

# BAYESIAN INFERENCE

In astronomy always have errors on everything

- Complete Data Likelihood

$$P(x_i, y_i | \xi_i, \gamma_i, G, \Psi) = P(x_i, y_i | \xi_i, \gamma_i) P(\gamma_i | \xi_i, \Theta) P(\xi_i, \Theta)$$

beat  $G, \Psi$  as *observers* we "wish we had"

- Observed data likelihood marginalise out

$$P(x_i, y_i | G, \Psi) = \iint \dots d\xi_i d\gamma_i$$

$$\Rightarrow P(x_i, y_i | \Theta, \Psi) = \prod P(x_i, y_i | \Theta, \Psi) \text{ assuming indep}$$

- In Frequentist stats parameters are not random, only data is

So what are  $(\xi_i, \gamma_i)$ ? Fill in conceptual gap "nuisance parameter"  
data or parameter? - not observed but have a prob distribution

## Bayesian viewpoint

- symmetry between data  $\leftrightarrow$  param  $\leftarrow$  RVs that are not observed  
all RVs described by  $P(D, G)$  can be observable  $D \leftrightarrow G$

$\rightarrow$  goal to infer unobserved parameters from observed data

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

prob of things I didn't observe  
given the things I did observe

Probability = degree of belief / uncertainty in hypotheses

$$P(G|D) = P(D|G) P(G) / P(D)$$

$\uparrow$   $\uparrow$   $\leftarrow$   $\leftarrow$

POSTERIOR  
degree of belief    LIKELIHOOD  
Sampling distribution    PRIOR  
degree of belief    NORMALISATION

## Bayesian credible interval

$$Y_1, \dots, Y_N \text{ iid } \sim N(\mu, \sigma^2)$$

$$\bar{Y} \sim N(\mu, \frac{\sigma^2}{N})$$

$$\text{flat prior } P(\mu) \propto 1 \quad p(\mu | \bar{Y} = \bar{y}_{\text{obs}}) \propto P(\mu) \cdot \prod_{i=1}^N (y_{\text{obs},i} | \mu, 1)$$

$$y_{\text{obs}} = (-0.64, -0.93, 0.16, -0.88)$$

product of gaussians

$$= N(\mu | \bar{y}, 1/N) = N(\mu | -0.57, 0.5^2)$$

degree of belief

$$P(\mu | D = y_{\text{obs}}) = 68\% \\ \in [-1.07, -0.07]$$

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

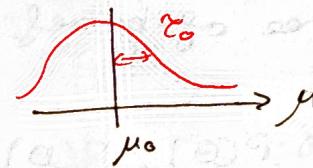
$$P(Y | \mu, \sigma^2) = \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^N \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right] \\ (y_i - \bar{y} + \bar{y} - \mu)^2 \\ = (y_i - \bar{y})^2 + 2(y_i - \bar{y})(\bar{y} - \mu) + (\bar{y} - \mu)^2 \\ = (y_i - \bar{y})^2 + 2(N-1)(\bar{y} - \mu)^2 + N(\bar{y} - \mu)^2 \\ = (6\sqrt{2\pi})^{-N} \exp \left\{ -\frac{1}{2\sigma^2} \left[ \sum_{i=1}^N (y_i - \bar{y})^2 + N(\bar{y} - \mu)^2 \right] \right\}$$

$$P(Y | \mu, \sigma^2) = (6\sqrt{2\pi})^{-N} \exp \left[ (N-1)s^2 + N(\bar{y} - \mu)^2 \right]$$

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

- Introduce prior suppose  $\sigma^2$  known ~~flat~~  $\mu \sim N(\mu_0, \tau_0^2)$

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



$$\Rightarrow \text{posterior } P(\mu | y) \propto P(y | \mu) P(\mu)$$

$$\propto (6\sqrt{2\pi})^{-N} \exp \left[ (N-1)s^2 + N(\bar{y} - \mu)^2 \right]$$

$$\cdot \frac{1}{\sigma \sqrt{2\pi}} \exp \left[ -\frac{1}{2\tau_0^2}(\mu - \mu_0)^2 \right]$$

again both gaussians in  $\mu$ , combine

$$\Rightarrow P(\mu | y) = N(\mu_0, \sigma_N^2)$$

$$\mu_N = \left( \frac{\mu_0}{\sigma_0^2} + \frac{N \bar{y}}{G^2} \right) / \left( \frac{1}{\sigma_0^2} + \frac{N}{G^2} \right) \quad \text{precision weighted mean}$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{G^2}$$

