

# (Predictive) model selection

Marilena Barbieri

Università Roma Tre

O'Bayes 2019: Objective Bayes Methodology Conference  
University of Warwick, 28 June - 02 July 2019

# (Predictive) model selection

Marilena Barbieri

Università Roma Tre

## Outline:

- ▶ Basics of Bayesian model selection
- ▶ Variable selection in linear model: the Median probability model

O'Bayes 2019: Objective Bayes Methodology Conference  
University of Warwick, 28 June - 02 July 2019

## General model selection: notation

Data  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  is assumed to have arisen from one of several **models**:

$M_1$ :  $\mathbf{y}$  has density  $f_1(\mathbf{y}|\beta_1)$

$M_2$ :  $\mathbf{y}$  has density  $f_2(\mathbf{y}|\beta_2)$

$\vdots$

$M_q$ :  $\mathbf{y}$  has density  $f_q(\mathbf{y}|\beta_q)$

$\beta_i$  ( $i = 1, \dots, q$ ) model specific parameter with dimension  $h_i$ .

Assign a **prior probability**  $P(M_i)$  to each model.

A common choice:  $P(M_i) = 1/q$

For **small**  $q$ ,  $P(M_i) = 1/q$  is sensible, although having decreasing  $P(M_i)$  as model complexity increases is often reasonable.

Care is needed if  $q$  is **huge** and multiplicity is an issue.

**Linear models** with  $k$  variables and a common constant:

$$P(M_i) = \frac{1}{k+1} \binom{k}{h_i}^{-1}$$

which is equivalent to assigning the set of all models having the same number of parameters probability  $1/(k+1)$  and dividing this probability equally among them.

Under model  $M_j$ :

**Prior density** of  $\beta_i$ :  $\pi_i(\beta_i)$

**Marginal density**:  $m_i(\mathbf{y}) = \int f_i(\mathbf{y}|\beta_i)\pi_i(\beta_i)d\beta_i$   
which measures “how likely is  $y$  under  $M_i$ ”

**Posterior density**:  $\pi_i(\beta_i|\mathbf{y}) = f_i(\mathbf{y}|\beta_i)\pi_i(\beta_i)/m_i(\mathbf{y})$

Bayes factor of  $M_j$  to  $M_i$ :

$$B_{ji} = \frac{m_j(\mathbf{y})}{m_i(\mathbf{y})}$$

Posterior probability of  $M_i$ :

$$P(M_i|\mathbf{y}) = \frac{P(M_i)m_i(\mathbf{y})}{\sum_{j=1}^q P(M_j)m_j(\mathbf{y})} = \left[ \sum_{j=1}^q \frac{P(M_j)}{P(M_i)} B_{ji} \right]^{-1}$$

# An important note of caution

- ▶ Improper priors are problematical, because they are unnormalized:

$$\pi_i(\beta_i) = c_i g_i(\beta_i)$$

Can not be used for parameters not occurring in all models since  $c_i$  is arbitrary and  $B_{ij}$  is arbitrary

$$B_{ji} = \frac{\int f_j(\mathbf{y}|\beta_j)\pi_j(\beta_j)d\beta_j}{\int f_i(\mathbf{y}|\beta_i)\pi_i(\beta_i)d\beta_i} = \frac{c_j \int f_j(\mathbf{y}|\beta_j)g_j(\beta_j)d\beta_j}{c_i \int f_i(\mathbf{y}|\beta_i)g_i(\beta_i)d\beta_i}$$

- ▶ *Vague* proper priors (arbitrarily large variances) may be even worse.

# Basic of Bayes prediction

For each model  $M_i$ , a future observation  $y^*$  has density  $g_i(y^*|\mathbf{y}, \beta_i)$ .

The **posterior predictive distribution** of  $y^*$ , once  $\mathbf{y}$  has been observed, is defined as

$$p_i(y^*|\mathbf{y}) = \int g_i(y^*|\mathbf{y}, \beta_i) \pi_i(\beta_i|\mathbf{y}) d\beta_i$$

A nice feature of the posterior predictive distribution is that the posterior predictive is proper, whenever the posterior distribution  $\pi_i(\beta_i|\mathbf{y})$  is proper.

