

Statistical Elucidation of Latent Structures via Copulas

Radu Craiu

Department of Statistical Sciences
University of Toronto

Copulas for serially correlated data with hidden structures

This project was done in collaboration with:

- ▶ Robert Zimmerman (Toronto)



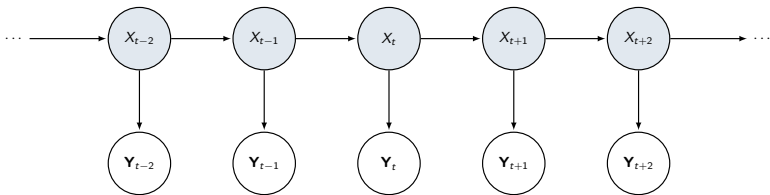
- ▶ Vianey Leos Barajas (Toronto)



Paper: *Copula Modelling of Serially Correlated Multivariate Data with Hidden Structures (JASA, 2024).*

Hidden Markov Models: A Primer

- ▶ A hidden Markov model (HMM) pairs an observed time series $\{\mathbf{Y}_t\}_{t \geq 1} \subseteq \mathbb{R}^d$ with a Markov chain $\{X_t\}_{t \geq 1}$ on some state space \mathcal{X} , such that the distribution of $\mathbf{Y}_s \mid X_s$ is independent of $\mathbf{Y}_t \mid X_t$ for $s \neq t$:



- ▶ $\mathbf{Y}_{t,h} \mid \{X_t = k\} \sim f_{k,h}(\cdot \mid \lambda_{k,h}) \quad \forall h = 1, \dots, d$
- ▶ $\{X_t\}$ is a Markov process (finite state space \mathcal{X}) with initial probability mass distribution $\{\pi_i\}_{i \in \mathcal{X}}$ and transition probabilities $\{\gamma_{i,j}\}_{i,j \in \mathcal{X}}$

Inferential aims for HMMs

- ▶ Typically, the chain $\{X_t\}_{t \geq 1}$ is partially or completely unobserved.
- ▶ The hidden states can correspond to a precise variable (occupancy data) or might be postulated (psychology, ecology, etc)
- ▶ **Aim 1:** Model the data generating mechanism (Nasri et al., 2020)
- ▶ **Aim 2:** Decode (i.e., classify) or predict the X_t 's from the observed data.

Examples

- ▶ A tri-axial accelerometer captures a shark's acceleration with respect to three positional axes depending on the shark's activity (resting, hunting, attacking).
- ▶ Hospitals repeatedly record a patient's vitals post-surgery to understand the state of their recovery (stable, complications, complete).
- ▶ Stock exchanges keep track of real-time prices for hundreds of stocks within an industry, depending on market conditions/states (stagnant, growing, shrinking).

Fusion of Multiple Data Sources

- ▶ In the real-world applications above, various sensors capture multiple streams of data, which are “fused” into a multivariate time series $\{\mathbf{Y}_t\}_{t \geq 1}$.
- ▶ In such situations, the components of any $\mathbf{Y}_t = (Y_{t,1}, \dots, Y_{t,d})$ cannot be assumed independent (even conditional on X_t).
- ▶ The corresponding assumption for HMMs – that of contemporaneous conditional independence (Zucchini et al, 2017) – is often violated.
- ▶ Instead, it is common to assume that \mathbf{Y}_t follows a multivariate Gaussian distribution, but this places limits on marginals and dependence structures.
- ▶ What if the strength/type of dependence between the components of \mathbf{Y}_t could be informative about the underlying state X_t ?

Occupancy Data

- ▶ The ability to detect whether a room is occupied using sensor data (such as temperature and CO₂ levels)
- ▶ Consider three publicly-available labelled datasets presented by Candanedo and Feldheim (2016) which contain multivariate time series of four environmental measurements (light, temperature, humidity, CO₂) and one derived metric (the humidity ratio).

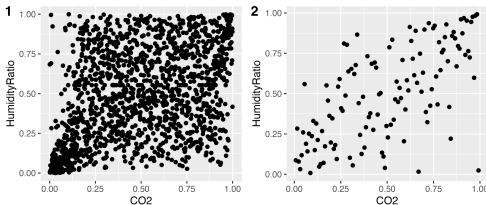


Figure: Pseudo-observations computed from unoccupied (Panel 1) and occupied (Panel 2) subsets.

Copulas: The Joys

- ▶ Copulas are mathematical devices used to model dependence between random variables regardless of their marginals.
- ▶ Copulas are useful for data fusion/integration as they lead to coherent joint models, even when the marginals are in different families or of different types.
- ▶ Copulas unlock information contained in the dependence part of the distribution (second-order) that complements the information in the marginals.
- ▶ Copulas extend statistical methods beyond the use of a multivariate Gaussian or Student.

