

- To get frequentist properties of estimator have to draw from a distribution, can test diff priors vs diff distrib.
- Using priors: encode background information
⇒ also want to test sensitivity of inferences to model

$$\text{POSTERIOR EXPECTATION } \mathbb{E}[f(\theta) | D] = \int f(\theta) P(\theta | D) d\theta$$

This is what we want to get, but computationally difficult

⇒ here Bayesian computation methods
"map out" / sample posterior density to compute expectation

Multi-parameter Bayesian inference

Gaussian example: Gelman BDA sec 3.2 - 3.3

$$Y_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2), i=1, \dots, n$$

$$\text{so } P(y | \mu, \sigma^2) = \prod_{i=1}^n N(y_i | \mu, \sigma^2) = (2\pi\sigma^2)^{-n} \exp \left\{ -\frac{1}{2\sigma^2} [n-1] s^2 + n(\bar{y} - \mu)^2 \right\}$$

↑
sample variance of y

- Prior: $P(\mu | G^2) = P(\mu | G^2) P(G^2)$

$$\mu | G^2 \sim N(\mu_0, \frac{G^2}{K_0})$$

↑
sample mean of y

how much information here about μ

and $G^2 \sim \text{Inv } \chi^2(\nu_0, G_0^2)$ inverse χ^2 distribution

$$P(\mu, G^2 | y) \propto P(y | \mu, G^2) P(\mu | G^2) P(G^2)$$

$$= N(\mu | \mu_n, \frac{G_n^2}{K_n}) \cdot \text{Inv } \chi^2(G^2 | \nu_n, G_n^2)$$

combine information of y and prior

where $K_n = K_0 + n$ $\mu_n = \frac{K_0}{K_0 + n} + \frac{n\bar{y}}{K_0 + n}$ $\nu_n G_n^2 = \nu_0 G_0^2 + (n-1)s^2 + \frac{K_0}{K_0 + n} (\bar{y} - \mu)^2$

$$\nu_n = \nu_0 + n$$

$$\underline{n \rightarrow \infty} \quad \mu_n \rightarrow \bar{y} \quad G_n^2 \rightarrow s^2 \quad K_n \rightarrow n \quad 2n \rightarrow n$$

$$\text{so } P(\mu, G^2 | y) = N(\mu | \bar{y}, G_n^2) \cdot \text{Inv}X^2(G^2 | n, s^2)$$

↑ similar to inverse P

now integrate out variables:

$$P(G^2 | y) = \int P(\mu | G^2, y) P(G^2 | y) d\mu = \text{Inv}X^2(G^2 | 2n, G_n^2)$$

$$\begin{aligned} P(\mu | y) &= \int P(\mu | t, G^2 | y) dG^2 = \dots \propto \left[1 + \frac{Kn(\bar{y} - \mu)^2}{2n G_n^2} \right]^{-(2n+1)/2} \\ &= t_{2n}(\mu | \bar{y}, G_n^2 / Kn) \quad t \text{ distribution!} \end{aligned}$$

$$\text{and as } n \rightarrow \infty \rightarrow t_n(\mu | \bar{y}, s^2/n) \rightarrow N(\mu | \bar{y}, s^2/n)$$

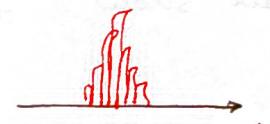
t distribution goes to gaussian

- What do we do if we can't derive this?

