

cvLME

A multi-language library to perform
cross-validated Bayesian model selection

Joram Soch, BCCN Berlin

Toolbox URL: <https://github.com/JoramSoch/cvLME>
Support Contact: Joram Soch <joram.soch@bccn-berlin.de>
Current Version: cvLME V0.3.4 [V0.languages.models]
Related Paper: Soch J, Allefeld C (2018). MACS – a new SPM toolbox for model assessment, comparison and selection. *Journal of Neuroscience Methods*, vol. 306, pp. 19-31; DOI: 10.1016/j.jneumeth.2018.05.017.

Contents

0	General remarks	1
0.1	Mathematics	1
0.2	Implementation	2
1	Model spaces and model selection	3
1.1	Log model evidence	3
1.2	Log Bayes factor	3
1.3	Posterior probabilities	3
1.4	Log family evidence	4
1.5	Implementation	5
2	Univariate General Linear Model	6
2.1	Likelihood function	6
2.2	Maximum likelihood	6
2.3	Prior distribution	7
2.4	Joint likelihood	7
2.5	Posterior distribution	8
2.6	Log model evidence	8
2.7	Cross-validated LME	10
2.8	Special cases	10
2.9	Implementation	10
3	Multivariate General Linear Model	11
3.1	Likelihood function	11
3.2	Maximum likelihood	11
3.3	Prior distribution	12
3.4	Joint likelihood	12
3.5	Posterior distribution	13
3.6	Log model evidence	14
3.7	Cross-validated LME	15
3.8	Special cases	15
3.9	Implementation	16
4	Poisson Distribution with Exposures	17
4.1	Likelihood function	17
4.2	Maximum likelihood	17
4.3	Prior distribution	17
4.4	Joint likelihood	17
4.5	Posterior distribution	18
4.6	Log model evidence	18
4.7	Cross-validated LME	19
4.8	Special cases	19
4.9	Implementation	20

0 General remarks

0.1 Mathematics

In the following sections, we are considering different model classes – such as linear regression for continuous data (see Section 2) or a Poisson model for count data (see Section 4). For each of these model classes – except for the section on model comparison where only some measures of model selection are introduced (see Section 1) –, we are going through several steps of mathematical derivation which are outlined here.

Step 1: First, the likelihood function is specified,

$$p(y|\theta, m) , \quad (0.1)$$

i.e. the probability of observing the data y , given a model m and parameters θ .

Step 2: Second, the maximum likelihood estimates are derived,

$$\hat{\theta} = \arg \max_{\theta} \log p(y|\theta, m) , \quad (0.2)$$

i.e. those model parameters that maximize the log-likelihood function $LL(\theta)$.

Step 3: Third, a (conjugate) prior distribution is specified,

$$p(\theta|m) , \quad (0.3)$$

i.e. a distribution over parameters that can be applied to the likelihood $p(y|\theta, m)$.

Step 4: Then, the joint likelihood function is calculated,

$$p(y, \theta|m) = p(y|\theta) p(\theta|m) , \quad (0.4)$$

i.e. the product of likelihood function and prior distribution over model parameters.

Step 5: Then, the posterior distribution is obtained,

$$p(\theta|y, m) = \frac{p(y|\theta, m) p(\theta|m)}{p(y|m)} \propto p(y, \theta|m) , \quad (0.5)$$

i.e. a distribution over model parameters that is proportional to the joint likelihood.

Step 6: Then, the log model evidence is derived,

$$\text{LME}(m) = \log \int p(y|\theta) p(\theta|m) d\theta = \log \frac{p(y|\theta, m) p(\theta|m)}{p(\theta|y, m)} , \quad (0.6)$$

i.e. the logarithm of the marginal likelihood, expected over the prior distribution.

Step 7: Then, the cross-validated log model evidence is presented,

$$\text{cvLME}(m) = \sum_{i=1}^S \log \int p(y_i|\theta, m) p(\theta| \cup_{j \neq i} y_j, m) d\theta \quad (0.7)$$

which basically consists in describing what kind of non-informative prior is used for obtaining the informative posterior from the training data in each cross-validation fold.

Steps 8 and 9: Finally, special cases of each model class are considered and implementation in the different programming languages is described.

0.2 Implementation

The cross-validated log model evidence (cvLME; see eq. 0.7) is implemented within several programming languages and the respective implementational details are described here. In what follows, `<name-of-the-model-class>` is either “GLM” or “Pois” (but more options are to follow in the future).

MATLAB: In MATLAB, the different methods for each model class (maximum likelihood estimation, Bayesian estimation, log model evidence, cross-validated log model evidence) are implemented as different functions called “`<name-of-the-model-class>_MLE/Bayes/LME/cvLME.m`” and these functions can be directly called with variables representing quantities in the model, e.g. measured data or experimental design information. Often, some quantities can be left empty.

Python: In Python, all methods for all model classes are implemented as a single module called “`cvBMS.py`” which can be simply imported at the beginning of a script using “`import cvBMS`”. Then, a particular model is initialized by calling “`model = cvBMS.<name-of-the-model-class>`” with variables representing model quantities. When a model has been initialized, statistical operations can be performed in object-oriented fashion by typing “`model.MLE/Bayes/LME/cvLME()`” where sometimes, some more input variables can or must be provided.

Further details are provided in the implementation subsections of each model class individually (see Sections X.9).

1 Model spaces and model selection

1.1 Log model evidence

A model space is defined as a set of models. In the context of these tools, a model space is always initialized with a set of *log model evidences* (LME)

$$\text{LME}(m) = \log p(y|m) = \log \int p(y|\theta, m) p(\theta|m) d\theta \quad (1.1)$$

or *cross-validated log model evidences* (cvLMEs)

$$\text{cvLME}(m) = \sum_{i=1}^S \log \int p(y_i|\theta, m) p(\theta | \cup_{j \neq i} y_j, m) d\theta \quad (1.2)$$

where S is the number of data subsets.

