

# **Example Class 4**

## Astrostatistics (Part III)

# Future

- Revision Class scheduled for Thu, 14 May @ 3:30pm (online, remote)
- basic online open book pass/fail exam with straightforward questions in June
- another one in September?
- Conventional classed (distinction/merit) exam when possible

# Example Sheet 4

## 1 Probabilistic Graphical Models and Gibbs Sampling

Consider linear regression of the quasars' X-ray spectral indices vs. bolometric luminosities in the presence of heteroskedastic measurement error in both quantities and intrinsic dispersion. Consider the probabilistic generative model described in class:

$$\xi_i | \mu, \tau^2 \sim N(\mu, \tau^2) \quad (1)$$

$$\eta_i | \xi_i; \alpha, \beta, \sigma^2 \sim N(\alpha + \beta \xi_i, \sigma^2) \quad (2)$$

$$x_i | \xi_i \sim N(\xi_i, \sigma_{x,i}^2) \quad (3)$$

$$y_i | \eta_i \sim N(\eta_i, \sigma_{y,i}^2) \quad (4)$$

The astronomer observes values  $\mathcal{D} = \{x_i, y_i\}$ , which are noisy measurements of the true luminosity  $\xi_i$  and the true spectral index  $\eta_i$  of each quasar. The measurement errors are independent and heteroskedastic with known variances  $\{\sigma_{x,i}^2, \sigma_{y,i}^2\}$ , for  $i = 1, \dots, N$  independent quasars.

# 1 Probabilistic Graphical Models and Gibbs Sampling

Consider linear regression of the quasars' X-ray spectral indices vs. bolometric luminosities in the presence of heteroskedastic measurement error in both quantities and intrinsic dispersion. Consider the probabilistic generative model described in class:

$$\xi_i | \mu, \tau^2 \sim N(\mu, \tau^2) \quad (1)$$

$$\eta_i | \xi_i; \alpha, \beta, \sigma^2 \sim N(\alpha + \beta \xi_i, \sigma^2) \quad (2)$$

$$x_i | \xi_i \sim N(\xi_i, \sigma_{x,i}^2) \quad (3)$$

$$y_i | \eta_i \sim N(\eta_i, \sigma_{y,i}^2) \quad (4)$$

The astronomer observes values  $\mathcal{D} = \{x_i, y_i\}$ , which are noisy measurements of the true luminosity  $\xi_i$  and the true spectral index  $\eta_i$  of each quasar. The measurement errors are independent and heteroskedastic with known variances  $\{\sigma_{x,i}^2, \sigma_{y,i}^2\}$ , for  $i = 1, \dots, N$  independent quasars.

1. Write down the joint distribution  $P(x_i, y_i, \xi_i, \eta_i | \alpha, \beta, \sigma^2, \mu, \tau^2)$  for a single quasar.

**Strategy: Factor joint pdf into conditional and marginal distributions and exploit conditional independence structure of model assumptions.**

**Solution:**

$$\begin{aligned} P(x_i, y_i, \xi_i, \eta_i | \alpha, \beta, \sigma^2, \mu, \tau^2) &= P(x_i, y_i | \xi_i, \eta_i)P(\xi_i, \eta_i | \alpha, \beta, \sigma^2, \mu, \tau^2) \\ &= P(x_i | \xi_i)P(y_i | \eta_i)P(\eta_i | \xi_i; \alpha, \beta, \sigma^2)P(\xi_i | \mu, \tau^2) \\ &= N(x_i | \xi_i, \sigma_{x,i}^2)N(y_i | \eta_i, \sigma_{y,i}^2)N(\eta_i | \alpha + \beta \xi_i, \sigma^2)N(\xi_i | \mu, \tau^2) \end{aligned} \quad (45)$$

Consider the probabilistic generative model described in class:

$$\xi_i | \mu, \tau^2 \sim N(\mu, \tau^2) \quad (1)$$

$$\eta_i | \xi_i; \alpha, \beta, \sigma^2 \sim N(\alpha + \beta \xi_i, \sigma^2) \quad (2)$$

$$x_i | \xi_i \sim N(\xi_i, \sigma_{x,i}^2) \quad (3)$$

$$y_i | \eta_i \sim N(\eta_i, \sigma_{y,i}^2) \quad (4)$$

The astronomer observes values  $\mathcal{D} = \{x_i, y_i\}$ , which are noisy measurements of the true luminosity  $\xi_i$  and the true spectral index  $\eta_i$  of each quasar. The measurement errors are independent and heteroskedastic with known variances  $\{\sigma_{x,i}^2, \sigma_{y,i}^2\}$ , for  $i = 1, \dots, N$  independent quasars.

2. Adopt “non-informative” hyperpriors on the hyperparameters: flat improper priors for each of  $P(\alpha)$ ,  $P(\beta)$ ,  $P(\mu)$ , and flat positive improper priors for each of  $P(\tau^2)$  and  $P(\sigma^2)$ . Write down the full joint distribution of all data  $\mathcal{D}$ , latent variables  $\{\xi_i, \eta_i\}$ , and hyperparameters  $\alpha, \beta, \sigma^2, \mu, \tau^2$ .

