

# Astrostatistics

modern statistical  
methods for astronomy

Lecture 2

27.1.25

+ data  $(x_i, y_i)$

① Where did it come from?

Invert question - if you knew "true" values  $(\eta_i, \xi_i)$ , how do you generate the data  
 $\leftarrow$  latent variables

Posit:  $x_i | \xi_i \sim N(\xi_i, \sigma_x^2)$  MEASUREMENT  
ERROR MODEL  
 $y_i | \eta_i \sim N(\eta_i, \sigma_y^2)$

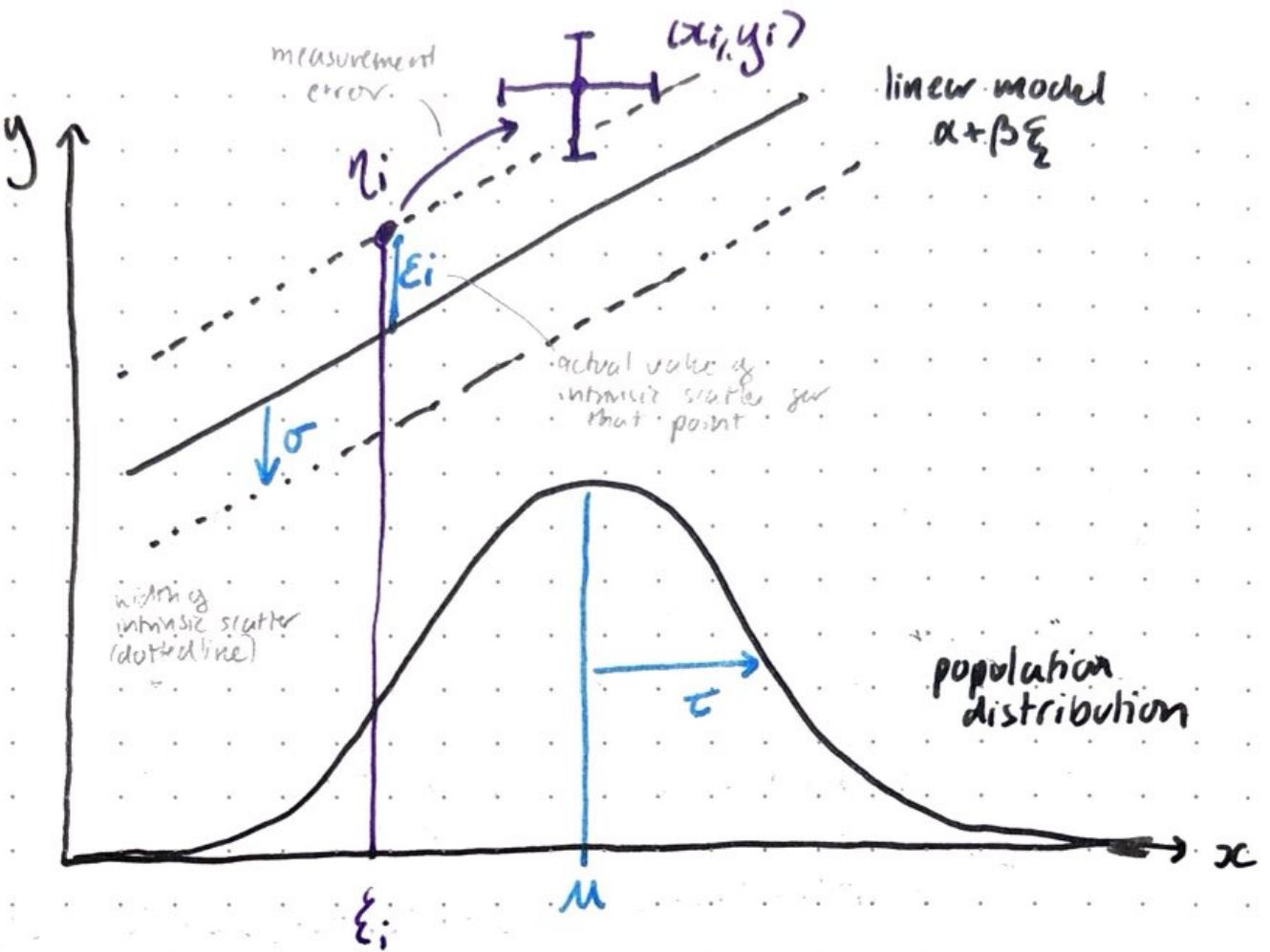
② Where do  $\eta_i, \xi_i$  come from?

③  $\eta_i$ ?

Posit:  $\eta_i = \alpha + \beta \xi_i + \epsilon_i$  intrinsic scatter  
 $\epsilon_i \sim N(0, \sigma^2)$   
linear model

④ Where does  $\xi_i$  come from?

Posit:  $\xi_i \sim N(\mu, \tau^2)$  POPULATION  
DISTRIBUTION  
MODEL  
"hyperparameters"



## GENERATIVE MODEL

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

Population Distribution

$$y_i | \xi_i \sim N(\alpha + \beta \xi_i, \sigma^2)$$

Regression Model

$$x_i | \xi_i \sim N(\xi_i, \sigma_x^2)$$

} Measurement Error

$$y_i | \eta_i \sim N(\eta_i, \sigma_y^2)$$

Knowns and Unknowns:

$$\Psi = (\mu, \tau)$$

Population Distribution Independent Variable "Hyperparameters"

$$\Theta = (\alpha, \beta, \sigma^2)$$

Regression Parameters

$$(\xi_i, \eta_i)$$

Latent (true) Variables  
((x,y) without measurement error)

$$(x_i, y_i)$$

Observed Data

$$(\sigma_{x,i}, \sigma_{y,i})$$

with measurement uncertainties

## Lecture 3

29.1.25

### FUNDAMENTALS AND NOTATION

(F&B ch. 2-3, Ivezic ch. 3-5)

measurement we take (some number)

possible outcome mass  
prob. distribution regarding our  
uncertainty in outcome

Data are realisations of random variables  
(outcomes of probabilistic experiments/ observations)

A random variable  $X$  can take on values in some domain (measurable space):

e.g. (discrete)  $\mathbb{Z}$ , (continuous)  $\mathbb{R}$ , (<sup>multivariate/</sup>vector continuous)  $\mathbb{R}^N$

#### Discrete case:

Define the probability distribution over possible outcomes with a probability mass function (PMF)

$$\Pr(X=k) = P_X(k)$$

prob (RV takes some value k)

$$\text{e.g. } \Pr(0 \leq X \leq 2) = \sum_{k=0}^2 P_X(k)$$

$$\text{Normalisation: } \sum_{k=0}^{\infty} P_X(k) = 1$$

#### Continuous case:

ie normally don't work w/ discrete e.g. 10C with photons  $\rightarrow$  continuous

Define the probability density function (PDF)

$$\Pr(x \leq X \leq x + dx) = P_X(x) dx$$

$$\text{e.g. } \Pr(0 \leq X \leq 2) = \int_0^2 P_X(x) dx$$

$$\text{Normalisation: } \int_{-\infty}^{\infty} P_X(x) dx = 1$$

and the cumulative distribution function (CDF)

$$\Pr(X \leq x_0) = \int_{-\infty}^{x_0} P_x(x) dx$$

Multivariate continuous case:

$\vec{x} \in \mathbb{R}^N$ ,  $S \subset \mathbb{R}^N$

$$\Pr(\vec{X} \in S) = \int_S P_{\vec{X}}(\vec{x}) d^N \vec{x}$$

volume element in  $\mathbb{R}^N$

(small  $\vec{x}$  is values R.V  $X$  takes on)

$$\text{Normalisation: } \int_{\mathbb{R}} P_{\vec{X}}(x) dx = 1$$

When there is no ambiguity, we may simplify notation,

$$\text{i.e. } P_x(x) \rightarrow P(x).$$

random variable

Given a distribution  $P(x)$  for RV  $x$ ,

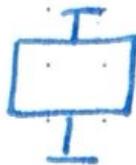
$X \sim P(x)$  means "distributed as", "is drawn from",  
"draw  $X$  from  $P(x)$ " depending on context  
e.g. in an algorithm means this  
not "approximation" sign

Astronomers' measurements are produced by physical processes that are inherently probabilistic (truly random processes thanks to quantum mechanics).

## Example 1: Photometry / Photon counting



x-ray



Chandrin  
x-ray  
observatory

$$\text{rate} = \frac{\text{photons arriving}}{\text{time}} = r$$

$$P_r(k \text{ photons in time } T) = \frac{e^{-\lambda} \lambda^k}{k!} \quad \text{where } \lambda = rT$$

$$P(k) = \text{Poisson}(k|\lambda)$$

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

Gaussian/Normal

$$\text{for } k \rightarrow \text{large}, P(k) \approx N(k|\lambda, \lambda) \quad \text{a limit theorem}$$

Define Gaussian Random Variable

$$X \sim N(\mu, \sigma^2), \quad X \in \mathbb{R}$$

$$\text{PDF: } P(x) = N(x|\mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2/\sigma^2}$$

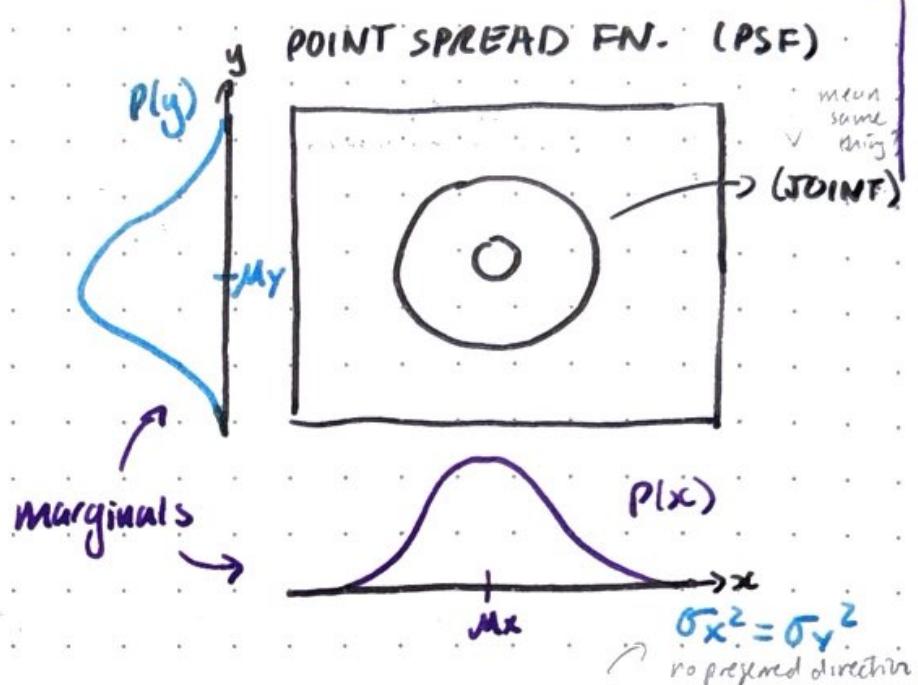
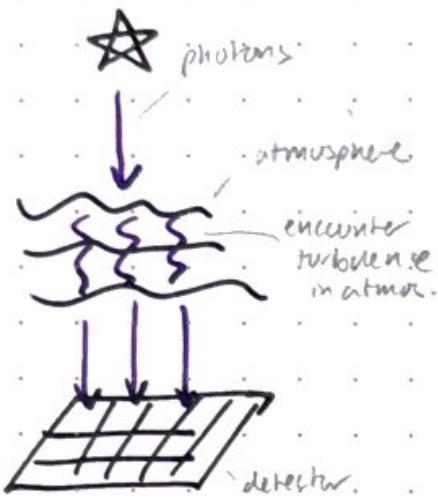
Generalise to Multivariate Gaussian Vector

$$\vec{v} \in \mathbb{R}^N \quad (\text{eg. } N=2, \vec{v} = \begin{pmatrix} x \\ y \end{pmatrix})$$

$$\text{PDF: } P(\vec{v}) = N(\vec{v}|\vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{\frac{N}{2}}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\vec{v}-\vec{\mu})^T \Sigma^{-1} (\vec{v}-\vec{\mu})}$$

$\vec{\mu} \in \mathbb{R}^N$ ,  $\Sigma$  is  $N \times N$  symmetric pos. def. matrix.

## Example: Astrometry



Joint distribution  $p(x,y) = N\left(\begin{pmatrix} x \\ y \end{pmatrix} \middle| \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{pmatrix}\right)$

simplifying (no preferred direction)

$$= N(x|\mu_x, \sigma_x^2) N(y|\mu_y, \sigma_y^2)$$

Law of total probability (LTP):

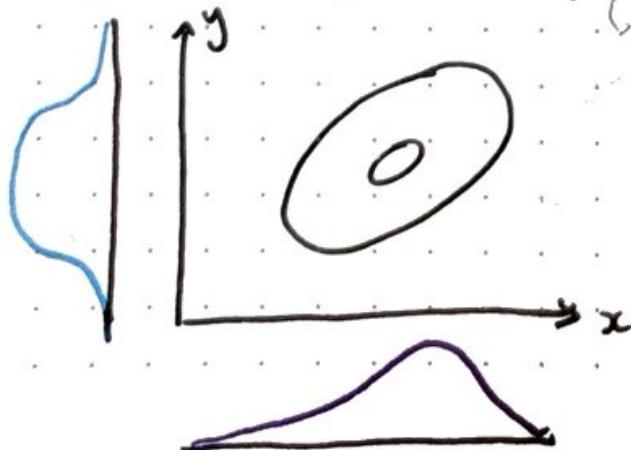
$$P(x) = \int_{-\infty}^{\infty} p(x,y) dy = N(x|\mu_x, \sigma_x^2)$$

$$P(y) = \int_{-\infty}^{\infty} p(x,y) dx = N(y|\mu_y, \sigma_y^2)$$

In this case,  $p(x,y) = P(x)P(y)$ ,  $x, y$  indep. RVs.

In general, this is not true:

$$p(x,y) = N\left(\begin{pmatrix} x \\ y \end{pmatrix} \middle| \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{pmatrix}\right) \quad \rho > 0$$



$$p(x,y) \neq P(x)P(y)$$

marginals same as above

## Conditional Probability

lower dim distributions  
obtained by integrating  
out other variables

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

Joint                      conditionals              marginal

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

BAYES' THEOREM

Note that LTP:  $P(x) = \int P(x,y) dy = \int P(x|y) p(y) dy$

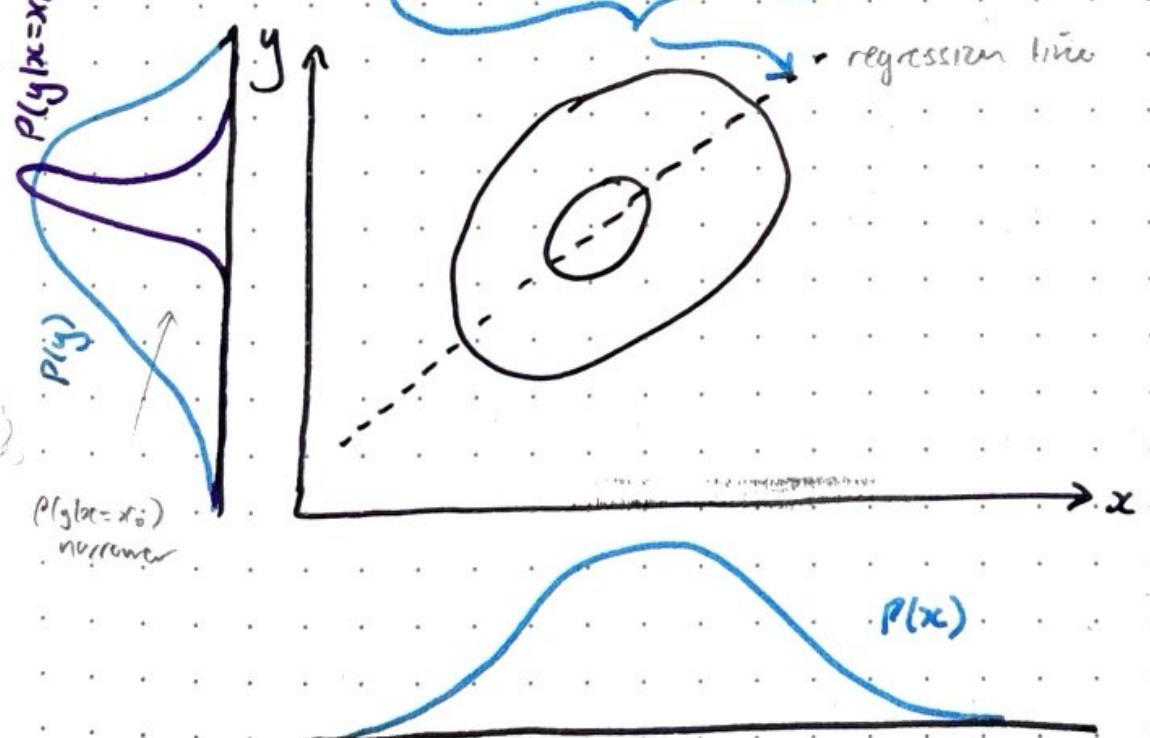
$$P(y|x) = \frac{P(x|y)P(y)}{\int P(x|y)P(y) dy}$$

In the bivariate case: Gaussian

$$P(y|x) = N(y | \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x), \sigma_y^2 (1 - \rho^2))$$

if i know  $x$ ; i can predict  $y$   
as it is bivariate gaussian  
b/c i know  $\sigma_y$ , variance in  $y$  decreases by  
factor  $(1-\rho^2)$

related to  
a property of gaussian  
(covars)



Observations  $x_0$  gives you information to predict  $y$   
 $P(y|x_0) \neq P(y)$ ,  $\rho > 0$ .

## Lecture 4

31.1.25

### MOMENTS & SUMMARIES OF PROB. DISTRIBUTIONS

$X \sim P(x) \rightarrow X$  is RV drawn from  $P(x)$

Expect values

$$\mathbb{E}[X] = \int x P(x) dx$$

$$\mathbb{E}[f(x)] = \int f(x) P(x) dx$$

Variance

$$\text{Var}(x) = \mathbb{E}[(x - \mathbb{E}(x))^2]$$

$$= \mathbb{E}[x^2] - (\mathbb{E}[x])^2$$

Multivariate  $(x, y)$ :

Covariance

$$\text{Cov}(x, y) = \mathbb{E}((x - \mathbb{E}[x])(y - \mathbb{E}[y]))$$

Multivariate Normal:

$$\begin{pmatrix} x \\ y \end{pmatrix} \sim N \left( \begin{pmatrix} M_x \\ M_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \sigma_x \sigma_y \rho \\ \sigma_x \sigma_y \rho & \sigma_y^2 \end{pmatrix} \right) \quad -1 < \rho < 1$$

\*  $X \perp Y \rightarrow \text{Cov}(X, Y) = 0$

\* general case  $\text{Cov}(X, Y) = \sigma_x \sigma_y \rho$

denotes independence  
of RVs

covariance  $\rho = \frac{\text{cov}}{\sigma_x \sigma_y}$   
normalization  
covariance

Useful Properties:

Collection of random variables

$$\{X_1, \dots, X_N\}, \{Y_1, \dots, Y_N\}$$

set of  $N \times N$  measurements  
labelled 1 to  $N$

$$\text{Def. } S = \sum_{i=1}^N a_i X_i, T = \sum_{j=1}^M b_j Y_j, \quad (a_i, b_j \in \mathbb{R})$$

Expectations are linear  $\rightarrow E[S] = \sum_{i=1}^N a_i E[X_i]$

expectation of sum is sum of expectation

Bilinearity of covariance: Cov of sum is sum of cov (bilinear in cov)

$$\text{Cov}(S, T) = \text{Cov}\left(\sum_{i=1}^N a_i X_i, \sum_{j=1}^M b_j Y_j\right) = \sum_i \sum_j a_i b_j \text{Cov}(X_i, Y_j)$$

$$\boxed{\text{Var}(S) = \text{Cov}(S, S) = \sum_{i=1}^N a_i^2 \text{Var}(X_i) + \sum_{i \neq j} a_i a_j \text{Cov}(X_i, X_j)}$$

$\therefore$  useful for summing multiple sources of error.

(variance of a sum of RV is covariance of sum w/ itself)

just this term  
is no errors correlated

some errors correlated

## TRANSFORMATIONS OF RVs

**IMPORTANT!**  
Variance of sum is sum of covariances:

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \sum_{j=1}^n \text{Cov}(X_i, X_j)$$

$$(\sum_{i,j} \text{Cov}(X_i, X_j) + \sum_i \text{Var}(X_i))$$

$X \sim P(x)$  invertible & differentiable transformation ( $y = \Phi(x)$ )

$Y = \Phi(X)$  what is  $P(y)$ ?

$$x = \Phi^{-1}(y)$$

$$P(y) = P(x) \left| \frac{dx}{dy} \right| = P(\Phi^{-1}(y)) \left| \frac{d\Phi^{-1}(y)}{dy} \right| \quad \text{jacobian}$$

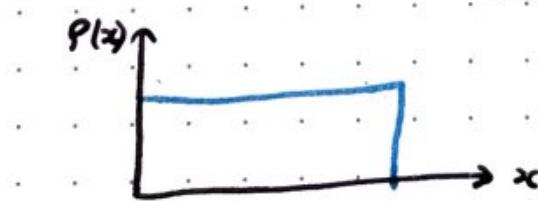
### Example:

$$X \sim N(0, 1)$$

$$Y = -\ln X$$

$$X = e^{-Y} \rightarrow P(y) = \begin{cases} e^{-y}, & 0 \leq y < \infty \\ 0, & \text{o/w} \end{cases}$$

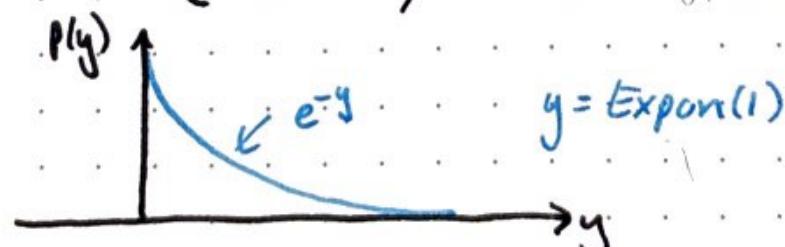
$$\frac{dx}{dy} = -e^{-y}$$



e.g. now that  $y \sim N(0, 1)$ .

then  $Y = Ax + b \sim N(b + \mu_A, A^T A)$ .

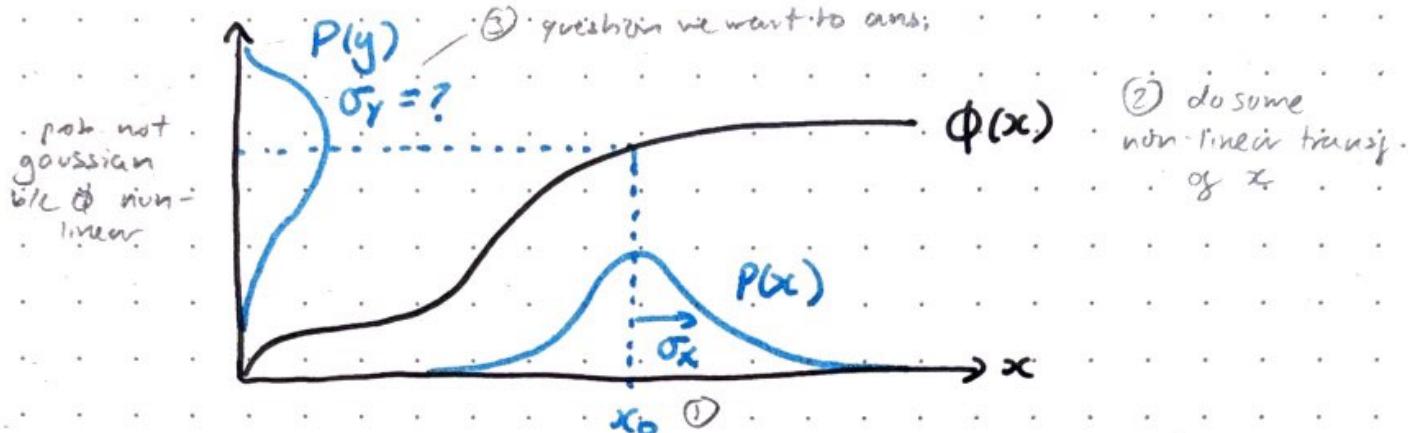
(clue:  $\left| \frac{dy}{dx} \right| = |A| = \det A$   
& properties of determinants in appendix)



can think of this as transforming RV  
but with variances

## Propagation of error

Measurement  $x_0$  w/ error variance  $\sigma_x^2$ , error bar  $\sigma_x$ .



$$y = \phi(x) = \phi(x_0) + \left. \frac{d\phi}{dx} \right|_{x=x_0} (x - x_0) + \dots$$

$$\text{Var}(y) \approx \left| \left. \frac{d\phi}{dx} \right|_{x=x_0} \right|^2 E((x-x_0)^2)$$

$$\sigma_y^2 \approx \left| \left. \frac{d\phi}{dx} \right|_{x=x_0} \right|^2 \sigma_x^2$$

$$\text{Var}(a+x) = \text{Var}(x)$$

$$\text{Var}(bx) = b^2 \text{Var}(x)$$

$$\begin{aligned} \text{Var}(y) &= \left| \left. \frac{d\phi}{dx} \right|_{x_0} \right|^2 \text{Var}(x-x_0) \\ &= \left| \left. \frac{d\phi}{dx} \right|_{x_0} \right|^2 \text{Var}(x). \end{aligned}$$

## Example:

Astronomers' flux  $f, \sigma_f$

magnitude  $m = -2.5 \log_{10}(f) + \text{const.}$

$$\sigma_m \approx \left( \frac{2.5}{\ln 10} \right) \frac{\sigma_f}{f} \approx \frac{\sigma_f}{f}$$

$$\log_{10} x = \frac{\ln x}{\ln 10}$$

goodness of approximation. Taylor exp. depends on how linear transformation is in the neighbourhood where we have uncertainty.  
(will break down if uncertainty large enough or transf. non-linear enough.)

MULTIVARIATE:  $\rightarrow y = g(x_1, \dots, x_n) \approx g(x^0) + \sum_i \left. \frac{dy}{dx_i} \right|_{x_i=x_i^0} (x_i - x_i^0)$  linear taylor exp.

$$\rightarrow \text{use rule: } \text{Var}(ax + by) = a^2 \text{Var}(x) + b^2 \text{Var}(y) + 2ab \text{Cov}(x, y)$$

$$\rightarrow \underbrace{\sigma_y^2}_{\sigma_y^2} \approx \sum_i \left( \frac{dy}{dx_i} \right)^2 \text{Var}(x_i) + \sum_{i \neq j} \left( \frac{dy}{dx_i} \right) \left( \frac{dy}{dx_j} \right) \text{Cov}(x_i, x_j) \rightarrow (\rho \sigma_x \sigma_y)$$

e.g. see wiki

intensity correlation

## Multivariate case:

just learn  
this?

$\rightarrow$  known Var/Cov

$$X_{\text{obs}} = (X_1, \dots, X_N), \quad Y = g(X_1, \dots, X_N)$$

$$\sigma_y^2 \triangleq \sum_{i=1}^N \left( \frac{dx}{dx_i} \right)^2 \Big|_{X_{\text{obs}}} \text{Var}(X_i) + \sum_{i \neq j} \left( \frac{dx}{dx_i} \right) \left( \frac{dx}{dx_j} \right) \Big|_{X_{\text{obs}}} \text{Cov}(X_i, X_j)$$

$$\sigma_y^2 \leq \sum_i \left( \frac{dx}{dx_i} \right)^2 \sigma_{x_i}^2 + \sum_{i \neq j} \left( \frac{dx}{dx_i} \right) \left( \frac{dx}{dx_j} \right) \text{Cov}(x_i, x_j). \quad \begin{matrix} \text{important term} \\ \text{accounts for correlated errors} \end{matrix}$$

## iid SEQUENCES OF RVs

why? e.g. in experiment do multiple observations, often model them as iid sequence (but not always!)

iid = independent & identically distributed  
(e.g. repeated measurement)

Suppose we have a sequence of indexed RVs:

$$X_1, X_2, X_3, \dots, X_N \stackrel{iid}{\sim} P(X) \quad \begin{matrix} \text{this means } (X_i) \text{ are independent!} \\ (\Rightarrow \text{Cov}(X_i, X_j) = 0) \end{matrix}$$

identically distrib.

$$X_i \sim P(X), \quad i \neq j \quad X_i \perp X_j$$

$$\text{Joint distribution} \sim P(X_1, X_2, X_3, \dots, X_N) = \prod_{i=1}^N P(X_i) \quad \begin{matrix} \text{indep} \\ \text{of each other condition?} \end{matrix}$$

Not example:  $\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N \left( \begin{pmatrix} \mu \\ \mu \end{pmatrix}, \begin{pmatrix} \sigma^2 & \sigma^2 \rho \\ \sigma^2 \rho & \sigma^2 \end{pmatrix} \right) \quad \rho \neq 0$

not. indep b/c correlated

May  $P(X_1) = \int P(X_1, X_2) dX_2 = N(X_1 | \mu, \sigma^2)$ ,  $P(X_2) = N(X_2 | \mu, \sigma^2)$

$X_1, X_2$  identically distributed but not independent (or as last lecture  $P(X_2 | X_1) \neq P(X_2)$ )

$$\text{Corr}(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1)} \sqrt{\text{Var}(X_2)}} = \rho \neq 0.$$

show not independent if  
 $\times P(a, b) \neq P(a)P(b)$   
or showing  
 $\star \text{Corr} = 0$   
(equivalently  $\text{Cov} = 0$ )

## Limit theorems for iid RVs

(see book for technicalities)

### ① LAW OF LARGE NUMBERS (LLN)

$$X_1, \dots, X_N \stackrel{iid}{\sim} P(x)$$

$$\mu = E[X_i] = \int x P(x) dx \rightarrow \left\{ \begin{array}{l} \text{theoretical model} \\ \text{population mean} \end{array} \right\}$$

$$\text{Def: Sample mean } \bar{X}_N \equiv \frac{1}{N} \sum_{i=1}^N X_i$$

as  $N \rightarrow \text{large}$ ,  $\bar{X}_N \rightarrow \mu$

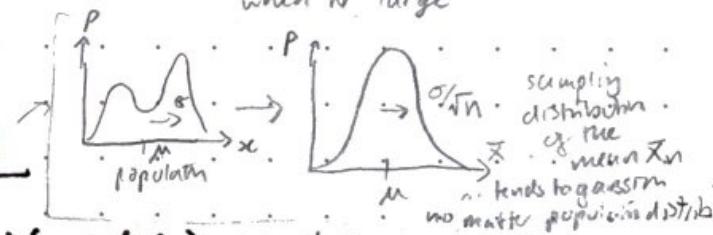
i.e. sample mean gives good approx. to population mean when  $N$  large

### ② CENTRAL LIMIT THEOREM

$$X_1, \dots, X_N \stackrel{iid}{\sim} P(x), \sigma^2 = \text{Var}(X_i) \text{ finite}$$

$$\text{As } N \rightarrow \text{large}, \frac{\bar{X}_N - \mu}{\sigma/\sqrt{N}} \xrightarrow{\text{approaches gaussian}} \text{unit gaussian}$$

$$\text{i.e. } \bar{X}_N \xrightarrow{} N(\mu, \sigma^2/N) \quad \text{error in the mean} = \sigma/\sqrt{N}$$



note: this is a general result, we haven't assumed  $P(x)$  is gaussian, it's just illustrative

## STATISTICAL ESTIMATORS

data, RV

Probability model  $P_\theta(x)$  or  $P(x|\theta)$

parameters

Suppose data are  $D = (X_1, \dots, X_N)$ ,  $X_1, \dots, X_N \stackrel{iid}{\sim} P_\theta(x)$

Estimate  $\theta$  from  $D$  by constructing an estimator  $\hat{\theta}(D)$  (fn. of  $D$ ) that yields a "good" value for the estimand  $\theta$ .

Thinking  
of estimating

data are drawn from prob. distribution

## Lecture 5

3.2.25

CMB  $\rightarrow$  radiation from ~300,000 years after big bang  
fluctuations in plasma temp of early universe  $\rightarrow$  map  
planck is a modern map of those fluctuations

distrib of temperatures  $(\text{at each pixel})$  measured is "the most perfect gaussian in nature"  
an example  $\rightarrow$  how might we find the mean of this distribution?

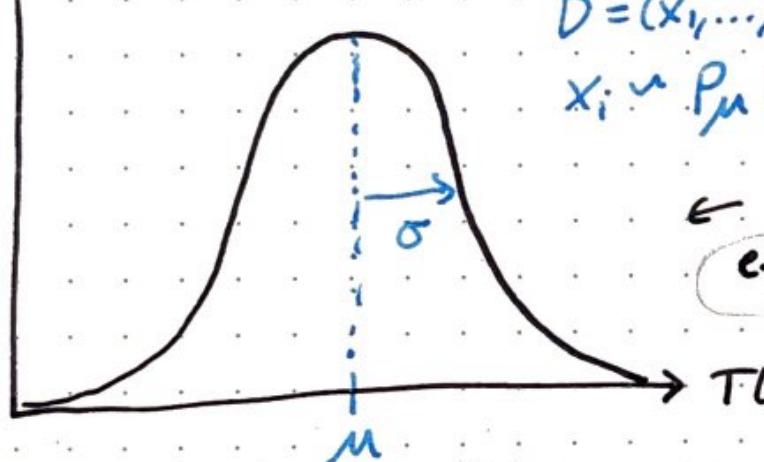
### STATISTICAL ESTIMATORS CONT.

Construct a "good" estimator  $\hat{\theta}(D) = \hat{\theta}(\vec{x})$  for estimand  $\theta$ .

Prob model  $P_\theta(x)$  or  $P(x|\theta)$

$D = (x_1, \dots, x_N) \stackrel{iid}{\sim} P(x|\theta)$

$x_i \sim P_\mu(x), \mu = E_p(x)$



$\leftarrow$  (CMB Example)

e.g. estimate  $\mu$  from temp measurement in each

pixel  $x_i \quad i=1, \dots, N$

Estimand  $\mu = E(x_i)$

Estimator  $\hat{\mu} = f(\vec{x})$

Can be many poss estimators for a particular estimand.

How do we choose?

① Sample mean

$$\frac{1}{N} \sum_{i=1}^N x_i$$

Sample mean of 1st  $K < N$  points

unbiased

truncated sample mean

② Take  $K < N$ ,  $\frac{1}{K} \sum_{i=1}^K x_i$

unbiased, inconsistent

only take K samples  
so doesn't have if  
large N

③  $\frac{1}{N-1} \sum_{i=1}^N x_i$

not unbiased  
but asymptotically unbiased

④ Midrange =  $\frac{1}{2} [\max(\vec{x}) + \min(\vec{x})]$

⑤ Median

⑥ Bin  $\rightarrow$  Histogram  $\rightarrow$  mode



7 Just report 3 Kelvin

How do you evaluate different possible estimators for a particular estimand?



$\hat{\mu}(\vec{x})$

### Criteria for Estimators $\hat{\theta}$ for $\theta$

$\hat{\theta} = h(\vec{x})$  is estimator for  $\theta$ ,

UNBIASEDNESS:  $\hat{\theta}$  is unbiased estimate for  $\theta$

if  $\boxed{E_p[\hat{\theta}] = \theta}$  — definition

e.g. ①, ②

bias:  $b(\hat{\theta}) = E[\hat{\theta}] - \theta$

over repeated experiments sampling the full data distribution  $P(\vec{x}| \theta)$

Imagine you did  $J$  experiments with fixed sample size  $N$ ,  $j=1, \dots, J$  — index experiments 1 to  $J$

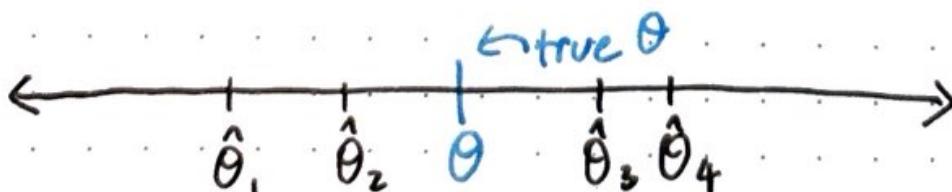
$\rightarrow \vec{x}_j = (x_{1,j}, \dots, x_{N,j})$  —  $j$  experiments

datasets:  $\vec{x}_j$

$$\vec{x}_1, \vec{x}_2, \dots, \vec{x}_J$$

$$\hat{\theta}_1 = \hat{\theta}(\vec{x}_1) \quad \hat{\theta}_2 = \hat{\theta}(\vec{x}_2) \quad \hat{\theta}_J = \hat{\theta}(\vec{x}_J)$$

Unbiased as  $J \rightarrow \infty$ ;  $\frac{1}{J} \sum_{j=1}^J \hat{\theta}_j = E[\hat{\theta}] = \theta$



as increase no. of experiments, avg tends to  $\theta$

N.B. You only really did one experiment!

ASYMPTOTICALLY UNBIASED:  $\boxed{\mathbb{E}[\hat{\theta}(x)] \rightarrow \theta \text{ as } N \rightarrow \infty}$

CONSISTENCY: As you gather more data (sample size  $N \rightarrow \text{large}$ ),  $\hat{\theta}$  converges to  $\theta$ ,

$$\boxed{\forall \varepsilon > 0, \Pr(|\hat{\theta} - \theta| > \varepsilon) \rightarrow 0 \text{ as } N \rightarrow \infty}$$

EFFICIENCY: Smallest mean squared error

$$\boxed{\text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta}(x) - \theta)^2]}$$

(for unbiased  $\rightarrow$  smallest variance)

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2] \\ &= \text{Var}(\hat{\theta}) + \text{Bias}^2(\hat{\theta}) \end{aligned}$$

BIAS-VARIANCE TRADE OFF

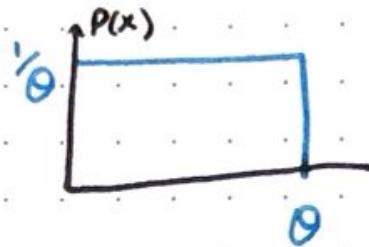
Sometimes the most efficient estimator is biased and the unbiased estimator is not the most efficient.

MINIMUM VARIANCE

UNBIASED ESTIMATORS (MVUE)

Caveat: Unbiased estimators can be wrong (obviously wrong)

Ex.  $X_i \stackrel{iid}{\sim} U(0, \theta)$   $i = 1, \dots, N$



Estimate  $\theta$ ?  
Unbiased Estimator?

Suppose we take sample mean:

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$

$$\mathbb{E}[\bar{X}] = \frac{\theta}{2}$$

$\hat{\theta}_u = 2\bar{X}$  unbiased estimator

since  $\mathbb{E}[\hat{\theta}_u] = 2(\frac{\theta}{2})$

Suppose instead always.

$$\hat{\theta}_m = \max(X_1, \dots, X_N) < \theta \quad \text{biased!}$$

generate sample  $\rightarrow$  RNG  $\bar{X} = (0.32, 0.46, 0.97, 2.77, 7.06)$

unbiased estimator  $\hat{\theta}_m = 4.63$  unbiased but clearly wrong

biased estimator  $\hat{\theta}_m = 7.06$  biased but much more sensible

Unbiasedness is not a property of any single experiment, only a property of averaging over all experiments that did not happen.

(cont'd) ... properties depend on underlying distribution

Ex. Estimator properties depend on data

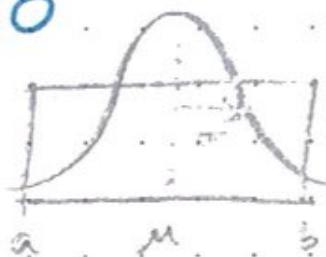
process  $P(X|\theta)$ ,  $N=10$   
(choose  $a, b$  s.t.  $M=(a+b)/2$ ,  $\text{Var}(X_i) = \frac{(b-a)^2}{12} = \sigma^2$ )

Estimator Variance	Gaussian $X_1, \dots, X_N \stackrel{iid}{\sim} N(\mu, \sigma^2=1)$	Uniform $X_1, \dots, X_N \stackrel{iid}{\sim} U(a, b)$
--------------------	---	---

(MSE)

Sample Mean  $\bar{X}$

$$\frac{\sigma^2}{N} = 0.10$$



$$\frac{\sigma^2}{N} = 0.10$$

Midrange  
 $\frac{\min(\bar{X}) + \max(\bar{X})}{2}$

$$\frac{\pi^2 \sigma^2}{24 \ln N} = 0.18$$

$$\frac{6\sigma^2}{(N+2)(N+1)} \approx \frac{6\sigma^2}{N^2} = 0.06$$

midrange less efficient

midrange more efficient

both sample mean & midrange are unbiased, but midrange  $\xrightarrow{\text{stdeviation}}$  goes down as  $1/\sqrt{N}$  whereas much slower for mean which goes down as  $1/\sqrt{N}$  (for uniform) see slides

$\rightarrow \theta^3 - \alpha^3 = (\theta - \alpha)(\theta^2 + \theta\alpha + \alpha^2)$   
 $\rightarrow$  much uniform & gaussian, because  $\theta^2$  is  $\propto \text{var}(\bar{X})$  with  $\text{var}(\bar{X}) \propto \sigma^2$

## LIKELIHOOD-BASED INFERENCE

Probability model  $P(\vec{D}|\theta)$  we call "sampling distribution"

Probability distribution for possible/potential datasets  $\vec{D}$ , for a given parameter value  $\theta$ .

We observe  $\vec{D}_{\text{obs}} : P(\vec{D} = \vec{D}_{\text{obs}} | \theta) = L(\theta)$

$\uparrow$   
Likelihood Function

Sampling dist. probability dist. over all possible outcomes of  $\vec{D}$  for a given  $\theta$ . ("before it is observed") N.B.  $\int P(\theta | \vec{D}) d\theta = 1$

## Lecture 6

another technique of point estimation:

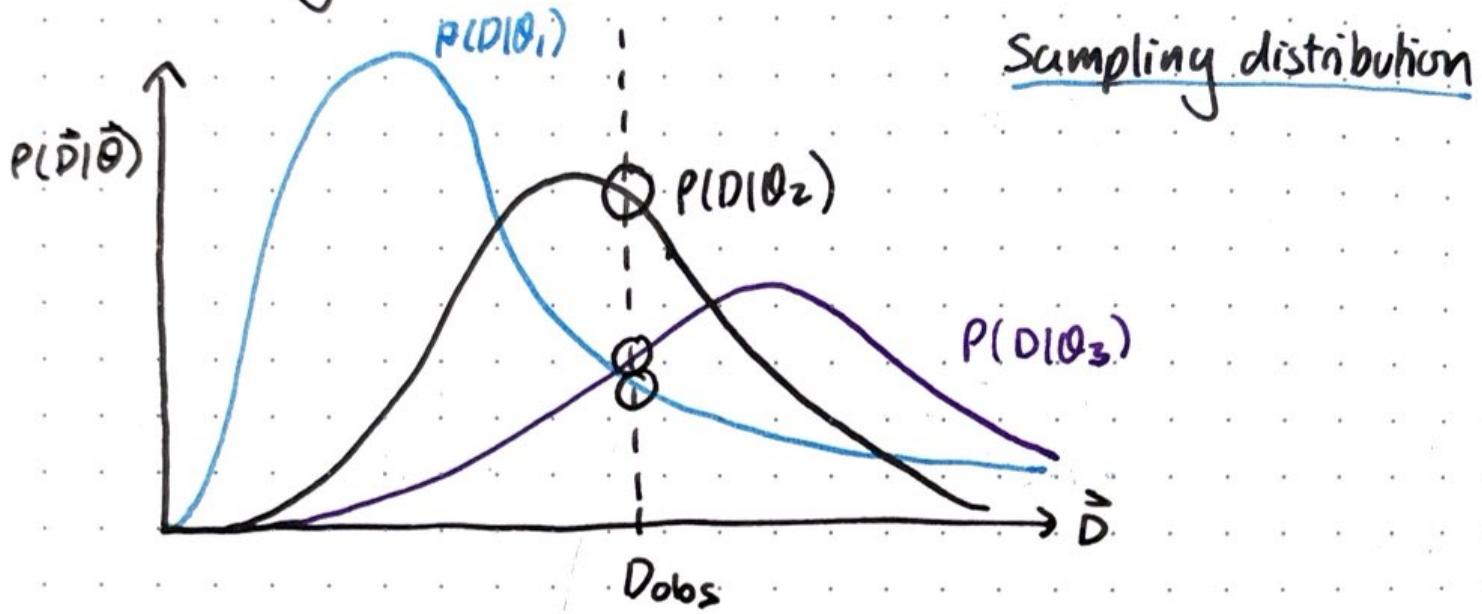
"change of given set of observations occurring"

$\hat{\theta}_{MLE}$

5.2.25

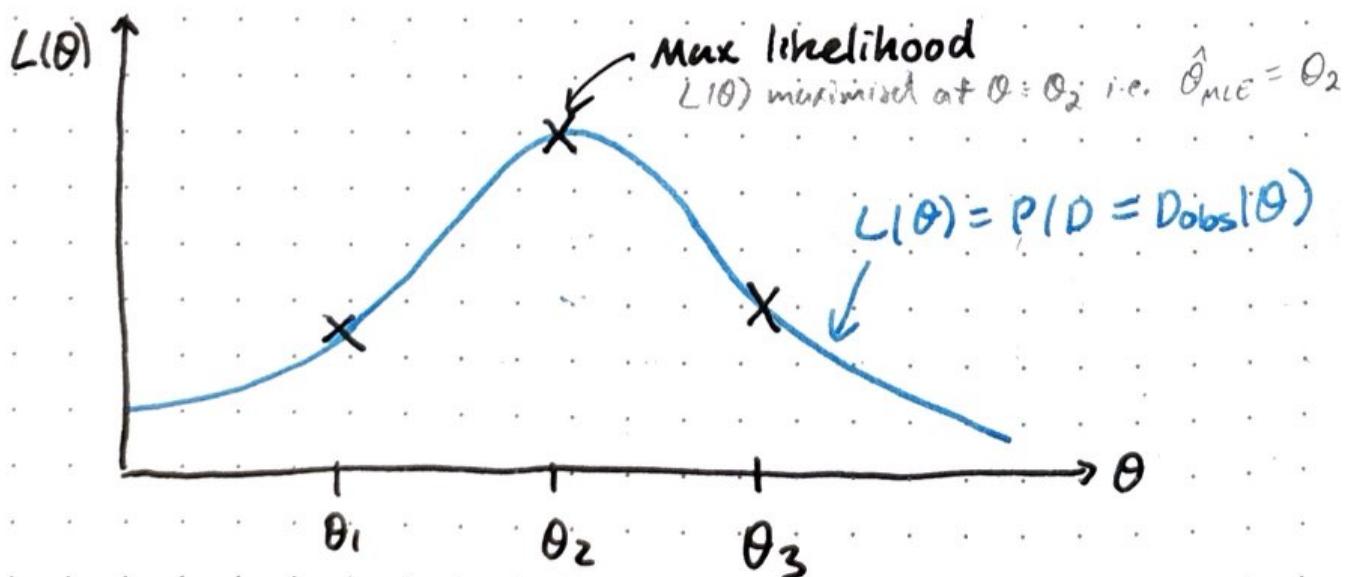
## LIKELIHOOD-BASED INFERENCE CONT.

