



# Data-analysis and Inverse Methods in Astronomy

Gibbs sampling

Antti Penttilä

Homepage <http://wiki.helsinki.fi/display/53834>



# Gibbs sampling

The algorithm “started” the ever-increasing popularity of MCMC methods when published by Geman and Geman at 1984. It was named after J. W. Gibbs who studied statistical physics.

Reasons behind popularity: Can solve (computationally) Bayesian problems that are analytically intractable; is (and was) fast enough for even early computers.

Idea in a nutshell – Joint distribution can be constructed using conditional distributions.



# ...Gibbs sampling

Suppose we want to know joint distribution  $D(\boldsymbol{\theta})$ , but we do not know it analytically. If we know all the conditional distributions

$$D(\theta_i | \boldsymbol{\theta}_{j \neq i})$$

The chain of simulated values for  $\theta_i$ 's from conditional distributions

$$\theta_1^{(1)} \sim D(\theta_1 | \boldsymbol{\theta}_{j \neq 1})$$

...

$$\theta_k^{(1)} \sim D(\theta_k | \boldsymbol{\theta}_{j \neq k})$$

$$\theta_1^{(2)} \sim D(\theta_1 | \boldsymbol{\theta}_{j \neq 1})$$

...

(when always using latest, updated values for  $\boldsymbol{\theta}_{j \neq i}$ ) will converge to the (unknown) distribution  $D(\boldsymbol{\theta})$



# ...Gibbs sampling

GS is a special case of Metropolis-Hastings, where update is done component-wise and proposed value is always accepted (therefore efficient).

The only drawback in GS is that conditionals must be known. In practice, hybrid methods are used. If conditionals are known for some block of parameters, pure GS is used there. For other blocks M-H could be used inside the GS step to approximate the unknown conditional.



# Gibbs sampling in regression models

GS suits especially well to Bayesian linear or nonlinear regression, if certain choices are made with model assumptions and a priori distributions. In regression the model is

$$\mathbf{y} = f(\mathbf{X}; \boldsymbol{\theta}) + \boldsymbol{\epsilon}$$

where  $\mathbf{y}$  are the observed data,  $f$  is the model,  $\mathbf{X}$  the explanatory variables,  $\boldsymbol{\theta}$  the model parameters, and  $\boldsymbol{\epsilon}$  the errors between model and data. Statistics come to play as we assume  $\boldsymbol{\epsilon}$  to be random having distribution  $D(\mathbf{0}, \boldsymbol{\Sigma})$ , thus  $\mathbf{y}$  has distribution  $D(f(\mathbf{X}; \boldsymbol{\theta}), \boldsymbol{\Sigma})$ . We want to estimate distribution for  $\boldsymbol{\theta}$ , which is written in Bayesian regression (BR) as

$$D(\boldsymbol{\theta}, \boldsymbol{\Sigma} | \mathbf{y}) = D_{pr}(\boldsymbol{\theta}, \boldsymbol{\Sigma}) D(\mathbf{y} | \boldsymbol{\theta}, \boldsymbol{\Sigma})$$



# Steps towards GS in BR

Usually it is assumed that the residual variance  $\Sigma$  and model parameters  $\theta$  are independent, thus

$$D_{pr}(\theta, \Sigma) = D_{pr}(\theta)D_{pr}(\Sigma)$$

and that observations are independent (and identical), thus

$$\Sigma = \sigma^2 \mathbf{I} \quad (\text{or } \sigma^2 \text{diag}(w_1 \dots w_n) \text{ if unequal errors})$$

Collecting these gives

$$D(\theta, \sigma^2 | \mathbf{y}) = D_{pr}(\theta)D_{pr}(\sigma^2)D(\mathbf{y} | \theta, \sigma^2)$$



# ...GS in BR

Now, conditional distributions are

$$D(\boldsymbol{\theta}|\sigma^2, \mathbf{y}) \propto D_{pr}(\boldsymbol{\theta})D(\mathbf{y}|\boldsymbol{\theta}, \sigma^2)$$

$$D(\sigma^2|\boldsymbol{\theta}, \mathbf{y}) \propto D_{pr}(\sigma^2)D(\mathbf{y}|\boldsymbol{\theta}, \sigma^2)$$

GS is efficient choice if the above distributions are known. Even if the choices above are suitable for only some of the parameters, hybrid GS combining 'pure GS' and M-H updates can be efficient.



