Summary of *Extending the Rank Likelihood for Semiparametric Copula Estimation*, by Peter Hoff

David Gerard

Department of Statistics University of Washington gerard2@uw.edu

May 2, 2013

- INC Respondent's Income
- DEG Highest degree obtained
- CHILD Number of children
- PINC Parent's income when respondent was 16
- PDEG Max(mother's degree, father's degree, na.rm = T)
- PCHILD Number of siblings + 1
- AGE Age of respondent

$INC_{i} = \beta_{0} + \beta_{1}CHILD_{i} + \beta_{2}DEG_{i} + \beta_{3}AGE_{i} + \beta_{4}PCHILD_{i} + \beta_{5}PINC_{i} + \beta_{6}PDEG_{i} + \epsilon_{i}$

or

 $CHILD_i \sim Pois(exp\{\beta_0 + \beta_1 INC_i + \beta_2 DEG_i + \beta_3 AGE_i + \beta_4 PCHILD_i + \beta_5 PINC_i + \beta_6 PDEG_i\})$

Response	INC	CHILD	DEG
INC	NA	1.103(0.112)	7.025(<0.001)
CHILD	0.005(<mark>0.00922</mark>)	NA	-0.068(0.056)

Response	AGE	PCHILD	PINC	PDEG
INC	0.335(<0.001)	0.284(0.407)	4.070(0.001)	1.399(0.115)
CHILD	0.037(<0.001)	0.021(0.080)	-0.063(0.195)	-0.051(0.204)

э

・ロト ・ 日 ト ・ 日 ト ・

Response	INC	CHILD	DEG
INC	NA	1.103(0.112)	7.025(<0.001)
CHILD	0.005(<mark>0.00922</mark>)	NA	-0.068(0.056)

Response	AGE	PCHILD	PINC	PDEG
INC	0.335(<0.001)	0.284(0.407)	4.070(0.001)	1.399(0.115)
CHILD	0.037(<0.001)	0.021(0.080)	-0.063(0.195)	-0.051(0.204)

Using a sandwich estimator we have:

Response	INC	CHILD	
INC	NA	1.103(<mark>0.108</mark>)	
CHILD	0.005(<mark>0.00939</mark>)	NA	

< ∃ > <

Let $y_{i,j}$ be the value of the j^{th} variable taken on by the i^{th} observational unit.

Gaussian Copula Sampling Model

$$m{z_1},\ldots,m{z_n}|m{C}\sim \textit{i.i.d.}$$
 multivariate normal $(m{0},m{C})$
 $y_{i,j}=F_j^{-1}[\Phi(z_{i,j})]$

- Estimate association parameters without having to estimate marginals.
- Do this by only using the partial ordering induced by the data: $y_{k,j} < y_{i,j} \Rightarrow z_{k,j} < z_{i,j}$.
- A partial ordering is a total ordering without the *totality* condition (i.e., some elements may be incomparable).

Let *D* be the set of $\mathbf{Z} := (z_{i,j})$'s that satisfy the partial ordering. Use the following "marginal likelihood" for inference:

$$Pr(\mathbf{Z} \in D | \mathbf{C}, F_1, \dots, F_p) = \int_D p(\mathbf{Z} | \mathbf{C}) d\mathbf{Z} = Pr(\mathbf{Z} \in D)$$

Full conditionals of the $z_{i,j}$'s are easy to derive, so we can implement a Gibb's sampler (using covariance matrix rather than correlation matrix, but doesn't matter for estimation).

- Sample z_{i,j} |Z_[-i,-j], V from a truncated normal with the bounds set by the partial ordering and the conditional mean and variance found in the usual way: σ_j² = V_[j,j] - V_[j,-j]V⁻¹_[-j,-j]V_[-j,j] and μ = V_[j,-j]V_[-j,-j]Z^T_[i,-j].
- Sample V from an inverse-Wishart distribution (if you use the conjugate prior).

3 Let
$$C_{[i,j]} = V_{[i,j]} / \sqrt{V_{[i,i]} V_{[j,j]}}$$



95% Posterior Credible Intervals



Visualizing Conditional Dependencies of GSS data.



Figure: Reduced conditional dependence graph for the General Social Survey data.

- We can also sample from a posterior predictive distribution to do inference on y's.
- Compare to empirical distributions for model checking.

• sample
$$\mathbf{C} \sim p(\mathbf{C} | \mathbf{Z} \in D)$$
;

- **2** sample $z \sim multivariate normal(0, C);$
- 3 set $y_j = \hat{F}_j^{-1}(\Phi(z_j))$.

Posterior Predictive and Empirical distributions of DEG given PINC



David Gerard (UW)

Posterior Predictive and Empirical distributions of INC given DEG and PINC



- The paper proves that ranks are "G-sufficient" and "L-Sufficient" when we have continuous marginals.
- However, when the data are discrete the partial ordering does not have either of these properties.

Definition

A collection $\mathcal G$ of 1-1 transformations of $\mathcal X$ (the sample space) is a group if

- $\bullet \quad \text{For all } g_1,g_2\in \mathcal{G},\ g_1g_2\in \mathcal{G}$
- 2 For all $g \in \mathcal{G}, g^{-1} \in \mathcal{G}$.

Definition

Two points $x_1, x_2 \in \mathcal{X}$ are equivalent if there exists a $g \in \mathcal{G}$ such that $x_1 = gx_2$. The sets of equivalent points are the orbits of \mathcal{G} .

