

コピュラに基づくダイバージェンスを用いた適合度検定

加藤 昇吾 統計基盤数理研究系 准教授

はじめに

コピュラ

d 次元コピュラとは、それぞれの1次元周辺分布が $[0, 1]$ 上の一様分布となる d 次元分布関数のことをいう。

d 次元コピュラを変換することにより、任意の1次元周辺分布関数を持つ d 次元分布関数を得ることが可能（スクラーの定理）。

⇒ コピュラにより、柔軟なモデリングが可能となる。

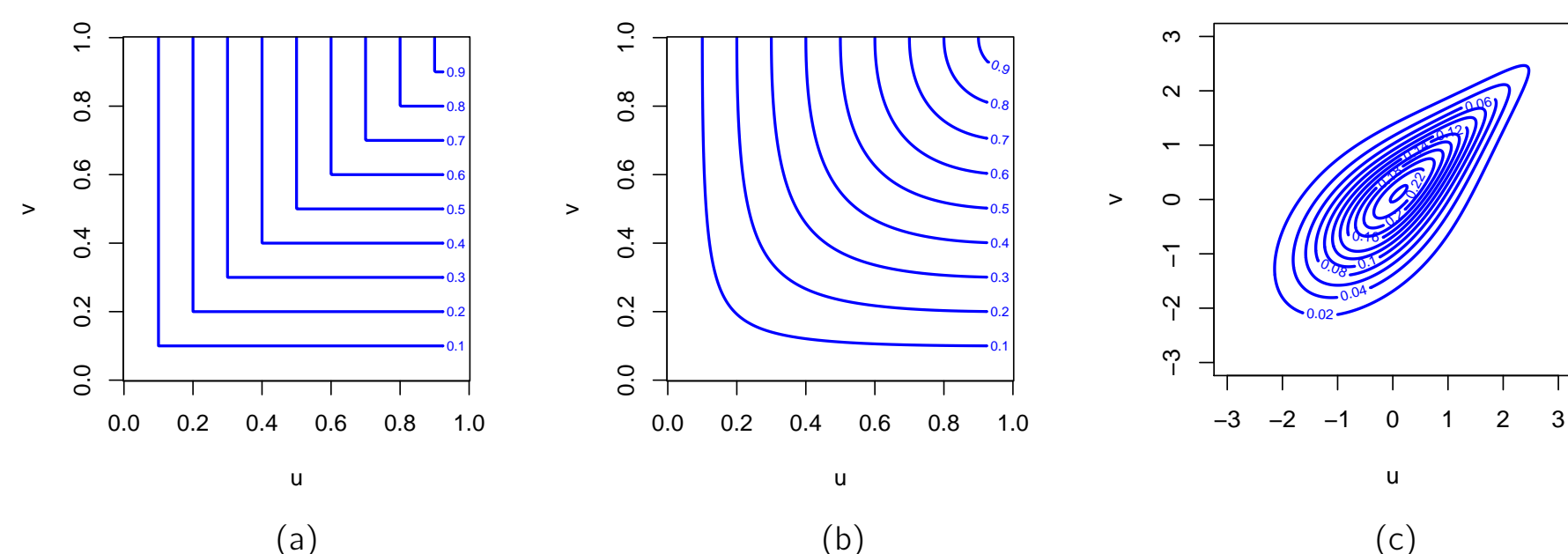


図1. (a),(b): 2次元コピュラの等高線プロットの例、
(c): 標準正規分布を周辺分布に持つコピュラの確率密度関数の等高線プロット。

研究の目的

コピュラのための適合度検定としては、Cramér-von Mises距離に基づく検定（Cramér-von Mises検定）がしばしば用いられる。

⇒ この検定は裾部分での感度が低い問題が知られている。

本研究では、裾部分での重みを調節可能な新たなダイバージェンスを定義し、それに基づく適合度検定を提案する。

なお本研究は、江口真透名誉教授（統計数理研究所）・吉羽要直教授（東京都立大学）との共同研究である。

コピュラに基づくダイバージェンスと適合度検定への応用

コピュラに基づくダイバージェンス

確率密度関数に基づく β ダイバージェンス(Basu et al., 1998; Minami & Eguchi, 2002)から着想を得て、以下のダイバージェンスを定義する。

