

Density Power Divergenceを用いた ロバストな情報量規準の提案とその応用

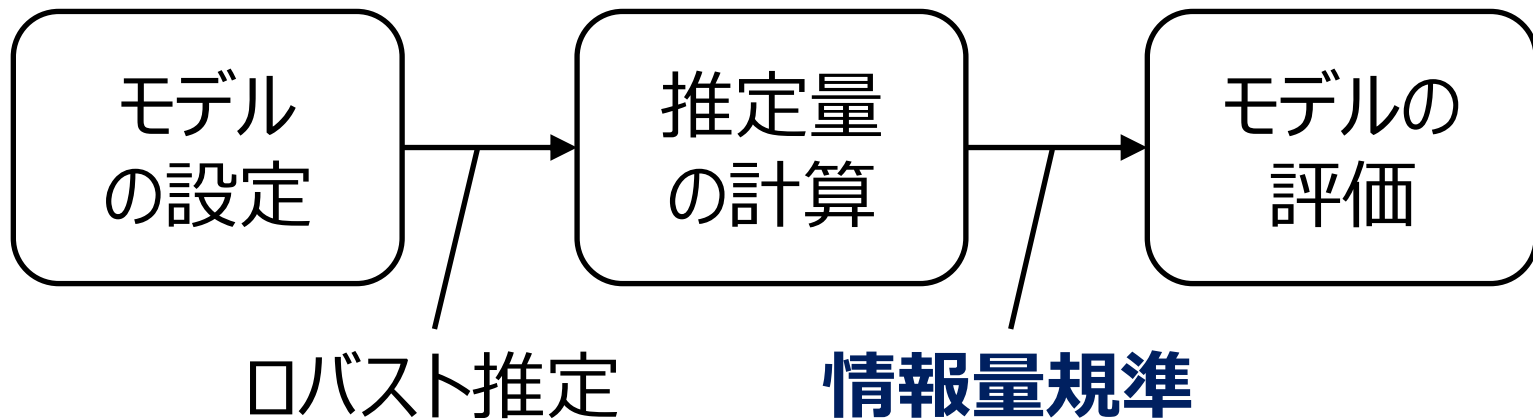
大阪大学大学院基礎工学研究科

奥野 彰文

下平 英寿

研究の目的

外れ値を含むデータを分析する



ロバストな情報量規準を導出する

Density Power Divergence (Basu et al. (1998))を用いる.

1. DPDと β -推定の紹介

Kullback Leibler Divergence (KL)

$$d(q, p_\theta) := - \int q(x) \log p(x | \theta) dx$$

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta \in \Theta} d(\hat{q}, p_\theta)$$

Density Power Divergence (DPD)

(Basu et al.(1998). 別名 β -divergence)

$$d_\beta(q, p_\theta) := - \int q(x) \frac{p(x | \theta)^\beta - 1}{\beta} dx + b_{1+\beta}(\theta) \quad (\beta > 0)$$

β -推定量 $\hat{\theta}_\beta := \arg \min_{\theta \in \Theta} d_\beta(\hat{q}, p_\theta)$

$\hat{\theta}_\beta \rightarrow \hat{\theta}_{\text{MLE}}, (\beta \downarrow 0)$ β -推定量はM-推定量である

$$\frac{\partial d_\beta(\hat{q}, p_\theta)}{\partial \theta} = -\frac{1}{n} \sum_{t=1}^n \overset{\text{重み}}{\boxed{p(x_t | \theta)^\beta}} \overset{\text{スコア関数}}{\boxed{\frac{\partial \log p(x_t | \theta)}{\partial \theta}}} + \frac{\partial b_{1+\beta}(\theta)}{\partial \theta} = 0$$