Probability model  $D \sim P(\vec{D}|\vec{\theta})$ ,  $\int P(\vec{D}|\vec{\theta})d\vec{D} = 1$



Likelihood:  $P(D=D_{obs}|\theta) = L(\theta)$

N.B.  $\int L(\theta)d\theta$  not necessarily = 1



Suppose we now observe data  $D_{obs}$  (outcome or realisation of  $\vec{D}$ ), the likelihood is the sampling dist. evaluated at  $D=D_{obs}$ , viewed as a fn. of the parameters i.e.  $L(\vec{\theta}) = P(\vec{D}=D_{obs}|\vec{\theta})$

## Notation:

Often we will elide the distinction between  $\vec{D}$  and  $\vec{D}_{\text{obs}}$ .

$P(\vec{D}|\vec{\theta})$  understood to mean  $P(\vec{D} = \vec{D}_{\text{obs}}|\vec{\theta})$ .

log likelihood  $\ell(\theta) = \ln L(\theta)$ .

## iid case:

$$x_i \stackrel{\text{iid}}{\sim} f(x|\theta) \quad i=1, \dots, N, \quad P(\vec{x} = (x_1, \dots, x_N)|\theta) = \prod_{i=1}^N f(x_i|\theta)$$

*individual likelihood*

$$L(\theta) = P(\vec{x}|\theta) = \prod_{i=1}^N f(x_i|\theta)$$

*full likelihood*

$$\ell(\theta) = \sum_{i=1}^N \ln f(x_i|\theta)$$

*log likelihood*

## Fisher Information $\longleftrightarrow$ Uncertainty

$$L(\theta) = P(x|\theta)$$

DEFINITION:

$$\text{Score} = S = \frac{\partial}{\partial \theta} \ln L(\theta) = \frac{\partial}{\partial \theta} \ln P(x|\theta)$$

$\rightarrow$  w.r.t.  $P(x|\theta)$

$$\mathbb{E}[S] = \mathbb{E}\left[\frac{\partial}{\partial \theta} \ln P(x|\theta)\right] = \int \left[\frac{\partial}{\partial \theta} \ln P(x|\theta)\right] P(x|\theta) dx$$

$$= \int \left[\frac{1}{P(x|\theta)} \frac{\partial}{\partial \theta} P(x|\theta)\right] P(x|\theta) dx$$

expectation  
over sample  
space

$$* = \frac{\partial}{\partial \theta} \int P(x|\theta) dx = \frac{\partial}{\partial \theta} [1] = 0$$

\* Under regularity conditions  $\frac{\partial}{\partial \theta} \leftrightarrow \int dx$

Fisher Information = Variance of the Score

$$(I(\theta) = \text{Var}(S) = E(S^2) - (E(S))^2)$$

= Ongoing  
process  
results

$$I(\theta) = E\left[\left(\frac{\partial}{\partial \theta} \ln P(x|\theta)\right)^2\right] - \left(E\left[\frac{\partial}{\partial \theta} \ln P(x|\theta)\right]\right)^2$$

$$\stackrel{?}{=} -E\left[\frac{\partial^2}{\partial \theta^2} \ln P(x|\theta)\right]$$

$$E\left[\frac{\partial^2}{\partial \theta^2} \ln P(x|\theta)\right] = E\left[\frac{\partial}{\partial \theta}\left[\frac{1}{P(x|\theta)} \frac{\partial}{\partial \theta} P(x|\theta)\right]\right]$$

$$= E\left[-\frac{1}{P(x|\theta)} \cdot \frac{\partial}{\partial \theta} P(x|\theta) \frac{\partial}{\partial \theta} P(x|\theta) + \frac{1}{P(x|\theta)} \frac{\partial^2}{\partial \theta^2} P(x|\theta)\right]$$

$$= E\left[-\left(\frac{\partial}{\partial \theta} \ln P(x|\theta)\right)^2 + \underbrace{\frac{1}{P(x|\theta)} \frac{\partial^2}{\partial \theta^2} P(x|\theta)}_{\int \left[\frac{1}{P(x|\theta)} \frac{\partial^2}{\partial \theta^2} P(x|\theta)\right] P(x|\theta) dx = \frac{\partial^2}{\partial \theta^2} (1) = 0}\right]$$

$$\int \left[\frac{1}{P(x|\theta)} \frac{\partial^2}{\partial \theta^2} P(x|\theta)\right] P(x|\theta) dx = \frac{\partial^2}{\partial \theta^2} (1) = 0$$

$$= -E\left[\left(\frac{\partial}{\partial \theta} \ln P(x|\theta)\right)^2\right]$$

$$= -I(\theta),$$

$$I(\theta) = E\left[\frac{\partial^2}{\partial \theta^2}(-\ln P(x|\theta))\right]$$

From Wiki:

- amount of information R.V.  $x$  carries about an unknown parameter  $\theta$  (upon which the probability  $P(x|\theta)$  depends)

$$I(\theta) \geq 0$$

F.I.

Formally:

- variance of score

or - repeated values of observed information

↳ used to calculate covariance matrices associated with MLE

bound on variance of estimators

"inverse of Fisher info. is a lower bound on the variance of any unbiased estimator"

## Cramer - Rao Lower Bound

Suppose we have estimator  $T(x)$  for  $\theta$ ,

$$\mathbb{E}[T] = \theta + b(\theta).$$

Consider  $\text{Cov}(S(x), T(x)) = \mathbb{E}[S(x)T(x)] - \mathbb{E}[S]\mathbb{E}[T]$

$\stackrel{\text{score}}{=}$

$$\begin{aligned} &= \mathbb{E}\left[T(x) \frac{1}{P(x|\theta)} \frac{\partial}{\partial \theta} P(x|\theta)\right] \\ &= \int \left[T(x) \frac{1}{P(x|\theta)} \frac{\partial}{\partial \theta} P(x|\theta)\right] P(x|\theta) dx \\ &= \int \frac{\partial}{\partial \theta} T(x) P(x|\theta) dx \\ * &= \frac{\partial}{\partial \theta} \int T(x) P(x|\theta) dx \\ &= \frac{\partial}{\partial \theta} \mathbb{E}[T(x)] = 1 + b'(\theta) \end{aligned}$$

(under regularity conditions)

$-1 \leq \text{corr} \leq 1$   
 $| \text{corr} | \leq 1$

$$\text{Cov}(S, T) = \sqrt{\text{Var}(S)} \sqrt{\text{Var}(T)} \text{Corr}(S, T)$$

$$| \text{Cov}(S, T) | \leq \sqrt{\text{Var}(S)} \sqrt{\text{Var}(T)}$$

$$(1 + b'(\theta))^2 \leq \text{Var}(S) \text{Var}(T)$$

$$\leq I(\theta) \text{Var}(T)$$

$$\boxed{\text{Var}(T) \geq \frac{(1 + b'(\theta))^2}{I(\theta)}}$$

$T$  unbiased  $\rightarrow b = 0$

$$\boxed{\text{Var}(T) \geq I(\theta)^{-1}}$$

(why relevant to MLE?)  
→ MLE is efficient  $\Rightarrow$  achieves CRLB as  $N \rightarrow \infty$   
→  $\Rightarrow$  no consistent estimator is more efficient than MLE  
(asymptotically?)

**CRLB**

gives lower bound on minimum poss. variance for an estimator.

if there are 2 or more unbiased estimators, the one with the lowest variance is often preferred.

useful property

## CRLB Multivariate case

$$\vec{X} \sim P(\vec{x} | \vec{\theta})$$

estimator for  
✓

and  $\vec{T}(\vec{x})$  is unbiased for  $\vec{\theta}$ ,

$$E[\vec{T}(\vec{x})] = \vec{\theta}, \quad E[T_j(\vec{x})] = \theta_j$$

Then  $\boxed{\text{Cov}(\vec{T}(\vec{x})) \geq \mathbb{I}^{-1}(\theta)}$

thus elements  
 $\text{Cov}(T_j, T_k)$

Fisher Matrix  $\mathbb{I}(\theta)$  has  
elements

$$I_{jk} = E\left[-\frac{\partial^2}{\partial \theta_j \partial \theta_k} \ln P(\vec{x} | \vec{\theta})\right]$$

Hessian

In particular:

$$\boxed{\text{Var}(T_i) \geq [I^{-1}(\theta)]_{ii}}$$

(diagonal terms)

$$H(\theta) = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} \\ \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_2^2} \end{pmatrix}$$

full eval at  
 $\theta = \hat{\theta}$

## Max Likelihood

$$X_i \stackrel{\text{iid}}{\sim} f(x|\theta) \quad i=1, \dots, N$$

$$L(\theta) = P(\vec{x}|\theta) = \prod_{i=1}^N f(x_i|\theta)$$

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

$$= \underset{\theta}{\operatorname{argmax}} \ln L(\theta)$$

MLE Properties (Model <sup>is</sup> true!)

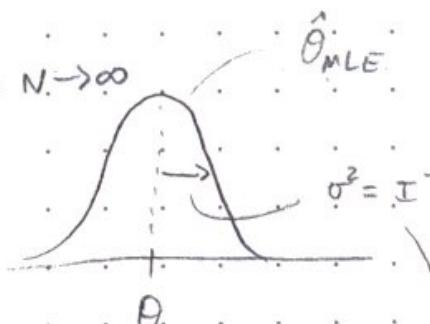
\* Consistent :  $\hat{\theta}_{\text{MLE}} \rightarrow \theta_{\text{true}}$  as  $N \rightarrow \infty$

\* Asymptotically unbiased:  $E[\hat{\theta}_{\text{MLE}}] \rightarrow \theta_{\text{true}}$  as  $N \rightarrow \infty$

\* Asymptotically Normal  $(\hat{\theta}_{\text{MLE}} - \theta_{\text{true}}) \xrightarrow{d} N(0, I^{-1})$

↳ Asymptotic normality: consistent estimators  $\hat{\theta}$  have dist around true  $\theta$  that approaches a normal dist. (variance decreasing as  $N \rightarrow \infty$ )

why gather info  
is useful: find  
covariance matrices  
associated w/ MLE



achieves  
CRLB as  $N \rightarrow \infty$   
(efficient!)  
no other estimator  
is more efficient

Max Likelihood Properties (cont.)

$$X_i \stackrel{iid}{\sim} f(x|\theta) \quad i=1, \dots, N$$

$$L(\theta) = P(\vec{x}|\theta) = \prod_{i=1}^N f(x_i|\theta)$$

$$\text{MLE} \rightarrow \hat{\theta}_{\text{MLE}} = \operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} \ln L(\theta)$$

\* Consistent:  $\hat{\theta}_{\text{MLE}} \rightarrow \theta_{\text{true}}$  as  $N \rightarrow \infty$

$$(\forall \epsilon > 0, \Pr(|\hat{\theta}_{\text{MLE}} - \theta_{\text{true}}| > \epsilon) \rightarrow 0 \text{ as } N \rightarrow \infty)$$

\* Asymptotically Unbiased:  $\mathbb{E}[\hat{\theta}_{\text{MLE}}] \rightarrow \theta_{\text{true}}$  as  $N \rightarrow \infty$   
(but not necessarily unbiased)

\* Asymptotically Normal:  $(\hat{\theta}_{\text{MLE}} - \theta_{\text{true}}) \xrightarrow{d} N(0, I^{-1})$

$$I(\theta) = -\mathbb{E}\left[\frac{\partial^2 \ell(\theta)}{\partial \theta^2}\right] = -N \mathbb{E}\left[\frac{\partial^2 \ln f}{\partial \theta^2}\right]$$

↑ Expected Fisher Information

$$\hat{I} = \text{"observed Fisher Info."} = -\left.\frac{\partial^2 \ell}{\partial \theta^2}\right|_{\theta=\hat{\theta}_{\text{MLE}}}$$

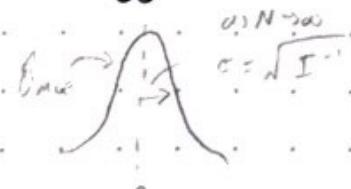
$$\approx I(\theta) \text{ as } N \rightarrow \infty$$

$\Rightarrow$  asymptotic normality

\* Efficient: Asymptotically achieves CRLB

$$\operatorname{Var}(\hat{\theta}_{\text{MLE}}) \rightarrow I^{-1} \text{ as } N \rightarrow \infty$$

"saturates"



with  
"equivariance property"

\* Functionally invariant:  $\alpha = g(\theta) \leftarrow$  (parameter transf.)  
 $\hat{\alpha}_{MLE} = g(\hat{\theta}_{MLE})$

if  $\hat{\theta}_{MLE}$  is MLE for  $\theta$   
 $g(\theta)$  is bijective transf.  
of  $\theta$  then MLE for  
 $\alpha = g(\theta)$  is  $\hat{\alpha}_{MLE} = g(\hat{\theta}_{MLE})$

## Multiparameter Case MLE

$$L(\vec{\theta}) = p(\vec{x}|\vec{\theta}), \quad \vec{\theta}_{MLE} = \underset{\vec{\theta}}{\operatorname{argmax}} L(\vec{\theta}) \\ = \prod_{i=1}^n p(x_i|\theta)$$

\* Asymptotically:  $\vec{\theta}_{MLE} - \vec{\theta}_{true} \xrightarrow{d} N(\vec{0}, \frac{1}{N} \mathbb{I}^{-1})$

Inverse Fisher Matrix

$I_{jk} = \text{Exp. Fisher Info Matrix}$

$$= \mathbb{E} \left[ -\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k} \right] = -N \mathbb{E} \left[ \frac{\partial^2 \ell_{\text{exp}}}{\partial \theta_j \partial \theta_k} \right]$$

similarly can define

Observed F.I. =  $\hat{I}_{jk} = -\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k} \Big|_{\vec{\theta} = \vec{\theta}_{MLE}}$

$\approx I_{jk}$  as  $N \rightarrow \infty$

\* Asymptotically Efficient

usually we want to know variance of particular components

$$\text{Cov}(\vec{\theta}_{MLE}) \rightarrow \frac{1}{N} \mathbb{I}^{-1}$$

In particular,

$$\text{Var}(\theta_{MLE,i}) \rightarrow (\mathbb{I}^{-1})_{ii} \quad (\text{diagonal})$$

Quick Example:

$$X_i \stackrel{iid}{\sim} N(\mu, \sigma^2), \quad i=1, \dots, N$$

$$\vec{\theta} = (\mu, \sigma^2)$$

$$L(\vec{\theta}) = \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}}$$

$$\rightarrow \ell(\vec{\theta}) = \sum_{i=1}^N -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}$$

(1st deriv.)

$$\frac{\partial \ell}{\partial \mu} = \sum_{i=1}^N \frac{(x_i - \mu)}{\sigma^2} = 0$$

$$\rightarrow \hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i = \bar{x} \rightarrow \text{Var}(\hat{\mu}_{MLE}) = \sigma^2/N$$

(2nd deriv.)

$$\frac{\partial^2 \ell}{\partial \mu^2} = \sum_{i=1}^N -\frac{1}{\sigma^2} = -N/\sigma^2$$

(1st deriv. with  $\sigma^2$ )

$$\frac{\partial \ell}{\partial \sigma^2} = \sum_{i=1}^N -\frac{1}{2} \frac{1}{\sigma^2} + \frac{1}{2} \frac{(x_i - \mu)^2}{(\sigma^2)^2} = 0$$

$$\rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

(N.B. this is biased!  $E[\hat{\sigma}_{MLE}^2] = \frac{N-1}{N} \sigma^2$ )

But sample variance  $s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$

is unbiased for  $\sigma^2$ .

(2nd deriv. with  $\sigma^2$ )

$$\frac{\partial^2 \ell}{\partial \sigma^2} = \sum_{i=1}^N \frac{1}{2\sigma^4} - \frac{(x_i - \mu)^2}{\sigma^6}$$

$$E\left[\frac{\partial^2 \ell}{\partial \sigma^2}\right] = -\frac{N}{2\sigma^4} \quad (\text{using } E[(x_i - \mu)^2] = \sigma^2)$$

check this yourself

Cross term  $\mathbb{E} \left[ \frac{\partial^2 \ell}{\partial \mu^2 \partial \sigma^2} \right] = 0$  (using  $\mathbb{E}[x_i - \mu] = 0$ )

$$\overline{I} = \begin{pmatrix} -\frac{\partial^2 \ell}{\partial \mu^2} & -\frac{\partial^2 \ell}{\partial \mu \partial \sigma^2} \\ -\frac{\partial^2 \ell}{\partial \sigma^2 \partial \mu} & -\frac{\partial^2 \ell}{\partial \sigma^4} \end{pmatrix} = \begin{pmatrix} N/\sigma^2 & 0 \\ 0 & N/2\sigma^4 \end{pmatrix}$$

(actually  $\mathbb{E}\left[\frac{\partial^2 \ell}{\partial \mu^2}\right]$  etc?)

$$\rightarrow \overline{I}^{-1} = \begin{pmatrix} \sigma^2/N & 0 \\ 0 & 2\sigma^4/N \end{pmatrix}$$

$$\begin{aligned} \text{Var}(\hat{\mu}_{MLE}) &= \sigma^2/N \stackrel{?}{\geq} (\overline{I}^{-1})_{11} \\ \text{Var}(\hat{\sigma}^2_{MLE}) &\stackrel{?}{\geq} (\overline{I}^{-1})_{22} \end{aligned}$$

check this  
yourself  
vs variance greater  
than our bound from NLE  
theory?

$$\text{Var}(\hat{\mu}_{MLE}) = \overline{I}^{-1}_{11} \quad \text{as } N \rightarrow \infty \quad \checkmark \quad \text{as expected}$$

# Laplace Approximation:

remember  
that:  
1) observed  
2) new info

recall  $\ell(\theta) = \ln L(\theta)$

Taylor approx around MLE:

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

score = 0 because  
we maximize to  
find MLE

$$-\hat{\mathbf{I}}$$

$$\underbrace{\quad}_{\text{+ ...}}$$

$$\frac{\partial^2 \ell}{\partial \theta^2} \Big|_{\hat{\theta}_{MLE}} (\theta - \hat{\theta}_{MLE})^2$$

Asymptotically,  $N \rightarrow \text{large}$

$$\ell(\theta) \approx \ell(\hat{\theta}_{MLE}) - \frac{1}{2} \hat{\mathbf{I}}^{-1} (\theta - \hat{\theta}_{MLE})^2$$

$$L(\theta) \approx L(\hat{\theta}_{MLE}) e^{-\frac{1}{2} (\theta - \hat{\theta}_{MLE})^2 / \hat{\mathbf{I}}^{-1}}$$

(Gaussian shape)

useful b/c can summarize  
in terms of  $N, \sigma$  or  $\mu, \Sigma$   
for MV.

Multiparameter case

$$L(\vec{\theta}) \approx L(\hat{\vec{\theta}}_{MLE}) \times e^{-\frac{1}{2} (\vec{\theta} - \hat{\vec{\theta}}_{MLE})^T \hat{\mathbf{I}}^{-1} (\vec{\theta} - \hat{\vec{\theta}}_{MLE})}$$

Reminder:

multidim Taylor expansion

$$f(\underline{x}) \approx f(\underline{a}) + \nabla f(\underline{a}) \cdot (\underline{x} - \underline{a}) + \frac{1}{2} (\underline{x} - \underline{a})^T H(\underline{a}) (\underline{x} - \underline{a})^T + \dots$$

grad f  
at a

hessian of f  
at a

$$\text{e.g. } f(x, y) \approx f(a, b) + \left. \frac{\partial f}{\partial x} \right|_{a,b} (x-a) + \left. \frac{\partial f}{\partial y} \right|_{a,b} (y-b) + \frac{1}{2} \left. \frac{\partial^2 f}{\partial x^2} \right|_{a,b} (x-a)^2 + \left. \frac{\partial^2 f}{\partial y^2} \right|_{a,b} (y-b)^2 + \frac{1}{2} \left. \frac{\partial^2 f}{\partial x \partial y} \right|_{a,b} (x-a)(y-b) + \frac{1}{2} \left. \frac{\partial^2 f}{\partial y \partial x} \right|_{a,b} (y-b)(x-a)$$

[Example (slides): Calibrating type Ia supernova absolute magnitudes with intrinsic dispersion and heteroskedastic measurement error. (Ex sheet 1)]

Determining Astronomical Distances using Standard Candles → key step in determining Hubble constant.

- ① Estimate or model luminosity  $L$  of a class of astronomical objects
- ② Measure apparent brightness or flux  $F$
- ③ Derive distance  $D$  to object using inverse square law

$$F = L / 4\pi D^2$$

- ④ Optical astronomer's units  $m = m - M$

$m$  = apparent magnitude (log apparent brightness, flux)

$M$  = absolute magnitude (log luminosity)

$M$  = distance modulus (log distance)

can use Hubble to find distances

Hubble:  $V = H_0 \times \text{distance}$

distance  $\propto$  velocity (redshift)

Hubbles law  
 $z \approx H_0 d$

Now use type Ia supernovae as standard candles

We want to measure the Hubble constant

→ can use (local) distance ladder → Hubble tension is dispersion between using this method & another one!

Can use cepheid stars: period of rotation linked to luminosity (pulsating)

1st tier: geometry, 2nd tier: cepheids, 3rd tier: supernovae

use inbetween tiers to calibrate e.g. same distance cepheids & supernovae

3. rungs of the distance ladder

hubble law - things that are further away tend to move away from us faster

calibrate when they are in the same galaxy

## Lecture 8

10.2.25

### Recap example (slides)

- Measure apparent brightness or flux  $F$
- Derive distance to object using inverse square law

$$F = L / 4\pi D^2$$

- key step in determining hubble constant ( $= \frac{\text{redshift}}{D}$ )
- 3 rungs to distance ladder: geometry, cepheids, supernovae (can calibrate between rungs)

we will focus on cepheids

back on board!

### Calibrating SNIa Abs Mags (board)

Calibrator sample we observe:  $N$  SNe,  $s=1, \dots, N$

$M_s$  = true apparent magnitude

$\hat{M}_s$  = measured apparent magnitude

$$\sim N(M_s, \sigma_{m,s}^2) \leftarrow \begin{array}{l} \text{given} \\ \text{in data table} \end{array}$$

equivalently,  $M_s = \hat{M}_s + E_{ms} \sim N(0, \sigma_{m,s}^2) \leftarrow \text{data}$

Cepheid distance estimates to SN galaxy:

$\mu_s$  = true distance modulus

$$= 25 + 5 \log_{10} \left( \frac{\text{distance}}{\text{Mpc}} \right)$$

again  
measured  
error in dist.  $\rightarrow \hat{\mu}_s \sim N(\mu_s, \sigma_{\mu,s}^2) \leftarrow \begin{array}{l} \text{given} \\ \text{in data table} \end{array}$

equiv.  $\hat{\mu}_s = \mu_s + E_{\mu s} \sim N(0, \sigma_{\mu,s}^2)$   
 "data"

## Abs Mag Distribution (Population):

$$M_s \sim N(M_0, \sigma_{int}^2)$$

↑  
unknown

what we want to estimate  
intrinsic error / not like  
measurement error  
intrinsic due to different mass &  
gas properties etc

$$\text{equiv. } M_s = M_0 + E_{int,s} \rightarrow N(0, \sigma_{int}^2)$$

like before, want to estimate  $\theta$ , so form likelihood eq:  
ie variable without bias  
 $M_s = M_s + \mu_s$  (Latent variable equation)  
relate the latent variables

Define:  $\hat{M}_s = \hat{M}_s - \bar{\mu}_s$  (Estimated Abs Mag)

$$= M_s + E_{int,s} - \mu_s - E_{\mu,s}$$

$$= M_s + E_{int,s} - E_{\mu,s}$$

$$= M_s + E_{int,s} \rightarrow N(0, \sigma_{int,s}^2 + \sigma_{\mu,s}^2)$$

use:  
variance of sum of  
gaussian rvs is  
sum of variance

combine errors

also can  
define

$$\hat{M}_s = M_0 + E_{int,s} + E_{err,s}$$

$$\hat{M}_s \sim N(M_0, \sigma_{int}^2 + \sigma_{err,s}^2)$$

old & new  
mean

combine variances  
( $M_0$  variance)

diff SN in  
diff galaxies  
is indep but not  
perpendicularly  
dist.

## Likelihood:

given independence  
between SNe

$$P(\hat{M}_s | M_0, \sigma_{int}^2) = N(\hat{M}_s | M_0, \sigma_{int}^2 + \sigma_{err,s}^2)$$

(independence)

$$L(\theta) = P(\hat{M}_s | M_0, \sigma_{int}^2) = \prod_{s=1}^N N(\hat{M}_s | M_0, \sigma_{int}^2 + \sigma_{err,s}^2)$$

example: heteroskedastic error is multiple sources of uncertainty.

$\sigma_{err,s}$  creates problem for MLE, cannot be solved analytically.  
so solve numerically

(see slides for code).  
example

cannot solve MLE  
analytically!

# Supernova Absolute Magnitude Distribution:

## Selection Effects

(board)

brighter SNs can only see  
brighter mag. miss dimmer things

same ex. but now

Assume measurement error =  $\sigma = \sigma_{m,s} = \sigma_{M,s}$

Population variability  $\sigma_{int} \rightarrow \sigma$

Only intrinsic population is source of variability  
( $\sigma^2 = \sigma_{int}^2$ )

Suppose distance is the same for an entire sample of SN

$$M_s = M$$

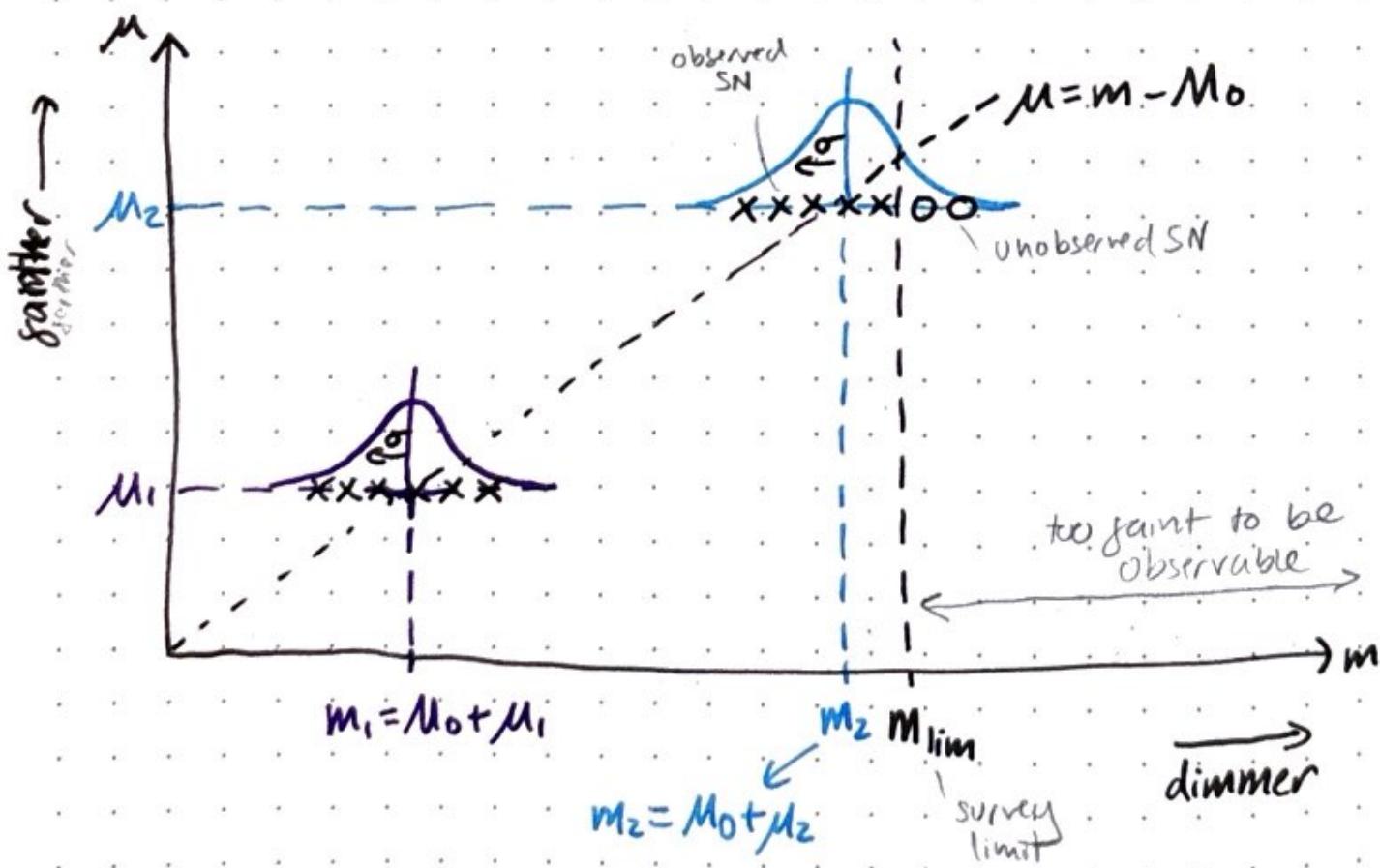
population dist

$$M_s \sim N(M_0, \sigma^2) \quad \text{iid}$$

$$(M_s = M_s + \mu_s) \quad \text{latent variable eq}$$

combine

$$\rightarrow M_s \sim N(M_0 + \mu, \sigma^2), \quad s=1, \dots, N$$



big topic)

Survey can only detect up to certain magnitude limit, can't see anything dimmer than this. Suppose we have population  $N(M_0, \sigma^2)$ . At  $m_2$  magnitudes look similar. But in the 2nd we don't see SN above some limit, what do we do? →

Data:  $\{m_s\} = \vec{m}$

fixed  $\mu$  for sample

Naive Likelihood

$$P(m_s | M_0, \sigma^2) = N(m_s | \mu + M_0, \sigma^2)$$

$$L(M_0, \sigma^2) = \prod_{s=1}^N P(m_s | M_0, \sigma^2)$$

↓ naively apply naive likelihood

If  $M_0 + \mu$  close to  $M_{lim}$ , MLE will be biased:

$$\left\{ \begin{array}{l} \hat{M}_0 \text{ too bright} \\ \hat{\sigma} \text{ too small} \end{array} \right. \quad \begin{array}{l} \text{(b/c didn't see} \\ \text{the dim covariance)} \end{array}$$

$\left( \begin{array}{l} \text{b/c tail of dist} \\ \text{is cut off} \\ \text{apparent width smaller than} \\ \text{true width} \end{array} \right)$

how do we solve this?

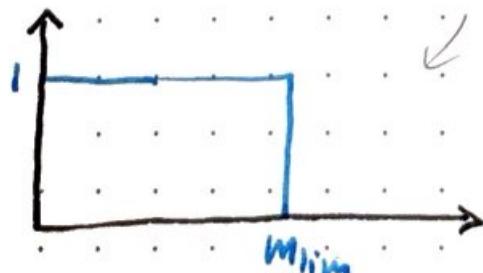
↓  
Accounting For Selection Effects (sheet 1.9.2)

Let  $I_s = \begin{cases} 1, & \text{if SN observed} \\ 0, & \text{if SN NOT observed} \end{cases}$

Define Selection Function:

prob indicator given data

$$P(I_s | m_s) = \begin{cases} 1, & m_s < M_{lim} \\ 0, & m_s \geq M_{lim} \end{cases}$$
$$= 1 - H(m_s - M_{lim})$$



↙ Heaviside step fun

## Lecture 9

12.2.25

(brief recap of SN problem setup on board - useful for sheet!)

### SELECTION EFFECTS cont.

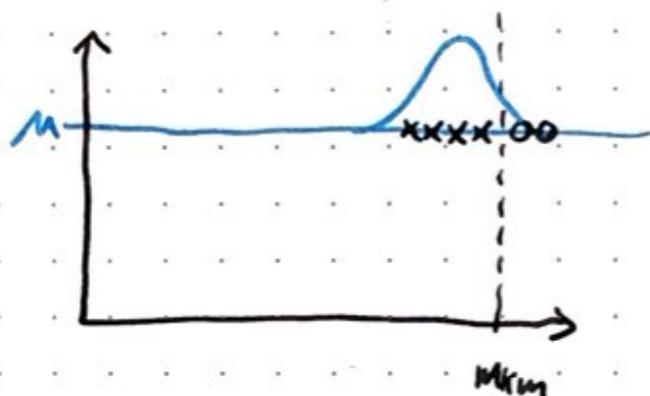
recap:

$$M_s = N(M_0, \sigma^2)$$

$$s = 1, \dots, N$$

$$M_s = M' \text{ (known)}$$

known  
object so  
no measurement  
error

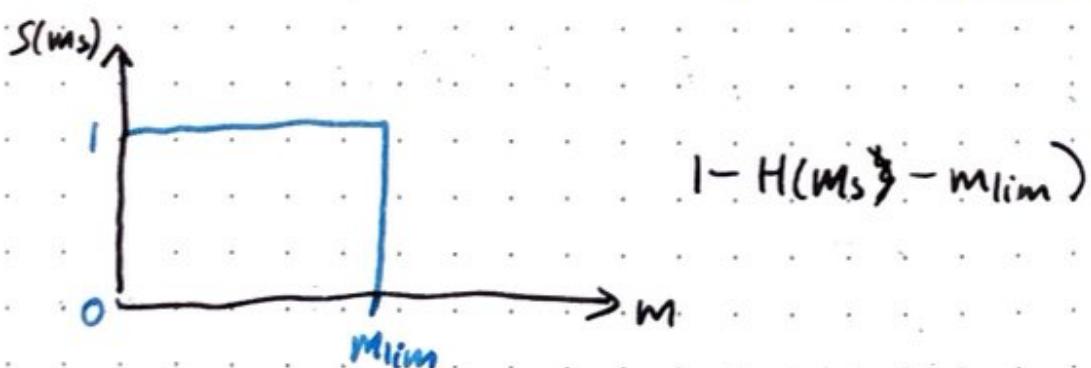


Formulate likelihood to account for selection effects:

$$\text{Let } I_s = \begin{cases} 1 & , \text{ if SN observed} \\ 0 & , \text{ if SN NOT observed} \end{cases}$$

Selection function:

$$S(M_s) = P(I_s | M_s) = \begin{cases} 1 & , M_s < M_{lim} \\ 0 & , M_s \geq M_{lim} \end{cases}$$

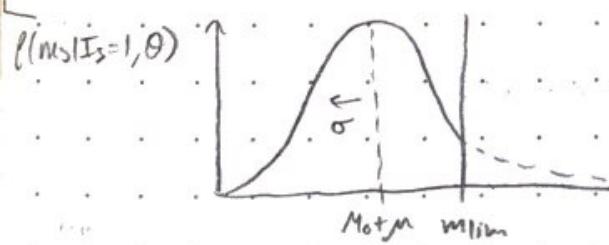


## Observed Data Likelihood:

$$\begin{aligned}
 P(m_s | I_s = 1, \theta) &= \frac{P(I_s = 1, m_s | \theta)}{P(I_s = 1 | \theta)} \\
 &= \frac{P(I_s = 1 | m_s, \theta) P(m_s | \theta)}{\int P(I_s = 1 | m_s, \theta) P(m_s | \theta) dm_s} \\
 &= \frac{S(m_s) N(m_s | M_0 + \mu, \sigma^2)}{\int S(m_s) N(m_s | M_0 + \mu, \sigma^2) dm_s}
 \end{aligned}$$

Gaussian CDF  $\Phi$   $\rightarrow \int_{-\infty}^{M_{\text{lim}}} N(m_s | M_0 + \mu, \sigma^2) dm_s$  normalize

(TRUNCATED)  
NORMAL



"look this up on wikipedia"

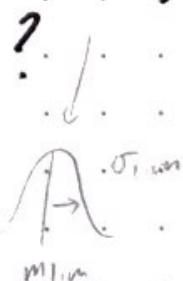
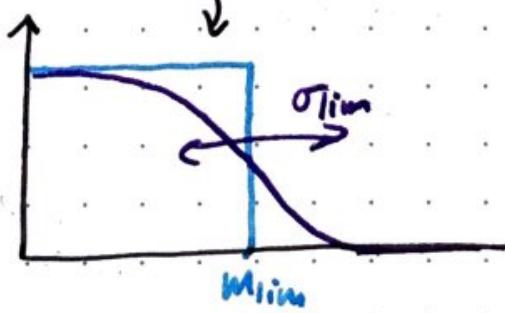
$$= TN(m_s | M_0 + \mu, \sigma^2, -\infty, M_{\text{lim}})$$

↑ untruncated mean & variance  
 ↑ lower truncation limit  
 ↑ upper truncation limit

Challenge: What if  $S(m_s) = P(I_s | m_s) = \Phi\left(\frac{M_{\text{lim}} - m_s}{\sigma_{\text{lim}}}\right)$

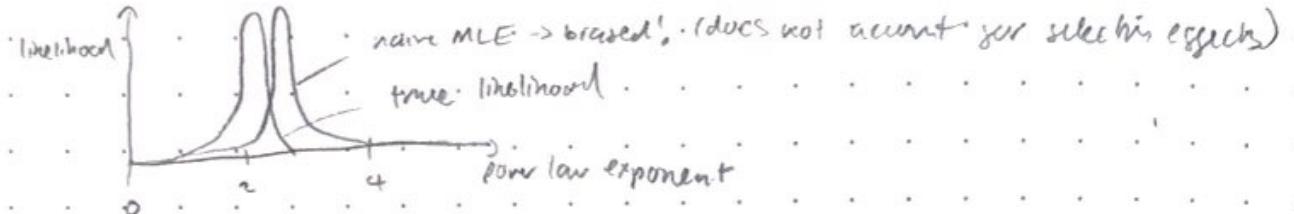
(inc regular selection  
mt to take into account  
survey limit is not the  
same every night)

limit not a fixed boundary.  
to take into account e.g.  
observing more / see further as  
clearer days



(slides)

- Ex sheet 1 q 2: star formation in Perseus (clouds of collapse under grav. & form stars here)  
→ in densest cloud regions stars are hidden  
→ densest regions will be where largest stars form → selection effect  
→ pareto or power law distribution.  $P(m) \propto m^{-\gamma}$ . ( $m > 60$ )



## QUANTIFYING UNCERTAINTY USING BOOTSTRAP

### Frequentist interpretation:

Consider variability of your estimator  $g(\bar{x})$  for  $\theta$  under (imaginary) repetitions of your experiment. (Random realisations of the potential data).

How does  $g(\bar{x})$  behave under the potential datasets you did not observe?

$$\text{e.g. } \text{Var}[g(\bar{x})] = \mathbb{E}^{\text{under } P(x|\theta)}[g(\bar{x}) - \mathbb{E}[g(\bar{x})]].$$

If  $g(\bar{x})$  is approximately Gaussian distributed

$$\rightarrow 68\% \text{ confidence interval } g(\bar{x}) \pm \sqrt{\text{Var}(g(\bar{x}))}$$

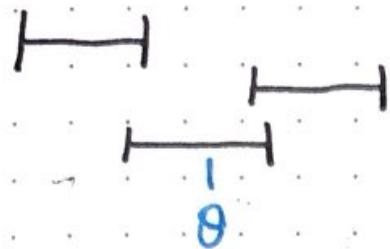
standard notation in science

$$[g(\bar{x}) - \sqrt{\text{Var}(g(\bar{x}))}, g(\bar{x}) + \sqrt{\text{Var}(g(\bar{x}))}]$$

$(1-\alpha)\%$  confidence interval  $[L(\bar{x}), U(\bar{x})]$  contains the true value  $\theta_{\text{true}}$  in at least  $(1-\alpha)\%$  of the realisations.

( $\approx$  repeat experiments, gathering intervals that contain true parameter)  
will approach 68%

NOT prob a particular interval contains the true parameter (common confusion)



get diff. interval. for each experiment  $\rightarrow$   
68% confidence interval  $\rightarrow$  68% of  
these realizations contain  $\theta$

$[L(x), U(x)]$  is a random interval vs.  $[L(x_{\text{obs}}), U(x_{\text{obs}})]$   
evaluated on observed numerical values  $\uparrow$ , only one  
dataset  $x_{\text{obs}}$ ! (either contains  $\theta$  true or doesn't)

## Bootstrap:

Use the observed dataset to simulate the variability of the unobserved (imaginary) data sets.

BOOTSTRAP SAMPLE = sample with replacement from the observed dataset to the sample size

Eq.  $x_1, \dots, x_5 \stackrel{\text{iid}}{\sim} \text{Poisson } (\lambda)$

$$\rightarrow P(x_i) = \frac{x_i^{\lambda} e^{-\lambda}}{x_i!}$$

could do w/ max. likelihood  
but suppose don't know dist.

Real data (observed)  $\vec{x}_{\text{obs}} = (3, 8, 2, 4, 5)$

Suppose you want to estimate the skewness of  $P(x)$  (asymmetry)

i.e. skewness =  $\frac{\mathbb{E}[(x-\mu)^3]}{(\sigma^2)^{3/2}}$

mean  $\rightarrow = \lambda$   
variance  $\rightarrow$

Sample skewness  $g(\bar{x}) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^3$

$$\left( \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} \right)^3$$

Bootstrap B "replicate" datasets from observed dataset: "sample w/ replacement" from original org. sample N times from 3, 8, 2, 4, 5 gives 2, 5, 4, 4, 4

$$\vec{X}^{\text{obs}} = (3, 8, 2, 4, 5)$$

$$\hat{g}^{\text{obs}} = g(\vec{X}^{\text{obs}}) = 0.6927$$

$$\vec{X}^{b=1} = (2, 5, 4, 4, 4)$$

$$\hat{g}_1 = g(\vec{X}^{b=1}) = -0.8675$$

$$\vec{X}^{b=2} = (2, 4, 2, 8, 8)$$

$$\hat{g}_2 = g(\vec{X}^{b=2}) = 0.2115$$

$$\vec{X}^{b=3} = (5, 2, 8, 2, 5)$$

$$\hat{g}_3 = g(\vec{X}^{b=3}) = 0.3436$$

⋮

⋮

$$\vec{X}^{b=B} =$$

$$\hat{g}_B =$$

Can now compute sample variance

$$\hat{\text{Var}}(\{\hat{g}_1, \dots, \hat{g}_B\})$$

$$\rightarrow \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (3)$$

$$\text{Standard error} = \sqrt{\hat{\text{Var}}}$$

$$\hat{g} = 0.6927 \pm 0.635 \approx 68\% \text{ C.I.} //$$

(slides)

back to stereogrammar peters example:

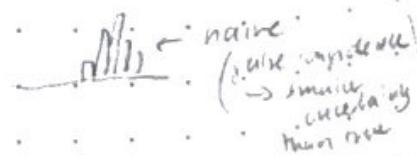
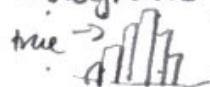
why do i get standard error on exponent?

running Fisher matrix approach and compare this against bootstraps

fit model to each bootstrapped realization



Plot Maximum likelihood estimates for each  
as histograms



## REGRESSION

(slides)

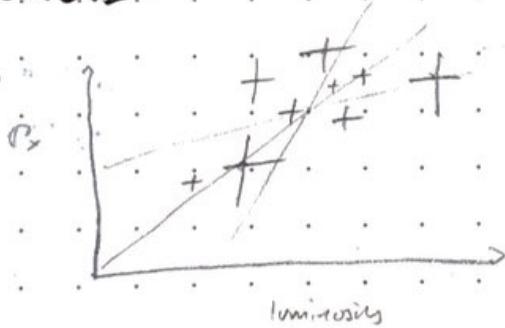
\* Fit a function  $E[y|x] = g(x; \theta)$  for the mean relation between  $y$  and  $x$

\* Basic approaches

- ordinary least squares (homoskedastic scatter)
- generalized least squares (heteroskedastic, correlated scatter)
- weighted least squares (minimum  $\chi^2$ , known variance)
- maximum likelihood

\* Real data problems require more complex modelling

- regression dilution from covariate measurement errors



e.g.  
each graser has  
diff measurement error  
regression dilution: namely apply  
OLS gives biased slope when  
there are errors in  $x$  (sheet 2)

## Ordinary least squares (OLS):

Linear model  $y_i = \beta_0 + \sum_{j=1}^{k-1} \beta_j x_{ij} + \epsilon_i$

$$Y = X\beta + \epsilon$$

$i=1, \dots, N$  objects,  $E[\epsilon_i] = 0$ , homoskedastic  
 $\text{Var}[\epsilon_i] = \sigma^2$  (known).

$$y_i = \beta_0 + \sum_{j=1}^{k-1} \beta_j x_{ij} + \epsilon_i$$

residual sum of squares

(don't think this is too important to mention?  
(slides))

argue invariance in y value even in no

minimum sum square of these distances

Minimise w.r.t.  $\beta$  (solve for gradient = 0):

$$RSS = \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^{k-1} \beta_j x_{ij})^2 = (Y - X\beta)^T (Y - X\beta)$$

("simple linear regression" von Wacker) per  
EZ PICO WURF

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y$$

just accept no?

$$\text{Var}(\hat{\beta}_{OLS}) = \sigma^2 (X^T X)^{-1}$$

unbiased  $\beta$  has minimum variance under these assumptions

$$\mathbb{E}(\hat{\beta}_{OLS}) = \beta \text{ (unbiased)}$$

BLUE

Estimate unknown variance

OLS estimator of variance is  $\hat{\sigma}^2 = \frac{RSS}{N-k}$   
(see 2 previous examples)

$$\hat{\sigma}^2 = \frac{1}{N-k} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$$

K=2 for simple linear y =  $a + b x$ ?