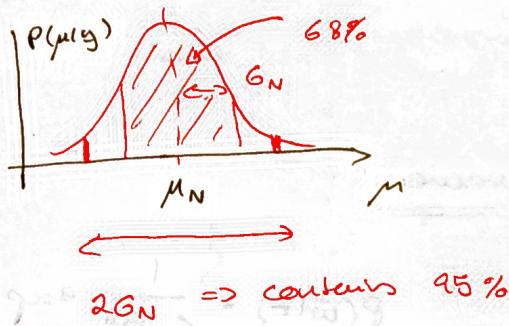
$\sigma_0 \rightarrow \infty$  flat prior so  $\mu_N \rightarrow \bar{y}$  and  $\sigma_N^2 \rightarrow G^2/N$

- even for  $N \rightarrow \infty$  get  $\mu_N \rightarrow \bar{y}$   $\sigma_N^2 \rightarrow G^2/N$

$\Rightarrow$  This is called CONJUGACY  $\mu$  is conjugate prior  
bc  $\mu$  is in the same family as  $y_i$

- ~~closed~~ credible interval

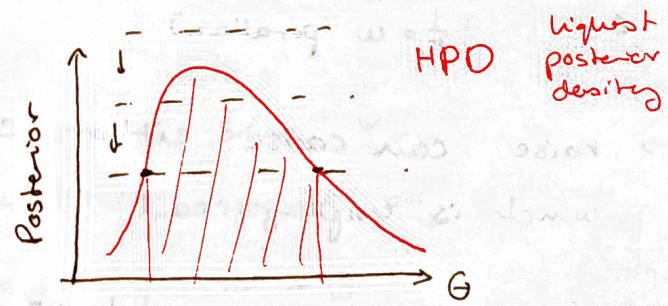
by convention report  
central interval



In general posterior could be quite complicated,  
want to communicate some information with only few numbers

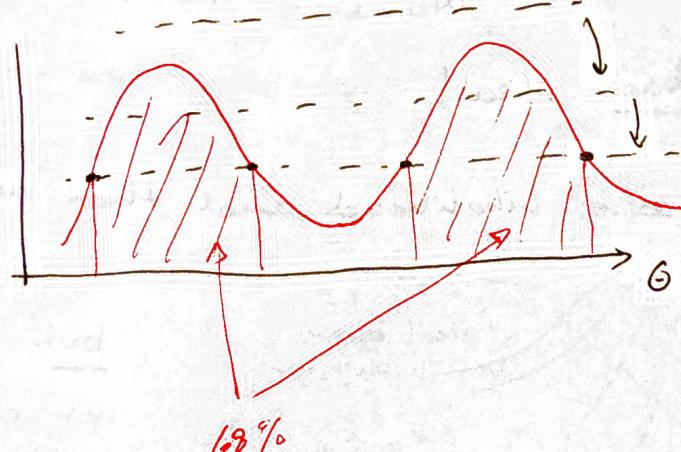
For a non-symm distribution

lower "bar" until interval  
between two points is 68%



For a bimodal distribution

lower "bar" → intersects  
at even # of points  
→ take intervals between  
pairs of points  
→ lower until intervals  
add up to 68%



- Frequentist  $\Leftrightarrow$  Bayesian

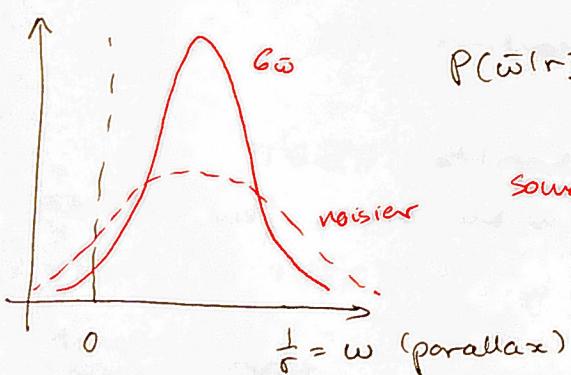
Bayesian answer is the full posterior  $P(\theta|D)$  "state of knowledge"  
actually get numerical estimate of posterior

