

Handout 15: Inference under model uncertainty^a

Lecturer: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim: To explain, design, and apply model selection and model determination Bayesian procedures.

References:

- Robert, C. (2007; Section 7.1 & 7.2). The Bayesian choice: from decision-theoretic foundations to computational implementation. Springer Science & Business Media.
- O'Hagan, A., & Forster, J. J. (2004; Section 7 Model comparison). Kendall's advanced theory of statistics, volume 2B: Bayesian inference (Vol. 2). Arnold.

^aAuthor: Georgios P. Karagiannis.

1 Model uncertainty problem

Note 1. All (statistical/mathematical) models are approximations/simplifications of the reality. They are arguable. Their purpose is to allow You draw conclusions about reality and take decisions.

- “All models are wrong, but some are useful”. –George Box’s quote

Note 2. Let $y = (y_1, \dots, y_n)$ be observables generated from the real but unknown data-generating process $y \sim R(y)$. Let $z = (y_{n+1}, \dots, y_{n+m})$ be future outcomes from the same process $R(\cdot)$.

Note 3. Statistician may be unsure of the proper way to formulate the statistical model $\mathcal{M} = \{F(y|\theta); \theta \in \Theta\}$ in order to acceptably represent the data-generating process $R(y)$. We call this situation ‘model uncertainty’.

Note 4. To address problems under model uncertainty, often

- statistician initially considers a restricted set $\{\mathcal{M}_k; k \in \mathcal{K}\}$ of available - reasonable - potential - competing - alternative statistical models $\mathcal{M}_k = \{F_k(y|\theta_k); \theta_k \in \Theta_k\}$. Then,
- interest lies in learning a model \mathcal{M} which ‘best’ represents $R(y)$, performing inference (predictive or parametric), and drawing conclusions under the pretense of model uncertainty.

Example 5. Consider the Normal linear regression problem. Let $\{\phi_0, \phi_1, \dots, \phi_{d-1}\}$ be a set of d possible regressor variables, which may affect the values of the response variable $y \in \mathcal{Y} \subseteq \mathbb{R}$. We are interested in learning the mapping

$$\{\phi_0, \dots, \phi_{d-1}\} \mapsto y$$

Assume there are available n pairs of observations: $(\Phi_1, y_1), (\Phi_2, y_2), \dots, (\Phi_n, y_n)$.

Consider the following approximations of the real mapping $\{\phi_0, \dots, \phi_{d-1}\} \mapsto y$: Assume that observables y_i are Normally distributed and independent for $i = 1, \dots, n$, with unknown variance σ^2 , and unknown mean which can be parametrised as a linear combination of some of the regressors $\{\phi_0, \dots, \phi_{d-1}\}$ with unknown coefficients β . Finally, consider, we are uncertain which set of $\{\phi_0, \dots, \phi_{d-1}\}$ actually affects values of y .

2 Standard framework

Note 6. Consider a collection of statistical models $\mathcal{M} = \{\mathcal{M}_k; k \in \mathcal{K}\}$, where each model

$$\mathcal{M}_k = \{F_k(y|\theta_k); \theta_k \in \Theta_k\},$$

is labeled by an index $k \in \mathcal{K}$, and associated with the sampling distribution $F_k(y|\theta_k)$ with pdf/pmf $f_k(y|\theta_k)$, unknown parameter $\theta_k \in \Theta_k$, and parametric space Θ_k .

Definition 7. The encompassing model \mathcal{M} is defined as the collection of statistical models

$$\begin{aligned} \mathcal{M} &= \{\mathcal{M}_k; k \in \mathcal{K}\} \\ &= \{\{F_k(y|\theta_k); \theta_k \in \Theta_k\}; k \in \mathcal{K}\} \end{aligned}$$

labeled by parameter $\vartheta = (k, \theta_k) \in \Theta$ which is defined on the joint parameter space

$$\Theta = \cup_{k \in \mathcal{K}} \{k\} \times \Theta_k$$

Example 8. (Cont...) We set-up the encompassing model $\mathcal{M} = \{\mathcal{M}_k; k \in \mathcal{K}\}$ as the collection of statistical models

$$\mathcal{M}_k = \{y|\beta_k, \sigma^2, k \sim N(\Phi_k \beta_k, I\sigma^2); \beta_k \in \mathbb{R}^{d_k}, \sigma^2 \in \mathbb{R}_+\}, \quad (1)$$

where $k \subseteq \{0, 1, \dots, d-1\}$ is the label of the model indicating the sub-set of the regressors included in the statistical model \mathcal{M}_k , and \mathcal{K} is the collection of all sub-sets of $\{0, 1, \dots, d-1\}$ so that $k \in \mathcal{K}$. Given \mathcal{M}_k , Φ_k is the design matrix consisting of the regressors with indexes in k , $\beta_k \in \mathbb{R}^{d_k}$ is the associated vector of regression coefficients, and d_k is the number of regressors. Notice that σ^2 is kept to be common among all the models, however one could have set it to be different between different models (e.g.; σ_k^2). The within model parameters are $\theta_k = (\beta_k, \sigma^2)$ defined on space $\Theta_k = \mathbb{R}^{d_k} \times \mathbb{R}$. The parameter of \mathcal{M} is $\vartheta = (k, \beta_k, \sigma^2)$ and defined on space $\Theta = \cup_{k \in \mathcal{K}} \{k\} \times \mathbb{R}^{d_k} \times \mathbb{R}$.

Note 9. To complete the Bayesian model we specify prior distributions on the joint parametric space Θ . Consider a collection of conditional priors on θ_k given model \mathcal{M}_k as

$$\theta_k | \mathcal{M}_k \sim \Pi(\theta_k | \mathcal{M}_k)$$

with pdf/pmf $\pi(\theta_k | \mathcal{M}_k)$; and consider a marginal prior distribution on \mathcal{M}_k

$$\mathcal{M}_k \sim \Pi(\mathcal{M}_k)$$

with pdf/pmf $\pi(\mathcal{M}_k)$, for all $k \in \mathcal{K}$. The encompassing prior model has distribution $\Pi(\mathcal{M}_k, \theta_k)$ with pdf/pmf

$$\pi(\mathcal{M}_k, \theta_k) = \pi(\theta_k | \mathcal{M}_k) \pi(\mathcal{M}_k), \quad \forall (k, \theta_k) \in \Theta$$

Note 10. Specifying the prior distribution $\Pi(\mathcal{M}_k, \theta_k)$ is a delicate matter.