(WEIGHTED LEAST SQUARES) — aka  $\chi^2$  minimisation:

$$\chi^2 = \sum_{i=1}^N \frac{(y_i - \beta_0 - \sum_{j=1}^{k-1} \beta_j x_{ij})^2}{\sigma_i^2}$$

$\chi^2$  r.v. = sum of squared Gaussian r.v.s

Minimise w.r.t.  $\beta$ :

If Gaussian errors, at  $\beta = \beta_{min}$ ,  $\chi^2 \sim \chi^2_{N-k}$

Model check:  $\mathbb{E}(\chi^2_{N-k}) = N-k \rightarrow$  ( $\chi^2$  assumptions (confidence) true)

$$\frac{\chi^2}{N-k} \approx 1 \quad (\text{for large } N-k)$$

$\chi^2$  statistic should have a  $\chi^2$  distribution with  $N-k$  deg. of freedom  
(if gaussian errors?)

reduced  $\chi^2$  statistic: a way to test model fit

These are special cases of generalised least squares.

Linear model  $y_i = \beta_0 + \sum_{j=1}^{k-1} \beta_j x_{ij} + \epsilon_i$ ,  $i=1, \dots, N$  objects  
 $\downarrow$   
 $Y = X\beta + \epsilon$

$$\mathbb{E}[\epsilon_i] = 0, \text{ correlated errors}$$

$$\text{Var}[\epsilon] = \text{Cov}[\epsilon, \epsilon^T] = W \text{ (known)}$$

Minimise w.r.t.  $\beta$

$$\text{RSS} = (Y - X\beta)^T W^{-1} (Y - X\beta)$$

$$\hat{\beta}_{\text{GLS}} = (X^T W^{-1} X)^{-1} X^T W^{-1} Y$$

$$\mathbb{E}[\hat{\beta}_{\text{GLS}}] = \beta \text{ (unbiased)}$$

$$\text{Var}[\hat{\beta}_{\text{GLS}}] = (X^T W X)^{-1}$$

They can also be thought of as maximum likelihood (assuming Gaussian errors)

$$Y = X\beta + \epsilon, \quad Y \sim N(X\beta, W)$$

Maximise w.r.t.  $\beta$  (solve for gradient = 0)

$$L(\beta) = P(Y|\beta, X) = N(Y|X\beta, W)$$

do we get  
1. likelihood eqn  
 $\beta$  (and also  
 $w$ )

(end of slides)

minimum  $\chi^2$  estimation: find values of  $\theta$  that make  $\chi^2(\theta)$  as small as possible

## $\chi^2$ MINIMISATION

(board)  
(vs max likelihood)

$(x_i, y_i), i=1, \dots, N$

Case: Variance unknown

$$y_i = g(x_i; \theta) + \epsilon_i, \quad \text{Var}(\epsilon_i) = \sigma_i^2 \text{ known}$$
$$= a + bx + \epsilon_i \text{ (linear)}$$

$$\chi^2(\theta) = \sum_{i=1}^N \frac{(y_i - g(x_i; \theta))^2}{\sigma_i^2}$$

$$\hat{\theta} = \arg \min_{\theta} \chi^2$$

Relation to maximum likelihood:

$$y_i = g(x_i; \theta) + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma_i^2)$$

$$P(y_i | x_i) = N(y_i | g(x_i; \theta), \sigma_i^2)$$

$$L(\theta) = \prod_{i=1}^N P(y_i | x_i; \theta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(y_i - g(x_i; \theta))^2}{2\sigma_i^2}}$$

$$-2 \ln L(\theta) = \sum_{i=1}^N \ln(2\pi\sigma_i^2) + \sum_{i=1}^N \frac{(y_i - g(x_i; \theta))^2}{\sigma_i^2}$$

$$= \sum_{i=1}^N \ln(2\pi\sigma_i^2) + \chi^2(\theta) \quad \chi^2(\theta) + \text{const.}$$

( $\min \chi^2(\theta)$  same as max likelihood)

→ If  $\sigma_i$  known:  $\hat{\theta}_{\min \chi^2} = \hat{\theta}_{MLE}$

→ If  $\sigma_i^2$  (or a component of it) unknown then NOT true!

(Note: you get even  $y_i$ ,  $m_i$  is same as  $y_i$ )

Case: Variance component unknowns: unknown

$$Y_i = g(x_{ij}|\theta) + \epsilon_{int} + \epsilon_m$$

$\epsilon_{int} \sim N(0, \sigma_{int}^2)$ $\epsilon_m \sim N(0, \sigma_{m,i}^2)$	<span style="color: blue;">known</span>
--	---

$$\chi^2(\theta, \sigma_{int}^2) = \sum_{i=1}^N \frac{(Y_i - g(x_{ij}|\theta))^2}{\sigma_{int}^2 + \sigma_{m,i}^2}$$

Problem:  $\chi^2$  is minimised when  $\sigma_{int} \rightarrow 00$  !!

Max likelihood:

$$P(Y_i | X_i, \theta, \sigma_{int}^2) = N(Y_i | g(x_{ij}|\theta), \sigma_{m,i}^2 + \sigma_{int}^2)$$

$$L(\theta, \sigma_{int}^2) = \prod_{i=1}^N \left[ 2\pi (\sigma_{int,i}^2 + \sigma_{m,i}^2) \right]^{-\frac{1}{2}} e^{-\frac{(Y_i - g(x_{ij}|\theta))^2}{\sigma_{int,i}^2 + \sigma_{m,i}^2}}$$

$$-2 \ln L(\theta, \sigma_{int}^2) = \sum_{i=1}^N \frac{(Y_i - g(x_{ij}|\theta))^2}{\sigma_{int}^2 + \sigma_{m,i}^2} + \ln [2\pi (\sigma_{int}^2 + \sigma_{m,i}^2)]$$

(1) (2)

When  $\sigma_{int} \rightarrow$  large, (1)  $\downarrow$  but (2)  $\uparrow$ .

Max likelihood  $\rightarrow$  finite estimate of  $\sigma_{int}$ !

→ assuming Gaussian, independent & known errors  
→ convenient but not robust  
→ greatest principle likelihood is not necessarily appropriate

Lesson:  $\chi^2$  is an ad-hoc prescription.  
Likelihood is derived from explicit modelling assumptions.

slides: OLS leads to regression dilution when there are errors in  $x$  (usually measurement errors) which causes slope to be shallower than true value aka regression dilution

OLS gives shallow slope if you have error on  $x$  data.  
(regression dilution)

# Probabilistic Generative Modelling

(slides)

$$P(D|\alpha, \theta)$$

- \* Forward model comprises series of probabilistic steps describing conceptually how the observed data was generated from the parameters of interest.
- \* Can introduce intermediate parameters / unobserved latent variables  $\alpha$  (e.g. true values corresponding to the observed data)
  - what data would be had if there been no measurement error.
- \* From forward model, derive the sampling distribution e.g.
$$P(D|\theta) = \int P(D|\alpha) P(\alpha|\theta) d\alpha$$
integrate out to leave params of interest
- \* Using observed data  $D$ , draw inference from likelihood function
$$L(\theta) = P(D|\theta)$$
- \* Or if Bayesian with prior  $P(\theta)$ : sample posterior
$$P(\theta|D) \propto P(D|\theta) P(\theta)$$

parameters:  
 $\theta$  want  
 $\alpha$  nuisance

$$P(\theta|D) = \int P(\theta, \alpha|D) d\alpha$$

use bayes:

$$\Rightarrow P(D|\theta) = \int p(D|\theta, \alpha) P(\alpha|\theta) d\alpha$$

show as before given w/  $P(\theta|D) = \int P(\theta, \alpha|D) d\alpha$

(how do I come up with model for this complex situation where I have measurement error in my x and y data)

## Generative Model

(slides) (same example from L3)

- Population distribution  $\xi \sim N(\mu, \tau^2)$
- Regression Model  $\eta_i | \xi_i \sim N(\alpha + \beta \xi_i, \sigma^2)$
- Measurement Error  $[x_i, y_i] | \xi_i, \eta_i \sim N([\xi_i, \eta_i], \Sigma)$

For this example let  $\Sigma = \begin{pmatrix} \sigma_{x,i}^2 & 0 \\ 0 & \sigma_{y,i}^2 \end{pmatrix}$  (for simplicity).

$$x_i = \xi_i + \epsilon_{x,i} \leftarrow \sigma_{x,i}^2$$

$$y_i = \eta_i + \epsilon_{y,i} \leftarrow \sigma_{y,i}^2$$

probabilistic  
generative  
model for linear  
regression

- Population dist. indep. variable  $\psi = (\mu, \tau)$
- Regression Parameters  $\theta = (\alpha, \beta, \sigma^2)$
- Latent (true) variables  $(\xi_i, \eta_i)$   
(i.e. no measurement error)
- Observed data with measurement  $(x_i, y_i)$
- uncertainties  $(\sigma_{x,i}, \sigma_{y,i})$

$$\Sigma = \begin{pmatrix} \sigma_{x,i}^2 & 0 \\ 0 & \sigma_{y,i}^2 \end{pmatrix}$$

Now can

Generate observed data from latent variables

(see slides -> same as L3)

derive this (Sheet 2a)  
from on slides

(now we've introduced all these new parameters how do I get rid of them?)  
Formulating likelihood Function: Marginalising latent variables

$$P(x_i, y_i | \theta, \psi) = \iint P(x_i, y_i, \xi_i, \eta_i | \theta, \psi) d\xi_i d\eta_i$$

"observed data likelihood" "complete data likelihood"

$$P(x_i, y_i | \theta, \psi) = \iint P(x_i, y_i | \xi_i, \eta_i) P(\eta_i | \xi_i, \theta) P(\xi_i | \psi) d\xi_i d\eta_i$$

measurement error regression ↑ population distribution of covariates

## Lecture 11

17.2.25

(slides)

- outcome from overhead: multiply marginal likelihood to get likelihood for all data
- Observed data likelihood:

$$P(x, y | \theta, \psi) = \prod_{i=1}^N P(x_i, y_i | \theta, \psi) \quad (\text{indep.})$$

In frequentist statistics: distinction between data and parameters: parameters are fixed and unknown, but not "random". Only "data" are random realisations of random variables.

### What is the nature of the latent variables $\xi_i, \eta_i$ ?

- They have distribution  $(\xi_i, \eta_i) \sim P(\xi_i, \eta_i | \theta, \psi) = P(\eta_i | \xi_i, \theta) P(\xi_i | \psi)$
- Often called "nuisance parameters"  $\rightarrow$  needed to complete the model but not the parameters of interest  $(\theta, \psi)$
- Are the latent variables "data" or "parameters"?  
 ("missing data"      "nuisance parameters")

### BAYESIAN VIEWPOINT

- \* There is a symmetry between data  $D$  and parameters  $\theta$  - both are r.v.s described by probability distributions
- \* Actually they are described by a joint probability  $P(D, \theta)$

(slides)

- \* Data are r.v.s whose realisations are observed, parameters are r.v.s not observed
- \* Goal is to infer the unobserved parameters from the observed data using the rules of probability

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

CONDITIONAL PROBABILITY

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

BAYES' THEOREM

- \* Probability interpreted as degree of belief / uncertainty in hypothesis  
how sure are you in the values of that parameter?

### Bayes' Theorem

Joint probability of data & parameters:  $P(D,\theta) = P(D|\theta)P(\theta)$   
 $= P(\theta|D)P(D)$

Probability of parameters given data:

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

prior prob.:  
degree of belief

posterior prob.:  
degree of belief

likelihood

sampling distribution

normalisation constant

Example: bayesian inference of the dimensionality of space-time (see lecture)

observe light & detect w/ interferometer  
signal from black hole merger (?)  
can use this info to infer dimensionality of spacetime

dimensionality of spacetime  
very signal vs. time to receive light  
time to compare long light  
with now  
time to reach us

## Simple Gaussian Example

(slide 3)

## Frequentist confidence vs. Bayesian credible intervals

①

### Frequentist confidence Interval

$$Y_1, \dots, Y_4 \stackrel{iid}{\sim} N(\mu, 1) \quad (\text{sample mean})$$

Sampling dist. of statistic  $\bar{Y} \sim N(\mu, \sigma^2/4)$

$[\bar{Y} - \sigma_{\bar{Y}}, \bar{Y} + \sigma_{\bar{Y}}]$  is a 68% confidence interval

- \* Under repeated experiments, 68% of the (random) confidence intervals constructed this way will contain (cover)  $\mu$

This does **NOT** mean that the probability is 68% that  $\mu$  lies within the interval (evaluated with  $y_{\text{obs}}$ ) frequentist view:  $\mu$  is a fixed number, it does not have a prob. A 68% interval will either cover that number or it doesn't

②

### Bayesian credible Interval

(board)

assume flat prior:  $p(\mu) \propto 1$

can derive posterior  $p(\mu | Y = y_{\text{obs}})$ :

$$Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

*unknown* *known*  
maybe

$$\begin{aligned} \text{Likelihood: } p(Y | \mu, \sigma^2) &= \prod_{i=1}^N N(Y_i | \mu, \sigma^2) \\ &= \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2} (Y_i - \mu)^2} \\ &= (2\pi\sigma^2)^{-N/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (Y_i - \mu)^2} \end{aligned}$$

every value  $\mu$  equally likely  
(not normalisable, cannot integrate so say  $\propto 1$  (!))

"improper prior"

Likelihood  $P(\vec{y} | \mu, \sigma^2)$  only depends on data  $\vec{y}$  through the sufficient statistics  $\bar{y}$  and  $S^2$ .

$$= (2\pi\sigma^2)^{-N/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{y})^2}$$

Any two datasets  $\vec{y}'$  and  $\vec{y}''$  with the same sufficient statistics should yield the same inferences.

$$= (2\pi\sigma^2)^{-N/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N [(y'_i - \bar{y}')^2 + 2(\bar{y}' - \mu)(y'_i - \bar{y}')] + (\bar{y}' - \mu)^2}$$

Likelihood principle: all info about parameters from data is located in likelihood function.

$$P(\bar{Y} | \mu, \sigma^2) = (2\pi\sigma^2)^{-N/2} e^{-\frac{(N-1)}{2\sigma^2} S^2} e^{-\frac{N}{2\sigma^2} (\bar{Y} - \mu)^2}$$

where  $\begin{cases} \bar{Y} = \frac{1}{N} \sum_{i=1}^N y_i & \xrightarrow{\text{sufficient statistics}} \\ S^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2. & \end{cases}$

"sample variance"  
 (unbiased estimator)  
 for  $\sigma^2$

**Case 1:**  $\sigma^2 = 1$  is known,  $P(\mu) \propto 1$  (improper prior)

Bayes' theorem  $P(\mu | \bar{Y}) \propto P(\bar{Y} | \mu) P(\mu)$   
 (ignore constant terms)

$$P(\mu | \bar{Y}) \propto e^{-\frac{N}{2\sigma^2} (\bar{Y} - \mu)^2} \quad (\text{unnormalised posterior})$$

notice looks like part of a gaussian, only not joint variable is  $\mu$   
 so must define probability dist. over  $\mu$

could:

$$P(\mu | \bar{Y}) = A e^{-\frac{N}{2\sigma^2} (\bar{Y} - \mu)^2} \rightarrow \int d\mu \rightarrow \text{Find } A$$

recognise gaussian

or:

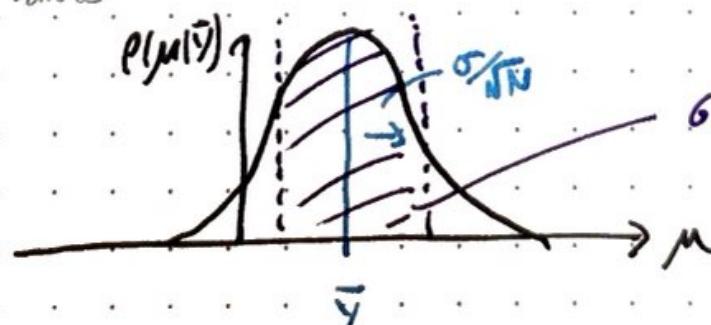
$$P(\mu | \bar{Y}) = \frac{1}{\sigma \sqrt{N} \sqrt{2\pi}} e^{-\frac{N}{2\sigma^2} \frac{(\bar{Y} - \mu)^2}{(1/\sigma^2 N)}}$$

gaussian on  $\mu$  tells us what degree of belief we have on various values of  $\mu$

68% prob.  $\mu$  lies between  $\bar{Y} \pm \sigma/\sqrt{N}$

$$= N(\bar{Y}, \sigma^2/N)$$

normalised posterior



68% posterior probability

aka  
"credible interval"

~~see slide: frequentist vs bayes for more detailed comparison~~

written over by:

## Lecture 12 (slides)

19.2.25

### Frequentist vs. Bayes

- \* Frequentists make statements about the data (or statistics or estimators = functions of the data) conditional on the parameter  $P(D|\theta)$  or  $P(f(D)|\theta)$
- \* Often goal is to get a "point estimate" or confidence intervals with good properties under repeated experiments  $\rightarrow \infty$  ("long run")
- \* Arguments are based on datasets that could have happened but didn't e.g. null hypothesis testing
- \* Bayesians make statements about the probability of parameters conditional on the dataset  $D = D_{obs}$  that you actually observed:  
 $P(\theta|D = D_{obs})$   
This requires an interpretation of the probability as quantifying a "degree of belief" in a hypothesis.
- \* Bayesian answer is the full posterior density  $P(\theta|D = D_{obs})$  quantifying the "state of knowledge" after seeing the data. Any numerical estimates are attempts to (imperfectly) summarise the posterior.

## Bayes Advantages

(slides)

- \* Ability to include prior information  $P(\theta)$ 
  - Incorporate info. from external datasets  
 $P(\theta)$  is the posterior from some other data  $P(\theta|\text{data})$ 
    - e.g. cosmology. diff datasets give diff ways to probe cosmological parameters  
( $\rightarrow$  encode external posterior given one analysis as prior to diff analysis). Natural way to incorporate information since one parameter often has multiple observables in cosmology
  - Regularisation: penalises overfitting with complex model. e.g. gaussian process prior
    - e.g. linear regression. ridge regression is linear regression w/ a penalty term; in Bayesian land we can incorporate such penalties. e.g. can say I don't think large slopes will occur and can incorporate this into prior. can also penalise complex models  $\rightarrow$  Occam's razor.
      - e.g. flat
      - e.g. wide gaussian vs. slightly more info than flat
  - "Noninformative"/weakly informative/definitive priors when you don't have much prior information
    - e.g. might only be able to guess param is between 0 and 1 but still information!  $\rightarrow$  cannot have prior without some information
- \* Likelihood (you're plugged in observed dataset and are viewing that as fn. of the parameters) is not a probability density in the parameters. But multiply by a prior (even flat) and the (normalised) posterior is a probability density, conditional/marginal probability can be computed.
  - (e.g. in multi-parameter case, useful to derive conditional densities later).
- \* Ability to deal with high dimensional parameter space. e.g. latent variables or nuisance parameters and marginalise them out (analytically/numerically)
  - (e.g. equivalent of MLE is called maximum a posteriori - can find place of posteriors and just use that as your point estimate)
- \* Note: Bayesian inference not necessarily completely opposed to frequentist statistics. Estimators derived from Bayesian arguments can still be evaluated in a frequentist basis. (e.g. James-Stein estimators)

think of Bayesian inference as just a way to generate estimators

or find credible intervals as intervals that contain some posterior prob 51

Last time: gaussian example of bayesian inference  
 (1D, flat prior  $\rightarrow$  like "hello world" of bayesian inference)  
 ↳ now slightly more involved example.

recap:

## Simple Gaussian Example w/ conjugate prior (board)

$$Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2) \quad (\mu \text{ unknown}, \sigma^2 \text{ known})$$

$$\begin{aligned} \text{Likelihood: } P(\vec{y} | \mu, \sigma^2) &= \prod_{i=1}^N N(y_i; \mu, \sigma^2) = \dots \\ &= (2\pi\sigma^2)^{-N/2} e^{-\frac{(N-1)s^2}{2\sigma^2}} e^{-\frac{N}{2\sigma^2}(\bar{y}-\mu)^2} \end{aligned}$$

Defined sufficient statistics:

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i \quad \text{"sample mean"}$$

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2 \quad \text{"sample variance"}$$

all info in data is contained in the likelihood & likelihood depends on sufficient statistics (summarizes info about data). e.g. if two diff. y vectors with same  $\bar{y}$  &  $s^2$ , world contains same info about parameters, even if y vector is different

$$\text{Conjugate Prior: } P(\mu) = N(\mu | \mu_0, \tau_0^2)$$

(equivalently can write  $\mu \sim N(\mu_0, \tau_0^2)$ )

$$P(\mu) = \frac{1}{\tau_0 \sqrt{2\pi}} e^{-\frac{1}{2\tau_0^2}(\mu - \mu_0)^2}$$

only goes  $\mu$   
 so disregard term  
 without  $\mu$

$$\text{Posterior: } P(\mu | \vec{y}) \propto P(\vec{y} | \mu) P(\mu)$$

$$\propto e^{-\frac{N}{2\sigma^2}(\bar{y} - \mu)^2} e^{-\frac{1}{2\tau_0^2}(\mu - \mu_0)^2}$$

recognise? [warm up ex1 - product of gaussian densities]

$$P(\mu | \vec{y}) = N(\mu | \mu_N, \tau_N^2)$$

↑  
 posterior  
 mean ↑  
 posterior  
 variance

think: realia  
 probability on  $\mu$   
 everything deterministic  
 w.p.  
 since is normalising const  
 i.e. when integrated over  
 N. we get const.  
 useful actually to see  
 it as prior w/ a  
 52 gaussian

(warm up sheet 1)

reminder: when you have product of gaussian densities,  
 the precision =  $\frac{1}{\text{variance}}$  of the resulting gaussian =  
 sum of precisions of individual gaussian.

$$\left[ \frac{1}{\tau_N^2} = \frac{1}{\tau_0^2} + \frac{N}{\sigma^2} \right] \quad \text{POSTERIOR PRECISION}$$

gaussian  
 $N(\mu, \sigma^2)$   
 $N(\mu_0, \sigma_0^2)$   
 add them:  
 $N(\mu + \mu_0, \sigma_0^2 + \sigma_0^2)$

Posterior precision = sum of precisions from likelihood  
 multiply them:  $N\left(\frac{\mu_0 \sigma_0^2 + \mu_1 \sigma_0^2}{\sigma_0^2 + \sigma_0^2}, \frac{1}{\sigma_0^2 + \sigma_0^2}\right) \rightarrow \text{prior}$

reminder:

$$\left[ M_N = \frac{\frac{1}{\tau_0^2} M_0 + \frac{N}{\sigma^2} \bar{y}}{\frac{1}{\tau_0^2} + \frac{N}{\sigma^2}} \right] \quad \text{POSTERIOR MEAN}$$

(divide by weights)  
 to get mean

Posterior mean = precision weighted average  
 of means from likelihood &  
 prior

(These formulas explicitly  
 tell us how our prior & likelihood  
 the bayesian inference

(mean & variance)

tends to

some as for  
 flat prior  
 likelihood

For fixed prior  $\tau_0$ , as  $N \rightarrow \text{large}$ ,  $M_N = \bar{y}$

$$(\text{plain}) = N(\mu_1, \frac{\frac{\sigma_0^2}{\tau_0^2} + \frac{N}{\sigma^2}}{\frac{\sigma_0^2}{\tau_0^2} + \frac{N}{\sigma^2}, \frac{\sigma_0^2 \tau_0^2}{\tau_0^2 + N \sigma^2}})$$

$$\tau_N^2 \rightarrow \sigma^2/N$$

(Data  $N \rightarrow \text{large}$ , likelihood dominates over prior)  
 basically never resulting likelihood pmf except w/ interpretation  
 of probability

Side note: What is conjugate prior?

for a given likelihood if you combine it  
 with a conjugate prior coming from 'nice' distribution  
 then the posterior is guaranteed to come  
 from the same class of 'nice' distributions  
 (e.g. gaussian, wishart, laplace)

(not sure how relevant/necessary this is)

## Astrophysics Example

(slides)

### BAYESIAN INFERENCE FOR PARALLAX

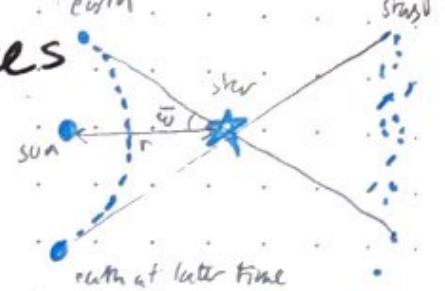
PAPER:  
 (Beamer Jones)  
 Estimating distances  
 from parallaxes

distant background stars

- parallax is a way to measure distances

$$\omega_{\text{true}} = \frac{\text{parsec}}{\text{arcsec}} \quad (\text{applying right})$$

$$\omega_{\text{true}} = \frac{1}{r} \quad (\text{in useful units})$$



- Gaia satellite makes parallax measurement to measure distances & map stars in our galaxy
- measurement error in parallax angle measured parallax is different to true parallax

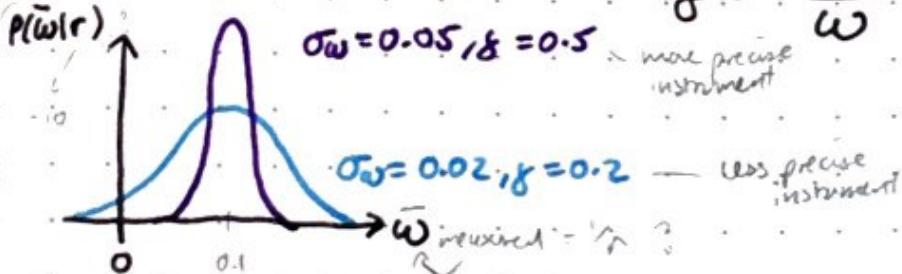
$$P(\bar{\omega}|r) = \frac{1}{\sigma_{\bar{\omega}} \sqrt{2\pi}} \exp \left[ -\frac{1}{2} \frac{(\bar{\omega} - \frac{1}{r})^2}{\sigma_{\bar{\omega}}^2} \right] \quad \sigma_{\bar{\omega}} > 0$$

measured parallax has some distribution;  
 assume gaussian distribution for simplicity

- fractional measurement error  $f = \frac{\sigma_{\bar{\omega}}}{\bar{\omega}}$

$$r = 10 \text{ pc}$$

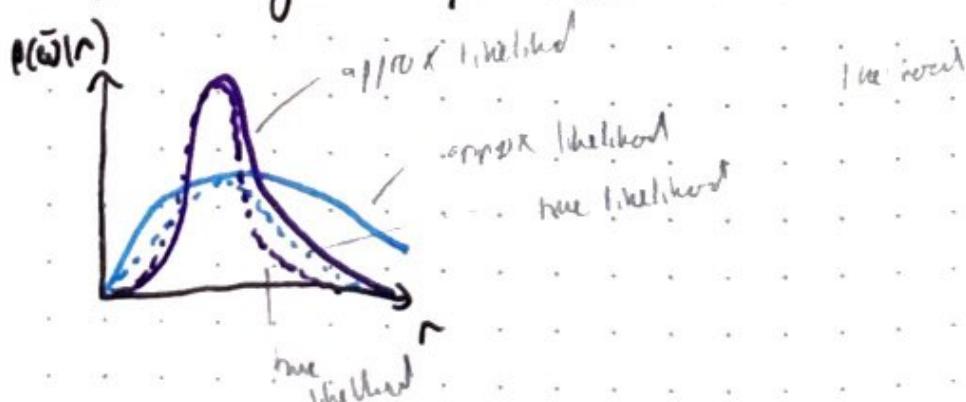
star at 10pc  
 actual measured parallax  
 would be tiny value given  
 your instrument's  
 precision



- possible to get negative parallax due to measurement error

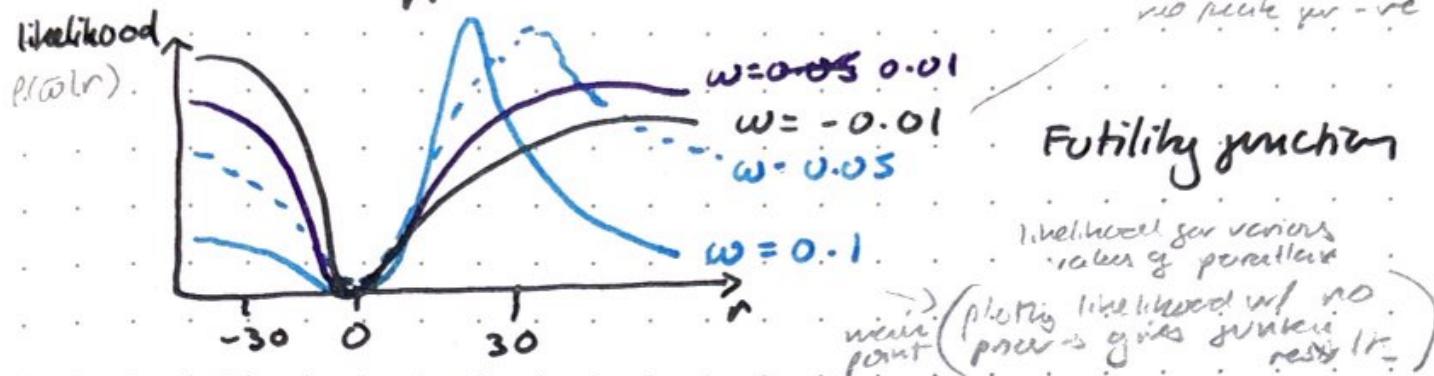
Likelihood

get skewed distributions



plot above? posterior vs prior in magnitude?  
 - Laplace approximation gives gaussian, incidentally this is mathematically equivalent to doing propagation of error (assume gaussian measurement error propagate to  $r_n \dots$ )

↳ this is an approximation (dotted vs solid line)?



### Futility function

likelihood for various values of parameter  
 $\rightarrow$  plotting likelihood w/ no point (prior is given symmetrically)

- likelihood is positive on negative values of distance (unphysical) for all values of measured parallax

- negative measurements have no mode (MLE)

like peak  $\rightarrow$  most likely distance peak for the but not for  $w = 0$

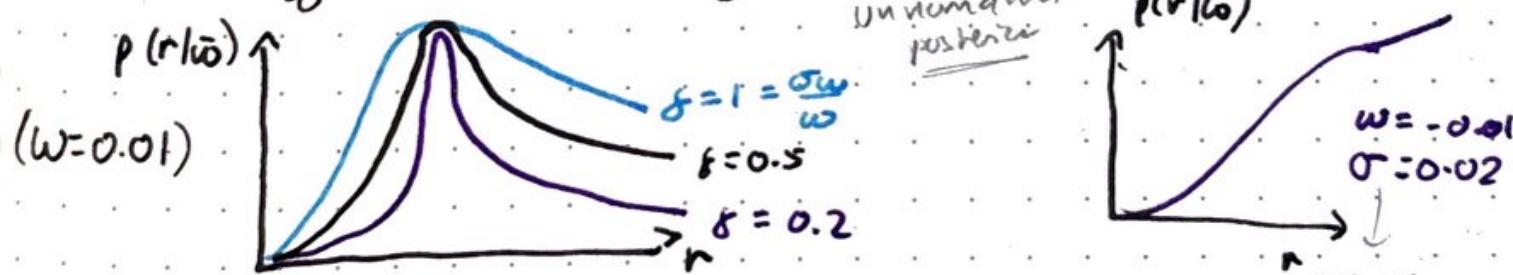
try to fix this:

- impose our knowledge that we know distance should be positive

$$P(r) \propto P(r|w) \propto P(w|r) P(r)$$

$$P(r) = \begin{cases} 1 & r > 0 \\ 0 & \text{otherwise} \end{cases} \rightarrow \begin{array}{l} \text{note:} \\ (\text{improper distribution} \\ = \text{not normalizable}) \end{array}$$

- this cuts off -ve values of  $r$  in our graph



- improper posterior: not-normalizable  
 no mean, variance etc.

- mode ( $r = w$ ) exists for the  $w$  but undefined for -ve  $w$  ( $r = \infty$ )

still have some problems

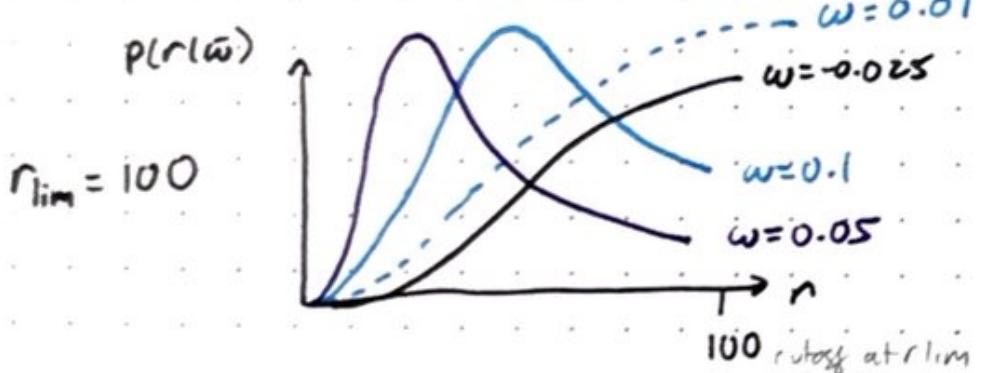
- impose a proper distance prior (impose limit to max. distance of a star in our survey)

$$P(r) = \begin{cases} \frac{1}{r_{\text{lim}}} & 0 < r \leq r_{\text{lim}} \\ 0 & \text{otherwise} \end{cases}$$

$$P(r|w) \propto P(w|r)P(r)$$

$$P_u^+(r|\bar{w}, \sigma_{\bar{w}}) = \begin{cases} \frac{1}{r_{\text{lim}}} P(\bar{w}|r, \sigma_{\bar{w}}) & 0 < r \leq r_{\text{lim}} \\ 0 & \text{otherwise} \end{cases}$$

numerical posterior



- implies unrealistic stellar density  $\sim \frac{1}{r^2}$
- cutoff of mode at edge  $r_{\text{lim}}$  for small or negative measured parallax ( $w < 0$ ) (but better than going to infinity)

$$\text{rest} = \begin{cases} \frac{1}{\bar{w}} & 0 < \frac{1}{\bar{w}} \leq r_{\text{lim}} \\ r_{\text{lim}} & \frac{1}{\bar{w}} > r_{\text{lim}} \\ r_{\text{lim}} & \bar{w} \leq 0 \end{cases} \quad (\text{extreme mode})$$

- still have some problems  $\rightarrow$  can go towards more & more astrophysically motivated priors

- problem w/ previous prior: volume density  $p(r)$   
 $P_r(\text{distance } \in [r, r+dr]) \propto p(r) 4\pi r^2 dr$

unphysical  $\rightarrow$

$$p(r) \propto \frac{1}{r^2}, r < r_{\text{lim}} \Rightarrow p(r) \propto \text{const.} \quad r < r_{\text{lim}} \\ 0 \quad \text{otherwise} \quad 0 \quad \text{otherwise}$$

- more physical: <sup>consider</sup> copernican principle  $p(r) \propto \text{const.}$ ,  $p(r) \propto \text{const.} \times 4\pi r^2$

Lecture 1321.2.25

- astrophysics example
- recap: understanding how priors impact inference
- keeps going as before: introduce more & more astrophysically motivated priors (see slides)

Conclude: Priors in Bayesian Inference

- \* Priors can be used to encode background info./ external knowledge about parameters
  - weak mathematical constraints on physical parameters e.g. positivity of distances
  - astrophysical info. e.g. distributions of stars
- \* test sensitivity of your inference to the priors (& likelihood) under various assumptions of the model

$$P(\bar{w}|r) \xrightarrow[\text{likelihood}]{\text{measurement}} P(r|w) \propto P(w|r)P(r) \xrightarrow[\text{posterior}]{\text{prior}}$$

Note: also assumed some measurement error e.g. Gaussian  
not just prior

# MULTIPARAMETER BAYESIAN MODELS (slides)

Example: gelman BDA 3.2 sheet 2

[analytic posterior for gaussian  $(\mu, \sigma^2)$  model with non-informative prior]

data generating process:  $y_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \quad i=1, \dots, n$

likelihood:  $p(y|\mu, \sigma^2) \propto (\sigma^2)^{-n/2} \exp\left(-\frac{(n-1)s^2}{2\sigma^2}\right) \exp\left(-\frac{n}{2\sigma^2}(\bar{y}-\mu)^2\right)$

(non informative)  $p(\mu) \propto 1$

improper prior:  $P(\log \sigma^2) \propto 1$  or  $P(\sigma^2) \propto \sigma^{-2}$  ( $\sigma^2 > 0$ )

even if prior improper, need proper posterior (corresponds to real observations)

joint posterior:  $P(\mu, \sigma^2 | y) \propto P(y | \mu, \sigma^2) \times P(\mu, \sigma^2)$

$P(\mu, \sigma^2 | \bar{y}) \propto (\sigma^2)^{-n/2-1} \exp\left(-\frac{(n-1)s^2}{2\sigma^2}\right) \exp\left(-\frac{n}{2\sigma^2}(\bar{y}-\mu)^2\right)$

joint is useful to find conditionals/marginals

conditional posterior:  $P(\mu | \sigma^2, y) = N(\mu | \bar{y}, \sigma^2/n)$

marginal posteriors:  $P(\sigma^2 | y) = \int P(\mu, \sigma^2 | y) d\mu = \text{Inv-}\chi^2(\sigma^2 | n-1, s^2)$

$$P(\mu | y) = \int P(\mu, \sigma^2 | y) d\sigma^2$$

$$\propto \left[ 1 + \frac{n(\mu - \bar{y})^2}{(n-1)s^2} \right]^{-n/2}$$

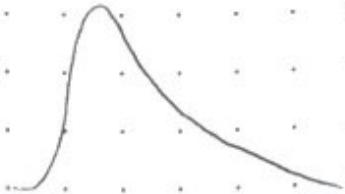
$$= t_{n-1}(\mu | \bar{y}, s^2/n)$$

t distribution

Note: ① inverse chi-squared distribution

$$\text{Inv-}\chi^2(\theta | v, s^2) \propto \theta^{-(v/2+1)} \exp\left(-\frac{vs^2}{2\theta}\right), \theta > 0$$

deg. of freedom  $v > 0$   
scale parameter  $s > 0$



② same as inverse-gamma

$$\text{Inv-gamma}(\theta|\alpha, \beta) \propto \theta^{-(\alpha+1)} \exp(-\beta/\theta), \theta > 0$$

shape parameter  $\alpha = \frac{\nu_0}{2}$   
scale parameter  $\beta = \frac{\nu_0}{2} \sigma_0^2$

useful normal distributions have random number generators in common python packages

③ student's t distribution

$$t_{v}(\theta|\mu, \sigma^2) \propto [1 + \frac{1}{v} (\frac{\theta - \mu}{\sigma})^2]^{-(v+1)/2}$$

deg. of freedom  $v > 0$   
location parameter  $\mu$   
scale parameter  $\sigma > 0$



details on distributions in book (gelman BDA) > appendix  
'your example sheet' "or just wikipedia"

What if you can't compute marginals/expectations analytically?  $\rightarrow$  Bayesian Computation

- \* Bayesian answer: full posterior  $p(\theta|D)$ , numerical estimates are attempts to (imperfectly) summarize the posterior e.g. mean, mode
- \* Often these are posterior expectations  $E[\delta(\theta|D)]$   
(computationally difficult)  $= \int \delta(\theta) p(\theta|D) d\theta$
- \* Bayesian computation: algorithms to "map out" / sample the posterior  $p(\theta|D)$  and compute expectations  $E[\delta(\theta|D)]$
- \* e.g. MCMC, nested sampling, importance sampling  
all models are wrong, some are useful. (even if your computation is 'right'!)

object of computation is to calculate some integral (posterior expectation)  $\rightarrow$  can do with monte carlo integration  
MONTE CARLO INTEGRATION (slides)

Typically we want to summarize the posterior & compute expectations of the form

$$I = \mathbb{E}[f(\theta)|D] = \int f(\theta) p(\theta|D) d\theta$$

Using  $m$  samples from the posterior

$$\theta_i \sim p(\theta|D)$$

$$\hat{I} = \frac{1}{m} \sum_{i=1}^m f(\theta_i) \rightarrow I \quad (\text{CLN for large } N)$$

/CLT

Monte Carlo Error (derive later on next page)

$$\text{Var}[\hat{I}] = \frac{1}{m^2} \sum_{i=1}^m \text{Var}[f(\theta)] = \frac{1}{m} \text{Var}[f(\theta)] \approx \frac{1}{m} \widehat{\text{Var}}[\{f(\theta_i)\}]$$

Monte carlo error  $\sim 1/m^{1/2}$  indep of dim  $[\theta]$   $\rightarrow$  convenient

Fundamental Theorem of Monte Carlo  $\rightarrow$  Bayesian computation using sampling

$$\mathbb{E}[f(\theta)|D] = \underbrace{\int f(\theta) p(\theta|D) d\theta}_{\text{posterior expectation}} \approx \frac{1}{m} \sum_{i=1}^m f(\theta_i) \underbrace{\approx}_{\text{sample average}}$$

E.g. Posterior mean  $\mu$   $f(\theta) = \theta$   
 Posterior variance  $f(\theta) = (\theta - \mu)^2$

Probability in an interval  $[a, b]$   $f(\theta) = I_{[a,b]}(\theta)$

$\uparrow$   
 (indicator function)  
 = one where variable between  
 $a$  &  $b$ .

# Monte Carlo Error

(board)

suppose:  $f$  is some invertible fn. of  $\theta$ .

$$\theta_i \stackrel{iid}{\sim} P(\theta) \quad i=1, \dots, m$$

monte carlo sample size (not same as data sample size)

→  $f(\theta_i)$  are iid

$$\text{Def: } \hat{I} = \frac{1}{m} \sum_{i=1}^m f(\theta_i)$$

$$\text{Var}(\hat{I}) = \text{Cov}(\hat{I}, \hat{I}) = \text{Cov}\left[\frac{1}{m} \sum_{i=1}^m f(\theta_i), \frac{1}{m} \sum_{j=1}^m f(\theta_j)\right]$$

$$= \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \text{Cov}[f(\theta_i), f(\theta_j)]$$

$\theta_i, \theta_j$  iid means ⇒ only non-zero contributions when  $j=i$

$$= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(f(\theta_i))$$

$$= \frac{1}{m} \text{Var}(f(\theta))$$

standard deviation of monte carlo error  
std MC error

$$\sqrt{\frac{1}{m} \text{Var}(f(\theta))} \times \frac{1}{\sqrt{m}}$$

$$\text{GLT: } m \rightarrow \text{large} \quad \hat{I} \stackrel{d}{\sim} N(\mathbb{E}[f(\theta)], \text{Var}(f(\theta))/m)$$

approximate variance

$$\text{Var}(\hat{I}) \approx \frac{1}{m} \widehat{\text{Var}}(\{f(\theta_i)\})$$



MC sample variance

$$\widehat{\text{Var}}(f(\theta_i)) = \frac{1}{m-1} \sum_{i=1}^m (f(\theta_i) - \hat{I})^2$$

slow convergence  
(increasing of MC error)  
but! nice: indep. of  
dim( $\theta$ )

so MC error rate  
is indep. of dimension  
useful!

- MC error scales as  $m^{-1/2}$
- Slow but indep. of  $\text{dim}(\theta)$ !
- can apprx.

## Lecture 14

notes on useful properties of  
gaussian on moodle:  
→ ex sheet 2

24.2.25

**Today:** Bayesian Computation, Direct sampling, Importance Sampling

### Monte carlo direct sampling (slides)

how do we sample from joint? → it factorises

example given before  
is gaussian  $(\mu, \sigma^2)$  model with  
non-informative prior

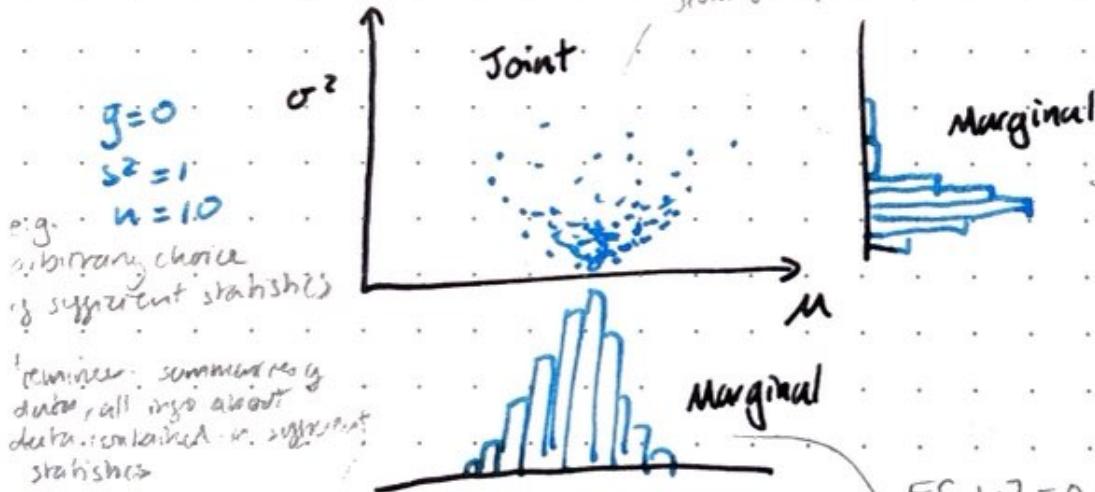
$\rightarrow P(\mu|\sigma^2)$

$$\text{Factorise posterior: } P(\mu, \sigma^2 | y) = P(\mu | \sigma^2, y) P(\sigma^2 | y)$$

$$\begin{aligned} 1. \sigma^2 &\sim P(\sigma^2 | y) & \text{Inv-}\chi^2 &\rightarrow (\mu_i, \sigma_i^2) \sim P(\mu, \sigma^2 | y) \\ 2. \mu_i | \sigma_i^2 &\sim P(\mu | \sigma_i^2, y) & \text{Normal} & \end{aligned}$$

drawn from posterior

Joint draw from posterior!



