

## Introduction:

- Bayes theorem,
- Predictive distribution
- Sequential learning

## Inference:

- Inference about functions of parameters
- Point & interval estimates
- Summary per hypothesis
- Statistical decision theory; optimal decision, when mean/mode/median
- Likelihood principle

## Priors:

- Conjugate priors; exponential family
- Non-informative; uniform & Jeffreys
- Hierarchical priors

## Graphs:

- DAG, moralising,  $X_1 \perp\!\!\!\perp X_2 \mid S$
- Factorisation theorem, Markov blanket, full conditional distribution

## Hierarchical models:

- Marginal prior
- exchangeability

## MCMC:

- Gibbs sampling; sampling from full cond distributions
- Burn-in; determine  $M$ , Gelman-Rubin diagnostic
- Determining  $N$ , MCSE,  $SE(\bar{f}_{MN})$  (batching)

## 1: Introduction

## Bayes Theorem

$$P(\theta|y) = \frac{p(\theta) p(y|\theta)}{p(y)} \propto \underset{\substack{\uparrow \\ \text{prior}}}{p(\theta)} \underset{\substack{\uparrow \\ \text{likelihood}}}{p(y|\theta)}$$

$\Theta$  is parameters,  $y$  is data  
(can also use on RVs,  $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ )

## Interpretation of Probability

• Frequentist:  $P(A) = \lim_{n \rightarrow \infty} \frac{m}{n}$    
  $\swarrow$  occurrence of A in n repeats

- Concat, but what about unobservable quantities?

- Bayesian:  $P(A)$  = subjective belief that A will occur

# Bayesian Inference

- Classical inference:  $\leftarrow$  what does  $y$  tell us about unknown  $\theta$

•  $y$  random,  $\theta$  fixed but unknown

Bayesian inference:  $y$  and  $\theta$  random  $\leftarrow$  how does  $y$  change our belief about  $\theta$

 $\gamma$  and  $\Theta$  random

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

## Predictive Distributions

$$p(\tilde{Y} = \tilde{y} | Y = y) = \int p(\tilde{y} | \theta) p(\theta | y) d\theta$$

(predict new  $\tilde{y}$  given already seen  $y$ )

- Sequential learning:  $\leftarrow$  'today's posterior is tomorrow's prior'

$$p(\theta|y_1, y_2) \propto p(y_2|\theta) \times p(\theta|y_1)$$

allows us to 'update' priors as we observe data

## 2: Inference

### Bayesian inference for the Normal distribution

$$Y_i | \theta, \gamma \sim N(\theta, \gamma^{-1}), i=1, \dots, n \quad (\gamma = \text{precision}, \sigma^2 = \gamma^{-1})$$

- $\theta$  and/or  $\gamma$  unknown, with priors

$$\text{eg } p(\theta, \gamma) \propto \gamma^{-1}$$

$$\Rightarrow p(\theta, \gamma | y) \propto p(\theta, \gamma) \cdot p(y | \theta, \gamma)$$

$$p(\gamma | y) = \int_{-\infty}^{\infty} p(\theta, \gamma | y) d\theta \quad (\text{similar for } \theta)$$

- Inference about functions of parameters:

$$\phi = g(\theta)$$

$$\Rightarrow p_{\phi|y}(\phi | y) = p_{\theta|y}(g^{-1}(\phi) | y) \cdot \left| \frac{d\theta}{d\phi} \right|$$

### Summarisation of posterior distributions

- 3 ways to summarise  $p(\theta | y)$ :

- Graphically:

- good for 1-2D; need to marginalise for higher dimensions (integrals often intractable; use MCMC)

- Quantitative:

- Point estimates: mean, mode, median

- Interval estimates:

- Credible interval:

$$P(\theta \in [c_1, c_2] | y) = 1 - \alpha \quad \left. \vphantom{P(\theta \in [c_1, c_2] | y)} \right\} \begin{array}{l} 100(1-\alpha)\% \\ \text{confidence} \end{array}$$