At the root of it all, a theorem

- ▶ Copulas are distribution functions on $[0, 1]^d$ that model dependence between continuous random variables.
- ▶ **Sklar's Theorem:** If Y_1, Y_2, \dots, Y_d are continuous r.v.'s with cdfs F_1, \dots, F_d , there exists an unique copula function $C : [0, 1]^d \rightarrow [0, 1]$ such that

$$H(t_1, \dots, t_d) = P(Y_1 \leq t_1, \dots, Y_d \leq t_d) = C(F_1(t_1), \dots, F_d(t_d)).$$

- ▶ The copula bridges the marginal distributions of Y_1, \dots, Y_d with the joint distribution. It corresponds to a distribution on $[0, 1]^d$ with uniform margins.
- ▶ This can be extended to conditional distributions and copulas:

$$P(Y_1 \leq t_1, \dots, Y_d \leq t_d | X) = C(F_1(t_1 | X), \dots, F_d(t_d | X) | X).$$

Copulas Within HMMs

- ▶ Our model consists of an HMM $\{(\mathbf{Y}_t, X_t)\}_{t \geq 1} \subseteq \mathbb{R}^d \times \mathcal{X}$ in which the state-dependent distributions are copulas:

$$\mathbf{Y}_t \mid (X_t = k) \sim H_k(\cdot) = \underbrace{C_k\left(F_{k,1}(\cdot; \lambda_{k,1}), \dots, F_{k,d}(\cdot; \lambda_{k,d})\right)}_{\text{depends on the hidden state value } k} \mid \theta_k.$$

- ▶ $C_k(\cdot, \dots, \cdot \mid \theta_k)$ is a d -dimensional parametric copula
- ▶ $\{X_t\}_{t \geq 1}$ is a Markov process on finite state space $\mathcal{X} = \{1, 2, \dots, K\}$ and K is known.
- ▶ In this model, virtually all aspects of the state-dependent distributions are allowed to vary between states

Information in the dependence

- For a range of $\theta \in [0, 100)$, we simulated a bivariate time series of length $T = 100$ from the 2-state HMM

$$Y_t | (X_t = k) \sim C_{Frank}(N(0, 1), N(0, 1) | (-1)^k \cdot |\theta|), \quad k = 1, 2$$

and then separately assessed the accuracy of a standard decoding algorithm, first assuming independent margins and then the true model:

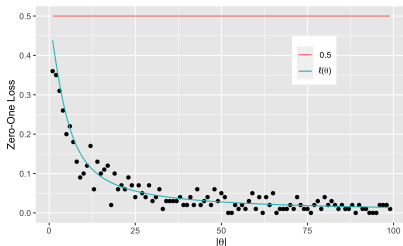


Figure: Zero-one losses for independent margins (red dots) and true model (blue dots)

Estimation with missing data

- ▶ Data consist in observed $\mathbf{Y}_{1:T}$ and missing $X_{1:T}$
- ▶ Parameters are $\eta = \{\lambda_{h,k}\}_{\substack{h=1:d \\ k=1:T}} \cup \{\theta_k\}_{k=1:T} \cup \{\gamma_{i,j}\}_{\substack{i=1:K \\ j=1:K}} \cup \{\pi_j\}_{j=1:K}$.
- ▶ The complete-data log-likelihood for one trajectory of the copula HMM is given by

$$\begin{aligned} \ell_{\text{com}}(\boldsymbol{\eta} \mid \mathbf{y}_{1:T}, X_{1:T}) &= \pi_{X_1} + \sum_{t=2}^T \log \gamma_{X_{t-1}, X_t} + \sum_{h=1}^d \log f_{X_t, h}(y_{t, h}; \lambda_{X_t, h}) \\ &+ \sum_{t=1}^T \log c_{X_t}(F_{X_t, 1}(y_{t, 1}; \lambda_{X_t, 1}), \dots, F_{X_t, 1}(y_{t, d}; \lambda_{X_t, d}) \mid \theta_{X_t}). \end{aligned} \quad (1)$$

Inference for HMMs Via the EM Algorithm

- ▶ Without copula, the estimation is done via the EM algorithm (aka Baum-Welch)

E-step Compute $Q(\eta|\eta^{(s)}) = E[l_{com}(\eta|\mathbf{Y}_{1:T}, X_{1:T})|\eta^{(s)}, \mathbf{Y}_{1:T}]$

M-step Set $\eta^{(s+1)} = \arg \max_{\eta} Q(\eta|\eta^{(s)})$

- ▶ The complete-data log-likelihood is written in terms of the state membership indicators $U_{k,t} = \mathbb{1}_{X_t=k}$ and $V_{j,k,t} = \mathbb{1}_{X_{t-1}=j, X_t=k}$
- ▶ In the **E-Step**, these indicators are estimated by the conditional probabilities $\hat{u}_{k,t} = \mathbb{P}(X_t = k | \mathbf{Y}_{1:T} = \mathbf{y}_{1:T})$ and $\hat{v}_{j,k,t} = \mathbb{P}(X_{t-1} = j, X_t = k | \mathbf{Y}_{1:T} = \mathbf{y}_{1:T})$, which are computed based on current parameter estimates
- ▶ This only requires evaluating the state-dependent densities at each of the observations $\mathbf{y}_1, \dots, \mathbf{y}_T$ (this is “OK”)