1. Within model priors $\theta_k | \mathcal{M}_k \sim \Pi(\theta_k | \mathcal{M}_k)$ should not be improper because the normalizing constants are not canceled in the joint posterior (7), and hence the Lindley-Jeffreys paradox may kick in.
 2. If some models are embedded (nested) into others, e.g. $\mathcal{M}_1 \subseteq \mathcal{M}_2$, possibly the marginal model prior should be $\pi(\mathcal{M}_1) \leq \pi(\mathcal{M}_2)$, in order to be coherent –this is not a panacea.
 3. Common parameters between models (such as σ^2 in Example) may be treated as separate parameters. However, it is common to be regarded as the same parameter and be assigned the same prior; this is an additional approximation/simplification that modeler does for computational convenience –less parameters to learn ;-).
- In the Example, σ^2 is common to all models. If we wanted to be in accordance with this note, we should have specified σ^2 separately for each model \mathcal{M}_k , as σ_k^2 , but we do not do this now (naughty).

59 *Note 11.* The full Bayesian model can be summarized as

$$60 \quad \begin{cases} y|\theta_k, \mathcal{M}_k \sim F_k(y|\theta) & , \text{ data generation from the sampling distribution} \\ \theta_k|\mathcal{M}_k \sim \Pi(\theta_k|\mathcal{M}_k) & , \text{ within model parameter generation from the conditional prior} \\ \mathcal{M}_k \sim \Pi(\mathcal{M}_k) & , \text{ model generation from the marginal model prior} \end{cases} \quad (2)$$

61 hence the joint distribution is such that $dP(y, \theta_k, \mathcal{M}_k) = dF_k(y|\theta)d\Pi(\theta_k|\mathcal{M}_k)d\Pi(\mathcal{M}_k)$.

62 **Example 12.** (Cont...) One can specify the following Bayesian linear regression model

$$63 \quad y|\beta_k, \sigma^2, \mathcal{M}_k \sim N(\Phi_k \beta_k, I\sigma^2) \quad , \text{ the sampling distribution} \quad (3)$$

$$64 \quad \beta_k|\sigma^2, \mathcal{M}_k \sim N(\mu_k, \sigma^2 V_k) \quad , \text{ conditional prior of the within model parameter} \quad (4)$$

$$65 \quad \sigma^2|\mathcal{M}_k \sim \text{IG}(a, \lambda) \quad \text{conditional prior of the within model parameter} \quad (5)$$

$$66 \quad \mathcal{M}_k \sim \pi(k) = \frac{1}{|\mathcal{K}|} = \frac{1}{2^d} \quad , \text{ marginal model prior} \quad (6)$$

69 Eq. 3 is the sampling distribution $F(y|\beta_k, \sigma^2, \mathcal{M}_k)$ for model \mathcal{M}_k . Eq. 4 and 5 are the conditional (or within) model
70 priors $\Pi(\beta_k, \sigma^2|\mathcal{M}_k)$; recall that they are conjugate to $F(y|\beta_k, \sigma^2, \mathcal{M}_k)$ and hence chosen for our computational
71 convenience. Notice that parameter σ^2 , which is common to all available models $\{\mathcal{M}_k\}$, is treated as obtaining the
72 same values among all models, and following the same prior. This is against the suggestion in Note 10(3)) but we
73 assume that the computational benefits of reducing the dimensionality of the parametric space dominate the losses
74 from the (possibly) worse approximation of the reality. In Eq. 6, the marginal model prior $\pi(k)$ is chosen to be
75 uniform across models (we have $|\mathcal{K}| = 2^d$ possible combinations or regressors).

76 *Notation 13.* To make the notation easier, we will denote \mathcal{M}_k as k and the probability measures as

- 77 • $F(y|\theta_k, k) := F_k(y|\theta_k)$ and $f(y|\theta_k, k) := f_k(y|\theta_k)$, as well as
- 78 • $\Pi(\theta_k|k) := \Pi(\theta_k|\mathcal{M}_k)$, $\pi(\theta_k|k) := \pi(\theta_k|\mathcal{M}_k)$, $\Pi(k) := \Pi(\mathcal{M}_k)$, and $\Pi(k, \theta_k) := \Pi(\mathcal{M}_k, \theta_k)$, etc...

79 3 Posterior distributions

80 According to the Bayesian, the following posteriors can be derived.

81 **Proposition 14.** *Given the Bayesian model (2), the joint posterior distribution $\Pi(k, \theta_k|y)$ has pdf/pmf*

$$82 \quad \pi(k, \theta_k|y) = \frac{f(y|\theta_k, k)\pi(\theta_k|k)\pi(k)}{\sum_{k \in \mathcal{K}} \int_{\Theta_k} f(y|\theta_k, k)\pi(\theta_k|k)d\theta_k \pi(k)}, \quad \forall (k, \theta_k) \in \Theta; \quad (7)$$

83 **Proposition 15.** *The joint posterior density (7) can be factorized as*

$$84 \quad \pi(k, \theta_k|y) = \pi(\theta_k|y, k)\pi(k|y)$$

85 where the first part is the pdf/pmf of the conditional posterior $\Pi(\theta_k|y, k)$ of θ_k given model \mathcal{M}_k i.e.

$$86 \quad \pi(\theta_k|y, k) = \frac{f(y|\theta_k, k)\pi(\theta_k|k)}{\int_{\Theta_k} f(y|\theta_k, k)\pi(\theta_k|k)d\theta_k} = \frac{f(y|\theta_k, k)\pi(\theta_k|k)}{f(y|k)}, \quad \forall \theta_k \in \Theta_k; \quad (8)$$

and the second term is the pdf/pmf marginal model posterior $\Pi(k|y)$ of model \mathcal{M}_k with pdf/pdf

$$\pi(k|y) = \frac{f(y|k)\pi(k)}{f(y)}, \forall k \in \mathcal{K}.$$

Here, $f(y|k)$ is the conditional marginal likelihood (or the model evidence) of model \mathcal{M}_k with pdf/pdf

$$f(y|k) = \int_{\Theta_k} f(y|\theta_k, k)\pi(\theta_k|k)d\theta_k$$

and $f(y)$ is the marginal likelihood with

$$f(y) = \sum_{k \in \mathcal{K}} \left[\int_{\Theta_k} f(y|\theta_k, k)\pi(\theta_k|k)d\theta_k \right] \pi(k) = \sum_{k \in \mathcal{K}} f(y|k)\pi(k)$$

