

重みつき最小二乗法で推定した 非線形ARXモデルの評価と選択

九州大学
マス・フォア・インダストリ研究所
秦 攀

発表内容

1. 研究動機
2. 基底関数に基づく非線形ARXモデルと重みつき最小二乗法
3. 一般化情報量基準(GIC)によるモデルの評価
4. ロバストなモデル選択法
5. LASSO型の選択法
6. 数値例

1

1. 研究動機

非線形システム同定



図1、非線形システム

- $u(t) \in R$: 外部の入力信号
- $w(t) \in R$: システムの出力信号
- $e(t) \in R$: 観測雑音
- $y(t) \in R$: は観測された出力信号
- $u(t)$ と $w(t)$ は非線形の関係を持つ
- 赤い信号 $w(t)$ と $e(t)$ は観測できないもの
- 離散時間信号 $\{u(t), y(t)\}$ のみに基づいて、システムをブラックボックスモデリング(Black-box modeling)する

2

1. 研究動機

非線形ARXモデル

- 非線形システム同定は離散時間信号 $\{u(t), y(t)\}$ に基づいて、非線形ARX(nonlinear autoregressive model with exogenous variables, NARX)モデル

$$M: y(t) = f(x(t)) + e(t)$$

を求めることである。ここで、

$$x(t) = [u(t-1), \dots, u(t-l_u), y(t-1), \dots, y(t-l_y)]^T$$

で定義する。

3

1. 研究動機

モデルの選択: 候補から、適切なモデルの構造を選択すること

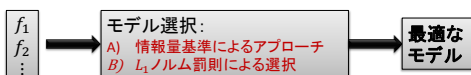


図2、モデル選択

本研究仮定

- $e(t) \sim N(0, \sigma^2)$ が独立同分布
- モデルは重みつき最小二乗法 (weighted least squares method, WLS) によって推定

4

2. 基底関数に基づくNARXモデルとRWLS

基底関数に基づくNARXモデル

$$y(t) = \sum_{k=1}^K \theta_k b_k(x(t)) + e(t) = \mathbf{b}(t)^T \boldsymbol{\theta} + e(t)$$

- $\theta_k(x(t)) \in R$ は $x(t)$ に関する基底関数(例えば、多項式、動径基底関数、スプラインなど)。
- $\theta_k \in R$ は線形的回帰係数。
- $\mathbf{b}(t) = [b_1(x(t)), \dots, b_K(x(t))]^T \in R^K$, $\boldsymbol{\theta}(t) = [\theta_1(x(t)), \dots, \theta_K(x(t))]^T \in R^K$,
- 全基底関数候補の集合 $\{\theta_k, k = 1, 2, \dots, K\}$ から適切な基底関数を選択する

基底関数の例

- 動径基底関数: $\exp\left\{-\frac{\|x-\mu\|^2}{2\sigma}\right\}$
- 3次スプライン: $\max\{0, (x-t_i)^3\}$

5

2.基底関数に基づくNARXモデルとWLS

重み付き最小二乗法 (weighted least squares method, WLS)

$$\hat{\theta}_{WLS} = \operatorname{argmin}_{\theta} \sum_{t=1}^N w_t (y(t) - b(t)^T \theta)^2 = \operatorname{argmin}_{\theta} (y - B\theta)^T W (y - B\theta)$$

$$\hat{\theta}_{WLS} = (B^T W B)^{-1} B^T W y$$

$W = I$ のとき、最小二乗法

- $w_t \geq 0$ は時刻 t の観測値につけられる **重み**
- $W = \operatorname{diag}(w_1, \dots, w_N)$ は **重み行列** (対角行列)
- $B = [b(1), \dots, b(N)]^T$ は基底関数による構成された **計画行列**

WLSを使う理由

- ある領域において、よりいい予測精度を期待するとき、その領域の観測値に大きい重みをつける。

例1: 計測装置によって、いい計測精度である領域の観測値に大きい重みをつける。

2.基底関数に基づくNARXモデルとWLS

WLSで推定されたNARXモデルの選択問題

- 最尤法に基づいたAICやBICなどの情報量基準が用いられなくなり、新しい情報量基準が必要になる。
- モデルの候補が大量であるとき、複雑な現象を統一的に説明でき、簡潔なモデルを構築できるロバストな選択手法を提案する。
- LASSO型の選択法をWLSに展開する。

