

Example Class 5

Department of Statistics and Actuarial Science,
The University of Hong Kong

October 25, 2022

Bayesian Model Averaging (BMA)

- Assume that we have K candidate models, M_1, \dots, M_K to fit the data D .
- $P(M_k)$ is the prior probability that M_k is the true model.
- For example, you can assign an equal weight to all candidate models, i.e., $P(M_k) = 1/K$ for $k = 1, \dots, K$.
- Each model M_k has its own parameter θ_k , and let $\pi(\theta_k|M_k)$ be the prior of θ_k .
- The posterior probability of M_k has the form,

$$P(M_k|D) = \frac{P(D|M_k)P(M_k)}{\sum_{j=1}^K P(D|M_j)P(M_j)}$$

$$P(D|M_k) = \int \pi(\theta_k|M_k)f(D|\theta_k, M_k)d\theta_k$$

BMA Estimator

- For the model parameter θ , the BMA estimator is given by

$$\bar{\theta} = \sum_{k=1}^K \hat{\theta}_k P(M_k|D)$$

- Here each $\hat{\theta}_k$ is the posterior mean of θ_k ,

$$\begin{aligned}\hat{\theta}_k &= \int \theta_k f(\theta_k|D, M_k) d\theta_k \\ &= \int \theta_k \frac{\pi(\theta_k|M_k) f(D|\theta_k, M_k)}{\int \pi(\theta_k|M_k) f(D|\theta_k, M_k) d\theta_k} d\theta_k\end{aligned}$$

- By assigning the posterior mean $\hat{\theta}_k$ a weight of $P(M_k|D)$, BMA automatically lean toward the best fitting model, and thus $\bar{\theta}$ will be close to the best parameter estimate
- If T is the quantity of interest,

$$f(T|D) = \sum_{k=1}^K f(T|D, M_k) P(M_k|D)$$

Why BMA?

- Frequentist
 - Model selection
 - Regularization

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$$\min_{\boldsymbol{\beta}} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}$$

- Bayesian
 - BMA
 - Marginalization

$$f(T|D) = \sum_{k=1}^K f(T|D, M_k) P(M_k|D)$$

$$P(D|M_k) = \int \pi(\boldsymbol{\theta}_k|M_k) f(D|\boldsymbol{\theta}_k, M_k) d\boldsymbol{\theta}_k$$

Why BMA?

- Averaging over all of the models provides better predictive ability, as measured by a logarithmic scoring rule, than using any single model M_j

$$-E \left[\log \left\{ \sum_{k=1}^K P(T|D, M_k) P(M_k|D) \right\} \right] \leq -E[\log P(T|D, M_j)]$$

- In decision theory, a score function, or scoring rule, measures the accuracy of probabilistic predictions.

- Simple linear regression

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \sigma^2 \mathbf{I}$$

- With $q + 1$ predictors (including the intercept term), overall you have $K = 2^{(q+1)}$ possible models $\{M_1, \dots, M_K\}$
- Simplified case including one predictor without intercept

$$y_i = \beta_\gamma x_{i\gamma} + \epsilon_i$$

BMA Example

- Linear regression with only one predictor

$$y_i = \beta_\gamma x_{i\gamma} + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

- Overall q predictors $\mathbf{X}_1, \dots, \mathbf{X}_q$
- The j -th Model M_j includes the j -th predictor \mathbf{X}_j
- Prior

$$\mathbf{Y} | M_j, \beta_j, \mathbf{X}_j, \sigma^2 \sim N(\mathbf{X}_j \beta_j, \sigma^2 \mathbf{I})$$

$$\beta_j | \sigma^2, \mu_0, \lambda_0 \sim N(\mu_0, \sigma^2 / \lambda_0)$$

$$\sigma^2 | a_0, b_0 \sim \text{InverseGamma}(a_0, b_0)$$

$$P(M_j) = 1/q$$

Normal-Inverse-Gamma Prior

- The normal-inverse-gamma distribution (or Gaussian-inverse-gamma distribution) is a four-parameter family of bivariate continuous probability distributions.
- It is the conjugate prior of a normal distribution with unknown mean and variance.
- (x, σ^2) has a normal-inverse-gamma distribution with parameter $(\mu, \lambda, \alpha, \beta)$ if

$$\begin{aligned}x|\sigma^2, \mu, \lambda &\sim N(\mu, \sigma^2/\lambda) \\ \sigma^2|\alpha, \beta &\sim IG(\alpha, \beta)\end{aligned}$$

$$f(x, \sigma^2|\mu, \lambda, \alpha, \beta) = \frac{\sqrt{\lambda}}{\sigma\sqrt{2\pi}} \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} e^{-\frac{2\beta+\lambda(x-\mu)^2}{2\sigma^2}}$$