**Solution:** We assume independent, improper flat priors on  $\alpha, \beta, \mu$ :  $P(\alpha) \propto 1$ ,  $P(\beta) \propto 1$ ,  $P(\mu) \propto 1$ . We assume flat positive priors on  $\sigma^2, \tau^2$ :  $P(\sigma^2) \propto 1, \sigma^2 > 0$ ,  $P(\tau^2) \propto 1, \tau^2 > 0$ . Then the joint distribution:

$$P(\{x_i, y_i\}, \{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2) = \left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i | \eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \times P(\alpha) P(\beta) P(\sigma^2) P(\mu) P(\tau^2) \quad (1)$$

**Utilise conditional independence of individual quasars**

$$\text{Let } H(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases} \quad \begin{aligned} P(\tau^2) &\propto H(z) \\ P(\sigma^2) &\propto H(z) \end{aligned}$$

Consider the probabilistic generative model described in class:

$$\xi_i | \mu, \tau^2 \sim N(\mu, \tau^2) \quad (1)$$

$$\eta_i | \xi_i; \alpha, \beta, \sigma^2 \sim N(\alpha + \beta \xi_i, \sigma^2) \quad (2)$$

$$x_i | \xi_i \sim N(\xi_i, \sigma_{x,i}^2) \quad (3)$$

$$y_i | \eta_i \sim N(\eta_i, \sigma_{y,i}^2) \quad (4)$$

The astronomer observes values  $\mathcal{D} = \{x_i, y_i\}$ , which are noisy measurements of the true luminosity  $\xi_i$  and the true spectral index  $\eta_i$  of each quasar. The measurement errors are independent and heteroskedastic with known variances  $\{\sigma_{x,i}^2, \sigma_{y,i}^2\}$ , for  $i = 1, \dots, N$  independent quasars.

2. Adopt “non-informative” hyperpriors on the hyperparameters: flat improper priors for each of  $P(\alpha)$ ,  $P(\beta)$ ,  $P(\mu)$ , and flat positive improper priors for each of  $P(\tau^2)$  and  $P(\sigma^2)$ . Write down the full joint distribution of all data  $\mathcal{D}$ , latent variables  $\{\xi_i, \eta_i\}$ , and hyperparameters  $\alpha, \beta, \sigma^2, \mu, \tau^2$ .

$$P(\{x_i, y_i\}, \{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i | \eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \times H(\sigma^2) H(\tau^2)$$

Consider the probabilistic generative model described in class:

$$\xi_i | \mu, \tau^2 \sim N(\mu, \tau^2) \quad (1)$$

$$\eta_i | \xi_i; \alpha, \beta, \sigma^2 \sim N(\alpha + \beta \xi_i, \sigma^2) \quad (2)$$

$$x_i | \xi_i \sim N(\xi_i, \sigma_{x,i}^2) \quad (3)$$

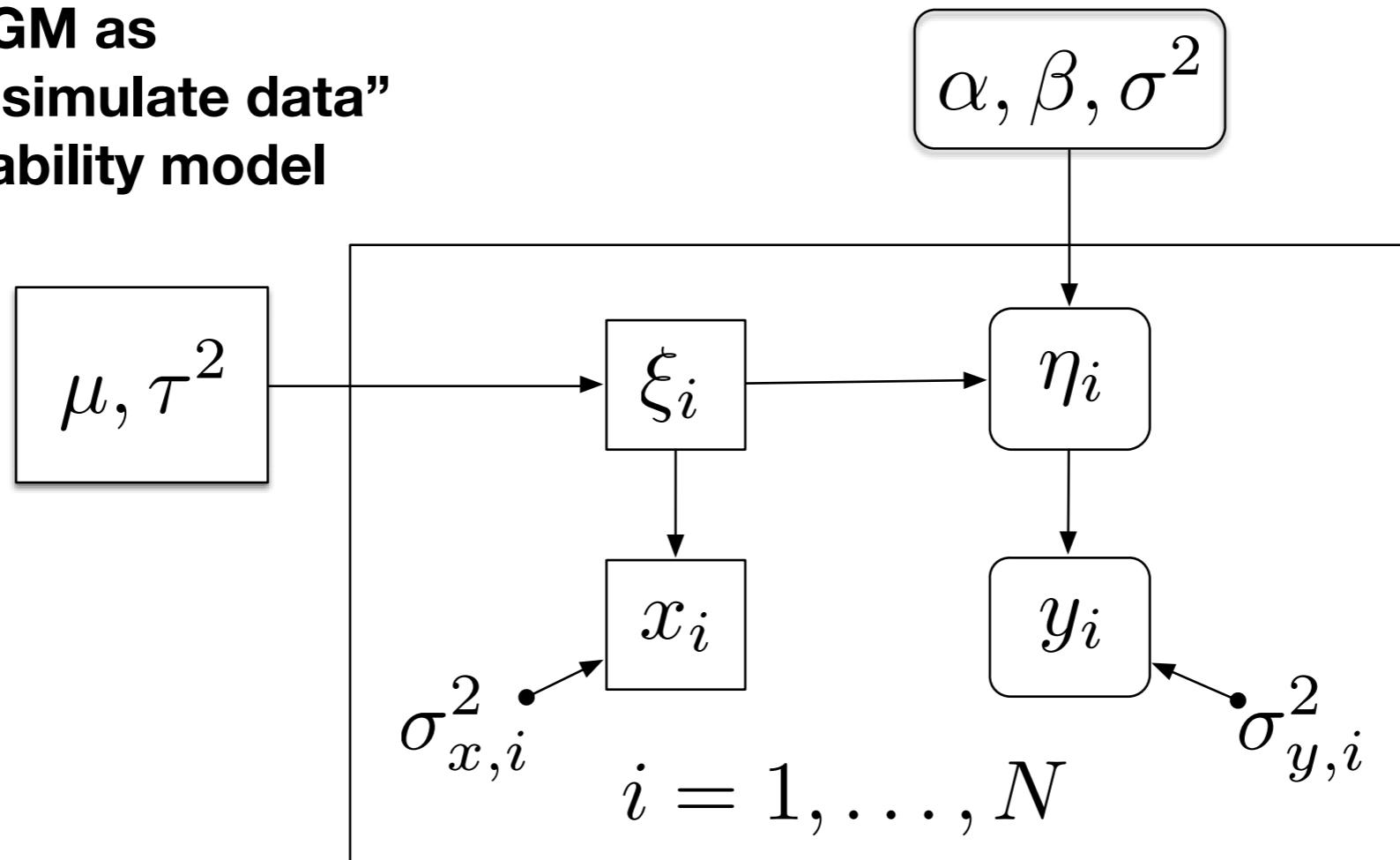
$$y_i | \eta_i \sim N(\eta_i, \sigma_{y,i}^2) \quad (4)$$

