

The case of linear regression models

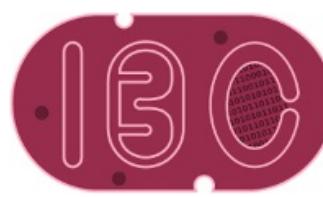
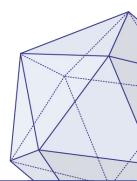
Jean-Michel Marin

Université de Montpellier

Institut Montpelliérain Alexander Grothendieck

Institut de Biologie Computationnelle (IBC)

Labex Numev



11ème école d'été de Peresq en traitement du signal et des images



Joint work with Christian Robert

Plan

The model

Natural conjugate prior family

Zellner's G -prior

Bayes factor

Prediction

The model

The dataset is then made up of the reunion of the vector of outcomes

$$\mathbf{y} = (y_1, \dots, y_n)$$

and the $n \times p$ matrix of explanatory variables

$$\mathbf{X} = [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_p] = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}.$$

The ordinary Gaussian linear regression model is such that:

$$\mathbf{y} | \alpha, \boldsymbol{\beta}, \sigma^2 \sim \mathcal{N}_n \left(\alpha \mathbf{1}_n + \mathbf{X} \boldsymbol{\beta}, \sigma^2 \mathbf{I}_n \right),$$

y_i 's are independent normal random variables with

$$\mathbb{E}[y_i | \alpha, \boldsymbol{\beta}, \sigma^2] = \alpha + \beta_1 x_{i1} + \dots + \beta_p x_{ip}, \quad \mathbb{V}[y_i | \alpha, \boldsymbol{\beta}, \sigma^2] = \sigma^2.$$

Given that the models studied in this section are all conditional on the regressors, we omit the conditioning on \mathbf{X} to simplify the notations.

A 1973 study on pine processionary caterpillars: y is the logarithmic transform of the average number of nests of caterpillars per tree in an area of 500 square meters, $n = 33$ and $p = 8$:

x_1 is the altitude (in meters),

x_2 is the slope (in degrees),

x_3 is the number of pine trees in the area,

x_4 is the height (in meters) of the tree sampled at the center of the area,

x_5 is the orientation of the area (from 1 if southbound to 2 otherwise),

x_6 is the height (in meters) of the dominant tree,

x_7 is the number of vegetation strata,

x_8 is the mix settlement index (from 1 if not mixed to 2 if mixed).

The goal of the regression analysis is to decide which explanatory variables have a strong influence on the number of nests.

We assume that $\text{rank} [\mathbf{1}_n \quad \mathbf{X}] = p + 1$.

$$\ell(\alpha, \boldsymbol{\beta}, \sigma^2 | \mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\beta}) \right\}.$$

Natural conjugate prior family

$$(\alpha, \beta) | \sigma^2 \sim \mathcal{N}_{p+1}((\tilde{\alpha}, \tilde{\beta}), \sigma^2 M^{-1}),$$

conditional on σ^2 and

$$\sigma^2 \sim \mathcal{IG}(a, b).$$

Even in the presence of genuine information on the parameters, the hyperparameters M , a and b are very difficult to specify and the posterior distributions.

Zellner's G -prior

$$\boldsymbol{\beta} | \alpha, \sigma^2 \sim \mathcal{N}_p \left(\tilde{\boldsymbol{\beta}}, g\sigma^2(\mathbf{X}^T \mathbf{X})^{-1} \right),$$

and a noninformative prior distribution is imposed on the pair (α, σ^2) ,

$$\pi(\alpha, \sigma^2) \propto \sigma^{-2}.$$

The factor g can be interpreted as being inversely proportional to the amount of information available in the prior relative to the sample.

For instance, setting $g = n$ gives the prior the same weight as one observation of the sample.

We will use this as our default value.