BMA Example

$$f(D|\beta_k, \sigma^2, M_k) \propto (\sigma^2)^{-\frac{n}{2}} e^{-\frac{\sum_{i=1}^n (y_i - \beta_k x_{ik})^2}{2\sigma^2}}$$

$$\begin{aligned} f(\beta_k, \sigma^2|D, M_k) &\propto (\sigma^2)^{-\frac{n}{2}} e^{-\frac{\sum_{i=1}^n (y_i - \beta_k x_{ik})^2}{2\sigma^2}} (\sigma^2)^{-\frac{1}{2}} e^{-\frac{\lambda_0(\beta_k - \mu_0)^2}{2\sigma^2}} (\sigma^2)^{-a_0-1} e^{-\frac{b_0}{\sigma^2}} \\ &\sim \text{NIG}\left(\frac{\lambda_0\mu_0 + \sum_{i=1}^n y_i x_{ik}}{\lambda_0 + \sum_{i=1}^n x_{ik}^2}, \lambda_0 + \sum_{i=1}^n x_{ik}^2, a_0 + \frac{n}{2}, b_{\text{new}}\right) \\ b_{\text{new}} &= b_0 + \frac{1}{2} \sum_{i=1}^n y_i^2 + \frac{\lambda_0\mu_0^2 \sum_{i=1}^n x_{ik}^2 - (\sum_{i=1}^n y_i x_{ik})^2 - 2\lambda_0\mu_0 \sum_{i=1}^n y_i x_{ik}}{2(\lambda_0 + \sum_{i=1}^n x_{ik}^2)} \end{aligned}$$

BMA Example

$$\begin{aligned}P(D|M_k) &= \int \pi(\beta_k, \sigma^2 | M_k) f(D | \beta_k, \sigma^2, M_k) d(\beta_k, \sigma^2) \\&= \int \frac{\sqrt{\lambda_0}}{\sigma \sqrt{2\pi}} \frac{b_0^{a_0}}{\Gamma(a_0)} (\sigma^2)^{-a_0-1} e^{-\frac{2b_0 + \lambda_0(\beta_k - \mu_0)^2}{2\sigma^2}} (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{\sum_{i=1}^n (y_i - x_{ik}\beta_k)^2}{2\sigma^2}} d(\beta_k, \sigma^2) \\&= \sqrt{\lambda_0} (2\pi)^{-\frac{n+1}{2}} \frac{b_0^{a_0}}{\Gamma(a_0)} \int (\sigma^2)^{-\frac{n+1}{2} - a_0 - 1} e^{-\frac{\sum_{i=1}^n (y_i - \beta_k x_{ik})^2}{2\sigma^2}} e^{-\frac{2b_0 + \lambda_0(\beta_k - \mu_0)^2}{2\sigma^2}} d(\beta_k, \sigma^2) \\&= \frac{\sqrt{\lambda_0}}{\sqrt{\lambda_0 + \sum_{i=1}^n x_{ik}^2}} (2\pi)^{-\frac{n}{2}} \frac{b_0^{a_0}}{\Gamma(a_0)} \frac{\Gamma(a_0 + n)}{b_{\text{new}}^{a_0 + \frac{n}{2}}}\end{aligned}$$

Or you can calculate $P(D|M_k)$ via the Monte Carlo method for integration.