The astronomer observes values  $\mathcal{D} = \{x_i, y_i\}$ , which are noisy measurements of the true luminosity  $\xi_i$  and the true spectral index  $\eta_i$  of each quasar. The measurement errors are independent and heteroskedastic with known variances  $\{\sigma_{x,i}^2, \sigma_{y,i}^2\}$ , for  $i = 1, \dots, N$  independent quasars.

3. Draw a probabilistic graphical model to represent this joint distribution.

$$P(\{x_i, y_i\}, \{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2) = \left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i | \eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \times P(\alpha) P(\beta) P(\sigma^2) P(\mu) P(\tau^2)$$

**Draw PGM as  
“how you would simulate data”  
from the probability model**



$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

4. Construct a Gibbs sampler for the posterior  $P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D})$  by deriving a sequence of proposed moves that are always accepted. Specify the order in which you run through your sequence. You have access to algorithms that generate random draws from univariate and multivariate Gaussian distributions, and scaled inverse- $\chi^2$  distributions.
- Derive a complete set of conditional posterior densities of each (subset of) parameters conditional on all the others (and the data)
  - For hierarchical models, the conditional independence structure enables a natural strategy for Gibbs sampling in **two stages**:
    - For each individual  $i$ , sample the latent variables  $(\xi_i, \eta_i)$  given the current values of the hyperparameters  $(\alpha, \beta, \sigma^2, \mu, \tau)$
    - Then sample the hyperparameters  $(\alpha, \beta, \sigma^2, \mu, \tau)$  given all the latent variables for all  $i$   $(\xi, \eta)$

$$\begin{aligned}
P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) &\propto \\
&\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\
&\times H(\sigma^2) H(\tau^2)
\end{aligned}$$

4. Construct a Gibbs sampler for the posterior  $P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D})$  by deriving a sequence of proposed moves that are always accepted. Specify the order in which you run through your sequence. You have access to algorithms that generate random draws from univariate and multivariate Gaussian distributions, and scaled inverse- $\chi^2$  distributions.
- Derive a complete set of conditional posterior densities of each (subset of) parameters conditional on all the others (and the data)
  - For hierarchical models, the conditional independence structure enables a natural strategy for Gibbs sampling in **two stages**:
    - For each individual  $i$ , sample the latent variables given the current values of the hyperparameters
$$P(\xi_i, \eta_i | \alpha, \beta, \sigma^2, \mu, \tau; \mathcal{D})$$
  - Then sample the hyperparameters given all the latent variables for all  $i$
- $$P(\alpha, \beta, \sigma^2, \mu, \tau | \boldsymbol{\xi}, \boldsymbol{\eta}; \mathcal{D})$$

**Each Stage can then be further broken up by inspecting conditional structure**

$$\begin{aligned}
P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) &\propto \\
&\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\
&\times H(\sigma^2) H(\tau^2)
\end{aligned}$$