1.2 Log Bayes factor

The *Bayes factor* (BF) is defined as the ratio of two model evidences,

$$\text{BF}_{12} = \frac{p(y|m_1)}{p(y|m_2)}, \quad (1.3)$$

such that the *log Bayes factor* (LBF) is the difference of two log model evidences,

$$\text{LBF}_{12} = \log \text{BF}_{12} = \log \frac{p(y|m_1)}{p(y|m_2)} = \text{LME}(m_1) - \text{LME}(m_2). \quad (1.4)$$

1.3 Posterior probabilities

Given more than two models, one can also calculate *posterior model probabilities* (PPs) by simply applying Bayes' theorem to the model evidences

$$p(m_i|y) = \frac{p(y|m_i) p(m_i)}{\sum_{j=1}^M p(y|m_j) p(m_j)} \quad (1.5)$$

or, equivalently, to the exponentiated log model evidences (LME)

$$p(m_i|y) = \frac{\exp[\text{LME}(m_i)] p(m_i)}{\sum_{j=1}^M \exp[\text{LME}(m_j)] p(m_j)} \quad (1.6)$$

where $p(m_i)$ are prior model probabilities and M is the number of models.

Note that posterior probabilities do not on depend on absolute LME values, but only on relative LME difference. For this reason, the mean LME over models is subtracted from all LMEs before PPs are calculated.

1.4 Log family evidence

The *family evidence* (FE) is obtained by marginalizing over “model” within “family”, i.e. as the marginal probability over the model evidences from all models within one family

$$p(y|f) = \sum_{m \in f} p(y|m) p(m|f) \quad (1.7)$$

and the *log family evidence* (LFE) is the natural logarithm of this quantity

$$\text{LFE}(f) = \log p(y|f) = \log \sum_{m \in f} p(y|m) p(m|f) \quad (1.8)$$

where $p(m|f)$ is a (most likely uniform) within-family prior distribution.

Note that, with a uniform within-family prior, the family evidence is the average of model evidences, but the log family evidence is not the average of the log model evidences! In particular, the problem is that we usually cannot access model evidences $p(y|m)$ directly, but only deal with log model evidences $\log p(y|m)$. LMEs are used to avoid computational problems with very small model evidences that could not be stored in standard computers, e.g. $p(y|m) = 10^{-100} \Rightarrow \log p(y|m) \approx -230$. However, just exponentiating LMEs does not work, because they often fall below a specific underflow threshold $-u$, e.g. $u = 745$, so that all model evidences would be 0.

The solution is to select the maximum LME within a family

$$L^*(f) = \max_{m \in f} [\text{LME}(m)] \quad (1.9)$$

and define differences between LMEs and maximum LME as

$$L'(m) = \text{LME}(m) - L^*(f) . \quad (1.10)$$

Then, the log family evidence can be written as

$$\text{LFE}(f) = \log p(y|f) = \log \left[\frac{1}{M_f} \sum_{i=1}^{M_f} \exp [\text{LME}(m_i)] \right] \quad (1.11)$$

which can be further developed in the following way:

$$\begin{aligned} \text{LFE}(f) &= \log \left[\frac{1}{M_f} \sum_{i=1}^{M_f} \exp [L'(m_i) + L^*(f)] \right] \\ &= \log \left[\frac{1}{M_f} \exp L^*(f) \sum_{i=1}^{M_f} \exp L'(m_i) \right] \\ &= L^*(f) + \log \sum_{i=1}^{M_f} \exp L'(m_i) - \log M_f . \end{aligned} \quad (1.12)$$

1.5 Implementation

In **MATLAB**, (log) Bayes factors, posterior model probabilities and log family evidences are implemented via the functions `MS_LBF`, `MS_PP` and `MS_LFE` which have to be called with an $M \times N$ matrix `LME` as input.

In **Python**, a model space object has to be initiated via `ms = cvBMS.MS(LME)` and (log) Bayes factors, posterior model probabilities and log family evidences are calculated via `ms.LBF`, `ms.BF`, `ms.PP`, and `ms.LFE`.

2 Univariate General Linear Model

2.1 Likelihood function

In the univariate general linear model (GLM), a single measured signal (y) is modelled as a linear combination (β) of predictor variables (X), where errors (ε) are assumed to be normally distributed around zero and to have a known covariance structure (V), but unknown variance factor (σ^2):

$$y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 V) . \quad (2.1)$$

In this equation, y is the $n \times 1$ measured signal, X is the $n \times p$ design matrix, β is a $p \times 1$ vector of regression coefficients, ε is an $n \times 1$ vector of errors, σ^2 is the variance of these errors and V is an $n \times n$ correlation matrix where n is the number of data points and p is the number of regressors.

The GLM equation (2.1) implies the following *likelihood function*

$$p(y|\beta, \sigma^2) = N(y; X\beta, \sigma^2 V) = \sqrt{\frac{1}{(2\pi)^n |\sigma^2 V|}} \exp \left[-\frac{1}{2\sigma^2} (y - X\beta)^T V^{-1} (y - X\beta) \right] \quad (2.2)$$

which, for mathematical convenience, can also be parametrized as

$$p(y|\beta, \tau) = N(y; X\beta, (\tau P)^{-1}) = \sqrt{\frac{|\tau P|}{(2\pi)^n}} \exp \left[-\frac{\tau}{2} (y - X\beta)^T P (y - X\beta) \right] \quad (2.3)$$

using the residual precision $\tau = 1/\sigma^2$ and the $n \times n$ precision matrix $P = V^{-1}$.