# Multiple models setting

When considering the set of parametric models  $\{M_i\}$ , the **posterior predictive distribution**, obtained by Bayesian Model Averaging, is

$$\begin{aligned} p(y^*|\mathbf{y}) &= \sum_i \int g_i(y^*|\mathbf{y}, \beta_i) \pi_i(\beta_i|\mathbf{y}) P(M_i|\mathbf{y}) d\beta_i = \\ &= \sum_i P(M_i|\mathbf{y}) p_i(y^*|\mathbf{y}) \end{aligned}$$

# Basic of predictive model selection

Comparing the predictive performances of several models is often (if not always explicitly) formulated as a decision problem. The optimal predictive model is the one that minimizes the expected loss

$$\int L(m_i, a_i^*, y^*) p(y^* | \mathbf{y}) dy^*$$

over  $i = 1, 2, \dots, q$ , where  $m_i$  corresponds to selecting  $M_i$  and  $a_i^*$  denotes the optimal subsequent decision regarding a future observation  $y^*$ .

# Basic of predictive model selection

The expected loss takes the form

$$- \int \log [p_i(y^*|\mathbf{y})] p(y^*|\mathbf{y}) dy^*$$

for a predictive distribution of  $y^*$  with **logarithmic score function** and

$$\int [y^* - \hat{y}_i^*]^2 p(y^*|\mathbf{y}) dy^*$$

for point prediction with **quadratic loss**, where  $\hat{y}_i^* = \mathbb{E}(y^*|\mathbf{y}, M_i)$  is the optimal prediction of a future observation  $y^*$ , given  $\mathbf{y}$  and assuming model  $M_i$ , that is the value which minimizes

$$\int (y^* - \hat{y}_i^*)^2 p_i(y^*|\mathbf{y}) dy^*.$$

# Motivation for the Bayesian approach

1. **Ease of interpretation:**

$P(M_i|\mathbf{y})$  reflects real expected error rates.

2. **Prior information can be incorporated**, if desired:

It is useful to separately report  $\{m_i(\mathbf{y})\}$  (or  $\{B_{ji}\}$ ) and  $\{P(M_i)\}$ . This allows computation of the  $P(M_i|\mathbf{y})$  for any prior probabilities.

3. **Consistency:**

If one of the  $M_i$  is true, then  $P(M_i|\mathbf{y}) \rightarrow 1$  as  $n \rightarrow \infty$ .

# Motivation for the Bayesian approach

## 4. Ockham's razor:

Bayes Factors automatically seek parsimony; no adhoc penalties for model complexity are needed. (Jefferys and Berger, 1992)

## 5. Generality of application:

The approach is viable for any models, regular or irregular, nested or not, large or small sample sizes, two or multiple models.

# Motivation for the Bayesian approach

## 6. Accounting for model uncertainty:

Selecting a single model and using it for inference ignores model uncertainty, resulting in inferior inferences, considerable overstatements of accuracy.

The Bayesian approach incorporates this uncertainty by model averaging: if inference concerning  $\theta$  (*same* meaning across models) is desired, it would be based on

$$\pi(\theta|\mathbf{y}) = \sum_i P(M_i|\mathbf{y}) \pi_i(\theta|\mathbf{y}).$$

Most frequent use is for prediction (Geisser 93, Draper 95, Raftery et al. 97, Clyde 99, Clyde and George, 04, . . . )

# Prediction with Normal linear models

Under the full model, the  $n \times 1$  observation vector would follow

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

where  $\mathbf{X}$  is the  $n \times k$  ( $k < n$ ) full rank design matrix of covariates,  $\boldsymbol{\beta}$  is a  $k \times 1$  vector of unknown coefficients and  $\boldsymbol{\varepsilon}$  is  $N(0, \sigma^2 I)$ .

The possible submodels are

$$M_I : \mathbf{y} = \mathbf{X}_I \boldsymbol{\beta}_I + \boldsymbol{\varepsilon}$$

where  $I = (I_1, I_2, \dots, I_k)$  is the model index,  $I_i$  being either 1 or 0 as covariate  $x_i$  is in or out of the model.

Assume that one of these models is true, and our goal is to predict a future observation at covariate values  $\mathbf{x}^*$

$$y^* = \mathbf{x}^* \boldsymbol{\beta} + \varepsilon$$

using square error loss

$$(y^* - \hat{y}^*)^2.$$

The best predictor of  $y^*$  is, via model averaging, the posterior mean of  $y^*$

$$\bar{y}^* = \mathbf{x}^* \bar{\boldsymbol{\beta}} \equiv \mathbf{x}^* \sum_I P(M_I | \mathbf{y}) \tilde{\boldsymbol{\beta}}_I$$

where  $\tilde{\boldsymbol{\beta}}_I$  is the posterior mean of  $\boldsymbol{\beta}$  under  $M_I$  (adding any needed zeroes to the  $\boldsymbol{\beta}$ )

# Selecting a single model

Often a single model  $M_I$  is desired for prediction, with the prediction then being  $\hat{y}_I^* = \mathbf{x}^* \tilde{\beta}_I$ .