- can include prior information - external dataset
- by choosing prior can regularise and penalise having too many free variables ie overfitting w complexity

$\rightarrow$  likelihood is not prob density in parameters but the posterior is so conditional, has marginal prob  
 $\rightarrow$  can deal w high dimensional param space by marginalising

NB: any Bayesian estimate can be evaluated as if frequentist

## Parallax Measurement



$$P(\bar{\omega}|r) = \frac{1}{6\omega \sqrt{2\pi}} \exp \left[ -\frac{1}{26\omega^2} (\bar{\omega} - \frac{1}{r})^2 \right]$$

sometimes say "measurement error" when mean  $6\omega$

fractional meas error  $f = \frac{6\omega}{\omega} = r \cdot 6\omega$  since  $\omega = \frac{1}{r}$

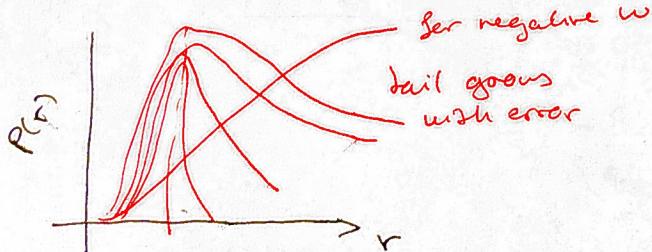
$\rightarrow$  noise can cause either 0 or negative  $\omega$ , which is unphysical how do we correct for this?

① could do:  $r = \frac{1}{\omega}$   $6r^2 \approx |\frac{1}{\omega^2}|^2 6\omega^2 \Rightarrow 6r = \frac{6\omega}{\omega^2}$

bad idea! get very skewed functions  $\rightarrow r$

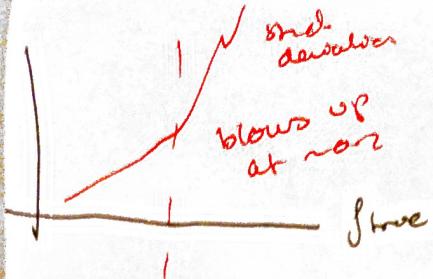


② could solve likelihood and then add  $P(r) = \begin{cases} 1 & r > 0 \\ 0 & \text{otherwise} \end{cases}$  so only positive r



but improper posterior  
not normalisable

→ if we compute frequentist properties see slides for method



problematic bc always blows up

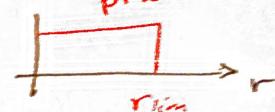
use diff distributions for r true to draw from to simulate properties of estimator

(3) put in a proper distance prior

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

cutoff limits damage of spurious measurements

$r_{\text{lim}}$  = max possible distance of star in survey



$$P(r|\omega) \propto P(\omega|r) \cdot P(r) \Rightarrow P(r|\bar{\omega}, G\bar{\omega}) = \begin{cases} \frac{1}{r_{\text{lim}}} P(\bar{\omega}|r, G\bar{\omega}) & r > 0 \\ 0 & \text{otherwise} \end{cases}$$

The cutoff means that for  $\bar{\omega} > r_{\text{lim}}$  and  $\bar{\omega} \leq 0$  we just set our r estimate to  $r_{\text{lim}}$  better than  $\infty$

→ doesn't blow up any more

$$(4) make prior physical P(r) = \begin{cases} \frac{3}{r_{\text{lim}}^3} r^2 & 0 < r \leq r_{\text{lim}} \\ 0 & \text{otherwise} \end{cases}$$

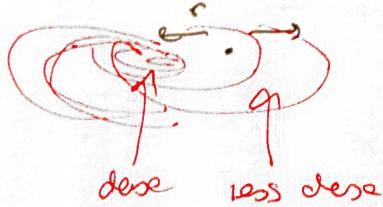


uniform density of stars up to  $r_{\text{lim}}$

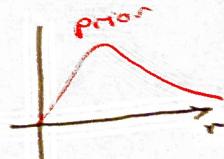


→ now freq pop very well behaved

(5) To make even more physical, include the fact that stars density decreases at outer parts of galaxy



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



no cutoff → no hard edge

again well behaved in freq properties