2.2 Maximum likelihood

Classical model estimation proceeds by maximizing the *log-likelihood* (LL)

$$LL(\beta, \sigma^2) = \log p(y|\beta, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\sigma^2 V| - \frac{1}{2\sigma^2} (y - X\beta)^T V^{-1} (y - X\beta) \quad (2.4)$$

which gives rise to *maximum-likelihood* (ML) parameter estimates

$$\begin{aligned} \hat{\beta} &= (X^T V^{-1} X)^{-1} X^T V^{-1} y \\ \hat{\sigma}^2 &= \frac{1}{n} (y - X\hat{\beta})^T V^{-1} (y - X\hat{\beta}) \end{aligned} \quad (2.5)$$

that can be used to form t - and F -statistics

$$\begin{aligned} t &= \frac{c^T \hat{\beta}}{\sqrt{\hat{\sigma}^2 c^T \text{cov}(\hat{\beta}) c}} \\ F &= (C^T \hat{\beta})^T (\hat{\sigma}^2 C^T \text{cov}(\hat{\beta}) C)^{-1} (C^T \hat{\beta}) \end{aligned} \quad (2.6)$$

where c is a $p \times 1$ *contrast vector*, C is a $p \times q$ *contrast matrix* and

$$\text{cov}(\hat{\beta}) = (X^T V^{-1} X)^{-1} . \quad (2.7)$$

2.3 Prior distribution

A conjugate prior distribution relative to the likelihood function given by (2.3) is the *normal-gamma distribution* over regression coefficients β and residual precision τ

$$p(\beta, \tau) = N(\beta; \mu_0, (\tau \Lambda_0)^{-1}) \cdot \text{Gam}(\tau; a_0, b_0) \quad (2.8)$$

which can be split into a conditional distribution and a marginal distribution

$$\begin{aligned} p(\beta|\tau) &= N(\beta; \mu_0, (\tau \Lambda_0)^{-1}) = \sqrt{\frac{|\tau \Lambda_0|}{(2\pi)^p}} \exp \left[-\frac{\tau}{2} (\beta - \mu_0)^T \Lambda_0 (\beta - \mu_0) \right] \\ p(\tau) &= \text{Gam}(\tau; a_0, b_0) = \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \end{aligned} \quad (2.9)$$

where μ_0 and Λ_0 are the prior mean and the prior precision of β and a_0 and b_0 are the prior shape and rate parameters for τ .

2.4 Joint likelihood

Combining the likelihood function (2.3) with the prior distribution (2.9), the *joint likelihood function* of the general linear model with normal-gamma priors (GLM-NG) becomes

$$\begin{aligned} p(y, \beta, \tau) &= p(y|\beta, \tau) p(\beta, \tau) = p(y|\beta, \tau) p(\beta|\tau) p(\tau) \\ &= \sqrt{\frac{|\tau P|}{(2\pi)^n}} \exp \left[-\frac{\tau}{2} (y - X\beta)^T P (y - X\beta) \right] \cdot \\ &\quad \sqrt{\frac{|\tau \Lambda_0|}{(2\pi)^p}} \exp \left[-\frac{\tau}{2} (\beta - \mu_0)^T \Lambda_0 (\beta - \mu_0) \right] \cdot \\ &\quad \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] . \end{aligned} \quad (2.10)$$

Collecting identical variables gives:

$$\begin{aligned} p(y, \beta, \tau) &= \sqrt{\frac{\tau^{n+p}}{(2\pi)^{n+p}} |P| |\Lambda_0|} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \\ &\quad \exp \left[-\frac{\tau}{2} ((y - X\beta)^T P (y - X\beta) + (\beta - \mu_0)^T \Lambda_0 (\beta - \mu_0)) \right] . \end{aligned} \quad (2.11)$$

Completing the square over β gives:

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^{n+p}}{(2\pi)^{n+p}}} |P| |\Lambda_0| \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} ((\beta - \mu_n)^T \Lambda_n (\beta - \mu_n) + (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n)) \right] . \quad (2.12)$$

2.5 Posterior distribution

The *posterior distribution* in the GLM-NG can be evaluated using Bayes' theorem:

$$p(\beta, \tau | y) = \frac{p(y | \beta, \tau) p(\beta, \tau)}{p(y)} . \quad (2.13)$$

Since $p(y)$ is just a normalization factor, the posterior is proportional to the joint:

$$p(\beta, \tau | y) \propto p(y | \beta, \tau) p(\beta, \tau) = p(y, \beta, \tau) . \quad (2.14)$$

From the term in (2.12), we can isolate the posterior distribution over β :

$$\begin{aligned} p(\beta | \tau, y) &= N(\beta; \mu_n, (\tau \Lambda_n)^{-1}) \\ \mu_n &= \Lambda_n^{-1} (X^T P y + \Lambda_0 \mu_0) \\ \Lambda_n &= X^T P X + \Lambda_0 . \end{aligned} \quad (2.15)$$

From the remaining term, we can isolate the posterior distribution over τ :

$$\begin{aligned} p(\tau | y) &= \text{Gam}(\tau; a_n, b_n) \\ a_n &= a_0 + \frac{n}{2} \\ b_n &= b_0 + \frac{1}{2} (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) . \end{aligned} \quad (2.16)$$

2.6 Log model evidence

According to the law of marginal probability, the *model evidence* of the GLM-NG is:

$$p(y | m) = \iint p(y | \beta, \tau) p(\beta | \tau) p(\tau) d\beta d\tau . \quad (2.17)$$