**Drawing new latent variables, conditional on hyperparameters:**

- (a) For each individual quasar  $i = 1, \dots, N$ , we update/sample the latent variables  $\xi_i, \eta_i$  from

$$\begin{aligned}
P(\xi_i, \eta_i | \alpha, \beta, \sigma^2, \mu, \tau^2) &\propto N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \\
&\propto N(\phi_i | \mathbf{z}_i, \Sigma_i) \times N(\phi_i | \boldsymbol{\gamma}, \mathbf{C})
\end{aligned} \tag{10}$$

where the second line is obtained using the properties of the multivariate Gaussian and  $\phi_i \equiv (\eta_i, \xi_i)^T$ ,  $\boldsymbol{\gamma} \equiv (\alpha + \beta \mu, \mu)^T$  and

$$\mathbf{C} \equiv \begin{pmatrix} \beta^2 \tau^2 + \sigma^2 & \beta \tau^2 \\ \beta \tau^2 & \tau^2 \end{pmatrix} \tag{11}$$

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

$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

## Drawing new latent variables, conditional on hyperparameters:

- (a) For each individual quasar  $i = 1, \dots, N$ , we update/sample the latent variables  $\xi_i, \eta_i$  from

$$P(\xi_i, \eta_i | \alpha, \beta, \sigma^2, \mu, \tau^2) \propto N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \\ \propto N(\phi_i | z_i, \Sigma_i) \times N(\phi_i | \gamma, C)$$

$\phi_i \equiv \begin{pmatrix} \eta_i \\ \xi_i \end{pmatrix}$

The product of two Gaussian densities is proportional to a single Gaussian density, whose mean is the precision-weighted average of the individual means, and whose precision (matrix) is the sum of the individual precision (matrices). Precision is the inverse of the variance, and a precision matrix is the inverse of a covariance matrix. Therefore let

$$\bar{\phi}_i \equiv (C^{-1} + \Sigma^{-1})^{-1} (C^{-1} \gamma + \Sigma_i^{-1} z_i) \quad (13)$$

$$V_\phi^i = (C^{-1} + \Sigma_i^{-1})^{-1} \quad (14)$$

Thus  $P(\phi | \dots, \mathcal{D}) = N(\phi | \bar{\phi}, V_\phi^i)$ . Therefore we make the following draw

$$\begin{pmatrix} \eta_i \\ \xi_i \end{pmatrix} \sim N(\bar{\phi}_i, V_\phi^i) \quad (15)$$

for each quasar  $i$ .

$$\begin{aligned}
P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) &\propto \\
&\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\
&\times H(\sigma^2) H(\tau^2)
\end{aligned}$$

**Drawing new hyperparameters, conditional on the updated latent variables.**  
**Break up into subsets:**

**Update Sequentially:**

$$(b) \sigma^2 \sim P(\sigma^2 | \dots; \mathcal{D}) = P(\sigma^2 | \alpha, \beta, \{\eta_i, \xi_i\})$$

$$(c) \alpha \sim P(\alpha | \dots; \mathcal{D}) = P(\alpha | \beta, \sigma^2, \{\eta_i, \xi_i\})$$

$$(d) \beta \sim P(\beta | \dots; \mathcal{D}) = P(\beta | \alpha, \sigma^2, \{\eta_i, \xi_i\})$$

$$(e) \mu, \tau^2 \sim P(\mu, \tau^2 | \dots; \mathcal{D}) = P(\mu | \tau^2 \{\xi_i\}) P(\tau^2 | \{\xi_i\})$$

**(draw a block of two parameters)**

$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

**Update:**

$$(b) \sigma^2 \sim P(\sigma^2 | \dots; \mathcal{D}) = P(\sigma^2 | \alpha, \beta, \{\eta_i, \xi_i\})$$

(b) Next we update  $\sigma^2$  from  $P(\sigma^2 | \dots \mathcal{D})$ .

$$P(\sigma^2 | \dots \mathcal{D}) \propto \prod_{i=1}^N N(\eta_i | \alpha + \beta \xi_i, \sigma^2) \\ \propto \sigma^{-N} \prod_{i=1}^N \exp \left( -\frac{1}{2} (\eta_i - \alpha - \beta \xi_i)^2 / \sigma^2 \right) \\ \propto (\sigma^2)^{-N/2} \exp \left( -\frac{1}{2} \text{SSR} / \sigma^2 \right) \\ = \text{Inv-}\chi^2(\sigma^2 | (N-2), \text{SSR}/(N-2)) \quad (16)$$

where  $\text{SSR} \equiv \sum_{i=1}^N (\eta_i - \alpha - \beta \xi_i)^2$ . Thus, we draw a new  $\sigma^2$  from the above scaled inverse  $\chi^2$  distribution.