When $p > 0$,

$$\begin{aligned}\alpha | \sigma^2, \mathbf{y} &\sim \mathcal{N}_1 (\bar{\mathbf{y}}, \sigma^2/n) , \\ \boldsymbol{\beta} | \mathbf{y}, \sigma^2 &\sim \mathcal{N}_p \left(\frac{g}{g+1} \left(\hat{\boldsymbol{\beta}} + \tilde{\boldsymbol{\beta}}/g \right), \frac{\sigma^2 g}{g+1} \{ \mathbf{X}^T \mathbf{X} \}^{-1} \right) ,\end{aligned}$$

where $\hat{\boldsymbol{\beta}} = \{ \mathbf{X}^T \mathbf{X} \}^{-1} \mathbf{X}^T \mathbf{y}$ is the maximum likelihood and least squares estimator of $\boldsymbol{\beta}$.

The posterior independence between α and $\boldsymbol{\beta}$ is due to the fact that \mathbf{X} is centered and that α and $\boldsymbol{\beta}$ are a priori independent.

$$\sigma^2 | \mathbf{y} \sim I\mathcal{G} \left[(n-1)/2, s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) / (g+1) \right]$$

$$\text{where } s^2 = (\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n - \mathbf{X}\hat{\boldsymbol{\beta}})$$

When $p = 0$,

$$\alpha | \mathbf{y}, \sigma^2 \sim \mathcal{N}(\bar{\mathbf{y}}, \sigma^2/n) ,$$

$$\sigma^2 | \mathbf{y} \sim I\mathcal{G} \left[(n-1)/2, (\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n)^T (\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n) / 2 \right] .$$

We can derive from the previous derivations that

$$\mathbb{E}^\pi [\alpha | \mathbf{y}] = \mathbb{E}^\pi [\mathbb{E}^\pi (\alpha | \sigma^2, \mathbf{y}) | \mathbf{y}] = \mathbb{E}^\pi [\bar{\mathbf{y}} | \mathbf{y}] = \bar{\mathbf{y}}$$

$$\begin{aligned}\mathbb{E}^\pi [\boldsymbol{\beta} | \mathbf{y}] &= \mathbb{E}^\pi [\mathbb{E}^\pi (\boldsymbol{\beta} | \sigma^2, \mathbf{y}) | \mathbf{y}] \\ &= \mathbb{E}^\pi \left[\frac{g}{g+1} (\hat{\boldsymbol{\beta}} + \tilde{\boldsymbol{\beta}}/g) | \mathbf{y} \right] \\ &= \frac{g}{g+1} (\hat{\boldsymbol{\beta}} + \tilde{\boldsymbol{\beta}}/g).\end{aligned}$$

This result gives its meaning to the above point relating g with the amount of information contained in the dataset.

$$\mathbb{E}^\pi [\boldsymbol{\beta} | \mathbf{y}] = \frac{g}{g+1} (\hat{\boldsymbol{\beta}} + \tilde{\boldsymbol{\beta}}/g)$$

When $g = 1$, the prior information has the same weight as this amount: the Bayesian estimate of $\boldsymbol{\beta}$ is the average between the least square estimator and the prior expectation.

The larger g is, the weaker the prior information and the closer the Bayesian estimator is to the least squares estimator.

When considering the marginal likelihood at the core of the Bayes factors, we have, if $p \neq 0$,

$$f(\mathbf{y}) = \int \left(\int \int f(\mathbf{y}|\alpha, \boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}|\alpha, \sigma^2) \pi(\sigma^2, \alpha) d\alpha d\boldsymbol{\beta} \right) d\sigma^2,$$

$$f(\mathbf{y}) = \frac{\delta \Gamma((n-1)/2)}{\pi^{(n-1)/2} n^{1/2}} (g+1)^{-p/2} \kappa^{-(n-1)/2}.$$

$$\begin{aligned} \kappa &= (\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n)^T (\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n) + \frac{1}{g+1} \left\{ -g\mathbf{y}^T \mathbf{P}\mathbf{y} + \tilde{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{P}\mathbf{X} \tilde{\boldsymbol{\beta}} - 2\mathbf{y}^T \mathbf{P}\mathbf{X} \tilde{\boldsymbol{\beta}} \right\} \\ &= s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})/(g+1). \end{aligned}$$

