

Øyvind Hoveid

Norwegian Institute of Bioeconomic Research,
Mail: oyvind.hoveid@nibio.no

AMLGM, Trondheim, September 17-18, 2015

Introduction

Gaussian distributed latent variables are unsurpassed in hierarchical/spatio/temporal models due to the equivalence of contingent independence and zero cross-precision (Rue and Held, 2005). However, many real world (e.g. economic) phenomena require other distributions. I will here point to a broader class of models in which the Gaussian latents depend on non-Gaussian shocks, and the predictors of observations are transformed Gaussian. Within the class it is necessary to keep determinants of Gaussian densities non-singular. Informative Inverted Wishart priors on variance matrices guard effectively against this situation. The idea stems from the term "Boundary avoiding priors" introduced by Gelman (2014). The priors are changed from Gamma to Inverse Wishart, though.

Informative priors need a justification. I will promote them as acceptable with a scaling factor $\omega \in (0, 1)$ for error variance determinants relative to the sum of latent and error variance determinants. With $\omega = 0$, error variance is singular. With $\omega = 1$, latent variance is singular. Estimation will take various fixed interior ω . A final value will be found according to a model selection criterion that rates models according to likelihood less a penalty for model complexity. Smaller ω means more complexity. Watanabe's information criterion (WAIC) (Watanabe 2010, 2013), (Gelman, Hwang, Vehtari 2013) works well for such models. With the framework on top, the priors are denoted *relatively informative priors* (RIP).

Model class

Consider models of the following class:

$$\phi(\mathbf{y} - \mathbf{B}\mathbf{x}, \mathbf{R}), \mathbf{x} \equiv \mathbf{F}(\mathbf{w}), \phi(\mathbf{w} - \mathbf{A}\mathbf{v}, \mathbf{Q}), \pi(\mathbf{v}|\mathbf{P}) \\ \mathcal{W}^{-1}(\mathbf{P}, \mathbf{Q}), \mathcal{W}^{-1}(\mathbf{R}|\omega, |\mathbf{P}|, |\mathbf{Q}|), \omega \in (0, 1)$$

where

- \mathbf{y} is observations with multivariate Gaussian distribution ϕ
- $\mathbf{B}\mathbf{x}$ is the predictor of \mathbf{y} selected from latents \mathbf{x} with selector matrix \mathbf{B} (with one non-zero element equal to 1 in each row)
- \mathbf{R} is noise variance
- \mathbf{F} is a 1-1 mapping from Gaussian latents \mathbf{w} to predictors \mathbf{x}
- $\mathbf{A}\mathbf{v}$ is the predictor of \mathbf{w} generated from the vector \mathbf{v} with a selector matrix \mathbf{A}
- \mathbf{Q} is the variance matrix of \mathbf{w} , with a dependence structure in terms of a sparse \mathbf{Q}^{-1} . \mathbf{w} is sorted to orient dependencies towards NW. Sparsity is then preserved in Cholesky factors with reversed Cholesky factorization: $\mathbf{Q}^{-1} = \mathbf{Q}^{-\frac{1}{2}}\mathbf{Q}^{-\frac{1}{2}}$
- \mathbf{v} is a vector of independent non-Gaussian variables with a diagonal scale matrix \mathbf{P}

Model selection within the fully Gaussian frame

Consider the total likelihood:

$$\mathcal{L}(\mathbf{y}, \mathbf{x}|\mathbf{Q}, \mathbf{R}, \mathbf{v}) = \phi(\mathbf{y} - \mathbf{B}\mathbf{x}, \mathbf{R})\phi(\mathbf{x} - \mathbf{A}\mathbf{v}, \mathbf{Q})$$

with posterior predictor of \mathbf{x} :

$$\hat{\mathbf{x}}(\mathbf{Q}, \mathbf{R}, \mathbf{v}, \mathbf{y}) = \mathbf{H}(\mathbf{B}^T\mathbf{R}^{-1}\mathbf{y} + \mathbf{Q}^{-1}\mathbf{A}\mathbf{v})$$

where

$$\mathbf{H} = (\mathbf{B}^T\mathbf{R}^{-1}\mathbf{B} + \mathbf{Q}^{-1})^{-1}$$

With flat priors, the hyper-parameters $(\mathbf{Q}, \mathbf{R}, \mathbf{v})$, can be varied until maximum likelihood. The likelihood can be stated as the total Gaussian one divided by the Laplace density of \mathbf{x} at $\hat{\mathbf{x}}$ (INLA). Analytically:

$$\log \pi(\mathbf{y}|\mathbf{Q}, \mathbf{R}, \mathbf{v}) = \log [\phi(\mathbf{y} - \mathbf{B}\hat{\mathbf{x}}, \mathbf{R})\phi(\hat{\mathbf{x}} - \mathbf{A}\mathbf{v}, \mathbf{Q}) / \phi(\mathbf{0}, \mathbf{H})] \\ = \log \phi(\mathbf{y} + \mathbf{C}\mathbf{A}\mathbf{v}, \mathbf{G}) + \frac{1}{2} [(\mathbf{A}\mathbf{v})^T (\mathbf{C}^T\mathbf{G}^{-1}\mathbf{C} - \mathbf{B}^T\mathbf{G}^{-1}\mathbf{B}) \mathbf{A}\mathbf{v}]$$

where

$$\mathbf{G} = \mathbf{R} + \mathbf{B}\mathbf{Q}\mathbf{B}^T, \quad \mathbf{C} = \mathbf{G}\mathbf{R}^{-1}\mathbf{B}\mathbf{H}^{-1}\mathbf{Q}^{-1}$$

There is no chance of hitting $|\mathbf{G}| = 0$ in this case and no informative prior is needed, but the model selection issue remains. The constraint:

$$\omega|\mathbf{Q}| = (1 - \omega)|\mathbf{R}| \quad (1)$$

with models solved for different ω , will span out a set from which model selection can take place.

Model selection with transformed Gaussian latents and flat priors

Transformed Gaussian latents, $\mathbf{x} = \mathbf{F}(\mathbf{w})$, brings in a Jacobian determinant to state the prior distribution of \mathbf{x} :

$$\mathcal{L}(\mathbf{y}, \mathbf{w}|\mathbf{Q}, \mathbf{R}, \mathbf{v}) = \phi(\mathbf{y} - \mathbf{B}\mathbf{F}(\mathbf{w}), \mathbf{R})\phi(\mathbf{w} - \mathbf{A}\mathbf{v}, \mathbf{Q}) / |\partial\mathbf{F}(\mathbf{w})|$$

Some structure, say

$$\mathbf{F}(\mathbf{w}) = (\mathbf{F}_1(\mathbf{w}_1), \dots, \mathbf{F}_t(\mathbf{w}_{1:t}))$$

will be helpful for the specification of the determinant by means of the Cholesky factor. Structures of contingent independence within \mathbf{w} will only be disturbed as far as $\mathbf{F}_k(\mathbf{w})$ depends on \mathbf{w}_j with $j < k$ while $\mathbf{Q}_{kj}^{-1} = 0$.

One might integrate $\mathcal{L}(\mathbf{y}, \mathbf{w}|\mathbf{Q}, \mathbf{R}, \mathbf{v})$ with respect to \mathbf{w} numerically, and search for a model of $\mathcal{L}(\mathbf{y}|\mathbf{Q}, \mathbf{R}, \mathbf{v})$. For sake of model selection one need not involve any constant ω , as $(\mathbf{Q}, \mathbf{R}, \mathbf{v})$ are constants themselves.

Model selection: transformed Gaussian latents and RIPs

Let \mathbf{Q} and \mathbf{R} be stochastic and introduce RIP-s, $\mathcal{W}^{-1}(\mathbf{Q}|\omega\mathbf{Q}, \mathbf{q})$ and $\mathcal{W}^{-1}(\mathbf{R}|\omega\mathbf{R}, \mathbf{r})$. The constraint on determinants has now been softened with scaling of Wishart parameters $\omega\mathbf{Q}$ and $(1 - \omega)\mathbf{R}$ respectively. \mathbf{q} and \mathbf{r} are related to dimensionality and degrees of freedom. \mathbf{Q} and \mathbf{R} are fixed matrices proportionate to $\mathbf{E}\mathbf{x} \mathbf{Q}$ and $\mathbf{E}\mathbf{x} \mathbf{R}$ respectively.

Total likelihood is:

$$\mathcal{L}(\mathbf{Q}, \mathbf{R}, \mathbf{w}, \mathbf{y}) = \phi(\mathbf{y} - \mathbf{B}\mathbf{F}(\mathbf{w}), \mathbf{R})\phi(\mathbf{w} - \mathbf{A}\mathbf{v}, \mathbf{Q})\mathcal{W}^{-1}(\mathbf{Q}|\omega\mathbf{Q}, \mathbf{q})\mathcal{W}^{-1}(\mathbf{R}|\omega\mathbf{R}, \mathbf{r}) / |\partial\mathbf{F}(\mathbf{w})|$$

and ignoring constants:

$$\log \mathcal{L}(\mathbf{Q}, \mathbf{R}, \mathbf{w}, \mathbf{y}) = \mathbf{q} \log |\mathbf{Q}^{-\frac{1}{2}}| + \mathbf{q} \log |\mathbf{R}^{-\frac{1}{2}}| - \log |\partial\mathbf{F}(\mathbf{w})| \\ - \frac{1}{2} [\text{tr}((\mathbf{y} - \mathbf{B}\mathbf{F}(\mathbf{w}))(\mathbf{y} - \mathbf{B}\mathbf{F}(\mathbf{w}))^T + (1 - \omega)\mathbf{R})\mathbf{R}^{-1} + ((\mathbf{w} - \mathbf{A}\mathbf{v})(\mathbf{w} - \mathbf{A}\mathbf{v})^T + \omega\mathbf{Q})\mathbf{Q}^{-1}]$$