$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

**Update:**

$$(c) \quad \alpha \sim P(\alpha | \dots; \mathcal{D}) = P(\alpha | \beta, \sigma^2, \{\eta_i, \xi_i\})$$

(c) Next we update  $\alpha$  from  $P(\alpha | \dots; \mathcal{D})$ .

$$\begin{aligned} P(\alpha | \dots; \mathcal{D}) &\propto \prod_{i=1}^N N(\eta_i | \alpha + \beta \xi_i, \sigma^2) \\ &\propto \prod_{i=1}^N \exp \left( -\frac{1}{2} (\eta_i - \alpha - \beta \xi_i)^2 / \sigma^2 \right) \\ &\propto \prod_{i=1}^N \exp \left( -\frac{1}{2} (\alpha + \beta \xi_i - \eta_i)^2 / \sigma^2 \right) \\ &\propto \prod_{i=1}^N N(\alpha | \eta_i - \beta \xi_i, \sigma^2). \end{aligned} \tag{17}$$

$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

$$(c) \quad \alpha \sim P(\alpha | \dots; \mathcal{D}) = P(\alpha | \beta, \sigma^2, \{\eta_i, \xi_i\})$$

$$\propto \prod_{i=1}^N N(\alpha | \eta_i - \beta \xi_i, \sigma^2)$$

Now we use the fact that the product of  $N$  Gaussian densities (in  $\alpha$ ) is proportional to a single Gaussian density (in  $\alpha$ ). The resulting mean is the precision-weighted average of the individual means, and the resulting precision (inverse variance) is the sum of the individual precisions. In the case, all the precision weights are equal (to  $\sigma^{-2}$ ).

$$P(\alpha | \dots; \mathcal{D}) = N\left(\alpha \middle| N^{-1} \sum_{i=1}^N (\eta_i - \beta \xi_i), \sigma^2/N\right) \quad (18)$$

Thus, we draw a new  $\alpha$  from this Gaussian distribution.

$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

**Update:** (d)  $\beta \sim P(\beta | \dots; \mathcal{D}) = P(\beta | \alpha, \sigma^2, \{\eta_i, \xi_i\})$

(d) Next we update  $\beta$  from  $P(\beta | \dots \mathcal{D})$ .

$$P(\beta | \dots \mathcal{D}) \propto \prod_{i=1}^N N(\eta_i | \alpha + \beta \xi_i, \sigma^2) \\ \propto \prod_{i=1}^N \exp \left( -\frac{1}{2} (\eta_i - \alpha - \beta \xi_i)^2 / \sigma^2 \right) \\ \propto \prod_{i=1}^N \exp \left( -\frac{1}{2} (\beta + \alpha/\xi_i - \eta_i/\xi_i)^2 / (\sigma/\xi_i)^2 \right) \\ \propto \prod_{i=1}^N N(\beta | (\eta_i - \alpha)/\xi_i, \delta_i^2 \equiv \sigma^2/\xi_i^2)$$

$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

$$(d) \quad \beta \sim P(\beta | \dots; \mathcal{D}) = P(\beta | \alpha, \sigma^2, \{\eta_i, \xi_i\})$$

$$\propto \prod_{i=1}^N N(\beta | (\eta_i - \alpha)/\xi_i, \delta_i^2 \equiv \sigma^2/\xi_i^2)$$

We can apply the “product-of-Gaussians” rule again to derive that this conditional posterior density of  $\beta$  is  $P(\beta | \dots, \mathcal{D}) = N(\beta | \bar{\beta}, V_\beta)$  where

$$\bar{\beta} = \frac{\sum_{i=1}^N \delta_i^{-2} (\eta_i - \alpha) / \xi_i}{\sum_{i=1}^N \delta_i^{-2}} \quad (20)$$

$$V_\beta = \left( \sum_{i=1}^N \delta_i^{-2} \right)^{-1} \quad (21)$$

Thus, we draw a new  $\beta$  from the above Gaussian density.

$$P(\{\xi_i, \eta_i\}, \alpha, \beta, \sigma^2, \mu, \tau^2 | \mathcal{D}) \propto$$

$$\left[ \prod_{i=1}^N N(x_i | \xi_i, \sigma_{x,i}^2) N(y_i, |\eta_i, \sigma_{y,i}^2) N(\eta_i | \alpha + \beta \xi_i, \sigma^2) N(\xi_i | \mu, \tau^2) \right] \\ \times H(\sigma^2) H(\tau^2)$$

$$(e) \mu, \tau^2 \sim P(\mu, \tau^2 | \dots; \mathcal{D}) = P(\mu | \tau^2 \{\xi_i\}) P(\tau^2 | \{\xi_i\})$$

(e) Finally, we update  $\mu, \tau^2$  from  $P(\mu, \tau^2 | \dots, \mathcal{D}) = P(\mu, \tau^2 | \{\xi_i\})$ . This is essentially the posterior of the mean and variance of a Gaussian given draws from that Gaussian. We can use the same approach as in Problem 1.3 of Example Sheet 2.

$$P(\mu, \tau^2 | \{\xi_i\}) \propto \prod_{i=1}^N N(\xi_i | \mu, \tau^2) \quad (22)$$