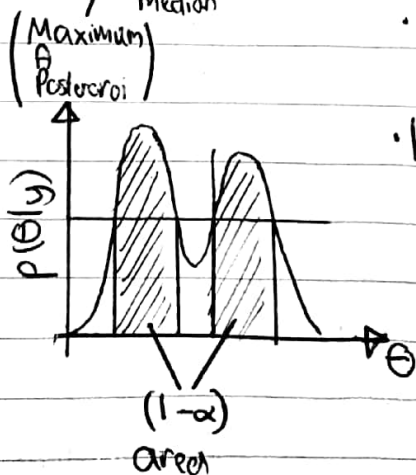
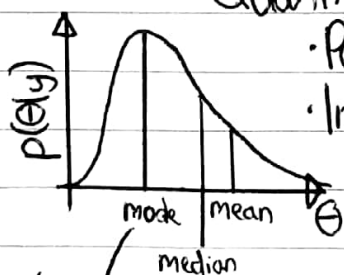
- Credible region:

$$P(\theta \in C | y) = 1 - \alpha$$

- Highest posterior density:

- interval s.t. any point inside has density  $>$  any point outside

- $100(1-\alpha)\%$  HPD interval is the shortest width  $100(1-\alpha)\%$  credible interval



• Summary for specific hypothesis:

eg  $H_0: \theta < \theta_0$  ( $H_1: \theta \geq \theta_0$ )

reject  $H_0$  if  $P(\theta < \theta_0 | y) < P(\theta \geq \theta_0 | y)$

(can scale by asymmetric losses too)

## Statistical Decision Theory

$$\left[ \underset{\substack{\uparrow \\ \text{(Type I loss)}}}{c_\alpha} P(\theta < \theta_0 | y) < \underset{\substack{\uparrow \\ \text{(Type II loss)}}}{c_\beta} P(\theta \geq \theta_0 | y) \right]$$

- parameter space  $\Theta$ , set of actions  $A$ , loss function  $L(\theta, a) \geq 0$
- decision function  $d(y)$  maps observed data  $y$  to action
- posterior expected loss:

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

- optimal decision minimises this

- Mean under square loss
- Median under absolute loss
- Mode under zero-one loss

## Comparison of Bayesian & classical inferences

• Point estimation:

expected  
= true  
↓

MSE  
↓

$n \rightarrow \infty$   
converge to true  
↓

not uniformly higher  
MSE than another  
estimator  
↓

• classical criteria: bias, efficiency, consistency, admissibility

• based around  $p(y|\theta)$ , as  $\theta$  fixed

• use asymptotic arguments; not used in Bayesian so much

• MAP = ML under uninformative prior / large samples

• Interval estimation:

• Bayesian:  $100(1-\alpha)\%$  chance  $\theta$  lies in interval

→ • Classical: Random interval containing  $\theta$  w/ probability  $100(1-\alpha)\%$

• In general can be quite different in practice Bayesian credible interval similar to classical confidence interval

• Likelihood principle:

•  $p(y|\theta)$  summarises all information about  $\theta$  provided by  $y$

• obeyed by Bayesian, violated by classical (as sampling distribution of  $Y$  can influence inference)

cannot be  
probability  
statement,  
as  $\theta$  not  
random

(ML estimator)  
obeys likelihood  
principle

## Pros / Cons of Bayesian Approaches

+ Theoretically sound, uses all available information

- where does prior come from (hard to justify), computational difficulties

### 3: Priors

#### Basic Considerations

$\Theta$ support	Suitable priors
$(-\infty, +\infty)$	Normal, t-dist
$(0, \infty)$	Gamma, log-normal
$(0, 1)$	Beta

functional form (support of  $\Theta$  must match) & parameter values

#### Conjugate Priors

Prior & Posterior in same class

class  $P$  of prior distributions  $p(\theta)$  is conjugate family for  $p(y|\theta)$  if  $p(\theta|y)$  is also in  $P$  for all  $y$

eg: likelihood | conjugate prior

Binomial

Beta

Normal (known)

Normal

exponential family

→

$f(y)g(\theta)e^{h(\theta)t(y)}$

$g(\theta)^v e^{h(\theta)\delta}$

$t(y)$  is sufficient statistic  
( $v$  = prior 'sample size'  
 $\delta$  = prior sufficient statistic)

Exponential family:  $p(y|\theta) = f(y)g(\theta)e^{h(\theta)t(y)}$

Distribution	$g(\theta)$	$h(\theta)$	$t(y)$
Normal ( $\theta \tau^{-1}$ )	$\exp[-\frac{\tau\theta^2}{2}]$	$\tau\theta$	$y_i$
Bin ( $m, \theta$ )	$(1-\theta)^m$	$\log \frac{\theta}{1-\theta}$	$y$

Choosing parameters:

choose  $v, \delta$  by roughly, or by knowing mean/variance of  $\theta$ , or confidence interval for  $\theta$  etc...