Example 16. To calculate the conditional posterior pdf $\pi(\beta_k, \sigma^2|y)$, we use the Bayes Theorem

$$\begin{aligned} \pi(\beta_k, \sigma^2|y, k) &\propto f(y|\theta_k, k)\pi(\theta_k|k) = \mathcal{N}(y|\Phi\beta_k, I\sigma^2) \mathcal{N}(\beta_k|\mu_k, V_k\sigma^2) \text{IG}(\sigma^2|a, \lambda) \\ &\propto \underbrace{\left(\frac{1}{\sigma^2}\right)^{\frac{d}{2}} \exp\left(-\frac{1}{2\sigma^2}(\beta_k - \mu_k^*)^\top (V_k^*)^{-1}(\mu_k - \beta_k^*)\right)}_{\propto \mathcal{N}(\beta_k|\mu_k^*, V_k^*\sigma^2)} \underbrace{\left(\frac{1}{\sigma^2}\right)^{\frac{n}{2}+a+1} \exp\left(-\frac{\lambda}{\sigma^2}\right)}_{\propto \text{IG}(\sigma^2|a^*, \lambda_k^*)} \end{aligned}$$

where

$$\begin{aligned} V_k^* &= (V_k^{-1} + \Phi_k^\top \Phi_k)^{-1}; \quad \mu_k^* = V_k^* (V_k^{-1} \mu_k + \Phi_k^\top y); \quad a^* = \frac{n}{2} + a; \quad \lambda_k^* = \lambda + \frac{1}{2} S_k^*; \\ S_k^* &= \mu_k^\top V_k^{-1} \mu_k - (\mu_k^*)^\top (V_k^*)^{-1} (\mu_k^*) + y^\top y \end{aligned}$$

Hence

$$\Rightarrow \begin{cases} \beta_k|y\sigma^2, k & \sim \mathcal{N}(\mu_k^*, V_k^*\sigma^2), \forall k \in \mathcal{K} \\ \sigma^2|y, k & \sim \text{IG}(a^*, \lambda_k^*) \end{cases}$$

Example 17. To calculate the model evidence / conditional marginal likelihood for \mathcal{M}_k , we integrate

$$f(y|k) = \int \overbrace{\mathcal{N}(y|\Phi\beta_k, I\sigma^2) \mathcal{N}(\beta_k|\mu, V_k\sigma^2)}^{f(y|\theta_k, k)} \overbrace{\text{IG}(\sigma^2|a, \lambda)}^{\pi(\theta_k|k)} d\beta_k d\sigma^2 = \dots \text{calc} \dots = \left(\frac{1}{2\pi}\right)^{\frac{n}{2}} \left(\frac{|V_k^*|}{|V_k|}\right)^{\frac{1}{2}} \frac{\lambda^a}{\Gamma(a)} \frac{\Gamma(a^*)}{(\lambda_k^*)^{a^*}}$$

Since \mathcal{K} is finite, the marginal likelihood $f(y)$ of y can be computed as

$$f(y) = \sum_{k \in \mathcal{K}} f(y|k)\pi(k) = \left(\frac{1}{2\pi}\right)^{\frac{n}{2}} \lambda^a \frac{\Gamma(a^*)}{\Gamma(a)} \frac{1}{2^d} \sum_{k \in \mathcal{K}} \left(\frac{|V_k^*|}{|V_k|}\right)^{\frac{1}{2}} \frac{1}{(\lambda_k^*)^{a^*}}$$

So the marginal model posterior probability $\pi(k|y)$ of \mathcal{M}_k is

$$\pi(k|y) = \frac{f(y|k)\pi(k)}{f(y)} = \left(\frac{|V_k^*|}{|V_k|}\right)^{\frac{1}{2}} \frac{1}{(\lambda_k^*)^{a^*}} \bigg/ \sum_{k' \in \mathcal{K}} \left(\frac{|V_{k'}^*|}{|V_{k'}|}\right)^{\frac{1}{2}} \frac{1}{(\lambda_{k'}^*)^{a^*}} \quad (9)$$

Proposition 18. Let $z = (y_{n+1}, \dots, y_{n+m})$ be a sequence of future outcomes. The conditional predictive distribution $G(z|y, k)$ of z given observables and model \mathcal{M}_k has pdf/pmf

$$g(z|y, k) = \int f(z|y, \theta_k, k) d\Pi(\theta_k|y, k) \stackrel{z \text{ indep. } y}{=} \int f(z|\theta_k, k) d\Pi(\theta_k|y, k) \quad (10)$$

Appendix
A

Appendix
B

The marginal predictive distribution $G(z|y)$ of z given the observables y has pdf/pmf

$$g(z|y) = \sum_{k \in \mathcal{K}} \int f(z|y, \theta_k, k) d\Pi(\theta_k|y, k) \pi(k|y) = \sum_{k \in \mathcal{K}} g(z|y, k) \pi(k|y)$$

Example 19. The predictive distribution $G(z|y, k)$ of future observables $z = (y_{n+1}, \dots, y_{n+m})$ at any new Φ_k^{new} given the observables y and the model \mathcal{M}_k ; i.e. $z = \Phi_k^{\text{new}} \beta_k + \epsilon$ and $\epsilon \sim N(0, I_m \sigma^2)$, has pdf

$$\begin{aligned} g(z|y, k) &= \int f(z|\theta_k, k) \pi(\theta_k|y, k) d\theta_k = \int N(z|\Phi_k^{\text{new}} \beta_k, I \sigma^2) N(\beta_k|\mu_k^*, V_k^* \sigma^2) \text{IG}(\sigma^2|a^*, \lambda_k^*) d\beta_k d\sigma^2 \\ &= \dots \text{calc} \dots = T_m \left(z|\Phi_k^{\text{new}} \mu_k^*, [I + V_k^*] \frac{\lambda_k^*}{a^*}, 2a^* \right) \end{aligned}$$