3.GICによるモデルを評価

情報量基準

$$IC = -2 \sum_{t=1}^N \log f(y(t) | \hat{\theta}) + 2[\text{corrected-bias}]$$

分布の種類や推定法などによって決める

一般化情報量基準 (generalized information criterion, GIC)
[Konishi and Kitagawa (1996)]

$$GIC = -2 \sum_{t=1}^N \log f(y(t) | \hat{\theta}_M) + 2\operatorname{tr}\{R^{-1}Q\}$$

$$R = -\frac{1}{N} \sum_{t=1}^N \frac{\partial \psi(y(t), \theta)}{\partial \theta} \bigg|_{\theta = \hat{\theta}_M} \in R^k$$

$$Q = -\frac{1}{N} \sum_{t=1}^N \psi(y(t), \hat{\theta}_M) \frac{\partial \log(y(t) | \theta)}{\partial \theta^T} \bigg|_{\theta = \hat{\theta}_M} \in R^k$$

$\hat{\theta}_M \in R^k$ は「M-estimator」であり、汎関数 $\sum_{t=1}^N \psi(y(t), \theta) = 0$ の解である。ここで、 $\psi \in R^k$ は ψ 関数である。

3.GICによるモデルを評価

WLSに対応する ψ 関数

$$\sum_{t=1}^N \psi(y(t), \theta) = \sum_{t=1}^N \frac{\partial}{\partial \theta} \{w_t (y(t) - b(t)^T \theta)^2\} = -2B^T W y + 2B^T W B \theta = 0$$

GIC_{WLS}

$$GIC_{WLS} = -2 \sum_{t=1}^N \log f(y(t) | \hat{\theta}_{WLS}) + 2\operatorname{tr} \left\{ (B^T W B)^{-1} \sum_{t=1}^N \frac{w_t \epsilon(t)^2}{\hat{\sigma}_{WLS}^2} b(t) b(t)^T \right\}$$

- 対数尤度: $\log f(y(t) | \hat{\theta}_{WLS}) = -\frac{1}{2} \left\{ N \log(2\pi \hat{\sigma}_{WLS}^2) + \sum_{t=1}^N \frac{\epsilon(t)^2}{\hat{\sigma}_{WLS}^2} \right\}$
- 推定されたノイズの分散: $\hat{\sigma}_{WLS}^2 = \frac{(y - B\hat{\theta}_{WLS})^T W (y - B\hat{\theta}_{WLS})}{\operatorname{tr}(W)}$
- 予測誤差: $\epsilon(t) = y(t) - b(t)^T \hat{\theta}$
- GIC_{WLS} を最小するモデルが望ましい

4.ロバストなモデル選択法

GIC_{WLS}に基づく増加法 (Forward stepwise method)

- モデル候補の数は 2^k である。 k が大きいとき、全部の候補を評価することによって最適なモデルを選択するのは困難になる。

```

graph TD
    A[Initial step: B^(0) = []] --> B["k-th Step: B^(k) = [b_1, ..., b_k]"]
    B --> C["b_{k+1} mostly reducing GIC_WLS"]
    C --> D["k + 1-th Step: B^(k+1) = [b^(k), b_{k+1}]"]
    D --> E["if k + 1 < K, k = k + 1"]
    E --> B
    D --> F["else if k + 1 = K, break"]
    F --> G["Nested models: B^(0) ⊂ B^(1) ⊂ ... ⊂ B^(K)"]
    G --> H[GIC_WLS が最小であるモデルを選択]
    
```

