min \left\{ 1, \frac{p(\vartheta^{(t+1)} | \mathbf{y})}{p(\theta^{(t)} | \mathbf{y})} \right\}.$$

Common choices for g include zero-mean Gaussian (my favorite) or uniform distributions!

The Metropolis-Hastings algorithm

- **Example (Gumble likelihood):** Consider a setting where y_1, \dots, y_n corresponds to a random sample from a Gumble distribution with location parameter θ , i.e.,

$$p(y_i | \theta) = \exp\{-(y_i - \theta) - \exp\{-(y_i - \theta)\}\} \quad y_i \in \mathbb{R}$$

and assume that we let $\theta \sim N(\xi, \kappa^2)$ a priori.

The posterior distribution associated with this model is intractable. However, creating a RWMH algorithm to sample from it is straightforward. Because θ can in principle take any real value, a Gaussian random walk seems appropriate,

$$q(\vartheta | \theta) = \frac{1}{\sqrt{2\pi}\tau} \exp\left\{-\frac{1}{2\tau^2}(\vartheta - \theta)^2\right\}$$

The Metropolis-Hastings algorithm

- **Example (Gumble likelihood, cont):** The corresponding acceptance probability becomes

$$\rho(\vartheta, \theta) =$$

$$\min \left\{ 1, \frac{\exp \left\{ - \sum_{i=1}^n (y_i - \vartheta) - \sum_{i=1}^n \exp \{ -(y_i - \vartheta) \} \right\}}{\exp \left\{ - \sum_{i=1}^n (y_i - \theta) - \sum_{i=1}^n \exp \{ -(y_i - \theta) \} \right\}} \frac{\exp \left\{ - \frac{(\vartheta - \xi)^2}{2\kappa^2} \right\}}{\exp \left\{ - \frac{(\theta - \xi)^2}{2\kappa^2} \right\}} \right\}$$

(Note that, as we discussed before, the ratio of the proposals cancels out because of the symmetry!)

This algorithm is implemented in the file `gumble.R`. Compare the results using $\tau^2 = 0.001$ (too small!), $\tau^2 = 0.07$ (about right, roughly 40% acceptance rate), and $\tau^2 = 5$ (too large!).

The Metropolis-Hastings algorithm

- **Example (Gaussian process):** Consider now a model where a vector $\mathbf{y} = (y_1, \dots, y_n)$ follows a joint multivariate distribution with mean 0 and variance-covariance matrix Σ with

$$[\Sigma(\lambda, \kappa)]_{i,j} = \kappa \exp\left\{-\frac{|x_i - x_j|}{\lambda}\right\}, \quad \lambda, \kappa > 0.$$

In the previous expression, λ and κ are unknown *positive* parameters, and the values x_1, \dots, x_n associated with each of the entries of \mathbf{y} are assumed to be known.

This problem is more complicated than the previous one in two ways: Now $\theta = (\lambda, \kappa)$ is two-dimensional, and the parameters are restricted to be positive!

The Metropolis-Hastings algorithm

- **Example (Gaussian process, cont):** Because the parameters must be positive, we use Gamma priors for both:

$$p(\lambda) = \frac{b_1^{a_1} \lambda^{a_1-1}}{\Gamma(a_1)} \exp\{-b_1\lambda\} \quad p(\kappa) = \frac{b_2^{a_2} \kappa^{a_2-1}}{\Gamma(a_2)} \exp\{-b_2\kappa\}$$

For the same reason a Gaussian random walk directly on (λ, κ) is not a great idea (you would be automatically rejecting every time you generate negative values).

However, a (bivariate!) Gaussian random walk for $(\log \lambda, \log \kappa)$ is feasible and does not lead to automatic rejections.

(Another alternative is to use a reflecting Gaussian random walk, but I will not pursue that idea further.)

The Metropolis-Hastings algorithm

- **Example (Gaussian process, cont):** Since we work with a transformation of the parameters, there are two ways to derive the algorithm (they lead to different formal expressions, but they are equivalent!):
 - ① Do a transformation so that the problem is reformulated as sampling from $p(z_1, z_2 | \mathbf{y})$ where $z_1 = \log \lambda$ and $z_2 = \log \kappa$ instead of $p(\lambda, \kappa | \mathbf{y})$. Once samples $z_1^{(1)}, \dots, z_1^{(B)}$ and $z_2^{(1)}, \dots, z_2^{(B)}$ have been generated, samples $\lambda^{(1)}, \dots, \lambda^{(B)}$ and $\kappa^{(1)}, \dots, \kappa^{(B)}$ can be constructed by letting $\lambda^{(b)} = \exp\{z_1^{(b)}\}$ and $\kappa^{(b)} = \exp\{z_2^{(b)}\}$
 - ② Keep the original target $p(\lambda, \kappa | \mathbf{y})$ but think of your proposal as being a bivariate log-Gaussian distribution rather than a Gaussian (and recognize that your proposal is almost symmetric, but not quite!).