⇒ Monte Carlo direct sampling

use factorising, ie that we know posterior is normal × something

draw from both factors of posterior



and make histogram of posterior samples

$$\underline{\theta_i} \sim P(\underline{\theta} | D) \quad \hat{I} = \frac{1}{m} \sum_{i=1}^m f(\underline{\theta}_i) \rightarrow \text{expectation for large } m$$

- f could be mean, variance, posterior prob over interval

NB ^{here} posterior expectation is not averaging over data
but over parameters

- For above: $\underline{G^2} \sim P(G^2 | y)$ } look at joint distribution
 $\mu | G^2 \sim P(\mu | G^2, y)$ } of samples

KDE kernel density estimate smooths density estimate

put a gaussian + bandwidth on each point ⇒ add them all up



- When if you can't sample posterior? MCMC, Nested, etc.
for long-run generation
- $(\mu_0, \sigma_0^2)^T \times \pi_0 = (\mu_0, \sigma_0^2) \pi_0(\mu_0, \sigma_0^2)$
- ### Importance sampling
- draw from a simpler distribution

- have probability density $P(G)$ — intractable can write down but can't easily draw sample

want to estimate $I = \mathbb{E}[f(G)] = \int f(G) P(G) dG$

- make $Q(G)$ instrumental distib. / importance function
easy to draw from $G_i \stackrel{iid}{\sim} Q(G)$

important: $Q(G) > 0$ whenever $P(G) > 0$

$$I = \mathbb{E}_P[f(G)] = \int f(G) \frac{P(G)}{Q(G)} Q(G) dG \approx \frac{1}{m} \sum_{i=1}^m \frac{P(G_i)}{Q(G_i)} f(G_i)$$

$$\hat{I} = \frac{1}{m} \sum_{i=1}^m W(G_i) f(G_i)$$

IMPORTANCE WIDTH $\frac{P(G_i)}{Q(G_i)}$

\mathbb{E}_Q expectation over $Q(G)$

- can show $\mathbb{E}_Q[\hat{I}] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_Q[f(G_i) W(G_i)]$

= $\frac{1}{m} \sum_{i=1}^m \underbrace{\mathbb{E}_P[f(G_i)]}_{I}$

① Self Normalised

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

$P(D)$ is called $Z_P = \int \tilde{P}(G|D) dG$

evidence / partition function

- often only have $\tilde{P}(G|D) = \underbrace{P(D|G) P(G)}_{\equiv \tilde{P}(G)} \quad \text{don't know } P(D)!$

- Similarly to above:

$$I = \mathbb{E}[f] = \int f(G) P(G|D) dG = \int f(G) \frac{\tilde{P}(G|D)}{\int \tilde{P}(G|D) dG} dG$$

$$= \frac{\int f(G) \tilde{P}(G)/Q(G) Q(G) dG}{\int \tilde{P}(G)/Q(G) Q(G) dG}$$

$$\Rightarrow \hat{I} = \frac{\sum_{i=1}^m f(G_i) \tilde{W}(G_i)}{\sum_{i=1}^m \tilde{W}(G_i)}$$

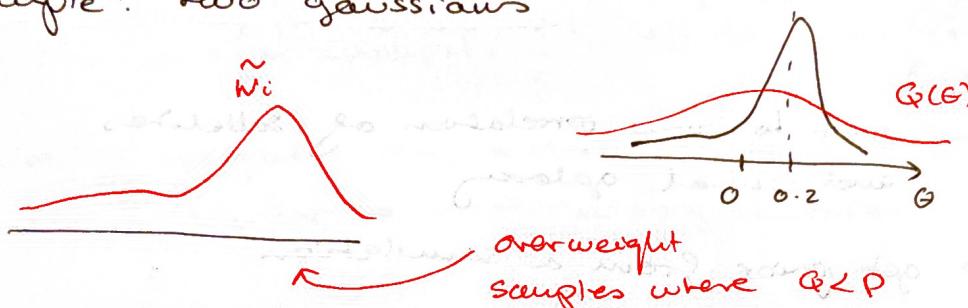
where $\tilde{W}_i = \frac{\tilde{P}(G_i)}{Q(G_i)}$ $G_i \stackrel{iid}{\sim} Q(G)$ but

$$= \sum_{i=1}^m f(G_i) W(G_i)$$

$$W_i = \frac{\tilde{W}_i}{\sum_{i=1}^m \tilde{W}_i} = \frac{\tilde{P}(G_i) / Q(G_i)}{\sum_{i=1}^m \tilde{P}(G_i) / Q(G_i)}$$