図3. 増加法の流れ図

4.ロバストなモデル選択法

増加法のモンテカルロ実験

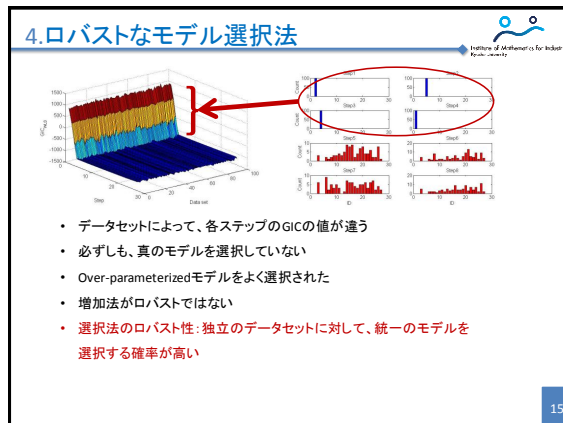
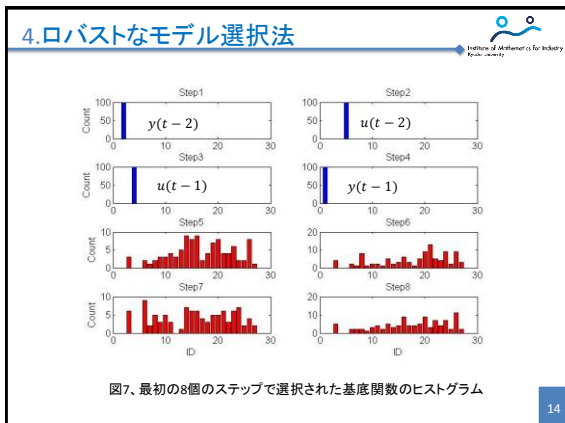
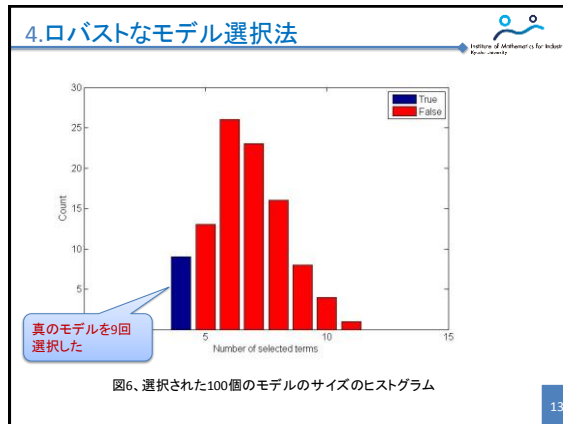
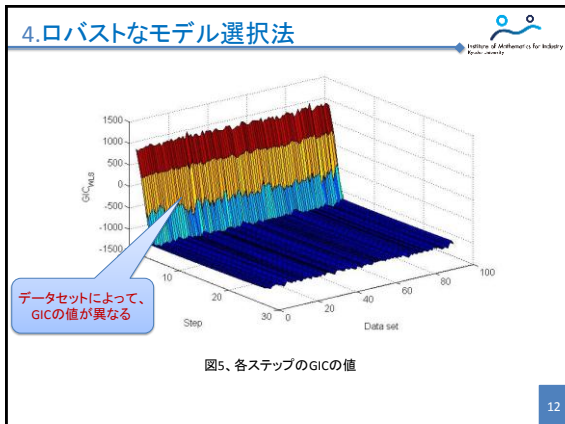
$$S_1: y(t) = -0.2y(t-1) + 0.75y(t-2) + 0.25u(t-1) + 0.3u(t-2) + e(t)$$

- $u(t) \sim N(0,1)$ $e(t) \sim N(0,0.01)$
- 100個の独立のデータセットを生成し、各データセットの長さが1000である
- そのデータセットを用い、GIC_{WLS}に基づく増加法を行う (即ち、100回のモンテカルロ実験)
- $\{y(t-1), y(t-2), u(t-1), u(t-2), u(t-3)\}$ の2次までの多項式を基底関数とする、全部は27個である。
- $w_t = |y(t)|$, 大きい観測値に大きい重みをつける

図4. 増加法のモンテカルロ実験

```

graph LR
    S1[S1] --> DS[Data set 1  
⋮  
Data set 100]
    DS --> Z[増加法]
    Z --> M[Model 1  
⋮  
Model 100]
    
```