定義 C_0, C_1 : d 次元コピュラとする。

このとき、 C_0 と C_1 に対して、新しいダイバージェンスを以下で定義する：

$$D_\beta(C_0, C_1) = \int_{\mathcal{H}} \left\{ \frac{C_0(\mathbf{u})^{\beta+1}}{\beta(\beta+1)} + \frac{C_1(\mathbf{u})^{\beta+1}}{\beta+1} - \frac{C_0(\mathbf{u})C_1(\mathbf{u})^\beta}{\beta} \right\} dW(\mathbf{u}), \quad (1)$$

ただし、 W : \mathcal{H} を台 (support) とする d 次元コピュラ測度、
 $\beta \in \mathbb{R} \setminus \{-1, 0\}$: チューニングパラメータである。

また、 $\beta = 0, -1$ のときは、以下のように定義する。

(i) $\beta = 0$ (累積Kullback-Leiblerダイバージェンス) :

$$D_0(C_0, C_1) = \lim_{\beta \rightarrow 0} D_\beta(C_0, C_1) = - \int_{\mathcal{H}} \left\{ C_0(\mathbf{u}) \log \left(\frac{C_1(\mathbf{u})}{C_0(\mathbf{u})} \right) + C_0(\mathbf{u}) - C_1(\mathbf{u}) \right\} dW(\mathbf{u}),$$

(ii) $\beta = -1$ (累積Itakura-Saitoダイバージェンス) :

$$D_{-1}(C_0, C_1) = \lim_{\beta \rightarrow -1} D_\beta(C_0, C_1) = \int_{\mathcal{H}} \left\{ -\log \left(\frac{C_1(\mathbf{u})}{C_0(\mathbf{u})} \right) + \frac{C_1(\mathbf{u})}{C_0(\mathbf{u})} - 1 \right\} dW(\mathbf{u}).$$

特別な場合 $\beta = 1$ のとき、提案したダイバージェンス(1)は、Cramér-von Mises距離と等しくなる：

$$D_1(C_0, C_1) = \frac{1}{2} \int_{\mathcal{H}} \{C_0(\mathbf{u}) - C_1(\mathbf{u})\}^2 dW(\mathbf{u}).$$

ダイバージェンスの性質 ダイバージェンス(1)は、以下を満たす：

- (i) 任意の $C_0, C_1 \in \mathcal{C}$ に対して、 $D_\beta(C_0, C_1) \geq 0$ 、
(ii) $D_\beta(C_0, C_1) = 0 \iff C_0(\mathbf{u}) = C_1(\mathbf{u})$ for all $\mathbf{u} \in \mathcal{H}$.

ただし、 \mathcal{C} は全ての d 次元コピュラの集合を表す。

適合度検定への応用

定義 $U_1, \dots, U_n \sim i.i.d.$ コピュラ G , C_0 : コピュラ, とする。

以下の仮説検定を考える：

$$H_0: G = C_0 \quad \text{vs.} \quad H_1: G \neq C_0.$$

この仮説検定のための検定統計量を以下で定義する：

$$T = 2nD_\beta(\hat{C}, C_0).$$

ここで、 \hat{C} は (U_1, \dots, U_n) の経験分布関数を表す。 T が大きいときに H_0 を棄却し、小さいときに H_0 を採択する。

検定統計量 T の性質

- $\beta = 1 \implies$ Cramér-von Mises検定量。
- 累積Kullback-Leiblerダイバージェンス($\beta = 0$)や累積Itakura-Saitoダイバージェンス($\beta = -1$)に基づく検定量も特別な場合として含む。

検定統計量 T の漸近分布 以下を仮定する：

$W(\mathbf{u})$: $\{\mathbf{v}_1, \dots, \mathbf{v}_k\} \subset \mathcal{H}$ を台とする分布関数、つまり、

$$W(\mathbf{u}) = \frac{1}{k} \sum_{j=1}^k \mathcal{I}(\mathbf{u} \leq \mathbf{v}_j).$$

ただし、 \mathcal{I} は定義関数を表す。このとき、帰無仮説 H_0 の下で、以下が検定統計量 T に対して成り立つ：

$$T \xrightarrow{d} \frac{1}{k} \sum_{j=1}^k \lambda_j \chi_j^2 \quad (n \rightarrow \infty).$$

ここに、 $\chi_1^2, \dots, \chi_k^2 \sim i.i.d.$ 自由度1の χ^2 分布、

$\lambda_1, \dots, \lambda_k$: $\Sigma_T = (\sigma_T(\mathbf{v}_i, \mathbf{v}_j))_{1 \leq i, j \leq k}$ の固有値、

$$\sigma_T(\mathbf{v}_i, \mathbf{v}_j) = \{C_0(\mathbf{v}_i)C_0(\mathbf{v}_j)\}^{(\beta-1)/2} \{C_0(\mathbf{v}_i \wedge \mathbf{v}_j) - C_0(\mathbf{v}_i)C_0(\mathbf{v}_j)\}.$$