According to the law of conditional probability, the integrand is equivalent to the joint:

$$p(y | m) = \iint p(y, \beta, \tau) d\beta d\tau . \quad (2.18)$$

In (2.12), we have already evaluated this term as

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^n |P|}{(2\pi)^n}} \sqrt{\frac{\tau^p |\Lambda_0|}{(2\pi)^p}} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} ((\beta - \mu_n)^T \Lambda_n (\beta - \mu_n) + (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n)) \right]. \quad (2.19)$$

Using the posterior distribution over β , we can rewrite this as

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^n |P|}{(2\pi)^n}} \sqrt{\frac{\tau^p |\Lambda_0|}{(2\pi)^p}} \sqrt{\frac{(2\pi)^p}{\tau^p |\Lambda_n|}} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \mathcal{N}(\beta; \mu_n, (\tau \Lambda_n)^{-1}) \exp \left[-\frac{\tau}{2} (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) \right]. \quad (2.20)$$

Now, β can be integrated out easily:

$$\int p(y, \beta, \tau) d\beta = \sqrt{\frac{\tau^n |P|}{(2\pi)^n}} \sqrt{\frac{|\Lambda_0|}{|\Lambda_n|}} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) \right]. \quad (2.21)$$

Using the posterior distribution over τ , we can rewrite this as

$$\int p(y, \beta, \tau) d\beta = \sqrt{\frac{|P|}{(2\pi)^n}} \sqrt{\frac{|\Lambda_0|}{|\Lambda_n|}} \frac{b_0^{a_0}}{\Gamma(a_0)} \frac{\Gamma(a_n)}{b_n^{a_n}} \text{Gam}(\tau; a_n, b_n). \quad (2.22)$$

Finally, τ can also be integrated out:

$$\iint p(y, \beta, \tau) d\beta d\tau = \sqrt{\frac{|P|}{(2\pi)^n}} \sqrt{\frac{|\Lambda_0|}{|\Lambda_n|}} \frac{\Gamma(a_n)}{\Gamma(a_0)} \frac{b_0^{a_0}}{b_n^{a_n}} = p(y|m). \quad (2.23)$$

Thus, the *log model evidence* of the GLM-NG is given by

$$\log p(y|m) = \frac{1}{2} \log |P| - \frac{n}{2} \log(2\pi) + \frac{1}{2} \log |\Lambda_0| - \frac{1}{2} \log |\Lambda_n| + \log \Gamma(a_n) - \log \Gamma(a_0) + a_0 \log b_0 - a_n \log b_n. \quad (2.24)$$

2.7 Cross-validated LME

For calculation of the *cross-validated log model evidence* (cvLME), the data are splitted into S subsets. In the training phase, all except one subset of the data are analyzed using a non-informative prior $p_{\text{ni}}(\beta, \tau)$ with the prior parameters

$$\mu_0 = 0_p, \Lambda_0 = 0_{pp} \quad \text{and} \quad a_0 = 0, b_0 = 0 \quad (2.25)$$

to obtain an informative posterior $p(\beta, \tau | \cup_{j \neq i} y_j)$ using equations (2.15) and (2.16). In the testing phase, this informative posterior is then applied as a prior distribution to obtain the out-of-sample log model evidence $\log p(y_i | \cup_{j \neq i} y_j)$ via equation (2.24). Summing up over data subsets yields the cvLME according to equation (1.2).

As one can see from equations (2.15) and (2.16), the priors in (2.25) are non-informative in the sense that only the data remain to influence the posteriors.

2.8 Special cases

The *univariate Gaussian with unknown variance* (UGuv) is a special case in which

$$X = 1_n, \quad \beta = \mu \quad \text{and} \quad V = I_n. \quad (2.26)$$

Furthermore, *simple linear regression* (SLR) is a special case of the GLM where

$$X = [1_n, x], \quad \beta = [\beta_0, \beta_1]^T \quad \text{and} \quad V = I_n. \quad (2.27)$$

The *one-sample t-test*, the *two-sample t-test*, the *paired t-test* and the *omnibus F-test* can all be emulated as comparisons of general linear models with specific design matrices.

2.9 Implementation

In **MATLAB**, maximum likelihood estimates and Bayesian posterior distributions can be obtained via the functions `GLM_MLE` and `GLM_Bayes` while log model evidence and cross-validated LME can be calculated using the functions `GLM_LME` and `GLM_cvLME`. Given an $n \times v$ data matrix Y , an $n \times p$ design matrix X , an $n \times n$ precision matrix P and a number of data subsets S , the cvLME for a GLM-NG is calculated as

$$\text{cvLME} = \text{GLM_cvLME}(Y, X, P, S); \quad (2.28)$$

In **Python**, a GLM object has to be initiated via `glm = cvBMS.GLM(Y, X, V)` and maximum likelihood estimates, Bayesian posterior distributions, log model evidence and cross-validated LME are evaluated via `glm.MLE`, `glm.Bayes`, `glm.LME`, and `glm.cvLME`. Given Y , X , V and S as above, the cvLME for a GLM-NG is calculated as

$$\text{cvLME} = \text{cvBMS.GLM}(Y, X, V).\text{cvLME}(S) \quad (2.29)$$

In all of the above, V and P default to I_n whereas S defaults to 2 when left empty.

3 Multivariate General Linear Model

3.1 Likelihood function

In the multivariate general linear model (MGLM), several measured signals (Y) are modelled as a linear combination (B) of predictor variables (X), where errors (E) are assumed to be normally distributed around zero and to have a known covariance across observations (V), but unknown covariance across measurements (E):

$$Y = XB + E, \quad E \sim \mathcal{MN}(0, V, \Sigma) . \quad (3.1)$$

In this equation, Y is the $n \times v$ data matrix, X is the $n \times p$ design matrix, B is a $p \times v$ matrix of regression coefficients, E is an $n \times v$ matrix of errors, V is an $n \times n$ correlation matrix and Σ is a $v \times v$ covariance matrix where n is the number of data points, v is the number of measured variables and p is the number of regressors.