Compute posterior  
summarises from  
monte carlo!

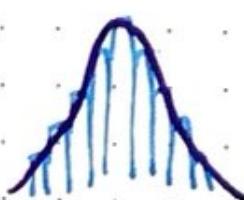
$$\begin{aligned} E[\sigma^2 | y] &= 1.10 \\ \text{std}[\sigma^2 | y] &= 0.25 \end{aligned}$$

$$\begin{aligned} E[\mu | y] &= 0.00 \\ \text{std}[\mu | y] &= 0.14 \end{aligned}$$

from  
posterior

(did not need to solve analytically to get marginals)

Kernel density estimate = estimate a smooth density from samples



can think  
of as a  
smooth  
histogram

## KERNEL DENSITY ESTIMATION

(lociord)

$$\theta_i \stackrel{iid}{\sim} P(\theta|D) \quad i=1, \dots, m$$

Approx  $P(\theta|D)$  using  $\hat{P}(\theta|D) \approx \frac{1}{m} \sum_{i=1}^m N(\theta|\theta_i; bw^2)$

$$\int \hat{P}(\theta|D) d\theta = 1 \quad \checkmark$$

nice  
normal-like  
gaussian  
peak

$\checkmark$   
select bandwidth

Lots of ways to choose bandwidth, nice one is

### Silverman's Rule of Thumb:

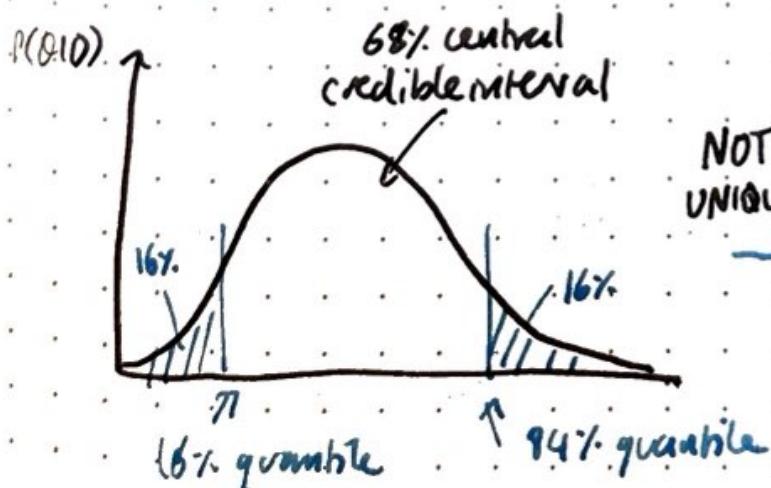
$$bw = \left( \frac{4\hat{\sigma}^5}{3m} \right)^{1/5} \quad \hat{\sigma}^2 = \text{Var}(\theta|D) \approx \widehat{\text{Var}}(\varepsilon\theta_i)$$

sample variance  
 $\frac{1}{m-1} \sum_{i=1}^m (\theta_i - \bar{\theta})^2$  ?

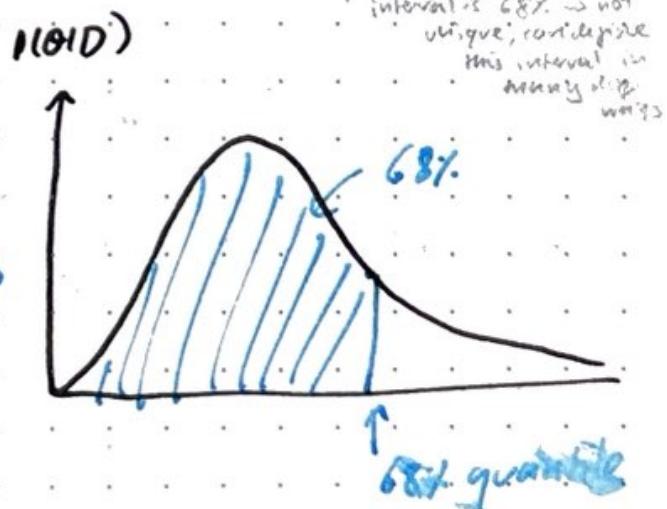
## SUMMARISING POSTERIOR UNCERTAINTIES

- \* Can compute posterior mean (or median) ± posterior standard deviation.
- \* If posterior  $P(\theta|D)$  is Gaussian, then this contains 68% posterior probability.
- \* But not necessarily for non-gaussian posteriors!

## CREDIBLE INTERVALS



NOT  
UNIQUE!



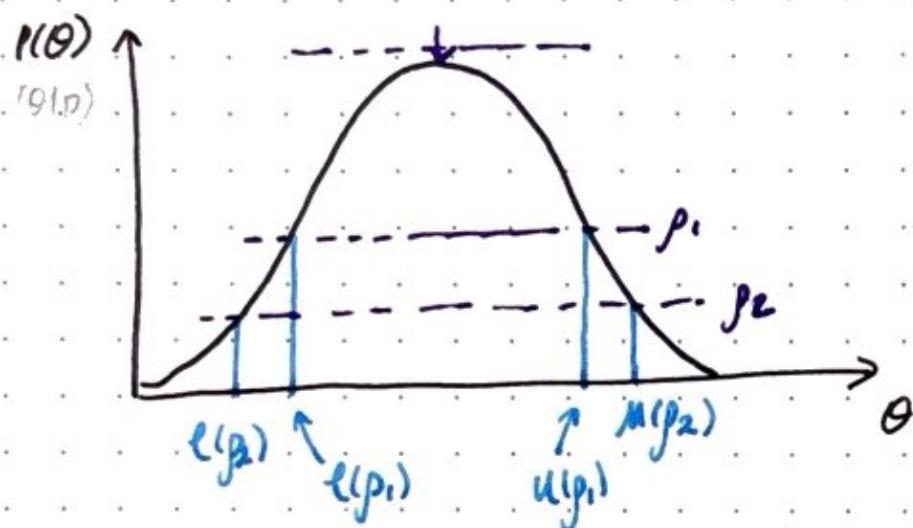
"In Bayesian terms  
there about credible intervals"

prob  $\theta$  lies in given  
interval is 68% → not  
unique; consider  
this interval in  
many ways

(HPD) (board)

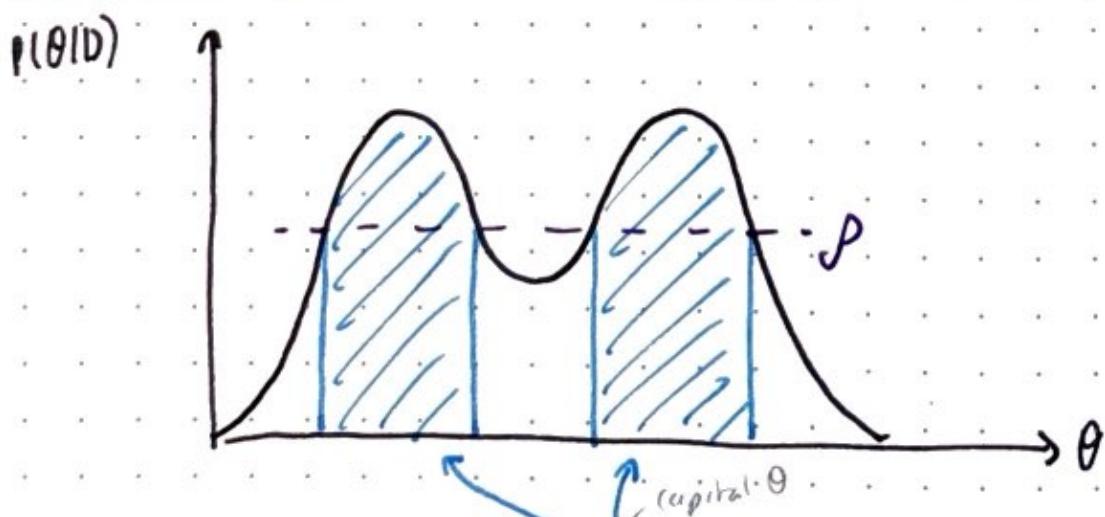
## HIGHEST POSTERIOR DENSITY CREDIBLE INTERVAL

- Unique
- Narrowest interval that contains  $x\%$  posterior probability



Find  $p$  s.t.  $[l(p), u(p)]$  contains  $x\%$  posterior probability.

HPD works for multi-modal posteriors!



$x\%$  HPD-C.I. = set  $\Theta$  (possibly disconnected)  
with highest  $p$  s.t.

$$\forall \theta \in \Theta, p(\theta|D) \geq p \text{ and } \int_{\Theta} p(\theta|D) = x$$

[What if you can't directly sample the posterior  
 $\theta_i \sim P(\theta|D)$ ?] (board)

$$\mathbb{E}[f(\theta)|D] = \int f(\theta) P(\theta|D) d\theta \approx \frac{1}{m} \sum_{i=1}^m f(\theta_i)$$

e.g. likelihood might be output of some computer program

### \* Importance sampling:

draw from an easier "tractable" distribution (importance function)  $\theta_i \sim Q(\theta)$  and weight the samples by  $w_i = P(\theta_i|D)/Q(\theta_i)$  to compute expectations

### \* Posterior simulation:

MCMC, Nested Sampling etc. generates draws from the posterior density iteratively in long-run

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

## IMPORTANCE SAMPLING (board)

Given a probability density  $P(\theta)$  [e.g. the posterior]

Want to estimate

$$I = \mathbb{E}_P[f(\theta)] = \int f(\theta) P(\theta) d\theta$$

However,  $P(\theta)$  is intractable, i.e. difficult to directly sample, e.g.  $\theta_i \sim P(\theta)$ , but can evaluate  $P(\theta_i)$  for a parameter value  $\theta_i$ .

Select an instrumental distribution (importance function) from which it is easy to draw samples:

$$\theta_i \sim Q(\theta) \rightarrow \begin{array}{l} \text{tractable} \\ \text{readily draw/sample} \\ \text{evaluate} \end{array}$$

and can evaluate  $Q(\theta_i)$ , and  $Q(\theta) > 0$  whenever  $P(\theta) > 0$ .

$$\begin{aligned} \text{Rewrite } I &= \mathbb{E}_P[f(\theta)] = \int f(\theta) \frac{P(\theta)}{Q(\theta)} Q(\theta) d\theta \\ &= \mathbb{E}_Q \left[ f(\theta) \frac{P(\theta)}{Q(\theta)} \right] \end{aligned}$$

Draw samples  $\theta_i \stackrel{\text{iid}}{\sim} Q(\theta)$  ( $i=1, \dots, m$ )

Approximate  $I$  with importance sampling estimate

$$\hat{I}^* = \frac{1}{m} \sum_{i=1}^m f(\theta_i) \frac{P(\theta_i)}{Q(\theta_i)} = \frac{1}{m} \sum_{i=1}^m f(\theta_i) w_i^*$$

where  $w_i^* = w^*(\theta_i) = \text{importance weights} = \frac{P(\theta_i)}{Q(\theta_i)}$

$$\begin{aligned}
 E_Q[\hat{I}^*] &= \frac{1}{m} \sum_{i=1}^m E_Q[f(\theta_i) w^*(\theta_i)] \\
 &= \frac{1}{m} \sum_{i=1}^m E_Q[f(\theta_i) \frac{P(\theta_i)}{Q(\theta_i)}] \\
 &= \frac{1}{m} \sum_{i=1}^m \int f(\theta) \frac{P(\theta)}{Q(\theta)} Q(\theta) d\theta \\
 &= \frac{1}{m} \sum_{i=1}^m \int f(\theta) P(\theta) d\theta \\
 &= \frac{1}{m} \sum_{i=1}^m E_P[f(\theta)] = E_P[f(\theta)] \\
 &= I \quad (\text{unbiased})
 \end{aligned}$$

$$\text{Var}_Q[\hat{I}^*] = \frac{1}{m^2} \sum_{i=1}^m \text{Var}(f(\theta) w^*(\theta))$$

$$= \frac{1}{m} \text{Var}(f(\theta) w^*(\theta))$$

$$\approx \frac{1}{m} \widehat{\text{Var}}(\{f(\theta_i) w^*(\theta_i)\})$$

$$\widehat{\text{Var}}(\{\cdot\}) = \frac{1}{m-1} \sum_{i=1}^m (f(\theta_i) w^*(\theta_i) - \hat{I}^*)^2$$

reminder:

$$\begin{aligned}
 &\text{Var}(x+bx) \\
 &= a^2 \text{Var}(x) + b^2 \text{Var}(y)
 \end{aligned}$$

## Lecture 15

useful properties of  
gaussian's noted on moodle  
see sheet 2

26.2.25

Today: case study "Importance sampling the mass of the milky way galaxy"  
Continuing: Bayesian computation, importance sampling, MCMC (Gelman, BDA ch10-12)

Recap: MC integration, want to summande posterior, estimate w/ samples from the posterior. (replace integral w/ MC sum) which converges to desired expectation for large  $N$ . Error on estimate  $\propto 1/N$  and is indep of dim. of parameter space!

### Importance sampling: recap (slides)

Objective: compute expectation w.r.t. distribution  $p(\theta)$

$$I = \mathbb{E}[f(\theta)] = \int f(\theta) p(\theta) d\theta$$

Example: Posterior (suppress "1D" here)

$$P(\theta) = \frac{L(\theta) \pi(\theta)}{\int L(\theta) \pi(\theta) d\theta}$$

← Likelihood      ← prior  
directly

Can evaluate  $p(\theta)$  but not sample from it. choose importance function  $Q(\theta)$  you can evaluate and sample!

Importance sampling estimate  $\theta_1, \dots, \theta_n \stackrel{iid}{\sim} Q(\theta)$

$$\hat{I}^* = \sum_{i=1}^n f(\theta_i) w^*(\theta_i)$$

$$w_i^* = w^*(\theta_i) = p(\theta_i) / Q(\theta_i)$$

## SELF-NORMALISED IMPORTANCE SAMPLING

(slides)

Objective: compute expectation

$$I = \mathbb{E}[f(\theta)] = \int f(\theta) P(\theta) d\theta$$

but can only evaluate unnormalised  
 $\tilde{P}(\theta)$ .

Example: unnormalised posterior

$$\tilde{P}(\theta) = L(\theta)\pi(\theta).$$

(evidence hard to evaluate)

$$\text{Normalised posterior is } P(\theta) = \frac{\tilde{P}(\theta)}{Z_p} = \frac{\tilde{P}(\theta)}{\int \tilde{P}(\theta) d\theta}$$

but, cannot easily calculate

$$Z_p = \int L(\theta)\pi(\theta) d\theta \text{ so cannot evaluate } P(\theta).$$

## Self-Normalised Importance Weighting

(board)

Often in Bayesian Analysis, you can only computer evaluate posterior  $P(\theta)$  up to a constant.

$$P(\theta) = \frac{L(\theta) \pi(\theta)}{Z_p} = \frac{\tilde{P}(\theta)}{Z_p}, Z_p = \int \tilde{P}(\theta) d\theta$$

↑ unnormalized  
normalized

**GOTL:** Estimate  $I = E_p[f(\theta)] = \int f(\theta) P(\theta) d\theta$

$$= \int f(\theta) \frac{\tilde{P}(\theta)}{\int \tilde{P}(\theta) d\theta} d\theta$$

$$= \frac{\int f(\theta) \frac{\tilde{P}(\theta)}{\int \tilde{P}(\theta) d\theta} Q(\theta) d\theta}{\int \frac{\tilde{P}(\theta)}{Q(\theta)} Q(\theta) d\theta}$$

instrumental  
distribution

preferably nice  
target distribution

"if you can't  
solve a hard problem  
solve an easier  
problem"

Choose  $Q$  s.t. drawing samples  $\theta_i \stackrel{iid}{\sim} Q(\theta)$  is easy ( $i=1, \dots, m$ ).

Approximate  $I$  with  $\hat{I}$  (FTMC)

$$I \approx \hat{I} = \frac{\sum_{i=1}^m f(\theta_i) \tilde{w}(\theta_i)}{\sum_{i=1}^m \tilde{w}(\theta_i)} = \sum_{i=1}^m f(\theta_i) w(\theta_i)$$

(Estimate normalization)  $\rightarrow \sum_{i=1}^m \tilde{w}(\theta_i)$   
using same sample from  $Q$

don't need full  
"in b/c top  
& bottom curve?"

where the self-normalised weights are,

$$w(\theta_i) = \frac{\tilde{w}(\theta_i)}{\sum_{i=1}^m \tilde{w}(\theta_i)} = \frac{\tilde{P}(\theta_i)/Q(\theta_i)}{\sum_{i=1}^m \tilde{P}(\theta_i)/Q(\theta_i)}$$

will work when  $P(\cdot)$  unnormalised  
we have done the bottom integral using importance sampling itself.

(Also works if  $Q$  is unnormalised, since const. factor in  $Q$  cancel above)

**COST:**  $\hat{I}$  is slightly biased

(whereas importance sampling is ~~biased~~ at. before unbiased)

$$\mathbb{E}[\hat{I}] = I + O(1/m)$$

b/c of denominator (work out yourself)

biased but consistent?  $\rightarrow$  as  $M \rightarrow \infty$   $P(\hat{I} = I) \rightarrow 1$

Bias arises from self-normalisation process

reminder:  $\mathbb{E}[AB] = \mathbb{E}[A]\mathbb{E}[B] + \text{cov}(A, B)$   
since  $\text{cov}(A, B) = \mathbb{E}[AB] - \mathbb{E}[A]\mathbb{E}[B]$

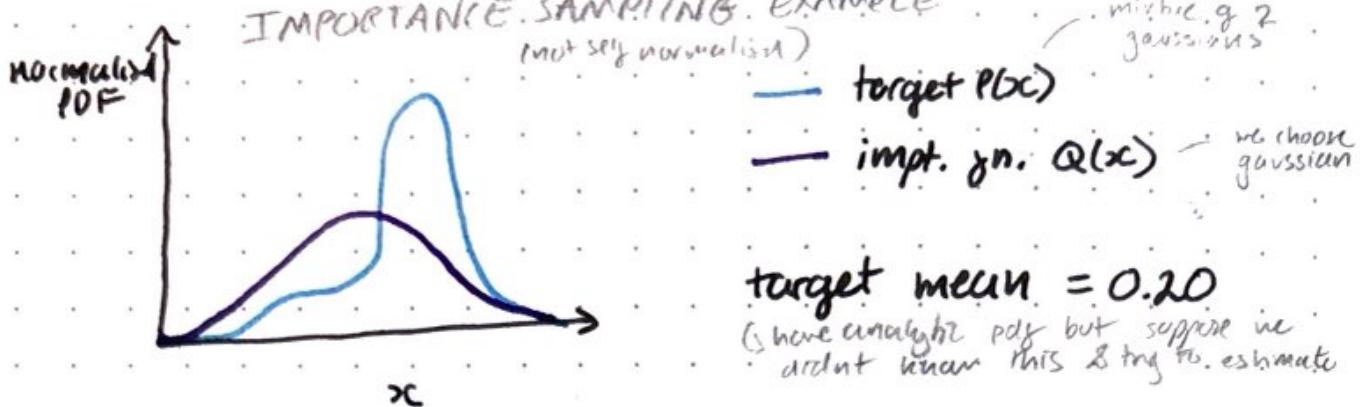
$$\mathbb{E}_Q[\hat{I}] = \sum_{i=1}^m \mathbb{E}_Q[\delta_i(\theta) w_i(\theta)] = \sum_{i=1}^m (\mathbb{E}_Q[\delta_i(\theta)] \mathbb{E}_Q[w_i(\theta)] + \text{cov}[\delta_i(\theta), w_i(\theta)])$$

unclear steps:  $\left\{ \begin{array}{l} \mathbb{E}_Q[w_i(\theta)] = 1/m + O(1/m^2) \\ \text{because } \mathbb{E}w_i = 1? \end{array} \right.$  (and  $\text{cov}[\delta_i(\theta), w_i(\theta)]$  ignorable?)  
 $\rightarrow \mathbb{E}_Q[\hat{I}] = \sum_{i=1}^m (\mathbb{E}_Q[\delta_i(\theta)] (1/m + O(1/m^2)))$  or  $\dots$ ?  
 $= \mathbb{E}_Q[\delta_i(\theta)] + O(1/m)$

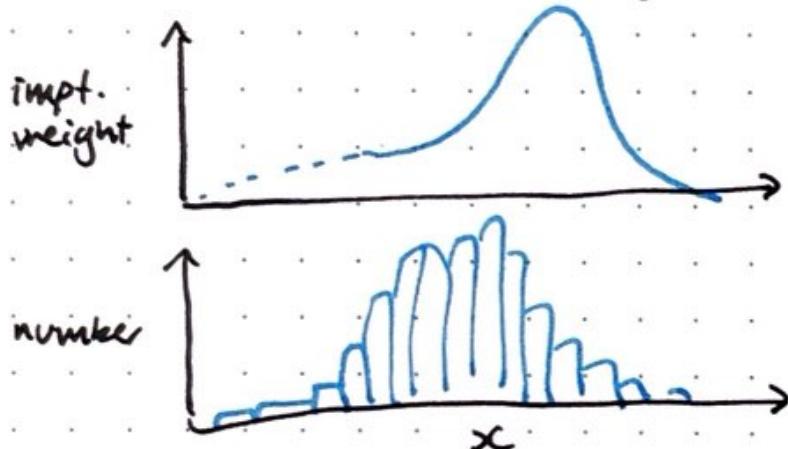
some kind of perturbation analysis/  
Taylor expansion?

$$\left( \text{Var}[\hat{I}] = \sum_{i=1}^m \text{Var}[\delta_i(\theta) w_i(\theta)] = m \text{Var}[\delta_i(\theta) w_i(\theta)] \right) ? \text{ like} \\ \approx m \widehat{\text{Var}}[\delta_i(\theta) w_i(\theta)]? \right)$$

## Contrived example: Gaussian Mixture, normalised PDF



Generate 1000 samples from impt. gn.



MC estimate of mean  
 $\sum_i w_i^* = 0.215 \pm 0.028$   
 $m$  within 10% target

$m=1000$   
draws from impt. gn.

could also do this for posterior variance etc.

$m=10,000$  draws: mc estimate of mean is

$$\sum_i w_i^* = 0.199 \pm 0.009$$

error goes down as  $1/\sqrt{m}$

## Parallax Example

(slides)

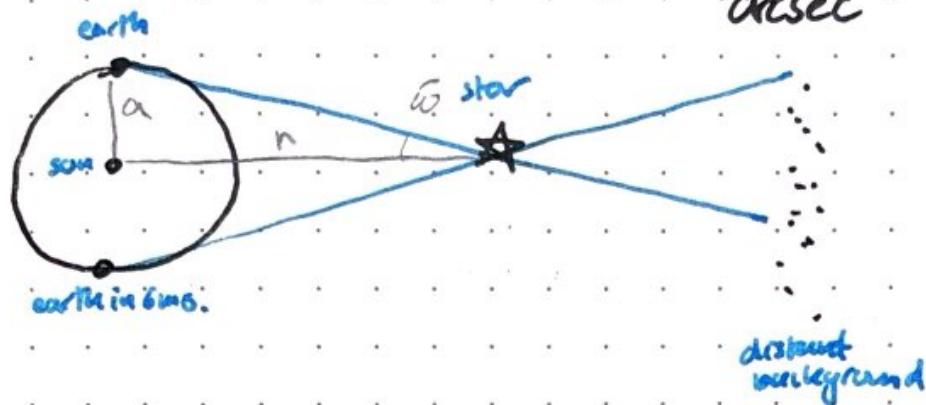
Likelihood  $w \sim N(\frac{1}{r}, \sigma_w^2)$

(gaussian measurement)  
error in  $w$ ?

$$P(\bar{\omega} | r) = \frac{1}{\sigma_{\bar{\omega}} \sqrt{2\pi}} \exp \left[ -\frac{1}{2\sigma_{\bar{\omega}}^2} (\bar{\omega} - \frac{1}{r})^2 \right], \quad \sigma_{\bar{\omega}} > 0$$

parallax  
recap

True relation (no errors):  $\frac{w}{\text{arcsec}} = \frac{\text{parsec}}{r}$



$w$  ← parallax angle

arcsec ←

parsec ←

$r$  ← distance

since  $\bar{\omega} \ll 1$

$\bar{\omega} = a/r$  to a good approx

when  $\bar{\omega}$  is 1 arcsecond,  
 $r$  is defined as the  
parsec

Introduce physically motivated prior

$$P(r) = \begin{cases} \frac{1}{2L^3} r^2 e^{-r/L} & \text{if } r > 0 \\ 0 & \text{otherwise} \end{cases}$$

as seen before  
some kind of

Exponential decrease in density of stars with galactic length scale  $L$

$$P(r|w) \propto P(w|r) P(r)$$

Unnormalised posterior

$$P_{\text{unnorm}}^*(r|\bar{\omega}, \sigma_{\bar{\omega}}) = \begin{cases} \frac{r^2 e^{-r/L}}{\sigma_{\bar{\omega}}} \exp \left[ -\frac{1}{2\sigma_{\bar{\omega}}^2} (\bar{\omega} - \frac{1}{r})^2 \right] & \text{if } r > 0 \\ 0 & \text{otherwise} \end{cases}$$

can't compute normalisation analytically

(in reality this is a 1D example so

we could actually evaluate this on a grid

and do integral numerically but

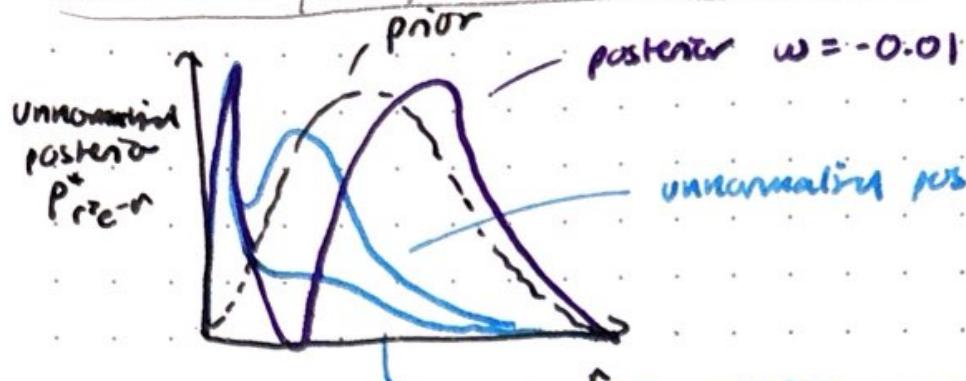
let's pretend we can't do that

want to compute posterior mean: consider posteriors of distance

blue?

posterior for  $w = 0.01$

with more values & error



unnormalized posterior  $g = \omega/\sigma = 0.31$

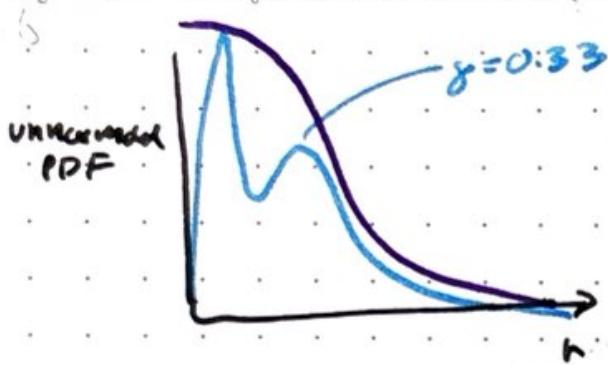
short peak at low r:  
ratio  $\hat{g}$ .

bimodal for high noise signal.

inhibiting  $\hat{g}$ ?

[e.g. try to get mean for  $w = 0.01$ , with error  $\hat{g} = 0.33$ ? need to use self-normalized importance sampling  $\Rightarrow$  need to choose imp. gn.]

bigger the var the  
more it looks like  
prior



$\hat{g} = 0.33$

$w = 0.01$

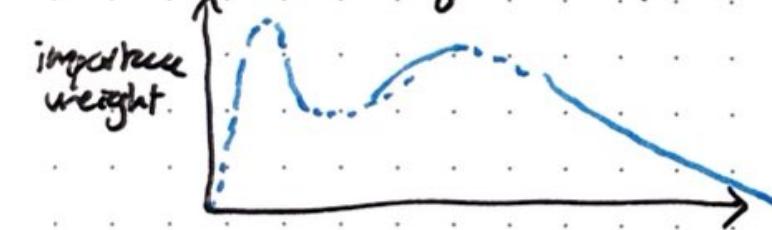
— target posterior  $p(r|w)$

— imp. gn.  $Q(r)$

choose half student t dist  
with  $V=3$  (low deg of freedom)  
choose because tail of student t  
is thicker than gaussian (gaussian  
drops off as  $e^{-x^2}$  but student t c  
want  $Q(r)$  to have fatter tail  
than the target, also choose b/c we  
want a named distribution

choose  $Q(r)$  that covers  $p(r)$  (target)  
↳ should be true everywhere the  
target is +ve  
For this case choose half student t  
with low deg. of freedom  $V=3$

$M = 10^4$  draws from  $Q(r)$ :

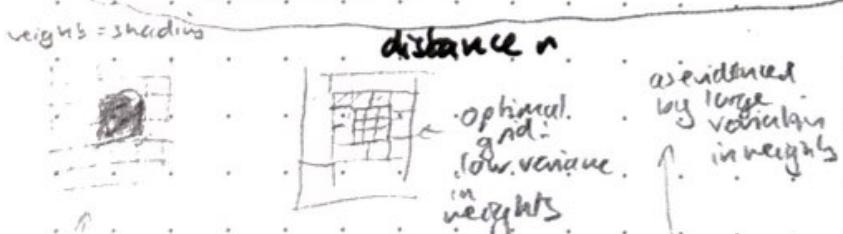


MC estimate of mean

$$\sum x_i w_i = 2391.263 \pm 3.918$$



$M = 10000$  draws from  
imp. gn.



unif.  
high var  
in weights  
of each square

idea of ESS,  
can get some info  
from uniform grid as:

non-uniform grid w/  
over samples but with  
more weights. so we get more  
effective info from each grid  
square  $\rightarrow$  i.e. the uniform grid has low  
ESS

think of weights like a  
resolution or grid / a

grain size  
or a coarsened resolution to  
MC

imagine approximating integral over  
prob dist as a sum over a grid of points  
grid squares in low prob areas do not  
contribute a lot to the sum. sum is  
in fact dominated by small no. of points  
with highest "weights" / probabilities. so  
we want our grid to be. finer here but  
coarser in regions that don't contain much  
 $\rightarrow$  equiv. to concept of having similar  
weights everywhere! so small variance in  
weights everywhere!

## CHOOSING A GOOD IMPORTANCE FN.

(slides)

(most important part: g. importance sampling)

within class, named distributions:

- \* Can be shown (sheet 2) that theoretically optimal (minimum variance) importance fn. is

$$Q^*(\theta) = \frac{|f(\theta)| P(\theta)}{\int |f(\theta)| P(\theta) d\theta}$$

high variance in importance weight (competing estimates  
 $\Rightarrow$  small effective sample size.)

- \* However, if we can't directly sample from  $P(\theta)$ , then we probably can't sample from  $|f(\theta)| P(\theta)$

- \* Want to keep the importance weights roughly const. otherwise large variations in  $P(\theta)/Q(\theta)$  will lead to high variance of estimate, smaller ESS

- \* Effective sample size

i.e. importance weights have similar shape to dist. we are trying to sample

$$ESS = \frac{m}{1 + \widehat{\text{Var}}[\sum w^*(Q_i)]}$$

e.g. this is in parallel example

- \* In practice, find thick-tailed distribution  $Q(\theta)$  that is positive everywhere and similar in shape to  $|f(\theta)| P(\theta)$

- \* Don't want  $Q(\theta)$  small when  $|f(\theta)| P(\theta)$  large!

problem: small sample size, high variance

→ ideal importance fn. is proportional to  $|f(\theta)| P(\theta)$

want for tail closely resembles target prior (i.e. importance weights are const  $\Rightarrow$  small variance in weight  $\Rightarrow$  large ESS)

→ in practice want positive, everywhere, covers the tail (thinnest tail), and (thick) may be helpful: named distributions

## Intro (slides)

satellite galaxies  
large magellanic cloud  
SMC  
Andromeda  
Milky way (central galaxy)

### Astrostatistics case study

#### "Bayesian estimates of the milky way and andromeda masses using high-precision astrometry & cosmological simulations"

- want to measure mass of milky way (constraint mass of our galaxy)
- satellite galaxies of our milky way (by LMX)  $\rightarrow$  use their orbits/motions to estimate mass of our own galaxy
- measure distance of satellite galaxy to central galaxy, angular momentum etc.
- use simulation to plot kinematic properties against mass of central galaxy (correlations between the two)  $\rightarrow p(\text{angular mom.} | \text{mass})$
- use bayesian inference to infer actual mass of our galaxy from measured kinematic properties & results of simulations

Illustris simulation: simulate evolution of galaxies  
 small density fluctuations over density regions attract mass, galaxies form & subtract gas which forms stars  $\rightarrow$  galaxies  
 (web)  
 can measure properties of satellite  
 galaxy properties vs mass of central galaxy  $\rightarrow$  get catalogue of properties  
 simulation generates samples (prior plasma)  $p(\text{angular mom.})$  is intractable in sense that we can't write it down, but we can draw from it

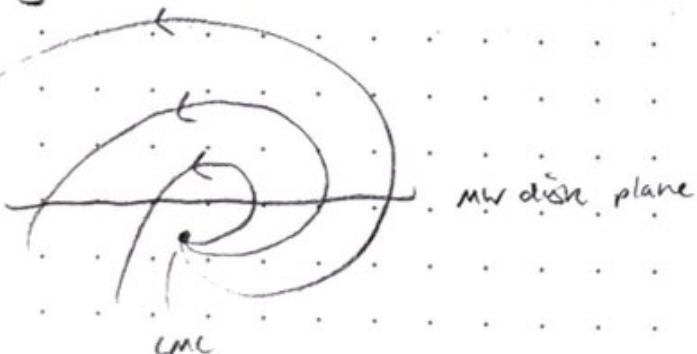
velocities ( $v$ ) positions ( $r$ ) momenta ( $p$ ) of satellites correlated w/ central galaxy mass via galaxy formation simulations (physical prior) (intractable but can sample from prior)  
 $x = \text{latent (true) values } v, r, p$        $M_{\text{vir}} = \text{mass of galaxy}$ ,  $\theta = \text{parameters } x, M_{\text{vir}}$   
 (joint prior  $p(v, r, p, M_{\text{vir}})$  from simulations)

measure LMC data  $v, r, p$ , we have analytic measurement likelihood  $p(x|d)$ , now use importance sampling to estimate posterior!

$$p(x|d) = L(x|d) = N(V_{obs}/V, \sigma_V^2) \times N(r_{obs}/r, \sigma_r^2) \times N(p_{obs}/p, \sigma_p^2)$$

- $\rightarrow$  Bayesian inference
- $\rightarrow$  Importance sampling
- $\rightarrow$  kernel density estimation

trajectories depend on mass of milky way



# IMPORTANCE SAMPLING THE MILKY WAY GALAXY (Aated+17)

$\vec{d}$  = (noisy) measurements of kinematic properties of our Milky way satellite's (e.g. Large Magellanic Cloud) properties (e.g. distance from central galaxy, velocity, angular momentum...)  $\vec{x} = v, r, j$ .

$\vec{x}$  = latent (true) values of satellite properties

$m = \log_{10} [\text{mass of central galaxy}]$

**GOAL:** Estimate  $E[g(m)|\vec{d}] = \int g(m) P(m|\vec{d}) dm$

(e.g. let  $g(m) = m$  for posterior mean or  $g(m) = (m - E[m|\vec{d}])^2$  for posterior variance)

marginalised posterior

$$\begin{aligned} \text{Posterior mean: } E[m|\vec{d}] &= \int m P(m|\vec{d}) dm \\ &= \int m \left[ \int P(\vec{x}, m|\vec{d}) d\vec{x} \right] dm \end{aligned}$$

$$P(\vec{x}, m|\vec{d}) = \frac{P(\vec{d}|\vec{x}) P(\vec{x}, m)}{\int \int P(\vec{d}|\vec{x}) P(\vec{x}, m) d\vec{x} dm}$$

(Gaussian)  
 measurement likelihood  
 ↓  
 Physical prior  
 (encoded in cosmology)

$P(d|x, m)$   
 why not  $P(x|m)$ ?  
 reasonable  
 assumption  
 if I know true  
 value e.g.  $v$  then  
 knows under not  
 provide any additional  
 information

$$E[m|\vec{d}] = \frac{\int \int m P(\vec{d}|\vec{x}) P(\vec{x}, m) d\vec{x} dm}{\int \int P(\vec{d}|\vec{x}) P(\vec{x}, m) d\vec{x} dm}$$

(Posterior  
is constant so  
don't need to integrate again)

Cosmo. sims generate galactic systems that can be thought of as giving samples  $(\vec{x}_j, m_j) \stackrel{iid}{\sim} P(\vec{x}, m)$   $j=1, \dots, m$  from the physical prior (intractable).

- ∴ want  $E[m|\vec{d}]$  so need  $P(m|\vec{d})$  which we get from marginalising  $P(\vec{x}, m|\vec{d})$ .
- we know  $P(\vec{x}, m)$  from sim so find  $P(m|\vec{d}) = P(\text{other}, m) P(\vec{x}|m) / P(\vec{d})$
- $P(\vec{d}|\vec{x})$  is just meas. error.
- then integrate to get  $E[m|\vec{d}]$  and integrate again to get  $E[E[m|\vec{d}]]$ .

→ use the prior itself as the importance function!

were we can get away with intractable (and prior?) since they cancel?

Apply Monte Carlo approx.:

$$\mathbb{E}[m|d] \approx \frac{\frac{1}{n} \sum_{j=1}^n m_j P(\vec{d}|\vec{x}_j)}{\frac{1}{n} \sum_{j=1}^n P(\vec{d}|\vec{x}_j)} = \sum_{j=1}^n m_j w_j$$

tractable  
prior?

from sim's we have samples from prior that not prior dist. itself.

now we only need samples from prior eval  $E[m]$  which we have?

Self-norm importance weight

$$w_j = \frac{P(\vec{d}|\vec{x}_j)}{\sum_{k=1}^n P(\vec{d}|\vec{x}_k)}$$

$$\sum_{j=1}^n w_j = 1$$

REMARKS

remove  $w_j$  ratio of normalized posterior to  $Q(m)$ ; so what happens to the  $Q$  in here? → effectively cancels out because we used prior as  $Q(m)$ .

$$w_j \propto \frac{P(d|\vec{x}_j) P(x_j|m)}{Q(m,m)}$$

cancels because encode prior as importance function itself

we're getting away without being able to evaluate density  $Q(m)$  itself (because we can only get samples from computer distribution) since it cancels with the prior!

$$P(d|x, m) = P(d|x)$$

$\vec{x}$  conditionally indep. of  $m$  given  $\vec{x}$ .

conditional independence assumption:  
measurements contain info about true dist.  $m$

- I know the  $x$  value
- can write down all the measurements and their
- no additional info contained in the mass

measuring distance or velocity  $y$  satellite galaxy does not probe mass directly  
→ the only info on mass comes from knowing  $m$  in the prior  $Q(m)$ .

## Lecture 16

28.2.25

Today: continuing Bayesian computation, importance sampling, MCMC

recommend: paper Speagle J "A Conceptual Introduction to Markov Chain Monte Carlo Methods"

why do we not just  
do grid sampling /  
importance sampling all the  
time?

### WEIGHTED KERNEL DENSITY ESTIMATE

(slides)

$$Wkde(\theta) = \sum_{s=1}^m w_s \times N(\theta | \theta_s, bw^2)$$

why do you Q?

$$w_s = \frac{p(\theta_s)}{\sum p(\theta_s)}$$

$w_s$  = normalised importance weights

bandwidth: Silverman's rule of thumb  $bw = \left(\frac{40^5}{3n}\right)^{1/5}$

Estimate  $\sigma^2$  from posterior  $\text{Var}[\theta]$  estimate.

Use importance sampling effective sample size (ESS) for  $n$ .

Ideal case:

if equal weights:  $w_i = \frac{1}{m}$  reduces to

$$kde(\theta) = \sum_{s=1}^m \frac{1}{m} \times N(\theta | \theta_s, bw^2)$$

$$\frac{m}{1 + \text{Var}[\ln w_s]}$$

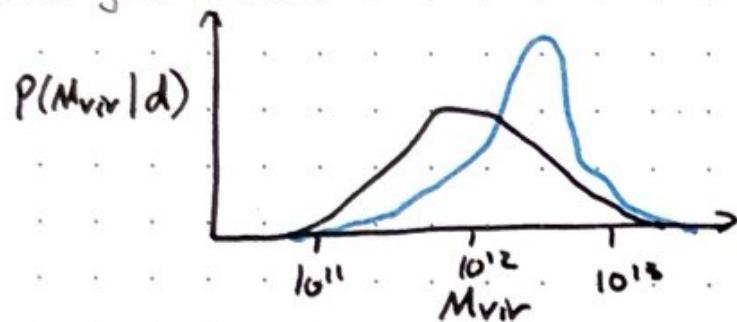
Continue from  
example before:

bisically  
"answer to  
sheet 2 q3"

(slides)

## Posterior of Milky Way w/ weighted KDE

posterior given LHC data

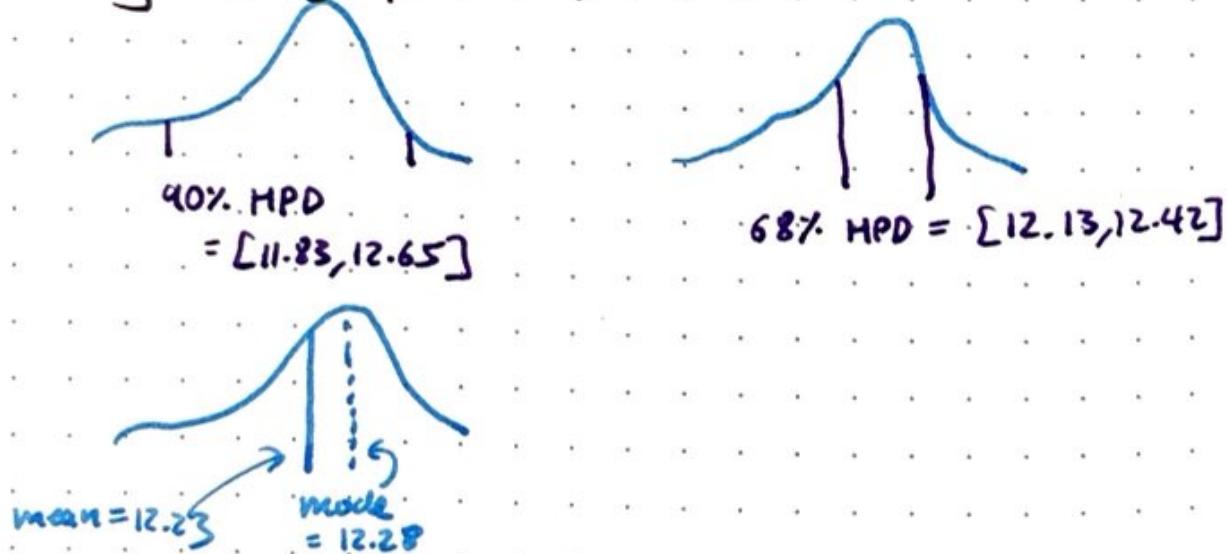


— prior  
— posterior  $P(M_{\text{vir}} | v_i, r_{ij})$   
MW  $M_{\text{vir}} = 1.7^{+1.33}_{-0.52} \times 10^{12}$

maximum weight?  
mean weight?

## Posterior HPD

X% HPD (X% credible interval(s)) with highest density  
(containing X% of posterior)



You make combined posterior from data from multiple satellite galaxies & do some calculation w/ weighted KDE

combine inferences from different sources to get more accurate estimate (This is an ideal case for bayesian inference)

From title: misconception that primary purpose of MCMC is to estimate the posterior, primary purpose is actually to estimate expectations! (i.e. integrals)

## MARKOV CHAIN MONTE CARLO MCMC

classical techniques for sampling from some probability dist.  
→ useful for dists that are too complex or high-dimensional to study  
with analytic techniques alone (e.g. dimensionality!).

Goal is to evaluate the posterior  $P(\theta|D)$

- \* Simple likelihoods/conjugate priors admit analytic solutions to the posterior
- \* Simple models may allow direct draws  $\theta_i \sim P(\theta|D)$   
i.e. "direct simulation"
- \* Small numbers of parameters  $p$  (small dimensionality)  
→ Evaluate posterior on a  $p$ -dimensional grid.  
Inefficient for  $p > 3$  (end up wasting compute time evaluating  $P(\theta|D)$  where it is close to zero.)
- \* Realistic models with many parameters  $p$   
(high dimensional parameter space) (and we have no analytic solution)  
↳ Markov Chain Monte Carlo is the gold standard!

whole

Point of MCMC is to generate samples from the posterior when we can't directly sample from it

$$E[S(\theta)|D] = \int s(\theta) P(\theta|D) d\theta \approx \frac{1}{m} \sum_{i=1}^m s(\theta_i)$$

\* Posterior simulation - MCMC

- \* How? - Generate a correlated sequence (chain) of random variates (monte carlo) that (in a limit) are draws from the posterior. The next value in the sequence only depends on the current values. (Markov)  
(explain how later)
- \* Algorithm cleverly constructed to ensure distribution of chain values → posterior dist. = stationary dist. in the long run.