Example 20. The marginal predictive distribution $G(z|k)$ of future observables $z = (y_{n+1}, \dots, y_{n+m})$ at any new Φ_k^{new} given the observables y has pdf

$$\begin{aligned} g(z|y) &= \sum_{k \in \mathcal{K}} g(z|y, k) \pi(k|y) \\ &= \sum_{k \in \mathcal{K}} T_m \left(z|\Phi_k^{\text{new}} \mu_k^*, V_k^* \frac{\lambda_k^*}{a^*}, 2a^* \right) \left(\frac{|V_k^*|}{|V_k|} \right)^{\frac{1}{2}} \frac{1}{(\lambda_k^*)^{a^*}} \bigg/ \sum_{k' \in \mathcal{K}} \left(\frac{|V_{k'}^*|}{|V_{k'}|} \right)^{\frac{1}{2}} \frac{1}{(\lambda_{k'}^*)^{a^*}} \end{aligned} \quad (11)$$

4 Inference based on a single Best model

For parametric or predictive inference, the statistician may proceed as follows:

1. select a single 'Best' model \mathcal{M}_{k^*} (optimal in some sense) from a set of available models $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \dots\}$ and then
2. perform Bayesian inference conditional on the selected model \mathcal{M}_{k^*} only and as if it had generated the data; i.e. based on

$$\begin{cases} y|\theta_{k^*}, \mathcal{M}_{k^*} \sim F(y|\theta_{k^*}, k^*) & , \text{ data generation from the sampling distribution} \\ \theta_{k^*}|\mathcal{M}_{k^*} \sim \Pi(\theta_{k^*}|k^*) & , \text{ within model parameter generation from the conditional prior} \end{cases}$$

4.1 Bayesian model selection: How to select the single 'Best' model \mathcal{M}_{k^*} .

Note 21. Model selection/choice is the selection of the "Best" statistical model from a set of competing statistical models $\mathcal{M} = \{\mathcal{M}_k; k \in \mathcal{K}\}$. It can be naturally performed as a parametric point estimation problem, statistical decision problem, or hypothesis test (all equivalent). See Criteria 22 and 23.

Criterion 22. We produce the Bayes estimator (of Bayes rule) $\hat{\delta}$ of parameter $k \in \mathcal{K}$ under a loss function $\ell(\delta, k, \theta_k)$ with decision space $\mathcal{D} = \mathcal{K}$ and posterior distribution $(k, \theta_k) \sim \Pi(k, \theta_k|y)$. I.e find $\hat{\delta}$:

$$\hat{\delta} = \arg \min_{\forall \delta \in \mathcal{K}} E_{\Pi}(\ell(\delta, k, \theta_k) | y) = \arg \min_{\forall \delta \in \mathcal{K}} \left(\int_{\Theta_k} \ell(\delta, k, \theta_k) d\Pi(k, \theta_k|y) \right)$$

Criterion 23. Under the 0 – 1 loss $\ell(\delta, k, \theta_k) = 1 (\delta \neq k)$ the best model corresponds to the MAP $\hat{\delta}$ that is model with the highest marginal model posterior probability $\pi(k|y)$, I.e $\hat{\delta}$ such as:

$$\hat{\delta} = \arg \max_{\forall k \in \mathcal{K}} \{\pi(k|y)\}$$

or equivalently $\hat{\delta}$ such that $\pi(\hat{\delta}|y) \geq \pi(k|y)$, for all $k \in \mathcal{K}$.

Bayes factors can be used for model selection. We re-define the Bayes factor to be suitable to the model selection framework.

Definition 24. Bayes factor $B_{k,j}(y)$ of model \mathcal{M}_k against model \mathcal{M}_j is defined as

$$B_{k,j}(y) = \frac{\pi(k|y)/\pi(k)}{\pi(j|y)/\pi(j)} = \frac{f(y|k)}{f(y|j)}; \quad \forall k, j \in \mathcal{K}$$

Example 25. (Cont.) The marginal model posterior $\pi(k|y)$ has been calculated in Example 17 in Eq. 9. Under the 0-1 loss $\ell(\delta, k, \theta_k) = 1$ ($\delta \neq k$), the best regression model is that with the largest marginal model posterior $\pi(k|y)$ (9). Consider data-set airquality {datasets} available from R repository, that involves $n = 153$ pairs of observations on the response variable y 'Ozone (ppb)' and the possible regressors

- ϕ_0 : Constant value 1
- ϕ_1 : Solar R (lang)
- ϕ_2 : Wind (mph)
- ϕ_3 : Temperature (degrees F)

The regressors have been standardized to adjust their unites in the model. The available models are: $\mathcal{M}_0 = \{1\}$, $\mathcal{M}_1 = \{1, 2\}$, $\mathcal{M}_2 = \{1, 3\}$, $\mathcal{M}_3 = \{1, 4\}$, $\mathcal{M}_4 = \{1, 2, 3\}$, $\mathcal{M}_5 = \{1, 2, 4\}$, $\mathcal{M}_6 = \{1, 3, 4\}$, and $\mathcal{M}_7 = \{1, 2, 3, 4\}$. The prior hyper-parameters are $\mu_k = 0$, $V_k = 100I_{d_k}$, $a_k = 1.0$, and $\lambda_k = 1.0$. The marginal model prior was $\pi(k|y) = 1/8$ for $k = \{0, \dots, 7\}$.

In Figure 1, we see the marginal model prior which is uniform, and the marginal model posterior which denotes that the MAP 'best' model is \mathcal{M}_6 .