4. ロバストなモデル選択法

ロバストなモデル選択基準

- $\Pr_j(b(t))$ はj番目のステップで基底関数 $b(t)$ が**選択された確率**を表す。この確率が $\Pr_j(b(t)) > \rho$ を満たすとき、 $b(t)$ を選ぶ。 $0 < \rho < 1$ はuser-definedパラメータである。

モンテカルロ法によって $\Pr_j(b(t))$ を近似

Data set 1
⋮
Data set m

増加法

$$\Pr_j^{MC}(b(t)) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}(b^{(i,j)}(t) = b(t))$$

図8、モンテカルロ法によって $\Pr_j(b(t))$ を近似

大量のデータセットを生成するのは困難である。

4. ロバストなモデルな選択法

Subsamplingによって $\Pr_j(b(t))$ を近似

$\{y(t_1), u(t_1)\}$
⋮
 $\{y(t_N), u(t_N)\}$

Sub-sequence 1
⋮
Sub-sequence m

増加法

$$\Pr_j^{SS}(b(t)) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}(b^{(i,j)}(t) = b(t))$$

図9、Subsamplingによって $\Pr_j(b(t))$ を近似

4.ロバストなモデルな選択法

Subsamplingの特徴

- 抽出:

$$X_{n_1}^{N_s} = \{y(t), u(t) | t = 1, \dots, N_s\}$$

$$X_{n_2}^{N_s} = \{y(t), u(t) | t = p + 1, \dots, N_s + p\}$$

$$\vdots$$

$$X_{n_m}^{N_s} = \{y(t), u(t) | t = (m-1)p + 1, \dots, N_s + (m-1)p\}$$
- $1 \leq N_s < N, p \geq 1$
 - $p \leq N_s$ のとき、Sub-sequenceが重なっているため、お互いに相関になる
 - モンテカルロ法に比べれば、一つのデータセットのみが必要
 - チューニングパラメータ N_s と p は結果に影響を与えられる

問題点:

- $e(t)$ の分布を仮定する
- パラメータの推定とモデル選択を分けて行う
- Subsamplingのため、計算量が増やす

18

5.LASSO型の選択法

重みつきLASSO(least absolute shrinkage and selection operator)推定 [Tibshirani (1996)]

$$\hat{\theta}_{LASSO} = \operatorname{argmin}_{\theta} (y - B\theta)^T W(y - B\theta) + \lambda_1 \sum_{k=1}^K |\theta_k|$$

(a) Ridge regression: $\theta_1^2 + \theta_2^2 < c_1$ (b) Lasso: $|\theta_1| + |\theta_2| < c_2$

図10. 2変数の縮小推定法

19

5.LASSO型の選択法

重みつきLASSO推定の計算

- 観測データと計画行列を重み行列で調節: $y^* = W^{\frac{1}{2}}y, B^* = W^{\frac{1}{2}}B$

$$\hat{\theta}_{LASSO} = \operatorname{argmin}_{\theta} (y^* - B^*\theta)^T (y^* - B^*\theta) + \lambda_1 \sum_{k=1}^K |\theta_k|$$
- Shooting, LARSなどの最適化アルゴリズムで推定

λ_1 の決め方

- $y^* = [y_1^*, \dots, y_p^*]^T, B^* = [B_1^{*T}, \dots, B_p^{*T}]^T$
- k-fold cross validation (CV)を定義する

$$CV(\lambda_1) = \sum_{p=1}^P \frac{1}{N_p} \left\| y_k^* - B^* \theta_{LASSO}^{(-p)} \right\|_2^2$$