## Simplest MCMC: METROPOLIS ALGORITHM

(slide)

- ① Choose random starting point  $\mu_0$
- ② At each step  $i=1, \dots, N_{MC}$  propose a new parameter value  $\mu_{\text{prop}}$  via  $J(\mu)$ . The proposal distribution  $J(\mu)$  is usually  $N(\mu_i, \tau^2)$   
The proposal scale  $\tau$  is chosen <sup>Carefully!</sup> for efficiency!  
drawn from proposal dist.  
I see jump  
hence like adding noise
- ③ Evaluate ratio of posteriors at proposed vs. current values  
Metropolis Ratio  $r = \frac{P(\mu_{\text{prop}}|y)}{P(\mu_i|y)}$
- ④ Is  $\mu_{\text{prop}}$  is a better solution (higher posterior),  $r > 1$ , accept the new value  
 $\mu_i = \mu_{\text{prop}}$   
Else accept with probability  $r$  (i.e. accept with probability  $\min(r, 1)$ ). Stay at same value  
 $\mu_i = \mu_i$   
(and include in chain - keeps repeated values)
- ⑤ Repeat 2-4 until reach some measure of convergence and gather enough samples to complete your inference.

MCMC = next value depends on current value  
Monte Carlo = randomness

**Example:** simple gaussian unknown  $\mu$  (slides)

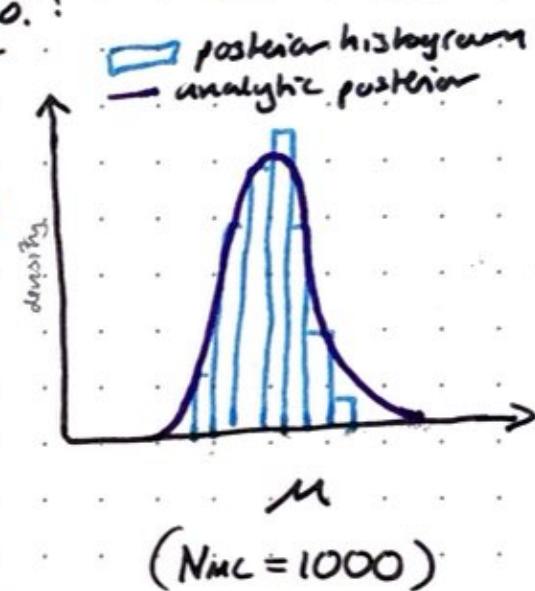
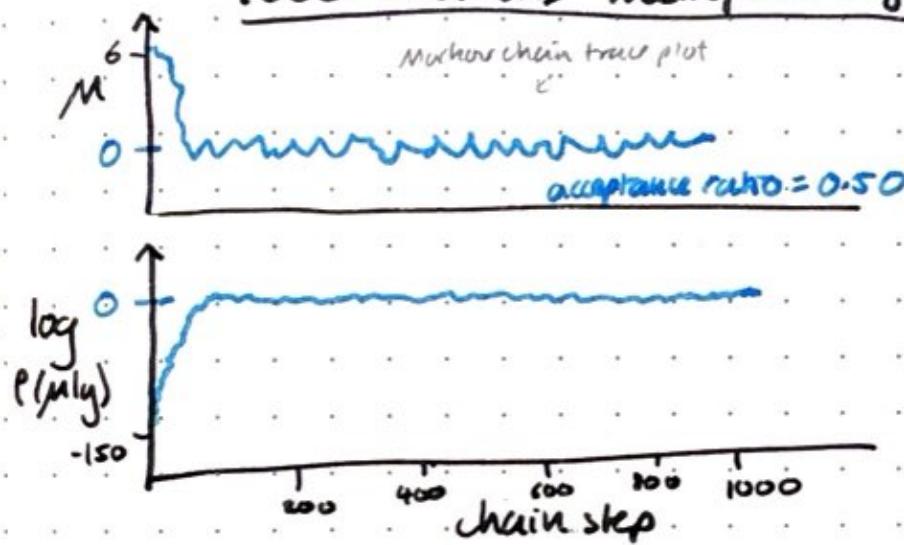
Likelihood:  $y_i \sim N(\mu, \sigma^2 = 1)$   $i=1, \dots, N$

Prior:  $P(\mu) \propto 1$

$\Rightarrow$  Posterior:  $P(\mu | y) = N(\mu | \bar{y}, \sigma^2/N)$  (analytic solution) (L17)

Do Metropolis step → (proposal scale 0.6) track no. of acceptances  
Code on moodle to try yourself (more  $\bar{y}$ -elements summing statistic various all info)  
accept data set by putting

1000 iterations metropolis algo.:

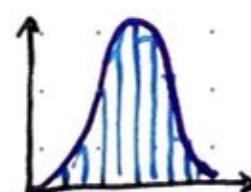
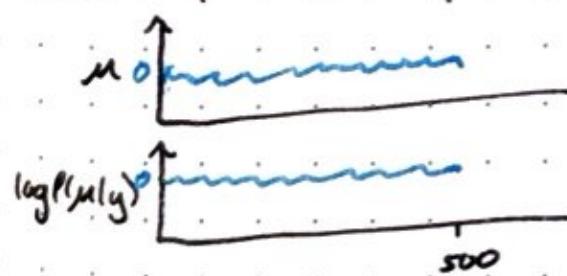


During burn-in e.g. first 50 iterations  
acceptance ratio = 0.42

acceptance ratio goes up as we make small steps ratio goes up so what he's saying

Can also

Cut off 1st half of chain to cut burn in (aggressive solution, more empirical cut-offs exist)



In reality, if we had 1D dist. probably wouldn't need to do MCMC, could just give CDF, invert & apply to uniform RV to get draw from it but show here as example.

Lecture 173.3.25

[Today: continuing Bayesian computation: MCMC]

Recap MCMC:

Generate a correlated sequence (chain) of random variates (Monte Carlo) that (in a limit) are draws from the posterior. The next value in the sequence only depends on current values (Markov).

ideally: want iid but in long run very fast we get enough indep steps that this works  
not iid (correlated) but looks like  
tends to large NMC

D-DIM METROPOLIS ALGORITHM: (slides)

Posterior  $P(\theta|D)$  where  $\dim(\theta) = d$

Symmetric proposal/jump dist.  $J(\theta^*|\theta) = J(\theta|\theta^*)$

- ① Choose random starting point  $\theta_0$
- ② At step  $i=1, \dots, N$  propose new parameter value

$$\theta^* \sim N(\theta_{i-1}, \Sigma_p)$$

The proposal distribution is

$$J(\theta^*|\theta_{i-1}) = N(\theta^*|\theta_{i-1}, \Sigma_p)$$

- ③ Evaluate ratio of posteriors at proposed vs current values

$$r = \frac{P(\theta^*|y)}{P(\theta_{i-1}|y)}$$

note to self:  
why do we metropolis? ->  
symmetric jump (not really)  
normal best-looking dist.  
usually normal so assume  
metropolis = normal  
for exam?

④ Accept  $\theta^*$  with probability

$$\min(r, 1) : \theta_i = \theta^*$$

If not accept, stay at same value  $\theta_i = \theta_{i-1}$  for the next step and include in chain.

⑤ Repeat 2-4 until reach some measure of convergence ( $\delta-R$ ) and gather enough independent samples to compute your inference (reduce Monte Carlo error).

widely used way to sample from MVG:

[HOW DO WE DRAW FROM MVG?]

(i.e. step 2 d-dim metropolis)

## Multivariate Gaussian Draws

$$\vec{\theta} \in \mathbb{R}^d$$

Metropolis proposal:  $\vec{\theta}^* \sim N(\vec{\theta}_{i-1}, \Sigma_p)$

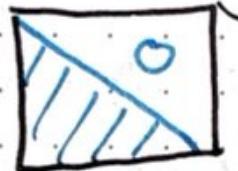
$\Sigma_p$  = real, symmetric, positive definite  $d \times d$  matrix

$\det(\Sigma_p) > 0$  and  $\Sigma_p^{-1}$  exists.

1) Cholesky decomposition, find  $L$  s.t.  $\Sigma_p = L L^T$

$L$  is lower triangular Cholesky factor

"matrix square root"



cholesky decomp expensive so don't want to do cholesky decomp at every point in your loops. so instead do this

2) Draw  $\vec{\theta}^* \sim N(\vec{\theta}_{i-1}, \Sigma_p)$

by  $\vec{\theta}^* = \vec{\theta}_{i-1} + L \vec{z}$

where  $\vec{z}$  is a  $d$ -vector of iid uniform unit normal r.v.s

$z_i \stackrel{iid}{\sim} N(0, 1) \quad i=1, \dots, d$

check that sampling  $\vec{\theta}^*$  like this is indeed the same as sampling from MVG:

Check:  $E[\vec{\theta}^*] = \vec{\theta}_{i-1}$

Wiki:  $\vec{\theta}_{i-1} + L \vec{z}$  has the desired dist. due to the affine transformation property: any affine transformation  $g(x) = Ax + b$  of a gaussian is a gaussian

$\text{Cov}[\vec{\theta}^*, \vec{\theta}^{*T}] = \text{Cov}[L \vec{z}, (L \vec{z})^T] \quad (\text{billinear})$

$$= L \text{Cov}[\vec{z}, \vec{z}^T L^T] = L \text{Cov}[\vec{z}, \vec{z}^T] L^T$$

$$= L I_d L^T = L L^T = \Sigma_p$$

note to self: SUMMARY  
1. drawing dist. on the  
multivariate gaussian  
due to affine invariance property  
2. can show resulting gaussian  
has preserved the particular  
band-diagonal covariance matrix //

Example:

(slides)

Recall: analytic posterior density for Gaussian  $(\mu, \sigma^2)$  model with non-informative priors  $P(\mu) \propto 1$ ,  $P(\sigma^2) \propto \sigma^{-2}$ ,  $\sigma^2 > 0$ .  
(multi-param. inference)

\* Joint posterior:  $(\sigma^2 > 0)$

$$P(\mu, \sigma^2 | y) \propto (\sigma^2)^{-n/2 - 1} \exp\left(-\frac{(n-1)s^2}{2\sigma^2}\right) \exp\left(-\frac{n}{2\sigma^2}(\bar{y} - \mu)^2\right)$$

\* Marginal of  $\mu$ : "marginal posterior" (1)

$$\begin{aligned} P(\mu | y) &= \int P(\mu, \sigma^2 | y) d\sigma^2 \propto \left[ 1 + \frac{n(\mu - \bar{y})^2}{(n-1)s^2} \right]^{-n/2} \\ &= t_{n-1}(\mu | \bar{y}, s^2/n) \end{aligned}$$

\* Marginal of  $\sigma^2$ : "marginal posterior" (2)

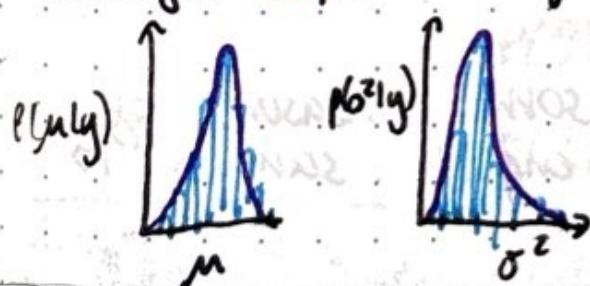
$$P(\sigma^2 | y) = \int P(\mu, \sigma^2 | y) d\mu = \text{Inv-}\chi^2(\sigma^2 | n-1, s^2)$$

\* Joint posterior

$$P(\mu, \sigma^2 | y) = P(\mu | \sigma^2, y) P(\sigma^2 | y) = N(\mu | \bar{y}, \sigma^2/n) \times \text{Inv-}\chi^2(\sigma^2 | n-1, s^2)$$



Metropolis 2D example  
- analytic  posterior histogram



MC chain  
 $N_{mc} = 2000$   
acc. ratio = 0.39

Reminder: if  $\underline{x} \sim N(\underline{0}, \Sigma)$  then  $\Sigma = \text{cov}(\underline{x}, \underline{x})$ . (some prefer to denote this  $\text{cov}(\underline{x}, \underline{x}^T)$ )

e.g.  $\underline{x}_j \sim N(\underline{0}, \Sigma)$   $\Sigma = \begin{pmatrix} \text{cov}_{11} & \text{cov}_{12} & \dots \\ \text{cov}_{21} & \text{cov}_{22} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} = \begin{pmatrix} \sigma_{11}^2 & \rho \sigma_{12} \sigma_{21} & \dots \\ \rho \sigma_{21} \sigma_{12} & \sigma_{22}^2 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$  but prefer to think components i.e.  $\Sigma_{ij} = \text{cov}(x_i, x_j)$

think in components!

$$\text{Cov}[\theta_i^k, \theta_j^k] = \text{Cov}[\theta_i + L_i z_k, \theta_j + L_j z_k] = \text{Cov}[L_i z_k, L_j z_k] = L_i L_j \text{Cov}[z_k, z_k] = L_i L_j$$

$\text{cov}(ax, b+cy) = \text{cov}(x, y)$   
since  $\theta_i$  const. ( $i$  is a number)

$\text{cov}(ax, by) = ab \text{cov}(x, y)$

since  $z_i$  iid  
 $\text{cov}(z_i, z_j) = \delta_{ij}$

(Metropolis in this case)  
algo generates a chain  
is "Markov process"  
combine this w/  
"monte carlo" nature  
of simulating new  
positions = MCMC → side note

## METROPOLIS - HASTINGS ALGORITHM

More General Jumping rule  $J(\theta^+ | \theta_i)$

(Need not be symmetric)  $J(\theta_a | \theta_b) \neq J(\theta_b | \theta_a)$

- ① Choose random starting point  $\theta_0$
- ② At step  $i=1, \dots, N$  propose a new parameter value

$$\theta^* \sim J(\theta^+ | \theta_{i-1})$$

- ③ Evaluate M-H ratio of posteriors at proposed vs current values.

$$r = \frac{P(\theta^* | y) / J(\theta^* | \theta_{i-1})}{P(\theta_{i-1} | y) / J(\theta_{i-1} | \theta^*)}$$

- ④ Accept  $\theta^*$  with probability

$$\min(r, 1): \theta_i = \theta^*$$

If not accept, stay at same value  $\theta_i = \theta_{i-1}$  and include in chain.

- ⑤ Repeat 2-4 until reach some measure of convergence and gather enough samples to compute your inference.

### METROPOLIS-HASTINGS

Theory: see wikipedia

Want markov process (chain) to reach stationary distribution which is equal to the posterior  
 now to achieve this  $\Rightarrow$  stationary condition for stationary process is condition of detailed balance (prob being in state  $\theta_i$  & trans.  
 to state  $\theta_{i+1}$  is equal to prob being in state  $\theta_{i+1}$  & transiting to state  $\theta_i$ ).  
 (not necessary) to state  $\theta_{i+1}$  is equal to prob being in state  $\theta_{i+1}$  & transiting to state  $\theta_i$ ).  
 Stationary dist. process, flow of sample between states as chain progresses so such that overall dist of samples is stationary.  
 also want this stationary dist to be unique (i.e. no existing markov chain converging to dist other than the same as  
 this is guaranteed by the ergodicity of the markov process (ergodicity = avg. behavior over time is the same as  
 avg. behaviour observed over many diff. instances ('estimates') of the system). So avg. trajectory.  
 for any starting point  $\Rightarrow$  we always converge to same stationary dist so it is unique?

design a markov process (by constraining transition probability) that fulfills these 2 conditions  
 88 and so stationary dist is the posterior  $\Rightarrow$  this is how metropolis hastings algo is

- \* d-dim Metropolis is just a special case, where  $J(\theta^* | \theta_i) = N(\theta^* | \theta_i, \Sigma_p) = N(\theta_i | \theta^*, \Sigma_p) = J(\theta_i | \theta^*)$  is a symmetric proposal dist.
- \* More general asymmetric proposals, allow "biased" proposals  $\rightarrow$  more probable to propose towards a certain direction.
- \* With some knowledge of structure of the posterior, can sometimes engineer a clever proposal  $J(\theta | \theta_i)$

$\downarrow$   
ex 10

## GIBBS SAMPLING

use when:  $\begin{cases} \text{can't sample from } P(\theta, \phi) \text{ directly but} \\ \text{can sample from } P(\theta|\phi), P(\phi|D) \text{ directly.} \\ \text{(conditional)} \end{cases}$

(cont.)

- \* Special case of Metropolis-Hastings (Gelman BDA ch 11)

Multi-dim. sampling, when you can utilize the set of conditional posterior distributions as proposal distns. for each parameter - Metropolis-Hastings ratio = 1 (always accept)

- \* If joint posterior is  $P(\theta, \phi | D)$

- \* And you can solve for tractable (you can draw) conditionals:

$$P(\theta | \phi, D)$$

$$P(\phi | \theta, D)$$

- \* Jump along each parameter-dimension one at a time

DERIVE: algos: metropolis, dist  $\pi_i$ , if Metropolis fails then starting dist exists, to construct algo, choose starting dist to be poster dist  $\pi_i$

- detected balance:  $P(\theta_i)P(\theta_{i+1} | \theta_i) = P(\theta_{i+1})P(\theta_i | \theta_{i+1})$  so  $\frac{P(\theta_{i+1} | \theta_i)}{P(\theta_i | \theta_{i+1})} = \frac{P(\theta_{i+1})}{P(\theta_i)}$
- prob transition  $\rightarrow P(\theta_{i+1} | \theta_i) = J(\theta_{i+1} | \theta_i) A(\theta_{i+1} | \theta_i)$  prob propose
- so  $A(\theta_{i+1} | \theta_i) = \frac{P(\theta_{i+1})}{P(\theta_i)} J(\theta_{i+1} | \theta_i)$  prob accept
- choose acceptance ratio that satis. this condition: Metropolis rule is  $A(\theta_{i+1} | \theta_i) = \min(1, \frac{P(\theta_{i+1})}{P(\theta_i)} J(\theta_{i+1} | \theta_i))$

works b/c either  $A(\theta_{i+1} | \theta_i) = 1$  or  $A(\theta_i | \theta_{i+1}) = 1$  so condition is always satisfied

## 2-dim GIBBS SAMPLER

cons/computational  
resources  $\rightarrow$  always  
accept

- ① Choose random starting point  $(\theta^0, \phi^0)$

- ② At cycle  $t$ , update

$$\theta^t \sim P(\theta | \phi^{t-1}, D)$$

- ③ Then update

$$\phi^t \sim P(\phi | \theta^t, D)$$

(each complete set (pair) of updates is called  
a Gibbs cycle)

- ④ Record current values of chain  $(\theta^t, \phi^t)$
- ⑤ Repeat 2-4 until reach some measure of convergence (G-R) and gather enough indep. samples to compute your inference (reduce Monte Carlo error).

## d-DIM GIBBS SAMPLER

Parameter vector  $\theta = (\theta_1, \dots, \theta_d)$

Current state at  $j$ -th update within cycle  $t$ :

$$\underbrace{(\theta_{+j}^t, \theta_{-j}^t)}_{\substack{\text{value of component } j \\ \text{at iteration } t}} \quad \underbrace{\theta_{-j}^t}_{\substack{\text{value of all components} \\ \text{except } j \text{ at iteration } t \\ \downarrow \\ \text{d-dim vector}}} \quad \begin{matrix} \text{value of all components} \\ \text{except } j \text{ at iteration } t \\ \downarrow \\ \text{d-dim vector} \end{matrix}$$

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

iteration  $t+1$  values for  
components already updated (?)

iteration  $t+1$  (parameters) values for  
components not yet updated (?)

- ① Choose random starting point  $\theta^0$
- ② At cycle  $t$ , update through the  $d$ -parameters:  
For each  $j=1, \dots, d$  move the  $j$ th parameter to  
 $\theta_j^{t+1} \sim p(\theta_j | \theta_{-j}^t, D)$   
(update  $\theta_j$  conditional on current values of  
all other parameters)
- ③ After updating all  $d$  parameters, record  
current state  $\theta^{t+1}$
- ④ Repeat 2-3 until reach convergence and  
enough samples

Gelman slightly diff notation:

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

Propose at iteration  $t$ :  $\theta_j^t \sim p(\theta_j | \theta_{-j}^{t-1})$

(e.g.  $d=3$ ) start  $\theta^0 = (0^0, 0^0, 0^0)$

$$\theta_1^1 \sim p(\theta_1 | 0^0, 0^0)$$

$$\theta_2^1 \sim p(\theta_2 | 0_1^1, 0^0)$$

$$\theta_3^1 \sim p(\theta_3 | 0_1^1, 0_2^1)$$

$$\theta_1^2 \sim p(\theta_1 | 0_1^1, 0_2^1)$$

$$\theta_2^2 \sim p(\theta_2 | 0_1^1, 0_2^1)$$

$$\theta_3^2 \sim p(\theta_3 | 0_1^1, 0_2^1)$$

## Gibbs Sampling: Example

(Wides)

(Gelman BDA section 11.1)

consider single observation  $y_1, y_2$

Likelihood:  $y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \sim N \left( \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$

( $\rho$  known)

Priors:  $P(\theta_1) = P(\theta_2) \propto 1$

Posterior:  $\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} | y \sim N \left( \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$

$P(\theta | y) = P(\theta_1, \theta_2 | y) P(\theta_2 | y) = P(\theta_2 | \theta_1, y) P(\theta_1 | y)$

Simple to draw directly from joint but for purpose of demonstrating gibbs sampling, find conditionals

Shrubby bivariate gaussian, joint can be decomposed into conditionals

→ Conditional Postiors (e.g. properties of MVGs): SEE MIDDLE

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

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

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

true

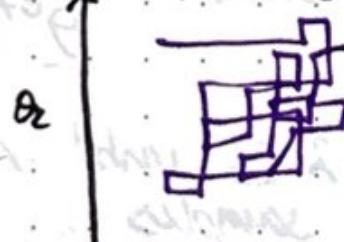
gaussian

we can draw from

$\theta_2$



Gibbs sampling trace path:  
 $P(\theta_1, \theta_2 | y)$



right angle!  
only update 1 parameter at a time

Reminder:

is plane

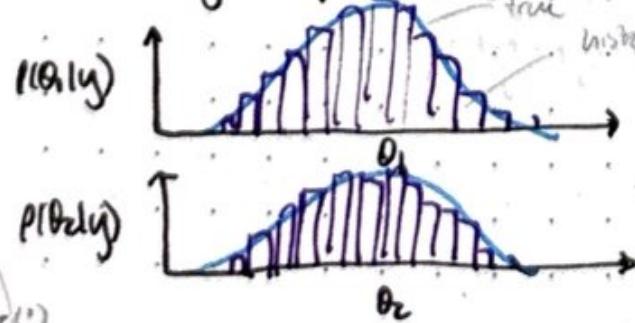
say  $\rho = 0.99$

gibbs sampler performs badly  
because gibbs sampler can only move in horizontal/vertical axis

bit plot is diagonal  
so hard for gibbs sampler to explore the space!

→ note will also be a problem  
for metropolis hastings  
but especially a problem here!?

Marginal posterior densities



## Lecture 18

ex sheet 3

(play w sampling  
only on moodle)

5.3.25

Today: Continuing Bayesian Computation:  
Metropolis-Hastings, Gibbs Sampling, applications  
to SN cosmology, Assessing convergence and  
Mixing.

GIBBS SAMPLING STEP AS A SPECIAL CASE OF  
METROPOLIS-HASTINGS (Gelman 11.3) (continued)

"Gibbs is a special case of metropolis hastings that always accepts."  
show this.

At the  $j$ th parameter update within iteration  $t$   
(target  $\pi(\vec{\theta}) = P(\vec{\theta} | \vec{D})$ )

$$\text{M-H ratio} \quad r = \frac{\pi(\vec{\theta}^*) / J(\vec{\theta}^* | \vec{\theta}^{t-1})}{\pi(\vec{\theta}^{t-1}) / J(\vec{\theta}^{t-1} | \vec{\theta}^*)}$$

Recall  $\vec{\theta}_j^{t-1} = (\underbrace{\theta_1^t, \dots, \theta_{j-1}^t}_{\text{have updated up to component } j-1}, \theta_j^{t-1}, \theta_{j+1}^{t-1}, \dots, \theta_d^{t-1})$   
current components of  $\vec{\theta}$  except  $\theta_j$

In Gibbs, we choose

can only move to  
new proposals where  
ones we are not updating  
have to stay the same

$$J(\vec{\theta}^* | \vec{\theta}^{t-1}) = \begin{cases} \pi(\vec{\theta}_j^* | \vec{\theta}_{-j}^{t-1}) & \text{if } \vec{\theta}_j^* = \vec{\theta}_j^{t-1} \\ 0 & \text{otherwise} \end{cases}$$

more detail in  
notes on moodle

$$r = \frac{\pi(\vec{\theta}^*) / \pi(\vec{\theta}_j^* | \vec{\theta}_{-j}^{t-1})}{\pi(\vec{\theta}^{t-1}) / \pi(\vec{\theta}_j^{t-1} | \vec{\theta}_{-j}^*)} \quad \vec{\theta}_j^* = \vec{\theta}_j^{t-1}$$

use this to  
expand terms into  
integrals.

$$= \frac{\pi(\vec{\theta}_j^* | \vec{\theta}_{-j}^*) \pi(\vec{\theta}_{-j}^*) / \pi(\vec{\theta}_j^* | \vec{\theta}_j^{t-1})}{\pi(\vec{\theta}_j^{t-1} | \vec{\theta}_{-j}^{t-1}) \pi(\vec{\theta}_{-j}^{t-1}) / \pi(\vec{\theta}_j^{t-1} | \vec{\theta}_{-j}^*)} = \frac{\pi(\vec{\theta}_j^*)}{\pi(\vec{\theta}_j^{t-1})} = 1$$

ALWAYS  
ACCEPT

saves on computation  $\rightarrow$  always accept, don't need to calc. ratio

another variations:

(slide >)

## Metropolis within Gibbs

Gibbs's problem: Parameters  $\theta$   $\rightarrow$  useful if natural posterior of parameters in problem, but difficult to solve for conditionals  $\rightarrow$  Metropolis substitution rule.

$$\theta = (\theta_1, \dots, \theta_d)$$

$$\theta_{-j} \equiv (\theta_1, \dots, \theta_{j-1}, \theta_{j+1}, \dots, \theta_d)$$

- \* When you can't solve for tractable conditional distributions for all  $\theta_j$ :

$$P(\theta_j | \theta_{-j}, D)$$

- \* Replace each step for updating each  $j$ th parameter  $\theta_j$  with a separate Metropolis rule, compute Metropolis ratio and accept/reject.

- \* Cycle through all parameters and repeat for all  $N$  MCMC steps

multiple accept/rejects within full scan of parameters

PREVIOUS PAGE

Gibbs as special case of MH:

$$MH: r = \frac{P(\theta^t)}{P(\theta^{t-1})} \frac{\pi(\theta^{t-1} | \theta^t)}{\pi(\theta^t | \theta^{t-1})}$$

(For Gibbs)

$$\text{Propose proposal distribution } \pi(\theta^t | \theta^{t-1}) = \begin{cases} \pi(\theta_j^t | \theta_{-j}^{t-1}) & \text{if } \theta_j^t = \theta_{-j}^{t-1} \\ 0 & \text{otherwise} \end{cases}$$

only propose new points where all components are the same except  $\theta_j^t$ , where  $\theta_j^t$  is picked according to  $P(\theta_j^t | \theta_{-j}^{t-1})$

$$\therefore r = P(\theta_j^t | \theta_{-j}^{t-1}) \cdot \frac{P(\theta_{-j}^{t-1} | \theta_j^t)}{P(\theta_j^{t-1} | \theta_{-j}^{t-1})}$$

$$\begin{aligned} r &= \frac{P(\theta_j^t | \theta_{-j}^{t-1}) P(\theta_{-j}^{t-1}) P(\theta_j^{t-1} | \theta_{-j}^{t-1})}{P(\theta_j^{t-1} | \theta_{-j}^{t-1}) P(\theta_{-j}^{t-1}) P(\theta_j^t | \theta_{-j}^{t-1})} \\ &\quad (\text{cancel since } \theta_j^t = \theta_{-j}^{t-1}) \end{aligned}$$

= 1  $\Rightarrow$  ALWAYS ACCEPT

so we directly propose what we need to make the cycle through parameters!

## d-dim Metropolis-within-Gibbs Sampler (slides)

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

① Choose a random starting point  $\theta_0$

② At cycle  $t=1, \dots, N$ , cycle through the d-parameters:

A For each  $j=1, \dots, d$ , propose a new  $j$ th parameter value from a 1-dimensional Gaussian

$$\theta_j^* \sim N(\theta_j^t, \tau_j^2)$$

proposals scale  $\tau^2$   
convergence of gaussian

B Evaluate ratio of posteriors at proposed vs current values:

$$r = P(\theta_j^*, \theta_{-j}^t | D) / P(\theta_j^t, \theta_{-j}^t | D)$$

$$= P(\theta_j^* | \theta_{-j}^t, D) / P(\theta_j^t | \theta_{-j}^t, D)$$

i.e.  $P(\text{new point}) / P(\text{old point})$   
(note: only component j changes!  $\theta_{-j}^t$  stays)  
Hence

C Accept  $\theta_j^{t+1} = \theta_j^*$  with prob  $\min(r, 1)$ , otherwise  $\theta_j^{t+1} = \theta_j^t$

③ After full cycle, record current values  $\theta^{t+1}$

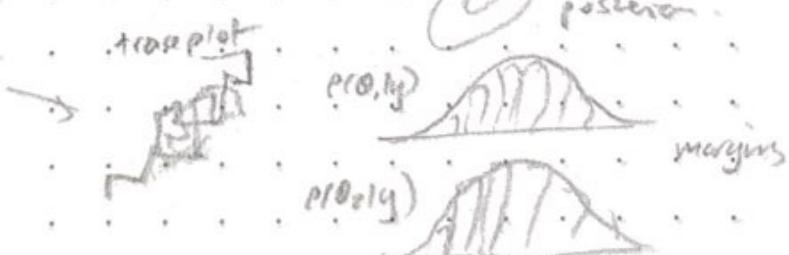
④ Repeat steps 2 for all parameters until convergence and enough samples

Note: code deposit  
example in lecture  
(same as last!)

adjust  $\tau^2$ :

empirically aim  
for 30% - 40%  
acceptance ratio

too small steps will accept every point  
but not explore parameter space, too large  
more anywhere



- advantage: → always accept
- disadvantage: → exact value very conditional derived from posterior gives algorithm implementation about which direction is likely to have more probability.
- disadvantage: → requires analytic properties,  
but if we can't derive them, can still use metropolis with Gibbs

## \* Mixed Gibbs Sampler:

→ Can replace sampling from conditionals with accept/reject for just the parameters with intractable conditionals

in cycle t  
for each j

if we do have analytic conditional  
just propose new point from that &  
always accept

if not

propose new point from normal  
and accept w/ prob  $\min(r, 1)$   
(aka metropolis)

$$r = \frac{p(\theta_j^{t+1}, \theta_{-j}^t)}{p(\theta_j^t, \theta_{-j}^t)}$$

## Tuning d-dim Metropolis

adjust proposal scale

- \*  $\theta^* \sim N(\theta_i, \Sigma_p)$ : If proposal scale  $\Sigma_p$  is too large, will get too many rejections and not go anywhere. If proposal scale too small, you will accept very many small moves: inefficient random walk.

### Laplace Approximation

$$P(\theta|D) \approx N(\theta|\theta_{MAP}, A^{-1})$$

$\theta_{MAP}$  = posterior mode

$$A_{ij} = -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln P(\theta|D) \Big|_{\theta=\theta_{MAP}}$$

✓ affine posterior cov matrix  $A$

- \* Choose  $\Sigma_p = c^2 A^{-1}$        $c \approx 2.4 / \sqrt{\dim(\theta)}$

- \* Scale Proposal to aim for acceptance ratio of 44% in 1D, 23% in  $\dim > 5$ .

notes will capture correlation  
of p. high

target highly  
concentrated

Laplace approx. favored  
proposal: strong correlation  
directions

## LAPLACE APPROXIMATION

Unnormalised Posterior:  $P(\theta|D) = P(D|\theta)P(\theta)$

Find MAP estimate: "max-a posterior"

$$\theta_{MAP} = \operatorname{argmax} \ln \tilde{P}(\theta)$$

$\leftarrow \tilde{P}(\theta|D)$  with Dobs  
plugged in?

Taylor expansion (first deriv. is zero):

$$\ln P(\theta|D) \approx \ln \tilde{P}(\theta_{MAP}|D) - \frac{1}{2}(\theta - \theta_{MAP})^T A (\theta - \theta_{MAP}) + \dots$$

$$A_{ij} = -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln \tilde{P}(\theta|D) \Big|_{\theta=\theta_{MAP}}$$

$$\hat{P}(\theta|D) \approx \hat{P}(\theta_{MAP}) \times \exp\left(-\frac{1}{2}(\theta - \theta_{MAP})^T A (\theta - \theta_{MAP})\right)$$

Gaussian Approximation:

(approximate  
posterior)

$$P(\theta|D) \approx N(\theta|\theta_{MAP}, A^{-1})$$

$A^{-1}$  is an approximate posterior covariance matrix.

reminder:

multidim. Taylor exp:

$$\cdot \text{let } \epsilon(\theta) = \ln \tilde{P}(\theta|D)$$

$$\cdot \epsilon(\theta) = \epsilon(\theta_{MAP}) + \nabla \epsilon(\theta_{MAP})^T (\theta - \theta_{MAP}) + \frac{1}{2} (\theta - \theta_{MAP})^T H(\theta_{MAP}) (\theta - \theta_{MAP}) + \dots$$

grad w.r.t.  
at  $\theta_{MAP}$

$$\therefore \epsilon(\theta) \approx \epsilon(\theta_{MAP}) - \frac{1}{2} (\theta - \theta_{MAP})^T A (\theta - \theta_{MAP})$$

(hesseinv. is taken matrix  
of 2nd order derivatives evaluated  
at  $\epsilon$ )

$A$  is w.r.t.  
 $H$  eval at  
 $\theta_{MAP}$

# Astrostatistics Case study: Supernova Cosmology

measuring brightness of SN:

Cosmology: previous control relationship between distance & redshift  
state equation EOS parameterized by  $w$ , ( $\rho = w\rho_0$ )

modern cosmology  $\rightarrow$  is  $w = -1$ ?

type Ia SN almost standard candles

$\hookrightarrow$  measure their magnitude & plot vs redshift

curve changes depending on how we tune

params.  $S_{R,\infty}, S_{R,0}$   $\rightarrow$  so observed data constrains

the parameters  $\rightarrow$  gamma analysis, told us  $\Delta n > 0$

comoving distance is integral from zero to redshift

distance probed by SN is "luminosity distance"

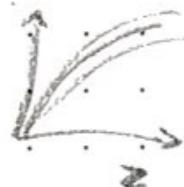
get  $dc$  by multiplying that by dimensions (of unit Mpc)

distance redshift also influenced by curvature of universe etc.

choose diff. param. plot

theoretical distance modulus  
vs. redshift

diff curve for diff galaxies etc.



data:



want to fit data to curve (fix parameter?)

(slides)

(lots of inconsistent notations here - ignore it)

For us: simplify, let's say SNe are standard candles.

$$M_s \sim N(M_0, \sigma_{int}^2)$$

$$M_s = M_s + \mu(z_s; \theta)$$

Population Distribution

(log) inverse square law

Assume data (apparent magnitudes and redshifts)  
 $\{M_s, z_s\}$  are measured perfectly

$$\theta = (H_0, \Omega_m, \Omega_c, w)$$

Cosmological Parameters

First assume  $w = -1$  (cosmological constant)

Derive model, likelihood, posterior

$$\tilde{\Sigma} = (\Omega_m, \Omega_c, w), h \equiv H_0 / (100 \text{ km s}^{-1} \text{ Mpc}^{-1})$$

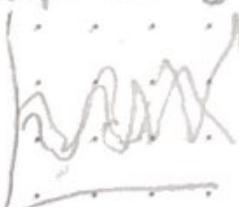
$$M(z_s; H_0, \tilde{\Sigma}) = 25 + \log_{10} \left( \frac{c}{H_0} d_L(z_s; \tilde{\Sigma}) \right)$$

$$M_s = M_s + \mu(z_s; H_0, \tilde{\Sigma})$$

$$p(M_s | z_s) = \dots$$

$M_0$  and  $\log_{10}(h)$  cannot be separately constrained!  
(not separately identifiable)

multiple ended chain  
overlap: "mixing"  $\rightarrow$  good sign, not getting stuck in place



ex sheet 3

true plot 4D metropolis  
use Laplace approx, metropolis  
to sample posterior

SN cosmology Reparametrisation

(40 min)

$$M_s \sim N(M_0, \sigma_{int}^2)$$

$$m_s = M_0 + 5 \log_{10} \left( \frac{d}{10 \text{pc}} \right)$$

$$M_s = M_0 + E_{int}$$

$$= M_0 + 5 \log_{10} \left( \frac{d}{\text{Mpc}} \right) + 25$$

$$m_s = M_0 + \mu(z_s | H_0, \vec{\Sigma}) + E_{int}$$

$$= M_0 + 5 \log_{10} \left( \frac{c}{H_0 d_L} \right) + 25$$

$$\vec{\Sigma} = (\Sigma_M, \Sigma_\Lambda, w)$$

$$M_0 + \mu(z | H_0, \vec{\Sigma}) = M_0 + 25 + 5 \log_{10} \left( \frac{c}{100} \frac{100}{H_0} d_L(z_s, \vec{\Sigma}) \right)$$

↑ 100 km s⁻¹ Mpc⁻¹

$$= M_0 - 5 \log_{10} \left( \frac{H_0}{100} \right) + 25 + 5 \log_{10} \left( \frac{c}{100} d_L(z_s, \vec{\Sigma}) \right)$$

$$M_0 - 5 \log_{10} h + \gamma(z_s, \vec{\Sigma})$$

$$M_0$$

important?  
 Why only three SN here (in sheet 1) had cepheids as well for which we knew their true distances) but there no cepheids to constraint also may so degeneracy between also may g. SN &  $H_0$ .  
 cannot distinguish between  $M_0$  &  $\gamma(z_s, \vec{\Sigma})$  alone.

So can only constrain the combination  
 is degenerate problem. b/c we can only constrain the linear combination of which one are really  $M_0$  &  $\log_{10} H_0$ .

$M_0$  and  $5 \log_{10} h$  are degenerate!

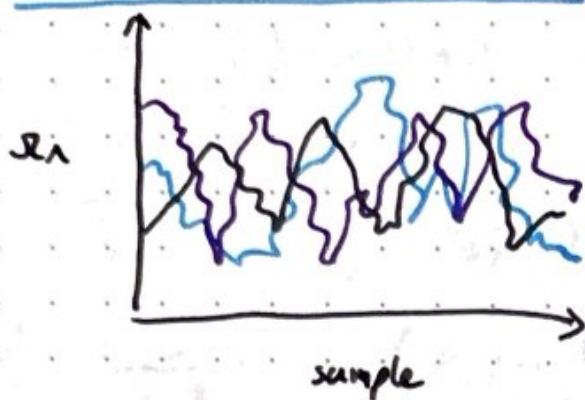
$$P(m_s | z_s; M_0, \vec{\Sigma}, \sigma_{int}^2) = N(m_s | M_0 + \gamma(z_s, \vec{\Sigma}), \sigma_{int}^2)$$

$$\text{Likelihood: } L(M_0, \vec{\Sigma}, \sigma_{int}^2) = \prod_{s=1}^N P(m_s | z_s; M_0, \vec{\Sigma}, \sigma_{int}^2)$$

$$(\text{Indep.}) \text{ priors: } M_0, \vec{\Sigma} \sim U(-\infty, \infty), \quad \sigma_{int}^2 \sim U(0, \infty)$$

$$\text{Posterior: } P(M_0, \vec{\Sigma}, \sigma_{int}^2 | \mathcal{E}_{M_s, z_s}) \propto L(M_0, \vec{\Sigma}, \sigma_{int}^2) \times P(M_0) P(\vec{\Sigma}) P(\sigma_{int}^2)$$

## Multiple Indep. Chains



chains overlap: "mixing"  
= good

## Assessing Convergence w/ multiple chains: Gelman - Rubin (G-R) ratio

Suppose we have simulated  $m$  parallel sequences, each of length  $n$  (after discarding the first half of the simulations). For each scalar estimated  $\theta$ , we label the simulation draws as  $\theta_{ij}$  ( $i = 1, \dots, n$ ;  $j = 1, \dots, m$ ), and we compute  $B$  and  $W$ , the between- and within-sequence variances:

$$B = \frac{n}{m-1} \sum_{j=1}^m (\bar{\theta}_{..j} - \bar{\theta}_{...})^2$$

(variance of the sum of all samples in each chain)

$$\bar{\theta}_{..j} = \frac{1}{n} \sum_{i=1}^n \theta_{ij}, \quad \bar{\theta}_{...} = \frac{1}{m} \sum_{j=1}^m \bar{\theta}_{..j}$$

avg over samples in given chain  
avg over every chain

$$W = \frac{1}{m} \sum_{j=1}^m s_j^2$$

(variance of each chain averaged over all chains)

$$s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\theta_{ij} - \bar{\theta}_{..j})^2$$

We can estimate  $\text{Var}(\theta|y)$ , the marginal posterior variance of the estimated estimand, by a weighted average of  $W$  and  $B$ , namely

$$\widehat{\text{Var}}^+(\theta|y) = \frac{n-1}{n} W + \frac{1}{n} B$$

G-R ratio:  $\hat{R} = \sqrt{\frac{\widehat{\text{Var}}^+(\theta|y)}{W}}$

102 To give measure of convergence/how well they have mixed

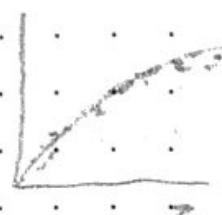
tune proposal scale to give acceptance ratio ~ 30%

run chains for 100000s  $\rightarrow$  want GE ratio  $\approx 1$

cutoffs  $\sim 20\sigma$  for burn-in (rule of thumb)

see slides sheet 3

SNRs may



Nobel prize 2011

MCMC  $\rightarrow$  can also print posterior on derived quantities (e.g.)  
see slides 1 (value)  $\rightarrow$  ex sheet 3

(1a)

$\rightarrow$  can compute posterior of young derived quantity e.g.  
deceleration parameter

gives high prob  $q < 0 \rightarrow$  universe is decelerating

separate chain

sums all  
samples in one  
chain

avg of sums of  
samples in each  
chain over  
all chains

$$\rightarrow B = \frac{1}{m-1} \sum_{j=1}^m \left( \sum_{i=1}^n \bar{x}_{ij} - \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n \bar{x}_{ij} \right)^2$$

think of  $B$  as variance in  
the mean of each chain times  $n$   
since  $\text{Var}(\bar{x}) = \frac{1}{n} \text{Var}(x)$

(chain mean is calculated from  $n$  samples)

this restores the scale to match  
the variance of individual samples  
(chart 6pt)

summary:  $\left( \begin{array}{l} W: \text{variance within chain} \\ B: \text{variance between chains scaled to} \\ \text{individual sample variance} \end{array} \right)$

not sure if this is correct? chart 6pt

$$\rightarrow \text{Var}(x) = E[\text{Var}(x|\text{chain})] + \text{Var}(E[x|\text{chain}]) \quad (\text{AHO TOTAL VARIANCE})$$

$$\rightarrow E[\text{Var}(x|\text{chain})] \quad \uparrow \text{expected within-chain variance}$$

$$\approx \frac{n-1}{n} W$$

$$\uparrow \text{variance of chain means}$$

$$\approx B/n$$

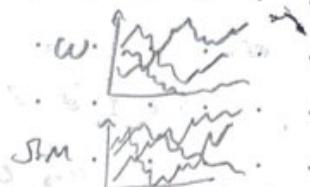
correcting for  
overdispersion,  
sample variance?

Today: continue MCMC - applications to SN cosmology,  
assessing convergence & mixing, comparison of algs.

(sheet 3)

### New case:

Now assume flat UNIV.  $\Omega_m = 1 - \Omega_m$   
but unknown  $M_{10}$ ,  $\sigma_8$ ,  $\Omega_m^2$ ,  $\Omega_m$ ,  $w$



Changing  $w \rightarrow v$ . small effect, so hard to constrain  $w$

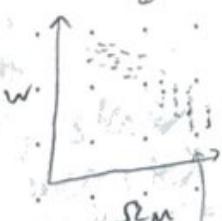
now instead use metropolis-within-gibbs 4x1D for each ~~each~~ <sup>sample w</sup> parameter (fix all other params to their current value and update one) tune 4x proposal scales & 4x acceptance ratios

(visually  $\rightarrow$  doesn't look like the best sampling on slides (not mixed much))

2D plot  $\rightarrow$  nice banded structure in samples: signaturing the algorithm due to rejection can use estimated cov. matrix (from initial MCMC or Laplace approx) for correlated proposal dist. in 4D metropolis!

$\rightarrow$  helps it know what direction to sample in

use this as proposal for 4D metropolis  $\rightarrow$  see much better movement/mixing, chains look more similar  $\rightarrow$  improvement!



better because we have tuned proposal jump dist to better match the true dat. so if ~~we~~ guides us to sample in the right direction (effectively explore sample space)

estimate cov matrix from samples from its initial run to use as our proposal dist. for new run metropolis

## How Many Iterations to get an $\approx$ independent sample? Autocorrelation Function

(slide)

For each scalar parameter  $\theta$

Chain:  $(\theta_1, \dots, \theta_N)$

Consider sample mean  $\bar{\theta} = \frac{1}{N} \sum_{i=1}^N \theta_i$

If chain were serially uncorrelated:  $\text{Cov}(\theta_i, \theta_j) = 0, (i \neq j)$

then  $\text{Var}(\bar{\theta}) = \text{Var}(\theta)/N$ .