The M-Step Is Hard

- ▶ In the **M-Step**, the resulting complete-data log-likelihood is maximized with respect to all parameters in the model simultaneously
 - ▶ Only for the simplest univariate models do the state-dependent MLEs exist in closed form; otherwise, one must resort to numerical methods (**this is hard and unstable!**)
 - ▶ Evaluating a copula density $c_k(\cdot, \dots, \cdot | \theta_k)$ in high dimensions is slow
 - ▶ When the state-dependent distributions in an HMM are copulas, performing the M-Step directly requires the evaluation of

$$\operatorname{argmax}_{\{\theta_k\}, \{\lambda_{k,h}\}} \left\{ \sum_{k=1}^K \sum_{t=1}^T \hat{u}_{k,t} \left[\log c_k \left(F_{k,1}(y_{t,1}; \lambda_{k,1}), \dots, F_{k,d}(y_{t,d}; \lambda_{k,d}) \mid \theta_k \right) + \sum_{h=1}^d \log f_{k,h}(y_{t,h}; \lambda_{k,h}) \right] \right\}$$

- ▶ This is very unstable (and slow)

Inference Functions for Margins

- ▶ Likelihood-based inference for copulas is easier when the goal is to estimate θ alone in the presence of known margins
- ▶ Why not perform inference on the marginal distributions first, and then on the copula itself?
- ▶ In the context of iid data, this is exactly the **inference functions for margins** (IFM) approach of Joe and Xu (1996):
 - ▶ First estimate each λ_h by its “marginal MLE” $\hat{\lambda}_h$ given $\{Y_{t,h}\}_{t \geq 1}$, for $h \in \{1, \dots, d\}$
 - ▶ Then estimate θ assuming fixed marginals $F_1(\cdot; \hat{\lambda}_1), \dots, F_d(\cdot; \hat{\lambda}_d)$
- ▶ One can show that the IFM estimator is consistent and asymptotically normal (although relatively less efficient than the MLE)

IFM Step

- ▶ For each $j \in \mathcal{X}$, estimate the initial distribution and transition probabilities:

$$\delta^{(s+1)} = (\hat{u}_{1,1}^{(s)}, \dots, \hat{u}_{K,1}^{(s)})$$

and

$$\gamma_{j,\cdot}^{(s+1)} = \left(\frac{\sum_{t=2}^T \hat{v}_{j,1,t}^{(s)}}{\sum_{k=1}^K \sum_{t=2}^T \hat{v}_{j,k,t}^{(s)}}, \dots, \frac{\sum_{t=2}^T \hat{v}_{j,K,t}^{(s)}}{\sum_{k=1}^K \sum_{t=2}^T \hat{v}_{j,k,t}^{(s)}} \right).$$

- ▶ For each $k \in \{1, \dots, K\}$ and $h \in \{1, \dots, d\}$, estimate the marginal parameters

$$\lambda_{k,h}^{(s+1)} = \arg \sup_{\lambda} \sum_{t=1}^T \hat{u}_{k,t}^{(s+1)} \cdot \log(f_{k,h}(y_{t,h}; \lambda)).$$

- ▶ For each $k \in \{1, \dots, K\}$, estimate the copula parameters

$$\tilde{\theta}_k^{(s+1)} = \arg \sup_{\theta} \sum_{t=1}^T \hat{u}_{k,t}^{(s+1)} \cdot \log \left(c_k \left(F_{k,1}(y_{t,1}; \lambda_{k,1}^{(s+1)}), \dots, F_{k,d}(y_{t,d}; \lambda_{k,d}^{(s+1)}) \mid \theta \right) \right).$$

New problems

- ▶ The EIFM algorithm is not an GEM algorithm

$$\sum_{t=1}^T \hat{u}_t \cdot \log \left(f_h(y_{t,h}; \lambda_h^{(s)}) \right) \leq \sum_{t=1}^T \hat{u}_t \cdot \log \left(f_h(y_{t,h}; \lambda_h^{(s+1)}) \right), \quad h \in \{1, \dots, d\} \quad (2)$$

does not imply

$$\begin{aligned} & \sum_{t=1}^T \hat{u}_t \cdot \log \left(c \left(F_1(y_{t,1}; \lambda_1^{(s)}), \dots, F_d(y_{t,d}; \lambda_d^{(s)}) \mid \theta^{(s)} \right) \right) \\ & \leq \sum_{t=1}^T \hat{u}_t \cdot \log \left(c \left(F_1(y_{t,1}; \lambda_1^{(s+1)}), \dots, F_d(y_{t,d}; \lambda_d^{(s+1)}) \mid \theta^{(s)} \right) \right). \end{aligned}$$

- ▶ The EIFM algorithm will converge (to a local or global maximum).
- ▶ The estimator is consistent and asymptotically normal (under mild regularity conditions).
- ▶ EIFM as a version of the ES algorithm of Elashoff and Ryan (2004).
- ▶ Use asymptotic theory of M-estimators for HMMs Jensen (2011).