- Five covariates:

- $(X_1, X_2) \sim N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix} \right)$
- $X_3 \sim N(0, 1)$
- $X_4 \sim \text{Bernoulli}(0.5) - 0.5$
- $X_5 = 2 \cdot X_1$

- True model: $y = \beta_{\text{true}} X_1 + \epsilon$, $\beta_{\text{true}} = 1$, $\epsilon \sim N(0, 1)$
- Prior: $P(M_k) = 1/5$, $\mu_0 = \beta_{\text{true}} = 1$, $\lambda_0 = 0.1$, $a_0 = b_0 = 0.1$
- Posterior Model Probability $P(M_k|D)$ for different sample sizes n

BMA Example

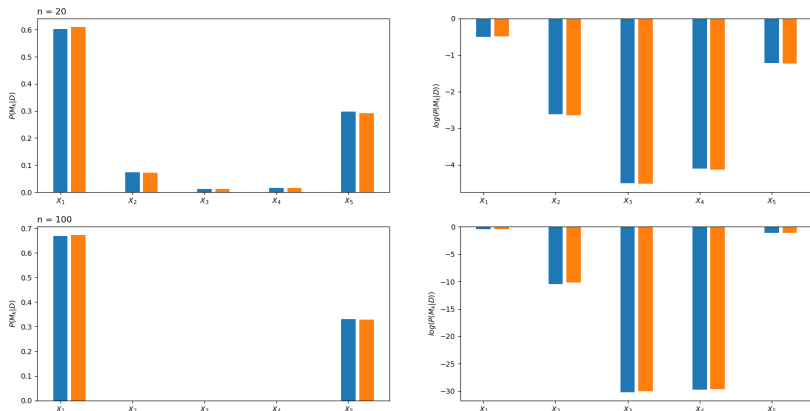


Figure 1: Barplots of $P(M_k|D)$ for $n = 20$ (left) and 100 (right).

Multivariate Linear Regression

- Multivariate Linear Regression

$$y_{i,1} = \mathbf{x}_i^T \boldsymbol{\beta}_1 + \epsilon_{i,1},$$

...

$$y_{i,m} = \mathbf{x}_i^T \boldsymbol{\beta}_m + \epsilon_{i,m}.$$

- The sets of errors $\boldsymbol{\epsilon}_i = \{\epsilon_{i,1}, \dots, \epsilon_{i,m}\}$ are correlated.
- $\mathbf{y}_i^T = \mathbf{x}_i^T \mathbf{B} + \boldsymbol{\epsilon}_i$; $\mathbf{y}_i = \{y_{i,1}, \dots, y_{i,m}\}$; $\mathbf{B} = [\beta_{i,j}]_{p \times m}$
- We assume that $\boldsymbol{\epsilon}_i$ is jointly normal, i.e., $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{\epsilon})$
- Regression problem in matrix form:

$$\begin{matrix} \mathbf{Y} \\ (n \times m) \end{matrix} = \begin{matrix} \mathbf{X} & \mathbf{B} \\ (n \times p) & (p \times m) \end{matrix} + \begin{matrix} \mathbf{E}, \\ (n \times m) \end{matrix}$$

Bayesian Multivariate Linear Regression

- Model: $\mathbf{Y} = \mathbf{XB} + \mathbf{E}$, $\epsilon_i \sim N(\mathbf{0}, \Sigma_\epsilon)$
- Observations: \mathbf{Y}, \mathbf{X}
- Parameters of interest: $\mathbf{B}, \Sigma_\epsilon$
- Likelihood:

$$\begin{aligned} p(\mathbf{Y}|\mathbf{X}, \mathbf{B}, \Sigma_\epsilon) &\propto |\Sigma_\epsilon|^{-\frac{n}{2}} e^{-\frac{1}{2}\{\sum_{i=1}^n (\mathbf{y}_i - \mathbf{x}_i^T \mathbf{B})^T \Sigma_\epsilon^{-1} (\mathbf{y}_i - \mathbf{x}_i^T \mathbf{B})\}} \\ &\propto |\Sigma_\epsilon|^{-\frac{n}{2}} e^{-\frac{1}{2}\text{tr}((\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB}) \Sigma_\epsilon^{-1})} \end{aligned}$$

- Priors:

$$\begin{aligned} \pi(\mathbf{B}, \Sigma_\epsilon) &= \pi(\mathbf{B}|\Sigma_\epsilon)\pi(\Sigma_\epsilon), \\ \pi(\mathbf{B}|\Sigma_\epsilon) &\sim N(\boldsymbol{\mu}_0, \Sigma_\epsilon/\lambda_0) \\ \pi(\Sigma_\epsilon) &\sim \mathcal{W}^{-1}(\boldsymbol{\Psi}_0, \nu_0) \end{aligned}$$