If $p = 0$, a similar expression emerges:

$$f(\mathbf{y}) = \int \left(\int f(\mathbf{y}|\alpha, \sigma^2) \pi(\alpha, \sigma^2) d\alpha \right) d\sigma^2,$$

$$f(\mathbf{y}) = \frac{\delta\Gamma((n-1)/2)}{\pi^{(n-1)/2} n^{1/2}} \left[(\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n)^T (\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n) \right]^{-(n-1)/2}$$

as the evidence associated with this “null” model.

Bayes factor

As explained, in the previous section, the computation of Bayes factors is plagued by the inability to include generic improper prior distributions.

In order to bypass this difficulty, we will assume that all the linear models under comparison do include the parameter α , which means that each regression model includes an intercept.

This assumption allows us to take the *same* improper prior on (α, σ^2) for all of those models.

Otherwise, the Bayes factors simply cannot be correctly defined.

When we compare two sets of regressors, we have to handle two regression matrices, \mathbf{X}^1 and \mathbf{X}^2 , with respective dimensions (n, p_1) and (n, p_2) , extracted from the original matrix \mathbf{X} by removing some columns.

$$B_{12}(\mathbf{y}) = \frac{(g_1 + 1)^{-p_1/2} \left[s_1^2 + (\tilde{\boldsymbol{\beta}}^1 - \hat{\boldsymbol{\beta}}^1)^T (\mathbf{X}^1)^T \mathbf{X}^1 (\tilde{\boldsymbol{\beta}}^1 - \hat{\boldsymbol{\beta}}^1) / (g_1 + 1) \right]^{-(n-1)/2}}{(g_2 + 1)^{-p_2/2} \left[s_2^2 + (\tilde{\boldsymbol{\beta}}^2 - \hat{\boldsymbol{\beta}}^2)^T (\mathbf{X}^2)^T \mathbf{X}^2 (\tilde{\boldsymbol{\beta}}^2 - \hat{\boldsymbol{\beta}}^2) / (g_2 + 1) \right]^{-(n-1)/2}}.$$

For the caterpillar dataset, if we have to test the null hypothesis H_0 : $\beta_8 = \beta_9 = 0$, using $\tilde{\beta}^1 = 0_8$, $\tilde{\beta}^2 = 0_6$, and an arbitrary $g_1 = g_2 = 100$, in Zellner's G -priors.

We obtain $B_{12}^\pi = 0.0165$ when model \mathfrak{M}_2 corresponds to H_0 .

Using Jeffreys' scale of evidence, this implies that $\log_{12}(B_{12}^\pi) = -1.78$, hence that the posterior distribution appears to strongly favor H_0 .

More generally, we can produce a Bayesian regression output:

	PostMean	PostStError	Log10bf	EvidAgaH0
Intercept	-0.8133	0.1407		
x1	-0.5039	0.1883	0.7224	(**)
x2	-0.3755	0.1508	0.5392	(**)
x3	0.6225	0.3436	-0.0443	
x4	-0.2776	0.2804	-0.5422	
x5	-0.2069	0.1499	-0.3378	
x6	0.2806	0.4760	-0.6857	
x7	-1.0420	0.4178	0.5435	(**)
x8	-0.0221	0.1531	-0.7609	

Posterior Mean of Sigma2: 0.6528

Posterior StError of Sigma2: 0.939

Prediction

The prediction of $m \geq 1$ future observations from units for which the explanatory variables $\tilde{\mathbf{X}}$ —but not the outcome variable $\tilde{\mathbf{y}}$ —have been observed or set is also based on the posterior distribution.

Logically enough, were α , β and σ^2 known quantities, the m -vector $\tilde{\mathbf{y}}$ would then have a Gaussian distribution with mean $\alpha\mathbf{1}_m + \tilde{\mathbf{X}}\beta$ and variance $\sigma^2\mathbf{I}_m$.