Implementation of EIFM

- ▶ K is assumed known
- ▶ An initial clustering algorithm may be used in which the observed multivariate data follow vine copulas (Sahin and Czado, 2022).
- ▶ We consider the k -means algorithm for clustering.
- ▶ Initial parameter values for the copula(s) are obtained using a Gaussian copula
- ▶ If marginals are Gaussian this means fitting a multivariate normal for each cluster.

Does This Work?

- ▶ For $T \in \{100, 1000, 5000\}$ and $d \in \{2, 5, 10\}$, we simulated a d -dimensional time series of length T from the 2-state HMM

$$\mathbf{Y}_t \mid (X_t = 1) \sim C_{\text{Frank}} \left(\left(\mathcal{N}(\mu_{1,h} = -h, 1) \right)_{h=1}^d \mid \theta_1 = 3 \right)$$

$$\mathbf{Y}_t \mid (X_t = 2) \sim C_{\text{Clayton}} \left(\left(\mathcal{N}(\mu_{2,h} = h, 1) \right)_{h=1}^d \mid \theta_2 = 3 \right)$$

and estimated $\boldsymbol{\eta} = (\mu_{1,1}, \dots, \mu_{2,d}, \theta_1, \theta_2)$ using both approaches

- ▶ Applied to the basic EM algorithm, R's `optim` with L-BFGS-B (i.e., quasi-Newton with box constraints) typically fails as soon as $d \geq 3$
 - ▶ The procedure is extremely sensitive to initial values and requires $\hat{\boldsymbol{\eta}}^{(0)} \approx \boldsymbol{\eta}$ just to avoid overflow
 - ▶ This kind of tuning is very tedious or impossible in high dimensions

Does This Work?

- ▶ We keep track of the **time** (in seconds) until the algorithm converges, and the **L_2 error** of the resulting estimate, $\epsilon = \|\boldsymbol{\eta} - \hat{\boldsymbol{\eta}}\|_2$
 - ▶ We used the `lbfgsb3c` package, which is more stable than `optim`

| | $d = 2$ | $d = 5$ | $d = 10$ |
|------------|----------------------------|-------------------------------|-------------------------------|
| $T = 100$ | 111.9 s, $\epsilon = 0.14$ | 123.4 s, $\epsilon = 299.98$ | 111.8 s, $\epsilon > 10^9$ |
| $T = 1000$ | 166.6 s, $\epsilon = 0.63$ | 169.5 s, $\epsilon > 10^{11}$ | 418.23 s, $\epsilon = 725.06$ |
| $T = 5000$ | ? | ? | ? |

Table: EM Algorithm

| | $d = 2$ | $d = 5$ | $d = 10$ |
|------------|----------------------------|----------------------------|----------------------------|
| $T = 100$ | 5.1 s, $\epsilon = 2.9$ | 3.0 s, $\epsilon = 0.94$ | 4.2 s, $\epsilon = 0.58$ |
| $T = 1000$ | 34.4 s, $\epsilon = 0.57$ | 22.9 s, $\epsilon = 0.60$ | 34.4 s, $\epsilon = 0.80$ |
| $T = 5000$ | 172.6 s, $\epsilon = 0.13$ | 106.2 s, $\epsilon = 0.12$ | 168.7 s, $\epsilon = 0.19$ |

Table: EFM Algorithm

Numerical Experiment I

- ▶ Generative model:

$$\mathbf{Y}_i \mid (X_i = k) \sim C_k (SN(\cdot; \xi_{k,1}, \omega_{k,1}, \alpha_{k,1}), SN(\cdot; \xi_{k,2}, \omega_{k,2}, \alpha_{k,2}) \mid \tau_k),$$

for $k \in \{1, \dots, 4\}$.

| State | Copula family | τ_k | $\xi_{k,1}$ | $\omega_{k,1}$ | $\alpha_{k,1}$ | $\xi_{k,2}$ | $\omega_{k,2}$ | $\alpha_{k,2}$ |
|-------|---------------|----------|-------------|----------------|----------------|-------------|----------------|----------------|
| 1 | Clayton | 0.2 | -4 | 1 | 5 | -1 | 1 | -3 |
| 2 | B4 | 0.4 | -2 | 1 | 3 | 2 | 1 | -3 |
| 3 | Gaussian | 0.6 | 0 | 1 | 5 | 3 | 1 | -5 |
| 4 | $t_{(\nu=5)}$ | 0.8 | 2 | 1 | 3 | 4 | 1 | -5 |

Table: True parameters for the state-dependent distributions.