It is commonly perceived that the best model will be that with the highest  $P(M_I|\mathbf{y})$ . This is true if:

- ▶ there are only two models;
- ▶  $X'X$  is diagonal,  $\sigma^2$  is known, and suitable priors are used.

## Selecting a single model

The best single model is the one that minimizes (where the expectation is with respect to the predictive distribution of  $y^*$  and ignoring irrelevant terms)

$$\begin{aligned}\mathbb{E}[(y^* - \hat{y}_I^*)^2] &= (\hat{y}_I^* - \bar{y}^*)^2 = \\ &= (\tilde{\beta}_I - \bar{\beta})' \mathbf{x}^{*'} \mathbf{x}^* (\tilde{\beta}_I - \bar{\beta})\end{aligned}$$

and will typically depend on  $\mathbf{x}^*$ .

## Selecting a single model

An important case is when the model will repeatedly be used to make future predictions, and covariates that will occur in the future are like those that occurred in the data, in the sense that  $\mathbb{E}[(\mathbf{x}^*)'(\mathbf{x}^*)] = \mathbf{X}'\mathbf{X}$

Then the averaged squared error predictive loss for future predictions, using  $M_I$  is

$$\mathbb{E}^{\mathbf{x}^*} (\hat{y}_I^* - \bar{y}^*)^2 = (\tilde{\beta}_I - \bar{\beta})' \mathbf{X}'\mathbf{X} (\tilde{\beta}_I - \bar{\beta})$$

# Posterior inclusion probabilities

The *posterior inclusion probability* for variable  $i$  is

$$p_i \equiv \sum_{I: I_i=1} P(M_I | \mathbf{y}),$$

that is, the overall posterior probability that variable  $i$  is in the model.

These are of interest in defining the (posterior) median probability model.

# The (posterior) median probability model

If it exists, the *median probability model*,  $M_{J^*}$ , is defined to be the model consisting of those variables whose posterior inclusion probability is at least  $1/2$ .

Formally,  $J^*$  is defined, coordinatewise, by

$$J_i^* = \begin{cases} 1, & \text{if } p_i \geq \frac{1}{2}, \\ 0, & \text{otherwise.} \end{cases}$$

Note: If computation is done by MCMC, the median probability model consists of those coordinates that were present in over half the iterations.

# Existence of the median probability model

The median probability model exists when the models under consideration follow a *graphical model structure*, including:

- ▶ when any subset of variables is allowed;
- ▶ the situation in which the allowed variables consist of main effects and interactions, but a higher order interaction is allowed only if lower order interactions are included;
- ▶ a sequence of nested models, such as arises in polynomial regression and autoregressive time series.

# Example: Polynomial Regression

Model  $M_j$  is

$$y = \sum_{i=0}^j \beta_i x^i + \varepsilon$$

(Model) $j$	0	1	2	3	4	5	6
$P(M_j y)$	$\approx 0$	0.06	0.22	0.29	0.38	0.05	$\approx 0$

(Covariate) $i$	0	1	2	3	4	5	6
$P(x^i \text{ is in model} y)$	1	1	0.94	0.72	0.43	0.05	0

Thus  $M_3$  is the median probability (optimal predictive) model, while  $M_4$  is the maximum probability model.

For nested models, the median probability model has a simpler representation as  $M_{I(j^*)}$ , where  $j^*$  is such that

$$\sum_{i=0}^{j^*-1} P(M_{I(i)}|\mathbf{y}) < \frac{1}{2} \quad \text{and} \quad \sum_{i=0}^{j^*} P(M_{I(i)}|\mathbf{y}) \geq \frac{1}{2}.$$

In other words, one just lists the sequence of posterior model probabilities and sums them up until the sum exceeds 1/2. The model at which the exceedance occurs is the median probability model.

# Optimality theorems

**Theorem 1.** If

- (i) the models under consideration have graphical structure,
- (ii)  $X'X$  is diagonal, and
- (iii) the posterior mean of  $\beta_I$  is simply the relevant coordinates of  $\tilde{\beta}$  (the posterior mean in the full model),

then the best predictive model is the median probability model.

Condition (iii) is satisfied under any mix of

constant priors for the  $\beta_i$ ;

independent  $N(0; \sigma^2 \lambda_i)$  priors for the  $\beta_i$ , with the  $\lambda_i$  given (objectively or subjectively specified, or estimated via empirical Bayes) and any prior for  $\sigma^2$ .