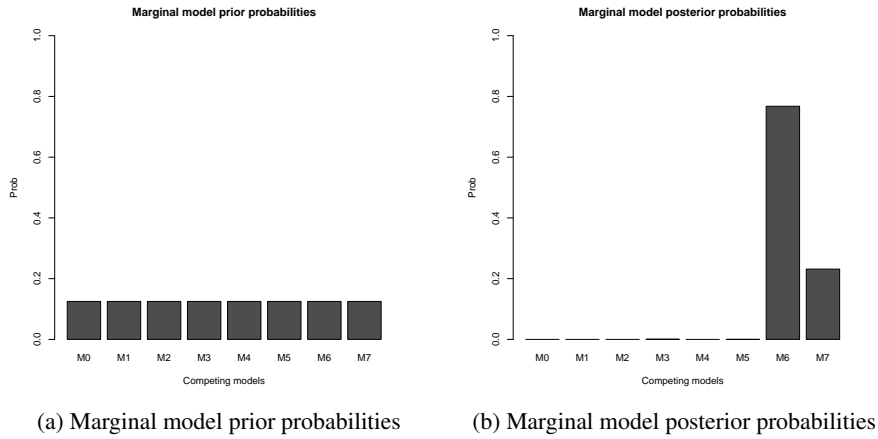


Figure 1: Marginal model prior and posterior probabilities produced in Example 25.
https://github.com/georgios-stats/Bayesian_Statistics/blob/master/HandoutsSupplementary/LinearRegressionModelUncertainty/MarginalModelPosteriorProbability.R

The Bayes factor is computed as

$$B_{k,j}(y) = \frac{f(y|k)}{f(y|j)} = \left(\frac{|V_j|}{|V_k|} \right)^{\frac{1}{2}} \left(\frac{|V_k^*|}{|V_j^*|} \right)^{\frac{1}{2}} \left(\frac{\lambda_j^*}{\lambda_k^*} \right)^{a^*} \quad (12)$$

4.2 Within model inference

Note 26. Given that \mathcal{M}_{k^*} is the selected “Best” model, interest may lie on learning the parameter $\theta_{k^*} \in \Theta_{k^*}$ (or any function of it), or future outputs $z = (y_{n+1}, \dots, y_{n+m})$. Then inference is performed based on the Bayesian model

$$\begin{cases} y|\theta_{k^*}, \mathcal{M}_{k^*} \sim F(y|\theta_{k^*}, k^*) & , \text{ data generation from the sampling distribution} \\ \theta_{k^*}|\mathcal{M}_{k^*} \sim \Pi(\theta_{k^*}|k^*) & , \text{ within model parameter generation from the conditional prior} \end{cases}$$

- Parametric inference about θ_{k^*} (or functions of it) is performed based on the conditional posterior in (8):

$$\theta_{k^*}|y, k^* \sim \Pi(\theta_{k^*}|y, k^*)$$

- Predictive inference about a sequence of future outcomes $z = (y_{n+1}, \dots, y_{n+m})$ is performed based on the conditional predictive distribution :

$$z|y, k^* \sim G(z|y, k^*)$$

in (10) with pdf

$$g(z|y, k^*) = \int f(z|y, \theta_{k^*}, k^*) d\Pi(\theta_{k^*}|y, k^*)$$

5 Bayesian model averaging (BMA)

Note 27. Selecting a single (best) model from a set of available ones and then proceeding as if the selected model had generated may not be optimal. It ignores the uncertainty due to model selection, which leads to over-confident inferences and decisions. BMA offers a coherent mechanism to account model uncertainty.

Definition 28. Model averaging is called the process where a posterior summary is obtained as a weighted mean of the summaries under each model with weights given by the marginal model posterior probability.

Parametric inference There may be certain model parameters or functions of model parameters which are present in all available models and retain a consistent interpretation across all models. Interesting may be to average out all the models.

Note 29. Let partition $\theta_k = (\psi, \phi_k)$, where $\psi \in \Psi$ is common in all models in \mathcal{M} . It is possible to obtain the marginal posterior distribution $\Pi(\psi|y)$ of ψ as a mixture of posterior distributions under each model with weights given by the marginal posterior probability of each model. i.e. $\Pi(\psi|y)$ has pdf/pmf

$$\pi(\psi|y) = \sum_{k \in \mathcal{K}} \pi(\psi|y, k) \pi(k|y); \quad \text{with} \quad \pi(\psi|y, k) = \int \pi(\psi, \phi_k|y, k) d\phi_k.$$

Note 30. Then any moment (expectation, variance, probability, etc...) of $\psi \in \Psi$ can be expressed in the aforesaid weighted average form. Let $h(\cdot)$ be a function defined on Ψ . Then

$$E_{\Pi}(h(\psi)|y) = \sum_{k \in \mathcal{K}} \int_{\Theta_k} h(\psi) d\Pi(\psi|y, k) \pi(k|y) = \sum_{k \in \mathcal{K}} E_{\Pi}(h(\psi)|y, k) \pi(k|y)$$

Example 31. The sampling distribution variance σ^2 is the common parameter across all models, so

$$\pi(\sigma^2|y) = \sum_{k \in \mathcal{K}} \frac{\pi(\sigma^2|y, k)}{\text{IG}(\sigma^2|a^*, \lambda_k^*)} \pi(k|y) = \frac{\sum_{k \in \mathcal{K}} \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2} + a^* + 1} \exp\left(-\frac{\lambda_k^*}{\sigma^2}\right) \left(\frac{|V_k^*|}{|V_k|}\right)^{\frac{1}{2}} / (\lambda_k^*)^{a^*}}{\sum_{k' \in \mathcal{K}} \left(\frac{|V_{k'}^*|}{|V_{k'}|}\right)^{\frac{1}{2}} / (\lambda_{k'}^*)^{a^*}} 1 \quad (\sigma^2 > 0)$$

188 and, because $E_{\sigma^2 \sim \text{IG}(a^*, \lambda_k^*)}(\sigma^2) = \frac{\lambda_k^*}{a^* - 1}$ for $a^* > 1$, we can get its marginal posterior expectation as

$$189 \quad E(\sigma^2|y) = \sum_{k \in \mathcal{K}} E_{\Pi}(\sigma^2|y, k) \pi(k|y) = \frac{\sum_{k \in \mathcal{K}} \frac{\lambda_k^*}{a^* - 1} \left(\frac{|V_k^*|}{|V_k|} \right)^{\frac{1}{2}} / (\lambda_k^*)^{a^*}}{\sum_{k' \in \mathcal{K}} \left(\frac{|V_{k'}^*|}{|V_{k'}|} \right)^{\frac{1}{2}} / (\lambda_{k'}^*)^{a^*}}$$

190 **Predictive inference** BMA plays an important role in the predictive inference.