Numerical Experiment I

| T | | 500 | 1000 | 2500 | 5000 |
|--------------------------|----------------|--------|--------|--------|--------|
| Stopping Rule Tolerance: | 0.01 | 14 | 24 | 23 | 15 |
| | 0.001 | 17 | 26 | 25 | 17 |
| | 0.0001 | 36 | 59 | 62 | 39 |
| | 0.00001 | 230 | 115 | 460 | 269 |
| Classifier: | k -means | 0.9020 | 0.9090 | 0.9200 | 0.9196 |
| | Local decoding | 0.9640 | 0.9640 | 0.9696 | 0.9732 |

For each $T \in \{500, 1000, 2500, 5000\}$: (Top rows) Number of iterations taken by the EIFM algorithm applied to $\{Y_t\}_{1 \leq t \leq T}$ before stopping using L_1 -norm tolerances in $\{0.01, 0.001, 0.0001, 0.00001\}$. (Bottom rows) Classification accuracy of initial k -means clustering and local decoding with parameter estimates provided by the EIFM algorithm.

Numerical Experiment I

| | | | | |
|--------------------------|-------|-------|--------|--------|
| T: | 500 | 1000 | 2500 | 5000 |
| Stopping Rule Tolerance: | | | | |
| 0.01 | 14 | 24 | 23 | 15 |
| 0.001 | 17 | 26 | 25 | 17 |
| 0.0001 | 36 | 59 | 62 | 39 |
| 0.00001 | 230 | 115 | 460 | 269 |
| Classifier: | | | | |
| k-means | 0.902 | 0.909 | 0.920 | 0.9196 |
| Local state decoding | 0.964 | 0.964 | 0.9696 | 0.9732 |

Numerical Experiment II

- The state-dependent distributions are Markov trees in which all conditional relationships are independent.



Figure: Markov trees for state 1 (left) and state 2 (right)

- The state-dependent distributions have densities supported on \mathbb{R}^5 given by

$$h_1(\mathbf{y}) = c_{1,12}(\Phi(y_1 - \mu_{1,1}), \Phi(y_2 - \mu_{1,2}) \mid \tau_{1,12}) \cdot c_{1,23}(\Phi(y_2 - \mu_{1,2}), \Phi(y_3 - \mu_{1,3}) \mid \tau_{1,23}) \\ \cdot c_{1,34}(\Phi(y_3 - \mu_{1,3}), \Phi(y_4 - \mu_{1,4}) \mid \tau_{1,34}) \cdot c_{1,45}(\Phi(y_4 - \mu_{1,4}), \Phi(y_5 - \mu_{1,5}) \mid \tau_{1,45}) \cdot \prod_{h=1}^5 \varphi(y_h - \mu_{1,h})$$

and

$$h_2(\mathbf{y}) = c_{2,14}(\Phi(y_1 - \mu_{2,1}), \Phi(y_4 - \mu_{2,4}) \mid \tau_{2,14}) \cdot c_{2,43}(\Phi(y_4 - \mu_{2,4}), \Phi(y_3 - \mu_{2,3}) \mid \tau_{2,43}) \\ \cdot c_{2,35}(\Phi(y_3 - \mu_{2,3}), \Phi(y_5 - \mu_{2,5}) \mid \tau_{2,35}) \cdot c_{2,52}(\Phi(y_5 - \mu_{2,5}), \Phi(y_2 - \mu_{2,2}) \mid \tau_{2,52}) \cdot \prod_{h=1}^5 \varphi(y_h - \mu_{2,h})$$

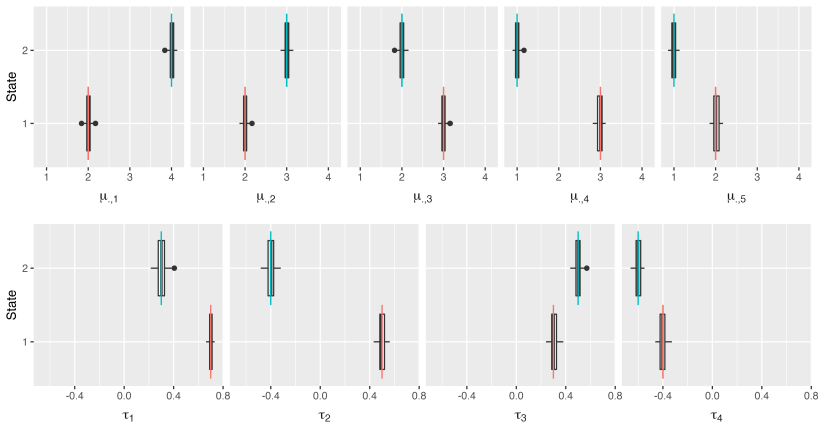


Figure: Parameter estimates based on 100 independent simulations and EIFM algorithm runs for the 5-dimensional 2-state HMMs

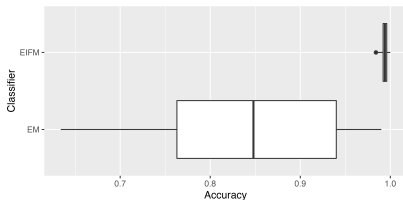


Figure: Accuracy across repetitions using initial independence model (bottom) and copula model (top)