The MGLM equation (3.1) implies the following *likelihood function*

$$\begin{aligned} p(Y|B, \Sigma) &= \mathcal{MN}(Y; XB, V, \Sigma) \\ &= \sqrt{\frac{1}{(2\pi)^{nv} |\Sigma|^n |V|^v}} \exp \left[-\frac{1}{2} \text{tr} (\Sigma^{-1} (Y - XB)^T V^{-1} (Y - XB)) \right] \end{aligned} \quad (3.2)$$

which, for mathematical convenience, can also be parametrized as

$$\begin{aligned} p(Y|B, T) &= \mathcal{MN}(Y; XB, P, T^{-1}) \\ &= \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \exp \left[-\frac{1}{2} \text{tr} (T (Y - XB)^T P (Y - XB)) \right] \end{aligned} \quad (3.3)$$

using the $v \times v$ precision matrix $T = \Sigma^{-1}$ and the $n \times n$ precision matrix $P = V^{-1}$.

3.2 Maximum likelihood

Classical model estimation proceeds by maximizing the *log-likelihood* (LL)

$$\begin{aligned} \text{LL}(B, \Sigma) &= \log p(Y|B, \Sigma) \\ &= -\frac{nv}{2} \log(2\pi) - \frac{n}{2} \log(|\Sigma|) - \frac{v}{2} \log(|V|) \\ &\quad - \frac{1}{2} \text{tr} [\Sigma^{-1} (Y - XB)^T V^{-1} (Y - XB)] . \end{aligned} \quad (3.4)$$

which gives rise to *maximum-likelihood* (ML) parameter estimates

$$\begin{aligned} \hat{B} &= (X^T V^{-1} X)^{-1} X^T V^{-1} Y \\ \hat{\Sigma} &= \frac{1}{n} (Y - X \hat{B})^T V^{-1} (Y - X \hat{B}) . \end{aligned} \quad (3.5)$$

3.3 Prior distribution

A conjugate prior distribution relative to the likelihood function given by (3.3) is the *normal-Wishart distribution* over regression coefficients B and noise precision T

$$p(B, T) = \mathcal{MN}(B; M_0, \Lambda_0^{-1}, T^{-1}) \cdot \mathcal{W}(T; \Omega_0^{-1}, \nu_0) \quad (3.6)$$

which can be split into a conditional distribution and a marginal distribution

$$\begin{aligned} p(B|T) &= \sqrt{\frac{|T|^p |\Lambda_0|^v}{(2\pi)^{pv}}} \exp \left[-\frac{1}{2} \text{tr} (T(B - M_0)^T \Lambda_0 (B - M_0)) \right] \\ p(T) &= \frac{1}{\Gamma_v \left(\frac{\nu_0}{2} \right)} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}}} |T|^{(\nu_0 - v - 1)/2} \exp \left[-\frac{1}{2} \text{tr} (\Omega_0 T) \right] \end{aligned} \quad (3.7)$$

where M_0 and Λ_0 are the prior mean and the prior precision of B and Ω_0 and ν_0 are the prior inverse scale matrix and degrees of freedom for T .

3.4 Joint likelihood

Combining the likelihood function (3.3) with the prior distribution (3.7), the *joint likelihood function* of the multivariate general linear model with normal-Wishart priors (MGLM-NW) becomes

$$\begin{aligned} p(Y, B, T) &= p(Y|B, T) p(B, T) = p(Y|B, T) p(B|T) p(T) \\ &= \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \exp \left[-\frac{1}{2} \text{tr} (T(Y - XB)^T P (Y - XB)) \right] \cdot \\ &\quad \sqrt{\frac{|T|^p |\Lambda_0|^v}{(2\pi)^{pv}}} \exp \left[-\frac{1}{2} \text{tr} (T(B - M_0)^T \Lambda_0 (B - M_0)) \right] \cdot \\ &\quad \frac{1}{\Gamma_v \left(\frac{\nu_0}{2} \right)} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}}} |T|^{(\nu_0 - v - 1)/2} \exp \left[-\frac{1}{2} \text{tr} (\Omega_0 T) \right] . \end{aligned} \quad (3.8)$$

Collecting identical variables gives:

$$\begin{aligned} p(Y, B, T) &= \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|T|^p |\Lambda_0|^v}{(2\pi)^{pv}}} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}}} \frac{1}{\Gamma_v \left(\frac{\nu_0}{2} \right)} \cdot |T|^{(\nu_0 - v - 1)/2} \exp \left[-\frac{1}{2} \text{tr} (\Omega_0 T) \right] \cdot \\ &\quad \exp \left[-\frac{1}{2} \text{tr} (T [(Y - XB)^T P (Y - XB) + (B - M_0)^T \Lambda_0 (B - M_0)]) \right] . \end{aligned} \quad (3.9)$$

Expanding the products in the exponent gives:

$$\begin{aligned}
p(Y, B, T) = & \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|T|^p |\Lambda_0|^v}{(2\pi)^{pv}}} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}} \frac{1}{\Gamma_v\left(\frac{\nu_0}{2}\right)}} \cdot |T|^{(\nu_0 - v - 1)/2} \exp \left[-\frac{1}{2} \text{tr}(\Omega_0 T) \right] \cdot \\
& \exp \left[-\frac{1}{2} \text{tr} \left(T \left[Y^T P Y - Y^T P X B - B^T X^T P Y + B^T X^T P X B + \right. \right. \right. \\
& \quad \left. \left. \left. B^T \Lambda_0 B - B^T \Lambda_0 M_0 - M_0^T \Lambda_0 B + M_0^T \Lambda_0 \mu_0 \right] \right) \right] .
\end{aligned} \tag{3.10}$$