β が

- 小さい \Rightarrow (データの)有効性が高い
- 大きい \Rightarrow ロバスト

β は通常 0~1 の値を取るが、
 外れ値によっては $\beta=2$ 程度が必要となる。

2. IC_B の提案

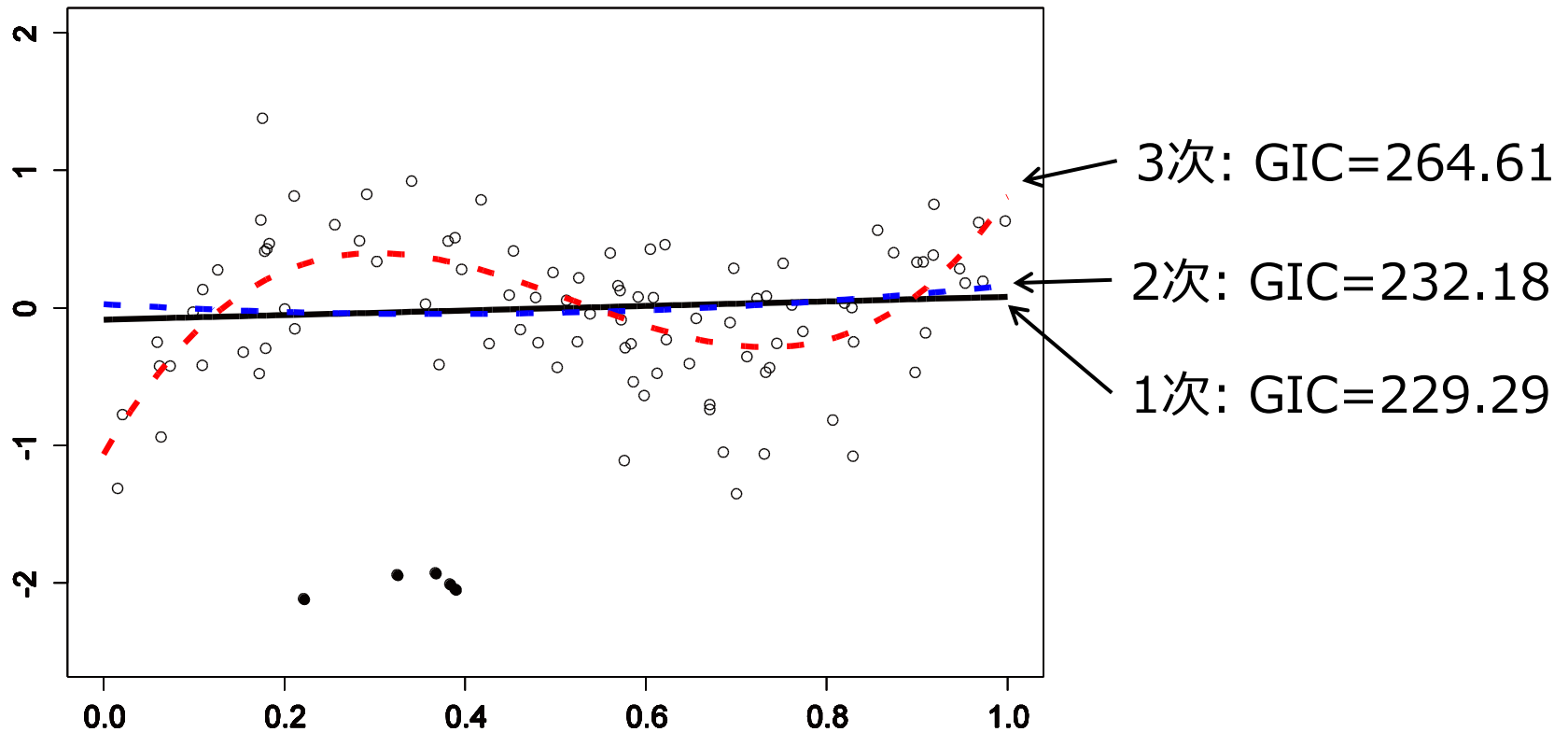
β -推定はM-推定である

GICはM-推定されたモデルを評価できる。

$GIC_{\beta} \iff d(q, \hat{p}_{\beta})$ β -推定量をKLで評価

GICは外れ値に弱い

ロバスト推定された1次,2次,3次のモデルをGICで評価



GICは**1次**のモデルを選択する。(当てはまりが良くない) 8

Konishi and Kitagawa(1996) $\text{GIC}_\beta \iff d(q, \hat{p}_\beta)$

KL情報量はロバストでない

Ronchetti(1982) $\text{AICR}_\beta \iff d_\beta(q, \hat{p}_\beta)$

β に依存したモデル評価

本研究 $\text{IC}_B \iff d_B(q, \hat{p}_\beta)$

ロバスト かつ β に依存しない モデル評価

今回提案する情報量規準 IC_B

$B > 0$: fixed,

$$IC_B := 2n d_B(\hat{q}, \hat{p}_\beta) + 2\text{tr}G_{\beta,B}H_\beta^{-1}$$

$$\rho_\beta(x | \theta) := -\frac{1}{\beta}p(x | \theta)^\beta + \frac{1}{1 + \beta} \int p(x | \theta)^{1+\beta} dx$$