- CVを最小する λ_1 を求める

20

5.LASSO型の選択法

図11. 重みつきLASSO推定の数値結果

21

5.LASSO型の選択法

Adaptive LASSO推定[Zou (2006)]

$$\hat{\theta}_{LASSO} = \operatorname{argmin}_{\theta} (y - B\theta)^T W(y - B\theta) + \lambda_1 \sum_{k=1}^K v_k |\theta_k|$$

adaptive LASSO推定の計算

- $\theta^* = [v_1 \theta_1, v_2 \theta_2, \dots, v_K \theta_K]^T, \Phi^* = \begin{bmatrix} \Phi_{11}^* & \Phi_{12}^* & \dots & \Phi_{1K}^* \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_{K1}^* & \Phi_{K2}^* & \dots & \Phi_{KK}^* \end{bmatrix}$
- $\hat{\theta}_{LASSO}^* = \operatorname{argmin}_{\theta} (y^* - B^*\theta)^T (y^* - B^*\theta) + \lambda_1 \sum_{k=1}^K |\theta_k^*|$
- $\theta_{adaLASSO}^* = \begin{bmatrix} \hat{\theta}_1^* & \hat{\theta}_2^* & \dots & \hat{\theta}_K^* \end{bmatrix}^T$

21

5.LASSO型の選択法

適応重み v_k のつけ方

- Ridge regression推定を計算

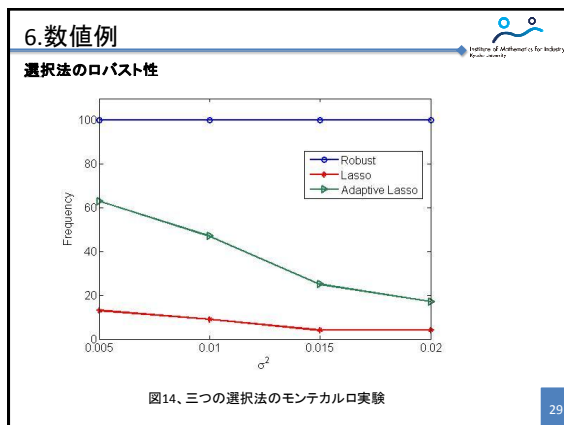
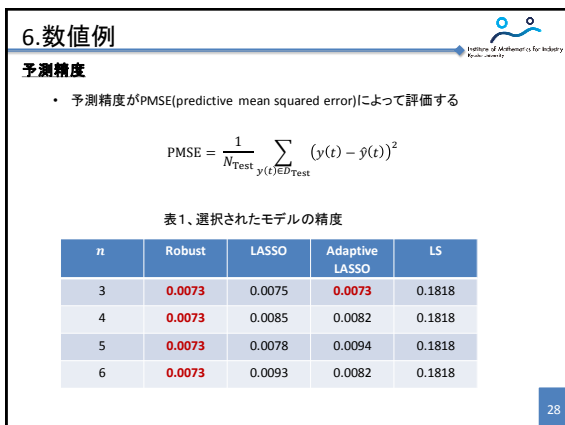
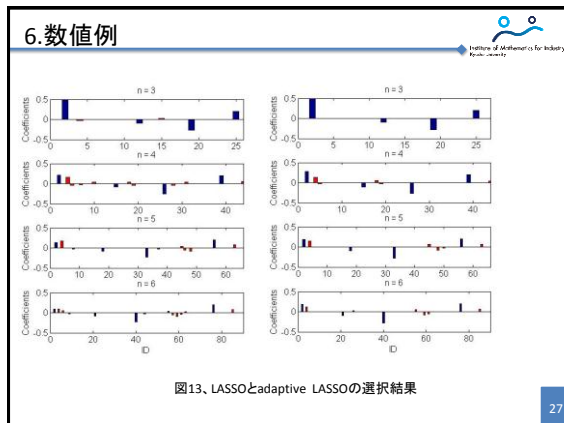
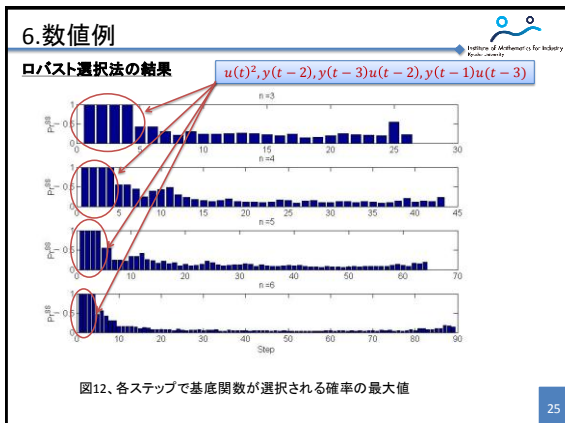
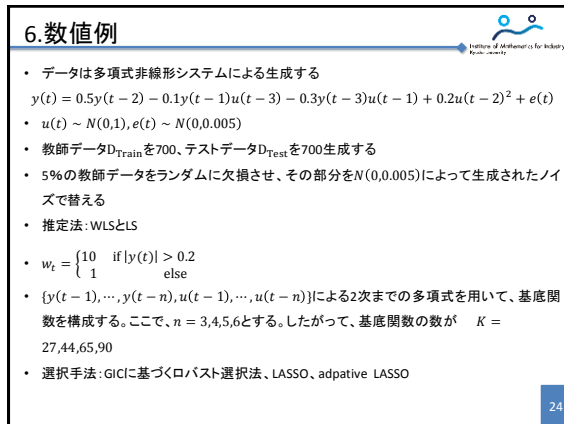
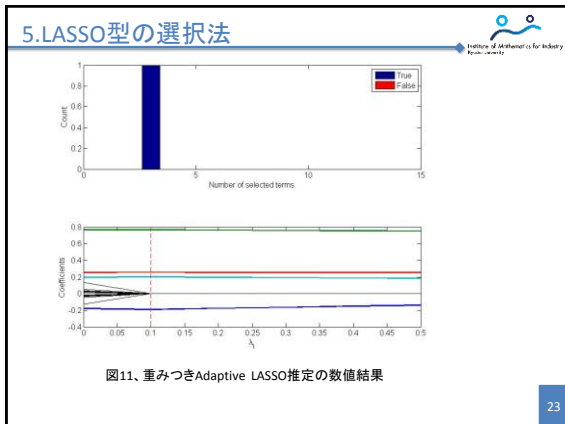
$$\hat{\theta}_{ridge} = \operatorname{argmin}_{\theta} (y - B\theta)^T W(y - B\theta) + \lambda_r \sum_{k=1}^K \theta_k^2$$
- $\hat{\theta}_{ridge}$ のk番目の要素 $\frac{1}{|\hat{\theta}_{ridge}|}$ を v_k とする