#### Non-informative Priors

Uniform: gives  $p(\theta|y) \propto p(y|\theta)$

if support is infinite,  $p(\theta)$  is improper;  $\int p(\theta)d\theta = \infty$   
⇒ may give improper posterior

Locally uniform:

'roughly' uniform where likelihood is non-negligible, dominated by likelihood

Beta (1,1) vs Beta (0,0):

Beta (1,1) uniform & proper, but influences posterior more than Beta (0,0) (which is improper)

eg large variance gaussian

→

maybe do Beta (0.001, 0.001) instead!

- Beta(0.001, 0.001) best of both, but if likelihood non-negligible at  $\theta=0$  or 1, then highly informative

• Jeffreys' prior:

- goal:  $p(\theta)$  remain invariant under transformation  $\phi = g(\theta)$

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

$$p(\theta) \propto \sqrt{I(\theta)} \quad (\Rightarrow p(\phi) = \sqrt{I(\phi)} \text{ for } \phi = g(\theta)) \quad \downarrow \text{proper}$$

Normal( $\theta, \gamma^{-1}$ ), known  $\theta \Rightarrow p(\gamma) \propto \frac{1}{\gamma}, \gamma \sim \text{Gamma}(0, 0)$   
 Bin( $n, \theta$ ),  $\theta \sim \text{Beta}(\frac{1}{2}, \frac{1}{2}) \leftarrow \text{improper}$

- often improper, violates likelihood principle, can lead to inconsistencies in multiparameter case ( $\theta_1$  can influence prior on  $\theta_2$  etc)

## Hierarchical priors

- divide model into stages; hyperpriors on priors
- eg  $Y \sim \text{Bin}(10, \theta)$   
 $\theta \sim \text{Beta}(\alpha, \beta)$   
 $\alpha \sim \text{Gamma}(4, 4), \beta \sim \text{Gamma}(5, 6)$
- often useful when  $\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$  &  $\otimes \theta$ s are exchangeable

## Summary

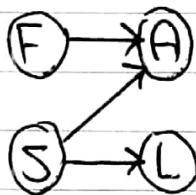
- Conjugate computationally convenient, but limiting
- Non-informative priors must be used carefully
- Sensitivity analysis of priors is important

## 4: Graphs

### Introduction

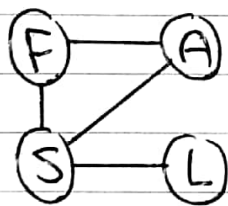
- Graphical models visual conditional independence structure
- $X \perp\!\!\!\perp Y$  if  $p(x,y) = p(x)p(y) \quad \forall x,y$
- $X \perp\!\!\!\perp Y|Z$  if  $p(x,y|z) = p(x|z)p(y|z)$

### Directed Acyclic Graphs



- Lines are indicated 'causation'; direct influence
- Node is conditionally independent of ancestors given parents
- Founders (no parents) are marginally independent
- Parents are conditionally dependent given shared child

### Moralising a DAG



ie  $X_3$  separates  $X_1$  &  $X_2$   
(this is global Markov property)

- 'marry' parents, drop directions
- to determine if  $X_1 \perp\!\!\!\perp X_2 | X_3$ ,
  - keep  $X_1, X_2, X_3$  & their ancestors
  - moralise this subgraph
  - true if  $X_3$  blocks all paths from  $X_1$  to  $X_2$  (else false)

### Factorisation theorem; Markov blankets & full conditional distributions

- Factorisation theorem:

$$p(X) = \prod_k p(X_k | \text{parents}[X_k])$$

(eg, above  
 $p(F,A,S,L) = p(F)p(S)p(A|F,S)p(L|S)$ )

- Markov blanket:

$$X_R \perp\!\!\!\perp X_{\setminus(X_R, \text{bl}[X_R])} | \text{bl}[X_R]$$

Useful when using MCMC methods to fit Bayesian models

where  $\text{bl}[X_R] = \text{parents}[X_R] \cup \text{children}[X_R] \cup \text{partners}[X_R]$   
(ie all nodes connected to  $X_R$  in moralised graph)

• Full-conditional distributions:

$$p(X_R | X_{\setminus X_R}) \propto p(X_R | \text{parents}[X_R]) \times \prod_{W \in \text{ch}[X_R]} p(W | \text{parents}[W])$$

• i.e. full-conditional dist of  $X_R$  depends only on variables in  $X_R$ 's Markov blanket