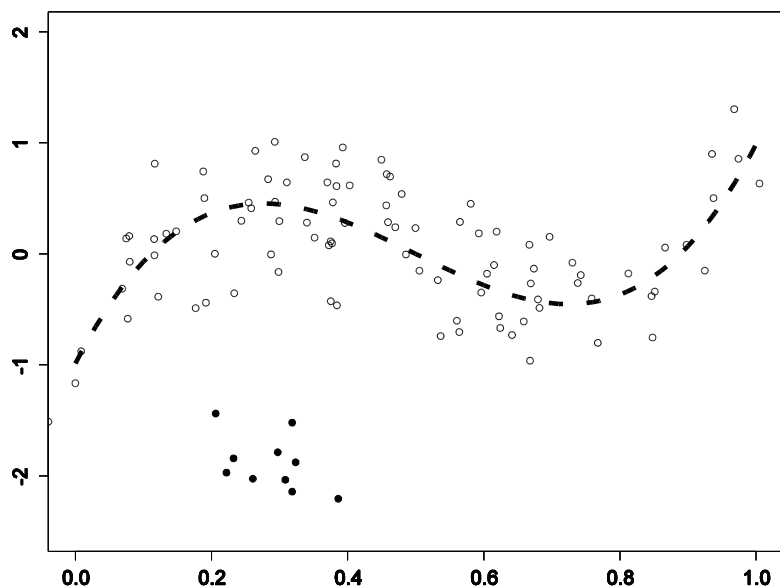
$$G_{\beta,B} := -\frac{1}{n} \sum_{t=1}^n \frac{\partial \rho_\beta(x_t | \theta)}{\partial \theta} \bigg|_{\theta=\hat{\theta}_\beta} \frac{\partial \rho_B(x_t | \theta)}{\partial \theta'} \bigg|_{\theta=\hat{\theta}_\beta}$$

$$H_\beta := -\frac{1}{n} \sum_{t=1}^n \frac{\partial^2 \rho_\beta(x_t | \theta)}{\partial \theta \partial \theta'} \bigg|_{\theta=\hat{\theta}_\beta}$$

3. 数値実験 と まとめ

多項式回帰によるモデル選択(平均予測二乗誤差)

真のモデル: 3次



モデルの候補

1. $y = a_0 + a_1x$
2. $y = a_0 + a_1x + a_2x^2$
- ⋮
5. $y = a_0 + a_1x + \dots + a_5x^5$

GIC, IC_B は $\beta = 0.1, 0.3, \dots, 1.9$ も選択する

結果

	既存の方法				新しい
	GIC_{β}	$AICR_{0.01}$	$AICR_{0.5}$	$AICR_{1.0}$	IC_B
平均予測 二乗誤差	0.1758	0.1951	0.0471	0.0518	0.0464

ただし $B=1.0$ とし, PMSEは1000回の平均を取った

まとめ

今回提案した IC_B は

ロバスト性 を持ち

β の選択 を含めて

良い予測を示す.

今後の展望:

DPDは外れ値が多いとバイアスが大きい

$q(x) := (1 - \varepsilon)f(x) + \varepsilon\xi(x)$ とすると
真の分布

$$\arg \min_{p \in \mathcal{D}} d_{\beta}(q; p) \approx \arg \min_{p \in \mathcal{D}} d_{\beta}((1 - \varepsilon)f, p) \neq f$$

Gamma-divergence (Fujisawa and Eguchi 2008) に置き換える

$$\arg \min_{p \in \mathcal{D}} d_{\gamma}(q; p) \approx \arg \min_{p \in \mathcal{D}} d_{\gamma}((1 - \varepsilon)f, p) \approx f$$