The *predictive distribution* on $\tilde{\mathbf{y}}$ is defined as the marginal in \mathbf{y} of the joint posterior distribution on $(\tilde{\mathbf{y}}, \alpha, \beta, \sigma^2)$.

Conditional on σ^2 , the vector $\tilde{\mathbf{y}}$ of future observations has a Gaussian distribution and we can derive its expectation—used as our Bayesian estimator—by averaging over α and β ,

$$\begin{aligned}\mathbb{E}^\pi[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}] &= \mathbb{E}^\pi[\mathbb{E}^\pi(\tilde{\mathbf{y}}|\alpha, \beta, \sigma^2, \mathbf{y})|\sigma^2, \mathbf{y}] \\ &= \mathbb{E}^\pi[\alpha \mathbf{1}_m + \tilde{\mathbf{X}}\beta|\sigma^2, \mathbf{y}] \\ &= \hat{\alpha} \mathbf{1}_m + \tilde{\mathbf{X}} \frac{\tilde{\beta} + g\hat{\beta}}{g+1},\end{aligned}$$

which is independent from σ^2 . This representation is quite intuitive, being the product of the matrix of explanatory variables $\tilde{\mathbf{X}}$ by the Bayesian estimator of β .

Similarly, we can compute

$$\begin{aligned}
 \mathbb{V}^\pi(\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}) &= \mathbb{E}^\pi[\mathbb{V}^\pi(\tilde{\mathbf{y}}|\alpha, \boldsymbol{\beta}, \sigma^2, \mathbf{y})|\sigma^2, \mathbf{y}] \\
 &\quad + \mathbb{V}^\pi(\mathbb{E}^\pi(\tilde{\mathbf{y}}|\alpha, \boldsymbol{\beta}, \sigma^2)|\sigma^2, \mathbf{y}) \\
 &= \mathbb{E}^\pi[\sigma^2 I_m|\sigma^2, \mathbf{y}] + \mathbb{V}^\pi(\alpha \mathbf{1}_m + \tilde{\mathbf{X}} \boldsymbol{\beta}|\sigma^2, \mathbf{y}) \\
 &= \sigma^2 \left(I_m + \frac{g}{g+1} \tilde{\mathbf{X}} (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T \right).
 \end{aligned}$$

Due to this factorisation, and the fact that the conditional expectation does not depend on σ^2 , we thus obtain

$$\mathbb{V}^\pi(\tilde{\mathbf{y}}|\mathbf{y}) = \hat{\sigma}^2 \left(I_m + \frac{g}{g+1} \tilde{\mathbf{X}} (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T \right).$$

This decomposition of the variance makes perfect sense: Conditionally on σ^2 , the posterior predictive variance has two terms, the first term being $\sigma^2 I_m$, which corresponds to the sampling variation, and the second one being $\sigma^2 \frac{g}{g+1} \tilde{\mathbf{X}}(\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T$, which corresponds to the uncertainty about β .

HPD credible regions and tests can then be conducted based on this conditional predictive distribution

$$\tilde{\mathbf{y}} | \sigma^2, \mathbf{y}, \sigma^2 \sim \mathcal{N} (\mathbb{E}^\pi[\tilde{\mathbf{y}}], \mathbb{V}^\pi(\tilde{\mathbf{y}} | \mathbf{y}, \sigma^2)) .$$

Integrating σ^2 out to produce the marginal distribution of $\tilde{\mathbf{y}}$ leads to a multivariate Student's t distribution

$$\begin{aligned}\tilde{\mathbf{y}}|\mathbf{y} &\sim \mathcal{T}_m \left(n, \hat{\alpha} \mathbf{1}_m + g\tilde{\boldsymbol{\beta}}/(g+1), \right. \\ &\quad \left. \frac{s^2 + \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}}}{n} \left\{ \mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T \right\} \right).\end{aligned}$$