Note that, in both cases, a Jacobian is going to be involved.

The Metropolis-Hastings algorithm

- **Example (Gaussian process, cont):** Taking the second route, the proposal distribution is

$$p(\vartheta_1, \vartheta_2 | \lambda, \kappa) = \frac{1}{2\pi} \frac{1}{\vartheta_1 \vartheta_2} |\boldsymbol{\Omega}|^{-1/2} \exp \left\{ -\frac{1}{2} \begin{pmatrix} \log \vartheta_1 - \log \lambda \\ \log \vartheta_2 - \log \kappa \end{pmatrix}^T \boldsymbol{\Omega}^{-1} \begin{pmatrix} \log \vartheta_1 - \log \lambda \\ \log \vartheta_2 - \log \kappa \end{pmatrix} \right\}$$

Note that $\boldsymbol{\Omega}$ does not have to be diagonal, so the moves could (and, turns out, should!) be correlated.

The Metropolis-Hastings algorithm

- **Example (Gaussian process, cont):** The corresponding acceptance probability for the algorithm becomes

$$\rho(\vartheta_1, \vartheta_2, \lambda, \kappa) = \min \left\{ 1, \frac{\left| \Sigma(\vartheta_1, \vartheta_2) \right|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \mathbf{y}^T [\Sigma(\vartheta_1, \vartheta_2)]^{-1} \mathbf{y} \right\} \vartheta_1^{a_1} \exp \{b_1 \vartheta_1\} \vartheta_2^{a_2} \exp \{b_2 \vartheta_2\} \vartheta_1 \vartheta_2}{\left| \Sigma(\lambda, \kappa) \right|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \mathbf{y}^T [\Sigma(\lambda, \kappa)]^{-1} \mathbf{y} \right\} \lambda^{a_1} \exp \{b_1 \lambda\} \kappa^{a_2} \exp \{b_2 \kappa\} \lambda \kappa} \right\}$$

(Note that the symmetric part of the proposal cancels out but the Jacobians are not symmetric, and they remain ...)

This algorithm is implemented in the file GP-MH.R. To tune the variance, run the algorithm once with $\Omega = \text{diag}\{0.3, 0.3\}$, compute $\text{Var}\{(\log \lambda, \log \kappa)^T | \mathbf{y}\}$, and then rerun it with $\Omega = \frac{2.38^2}{d} \text{Var}\{(\log \lambda, \log \kappa)^T | \mathbf{y}\}$. (In this case $d = 2$.)

The Metropolis-Hastings algorithm

Suggested exercise: Show that both approaches to constructing the random-walk Metropolis for the Gaussian process example lead to the same algorithm!

The Metropolis-Hastings algorithm

- The magic numbers $\frac{2.38^2}{d} \text{Var}(\theta|\mathbf{y})$ for proposals and 44% acceptance rates for unidimensional problems and 23% for multivariate problems comes from Gelman et al. (1996) Roberts et al. (1997) and Roberts et al. (2001), which derived general theoretical results and performed experiments on a Gaussian.
- Since, with enough data, most posteriors are approximately Gaussian, these are reasonable rules of thumb!
- Note that, because we are making two independent runs, the fact that we change the variance of the chain is not (conceptually) an issue as it does not invalidate the Markov property.
- If a good approximation for $\text{Var}(\theta|\mathbf{y})$ is available (rarely the case) then the preliminary run can be skipped.

①

Multivariate normal model

Multivariate normal distribution.

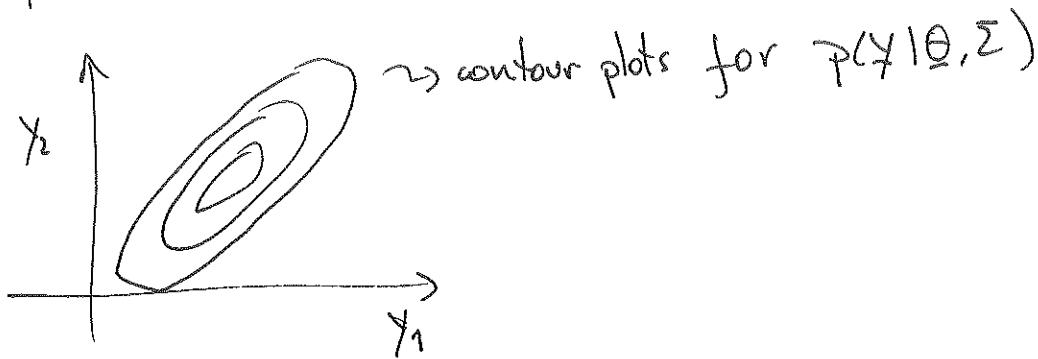
$$\mathbf{y} = (y_1, \dots, y_p)^T \quad \mathbf{y} \sim N_p(\underline{\theta}, \Sigma) \quad \underline{\theta} = (\theta_1, \dots, \theta_p)^T$$

if

$$p(\mathbf{y} | \underline{\theta}, \Sigma) = \left(\frac{1}{\sqrt{(2\pi)^p}} \right) |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \underline{\theta})^T \Sigma^{-1} (\mathbf{y} - \underline{\theta}) \right\} \quad \Sigma = \Sigma^T$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{pmatrix}$$