Definition

A function M is said to be maximally invariant if it is in 1-1 correspondence with the orbits of \mathcal{G} . i.e. M(gx) = M(x) for all $g \in \mathcal{G}$ and $M(x_1) = M(x_2) \Rightarrow x_2 = gx_1$ for some $g \in \mathcal{G}$.

<ロト </p>

- ${\cal G}$ induces a group $\bar{{\cal G}}$ over the parameter space.
- Let $\overline{\mathcal{G}} = \{\overline{g} : \Omega \to \Omega \text{ s.t. } \overline{g}\theta = \theta' \text{ if } X \sim P_{\theta} \text{ and } gX \sim P_{\theta'}\}.$
- There are also maximally invariant parameters for $\bar{\mathcal{G}}$.

- Under a group of transformations, *G*, the maximally invariant statistic is called "G-sufficient" for the maximally invariant parameter.
- The ranks are the maximally invariant statistics under the group of continuous strictly increasing functions [Lehmann and Romano, 2005, pp 215 - 216]. You can also put the correlation matrix in a 1-1 correspondence with the induced group \$\bar{\mathcal{G}}\$'s orbits (but only if the marginals are continuous).
- Hence, the ranks are "G-sufficient" for the correlation matrix.
- Intuition: if we assume the marginals are unknown, then applying strictly increasing continuous functions to the data should not change the estimation problem

Simple Example

- Let $X_1, \ldots, X_n \sim i.i.d.N(\mu, \sigma^2)$, $\Omega = \{(\mu, \sigma) : \mu \in \mathbb{R}, \sigma^2 > 0\}$ and $\mathcal{X} = \mathbb{R}$.
- Let $\mathcal{G} = \{g : g(\mathbf{x}) = \mathbf{x} + a\mathbf{1}, a \in \mathbb{R}\}$
- Then $g(\mathbf{X}) = (X_1 + a, \dots, X_n + a)$ s.t. $X_i \sim \text{i.i.d. } N(\mu + a, \sigma^2)$ and $\bar{\mathcal{G}} = \left\{ \bar{g} : \bar{g} \begin{pmatrix} \mu \\ \sigma \end{pmatrix} = \begin{pmatrix} \mu + a \\ \sigma \end{pmatrix} \right\}.$
- Orbits of Ω under $\bar{\mathcal{G}}$ are defined by $\sigma.$ So σ is the maximally invariant parameter.
- Orbits of \mathcal{X} under \mathcal{G} are defined by the differences from the mean $(X_1 \overline{X}, \dots, X_n \overline{X})$, so the differences are G-sufficient for σ and inference should only be based on the differences.
- Can extend these ideas to *minimal* G-sufficient ("maximal invariant reduction of minimal sufficient statistics") [Barnard, 1963]. Here, *s*².

- For {F₁,..., F_p} ∈ F the marginal distributions, a statistic t(Y) is L-sufficient for C if
 - $t(\mathbf{Y}_0) = t(\mathbf{Y}_1) \Rightarrow \sup_{\{F_1,\ldots,F_p\}\in\mathcal{F}} p(\mathbf{Y}_0|\mathbf{C},F_1,\ldots,F_p) = \sup_{\{F_1,\ldots,F_p\}\in\mathcal{F}} p(\mathbf{Y}_1|\mathbf{C},F_1,\ldots,F_p); \text{ and}$ $p(t(\mathbf{Y})|\mathbf{C},F_1,\ldots,F_p) = p(t(\mathbf{Y})|\mathbf{C})$
- If there are no nuisance parameters, L-sufficiency becomes "full sufficiency".
- Intuition: partition generated by L-sufficient statistic is "at least as fine" as the partition generated by the M.L.E. of C [Rémon, 1984] (i.e. MLE is function of ranks alone).

- Develop maximum likelihood approach using the extended rank likelihood.
 - Lots of trouble with this. [Pettitt, 1982] used an approximation for the rank likelihood in the standard multivariate normal setting.
- Run simulation study against the competitors [Genest et al., 1995] (just plugged in ECDF); [Olsson, 1979] (estimate threshholding values when number of categories is known)
- Develop technique for other copulas, then simulate and see what info is lost by using the wrong copula.

References



Barnard, G. A. (1963).

Some logical aspects of the fiducial argument. Journal of the Royal Statistical Society. Series B (Methodological), pages 111–114.



Genest, C., Ghoudi, K., and Rivest, L.-P. (1995).

A semiparametric estimation procedure of dependence parameters in multivariate families of distributions. *Biometrika*, 82(3):543–552.



Hoff, P. D. (2007).

Extending the rank likelihood for semiparametric copula estimation. *The Annals of Applied Statistics*, pages 265–283.



Lehmann, E. E. L. and Romano, J. P. (2005).

Testing statistical hypotheses. Springer Science+ Business Media.



Olsson, U. (1979).

Maximum likelihood estimation of the polychoric correlation coefficient. *Psychometrika*, 44(4):443–460.



Pettitt, A. (1982).

Inference for the linear model using a likelihood based on ranks. Journal of the Royal Statistical Society. Series B (Methodological), pages 234–243.



Rémon, M. (1984).

On a concept of partial sufficiency: L-sufficiency. International Statistical Review/Revue Internationale de Statistique, pages 127–135.

(日) (同) (三) (三)