Using the same derivation, we can factor

$$\begin{aligned} P(\mu, \tau^2 | \{\xi_i\}) &= P(\mu | \tau^2, \{\xi_i\}) \times P(\tau^2 | \{\xi_i\}) \\ &= N(\mu | \bar{\xi}, \tau^2/N) \times \text{Inv-}\chi^2 \left( \tau^2 \middle| (N-3), \frac{(N-1)}{(N-3)} S_\xi^2 \right) \end{aligned} \quad (23)$$

where  $\bar{\xi} = N^{-1} \sum_{i=1}^N \xi_i$

$$S_\xi^2 = \frac{1}{N-1} \sum_{i=1}^N (\xi_i - \bar{\xi})^2. \quad (24)$$

Hence we generate a joint draw of  $\mu, \tau^2$  from their conditional by first drawing  $\tau^2$  from the scaled inverse- $\chi^2$  and then  $\mu | \tau^2, \bar{\xi}$  from the Gaussian.

$$(e1) \tau^2 \sim \text{Inv-}\chi^2 \left( \tau^2 \middle| (N-3), \frac{(N-1)}{(N-3)} S_\xi^2 \right) \quad (e2) \mu | \tau^2 \sim N(\mu | \bar{\xi}, \tau^2/N)$$

# Summary

**Update Sequentially:**

**For every quasar i:** (a)<sub>i</sub>  $(\xi_i, \eta_i) \sim P(\xi_i, \eta_i | \alpha, \beta, \sigma^2, \mu, \tau^2; \mathcal{D})$

**Then update hyperparameters**

$$(b) \sigma^2 \sim P(\sigma^2 | \dots; \mathcal{D}) = P(\sigma^2 | \alpha, \beta, \{\eta_i, \xi_i\})$$

$$(c) \alpha \sim P(\alpha | \dots; \mathcal{D}) = P(\alpha | \beta, \sigma^2, \{\eta_i, \xi_i\})$$

$$(d) \beta \sim P(\beta | \dots; \mathcal{D}) = P(\beta | \alpha, \sigma^2, \{\eta_i, \xi_i\})$$

$$(e) \mu, \tau^2 \sim P(\mu, \tau^2 | \dots; \mathcal{D}) = P(\mu | \tau^2 | \{\xi_i\}) P(\tau^2 | \{\xi_i\})$$

$$(e1) \tau^2 \sim P(\tau^2 | \{\xi_i\})$$

$$(e2) \mu | \tau^2 \sim P(\mu | \tau^2, \{\xi_i\})$$

**Cycle through steps a-e many, many times through reaching convergence criteria  
and until reaching sufficient effective sample size**

# Implementation

- Initialisation
- Tuning?
- Assessing Convergence/Mixing - how long to run MCMC?
- Postprocessing
- Computing posterior inferences from MCMC samples

# Bayesian Model Comparison

1. Data points  $\{x_i\}$  come independently from a probability distribution  $P(x)$ . According to model  $H_0$ ,  $P(x)$  is a uniform distribution  $P(x|H_0) = \frac{1}{2}$  for  $x \in (-1, 1)$ . According to model  $H_1$ ,  $P(x)$  is a nonuniform distribution with an unknown parameter  $m \in (-1, 1)$ :

$$P(x|m, H_1) = \frac{1}{2}(1 + m x), \quad (25)$$

for  $x \in (-1, 1)$ . Given the data  $\mathcal{D} = \{0.3, 0.5, 0.7, 0.8, 0.9\}$ , what is the evidence for  $H_0$  and  $H_1$ ?

- Strategy is to compute the evidence (marginal likelihood) under each model  $H$ :

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

# Bayesian Model Comparison

1. Data points  $\{x_i\}$  come independently from a probability distribution  $P(x)$ . According to model  $H_0$ ,  $P(x)$  is a uniform distribution  $P(x|H_0) = \frac{1}{2}$  for  $x \in (-1, 1)$ . According to model  $H_1$ ,  $P(x)$  is a nonuniform distribution with an unknown parameter  $m \in (-1, 1)$ :

$$P(x|m, H_1) = \frac{1}{2}(1 + m x), \quad (25)$$

for  $x \in (-1, 1)$ . Given the data  $\mathcal{D} = \{0.3, 0.5, 0.7, 0.8, 0.9\}$ , what is the evidence for  $H_0$  and  $H_1$ ?

**Solution:** We note that all the observed data points are within the range of the sampling distribution.

Under model  $H_0$ , we have

$$P(\mathbf{x}|H_0) = \prod_{i=1}^N P(x_i|H_0) = \left(\frac{1}{2}\right)^N$$

**(no parameters to integrate over)**

# Bayesian Model Comparison

1. Data points  $\{x_i\}$  come independently from a probability distribution  $P(x)$ . According to model  $H_0$ ,  $P(x)$  is a uniform distribution  $P(x|H_0) = \frac{1}{2}$  for  $x \in (-1, 1)$ . According to model  $H_1$ ,  $P(x)$  is a nonuniform distribution with an unknown parameter  $m \in (-1, 1)$ :

$$P(x|m, H_1) = \frac{1}{2}(1 + mx), \quad (25)$$

for  $x \in (-1, 1)$ . Given the data  $\mathcal{D} = \{0.3, 0.5, 0.7, 0.8, 0.9\}$ , what is the evidence for  $H_0$  and  $H_1$ ? **Under model  $H_1$ ,**