For $p=2$



How do we learn the parameter $\underline{\theta}$ and Σ from a sample

y_1, \dots, y_n .

First, assume Σ is known but $\underline{\theta}$ unknown.

$$p(\underline{\theta}) = \left(\frac{1}{\sqrt{(2\pi)^p}} \right) |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\underline{\theta} - \underline{d})^T \Sigma^{-1} (\underline{\theta} - \underline{d}) \right\} \text{ is a conjugate prior}$$

$$\underline{\theta} \sim N_p(\underline{d}, \Sigma)$$

$$P(\theta | x_1, \dots, x_n) \propto \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \right)^p |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \sum_i (x_i - \theta)^T \Sigma^{-1} (x_i - \theta) \right\} \quad (2)$$

$$\left(\frac{1}{\sqrt{2\pi}} \right)^p |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\theta - d)^T \Sigma^{-1} (\theta - d) \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} (\theta - d)^T \Sigma^{-1} (\theta - d) - \frac{1}{2} \sum_{i=1}^n (x_i - \theta)^T \Sigma^{-1} (x_i - \theta) \right\}$$

$$= \exp \left\{ -\frac{1}{2} \left[\theta^T \Sigma^{-1} \theta - \underbrace{\theta^T \Sigma^{-1} d}_{\theta^T \Sigma^{-1} \bar{x}} - \underbrace{d^T \Sigma^{-1} \theta}_{d^T \Sigma^{-1} \bar{x}} + d^T \Sigma^{-1} d + \sum_{i=1}^n x_i^T \Sigma^{-1} x_i - \theta^T \Sigma^{-1} \sum_{i=1}^n x_i - \left(\sum_{i=1}^n x_i^T \right) \Sigma^{-1} \theta + n \theta^T \Sigma^{-1} \theta \right] \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} \left[\theta^T \Sigma^{-1} \theta + n \theta^T \Sigma^{-1} \theta - 2 \theta^T \Sigma^{-1} d - 2 \theta^T \Sigma^{-1} n \bar{x} \right] \right\}$$

~~$$\propto \exp \left\{ -\frac{1}{2} \left[\theta^T (\Sigma^{-1} + n \Sigma^{-1}) \theta - 2 \theta^T [\Sigma^{-1} d + n \Sigma^{-1} \bar{x}] \right] \right\}$$~~

$$\propto \exp \left\{ -\frac{1}{2} \left[\theta - (\Sigma^{-1} + n \Sigma^{-1})^{-1} (\Sigma^{-1} d + n \Sigma^{-1} \bar{x}) \right]^T (\Sigma^{-1} + n \Sigma^{-1}) \right. \\ \left. \left[\theta - (\Sigma^{-1} + n \Sigma^{-1})^{-1} (\Sigma^{-1} d + n \Sigma^{-1} \bar{x}) \right] \right\}$$

$$\theta \stackrel{\text{Data}}{\sim} N \left((\Sigma^{-1} + n \Sigma^{-1})^{-1} (\Sigma^{-1} d + n \Sigma^{-1} \bar{x}), (\Sigma^{-1} + n \Sigma^{-1})^{-1} \right)$$

We say that $\Omega (= \Sigma^{-1})$ has a Wishart distribution ③

If:

$$P(\Omega) = \frac{1}{2^{vp/2} |A|^{-v/2} \Gamma_p(v/2)} |\Omega|^{(v-p-1)/2} \exp\{-\text{tr}(A\Omega)\}$$

where $\text{tr}(B)$ = trace of the matrix

$$= \sum_{i=1}^p b_{ii} = \text{sum of diagonal elements}$$

$$\Gamma_p(v/2) = \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma\left(\frac{j}{2} + \frac{1-j}{2}\right)$$

↓
"multivariate" gamma function

regular gamma function

$\Omega \in \mathbb{R}^{p \times p}$ and symmetric + positive definite.