Occupancy Data

| Classifier | Train | Test 1 |
|---------------------------------|-------|--------|
| k-means clustering | 0.865 | 0.818 |
| Independence copulas within HMM | 0.895 | 0.846 |
| BB7/Tawn copulas within HMM | 0.900 | 0.852 |

Latent Variable Models with Copulas

This project is currently developed in collaboration with



Robert Zimmerman

Latent Variables (LV)

- ▶ The variable of interest W is sometimes impossible to measure directly
 - ▶ State of the economy
 - ▶ Traffic in a city
 - ▶ State of your health
 - ▶ State of a complex disease

- ▶ Instead, one measures
 - ▶ $\mathbf{Y} = (Y_1, \dots, Y_k)^T$ whose components are surrogates of W and each provide partial information about W
 - ▶ Covariate $\mathbf{X} \in \mathbb{R}^p$

- ▶ We are often interested in the explanatory power of \mathbf{X} for W .

An example

- ▶ Cardiocotography (CTG) is a medical procedure that monitors the fetal heart rate.
- ▶ The LV is the fetus' underlying state of health during birth, W .
- ▶ Our surrogate response is the bivariate vector (Q, Y) where
 - ▶ Q is the number of peaks (acceleration followed by a deceleration of heart beats) for the signal recorded by the CTG
 - ▶ Y is the log of mean short-term "beat-to-beat" variability (MSTV) where the short-term variability (STV) is obtained by measuring the time between successive R waves (cardiac systoles) of the fetus' electrocardiogram.
- ▶ The covariates are FM (fetal movement) and UC (uterine contraction), two continuous variables monitored during birth.

Conditional independence LV model

- ▶ A canonical LV model:

$$Y_i \perp Q_i | W_i$$

$$Y_i | W_i \sim N(\mu_c + \lambda_c W_i, \sigma^2)$$

$$Q_i | W_i \sim \text{Poisson}(\exp(\mu_d + \lambda_d W_i))$$

$$W_i | X_i \sim N(X_i \beta, \eta^2)$$

- ▶ This implies that the two marginal regressions share a common random effect so they are marginally dependent (and conditionally independent)
- ▶ The induced dependence is not analytically available.

Conditional independence is a Copula LV

- ▶ The copula alternative is, conditional on W_i ,

$$H(Y_i, Q_i|W_i) = C_{\theta_i}(F_Y(Y_i|W_i), F_Q(Q_i|W_i)), \quad \theta_i = \kappa^{-1}(\xi_0 + \xi_1 W_i)$$

$$Y_i \sim N(\mu_c + \lambda_c W_i, \sigma^2); \quad Q_i \sim \text{Poisson}(\exp(\mu_d + \lambda_d W_i))$$

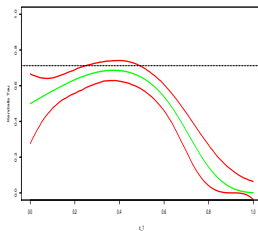
- ▶ The whole joint distribution of (Y, Q) is varying with W not just the marginals.
- ▶ The copula captures the residual dependence on W after the marginal effects have been accounted for.
- ▶ The previous model is obtained when the copula is the independence copula.

Why the Conditional Copula?

- ▶ $Y_i|x \sim N(f_i(x), \sigma_i) \quad x \in \mathbb{R}^2$
- ▶ True marginal means:
 - ▶ $f_1(x) = 0.6 \sin(5x_1) - 0.9 \sin(2x_2)$
 - ▶ $f_2(x) = 0.6 \sin(3x_1 + 5x_2)$
 - ▶ $\sigma_1 = \sigma_2 = 0.2, \mathbf{X}_1 \perp \mathbf{X}_2.$
- ▶ Copula : $\theta(x) = 0.71$
- ▶ Suppose x_2 is not observed so inference is based only on x_1

Why the Conditional Copula?

- ▶ $Y_i|x \sim N(f_i(x), \sigma_i) \quad x \in \mathbb{R}^2$
- ▶ True marginal means:
 - ▶ $f_1(x) = 0.6 \sin(5x_1) - 0.9 \sin(2x_2)$
 - ▶ $f_2(x) = 0.6 \sin(3x_1 + 5x_2)$
 - ▶ $\sigma_1 = \sigma_2 = 0.2, \mathbf{X}_1 \perp \mathbf{X}_2$.
- ▶ Copula : $\theta(x) = 0.71$
- ▶ Suppose x_2 is not observed so inference is based only on x_1



- ▶ Additional details in Levi and Craiu, 2018.