参考文献

1. Ronchetti, Elvezio. (1982). Robust alternatives to the F-test for the linear model. *Springer Netherlands*.
2. Konishi, Sadanori. and Genshiro, Kitagawa. (1996). Generalised information criteria in model selection. *Biometrika* **83,4** 875-890.
3. Shimodaira, H. (2000). Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*. **90,2** 227–244.
4. Basu, A. et al. (1998). Robust and efficient estimation by minimising a density power divergence. *Biometrika*. **85,3** 549-559
5. Mattheou, Kyriacos, and Alex Karagrignoriou. (2006). Robust model selection criteria with applications.
6. Fujisawa, H. and Eguchi, S. (2008). Robust parameter estimation with a small bias against heavy contamination. *J. Stat. Plan. And Inf.*
7. Yata, Kazuyoshi. et al. (2011). Note on robust model selection by density power divergence in a contaminated regression model. *数理解析研究所講究録*. **1758** 150-159

補足スライド

β -推定量は反復的に計算できる:

① $\hat{\theta}^{(0)} := \hat{\theta}_{\text{MLE}}$ とする.

$$\begin{aligned} \text{② } E_{\hat{q}}[p(X_t | \hat{\theta}^{(i)})^\beta S(X_t, \theta)] \\ = E_{p_\theta}[p(X | \theta)^\beta S(X, \theta)] \Big|_{\theta = \hat{\theta}^{(i)}} \end{aligned}$$

の解を $\hat{\theta}^{(i+1)}$ と定義する.

③ $\|\theta^{(i+1)} - \theta^{(i)}\| < \varepsilon$ となるまで②を繰り返す

重回帰での反復計算の例

$X_i \stackrel{\text{i.i.d.}}{\sim} U_d(0, 1)$ かつ $Y_i \stackrel{\text{i.i.d.}}{\sim} N(X_i^T \beta, 1^2)$ を仮定する.

確率モデルを $p(\{\mathbf{x}, y\}, \beta) := 1_{[0,1]^d}(\mathbf{x})N(y; \mathbf{x}^T \beta, 1^2)$

とすると、パラメータの更新式は陽的に解けて、

$$\beta^{(i+1)} = (X^T W^{(i)} X)^{-1} X^T W^{(i)} y$$

ただし $W^{(i)} = (p(x_1 | \theta^{(i)}), \dots, p(x_N | \theta^{(i)}))$

他の情報量規準との関係

$$\text{IC}_B = \text{AICR}_\beta, (B = \beta)$$

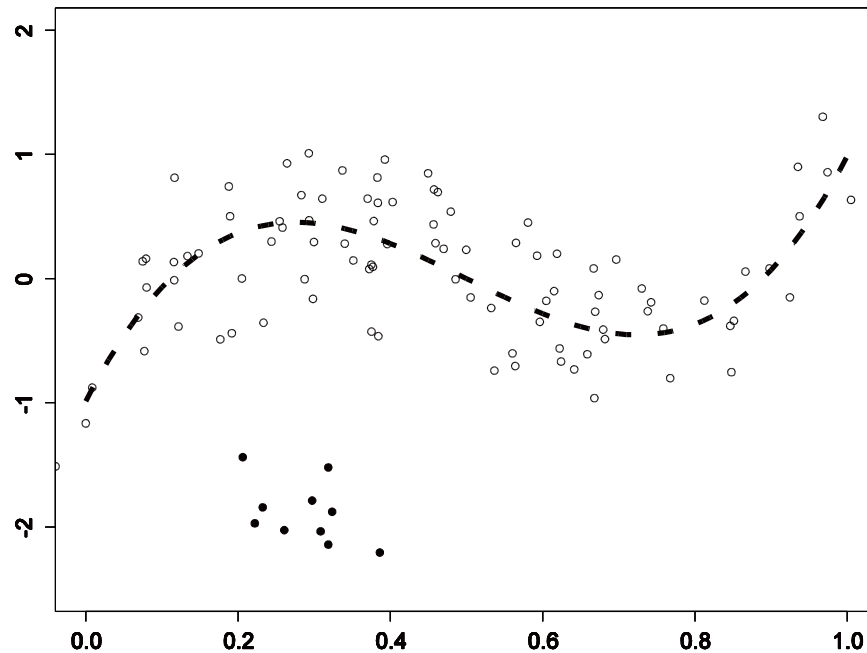
$$\text{IC}_B \rightarrow \text{GIC}_\beta, (B \rightarrow 0)$$

$$\text{IC}_B \rightarrow \text{TIC}, (B, \beta \rightarrow 0)$$

$q(x) = p(x | \bar{\theta}_\beta)$ の仮定の下で

$$\text{IC}_B \rightarrow \text{AIC}, (B, \beta \rightarrow 0)$$

モデル選択 (真のモデルを含まない)