$$\begin{aligned} P(\mathbf{x}|H_1) &= \int dm P(\mathbf{x}|m, H_1) P(m|H_1) \\ &= \int \left[ \prod_{i=1}^N P(x_i|m, H_1) P(m|H_1) \right] dm \\ &= \left( \frac{1}{2} \right)^N \int_{-1}^1 \prod_{i=1}^N (1 + mx_i) dm \\ &= \left( \frac{1}{2} \right)^4 \int_{-1}^1 (1 + mx_1)(1 + mx_2)(1 + mx_3)(1 + mx_4) dm \\ &= \left( \frac{1}{2} \right)^4 \int_{-1}^1 (1 + S_1 m + S_2 m^2 + S_3 m^3 + P m^4) dm \end{aligned}$$

$$S_2 \equiv x_1 x_2 + x_1 x_3 + x_1 x_4 + x_2 x_3 + x_2 x_4 + x_3 x_4 = 1.91$$

$$P \equiv x_1 x_2 x_3 x_4 = 0.0756 \quad (\mathbf{S}_1, \mathbf{S}_2 \text{ don't matter due to symmetry})$$

# Bayesian Model Comparison

1. Data points  $\{x_i\}$  come independently from a probability distribution  $P(x)$ . According to model  $H_0$ ,  $P(x)$  is a uniform distribution  $P(x|H_0) = \frac{1}{2}$  for  $x \in (-1, 1)$ . According to model  $H_1$ ,  $P(x)$  is a nonuniform distribution with an unknown parameter  $m \in (-1, 1)$ :

$$P(x|m, H_1) = \frac{1}{2}(1 + mx), \quad (25)$$

for  $x \in (-1, 1)$ . Given the data  $\mathcal{D} = \{0.3, 0.5, 0.7, 0.8, 0.9\}$ , what is the evidence for  $H_0$  and  $H_1$ ?

**Exploiting symmetries of the integrand, we find:**

$$\begin{aligned} P(\mathbf{x}|H_1) &= \left(\frac{1}{2}\right)^5 \int_{-1}^1 (1 + S_1m + S_2m^2 + S_3m^3 + Pm^4) dm \\ &= \frac{1}{16} [1 + S_2/3 + P/5] \end{aligned}$$

where  $S_2 \equiv x_1x_2 + x_1x_3 + x_1x_4 + x_2x_3 + x_2x_4 + x_3x_4 = 1.91$  and  $P \equiv x_1x_2x_3x_4 = 0.0756$ . Thus  $P(\mathbf{x}|H_1) = 0.0883$  and  $P(\mathbf{x}|H_0) = 0.0625$ , so the Bayes Factor  $B_{10} = P(\mathbf{x}|H_1)/P(\mathbf{x}|H_0) = 1.41$ , or  $\log B_{10} = 0.34$ . According to the Jeffrey's scale this is not significantly more evidence for model  $H_1$  vs. model  $H_0$ .

# Bayesian Model Comparison 2

Datapoints  $\mathcal{D} = \{x_i, y_i\}$  are believed to arisen from a straight line. The experimenter chooses  $x_i$ , and  $y_i$  is Gaussian-distributed around  $w_0 + w_1 x_i$  with variance  $\sigma^2$ . According to model  $H_0$ , the straight line is horizontal, so  $w_1 = 0$ . According to model  $H_1$ ,  $w_1$  is a parameter with prior distribution  $w_1 \sim N(0, \tau^2)$ . Both models assign a prior distribution  $w_0 \sim N(0, \tau^2)$ .

- (a) For each model, derive an expression for the posterior distribution of the regression parameter(s),  $P(w_0 | \mathcal{D}, H_0)$  and  $P(w_0, w_1 | \mathcal{D}, H_1)$ .
- (b) For each model, derive an expression for the evidence (or marginal likelihood).

**Solution:** For  $H_0$ , defining  $\bar{y} = (\sum_{i=1}^N y_i)/N$ , we have the posterior:

$$\begin{aligned} P(w_0 | \mathcal{D}) &\propto \left[ \prod_{i=1}^N N(y_i | w_0, \sigma^2) \right] N(w_0 | 0, \tau^2) \\ &\propto N(w_0 | \bar{y}, \sigma^2/N) N(w_0, \tau^2) \\ &= N\left[w_0 \middle| \frac{N\sigma^{-2} \bar{y}}{N\sigma^{-2} + \tau^{-2}}, (N\sigma^{-2} + \tau^{-2})^{-1}\right]. \end{aligned}$$

# Bayesian Model Comparison 2

Datapoints  $\mathcal{D} = \{x_i, y_i\}$  are believed to arisen from a straight line. The experimenter chooses  $x_i$ , and  $y_i$  is Gaussian-distributed around  $w_0 + w_1 x_i$  with variance  $\sigma^2$ . According to model  $H_0$ , the straight line is horizontal, so  $w_1 = 0$ . According to model  $H_1$ ,  $w_1$  is a parameter with prior distribution  $w_1 \sim N(0, \tau^2)$ . Both models assign a prior distribution  $w_0 \sim N(0, \tau^2)$ .