However, chain is typically serially correlated:

$$\text{Cov}(\theta_i, \theta_j) = C_{|i-j|} = \text{Var}(\theta) p_{|i-j|}$$

Then for large  $N$ :  $\text{Var}(\bar{\theta}) = \text{Var}(\theta) \times \tau_N = \text{Var}(\theta)/N_{\text{eff}}$

where  $\tau = 1 + 2 \sum_{t=1}^{\infty} p_t$  is the AUTOCORRELATION TIMESCALE

basically tells us

how many more steps it's good to take before your next "indep" sample?

and  $N_{\text{eff}} = N/\tau$  is the EFFECTIVE SAMPLE SIZE  
(ESS) (equivalent number of independent samples).

$$\text{Cov}(\theta_i, \theta_j) = \text{Var}(\theta) p_{|i-j|} \quad (\text{remind: } \text{cov}(x, y) = \sqrt{\text{var}(x)} \sqrt{\text{var}(y)})$$

serially correlated: errors at one point in chain are related to errors at another point in chain. each value in chain depends on values of other values in chain. as a fn. of distance between those points in chain!

$\Rightarrow$  Find  $\text{Var}(\bar{\theta})$ : reminder: variance of sum is sum of covariances!

$$\begin{aligned} \text{Var}(\bar{\theta}) &= \text{Var}\left(\frac{1}{N} \sum_{i=1}^N \theta_i\right) = \frac{1}{N^2} \text{Var}\left(\sum_{i=1}^N \theta_i\right) = \frac{1}{N^2} \left( \sum_{i=1}^N \sum_{j=1}^N \text{Cov}(\theta_i, \theta_j) \right) \\ &= \frac{1}{N^2} \left( \sum_{i,j} \text{cov}(\theta_i, \theta_j) + \sum_{i=1}^N \text{Var}(\theta_i) \right) = \frac{1}{N^2} \left( \sum_{i=1}^N \sum_{j \neq i}^N \text{Var}(\theta) p_{|i-j|} + N \text{Var}(\theta) \right) \\ &= \frac{1}{N} \text{Var}(\theta) \left[ \sum_{i=1}^N \sum_{j=1}^N p_{|i-j|} + 1 \right] \end{aligned}$$

for each point  
in chain

generalization of distance: for infinitely long chain imagine going in each direction & adding up all distances: get factor of 2 going left vs right & factor of  $N$  since do this for each of  $N$  values in chain. (which is apparently infinite... or?)

$$\text{so for large } N: \text{Var}(\bar{\theta}) = \frac{1}{N} \text{Var}(\theta) \left[ 2 \sum_{t=1}^{\infty} p_t + 1 \right]$$

do this in practice:

## Estimating the Autocorrelation / ESS

For each scalar parameter  $\theta$

Sample covariance of lag  $t$ :

$$\hat{C}_t = \frac{1}{N-t-1} \sum_{i=1}^{N-t} (\theta_i - \bar{\theta})(\theta_{i+t} - \bar{\theta})$$

To estimate covariance of distance  $t$  apart and  $t+1$  apart

$C_0$  = sample variance of  $\theta$

$$\hat{p}_t = \frac{\hat{C}_t}{C_0}$$
 Sample autocorrelation of lag  $t$

$$\hat{T} = 1 + 2 \sum_{t=1}^T \hat{p}_t \quad \text{Estimated autocorrelation time}$$

Truncate at  $T$  lags s.t.  $\hat{p}_T \approx 0.1$

$N_{eff} = N/\hat{T}$  Effective number of independent samples

Slowest parameter is the limiting one!

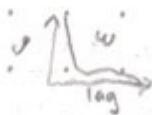
In practice, what matters is:

## Compute time / Effectively Independent Sample

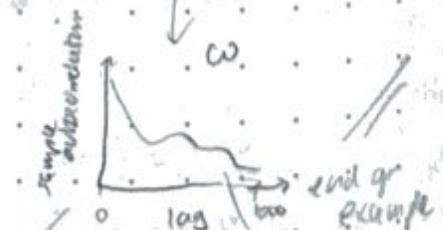
weigh up how fast algo is vs how many effective indep. samples  
i.e. how it takes to get 1 indep. sample! good to compare metric

See slides  $\rightarrow$  worst autocorr is  $w$ , compare to better sampling algo  $\rightarrow$  has larger effective sample number (faster rate q)

better sampling  $\rightarrow$  don't see banded structure in sampling.

  
try a diff. sampling method  
(1) 4D metropolis  
(2) v16x faster

better looking graph, faster rate q, indep. samples (larger  $N_{eff}$ )



for  $N=10,000$ ,

$\Rightarrow N_{eff} = 13$

not good!

$\Rightarrow 433$  to get one independent sample

still large p  
at distance four  
 $\rightarrow$  not good

can mix & match these algs: e.g.

## Mixed Samplers

(slides)

Target Posterior (intractable - cannot sample)  $\pi(\alpha, \beta, \theta)$

Suppose these

$$\pi(\alpha | \beta, \theta)$$

$$\pi(\beta | \alpha, \theta)$$

$$\pi(\alpha, \beta | \theta)$$

are tractable

but  $\pi(\theta | \alpha, \beta)$  is not tractable

### Mixed Sampler:

1) Sample  $\alpha^t \sim \pi(\alpha | \beta^{t-1}, \theta^{t-1})$

2) Sample  $\beta^t \sim \pi(\beta | \alpha^t, \theta^{t-1})$

3) Metropolis (proposal/acc./rej.) update  $\theta^t$

## Parameter Blocking

Target posterior (intractable - cannot sample):  $\pi(\kappa, \beta, \theta)$

As before suppose  $\pi(\theta | \alpha, \beta)$  not tractable.

(other conditionals  
tractable  
incl  $\pi(\alpha, \beta | \theta)$ )

### Mixed Sampler with Blocking: (more efficient)

1) Jointly sample  $\alpha^t, \beta^t \sim \pi(\alpha, \beta | \theta^{t-1})$

2) Metropolis (proposal/acc./rej.) update

usually more efficient

more shown how MCMC works  $\rightarrow$  now show why it works

(36/38) Run 10,000 chains:

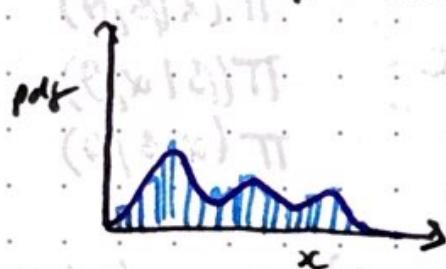
→ start all chains at some value

Histogram of initial values of chains ( $t=1$ )



histogram of ensemble of chains at time  $t$  (not single chain is used to evaluate it)

Histogram of chain values at step  $t=100$



### ERGODICITY:

Convergent distribution of chain ensemble at single time-slice = distribution-over-time of single chain  
= target (posterior) distribution (convergence to a stationary distribution)

Time average = Ensemble average

avg behavior over time is the same as behavior observed over many diff instances ('ensembles')

key q why we do MCMC  $\rightarrow$  want as fine grained sample to match true distribution, but we don't want to run loads of chains 100k times to reach that dist.  $\rightarrow$  we just want to run one chain & have it approximate the dist. at large times

(e.g.  $\pi(\theta)$  (symmetric)  $\Rightarrow$  unique stationary distribution?)

distribution  $(\alpha, \beta) \pi$  unique ergodic?

(ergodicity theorem): unique stationary distribution

( $\pi(\theta)$  and the ergodic limit (1. stationary (pd) and learning) ergodicism (2.

stationary (pd) and learning) ergodicism (2.

## A SKETCH OF MCMC THEORY

Markov chain: A sequence of random variables:

$$\theta^0, \theta^1, \dots, \theta^{t-1}, \theta^t \in \Theta$$

(statespace = parameter space)

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

Only dependent previous state, not all previous

Joint distribution:  $P(\theta^t, \theta^{t-1}, \dots, \theta^1, \theta^0)$

$$= P(\theta^t | \theta^{t-1}) \times P(\theta^{t-1} | \theta^{t-2}) \times \dots \times P(\theta^1 | \theta^0) P(\theta^0)$$

call  $\pi_0(\theta^0)$

(pdf of initial state)

$$= \pi_0(\theta^0) T_1(\theta^0 \rightarrow \theta^1) \times \dots \times T_{t-1}(\theta^{t-2} \rightarrow \theta^{t-1}) \times T_t(\theta^{t-1} \rightarrow \theta^t)$$

transition probabilities

$P(\theta^t | \theta^{t-1})$

If  $T_t(\theta^{t-1} \rightarrow \theta^t) = T(\theta^{t-1} \rightarrow \theta^t)$  not explicitly time-dependent  $\rightarrow$  "time homogeneous".

(e.g. prob transition state  $\theta^0$  to  $\theta^1$  is same whether at beginning or end of chain)

Marginal distribution:

$$P(\theta^t) = \int \pi_0(\theta^0) T(\theta^0 \rightarrow \theta^1) \times \dots \times T(\theta^{t-1} \rightarrow \theta^t) d\theta^0 \dots d\theta^{t-1}$$

Goal of MCMC: Want to engineer Markov chain s.t.

① Stationary distribution exists. ( $\Rightarrow$  unique)

① Limiting distribution

Stationary distribution  
aka invariant distribution

$$\pi_{\text{lim}}(\theta) = \lim_{t \rightarrow \infty} P(\theta^t = \theta)$$

$$= \pi_{\text{stat}}(\theta)$$

and it is equal to the target.

②  $\pi_{\text{lim}}(\theta) = \pi_{\text{stat}}(\theta)$  is equal to target posterior  
 $\pi_{\text{post}}(\theta)$

① system converges to unique stationary dist (?)  
② that dist is the target

## Good Properties for MCMC:

(board).  
can explore even  
space in good time

### 1) IRREDUCIBLE

MC can get to any state from any other state in a finite # of steps.

### 2) APERIODIC

[Periodic: MC has positive probability of returning to a state only with period  $k \geq 1$ ]

### 3) NOT TRANSIENT

[Transient: Non-zero probability of chain never returning to a particular state.]

$$\begin{aligned}
 & \text{Wikipedia: MCMC diag 3b} \\
 & \text{Markov process has a unique stationary dist - provided} \\
 & \quad \text{(implied by)} \\
 & \quad \text{① stationary dist exists (this is true b/c of condition of detailed balance)} \\
 & \quad \text{② stationary dist is unique. guaranteed by ergodicity which implies} \\
 & \quad \text{non-transient} \\
 & \quad \rightarrow \text{positive recurrent} \\
 & \text{positive recurrent} \\
 & \text{if we need irreducible? to converge?} \\
 & \text{if mixing dist is positive?} \\
 & \text{if } (\mathbb{Q})_{\text{mix}} \Pi = (\mathbb{Q})_{\text{mix}} \Pi \\
 & \text{if } (\mathbb{Q})_{\text{mix}} \Pi = (\mathbb{Q})_{\text{mix}} \Pi
 \end{aligned}$$

Lecture 2010.3.25

"Detailed balance": e.g. if oxygen molecule moves to other side of room another molecule will take its place & come of oxygen will stay the same, same identity dist. ~~independent of MCMC sample~~ will move around as chains sample diff. pos. but scroll overall.

[Today: finishing Bayesian computation (MCMC), soon: gaussian processes]

A sketch of MCMC theory cont.

(First, some jargon:)

**Positive Recurrent:** has finite expected time to return to each state.

**Ergodic:** Aperiodic & positive recurrent (time avg. = expectation over ensemble)

If we have:

Markov chain which is

- 1) Irreducible
  - 2) Aperiodic
  - 3) Not transient
- } MC converges to a unique, stationary distribution  $\pi_{\text{stat}}$  as the limiting  $\pi_{\text{lim}}(t \rightarrow \infty)$
- positive recurrent (chapt)  
(stronger condition)  $\Rightarrow$  not transient  
alone does not guarantee convergence

- We want to construct transition probabilities

s.t.  $\pi_{\text{stat}} = \pi_{\text{post}}$  (target)

- To do that, show that target distribution is invariant under transition probabilities, i.e. target is the unique stationary distribution.

## Stationary Distribution

Suppose  $\theta^{t-1} = a \sim \pi(\theta)$

suppose we start at position  
 $\theta^{t-1}$

what is probability next P(B)?  
 $\theta^t$

board)

but avg over all starting points a.

$$\int p(\theta^{t-1} = a) T(\theta^t = b | \theta^{t-1} = a) d\theta^{t-1}$$

$$\text{Next step: } P(\theta^t = b) = \int p(\theta^{t-1}) T(\theta^{t-1} = a \rightarrow \theta^t = b) d\theta^{t-1}$$

prob next point.  
is b regardless  
of history point

order of states  
from left  
to right  
(1) shift from  
(2) transition  
dist

iff  $P(\theta^t = b) = \pi(b)$  then  $\pi$  is the stationary distribution  $\pi_{\text{stat}}$

$\rightarrow$  (forget dist  $\pi$  is invariant under transition probabilities)

(transition probabilities do not induce any net change in pdf over all states)

## Detailed Balance:

sur stationary dist to exist

A sufficient condition - chain respects D.B. aka reversible, with  $\pi$  as stationary distribution

Reversible:  $P(\theta^{t-1} = a, \theta^t = b) = P(\theta^{t-1} = b, \theta^t = a)$

|| same may be true if  $\exists$  stationary dist  $\pi$  that satisfies D.B.  
process is reversible

Condition of D.B.:

ok  $\rightarrow$  production goes one way = prob it gives the other way

$$\pi(a) T(a \rightarrow b) = \pi(b) T(b \rightarrow a)$$

$\pi(a) T(b|a) = \pi(b) T(a|b)$  ) alternative notation

(flow of probability  $a \rightarrow b$  perfectly balanced by flow from  $b \rightarrow a$ )

$$\text{If } \theta^{t-1} = a \sim \pi, P(\theta^t = b) = \int \pi(a) T(a \rightarrow b) da$$

$$= \int \pi(b) T(b \rightarrow a) db$$

$$P(\theta^t = b) = \pi(b) \int P(a|b) da = \pi(b)$$

$\rightarrow \pi$  is the stationary distribution  $\pi_{\text{stat}}$ !

(board)

8 transition prob. to dist  $\Pi$  satisfy detailed balance there under condition  
1, 2, 3 if w.l.o.g. the stationary is  $\Pi_{\text{stat}}$ .  
Now we want to show under metropolis rule stationary detailed balance.

(SEE APPENDIX).

## Metropolis algorithm respects D.B. w/ $\Pi_{\text{stat}} = \Pi_{\text{post}}$

(shown that 1) transition prob. satisfies DB then stationary dist is  $\Pi_{\text{post}}$   
(show that our algo is constructed s.t. transition probabilities satisfying DB) is the stationary dist

Suppose w.l.o.g. we pick points  $a, b \in \Theta$  s.t.  
 $\Pi_{\text{post}}(a) < \Pi_{\text{post}}(b)$ , and  $a, b \in \Pi_{\text{post}}(\emptyset)$ , and  
proposal distribution  $J(b|a) = J(a|b)$  (symmetric). ( $r > 1$ ).

Forward:

(backward expression is  
symmetric in  $a, b$ )  
even though we never  
try to move from  $b$  to  $a$ ,  
this is true in general

prob make  
transition

prob propose  
transition

$$P(\theta^{t-1} = a, \theta^t = b) = \Pi_{\text{post}}(a) T(a \rightarrow b) = \Pi_{\text{post}}(a) J(b|a) \times \min\left(\frac{\Pi_{\text{post}}(b)}{\Pi_{\text{post}}(a)}, 1\right)$$

Backward:

$$\begin{aligned} P(\theta^{t-1} = b, \theta^t = a) &= \Pi_{\text{post}}(b) T(b \rightarrow a) \\ &= \Pi_{\text{post}}(b) \times J(a|b) \times \frac{\Pi_{\text{post}}(a)}{\Pi_{\text{post}}(b)} \\ &= J(a|b) \Pi_{\text{post}}(a) = J(b|a) \Pi_{\text{post}}(a) \end{aligned}$$

Forward = Backward  $\rightarrow$  transition probabilities

(reminder: just spring paper about Metropolis Hastings)  
respect D.B. w/  $\Pi_{\text{stat}} = \Pi_{\text{post}}$

(Exercise:)

(see appendix!)

① Show DB also holds for Metropolis Hastings w/ asymmetric proposal  $J(a|b) \neq J(b|a)$

② Gibbs Sampling: suppose  $\Pi_{\text{post}}(\emptyset, \emptyset)$  is the posterior over two parameters  $\emptyset, \emptyset$ . the conditionals  $\Pi(\emptyset|\emptyset)$  and  $\Pi(\emptyset|\emptyset)$  are tractable. By considering a full cycle of the Gibbs sampler, show that  $\Pi_{\text{post}}(\emptyset, \emptyset)$  is the stationary dist.

show that one step of a Gibbs sampler is equivalent to a M-H step where the acceptance ratio is always 1.

## MCMC in Practice

(slides)

posterior mode

MAP

(peaks)

1. Find the model(s) using optimisation, can use Laplace approx. to obtain a proposal cov. matrix  
alternative to replace it  
not tractable posterior  
tractable posterior tools  
very easily of course  
it is to do initial run  
& estimate cov matrix  
from sample as proposal  
dist per new run!
2. Begin multiple (4-8 parallel) chains at starting positions dispersed around the model(s)  
for Laplace approx  
need expression  
for posterior can  
distribute (not necessarily  
normalised?)  
can see particular chain  
gets stuck
3. Scale Metropolis proposals to tune 25-50% acceptance rate (depending on dimensionality of jump).
4. Use proposal cov. matrix that reflects the shape of the posterior.  
e.g. bimodal posterior  
example 678/100
5. After run, look at chains (if possible) to check for obvious problems.  
should look like noise
6. Compute Gelman-Rubin ratio comparing within-chain-variance to between-chain-variance to check that chains are well-mixed (should be very close to 1), and assess burn-in.
7. Compute autocorrelation timescale and effective sample size to make sure you have enough independent samples for inference.  
not good if say 3000 not mixed but with say 100 not mixed
8. If all checks out, remove burn-in, thin, and combine chains to compute posterior inferences.

## Overview of MCMC

### \* MCMC Algorithms

- Metropolis / M-H algorithms
- Gibbs Sampling
- Metropolis-within-Gibbs

most other algos are just special cases of M-H  
didn't cover modern MCMC use gradient (autodiff) - Hybrid / HMC / NUTS

### \* Assessing / Comparing performance of MCMC algs.

- Gelman-Rubin statistics for comparing mixing of multiple chains
- Autocorrelation time - how long to get an independent sample?
- Effective sample size - how many indep. samples do I have?

### \* Detailed balance & theoretical considerations

end of MCMC!

is next time gaussian processes

GAUSSIAN PROCESSES

(slide)

What is a Gaussian Process?

- \* A GP is a (possibly infinite) collection of RVs  $\{g_t\}$ , (typically indexed w/ some ordering in time, space or wavelength), such that any finite subset of RVs have a jointly multivariate gaussian dist.
- \* Any vector  $g = \{g_t : t=1, \dots, N\}$  of a finite subset is multivariate Gaussian, therefore it is completely described by a mean  $E[g]$  and cov. matrix  $\text{Var}[g] = \text{cov}[g, g^T]$ .
- \* Elements of the cov. matrix are determined by a function of the coordinates e.g.  $\text{cov}[g_t, g_{t'}] = k(t, t')$  called the covariance function or kernel.
- \* A G.P. w/ mean function  $m(t)$  and kernel  $k(t, t')$  is denoted
 
$$g(t) \sim GP(m(t), k(t, t'))$$
- \* A G.P. provides a distribution over functions

the dist of a GP  
is the joint of all  
these (ininitely many)  
variables i.e. it is a dist  
over fns with a continuous  
domain (i.e. time or space)

e.g.  $g = \{g_1, g_2\} \sim N\left(\begin{bmatrix} m(1) \\ m(2) \end{bmatrix}, \begin{bmatrix} k(1,1) & k(1,2) \\ k(2,1) & k(2,2) \end{bmatrix}\right)$

since gaussian, GP is completely defined  
by its mean  $m(t)$  and covariance  
 $(m, k(t, t'))$

Review: Properties of multivariate Gaussians

Full prob. density:

$$N(\mathbf{g}|\mu, \Sigma) = [\det(2\pi\Sigma)]^{-\frac{1}{2}} \exp(-\frac{1}{2}(\mathbf{g}-\mu)^T \Sigma^{-1} (\mathbf{g}-\mu)) \quad (\Sigma \text{ positive definite})$$

Joint dist. of components:

$$\mathbf{g} = \begin{pmatrix} \mathbf{U} \\ \mathbf{V} \end{pmatrix} \sim N\left(\begin{bmatrix} \mathbf{U}_0 \\ \mathbf{V}_0 \end{bmatrix}, \begin{bmatrix} \Sigma_U & \Sigma_{UV} \\ \Sigma_{UV} & \Sigma_V \end{bmatrix}\right)$$

If you observe/know/condition on  $\mathbf{V}$ :Conditional dist.  $\mathbf{U}|\mathbf{V} \sim N(\mathbb{E}[\mathbf{U}|\mathbf{V}], \text{Var}[\mathbf{U}|\mathbf{V}])$ Conditional mean  $\mathbb{E}[\mathbf{U}|\mathbf{V}] = \mathbf{U}_0 + \Sigma_{UV}\Sigma_V^{-1}(\mathbf{V} - \mathbf{V}_0)$ Conditional variance  $\text{Var}[\mathbf{U}|\mathbf{V}] = \Sigma_U - \Sigma_{UV}\Sigma_V^{-1}\Sigma_{VU}$ If  $\mathbf{V}$  = observed data,  $\mathbf{U}$  = unobserved parameters, then  $P(\mathbf{U}|\mathbf{V})$  is a posterior pdf!

reminder, marginals:  $P(\mathbf{V}) = \int P(\mathbf{U}, \mathbf{V}) d\mathbf{U} = N(\mathbf{V}|\mathbf{U}_v, \Sigma_v)$   
 $P(\mathbf{U}) = N(\mathbf{U}|\mathbf{U}_u, \Sigma_u)$

(simply just drop irrelevant variables from mean vector & cov matrix! e.g.  $\mathbf{M} = \begin{bmatrix} \mathbf{M}_v \\ \mathbf{M}_u \end{bmatrix} \rightarrow \mathbf{M}_u$ )  
 $\Sigma = \begin{bmatrix} \Sigma_U & \Sigma_{UV} \\ \Sigma_{UV} & \Sigma_V \end{bmatrix} \rightarrow \Sigma_U$

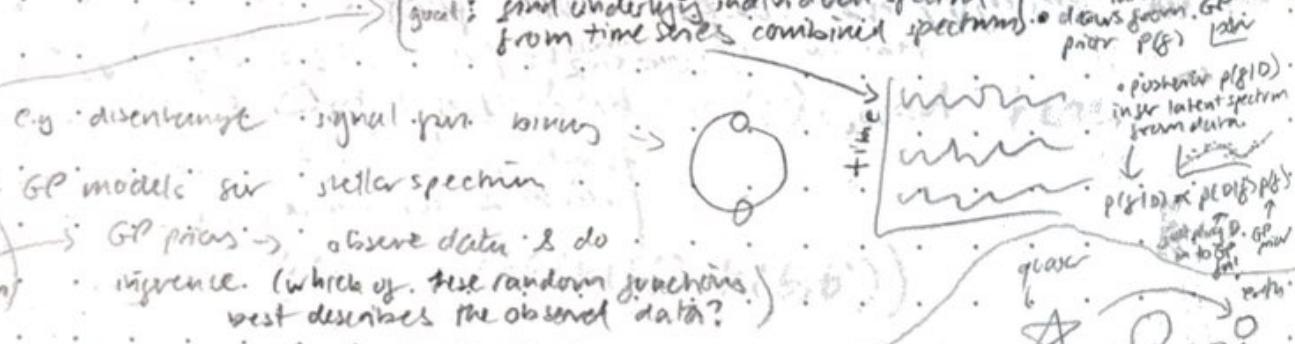
new conditionals come from  $p(\mathbf{U}|\mathbf{V}) = \frac{p(\mathbf{U}, \mathbf{V})}{p(\mathbf{V})}$ 

(don't know g. i need to know prior for this?)

# What are GPs used for in Astrophysics?

- Some physical models are very nearly gaussian  
e.g. CMB is GP on the sphere) or w/ gaussian  
(e.g. damped random Walks (quasar light))  
curves are approximately gaussian
- "Nonparametric" models: flexible sn.s to use when an accurate parametric astrophysically-motivated fn. is not available or is impossible.
- Nonparametric = no. of parameters / latent variables grows with the dataset.
- Interpolation / Emulation: to generate a smooth curve going through some observation or simulation points
- Correlated noise/error model: when you marginalise out the latent function, you are effectively accounting for noise/gluciations correlated over time/space/wavelength  $\rightarrow$  "nuisance function"

ex!



ex. 2

- lensed quasar: brightness randomly fluctuates in time
- but if, behind galaxy, it will be lensed & light split into two copies: fluctuations in one image will appear much later than in the other image
- use GP to describe underlying pr. so we can predict time delays
- $\hookrightarrow$  use GP to model underlying latent (true) light curve, as a damped random walk
- supervisive can also be gravitationally lensed  $\rightarrow$  get 4 copy's of the same SN!
- each copy is magnified by some value & also time shift between copies
- model light curves realising GP, likelihood programs to estimate  $\beta$ ,  $a$ ?

e.g. fit GP to SN light curve time series

## G.P. as a prior on functions (slides)

But we only ever evaluate the fn. on a finite set of points!

A grid of times:  $t = (t^1, \dots, t^i, \dots, t^n)^T$

A vector of fn. values evaluated on the time grid:

$$g = (g(t^1), \dots, g(t^i), \dots, g(t^n))^T$$

Assume a squared exponential kernel / cov. fn.

$$\text{Cov}[g(t), g(t')] = k(t, t') = A^2 \exp(-|t - t'|^2/\tau^2)$$

Assume a const. prior mean fn.:

$$\mathbb{E}[g(t)] = m(t) = c \quad (\text{often assume zero-}c \text{ mean})$$

Prior on function:  $P(g | A, \tau) = N(g | 1c, K)$

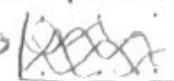
Drawing from prior:  $g | A, \tau \sim N(1c, K)$

Cov. matrix  $K$  is populated by evaluating the kernel:

$$\text{Cov}[g(t), g(t')] = k(t, t') = A^2 \exp(-|t - t'|^2/\tau^2)$$

For all pairs of points in  $t$ :

$$K_{ij} = k(t_i, t_j) = A^2 \exp(-|t_i - t_j|^2/\tau^2)$$

Non-stationary kernel, set  $A, \tau$  to particular values to generate say 20 random obs. from GP prior  $\rightarrow$  

(notice how pris change if we change  $\tau, A$ !)

hyperparameters  
kernel have  
significant influence  
on distribution

(Sides)

## Fitting a GP to data

### ① PREDICTION

If we knew the characteristic scales of the kernel (hyperparameters) ( $A, \tau^2$ ), then how do we fit the data at observed times to find the curve for unobserved times?

### ② MODEL SELECTION

Given the observed data, how do we fit for the characteristic scales of the kernel (hyperparameters)?

## Posterior inference w/ G.P.s, Estimating the underlying curve:

$f_0 = f_n$ . at observed times  $t_0$  (training set)

$f_* = f_n$ . at unobserved times  $t_*$  (prediction set)

Joint:  $\begin{pmatrix} f_0 \\ f_* \end{pmatrix} \sim N \left( \begin{bmatrix} 1_c \\ 1_c \end{bmatrix}, \begin{bmatrix} K(t_0, t_0) & K(t_0, t_*) \\ K(t_*, t_0) & K(t_*, t_*) \end{bmatrix} \right)$

Populating the cov. matrix:

$K(t, t')$  has  $i, j$ -th entry =  $k(t_i, t_j)$

Using the assumed kernel fn.

$$\text{Cov}[f(t), f(t')] = k(t, t') = A^2 \exp(-|t - t'|^2 / \tau^2)$$

now we have m.v. vectors  
partition

Posterior (conditional on observations, is also Gaussian):

$$f^*|y_0 \sim N(\mathbb{E}[f^*|y_0], \text{Var}[f^*|y_0])$$

Posterior predictive mean:

$$\mathbb{E}[f^*|y_0] = 1c + K(t_*, t_0)K(t_0, t_0)^{-1}(y_0 - 1c)$$

Posterior predictive variance / covariance:

$$\text{Var}[f^*|y_0] = K(t_*, t_*) - K(t_*, t_0)K(t_0, t_0)^{-1}K(t_0, t_*)$$

ignore measurement error  $\rightarrow$  bad

GP tries to fit by going through centre  
of every single data point

including measurement error: good

(~~measurement~~)

how do we account for this?

### Accounting for measurement error

$$y_0|f_0 \sim N(f_0, W)$$

"Model of  
measurement error  
as adding  
gaussian process  
to  
gaussian process"

$y_0$  are measured values of  $f_0$  at time  $t_0$

at each observation  $y_i|f_0^i \sim N(f_0^i, \sigma_i^2)$

$W$  is measurement cov. matrix.

Most common case: heteroskedastic uncorrelated measurement error:

$$\text{Cov}(e_i, e_j) \equiv W_{ij} = \delta_{ij}\sigma_i^2$$

Measurement error model:

data = latent value + meas. error

$$y_0 = f_0 + e$$

$$e \sim N(0, W)$$

(mean-zero gaussian noise)

W is cov matrix W

[Derivation as the sum of two GPs at the observed times: Probabilistic Generative Model]

### GP of intrinsic curve

$$f(t) \sim GP(m(t) = c, k(t, t'))$$

$f_0$  = latent fn. at observed times  $t_0$

GP Prior:  $f_0 \sim N(1c, K(t_0, t_0))$

### GP of measurement error:

$y_0$  = measurements (with error) of  $f_0$  at times  $t_0$

$$y_0 | f_0 \sim N(f_0, w)$$

Same as:

$$y_0 = f_0 + \epsilon \quad (\text{mean-zero error})$$

$$\epsilon \sim N(0, w)$$

GOAL: Need Joint Dist. of data and latent fn.

we know  
want to  
infer true curve  
over unobserved times

$$(y_0) \sim N\left(\begin{bmatrix} ? \\ ? \end{bmatrix}, \begin{bmatrix} ? & ? \\ ? & ? \end{bmatrix}\right)$$

Then can calc fn. (posterior) prediction at unobserved points

$$f_{*} | y_0 \sim N(E[f_* | y_0], \text{Var}[f_* | y_0])$$

Using conditional properties of MV Gaussian as before.

$$\begin{bmatrix} y_0 \\ f_* \end{bmatrix} \sim N\left(\begin{bmatrix} E[y_0] \\ E[f_*] \end{bmatrix}, \begin{bmatrix} \text{cov}[y_0, y_0] & \text{cov}[y_0, f_*] \\ \text{cov}[f_*, y_0] & \text{cov}[f_*, f_*] \end{bmatrix}\right)$$

just substitute  $y_0 = f_0 + \epsilon$  &  $\text{cov}[\epsilon, \epsilon] = 0$

Intrinsic/Latent Process  $f_0 \sim N(1c, k(t_0, t_0))$

Measurement Process  $y_0 | f_0 \sim N(f_0, w)$

$$y_0 = f_0 + \epsilon \quad \epsilon \sim N(0, w) \quad (\text{meas. error})$$

To find cov. submatrices, apply bilinearity of covariance

$$\text{Cov}(y_0, y_0) = \text{Cov}(f_0, f_0) + \text{Cov}(\epsilon, \epsilon) + 2\text{Cov}(f_0, \epsilon)$$

$$\text{Cov}(f_0, f_0) = k(t_0, t_0) \quad (\text{GP of intrinsic curve})$$

$$\text{Cov}(\epsilon, \epsilon) = w \quad (\text{measurement noise})$$

(The two processes are indep.  $\rightarrow$  uncorrelated)

$$2\text{Cov}(f_0, \epsilon) = 0$$

$$\therefore \text{Cov}[y_0, y_0] = k(t_0, t_0) + w$$

Use joint of latent  $g$  at observed and unobserved times

$$\begin{pmatrix} f_0 \\ f_{t_*} \end{pmatrix} \sim N \left( \begin{bmatrix} 1c \\ 1c \end{bmatrix}, \begin{bmatrix} k(t_0, t_0) & k(t_0, t_*) \\ k(t_*, t_0) & k(t_*, t_*) \end{bmatrix} \right)$$

To derive similar arguments for

$$\text{Cov}[y_0, f_{t_*}] = \text{Cov}[f_0, f_{t_*}] + \text{Cov}[\epsilon, f_{t_*}] = k(t_0, t_*) + 0$$

$$\text{Cov}[f_{t_*}, f_{t_*}] = k(t_*, t_*)$$

remember:  
 $k$  is a matrix

Fill out cov. matrix:

$$\begin{pmatrix} y_0 \\ f_{t_*} \end{pmatrix} \sim N \left( \begin{bmatrix} 1c \\ 1c \end{bmatrix}, \begin{bmatrix} k(t_0, t_0) + w & k(t_0, t_*) \\ k(t_*, t_0) & k(t_*, t_*) \end{bmatrix} \right)$$

now can compute fn. predictions at unobserved points:

Posterior Predictive (conditional on data  $y_0$ )

$$f_*|y_0 \sim N(\mathbb{E}[f_*|y_0], \text{Var}[f_*|y_0])$$

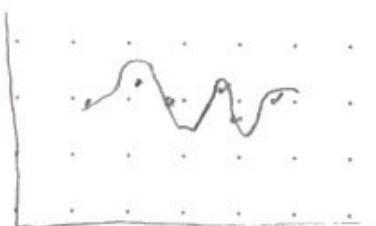
Using MV Gaussian conditional properties:

Posterior predictive mean

$$\mathbb{E}[g_*|y_0] = 1c + K(t_*, t_0)[K(t_0, t_0) + W]^{-1}(y_0 - 1c)$$

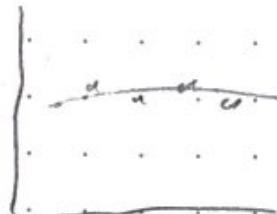
Posterior predictive variance/covariance

$$\text{Var}[f_*|y_0] = K(t_*, t_*) - K(t_*, t_0)[K(t_0, t_0) + W]^{-1}K(t, t_*)$$



ignore meas error  
overfitting

(see slides)



include meas  
error

→ looks more  
like an  
actual  
simpler  
curve

Now we know how to do ① PREDICTION,  
How do we do ② MODEL SELECTION?

## Tuning hyperparameters ( $A, \tau^2$ )

Recall our probabilistic generative model:

Intrinsic/latent GP process

(make explicit dependence on hyperparameters)

$$f(t) \sim GP(m(t) = c, K_{A, \tau^2}(t, t'))$$

$$\text{kernel: } K_{A, \tau^2}(t, t') = A^2 \exp(-|t - t'|^2 / \tau^2)$$

$$f_0 \sim N(1c, K_{A, \tau^2}(t_0, t_0))$$

$$\sim N(0, w)$$

Gaussian measurement process  $y_0 = f_0 + \epsilon$

$$\begin{aligned} & y \sim N(\mu_y, \Sigma_y), \quad \mu_y = \mu_x + \mu_f \\ & \text{then } \begin{bmatrix} x \\ y \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_y \end{bmatrix}\right) \end{aligned}$$

Integrating out the latent fn.  $f(t)$  (or use addition of MVN R.V.s) gives us the marginal likelihood:

$$\int p(y_0 | f_0, A, \tau^2) df_0 = p(y_0 | g_0; A, \tau^2) p(f_0 | A, \tau^2) df_0$$

$$\begin{aligned} p(y_0 | A, \tau^2) &= \int p(y_0 | f_0) \times p(f_0 | A, \tau^2) df_0 \\ &\text{meas. error} \quad \text{latent GP} \quad \begin{aligned} & \int p(y_0 | f_0) df_0 \\ & \text{joint normal} = \int N(y_0 | 1c, K_{A, \tau^2}(t_0, t_0)) df_0 \\ & \text{marginal} = N(y_0 | 1c, K_{A, \tau^2}(t_0, t_0)) \end{aligned} \\ &= \int N[y_0 | g_0, w] \times N[f_0 | 1c, K(t_0, t_0)] df_0 \end{aligned}$$

$$L(A, \tau^2) = p(y_0 | A, \tau^2) = N[y_0 | 1c, K_{A, \tau^2}(t_0, t_0) + w]$$

- $w$  = meas. error covariance
- $K_{A, \tau^2}$  = GP covariance

Which we can optimize (max likelihood) or specify a prior on  $(A, \tau)$  and sample from posterior.

results for this example on slides

(1) find for this case, it's not actually that sensitive to hyperparameters

remember we had 4 images shifted by a time delay each has indep. noise, how do we deal w/ this using GP?

## Case study: construct Bayesian Model

One latent fn. modelled as a draw from a GP

$$g(t) \sim GP(c, k(t, t')) \quad k(t, t') = A^2 \exp[-(t-t')^2/\tau^2]$$

Data of each image are measurements of  $(\Delta t_i, \Delta m_i)$ -shifted copies of this latent function

For image  $S_1$  w/ observations  $y_i$  at times  $t_i$ :

$$y_i(t_i) = g(t_i) + \epsilon_i, \quad \epsilon_{i,j} \sim N(0, \sigma_{i,j}^2)$$

$j$  indexes the observations at times  $t_{i,j}$

Images  $S \# i = 2, 3, 4$  w/ observations  $y_i$  at times  $t_i$  w/ (unknown) time delays and magnitude shifts (relative to  $S_1$ ):  $(\Delta t_i, \Delta m_i)$

$$i=2,3,4 \quad y_i(t_i) = \Delta m_i + g(t_i - \Delta t_i) + \epsilon_i \quad \epsilon_{i,j} \sim N(0, \sigma_{i,j}^2)$$

Sample Posterior:  $P(\Delta m, \Delta t, A, \tau^2 | D)$

slide results for sample w/ metropolis-within-gibbs  
time priors for sampling as needed

new  $\Delta t$  (left),  $\Delta m$  (right) = final best estimate  
for each time delay.

Shifting by estimated  $\Delta t, \Delta m$ , datasets line up pretty well as hoped!

## Lecture 22

14.3.25

gelman BDA ch 21, 5

[Today: Finishing Gaussian processes, starting  
hierarchical Bayes/ Probabilistic Graphical Models]

### Gaussian Processes cont.

Previously: used "squared exponential kernel"

$$\rightarrow \text{Cov}[g(t), g(t')] = k(t, t') = A^2 \exp(-|t - t'|^2/\tau^2)$$

$$g|A, \tau \sim N(1c, K)$$

(characteristic amplitude)

characteristic timescale of fluctuations

this kernel gives very smooth curves (infinitely differentiable)

False  
has  
properties:  
stationary to stationary

Stationary:  $k(t, t')$  time translation invariant to  
 $t \rightarrow t+k, t' \rightarrow t'+k$

Symmetric:  $k(t, t') = k(t', t)$

### OTHER COV FUNCTIONS:

"Ornstein Uhlenbeck Process" (damped random walk)  
Exponential covariance function

$$k(t, t') = A^2 \exp(-|t - t'|/\tau)$$

Long-run dist. of stochastic differential eq.

mean-reversion decay term random walk

$$dy(t) = \tau^{-1} [\mu - y(t)] dt + \sigma dW_t$$

long-term mean

$$A = \tau \sigma^2 / 2$$

mean-reversion timescale

wavy

this kernel is everywhere continuous but not differentiable

more jagged  $\Rightarrow$  smooth curve squared  
 $\Rightarrow$  want more jagged  $\Rightarrow$  this depends on what you want for your use case

### "Matern class" of cov functions

somewhere between smooth & jagged

$V = \frac{1}{2}$  special case of Matern kernel:

$$K_{\text{Matern}}(r) = \frac{2^{1-V}}{\Gamma(V)} \left(\frac{\sqrt{2V}r}{L}\right)^V K_V\left(\frac{\sqrt{2V}r}{L}\right)$$

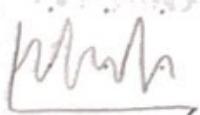
with free parameters  $V, L$ , where  $K_V$  is a modified Bessel fn.

When  $V = p + \frac{1}{2}$  is half-integer; exponential  $\times$   $p$ -polynomial

$$k_V = p + \frac{1}{2}(r) = \exp\left(-\frac{\sqrt{2V}r}{L}\right) \frac{\Gamma(p+1)}{\Gamma(2p+1)} \sum_{i=0}^p \frac{(p+i)!}{i!(p-i)!} \left(\frac{\sqrt{2V}r}{L}\right)^{p-i}$$

### "Periodic cov. functions"

$$\text{eg. } k(t, t') = A^2 \exp\left(-\frac{2r^2}{C^2} \sin^2(\pi(t-t')/T)\right)$$



$$(A^2 \sin^2(\pi(t-t')/T)) Q D^2 A = (A^2 D^2)$$

end of

$$Q D^2 A + A D^2 (-A)^{-1} J^{-1} J = (A) D^2$$

$$Q D^2 = A$$

# PROBABILISTIC GRAPHICAL MODELS

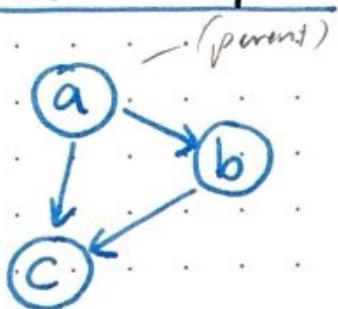
dependencies  
independence

(Depicting complex probability distributions.)

Suppose  $a, b, c$  are r.v.s with a joint pdf:

$$\begin{aligned} P(a, b, c) &= P(c|a, b)P(a, b) \\ &= P(c|a, b)P(b|a)P(a) \end{aligned} \quad \begin{matrix} \text{(general)} \\ \text{factorisation} \end{matrix}$$

Directed Acyclic Graph:



nodes are connected  
by one-way edges that  
do not form any cycles

DAG

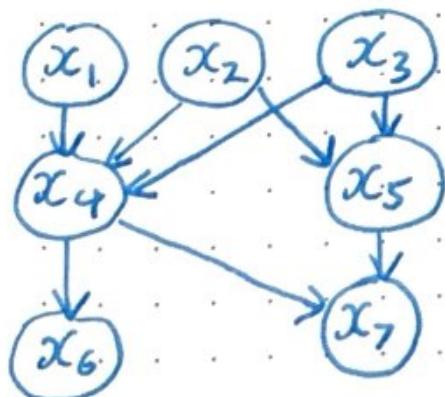
why we do  
digraphs:  
everything is Bayesian  
stochastic or not,  
or r.v.s can have  
a more structured joint  
posterior, etc.

Generally, if  $\vec{x}$  is a collection of  $k$  r.v.s, then given a DAG, the factorisation is written as

$$P(\vec{x}) = \prod_{k=1}^k P(x_k | \text{Parents of } x_k)$$

if  $x_k$  has no parents then just  $P(x_k)$

E.g.



$$\begin{aligned} P(x_1, \dots, x_7) &= P(x_1)P(x_2)P(x_3) \\ &\times P(x_4|x_1, x_2, x_3) \\ &\times P(x_5|x_3)P(x_6|x_4) \\ &\times P(x_7|x_4, x_5) \end{aligned}$$

we call this  
CONDITIONAL  
INDEPENDENCE

(the lack of connections implies some model structure!)

fact it isn't fully connected implies conditional independence (structure)

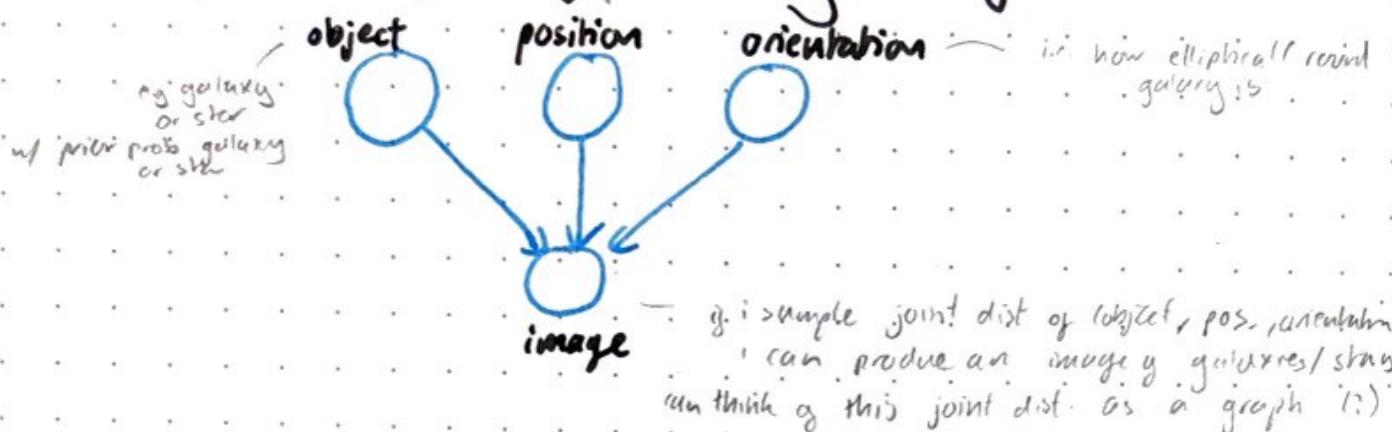
maybe you want to simulate an image with lots of galaxies in it

## Generative Models

(slides)

monday

### Causal process for generating images

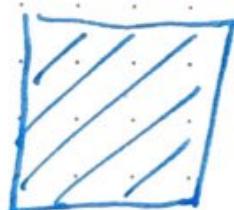


## Notation

(board)



open node = latent parameter / unobserved data



shaded node = observed parameter or data (conditioned on)



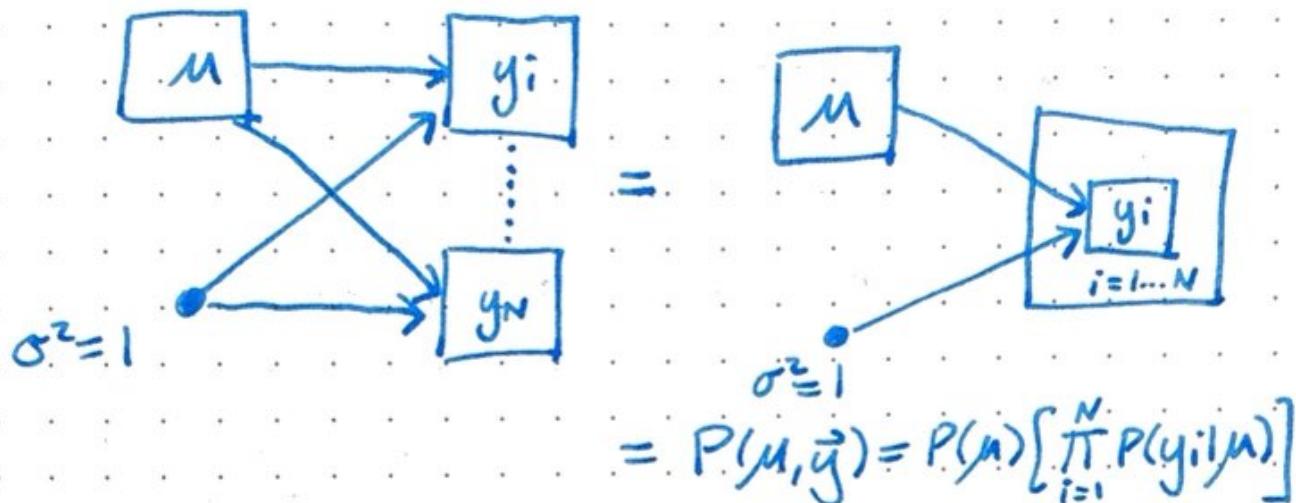
filled dot = fixed and known constant



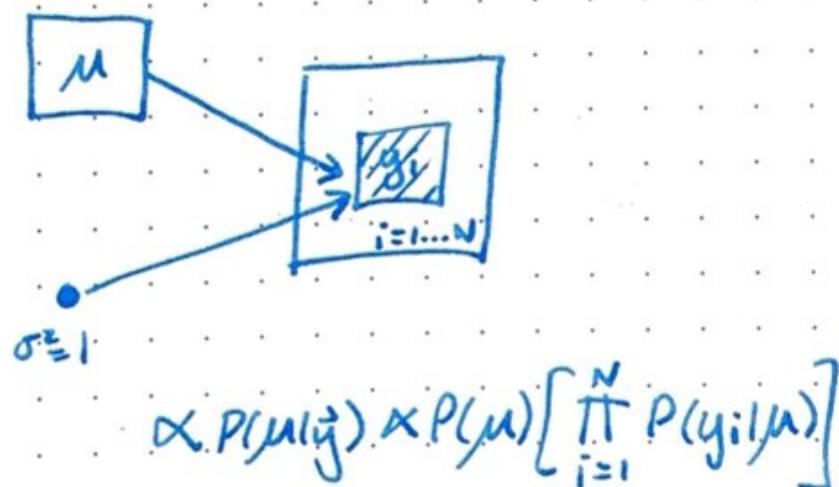
plate:  $N$  independent replications of what's inside

Eg.