**Theorem 2.** Suppose a sequence of nested linear models is under consideration. If

- (i) prediction is desired at 'future covariates like the past' and
- (ii) the posterior mean under  $M_I$  satisfies  $\tilde{\beta}_I = b\hat{\beta}_I$ , where  $b > 0$ ; that is, the posterior means are proportional to the least squares estimates, with the same proportionality constant across models

then the best predictive model is the median probability model.

Condition (ii) is satisfied if we use either

the objective priors for model parameters; or

*g*-type  $\mathcal{N}_{k_I}(\mathbf{0}, c\sigma^2(\mathbf{X}'_I\mathbf{X}_I)^{-1})$  prior with the same constant  $c > 0$  for each model and any prior for  $\sigma^2$ .

Theorems 1 and 2 essentially remain true even if there are non-orthogonal nuisance parameters (i.e., parameters common to all models) that are assigned the usual noninformative priors.

## Example.

Consider Hald's regression data [Draper and Smith (1981)], consisting of  $n = 13$  observations on a dependent variable  $y$ , with four (highly correlated) potential regressors:  $x_1, x_2, x_3, x_4$ . The full model is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \varepsilon$$

$\varepsilon \sim N(0, \sigma^2)$ ,  $\sigma^2$  unknown.

- ▶ All models that include the constant term are considered.
- ▶ Least squares estimates are used for parameters.
- ▶ Default posterior probabilities of each model are obtained using the Encompassing Arithmetic Intrinsic Bayes Factor (Berger and Pericchi, 1996) for linear models, together with equal prior model probabilities.

Model	$P(M_I \mathbf{y})$	$R(M_I)$	Model	$P(M_I \mathbf{y})$	$R(M_I)$
c	0.000003	2652.44	c,2,3	0.000229	353.72
c,1	0.000012	1207.04	c,2,4	0.000018	821.15
c,2	0.000026	854.85	c,3,4	0.003785	118.59
c,3	0.000002	1864.41	c,1,2,3	0.170990	1.21
c,4	0.000058	838.20	c,1,2,4	0.190720	0.18
c,1,2	0.275484	8.19	c,1,3,4	0.159959	1.71
c,1,3	0.000006	1174.14	c,2,3,4	0.041323	20.42
c,1,4	0.107798	29.73	c,1,2,3,4	0.049587	0.47

Computation of the posterior inclusion probabilities yields

$$p_1 = \sum_{I:l_1=1} P(M_I|\mathbf{y}) = 0.954556, p_2 = \sum_{I:l_2=1} P(M_I|\mathbf{y}) = 0.728377,$$

$$p_3 = \sum_{I:l_3=1} P(M_I|\mathbf{y}) = 0.425881, p_4 = \sum_{I:l_4=1} P(M_I|\mathbf{y}) = 0.553248.$$

Model	$P(M_I y)$	$R(M_I)$	Model	$P(M_I y)$	$R(M_I)$
c	0.000003	2652.44	c,2,3	0.000229	353.72
c,1	0.000012	1207.04	c,2,4	0.000018	821.15
c,2	0.000026	854.85	c,3,4	0.003785	118.59
c,3	0.000002	1864.41	c,1,2,3	0.170990	1.21
c,4	0.000058	838.20	c,1,2,4	0.190720	0.18
c,1,2	0.275484	8.19	c,1,3,4	0.159959	1.71
c,1,3	0.000006	1174.14	c,2,3,4	0.041323	20.42
c,1,4	0.107798	29.73	c,1,2,3,4	0.049587	0.47

- ▶ The median probability model is  $\{c, x_1, x_2, x_4\}$  which clearly coincides with the optimal predictive model.
- ▶ Note that the risk of the maximum probability model  $\{c, x_1, x_2\}$  is considerably higher than that of the median probability model.

## R package BayesVarSel (Garcia-Donato and Forte)

- ▶ freely available at CRAN
- ▶ a number of different priors on models and parameters
- ▶ in the output: posterior inclusion probabilities, the highest probability model, the median probability model, the Bayesian model average predictor of  $y^*$  at covariates  $x^*$

## When the median probability model can fail (Merlise Clyde)

Suppose that

- ▶ the only models under consideration are  $M_0$  with a constant term and the models  $M_i$  with the constant term and a single covariate  $x_i$ ,  $i = 1, \dots, k$ , with  $k \geq 3$ ;
- ▶ all models have equal prior probability of  $1/(k + 1)$ ;
- ▶ all covariates are nearly perfectly correlated, with each other and with  $y$ .

Then the posterior probability of  $M_0$  will be near zero, and that of each of the  $M_i$  will be approximately  $1/k$ .

Since these posterior inclusion probabilities are less than  $1/2$ , the median probability model will be  $M_0$ , which will have very poor predictive performance compared to any of the other models.

... more on this next time!

Thank you for your attention!