- (a) For each model, derive an expression for the posterior distribution of the regression parameter(s),  $P(w_0 | \mathcal{D}, H_0)$  and  $P(w_0, w_1 | \mathcal{D}, H_1)$ .
- (b) For each model, derive an expression for the evidence (or marginal likelihood).

The evidence  $Z_0$  can be computed by noting that the product of  $N$  independent 1D Gaussian sampling distributions can be re-expressed as a single  $N$ -dimensional multivariate Gaussian distribution with a diagonal covariance matrix.

**Use MV Gaussian marginalisation**

$$\begin{aligned} Z_0 &= \int \left[ \prod_{i=1}^N N(y_i | w_0, \sigma^2) \right] N(w_0 | 0, \tau^2) dw_0 \\ &= \int N(\mathbf{y} | \mathbf{1}_N w_0, \sigma^2 \mathbf{I}_N) N(w_0 | 0, \tau^2) dw_0 \\ &= N(\mathbf{y} | \mathbf{0}, \mathbf{C}_0 = \sigma^2 \mathbf{I}_N + \mathbf{1}_N \mathbf{1}_N^T \tau^2), \end{aligned}$$

where  $\mathbf{1}_N$  is an column vector of  $N$  ones, and  $\mathbf{I}_N$  is the  $N \times N$  dimensional identity matrix. A  $N \times N$  square matrix filled with ones results from  $\mathbf{1}_N \mathbf{1}_N^T$ .

# Bayesian Model Comparison 2

Datapoints  $\mathcal{D} = \{x_i, y_i\}$  are believed to arisen from a straight line. The experimenter chooses  $x_i$ , and  $y_i$  is Gaussian-distributed around  $w_0 + w_1 x_i$  with variance  $\sigma^2$ . According to model  $H_0$ , the straight line is horizontal, so  $w_1 = 0$ . According to model  $H_1$ ,  $w_1$  is a parameter with prior distribution  $w_1 \sim N(0, \tau^2)$ . Both models assign a prior distribution  $w_0 \sim N(0, \tau^2)$ .

- (a) For each model, derive an expression for the posterior distribution of the regression parameter(s),  $P(w_0 | \mathcal{D}, H_0)$  and  $P(w_0, w_1 | \mathcal{D}, H_1)$ .
- (b) For each model, derive an expression for the evidence (or marginal likelihood).

For  $H_1$ , let  $\mathbf{w} = (w_0, w_1)^T$  and  $\mathbf{D} = (\mathbf{1}_N, \mathbf{x})$

**Posterior:** 
$$P(\mathbf{w} | \mathcal{D}) \propto \left[ \prod_{i=1}^N N(y_i | w_0 + w_1 x_i, \sigma^2) \right] N(w_0 | 0, \tau^2)$$

$$\propto N(\mathbf{y} | \mathbf{D}\mathbf{w}, \sigma^2 \mathbf{I}_N) N(\mathbf{w} | \mathbf{0}, \tau^2 \mathbf{I}_2)$$

**Use result from OLS** 
$$\propto N(\mathbf{w} | (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{y}, \sigma^2 (\mathbf{D}^T \mathbf{D})^{-1}) \times N(\mathbf{w} | \mathbf{0}, \tau^2 \mathbf{I}_2)$$

**Product of Gaussian** 
$$\propto N(\mathbf{w} | \dots, \dots)$$

# Bayesian Model Comparison 2

Datapoints  $\mathcal{D} = \{x_i, y_i\}$  are believed to arisen from a straight line. The experimenter chooses  $x_i$ , and  $y_i$  is Gaussian-distributed around  $w_0 + w_1 x_i$  with variance  $\sigma^2$ . According to model  $H_0$ , the straight line is horizontal, so  $w_1 = 0$ . According to model  $H_1$ ,  $w_1$  is a parameter with prior distribution  $w_1 \sim N(0, \tau^2)$ . Both models assign a prior distribution  $w_0 \sim N(0, \tau^2)$ .

- (a) For each model, derive an expression for the posterior distribution of the regression parameter(s),  $P(w_0 | \mathcal{D}, H_0)$  and  $P(w_0, w_1 | \mathcal{D}, H_1)$ .
- (b) For each model, derive an expression for the evidence (or marginal likelihood).

The evidence  $Z_1$  is

**Using MV Gaussian  
marginalisation**

$$\begin{aligned} Z_1 &= \int N(\mathbf{y} | \mathbf{D}\mathbf{w}, \sigma^2 \mathbf{I}_N) N(\mathbf{w} | \mathbf{0}, \tau^2 \mathbf{I}_2) d\mathbf{w} \\ &= N(\mathbf{y} | \mathbf{0}, \mathbf{C}_1 = \sigma^2 \mathbf{I}_N + \mathbf{D}\mathbf{D}^T) \end{aligned}$$