Gaussian i.i.d data  $y_i \stackrel{iid}{\sim} N(\mu, \sigma^2 = 1)$   
 $i = 1, \dots, N$

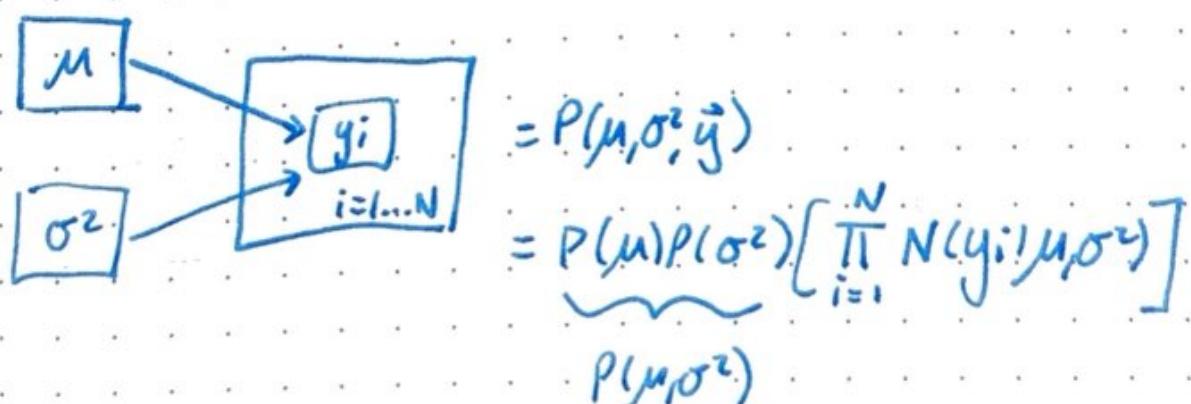


If  $y$  observed,  
 (shaded)



Eg.

Both  $\mu, \sigma^2$  unknown



# Conditional Independence and Probabilistic Graphical Models

Let  $a, b, c$  be r.v.s

(Marginal) Independence

$$P(a, b) = P(a)P(b) \rightarrow a \perp b$$

$$\rightarrow a \perp b | \emptyset$$

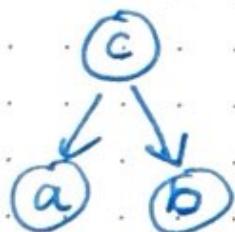
(not conditioning on anything)

Conditional Independence

$$P(a, b | c) = P(a | c)P(b | c) \rightarrow a \perp b | c$$

i.e.  $P(a | b, c) = P(a | c)$  ← if i know then  
b gives me no more info about a

Eg.

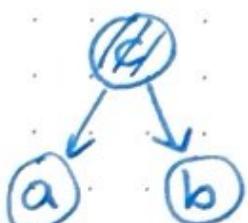


$$P(a, b, c) = P(c)P(a | c)P(b | c)$$

$$P(a, b) = \int P(c)P(a | c)P(b | c) dc \neq P(a)P(b)$$

$$\rightarrow a \not\perp b | \emptyset$$

"a is not indep. b marginally"

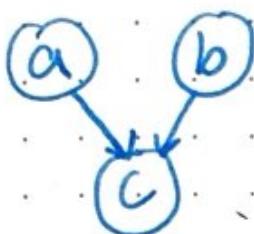


$$P(a, b | c) = \frac{P(a, b, c)}{P(c)} = P(a | c)P(b | c)$$

$$\rightarrow a \perp b | c$$

∴ a is indep. b  
conditioned on c

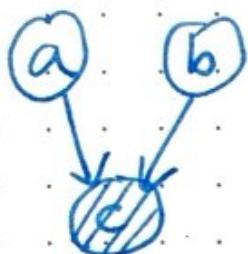
Eg



$$P(a, b, c) = P(a)P(b)P(c|a, b)$$

$$P(a, b) = P(a)P(b) \int P(c|a, b)P(c)dc \\ = P(a)P(b) \rightarrow a \perp b | \emptyset$$

"marginal independence between a & b"



$$P(a, b | c) = \frac{P(a, b, c)}{P(c)} \neq P(a|c)P(b|c) \\ \rightarrow a \not\perp b | c$$

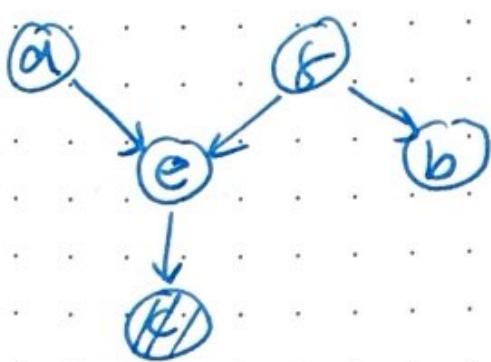
i.e.  
by observing c, the  
children a & b, must  
influence a, b are  
now correlated

"a is no longer indep. of b  
if you condition on c"

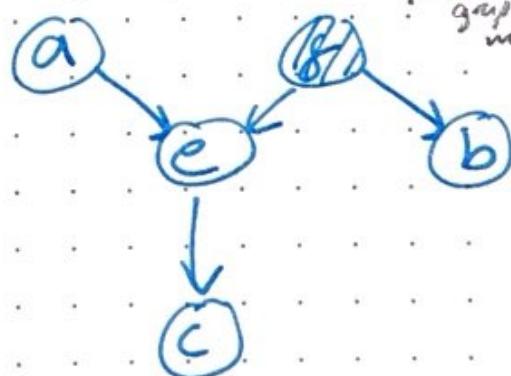
## D-separation Example

(stars)

want to test whether nodes are  
indep given other nodes  
read off independence properties from  
probabilistic graphical models



$$a \not\perp b | c$$



$$a \perp b | c$$

path a  $\rightarrow$  b is  
not blocked  
so  $a \not\perp b | c$

path a  $\rightarrow$  b  
blocked  
so  
 $a \perp b | c$

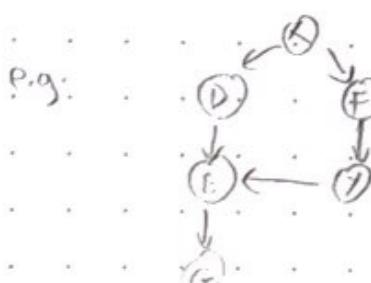
## D-separation

(slides)

- \* A, B and C are non-intersecting subsets of nodes in a directed graph
- \* A path from A to B is blocked if it contains a node such that either
  - a) The arrows on the path meet either head-to-tail or tail-to-tail at the node, and the node is in the set C, or
  - b) The arrows meet head-to-head at the node, and neither the node, nor any of its descendants are in the set C.
- \* If all paths from A to B are blocked, A is said to be d-separated from B by C.
- \* If A is d-separated from B by C, the joint over all variables in the graph satisfies  $A \perp B | C$

example → see slides

(very brief, didn't draw on my slide)

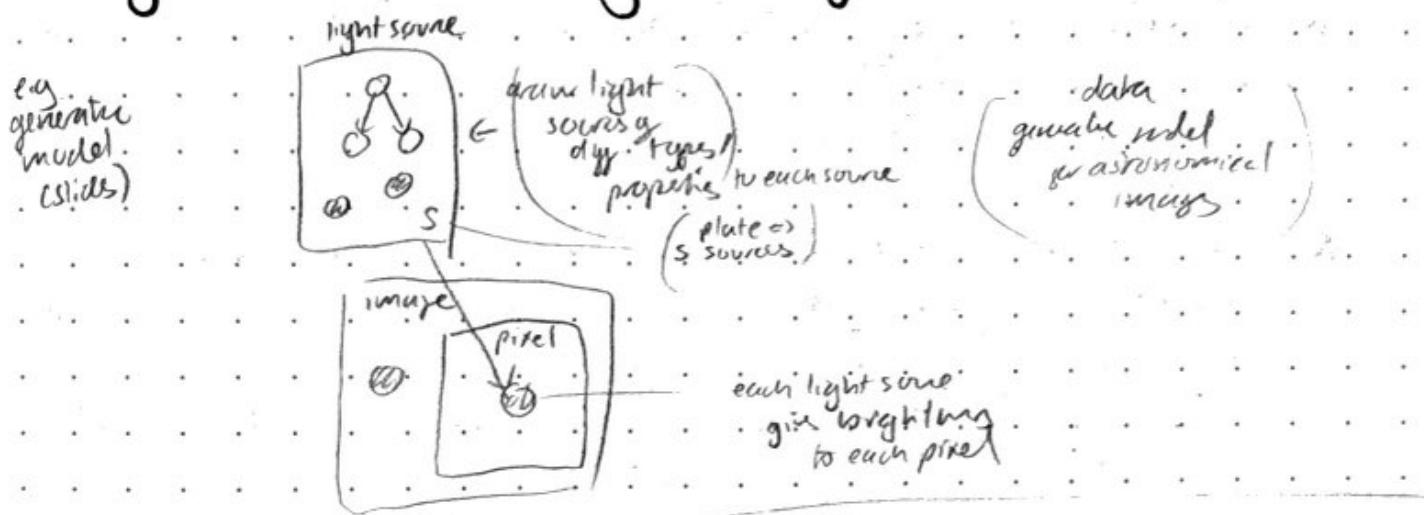


$X \perp Y$   
 $X \perp Y \text{ IF}$   
 $X \perp Y \text{ IF}, G$

way to think about statistical model  
 data generating process:  
 how do I generate samples from parameters?  
 first step, this is what our graph can describe  
 → strong start w DGP  
 gelman BDA (Why we like models? why they are so useful, fundamental to our inference process.)  
 hierarchical → joint prior  
 inference

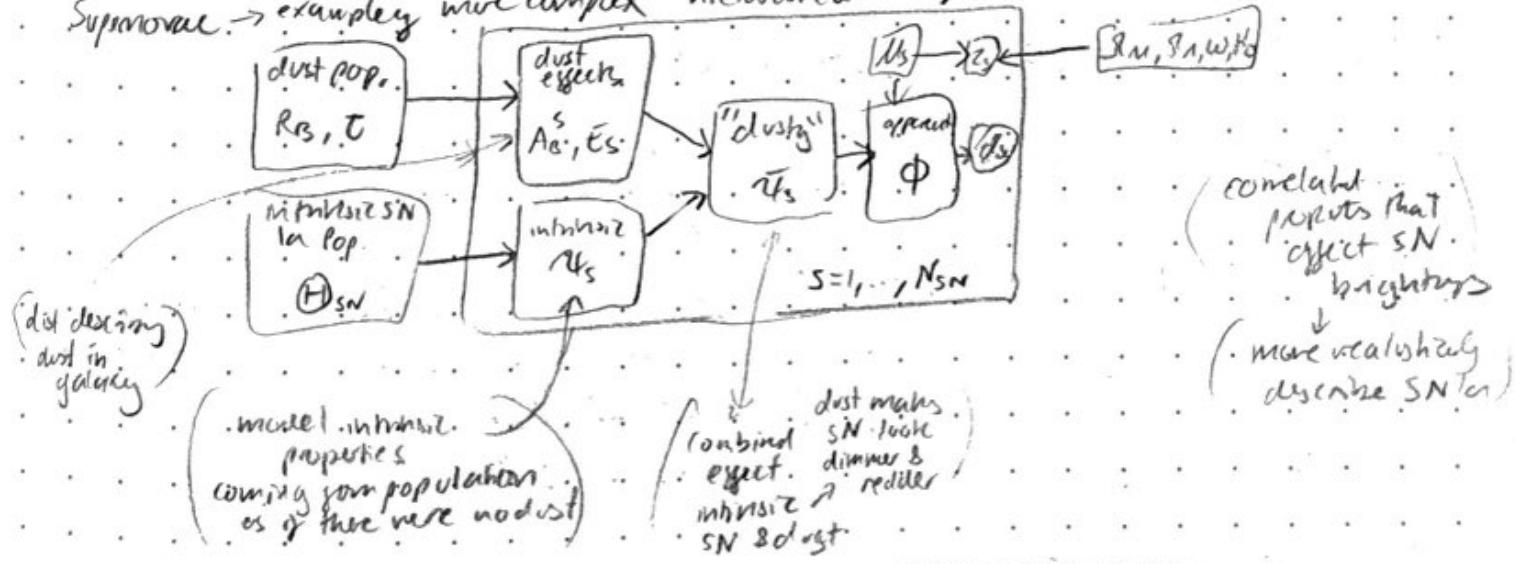
17.3.25

## Today: Hierarchical Bayes/ Bayesian Model Selection



data  
 generative model  
 astronomical images

'PGM linked to hierarchical bayesian models  
 (all g. research is basically hierarchical bayes)  
 Supernovae → example of more complex hierarchical bayesian model



model spectral energy distribution: SED = model for brightness (flux) of sn vs time & wavelength

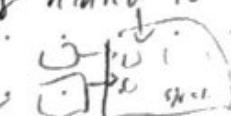
Now diff SN can have diff time series as a fn of wavelength - trace it to baryons in the source or a fn of time



slide 4 derives 1 SN, slide 5 deriving N SN (plate)

of params.

stuff outside plate describes post distribution / cosmological parameters (things not related to individual SN) and stuff inside plate describes stuff related to individual supernovae



while drawing example: big complicated hierarchical Bayesian model.

(Not quite important)  
see slides if you are interested

## Common Problems in Astronomy (slides)

- \* Want to learn about a population of objects from a finite sample of individuals, each measured with error
- \* Observed data is actually a combination of uncertain astrophysical & instrumental & selection effects. Need to model them to infer the "intrinsic" properties of the object or population of objects ("deconvolve")

## What is Hierarchical Bayes? (slides)

"simple" Bayes:  $D|O \sim \text{Model}(O)$   $\boxed{\theta} \rightarrow \boxed{D|O}$

Posterior (Bayes' theorem):  $P(O|D) \propto P(D|O)P(O)$

Hierarchical Bayes:  
 $O_i$ : parameter of individual  
 $\alpha, \beta$ : hyperparameters of population  
 $D_i|O_i \sim \text{Model}(O_i)$   
 $O_i|\alpha, \beta \sim \text{Pop Model}(\alpha, \beta)$

Joint posterior:

$$P(\epsilon O; \beta, \alpha | \epsilon D; \beta) \propto \left[ \prod_{i=1}^N P(D_i|O_i) P(O_i|\alpha, \beta) \right] P(\alpha, \beta)$$

build up complexity by layering conditional probabilities

diff hierarchies e.g. population level parameters (pop. dist. of host galaxies) <sup>dist</sup>  
obj params. object level parameters e.g. property of individual SN

Forward model:

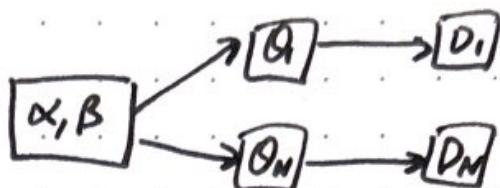
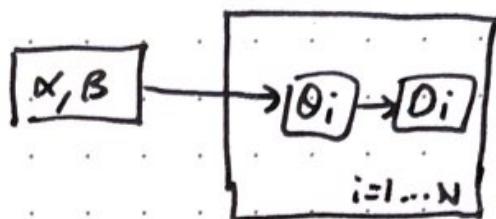


Plate notation:  
(loop over individuals  
in sample)



## Advantages of Hierarchical Bayesian Models

- Common problem in astronomy: infer properties of population from finite sample of individuals with noisy measurements
- Incorporate multiple sources of randomness & uncertainty as "latent variables" with distributions underlying the data
- Express structured probability models adapted to data-generating process ("forward model")
- Bayesian: Full (non-gaussian) probability distribution = global, coherent quantification of uncertainties
- Completely explore & marginalize posterior trade-offs / degeneracies between parameters / hyperparameters

# Simplest Hierarchical Bayesian / Multi-level Model: "Normal-Normal" for standard candle Mag.

$s=1, \dots, N$

**Level 1:** Pop. distribution  
of latent variables  
(absolute mags.)  
"population dist/prior"

$$M_s \sim N(M_0, \tau^2)$$

latent variables

hyperparams  
(pop. means  
variance)

**Level 2:** Measurement  
error process

"measurement  
likelihood"

$$D_s | M_s \sim N(M_s, \sigma_s^2)$$

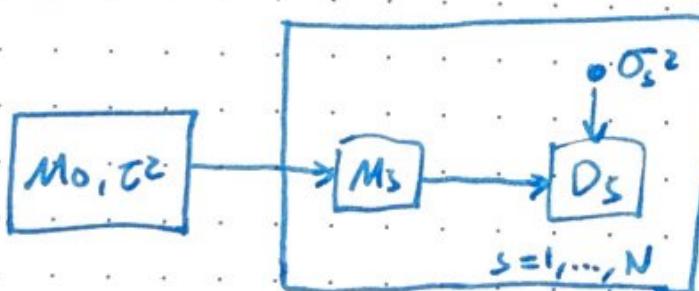
measurements  
(data)

heteroskedastic  
meas. error  
variance (known)

Joint probability density of data, latent variables,  
hyperparameters

"backwardly  
PGM for  
this joint"

$$P(\{D_s\}, \{M_s\}, H = \{M_0, \tau^2\})$$



joint factors  
into conditional  
and marginal pdfs  
based on model  
assumptions

Joint probability density of all the things: data,  
latent variables, hyperparameters

$$P(\{D_s\}, \{M_s\}, H) = \left[ \prod_{s=1}^N P(D_s | M_s) P(M_s | M_0, \tau^2) \right] \times P(H)$$

measurement  
likelihood

population  
distribution/  
prior

hyperprior

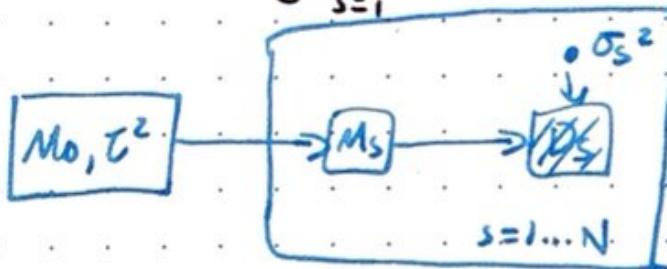
## Putting the Bayesian in Hierarchical Bayesian

Joint posterior of all unknowns given the data

$$P(\{\underline{M}_s\}, H | \underline{D}_s) = \frac{P(\{\underline{M}_s\}, H, \underline{D}_s)}{P(\underline{D}_s)} \quad \begin{matrix} \text{ignorable} \\ \text{normalization const.} \\ \text{"ignore until next time"} \end{matrix}$$

(posterior on  $N+2$  dim. parameter space)

$$P(\{\underline{M}_s\}, H | \underline{D}_s) \propto \left[ \prod_{s=1}^N P(D_s | M_s) P(M_s | M_0, \tau^2) \right] \times P(H)$$



## Hierarchical vs Regular Bayes

\* Could regard as just a general Bayesian inference problem in a very high dim. parameter space, e.g.

$$\Theta = \{\underline{M}_1, \dots, \underline{M}_N, M_0, \tau^2\} = \{\underline{M}; M_0, \tau^2\}$$

$$P(\underline{\Theta} | D) \propto P(D | \underline{\Theta}) P(\underline{\Theta})$$

$$P(\underline{\Theta} | D) \propto P(D | \underline{M}) P(\underline{M} | M_0, \tau^2) P(M_0, \tau^2)$$

$$P(\underline{\Theta} | D) \propto \left[ \prod_{s=1}^N P(D_s | M_s) P(M_s | M_0, \tau^2) \right] P(M_0, \tau^2)$$

helpful in  
bayesian computation  
if you can take  
advantage of  
the structure of  
the problem

- \* However, special hierarchical structure is useful for modelling, estimation & computation
- \* For large  $N$ , wouldn't want to do  $N+2$  dim. Metropolis MCMC!

# Gibbs Sampling & Hierarchical Bayes

can develop a ~~algorithm~~ for  
gibbs jumping that uses that starts  
in an efficient way

Utilises conditional independence structure of PGM/posterior to derive conditional posterior densities

$$P(\{\mathbf{E}M_s\}, H | \{\mathbf{D}_s\}) \propto \left[ \prod_{s=1}^N P(D_s | M_s) P(M_s | M_0, \tau^2) \right] \times P(H)$$

Gibbs: use  $P(\mathbf{E}M_s | H, \{\mathbf{D}_s\}) \propto P(H | \mathbf{E}M_s, \{\mathbf{D}_s\})$  instead to sample  
(need to plug in initial guesses first)

- ① For  $s = 1, \dots, N$ : Sample latent variables conditional on data and hyperparameters

$$P(M_s | H, D_s) \propto P(D_s | M_s) \times P(M_s | M_0, \tau^2)$$

(conditional independence)      (gaussian)      (gaussian)

- ② Sample hyperparameters from conditional on data and latent variables

$$P(M_0, \tau^2 | \{\mathbf{E}M_s\}; \{\mathbf{D}_s\}) = P(M_0, \tau^2 | \{\mathbf{E}M_s\})$$

(conditional independence)

Hyperparameters indep of  $\{\mathbf{D}_s\}$   
given  $\{\mathbf{E}M_s\}$  - can read  
from graph!

$$= P(M_0 | \tau^2, \{\mathbf{E}M_s\}) P(\tau^2 | \{\mathbf{E}M_s\})$$

(gaussian)      (inv- $\chi^2$ )

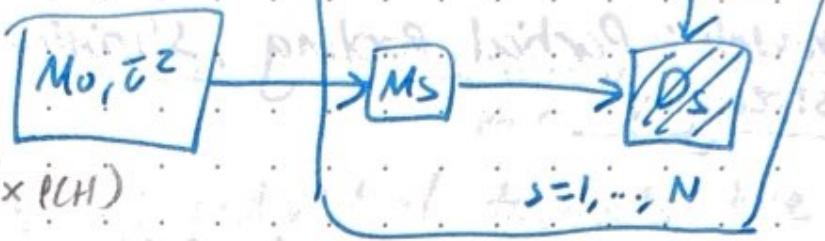
Reduces to the familiar posterior for unknown mean and variance of gaussian data (sheet 2 prob. 1 BDA 3.2.3.3)

- fix data (condition on D always)
- input some guess for  $M_0, \tau^2$  and condition on that, sample latent variables
- turns out that if we condition on H,  $M_0$  of SN 1 is indep to  $M_0$  SN 2, 3 etc.
- ① so simplifies to two terms rather than  $N+2$  terms which depend on  $s$   
so can loop through each  $s$  which is easy to sample given b/c gaussian
- ② turn out we read conditioning relationships of graph  
we don't have to worry about data, only these  $N$  factors but those look like gaussian which nicely factors into a gaussian &  $\text{inv-}\chi^2$   $\rightarrow$  much easier to sample from  
(w unknown  $M_0$ ?)

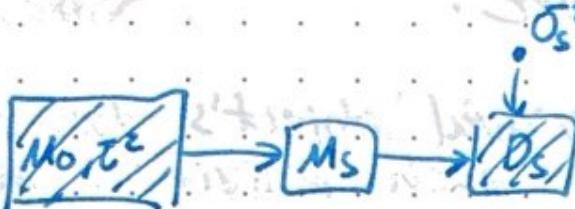
Full posterior:

$$P(\mathbf{EM}_S, \mathbf{H} | \mathbf{ED}_S)$$

$$\propto \left[ \prod_{s=1}^N P(D_s | M_s) P(M_s | M_0, \tau^2) \right] \times P(\mathbf{H})$$

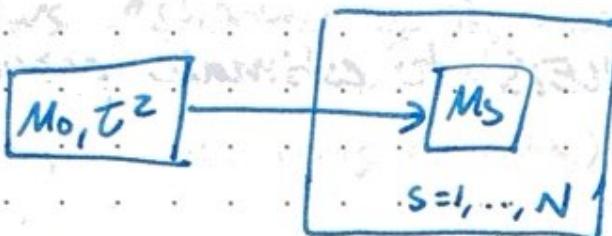


①



$$P(M_s | M_0, \tau^2, \mathbf{ED}_S) \propto P(D_s | M_s) \propto N(D_s | M_s, \sigma_s^2) \propto N(M_s | M_0, \tau^2)$$

②



$$P(M_0, \tau^2 | \mathbf{EM}_S, \mathbf{ED}_S) = P(M_0, \tau^2 | \mathbf{EM}_S) = P(M_0 | \tau^2, \mathbf{EM}_S) P(\tau^2 | \mathbf{EM}_S)$$

gaussian prior

(similar to sheet 2 problem 1)

"think about this on your own"

Hyperprior  $P(M_0, \tau^2) \propto 1$

Draw from  $\tau^2 | \mathbf{EM}_S \sim \text{Inv-}\chi^2(N-3, \frac{(N-1)}{(N-3)} S^2)$

$$M_0 | \tau^2, \mathbf{EM}_S \sim N(\bar{M}, \tau^2/N)$$

$$\bar{M} = \frac{1}{N} \sum_{s=1}^N M_s$$

$$S^2 = \frac{1}{N-1} \sum_{s=1}^N (M_s - \bar{M})^2$$

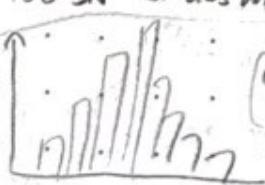
e.g. 100 SN w/ abs mag as distinct measurements

$$D_S = \hat{m}_S - \mu_S$$

see slide

{run male histogram  
of obs data Ds}

data with same error



obs abs mag  $D_S$

per individual SN calc  
posterior abs mag

observed

gibbs sampling sample  $M_s | M_0, \tau^2$   
(in 102 dim, marginalized so we have 1 dim).



pop mean  $M_0$ ,  $\tau^2$  width

Marginal posterior estimate

Scatterplot  
FIT

pooling evidence  
close to data  
by doing hierarchical bayes,  
the posterior estimate  
clustering truth insipred  
from original data contains  
measurement error (large spread than true)

individual estimates have  
intrinsic spread from intrinsic  
errors (have learnt to  
determine out measurement  
noise)

"partial"  
pooling phenomenon

measurement error (large spread than true)

# HB Models: Partial Pooling, Shrinkage & "Borrowing of strength"

HB models implement partial pooling  $\rightarrow$  strengthen

(3/3)

each  $i$  has objective  $D_i$   
individual prior  
first. Once we have

## \* Common sense procedure:

- Analyse each individual object's data  $D_i$  separately and get each individual MLE estimate (with error)
- Plug in all  $\hat{M}_{MLE,i}$ 's to estimate population hyperparameters

\* **PROBLEM**: Each individual  $D_i$  estimate may be unbiased but collectively give a biased estimate of population (e.g. variance) because of errors

\* **SOLUTION**: Use HB to model  $N$  individuals & population simultaneously and get better estimates of both

notes to self: (attempt to develop intuition).

$$\text{HB full posterior} \rightarrow P(M_s | M_0, \tau, D) \propto P(D|M_s)P(M_s | M_0, \tau^2)$$

individual estimate  $\uparrow$  pop dist acts.  
as prior

say we want to infer  $M_s$  after and take peak as our best estimate of  $M_s$ :

- \* individual estimate just finds whatever  $M_s$  maximises prob. of observing data (MLE)
- \* but HB model still posterior incorporates prior on  $M_s$  so e.g.  $M_s$  may maximise  $P(D|M_s)$  but if  $M_s$  is very unlikely this will sway our estimate. e.g.  $M_s$  through the prior  $\rightarrow$  balance between maximizing prob. of data but also probability of intrinsic abs may  $M_s$  existing in the population in the first place

How is this diff to MAP?

$\rightarrow$  point is HB gives  $P(D|M_s)$ , so we can infer those hyper params instead of merely fit those them ourselves (MAP)

(slides)

"often not unbiased but their biases helps you in away"

## Shrinkage Estimators

Closer to truth on avg than  
considering individual estimate  
for each individual in isolation

- \* Bias estimator of individual towards the population
- \* Leads to overall lower MSE than individual unbiased estimators
- \* Allows "sharing of information" between individuals to improve overall estimation
  - example: sticks (brief explanation)

## Shrinkage with Hierarchical Model

Return to our previous example:

$$\begin{array}{lll} \text{level 1: pop. dist. of latent variables} & M_s \sim N(M_0, \tau^2) & \text{"pop. dist." / "prior"} \\ \text{Level 2: measurement error process} & D_s | M_s \sim N(M_s, \sigma^2) & \text{"measurement likelihood"} \\ & & \text{(caused by shrinkage to form likelihood } P(D_s|M_s)) \end{array}$$

Individual MLE of SN's alone:  $\hat{M}_s = D_s$  (unbiased)

Population dist. acts as prior

$$P(M_s | M_0, \tau^2, D_s) \propto P(D_s | M_s) P(M_s | M_0, \tau^2)$$

and pulls posterior estimate of individual closer to population mean estimate

Pull/shrinkage controlled by population variance  $\tau^2$

This comes from full posterior:

$$P(M_s, M_0, \tau^2 | D_s) \propto P(H) \left[ \prod_{s=1}^N P(D_s | M_s) P(M_s | M_0, \tau^2) \right]$$

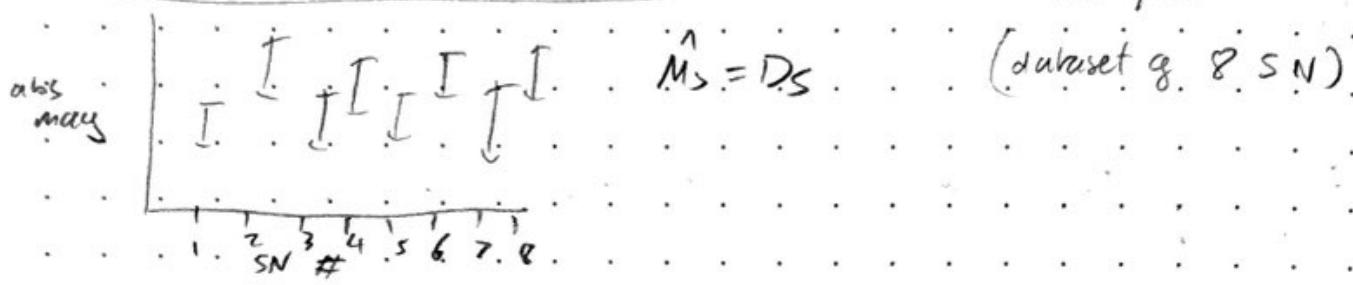
just ignore terms not in  $M_s$  like in simple geometry principle

$$P(M_s | M_0, \tau^2, D_s) \propto P(D_s | M_s) P(M_s | M_0, \tau^2)$$

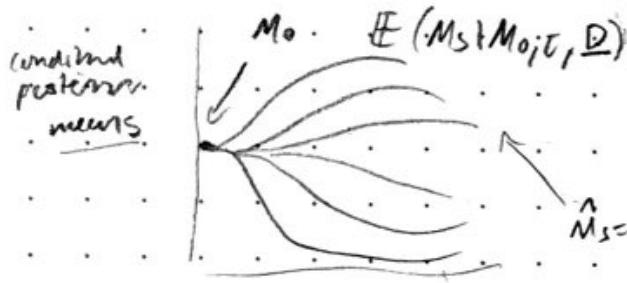
T controls how much we shrink our poster estimate towards pop. mean estimate

How does HB implement shrinkage?

example



If you know  $M_0, \tau$  of individual  $M_s$  toward it depends on  $\tau$ .



if we can take  $\tau$  to zero; converge to  $M_0$

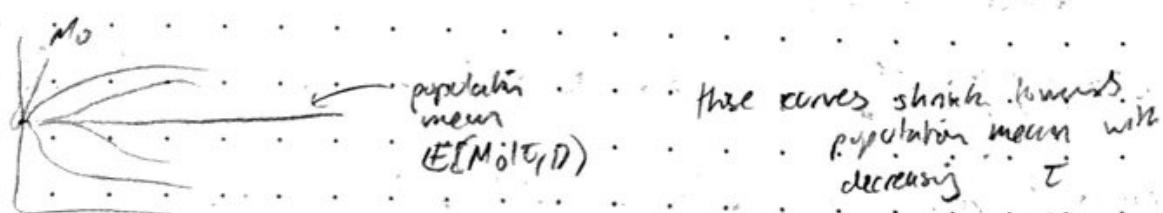
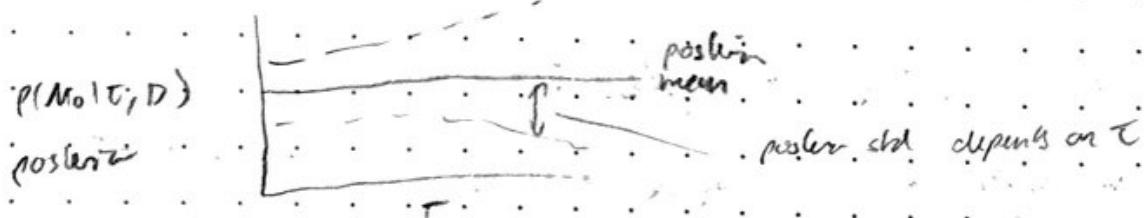
$\tau$  controls degree of shrinkage toward  $M_0$   
(increase  $\tau$ ; wider range of  $M_s$  probable  
so more spread)

$\tau$  (one track for each SN,  
posterior mean as most)

but we don't know  $M_0 \rightarrow$  how to estimate that from marginal posterior

$$p(M_0, \tau | D) = p(M_0 | \tau, D) p(\tau | D)$$

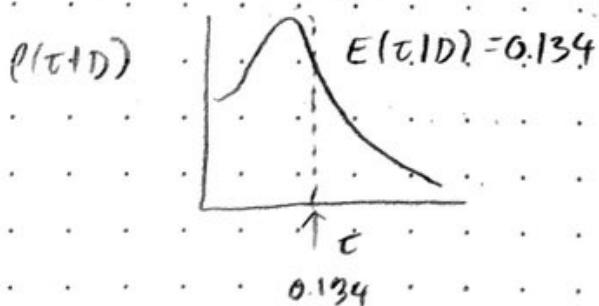
marginal out  
individual  $M_s$



these curves shrink toward  
population mean with  
decreasing  $\tau$

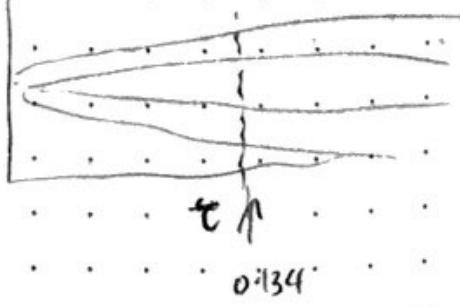
hierarchical bayes gives the  
posterior density  $p(\tau | D)$   $\rightarrow$  way to estimate  $\tau$

(don't necessarily have to take  
point estimate of  $\tau$ , can do  
estimates of  $\tau$ )



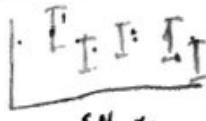
how can work out best trade-off  
between shrinking toward same pop. value vs  
allowing them to be independently estimated  
 $\Rightarrow$  this is at about  $\tau = 0.134$

$$E[M_s | M_0, \tau, D]$$



on avg HB estimate  
closer to truth than  
individual MLE

abs max



## Lecture 24

recommended reading (will's book!)

→ Mackay "Information theory  
inference & learning algorithms" ch. 28

19.3.25

## Today: Bayesian Model selection

### Model comparison & selection

Mackay:  
"probabilistic inference"  
1st level: assume model true  
→ fit that model to  
the data ( $P(D|M)$ )  
2nd level  
→ to  
compare  
models

- No. of spectral lines in a (noisy) spectrum?
- Clustering / mixture models - how many clumps?
- Time series - curve fitting
  - Is there a trend?
  - Complexity / order / degree of best model
  - which GP kernel best explains the data?
- Cosmology - standard (8-parameter) cosmological model vs more exotic (more parameters) models?
  - $\Lambda$  in standard model ( $\Lambda$ CDM)  
but people looking for evidence of more exotic varieties of the model! (to try and figure out dark energy)  
but then we need to introduce more parameters  
• which extra parameters are warranted by data  
adding more params will always give better fit but how do we know if adding extra params is just fitting to noise or actually warranted by data? → hope to answer this in bayesian model comparison

## Bayesian Model Comparison

(Odds?)

Parameter estimation: posterior on parameters

$$P(\theta_1 | D, M_1) = \frac{P(D | \theta_1, M_1) P(\theta_1 | M_1)}{P(D | M_1)}$$

Model selection: posterior on models

$$P(D | M_1) = \int P(D | \theta_1, M_1) P(\theta_1 | M_1) d\theta_1$$

needs proper prior!

evidence or  
marginal likelihood

between two competing models

$$\frac{P(M_1 | D)}{P(M_2 | D)} = \underbrace{\frac{P(D | M_1)}{P(D | M_2)}}_{BF} \times \underbrace{\frac{P(M_1)}{P(M_2)}}_{\substack{\text{Prior odds} \\ \text{or} \\ \text{improper}}}$$

prior probability given  
model 1 vs model 2  
two iscientifically  
based on if you're  
agnostic can set prior  
equal probabilities

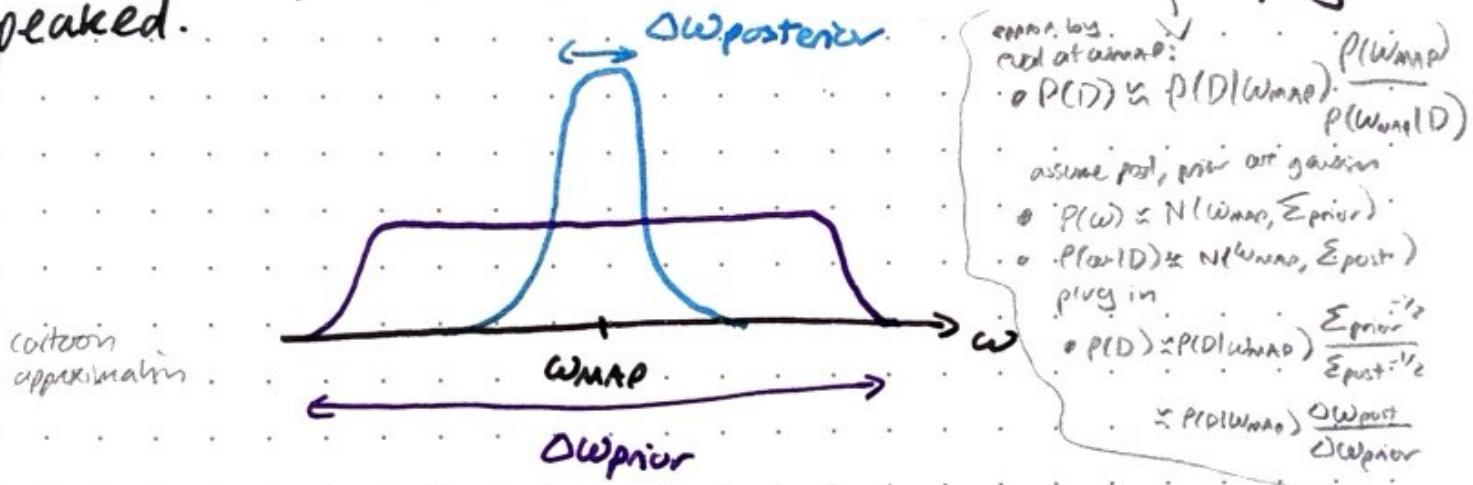
$$\text{Bayes Factor}_{12} = \frac{P(D | M_1)}{P(D | M_2)}$$

(let's develop some intuition)

For a given model with a single parameter,  $w$ , consider the approximation.

$$p(D) = \int P(D|w)p(w)dw \approx P(D|w_{MAP}) \frac{\Delta w_{posterior}}{\Delta w_{prior}}$$

Where the posterior is assumed to be sharply peaked.