Two parameter v, A

$$E(\Omega) = A^{-1} v$$

$$E(\Omega^{-1}) = E(\Sigma) = \frac{A}{v-p-1}$$

How to sample a Wishart (Ω, I)

(4)

$$\Omega = \begin{pmatrix} \omega_{11} & \omega_{12} \\ \omega_{21} & \omega_{22} \end{pmatrix}$$

~~$$\begin{matrix} \omega_{11} \sim \chi^2_1 \\ \omega_{12} \sim N(0, 1) \\ \omega_{21} \sim N(0, 1) \\ \omega_{22} \sim \chi^2_{n-1} \end{matrix}$$~~

check!!!

$$\Omega = BB^T$$

$$B = \begin{pmatrix} b_{11} & 0 & 0 & \cdots \\ b_{21} & b_{22} & 0 & \cdots \\ b_{31} & b_{32} & b_{33} & \cdots \end{pmatrix}$$

$$\begin{aligned} b_{11} &\sim \sqrt{\chi^2_D} & \omega_{11} &\sim \chi^2_D \\ b_{22} &\sim \sqrt{\chi^2_{D-1}} & \omega_{22} &\sim \chi^2_{D-1} \\ &\vdots && \vdots \\ b_{21} &\sim N(0, 1) & & \end{aligned}$$

The Wishart is a conjugate prior for the precision matrix $\Omega = \Sigma^{-1}$. Assume θ known.

$$P(y_1, \dots, y_n | \theta, \Omega) \propto |\Omega|^{n/2} \exp\left\{-\frac{1}{2}(\underline{y}_i - \underline{\theta})^T \Omega (\underline{y}_i - \underline{\theta})\right\}$$

$$P(\Omega | y_1, \dots, y_n) \propto |\Omega|^{(\frac{n-p-1}{2})/2} \exp\left\{-\frac{1}{2}\text{tr}(A\Omega)\right\} |\Omega|^{n/2}$$

$$\exp\left\{-\frac{1}{2} \sum_{i=1}^n \text{tr}((\underline{y}_i - \underline{\theta})^T \Omega (\underline{y}_i - \underline{\theta}))\right\}$$

If the dimensions allow it, then (5)

$$\text{trace}(ABC) = \text{trace}(BCA) = \text{trace}(CAB)$$

$$\Rightarrow P(\Omega | \text{data}) \propto |\Omega|^{(v+n-p-1)/2} \exp\left\{-\frac{1}{2}\left[\text{tr}(A\Omega) + \text{tr}((\hat{y}_i - \theta)(\hat{y}_i - \theta)^T \Omega)\right]\right\}$$

$$\Rightarrow \Omega | \text{data} \sim \text{Wishart}\left(v+n, A + \underbrace{\sum_{i=1}^n (\hat{y}_i - \theta)(\hat{y}_i - \theta)^T}_{p \times p \text{ matrix}}\right)$$

Inference for multivariate normal distributions.

①

$$\underline{Y}_1, \dots, \underline{Y}_n \stackrel{\text{iid}}{\sim} N_p(\underline{\theta}, \Sigma) \quad \Omega = \Sigma^{-1} \quad \text{both } \underline{\theta} \text{ and } \Sigma \text{ are unknown.}$$

Two options:

* Fully conjugate Normal-Wishart prior

$$\underline{\theta} | \Sigma \sim N_p(\underline{d}, \frac{1}{k} \Sigma)$$

$$\Omega = \Sigma^{-1} \sim \text{Wishart}(r, A)$$

* Conditionally conjugate prior

$$\underline{\theta} \sim N_p(\underline{d}, D)$$

$$\Omega = \Sigma^{-1} \sim \text{Wishart}(r, A)$$

We will focus on the second!!

In that case computation proceed using "Gibbs sampling"

Posterior:

$$P(\underline{\theta}, \Sigma | \underline{Y}_1, \dots, \underline{Y}_n) \propto \left(\frac{1}{\sqrt{2\pi}} \right)^{np} |\Sigma|^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\underline{y}_i - \underline{\theta})^\top \Sigma (\underline{y}_i - \underline{\theta}) \right\}$$

$$\left(\frac{1}{\sqrt{2\pi}} \right)^p |D|^{-\frac{r}{2}} \exp \left\{ -\frac{1}{2} (\underline{\theta} - \underline{d})^\top D^{-1} (\underline{\theta} - \underline{d}) \right\}$$

$$|\Sigma|^{-\frac{r-p-1}{2}} \exp \left\{ -\frac{1}{2} \text{trace}\{A \Sigma\} \right\}$$