CTG: The LV Copula Model

- ▶ $(Q_i, Y_i) | W_i$ has joint density

$$f_{(Q,Y)}(q, y) = f_c(y) \cdot [C_{d|c}(F_d(q), F_c(y)) - C_{d|c}(F_d(q-), F_c(y))],$$

where

$$C_{d|c}(u_d, u_c) = \frac{\partial}{\partial u_c} C(u_d, u_c).$$

- ▶ Data Augmentation: Introduce latent variable Z such that

$$Q \stackrel{d}{=} F_d^-(F_Z(Z)),$$

- ▶ The copula between (Y, Z) is the same as the copula between (Y, Q)
- ▶ We can choose the distribution of Z to help the computation.
- ▶ For instance if we use a Gaussian copula, it helps to have $Z \sim N(0, 1)$
- ▶ Craiu and Sabeti (2012); Smith and Khaled (2012).

The Augmented LV Copula Model for the CTG Example

- ▶ The augmented model for CTG data is

$$W_i \sim \mathcal{N}(x_i\beta, 1),$$

$$Z_i \sim \mathcal{N}(0, 1),$$

$$Y_i | W_i \sim \mathcal{N}(\mu_c + \lambda_c W_i, \sigma^2)$$

$$(Z_i, Y_i) | W_i \sim C^{\text{Gauss}} \left(\Phi(Z_i), \Phi \left(\frac{Y_i - \mu_c - \lambda_c W_i}{\sigma} \right) \mid \theta(W_i, \xi) \right)$$

$$P(Q_i = q \mid W_i, Z_i) \propto \mathbb{1}_{F_Z^{-1}(F_d(q - |\mu_d, \lambda_d, W_i)) \leq z < F_Z^{-1}(F_d(q | \mu_d, \lambda_d, W_i))}.$$

CTG: The Augmented LV Copula Model

- ▶ Let $\xi = (\xi_0, \xi_1) \in \mathbb{R}^2$ and $A(w) = \xi_0 + \xi_1 \cdot w$. Then we set

$$\theta(w, \xi) = \frac{e^{A(w)} - e^{-A(w)}}{e^{A(w)} - e^{-A(w)}}$$

as the correlation parameter of the bivariate Gaussian conditional copula of $(Y, Z) | W = w$.

- ▶ Parameters are a priori independent.

Model Selection: WAIC

- ▶ The WAIC is defined as

$$\text{WAIC}(\mathcal{M}) = -2\text{fit}(\mathcal{M}) + 2\text{p}(\mathcal{M}), \quad (3)$$

where the model fitness is

$$\text{fit}(\mathcal{M}) = \sum_{i=1}^n \log(\mathbb{E}[\Pr(y_i, \mathbf{q}_i | \omega, \mathcal{M})]) \quad (4)$$

and the penalty

$$\text{p}(\mathcal{M}) = \sum_{i=1}^n \text{Var}(\log(\Pr(y_i, \mathbf{q}_i | \omega, \mathcal{M}))), \quad (5)$$

where ω contains all the parameters and latent variables in the model.

Spotlight on dependence: A conditional WAIC

- ▶ We use the following two conditional WAICs (Levi and Craiu, 2018)

$$\begin{aligned} \text{CWAIC}_{Y|Q}(\mathcal{M}) &= -2 \sum_{i=1}^n \log(\mathbb{E}[\Pr(y_i|q_i, \omega, \mathcal{M})]) + \\ &\quad + 2 \sum_{i=1}^n \text{Var}(\log(\Pr(y_i|q_i, \omega, \mathcal{M}))), \end{aligned}$$

$$\begin{aligned} \text{CWAIC}_{Q|Y}(\mathcal{M}) &= -2 \sum_{i=1}^n \log(\mathbb{E}[\Pr(q_i|y_i, \omega, \mathcal{M})]) + \\ &\quad + 2 \sum_{i=1}^n \text{Var}(\log(\Pr(q_i|y_i, \omega, \mathcal{M}))), \end{aligned}$$

- ▶ $\frac{1}{2}(\text{CWAIC}_{1|2} + \text{CWAIC}_{2|1})$ is asymptotically equivalent to CCV for the marginal likelihood

$$\text{CCV}(\mathcal{M}) = \frac{1}{2} \left\{ \sum_{i=1}^n \log(\Pr(y_i|q_i, \mathcal{D}_{-i}, \mathcal{M})) + \sum_{i=1}^n \log(\Pr(q_i|y_i, \mathcal{D}_{-i}, \mathcal{M})) \right\}.$$

Simulation Experiment

- ▶ Generate data using a Gaussian copula

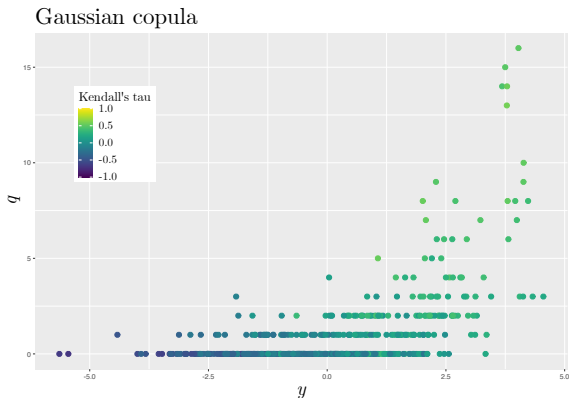


Figure: Bivariate scatterplot of the generated data with Gaussian copula, and Poisson and normal marginals