Completing the square over B gives:

$$\begin{aligned}
p(Y, B, T) = & \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|T|^p |\Lambda_0|^v}{(2\pi)^{pv}}} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}} \frac{1}{\Gamma_v\left(\frac{\nu_0}{2}\right)}} \cdot |T|^{(\nu_0 - v - 1)/2} \exp \left[-\frac{1}{2} \text{tr}(\Omega_0 T) \right] \cdot \\
& \exp \left[-\frac{1}{2} \text{tr} \left(T \left[(B - M_n)^T \Lambda_n (B - M_n) + (Y^T P Y + M_0^T \Lambda_0 M_0 - M_n^T \Lambda_n M_n) \right] \right) \right] .
\end{aligned} \tag{3.11}$$

3.5 Posterior distribution

The *posterior distribution* in the MGLM-NW can be evaluated using Bayes' theorem:

$$p(B, T|Y) = \frac{p(Y|B, T) p(B, T)}{p(Y)} . \tag{3.12}$$

Since $p(Y)$ is just a normalization factor, the posterior is proportional to the joint:

$$p(B, T|Y) \propto p(Y|B, T) p(B, T) = p(Y, B, T) . \tag{3.13}$$

From the term in (3.11), we can isolate the posterior distribution over B :

$$\begin{aligned}
p(B|T, Y) = & \mathcal{MN}(B; M_n, \Lambda_n^{-1}, T^{-1}) \\
M_n = & \Lambda_n^{-1} (X^T P Y + \Lambda_0 M_0) \\
\Lambda_n = & X^T P X + \Lambda_0 .
\end{aligned} \tag{3.14}$$

From the remaining term, we can isolate the posterior distribution over T :

$$\begin{aligned}
p(T|Y) = & \mathcal{W}(T; \Omega_n^{-1}, \nu_n) \\
\Omega_n = & \Omega_0 + Y^T P Y + M_0^T \Lambda_0 M_0 - M_n^T \Lambda_n M_n \\
\nu_n = & \nu_0 + n .
\end{aligned} \tag{3.15}$$

3.6 Log model evidence

According to the law of marginal probability, the *model evidence* of the MGLM-NW is:

$$p(Y|m) = \iint p(Y|B, T) p(B, T) dB dT . \quad (3.16)$$

According to the law of conditional probability, the integrand is equivalent to the joint:

$$p(Y|m) = \iint p(Y, B, T) dB dT . \quad (3.17)$$

In (3.11), we have already evaluated this term as

$$p(Y, B, T) = \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|T|^p |\Lambda_0|^v}{(2\pi)^{pv}}} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}} \frac{1}{\Gamma_v\left(\frac{\nu_0}{2}\right)}} \cdot |T|^{(\nu_0 - v - 1)/2} \exp\left[-\frac{1}{2} \text{tr}(\Omega_0 T)\right] \cdot \exp\left[-\frac{1}{2} \text{tr}\left(T \left[(B - M_n)^T \Lambda_n (B - M_n) + (Y^T P Y + M_0^T \Lambda_0 M_0 - M_n^T \Lambda_n M_n)\right]\right)\right] . \quad (3.18)$$

Using the posterior distribution over B , we can rewrite this as

$$p(Y, B, T) = \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|T|^p |\Lambda_0|^v}{(2\pi)^{pv}}} \sqrt{\frac{(2\pi)^{pv}}{|T|^p |\Lambda_n|^v}} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}} \frac{1}{\Gamma_v\left(\frac{\nu_0}{2}\right)}} \cdot |T|^{(\nu_0 - v - 1)/2} \exp\left[-\frac{1}{2} \text{tr}(\Omega_0 T)\right] \cdot \mathcal{MN}(B; M_n, \Lambda_n^{-1}, T^{-1}) \cdot \exp\left[-\frac{1}{2} \text{tr}\left(T \left[Y^T P Y + M_0^T \Lambda_0 M_0 - M_n^T \Lambda_n M_n\right]\right)\right] . \quad (3.19)$$

Now, B can be integrated out easily:

$$\int p(Y, B, T) dB = \sqrt{\frac{|T|^n |P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|\Lambda_0|^v}{|\Lambda_n|^v}} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}} \frac{1}{\Gamma_v\left(\frac{\nu_0}{2}\right)}} \cdot |T|^{(\nu_0 - v - 1)/2} \cdot \exp\left[-\frac{1}{2} \text{tr}\left(T \left[\Omega_0 + Y^T P Y + M_0^T \Lambda_0 M_0 - M_n^T \Lambda_n M_n\right]\right)\right] . \quad (3.20)$$

Using the posterior distribution over T , we can rewrite this as

$$\int p(Y, B, T) dB = \sqrt{\frac{|P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|\Lambda_0|^v}{|\Lambda_n|^v}} \sqrt{\frac{|\Omega_0|^{\nu_0}}{2^{\nu_0 v}}} \sqrt{\frac{2^{\nu_n v}}{|\Omega_n|^{\nu_n}} \frac{\Gamma_v\left(\frac{\nu_n}{2}\right)}{\Gamma_v\left(\frac{\nu_0}{2}\right)}} \cdot \mathcal{W}(T; \Omega_n^{-1}, \nu_n) . \quad (3.21)$$