(2)

$$P(\underline{\theta} | \Omega, \gamma) = P(\underline{\theta} | \dots)$$

$$\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\gamma_i - \underline{\theta})^\top \Omega (\gamma_i - \underline{\theta}) \right\}$$

$$\exp \left\{ -\frac{1}{2} (\underline{\theta} - \underline{\alpha})^\top D^{-1} (\underline{\theta} - \underline{\alpha}) \right\}$$

$$\underline{\theta} | \Omega, \gamma \sim N \left((n \Omega + D^{-1})^{-1} (n \Omega \gamma + D^{-1} \alpha), (n \Omega + D^{-1})^{-1} \right)$$

$$P(\Omega | \theta, \gamma) \propto |\Omega|^{n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\gamma_i - \theta)^\top \Omega (\gamma_i - \theta) \right\}$$

$$|\Omega|^{\frac{D-1}{2}} \exp \left\{ -\frac{1}{2} \text{tr}(A \Omega) \right\}$$

$$\Omega | \theta, \gamma \sim \text{Wish} \left(p+n, A + \sum_{i=1}^n (\gamma_i - \theta) (\gamma_i - \theta)^\top \right)$$

(3)

A slightly fancier version deals with some observations being missing.

$$\mathbf{y}_i = \begin{pmatrix} 2.03 \\ 1.82 \\ -1.07 \\ \vdots \\ -0.53 \end{pmatrix}$$

\Downarrow
fully observed

$$\mathbf{y}_i^* = \begin{pmatrix} 2.17 \\ \text{NA} \\ 0.17 \\ \vdots \\ -0.16 \end{pmatrix}$$

\Downarrow
partially observed

\mathbf{Y} = all observations.

~~Partially Obs.~~

$\mathbf{Y} = (\mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}})$ $\xrightarrow{\text{the value of the observations that I did not get}}$

The previous model can still be used for \mathbf{Y} .

$$P(\mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}} | \theta, \Omega) \sim \prod_{i=1}^n N(y_i | \theta, \Omega).$$

$$P(\mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}}, \theta, \Omega) = P(\mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}} | \theta, \Omega) P(\theta, \Omega)$$

We are interested in

$$P(\theta, \Omega | \mathbf{y}_{\text{obs}}) = \int P(\theta, \Omega, \mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}) d\mathbf{y}_{\text{mis}}.$$

(4)

We can sample from $P(Y_{\text{mis}}, \theta, \Sigma | Y_{\text{obs}})$
using Gibbs sampling.

$$P(\theta | \Sigma, Y_{\text{mis}}, Y_{\text{obs}})$$

$$P(\Sigma | \theta, Y_{\text{mis}}, Y_{\text{obs}}) \quad \text{Three full conditional distributions}$$

$$P(Y_{\text{mis}} | \theta, \Sigma, Y_{\text{obs}})$$

* The first two reduce to the case where there is no missing values.

* For the third, each missing value comes from a normal distribution that is the conditional distribution of the missing entry given the observed ones:

$$Y_i = \begin{pmatrix} Y_{\text{obs},i} \\ Y_{\text{mis},i} \end{pmatrix} \mid \begin{pmatrix} \theta_{\text{obs}} \\ \theta_{\text{mis}} \end{pmatrix}, \begin{pmatrix} \Sigma_{00} & \Sigma_{0m} \\ \Sigma_{m0} & \Sigma_{mm} \end{pmatrix} \sim N(Y_i | \theta, \Sigma)$$

$$Y_{\text{mis},i} \mid Y_{\text{obs},i}, \theta, \Sigma \sim N(\underline{\theta}_{\text{mis}} + \sum_{m0}^{-1} (\underline{Y}_{\text{obs},i} - \theta_{\text{obs}}), \Sigma_{mm} - \sum_{m0}^{-1} \sum_{0m})$$

Hierarchical modeling.

D

y_{ij} = SAT of the j -th student in the i -th school of a given school district.

Simplest model:

$$y_{ij} | \theta, \sigma^2 \sim N(\theta, \sigma^2) \quad i=1, \dots I \quad \text{and} \quad j=1, \dots J_i$$

$$\theta \sim N(\mu, \tau^2)$$

$$\sigma^2 \sim IG(a, b)$$

Does not account for the nested structure in the data

A better model:

$$y_{ij} | \theta_i, \sigma^2 \sim N(\theta_i, \sigma^2) \quad \sigma^2 \sim IG(a, b)$$

$$\theta_i \sim N(\mu, \tau^2)$$

If μ and τ^2 are fixed, then this is ^{almost} equivalent to fitting a separate model for each school.

Instead use priors for μ, τ^2

$$\Rightarrow \theta_i | N, \tau^2 \sim N(\mu, \tau^2)$$

$$\mu \sim N(m, K^2) \quad \tau^2 \sim IG(c, d)$$

(2)

- Repeated measurements model
- Random effects model
- Hierarchical model

How do we perform computation for this model?

$$P(\theta_1, \dots, \theta_I, \sigma^2, \mu, \tau^2 | \{y_{ij}\})$$

$$= \left[\prod_{i=1}^I \prod_{j=1}^{J_i} P(y_{ij} | \theta_i, \sigma^2) \right] \left[\prod_{i=1}^I P(\theta_i | \mu, \tau^2) \right] P(\mu) P(\tau^2) P(\sigma^2)$$