191 *Note 32.* Let $z = (y_{n+1}, \dots, y_{n+m})$ be a sequence of future outcomes. Then the predictive distribution $dG(z|y)$ of z
192 given the observables y has pdf/pmf

$$193 \quad g(z|y) = \sum_{k \in \mathcal{K}} g(z|y, k) \pi(k|y)$$

194 where

$$195 \quad g(z|y, k) = \int f(z|y, \theta_k, k) \pi(\theta_k|y, k) d\theta_k$$

196 is the predictive pdf of z given model \mathcal{M}_k .

197 **Example 33.** The BMA predictive distribution for the Normal linear regression model has been calculated in Example
198 19, Eq. c11.

199 *Note 34.* By averaging out all models in this fashion, BMA provides better predictive ability as measured by the
200 logarithmic scoring rule $-\log(\pi(A))$, than using any single model \mathcal{M}_k , conditional on \mathcal{M} . See the exercise below.

201 6 Discussion

202 In model uncertainty framework, we often meet two situations;

203 **\mathcal{M} -close perspective:** the encompassing Bayesian model includes the real data-generation process $R(\cdot)$. It is argued
204 that this is not a realistic situation. The introduced model uncertainty framework is suitable in this situation.

205 **\mathcal{M} -open perspective:** the encompassing Bayesian model does not include the real data-generation process $R(\cdot)$. This
206 is a more realistic situation. It is argued that the introduced model uncertainty framework is not suitable in
207 this situation in the sense that $\pi(\mathcal{M}_k)$ does not make sense to exist.

208 What do we do? This is a rather open problem, and may depends on the application under consideration. Advise; Try
209 to set-up an encompassing Bayesian model rich enough to promise, in some sense, that it acceptably approximate (at
210 least) the real data generating process.

212 **Question 35.** For practice address the Exercises 69, 70, and 71, from the Exercise sheet.

213 **Exercise 36.** Show that for any $k \in \mathcal{K}$

$$214 \quad -E \left(\log \left(\sum_{k \in \mathcal{K}} \pi(z|y, k) \pi(k|y) \right) \right) \leq -E(\log(\pi(z|y, k)))$$

215 where the expectations are with respect to $\sum_{k \in \mathcal{K}} \pi(z|y, k) \pi(k|y)$.

216 **Solution.** It is

$$217 \quad -E \left(\log \left(\sum_{k \in \mathcal{K}} \pi(z|y, k) \pi(k|y) \right) \right) \leq -E(\log(\pi(z|y, k))) \iff 0 \leq \underbrace{-E \left(\log \left(\frac{\pi(z|y, k)}{\sum_{k \in \mathcal{K}} \pi(z|y, k) \pi(k|y)} \right) \right)}_{= \text{KL}(\pi(z|y) || \pi(z|y, k))}$$

A Solution to Example 16

Use the following identity

$$(y - \Phi\beta)^\top (y - \Phi\beta) + (\beta - \mu)^\top V^{-1}(\beta - \mu) = (\beta - \mu^*)^\top (V^*)^{-1}(\beta - \mu^*) + S^*; \\ S^* = \mu^\top V^{-1}\mu - (\mu^*)^\top (V^*)^{-1}(\mu^*) + y^\top y; \quad V^* = (V^{-1} + \Phi^\top \Phi)^{-1}; \quad \mu^* = V^* (V^{-1}\mu + \Phi^\top y)$$

Solution. For simplicity, we suppress the indexing \cdot_k denoting the sub-set of the regressors. It is

$$\begin{aligned} \pi(\beta_k, \sigma^2 | y, k) &\propto f(y | \theta_k, k) \pi(\theta_k | k) = N(y | \Phi\beta_k, I\sigma^2) N(\beta_k | \mu_k, V_k \sigma^2) \text{IG}(\sigma^2 | a, \lambda) \\ &\propto \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2}(y - \Phi\beta)^\top (y - \Phi\beta)\right) \times \left(\frac{1}{\sigma^2}\right)^{\frac{d}{2}} \exp\left(-\frac{1}{2\sigma^2}(\beta - \mu)^\top V^{-1}(\beta - \mu)\right) \\ &\quad \times \left(\frac{1}{\sigma^2}\right)^{a+1} \exp\left(-\frac{\lambda}{\sigma^2}\right) \\ &= \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2} + \frac{d}{2} + a + 1} \exp\left(-\frac{1}{2\sigma^2}(y - \Phi\beta)^\top (y - \Phi\beta) - \frac{1}{2\sigma^2}(\beta - \mu)^\top V^{-1}(\beta - \mu)\right) \exp\left(-\frac{\lambda}{\sigma^2}\right) \\ &= \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2} + \frac{d}{2} + a + 1} \exp\left(-\frac{1}{2\sigma^2}[(y - \Phi\beta)^\top (y - \Phi\beta) + (\beta - \mu)^\top V^{-1}(\beta - \mu)]\right) \exp\left(-\frac{\lambda}{\sigma^2}\right) \\ &= \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2} + \frac{d}{2} + a + 1} \exp\left(-\frac{1}{2\sigma^2}(\beta - \mu^*)^\top (V^*)^{-1}(\beta - \mu^*) - \frac{1}{2\sigma^2}S^*\right) \exp\left(-\frac{\lambda}{\sigma^2}\right) \\ &= \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2} + \frac{d}{2} + a + 1} \exp\left(-\frac{1}{2\sigma^2}(\beta - \mu^*)^\top (V^*)^{-1}(\beta - \mu^*)\right) \exp\left(-\frac{1}{\sigma^2}\left(\lambda + \frac{1}{2}S^*\right)\right) \\ &\propto \underbrace{\left(\frac{1}{\sigma^2}\right)^{\frac{d}{2}} \exp\left(-\frac{1}{2\sigma^2}(\beta - \mu^*)^\top (V_k^*)^{-1}(\beta - \mu^*)\right)}_{\propto N(\beta_k | \mu^*, V^* \sigma^2)} \times \underbrace{\left(\frac{1}{\sigma^2}\right)^{\frac{n}{2} + a + 1} \exp\left(-\frac{1}{\sigma^2}\left(\lambda + \frac{1}{2}S^*\right)\right)}_{\propto \text{IG}(\sigma^2 | a^*, \lambda^*)} \\ &\propto N(\beta | \mu^*, V^* \sigma^2) \text{IG}(\sigma^2 | a^*, \lambda^*) \end{aligned}$$

where

$$S^* = \mu^\top V^{-1}\mu - (\mu^*)^\top (V^*)^{-1}(\mu^*) + y^\top y \\ a^* = \frac{n}{2} + a; \quad \lambda^* = \lambda + \frac{1}{2}S^*; \quad V^* = (V^{-1} + \Phi^\top \Phi)^{-1}; \quad \mu^* = V^* (V^{-1}\mu + \Phi^\top y)$$