## ...GS in BR

To my knowledge, there are at least the following cases where conditionals in previous slide are known. The error term  $\varepsilon$  is assumed Gaussian, the model linear, and

$$D_{pr}(\boldsymbol{\theta}) \sim 1 \text{ (improper even dist.)}, \text{ shifted Exp}(\lambda_0), \text{ or } \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

$$D_{pr}(\sigma^2) \sim \text{scaled-inverse-}\chi^2(\alpha, \beta), \text{ or } \mathcal{IG}(\alpha, \beta) \text{ (inverse-gamma)}$$

For example, with even distribution for  $\boldsymbol{\theta}$  and inverse-gamma for  $\sigma^2$  the conditionals for Gibbs sampling step are

$$D(\boldsymbol{\theta}|\sigma^2, \mathbf{y}) \sim \mathcal{N}(X^T X)^{-1} X^T \mathbf{y}, \sigma^2 (X^T X)^{-1}$$

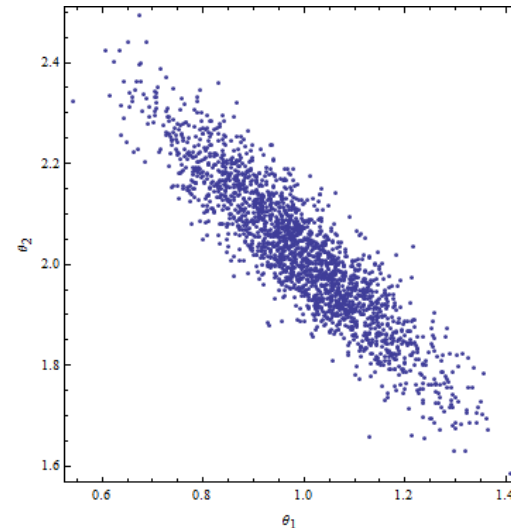
$$D(\sigma^2|\boldsymbol{\theta}, \mathbf{y}) \sim \mathcal{IG}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2}(\mathbf{y} - X\boldsymbol{\theta})^T(\mathbf{y} - X\boldsymbol{\theta})\right)$$





# Computer task

- Use information on previous slide, and write a program that uses Gibbs sampling to estimate model fit and parameter distribution
  - Model is  $y = \theta_1 x_1 + \theta_2 x_2$
  - Data can be found on course webpage
- Note:  $(X^T X)^{-1} X^T \mathbf{y} = \hat{\boldsymbol{\theta}}$ , a.k.a. the model fit
- $(\mathbf{y} - X\boldsymbol{\theta})^T (\mathbf{y} - X\boldsymbol{\theta})$  is the sum of squared residuals
- produce scatterplot of possible parameter values (see fig.) and estimates for parameter values
- (this example is a bit silly, because posterior is also known, and if noninformative priors are used, it does not offer much in addition to regular linear regression...)





# Tips

- Random number from inverse gamma
  - If  $x \sim \text{Gamma}(\alpha, 1/\beta)$  then  $1/x \sim \text{inv-Gamma}(\alpha, \beta)$
  - Suitable hyperparameters for noninformative inv-gamma would be  $\alpha = \beta = 0.001$
- In a nutshell
  - load data
  - Simulate parameters  $(\theta_1, \theta_2)$ , then  $\sigma^2$  over and over again using distributions in page 8. You will need a starting value for the first  $\sigma^2$ .
    - If you can simulate only one-dimensional normal distribution, see next slide
  - Throw off some of the first data, plot rest of  $(\theta_1, \theta_2)$ -pairs and compute estimates as means or medians of the simulated values.



## ...tips

If you can draw random numbers only from 1-dimensional normal distribution, you must alter the Gibbs chain accordingly:

- draw  $\theta_1$  from

$$\mathcal{N}\left(\hat{\theta}_1 + \frac{c_{12}(\theta_2 - \hat{\theta}_2)}{c_{22}}, \frac{(c_{11}c_{22} - c_{12}^2)\sigma^2}{c_{22}}\right)$$

- and  $\theta_2$  from

$$\mathcal{N}\left(\hat{\theta}_2 + \frac{c_{12}(\theta_1 - \hat{\theta}_1)}{c_{11}}, \frac{(c_{11}c_{22} - c_{12}^2)\sigma^2}{c_{11}}\right)$$

- and  $\sigma^2$  from

$$\mathcal{IG}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2}(\mathbf{y} - X\boldsymbol{\theta})^T(\mathbf{y} - X\boldsymbol{\theta}))\right)$$

- where [see next slide], and iterate this



## ...tips

$$(\hat{\theta}_1, \hat{\theta}_2) = (X^T X)^{-1} X^T \mathbf{y}$$

$$(X^T X)^{-1} = \begin{bmatrix} c_{11} & c_{12} \\ c_{12} & c_{22} \end{bmatrix}$$