In terms of *all parameters* $\mathbf{u} = (\mathbf{Q}, \mathbf{R}, \mathbf{w})$, $\mathcal{L}(\mathbf{u}, \mathbf{y})$ is bounded and one can find the mode $\hat{\mathbf{u}}$. The optimizing routine will deliver the negative Hessian matrix \mathbf{H}^{-1} at the mode, from which the Laplace approximation is: $\pi(\mathbf{u}) = \phi(\mathbf{u} - \hat{\mathbf{u}}, \mathbf{H})$. The Hessian matrix is expected to have a sparse factorization, $\mathbf{H}^{-1} = \mathbf{H}^{-\frac{1}{2}}\mathbf{H}^{-\frac{1}{2}}$. In order to sample from the Laplace distribution take $\mathbf{s} \in (0, 1)^h$, with h the dimension of \mathbf{H} . Denoting Φ as the cumulative standard normal distribution, $\Phi^{-1}(\mathbf{s})$ is then a h -dimensional standard normal, and the equation

$$\mathbf{H}^{-\frac{1}{2}}(\mathbf{u} - \hat{\mathbf{u}}) = \Phi^{-1}(\mathbf{s})$$

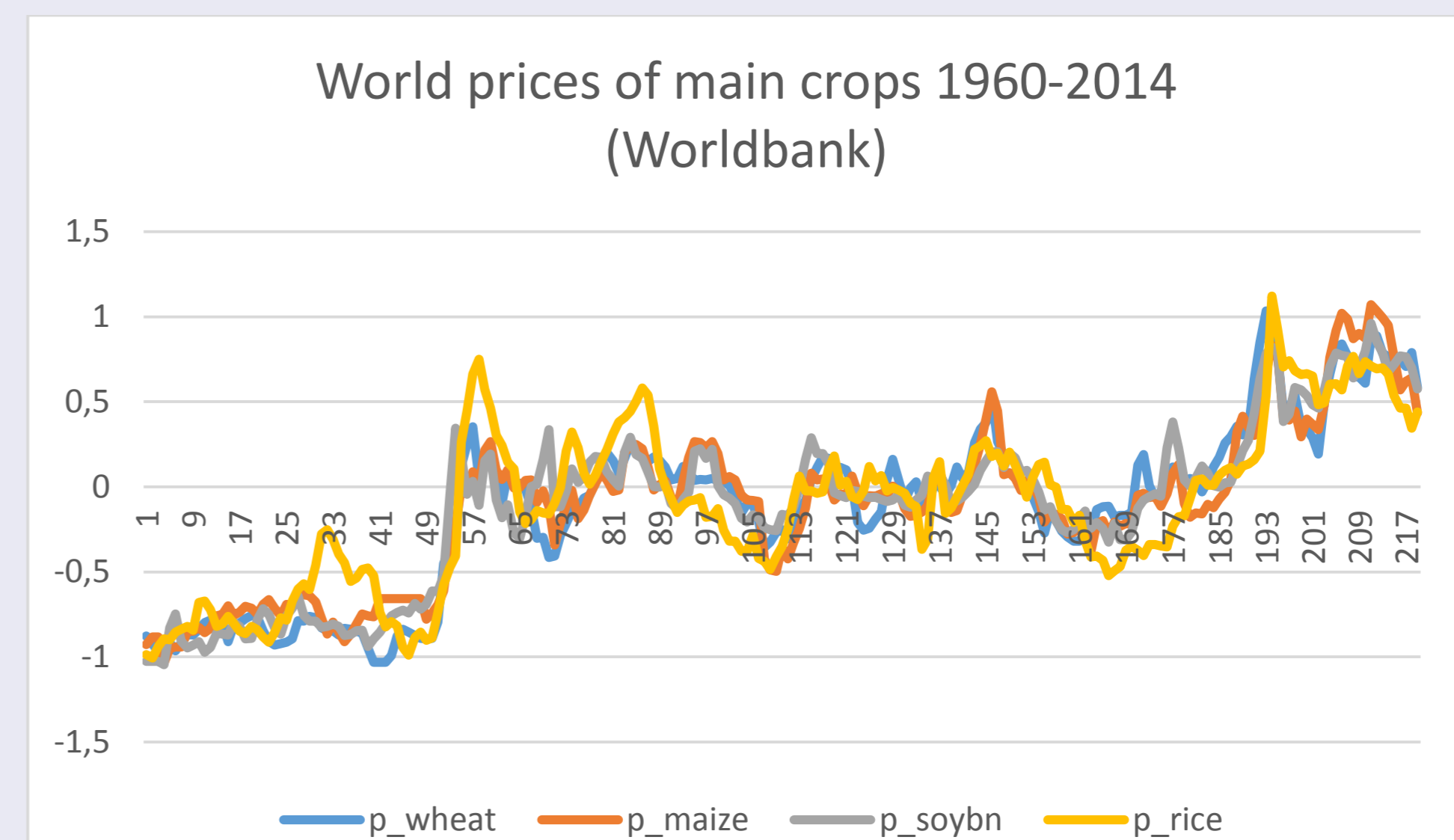
determines a sampled point. The log-probability of that selection is simply:

$$\log \pi(\mathbf{u}) = \log |\mathbf{H}^{-\frac{1}{2}}| - \frac{1}{2}(\Phi^{-1}(\mathbf{s}))^2$$

A weighted sample from the posterior is obtained by using likelihoods divided by selection probabilities as weights (importance sampling).

An economic model

Modelling starts out from observations of food prices in the world market. Logged and scaled prices of wheat, maize, rice and soybeans are displayed below.



Latent variables provide an opportunity to recover some of the underlying economic mechanisms. Food prices are assumed determined by temporary equilibria where demand equals supply. Demand behaviour of a representative consumer is specified in terms of a cost function, $\mathbf{C}(\mathbf{p}_t, \mathbf{u}_t)$, specifying the minimum income to obtain utility level \mathbf{u}_t at market prices \mathbf{p}_t . The structure of cost minimization then leads to the equations,

$$\mathbf{c}_t = \partial_{\mathbf{p}} \mathbf{C}(\mathbf{p}_t, \mathbf{u}_t)$$

$$\mathbf{i}_t = \mathbf{C}(\mathbf{p}_t, \mathbf{u}_t)$$

where \mathbf{c}_t is consumed quantity and \mathbf{i}_t is income. The functional specification of \mathbf{C} is a CES (constant elasticity of substitution) function:

$$\mathbf{C}(\mathbf{p}, \mathbf{u}) = \left[\sum_k (\beta_k(\mathbf{u}) \mathbf{p}_k)^\rho \right]^{\frac{1}{\rho}}, \quad \partial_{\mathbf{p}} \mathbf{C}(\mathbf{p}, \mathbf{u}) = \beta(\mathbf{u})^T \cdot \left[\frac{\beta(\mathbf{u})^T \cdot \mathbf{p}^T}{\mathbf{C}(\mathbf{p}, \mathbf{u})} \right]^{\rho-1}$$

With respect to production: it takes a long time to produce food, even longer times to increase capacity in periods of high prices, and the outcome is always uncertain. The production is consequently modelled as a long-lagged temporary process in which nature provide independent shocks, \mathbf{v}_t .

Translation to the statistical model:

- Observations \mathbf{y}_{tk} : logged prices and income at t , = $\log \mathbf{p}_{tk}$ with $\mathbf{k} \in \{\text{wheat, sugar, soybeans, maize, other}\}$, and also $\mathbf{y}_{t,i'} = \log \mathbf{i}_t$
- Predictors \mathbf{x}_{tk} : of all observations \mathbf{y}_{tk} , and also, $\mathbf{x}_{t,i'} = \mathbf{u}_t$
- Gaussian variables $\mathbf{w}_{tk} = \log \mathbf{c}_{tk}$ with temporal structure, \mathbf{Q} and expectation \mathbf{v}_{tk} , and also independent parameters: $(\mathbf{w}_{0k} = \log \beta_k, \mathbf{w}_{00} = \log \rho)$
- Transformation $\mathbf{x} = \mathbf{F}(\mathbf{w})$ defined implicitly by $\exp(\mathbf{w}_t) = \partial_{\mathbf{p}} \mathbf{C}(\mathbf{x}_t)$ and $\exp(\mathbf{x}_{t,i'}) = \mathbf{C}(\mathbf{x}_t)$
- Independent Cauchy shocks: \mathbf{v}_{tk} for each crop \mathbf{k}

Relevant issues:

- Which time dependencies dominate \mathbf{Q}^{-1} ?
- What is the better structure of the cost function, $\mathbf{C}(\mathbf{p}, \mathbf{u})$?
- What is the risk of the consumer due to food price fluctuations?