シミュレーション

$U_1, \dots, U_n \sim i.i.d.$ パラメータ θ のFrankコピュラ。

$$H_0: \theta = 10 \quad \text{vs.} \quad H_1: \theta \neq 10.$$

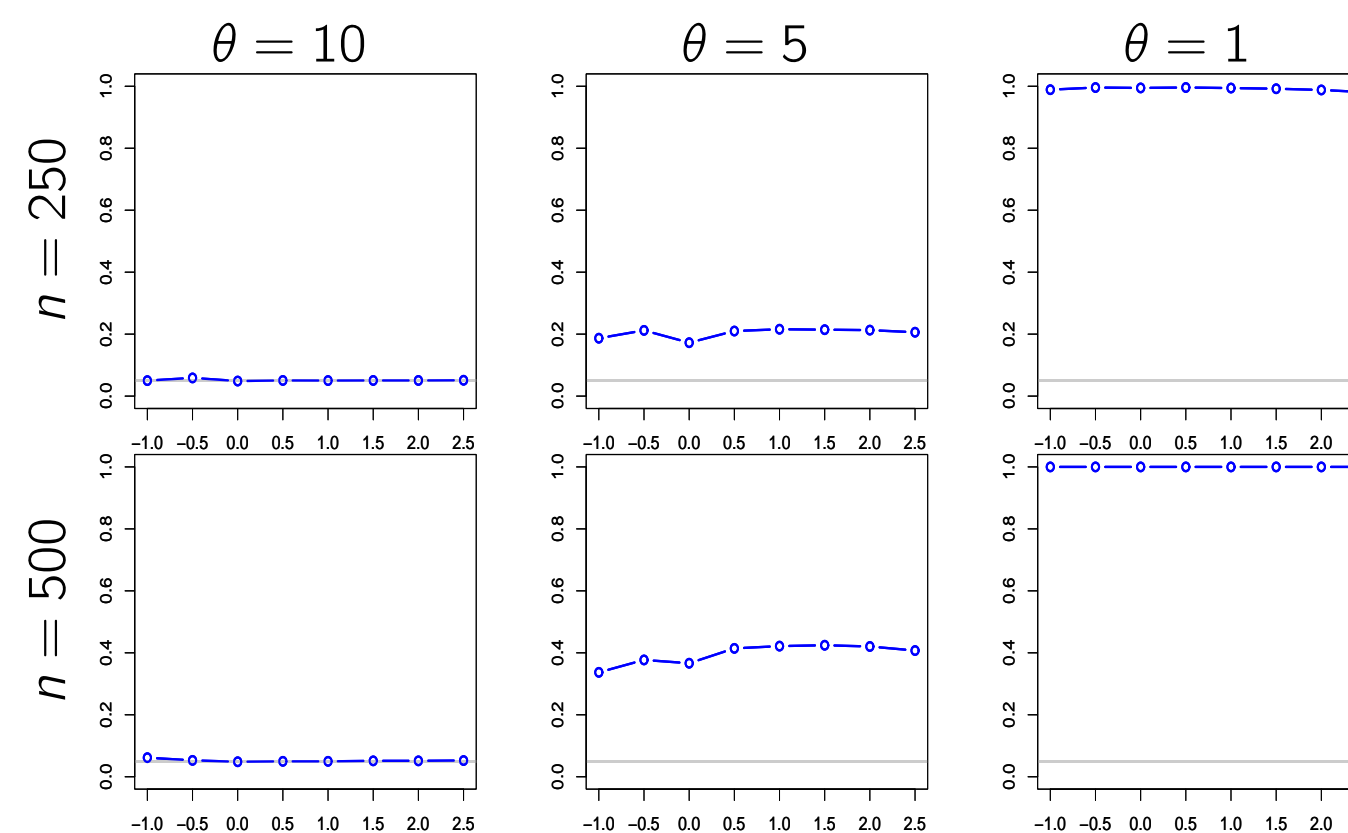


図2. 検定統計量 T に基づく検定の検出力のプロット。有意水準は0.05, $k = 100$ とし、検出力は10000個のサンプルから計算されている。(x軸は β , y軸は検出力を表す。)