Taking logarithms, we obtain

$$\ln P(D) \approx \ln P(D|w_{MAP}) + \ln \left( \frac{\Delta w_{posterior}}{\Delta w_{prior}} \right)$$

negative

since we hope  
prior width  
always  
larger than  
posterior width

generalise:

with  $M$  parameters, all assumed to have the same ratio  $\Delta w_{posterior}/\Delta w_{prior}$ , we get

$$\ln P(D) \approx \ln P(D|w_{MAP}) + M \ln \left( \frac{\Delta w_{posterior}}{\Delta w_{prior}} \right)$$

negative

term that is  
negative, proportional  
to no of parameters  
& relatively  
unstable

(OCCAM)  
FACTOR

gives simplest  
model that  
explains the  
data

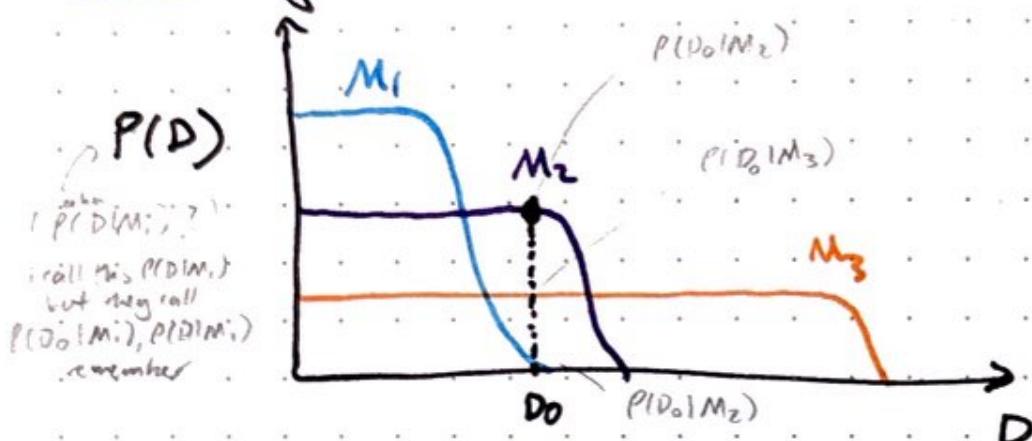
but this term  
penalizes you  
for adding more  
parameters

it involves adding  
more parameters will  
always increase likelihood of data  
but model is not this  
more likely is not always  
warranted because we could just be  
overfitting to noise. Occam's  
factor (2nd term) penalizes  
overfitting

this term  
increases w  
posterior width

want to find sweet spot  
enough parameters without overfitting 147

## Matching data and model complexity



"more complex models can predict a greater variety of data sets"

→ M<sub>1</sub> can only predict limited amount of data (e.g.: linear model)

→ M<sub>2</sub> can predict slightly more complex data sets (e.g. might be a quadratic model)

→ M<sub>3</sub> e.g. if a neural network may be able to predict arbitrarily complex functions

but remember:

$p(D)$  is itself a distribution over potential data sets

the more complex datasets a model can predict, the more spread out  $p(D)$  will be. If it is normalized  $\Rightarrow$  for any particular dataset it will have a lower probability overall

→ so simpler models will have more limited range of data sets they can predict, but because of that, overall  $p(D)$  is higher

→ whereas complex models predict greater range but  $p(D)$  lower b/c of that normalisation factor (pales any given one)

so this gives a good way to weigh fit of models to dataset vs their overall complexity

e.g. for dataset D<sub>0</sub>, M<sub>1</sub> is v. low probability. M<sub>1</sub> deserves its dataset. M<sub>3</sub> slightly more probable but M<sub>2</sub> has best balance between complexity of model (more complex than M<sub>1</sub>) but also not too complex (M<sub>2</sub>)

recommend read ch 28 of Mackay

## Marginalising Joint Distribution of Parameters and Data

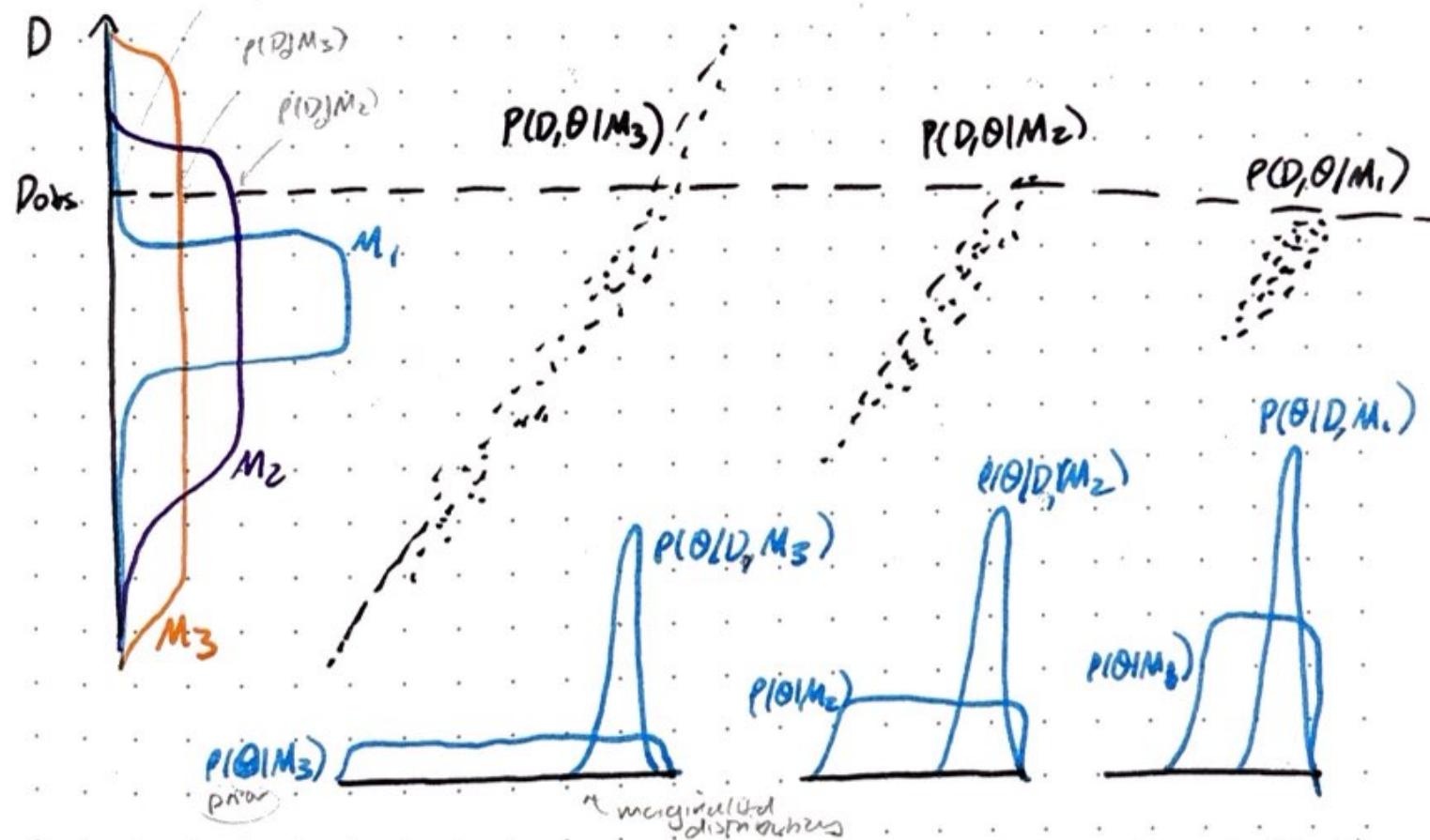
hypothesis  $H$

$$(w \sim P(w|H) \text{ & } D \sim P(D|w, H) \rightarrow (D, w) \sim P(D, w|H))$$

I'm going to call this model,  $M$ .

$$\theta \sim P(\theta|M) \text{ & } D \sim P(D|\theta, M) \rightarrow (D, \theta) \sim P(D, \theta|M)$$

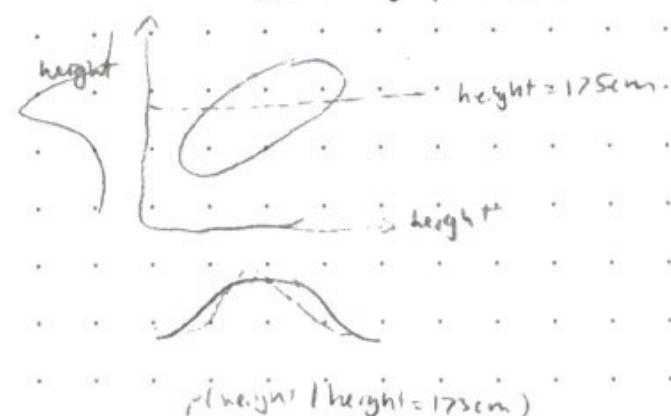
$P(D|M_i)$  (also  $P(D|M_i)$ )



draw  $\theta$  from prior then given  $\theta$  draw  $D$  from sampling distribution (what I called in serial joint)  $\rightarrow$  this gives  $\theta$  and  $D$  pairs from the joint distribution e.g.  $P(D, \theta|M_i)$   
These are not data points they are samples from joint 'observed' data set (there is just a single value shown by dotted line)

(draw from prior  $\theta \sim P(\theta|M)$ )  
given me to draw from joint  $D \sim P(D|M, \theta)$   
sampling dist.  $\theta, D$  from joint!

just like graph in mode



# Interpreting the Bayes Factor / Evidence Ratio: The Jeffreys Scale

$$\Delta \ln E = \ln BF$$

$$\begin{array}{l} \Delta \ln E < 1 \\ 1 < \Delta \ln E < 2.5 \\ 2.5 < \Delta \ln E < 5 \\ 5 < \Delta \ln E \end{array}$$

"not worth more than a bare mention"  
"significant"  
"strong to very strong"  
"decisive"

These categories  
are up for debate  
(not rigorous)

## How to calculate Evidence

$$P(D) = \int p(D|\theta) p(\theta) d\theta$$

\* Analytic (really simple, nice problems)

(analytic  
integral)

\* Laplace approximation

(calculator)

\* Savage-Dickey Ratio (Nested Models)

likely slow;  
almost certainly  
strongly peaked  
if it is at all informative  
about the data

\* Monte-Carlo: probably slow if likelihood is peaked

fundamentally  
never MC approx  
integral by sum

$$\theta_i \sim p(\theta_i)$$

$$P(D) \approx \frac{1}{M} \sum_{i=1}^M P(D|\theta_i)$$

approximate with  
MC avg of  
likelihood

\* Harmonic Mean Estimator (unstable)

$$\theta_i \sim p(\theta|D)$$

$$P(D) \approx \left[ \frac{1}{M} \sum_{i=1}^M P(D|\theta_i) \right]^{-1}$$

MC avg. vs  
inverse  
of the pdf  
inversely that

\* Nested Sampling

Calc. evidence via Laplace approx:

## (Recall: The Laplace Approximation)

Evidence  $P(D) = \int P^*(\theta|D) d\theta$

Unnormalized posterior  $P^*(\theta|D) = P(D|\theta)P(\theta)$

Find MAP estimate:  $\theta_0 = \operatorname{argmax}_{\theta} \ln P^*(\theta|D)$

Taylor expansion:

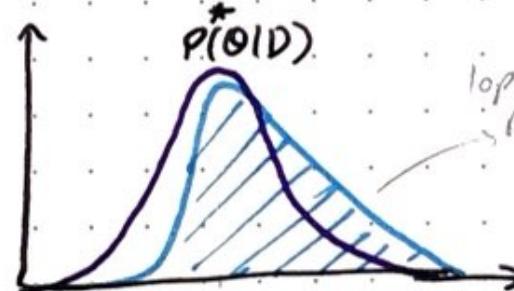
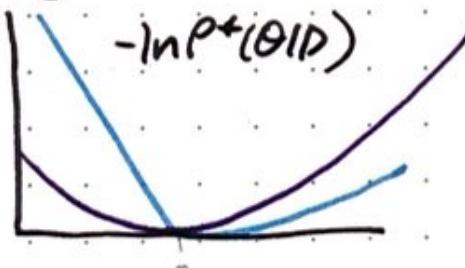
$$\ln P^*(\theta|D) \approx \ln P^*(\theta_0|D) - \frac{1}{2} (\theta - \theta_0)^T A (\theta - \theta_0) + \dots$$

Hessian at mode:  $A_{ij} = -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln P^*(\theta|D) \Big|_{\theta=\theta_0}$

$$P^*(\theta|D) \approx P^*(\theta_0|D) \times \exp(-\frac{1}{2} (\theta - \theta_0)^T A (\theta - \theta_0))$$

Not always accurate:

e.g.



e.g.  
top-sided unnormalized  
posterior  
Laplace  
upper bound  
approx

↓  
but if you trust your approximations you can do an integral over priors & get this:

### \*Evidence from Laplace Approximation:

$$P(D) = \int P^*(\theta|D) d\theta \approx P^*(\theta_0|D) \times 12\pi |A^{-1}|^{1/2} \int N(\theta|\theta_0, A^{-1}) d\theta$$

$$P(D) \approx P^*(\theta_0|D) \times \det(A/2\pi)^{-1/2}$$

$$P(D) \approx P^*(\theta_0|D) \times \det(2\pi A^{-1})^{-1/2}$$

$$P(D) \approx P(D|\theta_0) \times P(\theta_0) \times \det(A/2\pi)^{-1/2}$$

hessian  
recovery  
posterior

$$A = -\nabla \nabla \ln P^*(\theta|D)$$

?  
unnormalized?

need to be able to  
differentiate posterior at least  
2x to do this (?)  
(so only works for tractable posterior!)

this implements Occam's razor

## Evidence implements Occam's Razor

some sort  
similar to  
analysis  
before

$$P(D) \approx P(D|\theta_0) \times P(\theta_0) \times \det(A/2\pi)^{-1/2}$$

Evidence

best fit  
likelihood

Occam factor

$$\Sigma_{\text{post}} = A^{-1} \quad \text{roughly}$$

how can  
we see this?

$$\text{Suppose: } P(\theta) = N(\theta | \theta_{\text{prior}}, \Sigma_{\text{prior}})$$

Suppose, prior is  
a gaussian

Occam Factor:

$$\propto 12\pi \Sigma_{\text{prior}}^{-1/2} \times 12\pi \Sigma_{\text{post}}^{-1/2}$$

$$\propto \frac{1 \Sigma_{\text{post}}^{1/2}}{1 \Sigma_{\text{prior}}^{1/2}}$$

From inference book:

- factor by which our model's hypothesis space collapses when the data arrive
- magnitude of occam factor is a measure of complexity of the model
- depends on both
  - no of picums in model
  - prior prob model assigns to those picums

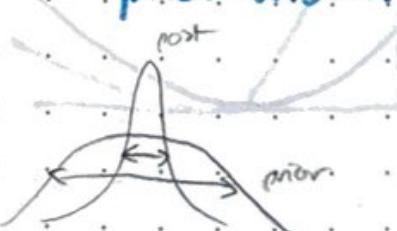
more picums → multiply prior by IV  
(wide prior, complex model)  
then smaller occam rule

$$\frac{P(D|M_1)}{P(D|M_2)}$$

notes to self:

implements Occam's razor for each  
tradeoff between maximizing data fit &  
minimizing complexity of model  
so get more accurate comparison?

(ratio of posterior to  
prior width)



Simple models concentrate probability around limited no of datasets  
complex models cover wide range of datasets but probability of a given dataset will be smaller i.e. even if 'boring' simple & complex model run predict data, complex will have much smaller prob of

## Savage - Dickey Ratio

for nested models?  
useful if we have simple  
model nested within  
complex model.

Suppose the parameters are  $\phi, \psi$  and complex  
model  $M_1$  reduces to simpler model  $M_0$  when  
 $\psi = 0$ . (Nested)

And the prior is separable:

$$P(\phi, \psi | M_1) = P(\psi | M_1) P(\phi | M_1)$$

nuisance  
parameters  
(uninteresting)

models have  
free parameters in common  
(not of direct interest so  
not the focus of our model  
(compensation))

or e.g. linear  
as a special  
case of  
polynomial  
constants  
= 0

And the prior on  $\phi$  is the same for each  
model:

invariant  
nuisance  
parameters

$$P(\phi | M_1) = P(\phi | M_0)$$

Then the Bayes Factor reduces to:

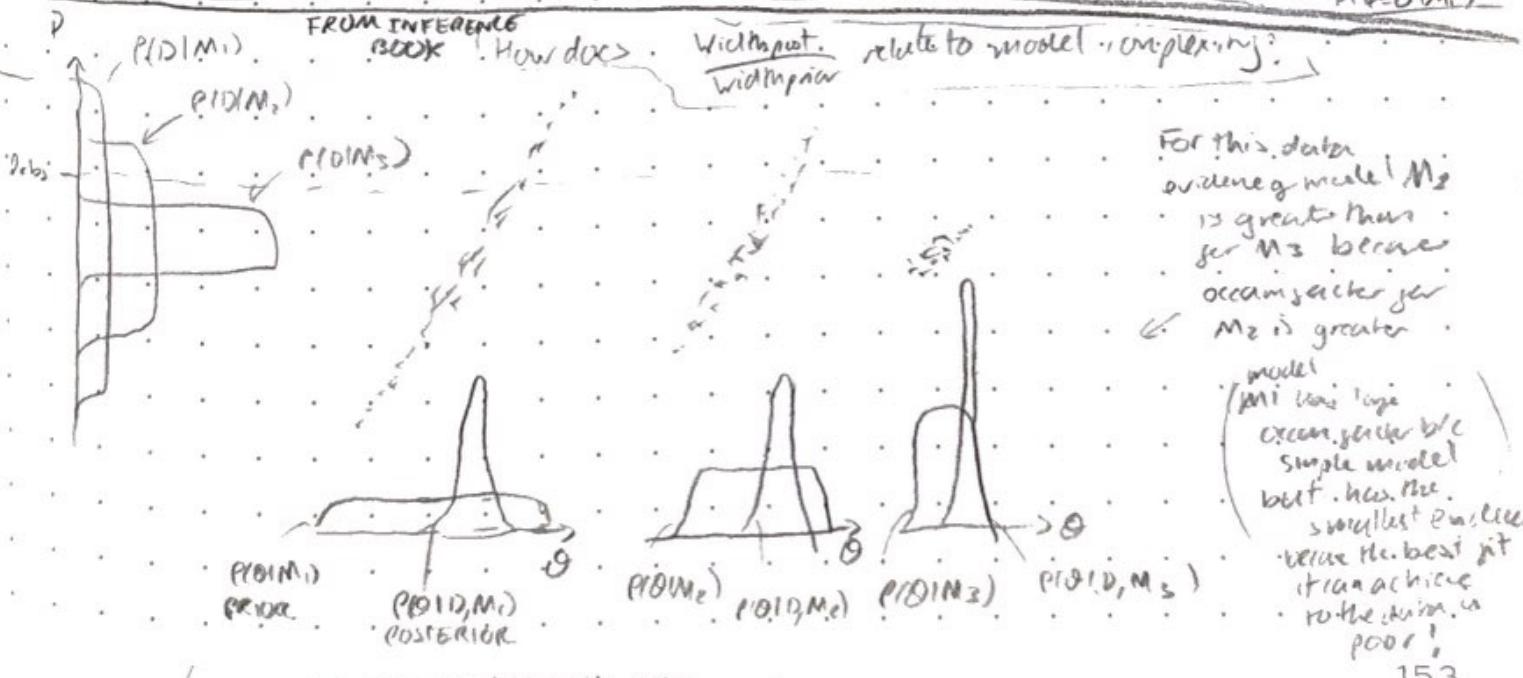
$$B_{01} = \frac{P(D|M_0)}{P(D|M_1)} = \frac{P(\psi | D, M_1)}{P(\psi | M_1)} \Big|_{\psi=0}$$

Show this:  
a. consider numerator.  
 $P(D|M_0)$

$$\begin{aligned} &= \int P(D|\phi, \psi=0, M_0) P(\phi|M_0) d\phi \\ &= \int P(D|\phi, \psi=0, M_1) P(\phi|\psi=0, M_1) d\phi \\ &= P(D|\psi=0, M_1) \\ &= P(\psi=0 | D, M_1) / P(D|M_1) \end{aligned}$$

Bayes  
(logarithms)

b. now ratio becomes  
 $B_{01} = P(\psi=0 | D, M_1) / P(\psi=0 | M_1)$



## NESTED SAMPLING

for evidence calculation  
but also for posterior sampling  
(use get logproduct!)

- An algorithm designed to compute Bayesian evidence (John Skilling, 2004)
- Can be better than MCMC for multi-modal distributions
- Get weighted samples from posterior for free
- Evolve an ensemble of "live points" successively sampling prior volume above a likelihood-level constraint
- Perform evidence integral over contours of equal likelihood
- Uses statistical estimate of prior mass within likelihood-level

Want to calculate multi-dim. evidence integral

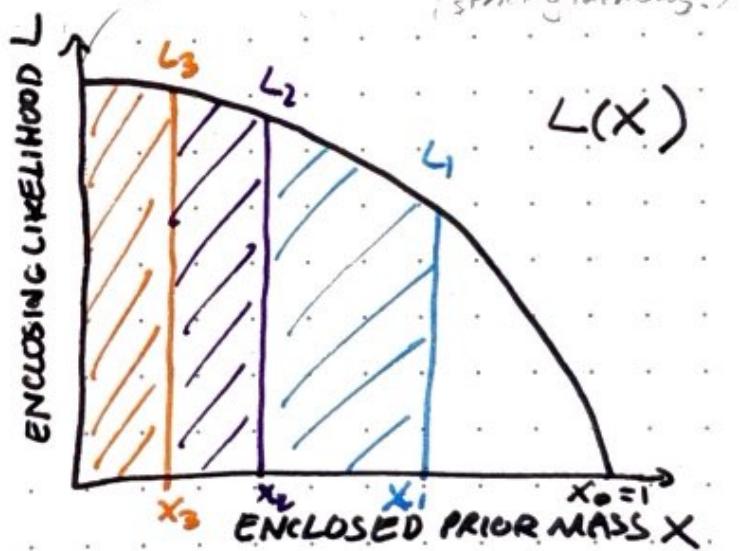
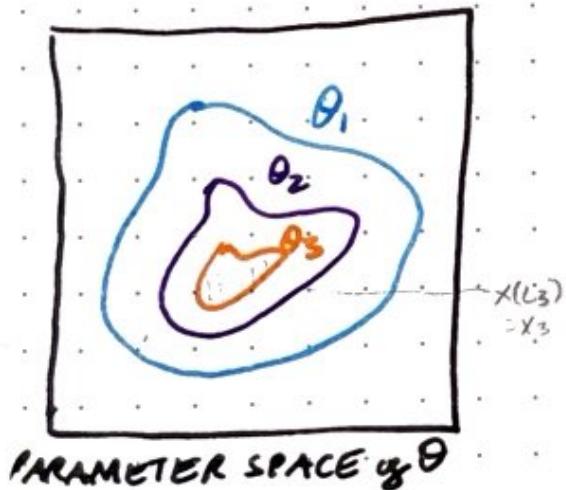
$$Z = \int L(\underline{\theta}) \pi(\underline{\theta}) d\underline{\theta}$$

Define "prior mass enclosed above likelihood  $L^*$ "

$$X(L^*) = \int_{L(\underline{\theta}) > L^*} \pi(\underline{\theta}) d\underline{\theta}$$

Inverse: Likelihood at enclosed prior mass:  $L(x)$

1D evidence integral:  $Z = \int L(x) dx$



Nested likelihood contours sort to enclosed prior mass  $X$ .

e.g. 2D example: take 3 random points from prior  $\theta_1, \theta_2, \theta_3$ , and calc. their likelihood, each will be associated with some likelihood contour. Each contour is associated with some percentage of the prior volume  $X(L^*)$ .

so can associate with each point a likelihood contour and also a prior mass enclosed above that likelihood.

if i can map from points in param. space to  $L(x)$  I can turn a previously multi-dim. integral in param space to a 1-Dim integral in  $x$  from 0 to 1.

how do we  
actually calc. this?

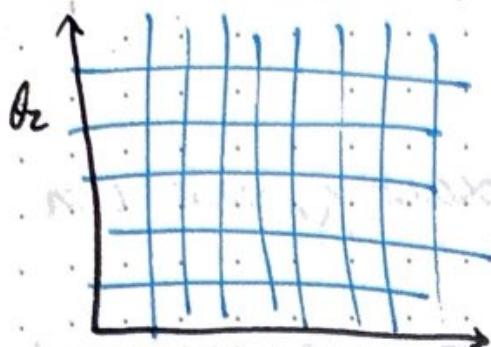
## Nested Sampling : How?

[NB: For the rest of this lecture assume prior is uniform box in parameter space (w/ total integral 1)]

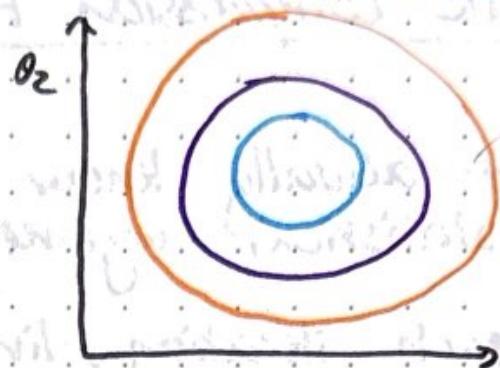
- ① Begin by sampling  $N$  live points  $\{\theta_i\}$  from the prior ( $x_0 = 1$ ), evaluate likelihood for each live point, initialise evidence  $Z = 0$ .
- ② Find the live point with smallest likelihood  $L^*$ .
- ③ Estimate compression factor  $t_i = x_i / x_{i-1}$ .  
volumes  $t_i$ : estimated statistically  
 $t_i \propto e^{-N \ln t_i}$
- ④ Accumulate evidence as a 1D integral  $\Delta Z = L^* \times (x_{i-1} - x_i) = L^* \times (1 - t_i)x_i$ .
- ⑤ Sample a new live point from prior constrained to  $L > L^*$ , compute its likelihood value.  
(e.g. with MCMC or otherwise)  
reject owing not higher  $L$  (diff to dead points)  
sample is w/o replacement  
reject if  $L \leq L^*$  (?)
- ⑥ Repeat steps 2-5 until convergence (i.e. evidence  $Z$  stops changing within some tolerance).
- ⑦ Add final evidence estimate of remaining live points  $\Delta Z = \bar{L} \times x_{\text{end}}$ .
- ⑧ List of dead points give weighted posterior samples  $\{\theta_i\}$  with  $w_i \propto (1 - t_i)x_i$ .

want to find  $Z$ : instead of dividing into cubes & summing over  $L$  for each cube, combine cubes w similar likelihood & perform easier 1D sum (only care about difference in overall their contents?)

$$Z \approx \sum L(\theta_1, \theta_2) d\theta_1 d\theta_2$$



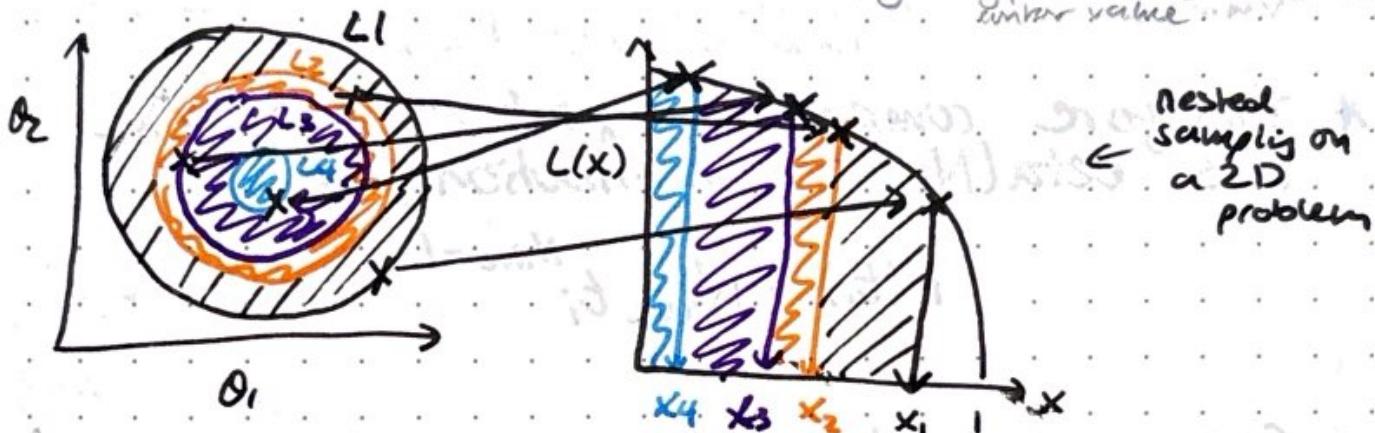
$$Z \approx \sum L(x) dx$$



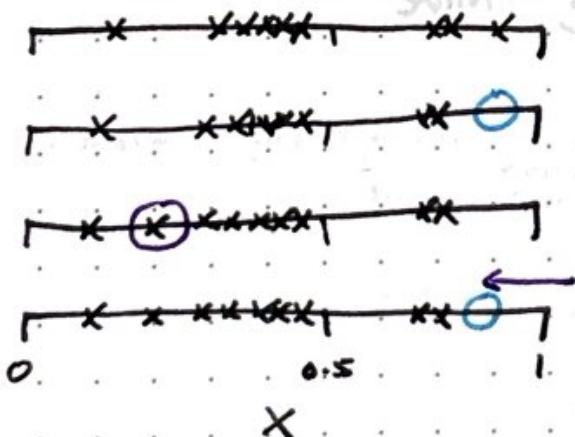
naive: sum over cubes and add them up

but in huge parameter spaces this is very computationally expensive

nested sampling: "combining" cubes by  
similar likelihood is summing over these  
in a prior sample integral by integrating  
over contours of equal likelihood &  
multiplying by prior volume "under" each  
contour value.



compression in 1. iterate of nested sampling:



uniformly distributed tie points

remove worst

draw replacement

compression,  $\epsilon = \beta N^{1/(n+1)}$

## Estimate compression Factor

tricky bit → can't actually see  $\mathbb{E}[X_i]$   
contours, esp in high dim  
spaces → so use statistical  
argument to estimate  
what they are

- \* Don't actually know  $\mathbb{E}[X_i]$  exactly, but can use statistical argument (to estimate it)
- \* At each iteration, live points are uniformly distributed between  $(0, X_{i-1})$ .
- \*  $X_i$  (lowest likelihood) in current iteration, is highest  $X$  of  $N_{\text{live}}$  uniform random values
- \* Therefore compression factor  $t_i = X_i / X_{i-1}$  has beta( $N_{\text{live}}, 1$ ) distribution, i.e.

$$P(t_i) = N_{\text{live}} t_i^{N_{\text{live}}-1}$$

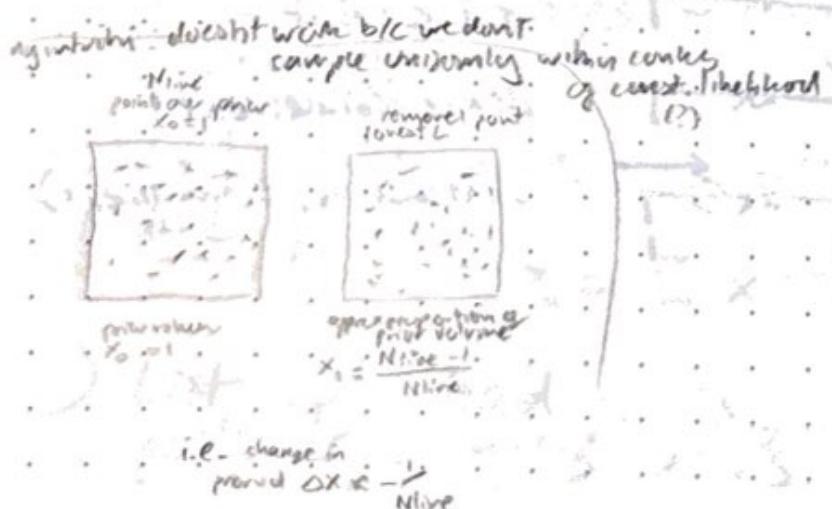
highest point of  
these  $N_{\text{live}}$  simply  
was beta dist  
so  $t$  also has  
beta dist

- \* Can approximate with mean, i.e.

$$\ln t_i \approx \langle \ln t_i \rangle = -1/N_{\text{live}}$$

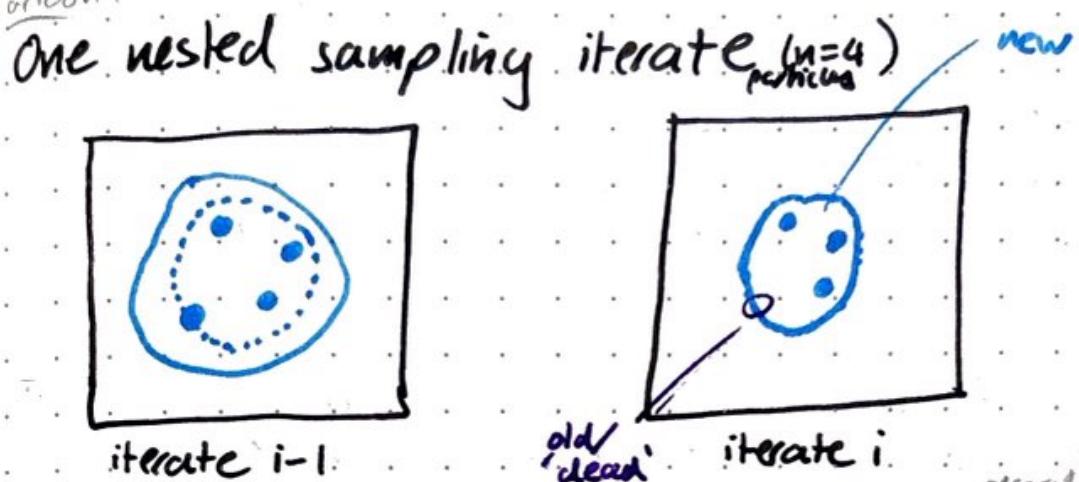
$$\text{or } t_i \approx e^{-1/N_{\text{live}}}$$

$$\begin{aligned} E[t] &= \int_0^1 N t^{N-1} dt \\ &= \frac{N}{N-1} ?? \end{aligned}$$



# Nested sampling Demo:

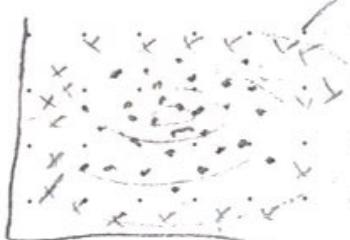
Cartoon:



Dep:

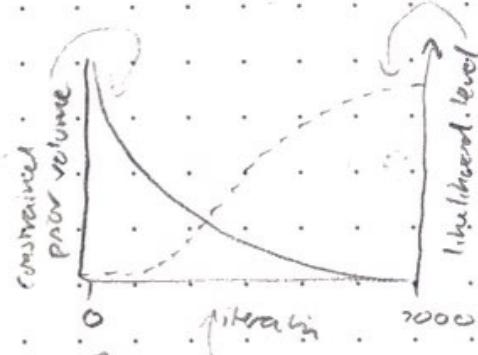
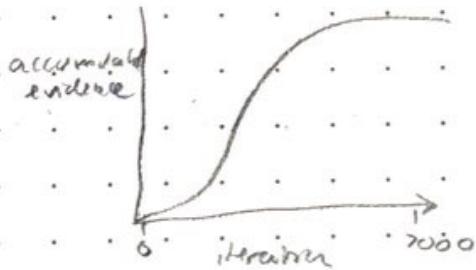
- 2D gaussians
- uniform prior
- over some range
- e.g.  $-5, 5$
- (basically covers gaussian)

Sample  
1000  
line  
points



record  
dead point 15

likelihood contours



know  
analytic  
 $2 - \alpha(1 - e^{-x})$

see lecture

- ↳ 5 nice gaussians
- we know analytic
- ↳ can see nested sampling approaches

correct?

(do lots of iterations)  
likelihood of live points  
that remain should  
basically be same

likelihood

( $\chi^2$  of that is ~0?)

only or not accumulate  
much but accumulate  
must where prior volume  
& likelihood have  
similar values (i.e. middle  
(where at either end  
both trend to zero))  
so not much  
accumulation  
=> S-shape  
(curve)

since  
accumulated  
evidence is  
more likelihood

use sticks because  
example, sample grows  
but no new dist.  
have a go if you  
want

~~basically~~ were considered the most basic nested sampling algorithm, but there are more!  
our way slows down w/ iteration  $\rightarrow$  smaller & smaller regions trying to sample you.

## Nested Sampling (Other Algorithms)

### State-of-the-Art Nested Sampling

- \* Smarter ways to sample from likelihood-constrained prior, e.g. MCMC, slice sampling, sequential Monte Carlo
- \* Stochastically sample a set of compression factors  $\{t_i\}$  rather than assume a constant mean value  $\rightarrow$  gives a measure of uncertainty in evidence estimate
- \* Better strategic allocation of live points e.g. dynamic nested sampling
- \* Codes: Multinest, UltraNest, Polychord, dynesty

yay done!

useful properties of determinants:

- $\det(A \underline{\underline{B}}) = \det(A) \det(\underline{\underline{B}})$  ( $= \det(\underline{\underline{A}}) \det(\underline{\underline{B}})$ ) ; rule
- $\det(\underline{\underline{A}}^T) = \det(\underline{\underline{A}})$
- $\det(\underline{\underline{A}}^{-1}) = \frac{1}{\det(\underline{\underline{A}})} = \det(\underline{\underline{A}})^{-1}$

non matrix  $\underline{\underline{A}}$ :  $\det(k\underline{\underline{A}}) = k^n \det(\underline{\underline{A}})$

$$\hookrightarrow \text{e.g. } 12\pi \underline{\underline{E}}^{-1/2} = (12\pi)^{-n/2} |\underline{\underline{E}}|^{-1/2}$$

# MCMC Theory

PREMISE:

If  $\pi(a)$  satisfies D.B.:  $\pi(a) T(b|a) = \pi(b) T(a|b)$ ,  
then stationary distribution  $\pi_{\text{stat}}$  exists! (sufficient but not necessary condition)

And additionally under the conditions 1) irreducible  
2) aperiodic and 3) positive recurrent, we can  
guarantee that our MC will converge to the  
unique stationary distribution  $\pi_{\text{stat}}$ .

SHOW THAT OUR ALGS ARE CONSTRUCTED S.T. DB SATISFIED:

We want to construct our algorithm s.t. DB is satisfied and  $\pi_{\text{stat}}(a) = \pi_{\text{post}}(a)$ , i.e. distribution we converge to is the target distribution (posterior).

$$\therefore \text{require } \pi_{\text{post}}(a) T(b|a) = \pi_{\text{post}}(b) T(a|b) \quad \nwarrow (\text{DB})$$

$$T(b|a) = J(b|a) A(b|a)$$

↑  
prob make the transition

↑  
prob propose the jump

↑  
prob accept the jump

$\Rightarrow$  we want to choose  $J, A$  s.t. they satisfy

$$\pi_{\text{post}}(a) J(b|a) A(b|a) = \pi_{\text{post}}(b) J(a|b) A(a|b)$$

$$\left( \frac{A(b|a)}{A(a|b)} = \frac{\pi_{\text{post}}(b)}{\pi_{\text{post}}(a)} \frac{J(a|b)}{J(b|a)} \right)$$

METROPOLIS SATISFIES DB:  $J(a|b) = J(b|a)$ ,  $A(b|a) = \min\left(\frac{\pi_{\text{post}}(b)}{\pi_{\text{post}}(a)}, 1\right)$

- suppose wlog.  $\pi_{\text{post}}(b) > \pi_{\text{post}}(a)$   $(LHS = RHS \text{ (should really do this separately)})$
- $\pi_{\text{post}}(a) J(b|a) = \pi_{\text{post}}(b) J(a|b) \frac{\pi_{\text{post}}(a)}{\pi_{\text{post}}(b)}$
- since  $J(b|a) = J(a|b)$ ,  $RHS = LHS$ !

SHOW THAT METROPOLIS-HASTINGS SATISFIES DB: (more general)

M-H:  $J(a|b) \neq J(b|a)$  necessarily,  $A(b|a) = \min\left(\frac{\pi_{\text{post}}(b)}{\pi_{\text{post}}(a)}, \frac{J(a|b)}{J(b|a)}, 1\right)$

$$\text{wts. } \pi_{\text{post}}(a) J(b|a) A(b|a) = \pi_{\text{post}}(b) J(a|b) A(a|b)$$

$\rightarrow$  since either  $A(b|a) = 1$  or  $A(a|b) = 1$ , easy to show  
 $LHS = RHS$ ,

GIBBS

$$\text{w.t.s. } \pi(\theta_1, \phi_1) T(\theta_2, \phi_2 | \theta_1, \phi_1) = \pi(\theta_2, \phi_2) T(\theta_1, \phi_1 | \theta_2, \phi_2) \quad A=1$$

$$\text{since for gibbs always accept, } T(\theta_2, \phi_2 | \theta_1, \phi_1) = J(\theta_2, \phi_2 | \theta_1, \phi_1) \times 1$$

For gibbs start at  $\theta_1, \phi_1$ , then propose move  $\theta$  then propose move  $\phi$   
so  $J(\theta_2, \phi_2 | \theta_1, \phi_1) = \pi(\theta_2 | \phi_1) \pi(\phi_2 | \theta_2)$

and use fact that going  $\theta_1, \phi_1 \rightarrow \theta_1, \phi_2 \rightarrow \theta_2, \phi_2$

$$\text{so LHS: } \pi(\theta_1, \phi_1) \pi(\theta_2 | \phi_1) \pi(\phi_2 | \theta_2) = \pi(\theta_1, \phi_1) \pi(\theta_2, \phi_2) \frac{\pi(\theta_2, \phi_1)}{\pi(\theta_1) \pi(\phi_1)}$$

$$\text{RHS: } \pi(\theta_2, \phi_2) \pi(\theta_1 | \phi_2) \pi(\phi_1 | \theta_1) = \pi(\theta_1, \phi_1) \pi(\theta_2, \phi_2) \frac{\pi(\theta_1, \phi_2)}{\pi(\theta_1) \pi(\phi_2)}$$

start  $\theta_1, \phi_1$ , update to  $\theta_1$  as f'n of  $\phi_2$  then update to  $\phi_2$  as f'n of  $\theta_1$ ,  $\theta_1$  is independent of  $\phi_2$ !

but since  $\phi_1$  is independent of  $\theta_2$  (per process on RHS)  
 $\pi(\theta_1) \pi(\phi_2) = \pi(\theta_1, \phi_2)$

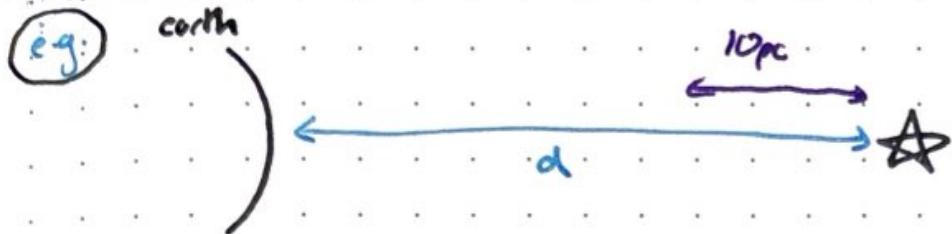
doing similarly for LHS  $\Rightarrow$  LHS = RHS

$\Rightarrow$  DB satisfied

$\Rightarrow$  stationary dist  $\pi_{\text{stat}} = \pi$  exists

# Brief Recap on Magnitude Scale

(eg) earth



star with intrinsic luminosity  $L_*$

how bright does it appear to us?

(Reminder: Flux)

energy time<sup>-1</sup> area<sup>-1</sup>

Luminosity)

energy time<sup>-1</sup>)

## APPARENT MAGNITUDE

apparent mag at distance  $d$

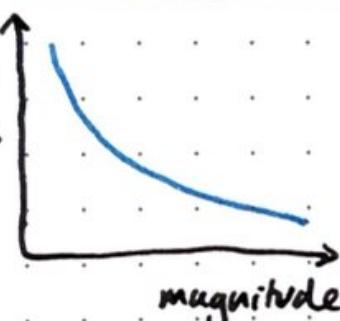
define apparent mag at distance  $d$

$$M_d = -2.5 \log_{10} \left( \frac{F_d}{F_0} \right)$$

flux at distance  $d$

$$= -2.5 \log_{10} F_d + \text{const.}$$

apparent brightness



some  
residue  
flux / zero  
point on  
filter

some  
calibration  
const. to do  
with residue  
flux & filter

e.g. decrease <sup>apparent</sup> mag 5 is  $10^{2.5} \times 2.5$   
star mag 5 is 10<sup>2.5</sup> times brighter than star mag 6

lower mag = brighter

## ABSOLUTE MAGNITUDE

defined as magnitude object would have if it were viewed at a distance of 10pc

$$\begin{aligned} M = M_{10pc} &= -2.5 \log_{10} \left( \frac{F_{10pc}}{F_0} \right) = -2.5 \log_{10} \left( \frac{F_d}{F_0} \frac{d^2}{(10pc)^2} \right) \\ &= -2.5 \log_{10} \left( \frac{F_d}{F_0} \right) - 5 \log_{10} \left( \frac{d}{10pc} \right) \end{aligned}$$

(luminosity distance is distance measured from luminosity  
(not necessarily "true" distance))

$$M = m - 5 \log_{10} \left( \frac{d}{10pc} \right) = m - 5 \log_{10} \left( \frac{d}{pc} \right) + 5$$

Distance modulus  $\mu$  def s.t.  $M = m - \mu$  (i.e.  $\mu = 5 \log_{10} \left( \frac{d}{10pc} \right)$ )  
(logarithmic distance)  $\rightarrow$  star at  $\mu=2$  is  $\sqrt{1.6}$  times further than star at  $\mu=1$  247

# Regression

(1)

dependent variable

independent variable

- special case of model fitting

- model describes  $E[y|x]$

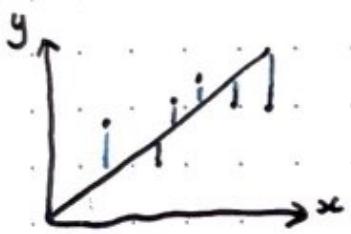
note: linear model = linear in params, e.g.  $f(x_i; \theta) = \theta_0 + \theta_1 x_i$   
 $\theta_0 + \theta_1 \sin(\theta_2 x)$  not linear model

## OLS: ordinary least squares

(we focus on linear model)

(assume no / negligible error in x)

e.g. simplest example: noisy data  $y_i = f(x_i; \theta) + \epsilon_i$   
want to fit model  $f(x_i; \theta) = \theta_0 + \theta_1 x_i$



method: find params which minimize

$$RSS = \sum_i (y_i - (\theta_0 + \theta_1 x_i))^2$$

(multiply out and find)  
 $\frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_0}$

$$\Rightarrow \begin{cases} \theta_0 = \bar{y} - \theta_1 \bar{x} \\ \theta_1 = \frac{\sum x_i y_i - N \bar{x} \bar{y}}{\sum x_i^2 - N \bar{x}^2} \\ = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \end{cases}$$

[also, OLS estimator of variance:

$$\hat{\sigma}^2 = \frac{RSS}{n-2}$$

## OLS relationship to max likelihood

→ (why did we square residual and not e.g. take magnitude?)

assume gaussian errors in y:  $P(y_i | \theta, M) = N(y_i | f(x_i; \theta), \sigma_i^2)$

Probability of observing dataset:

$$P(D|\theta, M) = \prod_i \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{1}{2} \frac{(x_i - f(x_i; \theta))^2}{\sigma_i^2}} = L(\theta)$$

$$\Rightarrow \ell(\theta) = -\frac{1}{2} \ln(2\pi\sigma_i^2) - \frac{1}{2} \sum_i \frac{(x_i - f(x_i; \theta))^2}{\sigma_i^2}$$

[ "least squares" = "max likelihood" ]

if data are gaussian

(2)

Generalise:

First introduce matrix notation  $Y = X\beta + \epsilon$   
 e.g. for previous example

$$Y = \begin{pmatrix} y_0 \\ y_1 \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_0 \\ 1 & x_1 \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

OLS:

$$\text{RSS} = (Y - X\beta)^T (Y - X\beta)$$

(check using Einstein summation  
convention notation)

$$\hat{\beta}_{\text{OLS}} = (X^T X)^{-1} X^T Y$$

(minimize  $y_i^2 - 2y_i x_{ij} \beta_j + x_{ij}^2 \beta_j^2$ )

unbiased!

- weighted least squares aka  $\chi^2$  minimisation is, like OLS, another special case of generalised least squares

$$\text{minimize } \chi^2 = \sum_i \frac{(y_i - x_{ij} \beta_j)^2}{\sigma_i^2}$$

with  $W = \begin{pmatrix} \sigma_1^2 & 0 & \dots \\ 0 & \sigma_2^2 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$ :  
diagonal cov matrix

Generalised least squares

linear model  $Y = X\beta + \epsilon$

correlated errors  $W = \text{Var}(\epsilon) = \text{Cov}(\epsilon, \epsilon^T)$  (known)

GLS:

$$\text{RSS} = (Y - X\beta)^T W^{-1} (Y - X\beta)$$

$$\hat{\beta}_{\text{GLS}} = (X^T W^{-1} X)^{-1} X^T W^{-1} Y$$

note: "generalised least squares" = "max likelihood"

holds in general for gaussian errors (so for OLS &  $\chi^2$  minimisation alike!)

OLS estimate of unknown function  
 $\hat{\sigma}^2 = \frac{1}{N-k} (Y - \hat{X}\hat{\beta})^T (Y - \hat{X}\hat{\beta})$   
 no. of estimated parameters