Full conditionals:

$$\theta_i | \dots \sim \sigma^2 | \dots$$

$$\mu | \dots \sim \tau^2 | \dots$$

$$\Rightarrow P(\theta_i | \dots) \propto \left[\prod_{j=1}^{J_i} P(y_{ij} | \theta_i, \sigma^2) \right] P(\theta_i | \mu, \tau^2)$$

$$\propto \left(\frac{1}{2\pi\sigma} \right)^{J_i} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^{J_i} (y_{ij} - \theta_i)^2 \right\} \sqrt{\frac{1}{2\pi} + \exp \left\{ -\frac{1}{2\tau^2} (\theta_i - \mu)^2 \right\}}$$

$$\Rightarrow \theta_i | \dots \sim N \left(\frac{\frac{J_i \bar{y}_{i:}}{\sigma^2} + \frac{\mu}{\tau^2}}{\frac{J_i}{\sigma^2} + \frac{1}{\tau^2}}, \frac{1}{\frac{J_i}{\sigma^2} + \frac{1}{\tau^2}} \right)$$

$$P(\mu | \dots) \propto \left[\prod_{i=1}^I P(\theta_i | \mu, \tau^2) \right] P(\mu) \quad (3)$$

$$\propto \left(\frac{1}{\sqrt{2\pi}\tau} \right)^I \exp \left\{ -\frac{1}{2\tau^2} \sum_{i=1}^I (\theta_i - \mu)^2 \right\} \frac{1}{\sqrt{2\pi/\kappa}} \exp \left\{ -\frac{1}{2\kappa^2} (\mu - \eta)^2 \right\}$$

$$\mu | \dots \sim N \left(\frac{\frac{I\bar{\theta}}{\tau^2} + \frac{\eta}{\kappa^2}}{\frac{I}{\tau^2} + \frac{1}{\kappa^2}}, \frac{1}{\frac{I}{\tau^2} + \frac{1}{\kappa^2}} \right)$$

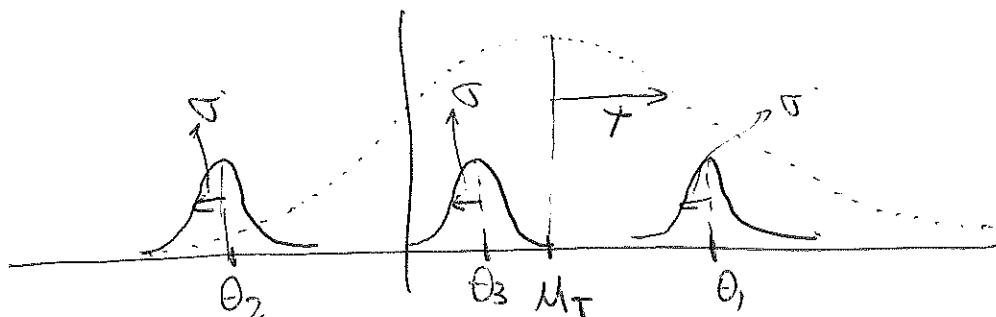
ANNEALING

$$P(\sigma^2 | \dots) \propto \prod_{i=1}^I \left(\frac{1}{\sigma^2} \right)^{J_i/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^{J_i} (y_{ij} - \theta_i)^2 \right\} \left(\frac{1}{\sigma^2} \right)^{a+1} \exp \left\{ -\frac{b}{\sigma^2} \right\}$$

$$\left(\frac{1}{\sigma^2} \right)^{a + \frac{\sum J_i}{2} + 1} \exp \left\{ -\frac{1}{\sigma^2} \left(b + \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^{J_i} (y_{ij} - \theta_i)^2 \right) \right\}$$

$$\sigma^2 \sim \text{IGam} \left(a + \frac{\sum J_i}{2}, b + \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^{J_i} (y_{ij} - \theta_i)^2 \right)$$

$$\tau^2 | \dots \sim \text{IGam} \left(c + \frac{I}{2}, d + \frac{1}{2} \sum_{i=1}^I (\theta_i - \mu)^2 \right)$$



Alternative model

(4)

$$Y_{ij} | \theta_i, \sigma_i^2 \sim N(\theta_i, \sigma_i^2)$$

$$\theta_i | \mu, \tau^2 \sim N(\mu, \tau^2) \quad \sigma_i^2 \sim \text{IGam}(a, b)$$

$$\mu \sim N(\eta, k^2)$$

$$b \sim \text{Gam}(\xi_1, \xi_2)$$

$$\tau^2 \sim \text{IGam}\left(\frac{\alpha}{c}, \frac{\beta}{d}\right)$$

NIMBLE

$$x_i | \lambda_i \stackrel{iid}{\sim} \text{Poisson}(\lambda_i)$$

$$\lambda_i = \theta_i t_i \quad \text{where } t_i \text{ is known.}$$

$$\theta_i | \alpha, \beta \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta)$$