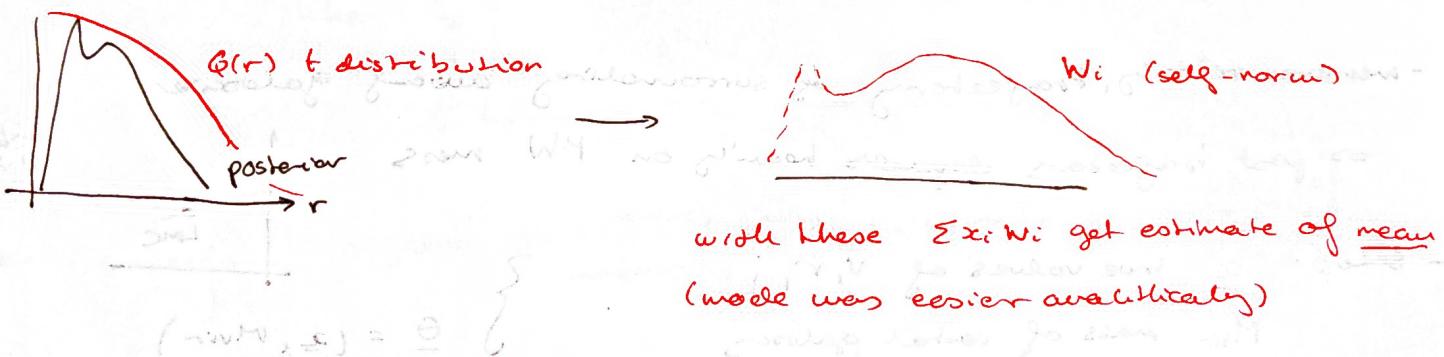
SELF-NORMALISED IMPORTANCE WEIGHT

- Example: two gaussians



~~from Q too many samples of $Q > P$, too little samples $Q < P$~~

- Example: in parallax example can't compute physical prior can't normalise posterior



NB error doesn't go down as \sqrt{N} bc samples are weighted. (from Q , not P)

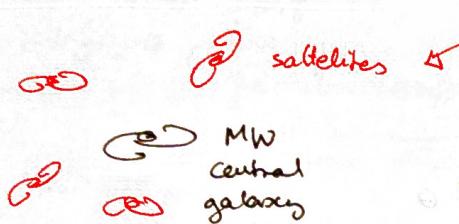
- Choosing G : theoretically want $\Phi = \frac{\int f(G) P(G)}{\int f(G) P(G) dG}$ but usually not possible

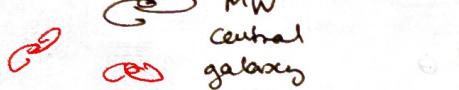
- \approx const W_i otherwise get too large variance
 - sample size is also diff bc of weighting!
 - increased by distribution of weighting
- find thick-tailed distribution infinities make variances large

$$\hat{I} = \frac{1}{m} \sum f(\theta_i) \longrightarrow I \quad \text{by LLN}$$

$$\text{error: } \text{Var}[\hat{I}] = \frac{1}{m} \text{Var}[f(\theta)] \propto \frac{1}{m} \text{Var}[\sum f(\theta_i)]$$