Simulation Experiment

- CWAIC $_{Y|Q}$ and CWAIC $_{Q|Y}$ selection criteria

| Criteria\Copula | Gaussian | Frank | Gumbel | Clayton | Indep |
|-----------------|----------------|---------|---------|---------|----------------|
| CWAIC $_{Y Q}$ | 1627.36 | 1642.36 | 2395.17 | 1637.17 | 1606.31 |
| CWAIC $_{Q Y}$ | 950.71 | 982.42 | 1673.57 | 976.05 | 997.43 |
| Average | 1289.04 | 1312.39 | 2034.37 | 1306.61 | 1301.87 |

Simulation Experiment

Gaussian copula

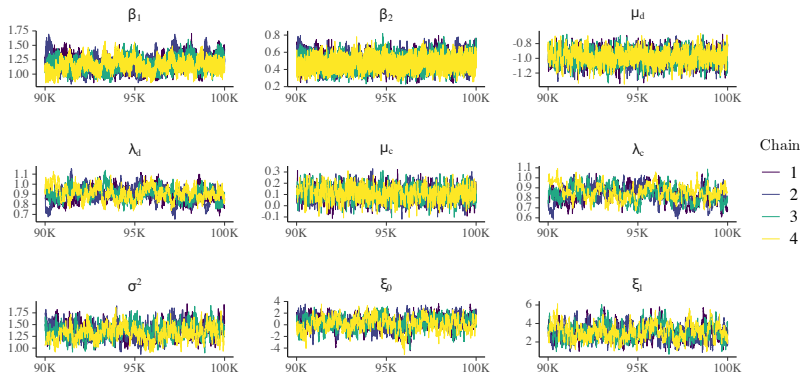
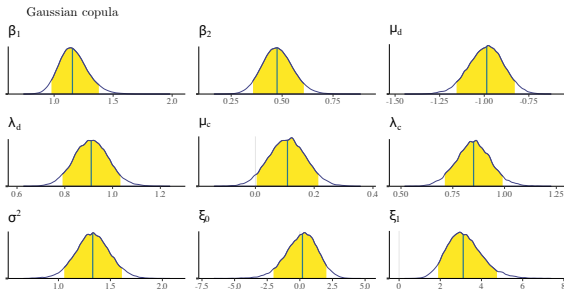


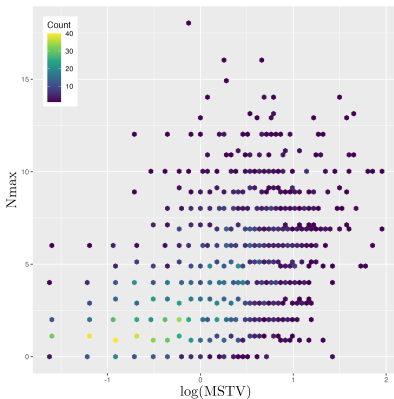
Figure: Traceplots for η 's components.

Simulation Experiment



| | β_1 | β_2 | λ_d | λ_c | ξ_1 |
|------|-----------|-----------|-------------|-------------|---------|
| Mean | 1.18 | 0.48 | 0.90 | 0.84 | 3.10 |
| True | 1 | 0.5 | 1 | 1 | 3 |

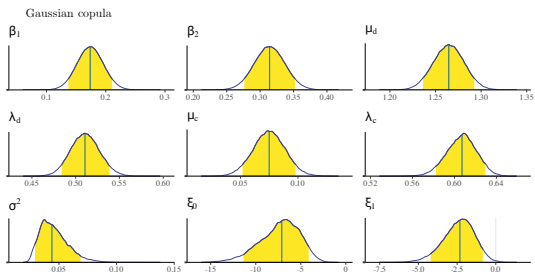
CTG: The data



CTG: Estimates

- ▶ $WAIC$, $WAIC_{Y|Q}$ and $WAIC_{Q|Y}$ all point to the Gaussian copula (over Gumbel, Frank, Clayton, Independence).
- ▶ The posterior means

| | β_1 (FM) | β_2 (UC) | λ_d | λ_c | ξ_1 |
|------|----------------|----------------|-------------|-------------|---------|
| Mean | 0.1744 | 0.3147 | 0.5101 | 0.6038 | -2.3401 |



The Past & Future

- ▶ Copulas offer a way to bypass the paucity of available joint distributions.
- ▶ Copulas allow the integration of multiple (dependent) sources of information/data via joint modeling
- ▶ Joint models can be used for prediction/imputation of an expensive variable given values for cheaper ones.
- ▶ So far have been used to further empower multivariate regressions, time series, HMMs, LV models, etc
- ▶ Computational challenges, especially in higher dimensions
- ▶ Papers available here: raducraiu.com/publications/