データの設定

X: $N(0.5, 0.3^2)$

Y: Xの3次曲線

N=100 (外れ値10%)

誤差: $Y - X^T \beta \sim N(0, 0.2^2)$

確率モデルの設定

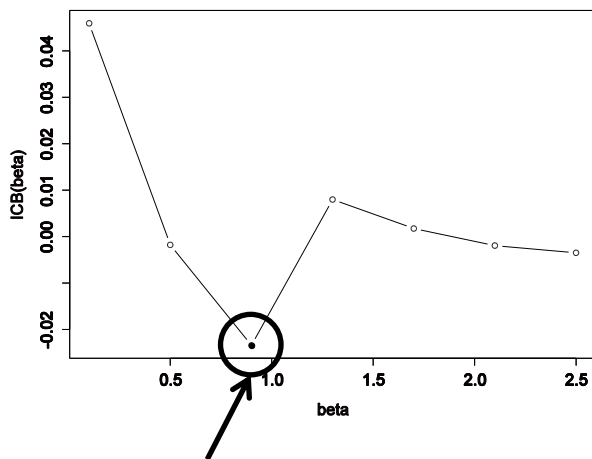
X: 一様に分布

Y: Xのd次曲線

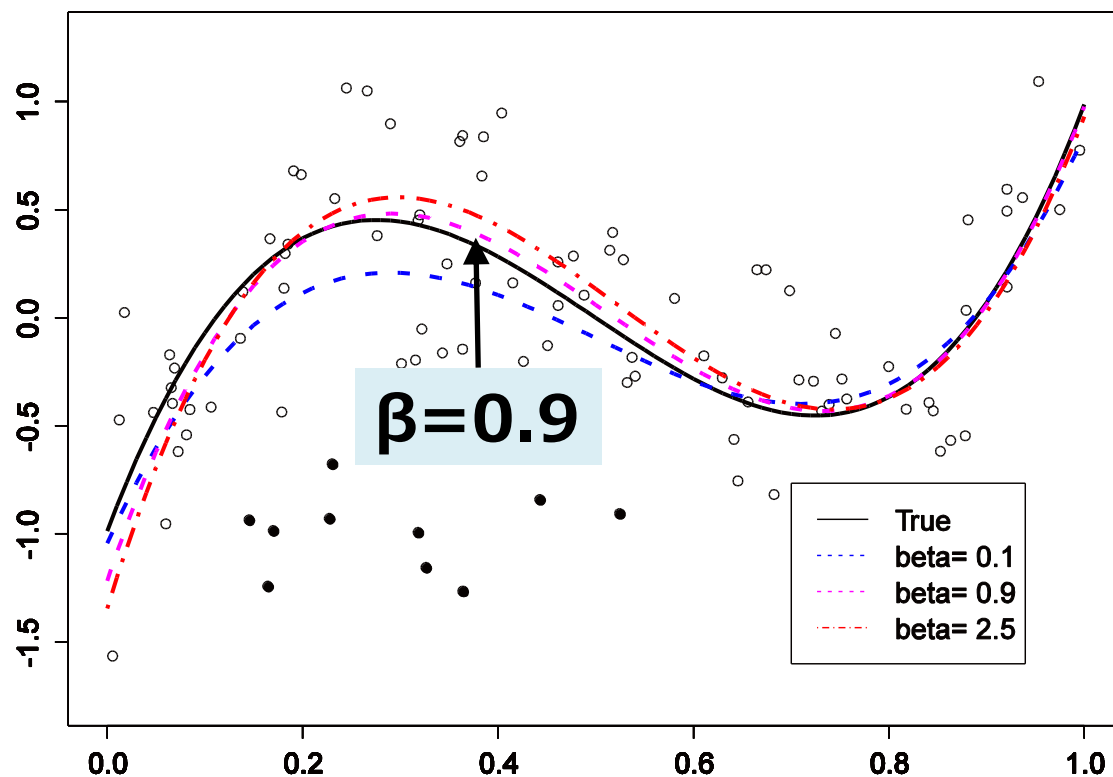
誤差: $Y - X^T \beta \sim N(0, \sigma^2)$

IC_Bを使うと, β を選択できる.

例) 3次多項式に限定, IC_Bが最小となる β を選択する.



$\beta = 0.9$ で最小



beta=* : 外れ値を含むデータで, beta=*としてbeta推定.