Milky Way Mass

 dwarf satellite galaxies
look at their mass, momentum, position

 want to know correlation of satellites
and central galaxy

\Rightarrow get prior from a simulation

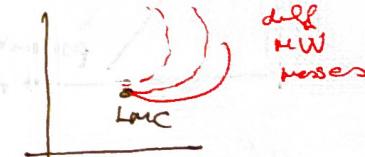
Illustris  simulates a box of the universe
with Dark Matter + normal matter
 \Rightarrow clumps into galaxies

 look at correlation
for all these simulated

- measure velocity, trajectory of surrounding dwarf galaxies

\Rightarrow past trajectory depends heavily on MW mass

- set up: \underline{x} true values of v, r, j $\underline{\theta}$ $\{$ $\underline{v}_{\text{obs}}$ $\underline{r}_{\text{obs}}$ $\underline{j}_{\text{obs}}$ $\}$
 M_{vir} mass of central galaxy $\} \quad \underline{\theta} = (\underline{z}, M_{\text{vir}})$



our prior can't be evaluated! only simulated

- measurement error likelihood: $L(\underline{z} | \underline{d})$ some normal error

\underline{d} are measurements of \underline{z}

- $P(\underline{z}, M_{\text{vir}} | \underline{d}) \propto P(\underline{d} | \underline{z}) \cdot P(\underline{z}, M_{\text{vir}})$

\uparrow observations \uparrow prior

have samples from prior as points in simulation

$$\int f(\underline{z}) P(\underline{z}, M_{\text{vir}} | \underline{d}) d\underline{z} \propto \frac{\sum f(\underline{z}_j) P(\underline{d} | \underline{z}_j)}{\sum P(\underline{d} | \underline{z}_j)}$$

- In this case f only depends on M_{vir}

$$\int f(M_{vir}) P(M_{vir} | \underline{z}) dM_{vir} \approx \sum f(M_{vir,j}) w_j$$

$$w_j = P(\underline{z} | \underline{x}_j) / \sum P(\underline{z} | \underline{x}_j) \quad \text{proportional to likelihood of each } \cancel{\text{sample}}$$

here no choice for Φ , given by simulation

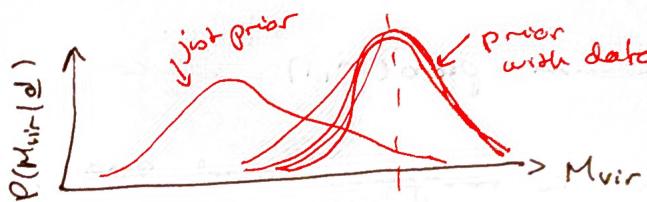
\Rightarrow never have to evaluate it because it cancels out!

- importance weight: only a few values contribute! systems similar to our galaxy contribute, others not much

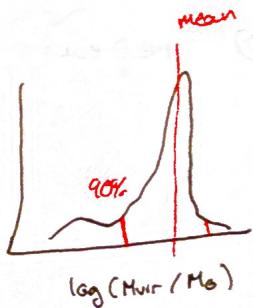
\Rightarrow can use these to calculate a weighted KDE

$$w\text{hole}(G) = \sum w_s \cdot N(G | G_s, bw^2) \quad \begin{matrix} \text{some} \\ \text{role of} \\ \text{bandwidth} \end{matrix}$$

make graphs like:



get μ_{MW} :
 $M_{vir} = 1.7^{+1.33}_{-0.52} \cdot 10^{12}$



- can calculate highest posterior density intervals

NB in paper don't use log so by Jensen's inequality
 get diff value for $\log M_{vir}$ than from us

- using multiple satellites get a better estimate

Markov Chain Monte Carlo

need to use this for higher number of parameters

- need to generate samples from a posterior without Q
- generate a chain of RVs that in a limit are draws from posterior
next value depends on current value! correlated

Metropolis

- ① choose μ_0 initial value choose directly
- ② propose $\mu_{\text{prop}} \sim N(\mu_i; \Sigma^2)$
- ③ evaluate $r = \frac{P(\mu_{\text{prop}} | y)}{P(\mu_i | y)}$
- ④ $r \geq 1$ μ_{prop} better $\mu_{i+1} = \mu_{\text{prop}}$
else accept with $\text{prob}(r, 1)$ so could accept even if not better
- ⑤ repeat until convergence and enough data