## Summary

- DAGs summarise statistical models as factorisation of simple relationships
- Moralising shows conditional independences



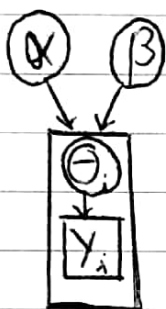
## 5: Hierarchical Models

### Non-hierarchical models



- must fully share ~~or have a copy of~~ parameter, or have all  $\Theta_i$  independent
- Want to encode some dependence between  $\Theta_i$ 's

### Hierarchical Models



$$p(\underline{\Theta}, \alpha, \beta) = \left[ \prod_i p(\Theta_i | \alpha, \beta) \right] p(\alpha) p(\beta)$$

- Can give weak priors on  $\alpha, \beta$  (eg  $\text{Exp}(0.01)$ )
- ensures  $\Theta_i$ 's weakly influence each other; ie assume  $\Theta_i$ 's similar but not identical
- shifts  $\Theta_i$ 's towards overall mean, and improves reliability when  $n_i$  is small



- in general, marginal prior

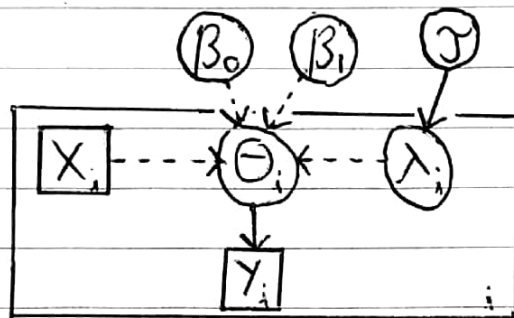
$$p(\underline{\Theta}) = \int p(\underline{\Theta} | \phi_2) p(\phi_2) d\phi_2$$

#### • Exchangeability:

- $\Theta_1, \dots, \Theta_n$  exchangeable if any permutation of these RVs has the same joint distribution
- Marginal independence  $\Rightarrow$  exchangeability (but not other way)  
& same marginal distribution
- Exchangeability implies a hierarchical model with some prior  $\phi$ , such that  $\Theta_i \perp \Theta_j | \phi$

# DAGs for hierarchical models

- circle nodes for unknowns
- Square nodes for observed RVs
- rectangular boxes for repetitive structure
- GLMM: Generalized Linear Mixed model



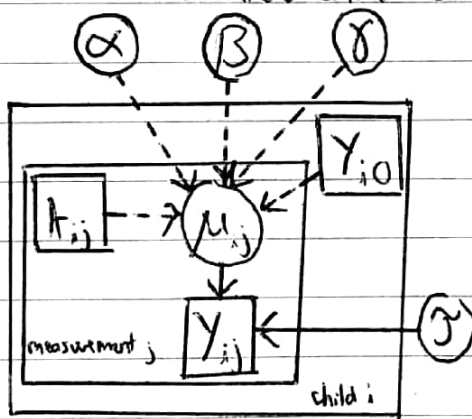
$$Y_i | \Theta_i \sim \text{Poisson}(C_i \Theta_i)$$

$$\log \Theta_i = \beta_0 + \beta_1 X_i + \lambda_i$$

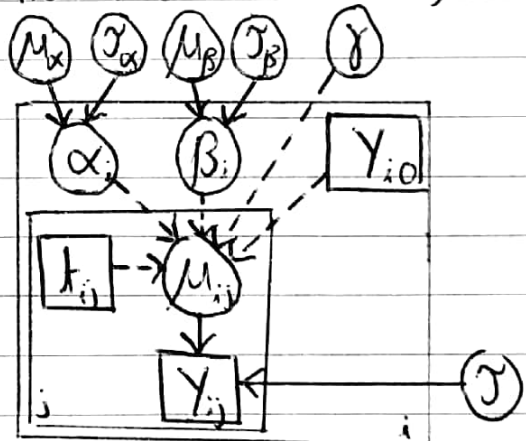
$$\lambda_i | \gamma \sim N(0, \gamma^{-1})$$

$$\beta_0, \beta_1, \gamma \sim \text{non-informative}$$

• Non-hierarchical LM:



Hierarchical LM (LMM):



$$Y_{ij} | \mu_{ij}, \gamma \sim N(\mu_{ij}, \gamma^{-1})$$

$$\mu_{ij} = \alpha + \beta(t_{ij} - \bar{t}) + \delta(Y_{i0} - \bar{Y}_0)$$

$$\alpha \sim N(0, 10000)$$

$$\beta \sim N(0, 10000)$$

$$\delta \sim N(0, 10000)$$

$$\gamma \sim \text{Gamma}(0.001, 0.001)$$

$$Y_{ij} | \mu_{ij}, \gamma \sim N(\mu_{ij}, \gamma^{-1})$$

$$\mu_{ij} = \alpha_i + \beta_i(t_{ij} - \bar{t}) + \delta(Y_{i0} - \bar{Y}_0)$$

$$\alpha_i | \mu_\alpha, \gamma_\alpha \sim N(\mu_\alpha, \gamma_\alpha^{-1})$$

$$\beta_i | \mu_\beta, \gamma_\beta \sim N(\mu_\beta, \gamma_\beta^{-1})$$

$$\mu_\mu, \mu_\alpha \sim N(0, 10000)$$

$$\gamma_\alpha, \gamma_\beta \sim \text{Gamma}(0.001, 0.001)$$

## Summary

- Breaking model into layers; parameters & hyper-parameters
- useful when data from similar but not identical 'units' which are exchangeable

hard to specify informative hyperpriors  
↓

# 6: Markov chain Monte Carlo

## Motivation

• Monte Carlo integration:

$$E[f(\theta) | x] = \int f(\theta) p(\theta | x) d\theta \approx \frac{1}{N} \sum_{i=1}^N f(\theta^{(i)}), \quad \theta \sim p$$

•  $\theta^{(i)}$  sampled from  $p(\theta | x)$

• Markov chain:

• used to sample from  $p(\theta | x)$ ; sample from an MC w/  $p(\theta | x)$  as its equilibrium distribution

## MCMC

• Metropolis-Hastings algorithm:

• general framework for MCMC; focus on special case; Gibbs sampling:

• Split  $\theta$  into  $K$  components

• Sample  $\theta_1^{(k+1)}$  from  $p(\theta_1 | \theta_{2:K}^{(k)})$  ← full conditional distribution

• repeat for  $\theta_2$  to  $\theta_K$  (always using most up to date  $\theta_j$ )

• now have  $\theta^{(k+1)}$ , repeat

• Sampling from full conditional distributions:

• may have closed form, if not use, eg, rejection sampling or ratio-of-uniforms method

## Convergence & Monte Carlo standard errors

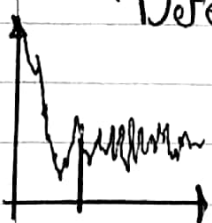
$$\hat{f}_{MN} - E[f(\theta) | x] \approx \frac{1}{N-M} \sum_{i=M+1}^N f(\theta^{(i)}) = \bar{f}_{MN}$$

• early iterations, 1 to  $M$ , are burn-in;  $\theta$  chain hasn't converged well yet, so don't use these

• Determining  $M$ :

• only truly converge at  $M = \infty$

• use trace plots; burnt in one sample look like random scatter about some value



• Convergence diagnostics:

• Gelman-Rubin diagnostic:

• run  $R$  chains, measure difference in behaviour ~~of~~ (when converge, should behave similarly)

$$R = \sqrt{\frac{V}{W}} \quad \begin{array}{l} V = \text{under-estimate of } \sigma_R^2 = \text{Var}(\theta_R | x) \\ W = \text{over-estimate of } \sigma_R^2 \\ (W, V \rightarrow \sigma_R^2 \text{ as } n \rightarrow \infty) \end{array}$$

•  $R < 1.05 \Rightarrow$  'practical' convergence  
(calculate  $R$  for all or several  $\theta$ )

• Determining  $N$ :

• run chain until Monte Carlo standard error (MCSE), less than 5% of parameters posterior standard deviation, for all parameters

• estimating  $SE(\bar{f}_{MN})$ :

• Batching:

$$\left( \begin{array}{l} b_q = \frac{1}{L} \sum_{i \in \text{batch}(q)} f(\theta^{(i)}) \\ \bar{b} = \frac{1}{Q} \sum_{q=1}^Q b_q \end{array} \right)$$

• divide sequence into  $Q$  uncorrelated batches of length  $L$

$$SE(\bar{f}_{MN}) = \sqrt{\frac{1}{Q(Q-1)} \sum_{q=1}^Q (b_q - \bar{b})^2}$$

## Pros & Cons of MCMC

- + freedom in modelling & inference
- + only available method for complex problems
- slow & computationally expensive
- hard to diagnose (lack of) convergence