$$\alpha \sim \text{Gamma}(1, 1)$$

$$\beta \sim \text{Gamma}(0.1, 1)$$

Full conditionals:

$$\theta_i | \text{data}, \alpha, \beta$$

$$\alpha | \text{data}, \{\theta_i\}$$

$$\beta | \text{data}, \cdot | \{\theta_i\}$$

$$p(\theta_i | \dots) \propto \frac{e^{-\theta_i t_i} (\theta_i t_i)^{x_i}}{x_i!} \left[\frac{\beta^\alpha}{\Gamma(\alpha)} \theta_i^{\alpha-1} \exp\{-\beta \theta_i\} \right]$$

$$\propto \exp\{-\theta_i(\beta + t_i)\} \theta_i^{x_i + \alpha - 1}$$

$$\theta_i | \dots \sim \text{Gamma}(\alpha + x_i, \beta + t_i)$$

(2)

$$P(\beta | \dots) \propto \left[\prod_{i=1}^n \frac{\beta^\alpha}{\Sigma(\alpha)} \theta_i^{\alpha-1} \exp\{-\beta \theta_i\} \right] \cancel{\exp\{-\beta\}} \\ \beta^{\alpha-1} \exp\{-\beta\}$$

$$\propto \beta^{n\alpha} \exp\{-\beta \sum_{i=1}^n \theta_i\} \beta^{-0.9} \exp\{-\beta\}$$

$$= \beta^{n\alpha - 0.9} \exp\{-\beta \left[1 + \sum_{i=1}^n \theta_i \right]\}$$

$$\beta | \dots \sim \text{Gamma} \left(n \alpha \cancel{+ 0.1}, 1 + \sum_{i=1}^n \theta_i \right)$$

$$P(\alpha | \dots) \propto \left[\prod_{i=1}^n \frac{\beta^\alpha}{\Sigma(\alpha)} \theta_i^{\alpha-1} \exp\{-\beta \theta_i\} \right] \exp\{-\alpha\}$$

$$\propto \frac{\beta^{n\alpha}}{\left[\Sigma(\alpha)\right]^n} \left[\prod_{i=1}^n \theta_i \right]^{\alpha-1} \exp\{-\alpha\}$$

Needs a random walk.

Linear models:

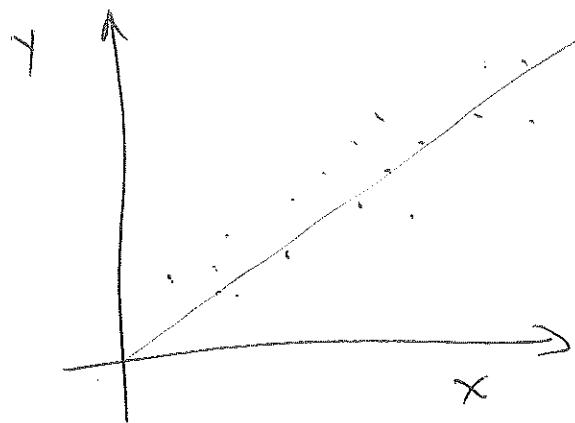
$$Y_i = \beta_0 + \sum_{k=1}^p x_{ik} \beta_k + \varepsilon_i$$

(1)

Known predictor
random error
observed values
Unknown parameters

Simple linear regression:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad (p=1)$$



Vector form:

$$Y_i = \underline{x}_i^T \underline{\beta} + \varepsilon_i \quad \underline{x}_i^T = (1, x_{i1}, x_{i2}, \dots, x_{ip})$$

$$\underline{\beta}^T = (\beta_0, \beta_1, \beta_2, \dots, \beta_p)$$

Matrix form:

$$\underline{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \underline{X} \underline{\beta} + \underline{\varepsilon}$$

$$\underline{X} = \begin{pmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{pmatrix} \quad \underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

First you need a likelihood ~~for~~ if you are going to do any kind of (Bayesian) inference ②

$$\varepsilon \sim N_n(0, \Sigma) \quad \text{simplest case } \Sigma = \sigma^2 I$$

but you could have more sophisticated scenarios

$$\Sigma = \sigma^2 \begin{pmatrix} 1 & p & 0 & \cdots & 0 \\ p & 1 & p & \cdots & 0 \\ 0 & p & 1 & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & 1 \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}^{-1} \rightarrow \text{AR}$$

Yet another way to write the model is

$$y \sim N_n(X\beta, \Sigma)$$

How do we do Bayesian inference for this model.