Finally, T can also be integrated out:

$$\iint p(Y, B, T) dB dT = \sqrt{\frac{|P|^v}{(2\pi)^{nv}}} \sqrt{\frac{|\Lambda_0|^v}{|\Lambda_n|^v}} \sqrt{\frac{\frac{1}{2} |\Omega_0|^{\nu_0}}{\frac{1}{2} |\Omega_n|^{\nu_n}} \frac{\Gamma_v\left(\frac{\nu_n}{2}\right)}{\Gamma_v\left(\frac{\nu_0}{2}\right)}} = p(Y|m) . \quad (3.22)$$

Thus, the *log model evidence* of the MGLM-NW is given by

$$\begin{aligned} \log p(Y|m) = & \frac{v}{2} \log |P| - \frac{nv}{2} \log(2\pi) + \frac{v}{2} \log |\Lambda_0| - \frac{v}{2} \log |\Lambda_n| + \\ & \frac{\nu_0}{2} \log \left| \frac{1}{2} \Omega_0 \right| - \frac{\nu_n}{2} \log \left| \frac{1}{2} \Omega_n \right| + \log \Gamma_v \left(\frac{\nu_n}{2} \right) - \log \Gamma_v \left(\frac{\nu_0}{2} \right) . \end{aligned} \quad (3.23)$$

3.7 Cross-validated LME

For calculation of the *cross-validated log model evidence* (cvLME), the data are splitted into S subsets. In the training phase, all except one subset of the data are analyzed using a non-informative prior $p_{\text{ni}}(B, T)$ with the prior parameters

$$M_0 = 0_{pv}, \quad \Lambda_0 = 0_{pp} \quad \text{and} \quad \Omega_0 = 0_{vv}, \quad \nu_0 = 0 \quad (3.24)$$

to obtain an informative posterior $p(B, T | \cup_{j \neq i} y_j)$ using equations (3.14) and (3.15). In the testing phase, this informative posterior is then applied as a prior distribution to obtain the out-of-sample log model evidence $\log p(y_i | \cup_{j \neq i} y_j)$ via equation (3.23). Summing up over data subsets yields the cvLME according to equation (1.2).

As one can see from equations (3.14) and (3.15), the priors in (3.24) are non-informative in the sense that only the data remain to influence the posteriors.

3.8 Special cases

The *univariate Gaussian with unknown variance* (UGuv) is a special case in which

$$Y = y, \quad X = 1_n, \quad B = \mu, \quad \Sigma = \sigma^2 \quad \text{and} \quad V = I_n . \quad (3.25)$$

Furthermore, *multiple linear regression* (MLR) is a special case of the MGLM where

$$Y = y, \quad B = \beta \quad \text{and} \quad \Sigma = \sigma^2 . \quad (3.26)$$

The *one-sample t-test*, the *two-sample t-test*, the *paired t-test* and the *omnibus F-test* can all be emulated as comparisons of general linear models with specific design matrices.

3.9 Implementation

In **MATLAB**, maximum likelihood estimates and Bayesian posterior distributions can be obtained via the functions `MGLM_MLE` and `MGLM_Bayes` while log model evidence and cross-validated LME can be calculated using the functions `MGLM_LME` and `MGLM_cvLME`. Given an $n \times v$ data matrix Y , an $n \times p$ design matrix X , an $n \times n$ precision matrix P and a number of data subsets S , the cvLME for a MGLM-NW is calculated as

$$\text{cvLME} = \text{MGLM_cvLME}(Y, X, P, S); \quad (3.27)$$

In **Python**, an MGLM object has to be initiated via `mglm = cvBMS.MGLM(Y, X, V)` and maximum likelihood estimates, Bayesian posterior distributions, log model evidence and cross-validated LME are evaluated via `mglm.MLE`, `mglm.Bayes`, `mglm.LME`, and `mglm.cvLME`. Given Y , X , V and S as above, the cvLME for a MGLM-NW is calculated as

$$\text{cvLME} = \text{cvBMS.MGLM}(Y, X, V).\text{cvLME}(S) \quad (3.28)$$

In all of the above, V and P default to I_n whereas S defaults to 2 when left empty.

4 Poisson Distribution with Exposures

4.1 Likelihood function

Let $y = \{y_1, \dots, y_n\}$ with $y_i \in \mathbb{N}$ be a series of observed *counts* and let $x = \{x_1, \dots, x_n\}$ with $x_i \in \mathbb{R}$ be a series of concurrent *exposures*, some quantity that might or might not influence the measured counts. Then, according to a relatively simple model, each observation (y) would be Poisson-distributed with the Poisson rate being a product of the concurrent exposure (x) and some unknown constant (λ):

$$p(y_i|\lambda) = \text{Poiss}(y_i; \lambda x_i) = \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!}. \quad (4.1)$$

Assuming independence between individual observations, i.e. factorization of individual likelihoods, this would imply the following *likelihood function*:

$$p(y|\lambda) = \prod_{i=1}^n p(y_i|\lambda) = \prod_{i=1}^n \text{Poiss}(y_i; \lambda x_i) = \prod_{i=1}^n \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!}. \quad (4.2)$$

4.2 Maximum likelihood

Classical model estimation proceeds by maximizing the *log-likelihood* (LL)

$$\text{LL}(\lambda) = \log p(y|\lambda) = \sum_{i=1}^n [y_i \log(\lambda x_i) - \lambda x_i - \log \Gamma(y_i + 1)] \quad (4.3)$$

which gives rise to *maximum-likelihood* (ML) parameter estimates

$$\hat{\lambda} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i} = \frac{n\bar{y}}{n\bar{x}} = \frac{\bar{y}}{\bar{x}}. \quad (4.4)$$

4.3 Prior distribution

A conjugate prior distribution relative to the likelihood function given by (4.2) is the *gamma distribution* over the Poisson rate λ which is given by

$$p(\lambda) = \text{Gam}(\lambda; a_0, b_0) = \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] \quad (4.5)$$

where a_0 and b_0 are the prior shape and rate parameters for λ .