The result in Example 16 results by indexing with \cdot_k the quantities that depend on it. I.e.

$$\begin{aligned} \pi(\beta_k, \sigma^2 | y, k) &= N(\beta_k | \mu_k^*, V_k^* \sigma^2) \text{IG}(\sigma^2 | a^*, \lambda_k^*) \\ S_k^* &= \mu_k^\top V_k^{-1}\mu_k - (\mu_k^*)^\top (V_k^*)^{-1}(\mu_k^*) + y^\top y \\ a^* &= \frac{n}{2} + a; \quad \lambda_k^* = \lambda + \frac{1}{2}S_k^*; \quad V_k^* = (V_k^{-1} + \Phi_k^\top \Phi_k)^{-1}; \quad \mu_k^* = V_k^* (V_k^{-1}\mu_k + \Phi_k^\top y) \end{aligned}$$

Hence

$$\Rightarrow \begin{cases} \beta_k | y \sigma^2, k & \sim N(\mu_k^*, V_k^* \sigma^2) \\ \sigma^2 | y, k & \sim \text{IG}(a^*, \lambda^*) \end{cases}, \quad \forall k \in \mathcal{K}$$

B Solution to Example 17

Use the following identity

$$(y - \Phi\beta)^\top (y - \Phi\beta) + (\beta - \mu)^\top V^{-1}(\beta - \mu) = (\beta - \mu^*)^\top (V^*)^{-1}(\beta - \mu^*) + S^*; \\ S^* = \mu^\top V^{-1}\mu - (\mu^*)^\top (V^*)^{-1}(\mu^*) + y^\top y; \quad V^* = (V^{-1} + \Phi^\top \Phi)^{-1}; \quad \mu^* = V^* (V^{-1}\mu + \Phi^\top y)$$

Solution. For simplicity, we suppress the indexing \cdot_k denoting the sub-set of the regressors.

The conditional marginal likelihood $f(y|k)$ (or model evidence of \mathcal{M}_k) is

$$\begin{aligned} f(y|k) &= \int_{\Theta_k} f(y|\theta_k, k) \pi(\theta_k|k) d\theta_k = \int \mathbf{N}(y|\Phi\beta, I\sigma^2) \mathbf{N}(\beta|\mu, V\sigma^2) \text{IG}(\sigma^2|a, \lambda) d\beta d\sigma^2 \\ &= \int \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} \exp \left(-\frac{1}{2\sigma^2} (y - \Phi\beta)^\top (y - \Phi\beta) \right) \times \left(\frac{1}{\sigma^2} \right)^{\frac{n}{2}} \left(\frac{1}{|V|} \right)^{\frac{1}{2}} \exp \left(-\frac{1}{2\sigma^2} (\beta - \mu)^\top V^{-1}(\beta - \mu) \right) \\ &\quad \times \frac{\lambda^a}{\Gamma(a)} \exp \left(-\frac{\lambda}{\sigma^2} \right) d\beta d\sigma^2 \\ &= \left(\frac{1}{2\pi} \right)^{\frac{n}{2} + \frac{d}{2}} \left(\frac{1}{|V|} \right)^{\frac{1}{2}} \frac{\lambda^a}{\Gamma(a)} \int \left(\frac{1}{\sigma^2} \right)^{\frac{n}{2} + \frac{d}{2} + a + 1} \exp \left(-\frac{1}{\sigma^2} \left(\frac{(y - \Phi\beta)^\top (y - \Phi\beta) + (\beta - \mu)^\top V^{-1}(\beta - \mu)}{2} + \lambda \right) \right) \\ &= \left(\frac{1}{2\pi} \right)^{\frac{n}{2} + \frac{d}{2}} \left(\frac{1}{|V|} \right)^{\frac{1}{2}} \frac{\lambda^a}{\Gamma(a)} \int \left(\frac{1}{\sigma^2} \right)^{\frac{n}{2} + \frac{d}{2} + a + 1} \exp \left(-\frac{1}{\sigma^2} \left(\frac{1}{2} S^* + \lambda \right) \right) \\ &\quad \times \int \left[\exp \left(-\frac{1}{2} \frac{1}{\sigma^2} (\beta - \beta^*)^\top (V^*)^{-1}(\beta - \beta^*) \right) d\beta \right] d\sigma^2 \\ &= \left(\frac{1}{2\pi} \right)^{\frac{n}{2} + \frac{d}{2}} \left(\frac{1}{|V|} \right)^{\frac{1}{2}} \frac{\lambda^a}{\Gamma(a)} \int \left(\frac{1}{\sigma^2} \right)^{\frac{n}{2} + \frac{d}{2} + a + 1} \exp \left(-\frac{1}{\sigma^2} \left(\frac{1}{2} S^* + \lambda \right) \right) \times \left[(2\pi)^{\frac{d}{2}} (\sigma^2)^{\frac{d}{2}} |V^*|^{\frac{1}{2}} \right] d\sigma^2 \\ &= \left(\frac{1}{2\pi} \right)^{\frac{n}{2}} \left(\frac{|V^*|}{|V|} \right)^{\frac{1}{2}} \frac{k^a}{\Gamma(a)} \frac{\Gamma(\frac{n}{2} + a)}{(\frac{1}{2} S^* + \lambda)^{\frac{n}{2} + a}} = \left(\frac{1}{2\pi} \right)^{\frac{n}{2}} \left(\frac{|V^*|}{|V|} \right)^{\frac{1}{2}} \frac{k^a}{\Gamma(a)} \frac{\Gamma(a^*)}{(\lambda^*)^{a^*}} \end{aligned}$$

$$\text{where} \quad S^* = \mu^\top V^{-1}\mu - (\mu^*)^\top (V^*)^{-1}(\mu^*) + y^\top y$$

$$a^* = \frac{n}{2} + a; \quad \lambda^* = \lambda + \frac{1}{2} S^*; \quad V^* = (V^{-1} + \Phi^\top \Phi)^{-1}; \quad \mu^* = V^* (V^{-1}\mu + \Phi^\top y)$$