Inverse Wishart Distribution

- The inverse Wishart distribution (\mathcal{W}^{-1}) is used as the conjugate prior for the covariance matrix of a multivariate normal distribution.
- For $\Sigma_\epsilon \sim \mathcal{W}^{-1}(\Psi, \nu)$, its pdf is,

$$f_p(\Sigma_\epsilon | \Psi, \nu) = \frac{|\Psi|^{\nu/2}}{2^{\nu p/2} \Gamma_p(\nu/2)} |\Sigma_\epsilon|^{-(\nu+p+1)/2} e^{-\frac{1}{2} \text{tr}(\Psi \Sigma_\epsilon^{-1})}$$

where Σ_ϵ and Ψ are $p \times p$ positive definite matrices, $|\cdot|$ is the determinant, and Γ_p is the multivariate gamma function.

- If $\Sigma_\epsilon \sim \mathcal{W}^{-1}(\Psi, \nu)$, its inverse Σ_ϵ^{-1} has a Wishart distribution $\mathcal{W}(\Psi^{-1}, \nu)$.

Normal-Inverse-Wishart (NIW) Distribution

- The normal-inverse-Wishart distribution is a multivariate four-parameter family of continuous probability distributions.
- We say $(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \sim NIW(\boldsymbol{\mu}_0, \lambda, \boldsymbol{\Psi}, \nu)$ if

$$f(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{\mu}_0, \lambda, \boldsymbol{\Psi}, \nu) = N(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \boldsymbol{\Sigma} / \lambda) \mathcal{W}^{-1}(\boldsymbol{\Sigma} | \boldsymbol{\Psi}, \nu)$$

$$\boldsymbol{\mu} | \boldsymbol{\mu}_0, \lambda, \boldsymbol{\Psi} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma} / \lambda)$$

$$\boldsymbol{\Sigma} | \boldsymbol{\Psi}, \nu \sim \mathcal{W}^{-1}(\boldsymbol{\Psi}, \nu)$$

- It is the conjugate prior of a multivariate normal distribution with unknown mean and covariance matrix.

Bayesian Multivariate Linear Regression

$$\pi(\mathbf{B}, \Sigma_\epsilon | \mathbf{Y}, \mathbf{X}) \quad (1)$$

$$\propto |\Sigma_\epsilon|^{-\frac{n}{2}} e^{-\frac{1}{2} \text{tr}((\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB}) \Sigma_\epsilon^{-1})} \times \quad \text{likelihood} \quad (2)$$

$$|\Sigma_\epsilon|^{-\frac{p}{2}} e^{-\frac{\lambda_0}{2} \text{tr}((\mathbf{B} - \mathbf{B}_0)^T (\mathbf{B} - \mathbf{B}_0) \Sigma_\epsilon^{-1})} \times \quad \text{prior for } \mathbf{B}, \mathbf{B}_0 \text{ is } \mu_0 \quad (3)$$

$$|\Sigma_\epsilon|^{-\frac{\nu_0 + m + 1}{2}} e^{-\frac{1}{2} \text{tr}(\Psi_0 \Sigma_\epsilon^{-1})} \quad \text{prior for } \Sigma_\epsilon \quad (4)$$

$$\propto |\Sigma_\epsilon|^{-\frac{p+m}{2}} e^{-\frac{1}{2} \text{tr}((\mathbf{B} - \mathbf{B}_n)^T (\lambda_0 \mathbf{I} + \mathbf{X}^T \mathbf{X}) (\mathbf{B} - \mathbf{B}_n) \Sigma_\epsilon^{-1})} \times \quad \text{MatrixNormal} \quad (5)$$

$$|\Sigma_\epsilon|^{-\frac{\nu_0 + n + 1}{2}} e^{-\frac{1}{2} \text{tr}((\Psi_0 + (\mathbf{Y} - \mathbf{XB}_n)^T (\mathbf{Y} - \mathbf{XB}_n) + \lambda_0 (\mathbf{B}_0 - \mathbf{B}_n)^T (\mathbf{B}_0 - \mathbf{B}_n)) \Sigma_\epsilon^{-1})} \quad \text{inverse-Wishart} \quad (6)$$

$$\text{where } \mathbf{B}_n = (\lambda_0 \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} (\lambda_0 \mathbf{B}_0 + \mathbf{X}^T \mathbf{Y})$$