4.4 Joint likelihood

Combining the likelihood function (4.2) with the prior distribution (4.5), the *joint likelihood function* of the Poisson distribution with exposures (Poiss-exp) becomes

$$\begin{aligned} p(y, \lambda) &= p(y|\lambda) p(\lambda) \\ &= \prod_{i=1}^n \left(\frac{(\lambda x_i)^{y_i} \exp[-\lambda x_i]}{y_i!} \right) \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda]. \end{aligned} \quad (4.6)$$

Multiplying out the product gives:

$$p(y, \lambda) = \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \lambda^{n\bar{y}} \exp[-n\bar{x}\lambda] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0\lambda] . \quad (4.7)$$

Collecting identical variables gives:

$$p(y, \lambda) = \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0+n\bar{y}-1} \exp[-(b_0 + n\bar{x})\lambda] . \quad (4.8)$$

4.5 Posterior distribution

The *posterior distribution* of the Poisson can be evaluated using Bayes' theorem:

$$p(\lambda|y) = \frac{p(y|\lambda) p(\lambda)}{p(y)} . \quad (4.9)$$

Since $p(y)$ is just a normalization factor, the posterior is proportional to the joint:

$$p(\lambda|y) \propto p(y|\lambda) p(\lambda) = p(y, \lambda) . \quad (4.10)$$

From the term in (4.8), we can isolate the posterior distribution over λ :

$$\begin{aligned} p(\lambda|y) &= \text{Gam}(\lambda; a_n, b_n) \\ a_n &= a_0 + n\bar{y} \\ b_n &= b_0 + n\bar{x} . \end{aligned} \quad (4.11)$$

Note that \bar{y} and \bar{x} are the averages of y and x and therefore $n\bar{y}$ and $n\bar{x}$ are the sums of all elements in y and x , respectively.

4.6 Log model evidence

According to the law of marginal probability, the *model evidence* of the Poisson is:

$$p(y|m) = \int p(y|\lambda) p(\lambda) d\lambda . \quad (4.12)$$

According to the law of conditional probability, the integrand is equivalent to the joint:

$$p(y|m) = \int p(y, \lambda) d\lambda . \quad (4.13)$$

In (4.8), we have already evaluated this term as

$$p(y, \lambda) = \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0+n\bar{y}-1} \exp[-(b_0 + n\bar{x})\lambda] . \quad (4.14)$$

Using the posterior distribution over λ , we can rewrite this as

$$p(y, \lambda) = \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} \frac{\Gamma(a_n)}{b_n^{a_n}} \text{Gam}(\lambda; a_n, b_n) . \quad (4.15)$$

Now, λ can be integrated out easily:

$$\int p(y, \lambda) d\lambda = \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} = p(y|m) . \quad (4.16)$$

Thus, the *log model evidence* of the Poisson is given by

$$\begin{aligned} \log p(y|m) &= \sum_{i=1}^n y_i \log(x_i) - \sum_{i=1}^n \log \Gamma(y_i + 1) + \\ &\quad \log \Gamma(a_n) - \log \Gamma(a_0) + a_0 \log b_0 - a_n \log b_n . \end{aligned} \quad (4.17)$$

4.7 Cross-validated LME

For calculation of the *cross-validated log model evidence* (cvLME), the data are splitted into S subsets. In the training phase, all except one subset of the data are analyzed using a non-informative prior $p_{\text{ni}}(\lambda)$ with the prior parameters

$$a_0 = 0 \quad \text{and} \quad b_0 = 0 \quad (4.18)$$

to obtain an informative posterior $p(\lambda | \cup_{j \neq i} y_j)$ using equation (4.11). In the testing phase, this informative posterior is then applied as a prior distribution to obtain the out-of-sample log model evidence $\log p(y_i | \cup_{j \neq i} y_j)$ via equation (4.17). Summing up over data subsets yields the cvLME according to equation (1.2).

As one can see from equation (4.11), the priors in (4.18) are non-informative in the sense that only the data remain to influence the posteriors.

4.8 Special cases

The *Poisson distribution without exposures* (Poiss) is a special case in which

$$x = 1_n , \quad (4.19)$$

i.e. the exposures x are constant and one, such that $\bar{x} = 1$ and $n\bar{x} = n$.

4.9 Implementation

In **MATLAB**, maximum likelihood estimates and Bayesian posterior distributions can be obtained via the functions `Poiss_MLE` and `Poiss_Bayes` while log model evidence and cross-validated LME can be calculated using the functions `Poiss_LME` and `Poiss_cvLME`. Given an $n \times v$ data matrix Y , an $n \times 1$ design vector x and a number of data subsets S , the cvLME for the Poisson is calculated as

$$\text{cvLME} = \text{Poiss_cvLME}(Y, x, S); \quad (4.20)$$

In **Python**, a Poisson object has to be initiated via `poiss = cvBMS.Poiss(Y, x)` and maximum likelihood estimates, Bayesian posterior distributions, log model evidence and cross-validated LME are evaluated via `poiss.MLE`, `poiss.Bayes`, `poiss.LME`, and `poiss.cvLME`. Given Y , x and S as above, the cvLME for the Poisson is calculated as

$$\text{cvLME} = \text{cvBMS.Poiss}(Y, x).\text{cvLME}(S); \quad (4.21)$$

In all of the above, x defaults to 1_n whereas S defaults to 2 when left empty.