正規化パラメータの決め方

$$CV(\lambda_1, \lambda_r) = \sum_{p=1}^P \frac{1}{N_p} \left\| y_k^* - B^* \theta_{adaLASSO}^{(-p)} \right\|_2^2$$

を最小する λ_1, λ_r

22



7. まとめ



- WLSで推定されたNARXモデルに対して、GICに基づくロバスト選択法、LASSO型の縮小選択法を提案した。
- GICに基づくロバスト選択法が予想のように、統一及び簡潔なモデルを選択できる。しかし、Subsamplingを用いるため、計算は複雑である。
- GICに基づくロバスト選択法に比べれば、LASSO型の選択法がモデルの推定と選択を同時に行える。
- Adaptive LASSOは確実にLASSOの選択パフォーマンスを改善できる。
- ノイズの分散を大きくすると、(すなわち、SN比(signal-noise-ratio)を小さくすると)、LASSO型のロバスト性が落ちる。それに比べれば、GICに基づくロバスト選択法がロバスト性を保てる。

30

参考文献



1. Konishi, S., Kitakawa, G.: 'Generalised information criteria in model selection' *Biometrika*, **83**(4), 875-890 (1996)
2. Garatti, S., Bittard, R.R.: 'On resampling and uncertainty estimation in linear system identification', *Automatica*, **46**(5), 785-795 (2010)
3. Tibshirani, R.: 'Regression Shrinkage and Selection via the Lasso', *J. Roy. Stat. Soc. B.*, **58**(1), 267-288 (1996)
4. Zou, H.: 'The adaptive Lasso and its oracle properties', *J. Am. Stat. Assoc.*, **101**(476), 1418-1429 (2006)

30