The result in Example 17 results by indexing with \cdot_k the quantities that depend on it. I.e.

$$\begin{aligned} f(y|k) &= \left(\frac{1}{2\pi} \right)^{\frac{n}{2}} \left(\frac{|V_k^*|}{|V_k|} \right)^{\frac{1}{2}} \frac{\lambda^a}{\Gamma(a)} \frac{\Gamma(a^*)}{(\lambda_k^*)^{a^*}} \\ S_k^* &= \mu_k^\top V_k^{-1}\mu_k - (\mu_k^*)^\top (V_k^*)^{-1}(\mu_k^*) + y^\top y \\ a^* &= \frac{n}{2} + a; \quad \lambda_k^* = \lambda + \frac{1}{2} S_k^*; \quad V_k^* = (V_k^{-1} + \Phi_k^\top \Phi_k)^{-1}; \quad \mu_k^* = V_k^* (V_k^{-1}\mu_k + \Phi_k^\top y) \end{aligned}$$

The marginal likelihood $f(y)$ is

$$f(y) = \sum_{k \in \mathcal{K}} f(y|k) \pi(k) = \left(\frac{1}{2\pi} \right)^{\frac{n}{2}} \lambda^a \frac{\Gamma(a^*)}{\Gamma(a)} \frac{1}{2^d} \sum_{k \in \mathcal{K}} \left(\frac{|V_k^*|}{|V_k|} \right)^{\frac{1}{2}} \frac{1}{(\lambda_k^*)^{a^*}}$$

The marginal model posterior probability $\pi(k|y)$ of \mathcal{M}_k is

$$\pi(k|y) = \frac{f(y|k) \pi(k)}{f(y)} = \left(\frac{|V_k^*|}{|V_k|} \right)^{\frac{1}{2}} \frac{1}{(\lambda_k^*)^{a^*}} \bigg/ \sum_{k' \in \mathcal{K}} \left(\frac{|V_{k'}^*|}{|V_{k'}|} \right)^{\frac{1}{2}} \frac{1}{(\lambda_{k'}^*)^{a^*}}$$

C Solution to Example 19

Use the following identity:

$$\begin{cases} y|\beta \sim N_m(\Phi\beta + c, \Sigma) \\ \beta \sim N_d(\mu, V) \end{cases} \implies y \sim N_m(\Phi\mu + c, \Sigma + \Phi V \Phi^\top)$$

Use the following property of the Student T distribution

$$\begin{cases} x|\xi \sim N_m(\mu, \Sigma\xi v) \\ \xi \sim \text{IG}(\frac{v}{2}, \frac{1}{2}) \end{cases} \implies x \sim T_m(\mu, \Sigma, v)$$

Solution. For simplicity, we suppress the indexing \cdot_k denoting the sub-set of the regressors. The conditional predictive distribution $G(z|y, k)$ given model \mathcal{M}_k and data y has pdf

$$\begin{aligned} g(z|y, k) &= \int f(z|\theta_k, k) \pi(\theta_k|y, k) d\theta_k = \int N_m(z|\Phi^{\text{new}}\beta, I\sigma^2) N_d(\beta|\mu^*, V^*\sigma^2) \text{IG}(\sigma^2|a^*, \lambda^*) d\beta d\sigma^2 \\ &= \int \left[\int N_m(z|\Phi^{\text{new}}\beta, I\sigma^2) N_d(\beta|\mu^*, V^*\sigma^2) d\beta \right] \text{IG}(\sigma^2|a^*, \lambda^*) d\sigma^2 \\ &= \int N_m(z|\Phi^{\text{new}}\mu^*, \sigma^2[I + V^*]) \text{IG}(\sigma^2|a^*, \lambda^*) d\sigma^2 \\ &= \int N_m(z|\Phi^{\text{new}}\mu^*, \sigma^2[I + V^*]) \text{IG}(\sigma^2|a^*, \lambda^*) d\sigma^2 \\ &= \int N_m\left(z|\Phi^{\text{new}}\mu^*, \sigma^2 \frac{\lambda^*}{\lambda^*} \frac{2a^*}{2a^*} [I + V^*]\right) \text{IG}(\sigma^2|a^*, \lambda^*) d\sigma^2 \\ &\stackrel{\xi = \frac{\sigma^2}{2\lambda^*} \sim \text{IG}(a^*, \frac{1}{2})}{=} \int N_m\left(z|\Phi^{\text{new}}\mu^*, \xi \lambda^* \frac{2a^*}{a^*} [I + V^*]\right) \text{IG}\left(\xi|a^*, \frac{1}{2}\right) d\xi \\ &= \int N_m\left(z|\Phi^{\text{new}}\mu^*, \xi 2a^* \frac{\lambda^*}{a^*} [I + V^*]\right) \text{IG}\left(\xi|\frac{a^*}{2}, \frac{1}{2}\right) d\xi \\ &= T_m\left(z|\Phi^{\text{new}}\mu^*, [I + V^*] \frac{\lambda^*}{a^*}, 2a^*\right) \end{aligned}$$

The result in Example 19 results by indexing with \cdot_k the quantities that depend on it. I.e.

$$g(z|y, k) = T_m\left(z|\Phi^{\text{new}}\mu_k^*, [I + V_k^*] \frac{\lambda_k^*}{a^*}, 2a^*\right)$$

The marginal predictive distribution $G(z|y)$ given data y has pdf/pmf

$$\begin{aligned} g(z|y) &= \sum_{k \in \mathcal{K}} g(z|y, k) \pi(k|y) \\ &= \sum_{k \in \mathcal{K}} T_m\left(z|\Phi^{\text{new}}\mu_k^*, V_k^* \frac{\lambda_k^*}{a^*}, 2a^*\right) \left(\frac{|V_k^*|}{|V_k|}\right)^{\frac{1}{2}} \frac{1}{(\lambda_k^*)^{a^*}} \bigg/ \sum_{k' \in \mathcal{K}} \left(\frac{|V_{k'}^*|}{|V_{k'}|}\right)^{\frac{1}{2}} \frac{1}{(\lambda_{k'}^*)^{a^*}} \end{aligned}$$