Consider first

$$y \sim N_n(X\beta, \sigma^2 D) \quad \text{with } D \text{ known and } \beta, \sigma^2 \text{ unknown.}$$

Priors:

$$\beta \sim N(\gamma, C)$$

$$\sigma^2 \sim \text{IGam}(a, b)$$

Computation with a Gibbs sampler:

(3)

$$P(\beta, \sigma^2 | \text{Data}) \propto \left(\frac{1}{\sqrt{2\pi}} \right) \left(\frac{1}{\sigma^2} \right)^{n/2} \left\{ D^{-1/2} \exp \left\{ -\frac{1}{2\sigma^2} (\gamma - X\beta)^T D^{-1} (\gamma - X\beta) \right\} \right.$$

$$\left. \left(\frac{1}{\sqrt{2\pi}} \right) \left(\frac{1}{\sigma^2} \right)^{n/2} \exp \left\{ -\frac{1}{2} (\beta - \eta)^T C^{-1} (\beta - \eta) \right\} \right\}$$

$$\left. \left(\frac{1}{\sigma^2} \right)^{a+1} \exp \left\{ -\frac{b}{\sigma^2} \right\} \right\}$$

$$P(\beta | \sigma^2, \text{data}) \propto \exp \left\{ -\frac{1}{2} \left[\beta^T X^T D^{-1} X + \frac{2}{\sigma^2} \beta^T X^T D^{-1} \gamma + \beta^T C^{-1} \beta - 2 \beta^T C^{-1} \eta \right] \right\}$$

$$= \exp \left\{ -\frac{1}{2} \left[\beta^T \left[\frac{1}{\sigma^2} X^T D^{-1} X + C^{-1} \right] \beta - 2 \beta^T \left(\frac{1}{\sigma^2} X^T D^{-1} \gamma + C^{-1} \eta \right) \right] \right\}$$

$$\beta | \sigma^2 \sim N \left(\left(\frac{1}{\sigma^2} X^T D^{-1} X + C^{-1} \right)^{-1} \left(\frac{1}{\sigma^2} X^T D^{-1} \gamma + C^{-1} \eta \right), \left(\frac{1}{\sigma^2} X^T D^{-1} X + C^{-1} \right)^{-1} \right)$$

$$P(\sigma^2 | \beta, \text{data}) \propto \left(\frac{1}{\sigma^2} \right)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} (\gamma - X\beta)^T D^{-1} (\gamma - X\beta) \right\}$$

$$\left. \left(\frac{1}{\sigma^2} \right)^{a+1} \exp \left\{ -\frac{b}{\sigma^2} \right\} \right\}$$

$$P(\sigma^2 | \beta, \text{data}) \propto \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2} + a + 1} \exp\left\{-\frac{1}{\sigma^2} \left(b + (\gamma - X\beta)^T D^{-1}(\gamma - X\beta)\right)\right\} \quad (4)$$

$$\sigma^2 \sim \text{IGam}\left(a + \frac{n}{2}, b + \frac{1}{2}(\gamma - X\beta)^T D^{-1}(\gamma - X\beta)\right)$$

If $D = I \Rightarrow$ independence of the observations
 and $X = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$ and $\beta = \beta_0$

Then

$$(\gamma - X\beta)^T D^{-1} (\gamma - X\beta) = \sum_{i=1}^n (\gamma_i - \beta_0)^2.$$

Gibbs sampler is relatively easy.

(5)

(More) Special cases of the linear model:

ANOVA model (one-way)

$$Y_{ij} = \theta_i + \varepsilon_{ij} \quad \varepsilon_{ij} \sim N(0, \sigma^2) \quad \begin{matrix} i=1, \dots, I \\ j=1, \dots, J_i \end{matrix}$$

$$\theta_i \sim N(\mu_i, \tau^2)$$

This can be rewritten in matrix form:

$$Y = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1,J_1} \\ Y_{2,1} \\ \vdots \\ Y_{2,J_2} \\ \vdots \\ Y_{I,1} \\ \vdots \\ Y_{I,J_I} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_I \end{pmatrix} + \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1,J_1} \\ \varepsilon_{2,1} \\ \vdots \\ \varepsilon_{2,J_2} \\ \vdots \\ \varepsilon_{I,1} \\ \vdots \\ \varepsilon_{I,J_I} \end{pmatrix}$$

X

① Fourier Regression:



$$y_i = f(x_i) + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$f(x) = \sum_{k=1}^B b_k(x) \theta_k$$

$$\Downarrow \cos\left(\frac{\omega k}{x}\right)$$

$$\Rightarrow y_i = \sum_{k=1}^B b_k(x) \theta_k + \varepsilon_i$$

$$Y = \begin{pmatrix} b_1(x_1) & \cdots & b_B(x_1) \\ b_1(x_2) & \cdots & b_B(x_2) \\ \vdots & & \vdots \\ b_1(x_n) & \cdots & b_B(x_n) \end{pmatrix} \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_B \end{pmatrix} + \varepsilon$$