$$\text{Psi_new} = \Psi_0 + (\mathbf{Y} - \mathbf{XB}_n)^T (\mathbf{Y} - \mathbf{XB}_n) + \lambda_0 (\mathbf{B}_0 - \mathbf{B}_n)^T (\mathbf{B}_0 - \mathbf{B}_n)$$

$$\text{MatrixNormal}(\mathbf{X} | \mathbf{X}_0, \mathbf{U}, \mathbf{V}) \propto |\mathbf{U}|^{-d_V/2} |\mathbf{V}|^{-d_U/2} e^{-\frac{1}{2} \text{tr}(\mathbf{V}^{-1} (\mathbf{X} - \mathbf{X}_0)^T \mathbf{U}^{-1} (\mathbf{X} - \mathbf{X}_0))}$$

$$V \text{ shape} : (d_V, d_V) \quad U \text{ shape} : (d_U, d_U)$$

Compare formula (5),(6) with the definition of MatrixNormal and \mathcal{W}^{-1} we get

$$\pi(\mathbf{B}, \Sigma_\epsilon | \mathbf{Y}, \mathbf{X}) \propto \mathcal{W}^{-1}(\nu_0 + n, \text{Psi_new}) \cdot \text{MatrixNormal}(\mathbf{B}_n, (\lambda_0 \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1}, \Sigma_\epsilon)$$

BMLR (Improper Prior)

- So

$$\Sigma_{\epsilon} | \mathbf{Y}, \mathbf{X} \sim \mathcal{W}^{-1}(\nu_0 + n, \text{Psi_new})$$

$$\mathbf{B} | \mathbf{Y}, \mathbf{X} \sim \text{MatrixNormal}(\mathbf{B}_n, (\lambda_0 \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1}, \Sigma_{\epsilon})$$

In the code, we will draw samples for Σ_{ϵ} and \mathbf{B} from this two posterior distribution, and get the estimated value of Σ_{ϵ} and \mathbf{B} .

- For simplicity, we can consider an improper prior for $(\mathbf{B}, \Sigma_{\epsilon})$

$$\pi(\mathbf{B}, \Sigma_{\epsilon}) \propto \Sigma_{\epsilon}^{-\frac{m+p+1}{2}}$$

$$\pi(\mathbf{B}, \Sigma_{\epsilon} | \mathbf{Y}, \mathbf{X}) \propto |\Sigma_{\epsilon}|^{-\frac{n}{2}} e^{-\frac{1}{2} \text{tr}((\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB}) \Sigma_{\epsilon}^{-1})} |\Sigma_{\epsilon}|^{-\frac{m+p+1}{2}}$$

BMLR (Example)

- Model:

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E}, \quad \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 0.5 \end{bmatrix}$$
$$\mathbf{E} \sim N(\mathbf{0}, \Sigma_{\epsilon}), \quad \Sigma_{\epsilon} = \begin{bmatrix} 1 & 0.3 \\ 0.3 & 1 \end{bmatrix}$$

- Prior:

$$\mu_0 = \mathbf{0}, \quad \lambda_0 = 0.01,$$
$$\Psi_0 = \mathbf{0}, \quad \nu_0 = 0.01.$$

BMLR Example

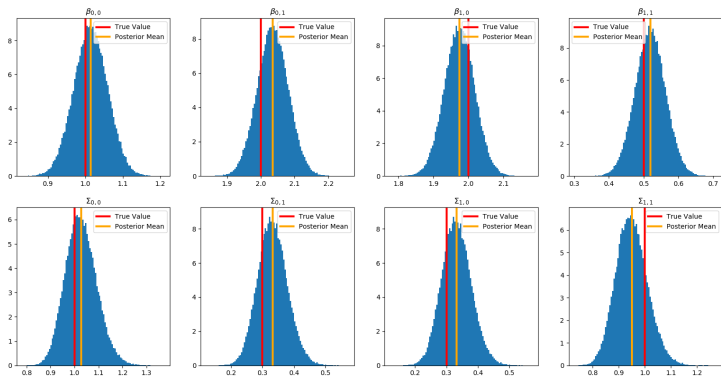


Figure 2: Histograms of posterior samples of \mathbf{B} (upper) and Σ_ϵ (lower).