

# ウミガメの遊泳を捉えた時系列の ARIMA モデルに基づく分割とクラスタリング

Naoto K. Inoue\*

## Abstract

ウミガメの福祉評価を行うためには、遊泳パターンの継続時間と変化を定量化する必要がある。この目的で、これまでに筆者は形状解析によってウミガメの輪郭形状を二次元時系列に圧縮した<sup>1</sup>。本論では、取得した二次元時系列に対し時系列解析とクラスタリング手法を用いることで、遊泳パターンの客観的な分割と分類を試みた。このとき、計算が軽く柔軟なモデルである AR モデルを仮定するが、クラスタリングで用いられる時系列間の非類似度のうち AR モデルをベースとしたものには、Piccolo (1990), Maharaj (1996), そして Kalpakis and Putagunta (2001) によって提案された指標が存在する。そこで、まずこれらの指標のレビューを行った。そして、ウミガメの輪郭形状から取得した時系列データに対して、AR モデルを仮定した分割・分類を行い、形成されるクラスタの特徴を調べた。レビューの結果、クラスタリング手法と相性が良く、結果の解釈性も高い piccolo と Kalpakis and Putagunta の 2 つの指標を用いることに決定した。結果として形成されたクラスタに明確な違いは見られず、現段階では分類精度が高くなかったが、これは時系列間の相関を考慮していないことが原因であると考えられる。

---

\*Department of Science, Kobe University, Hyogo, Japan. Email:226s405s@stu.kobe-u.ac.jp

<sup>1</sup>この内容はまとめて論文にしているので、近日中に発表します

# 1 導入

時系列クラスタリングは、パターンや傾向の似た時系列を分類する手法であり、株取引や売上の動向を調べる、血液分析の結果から健康状態を判断するといった目的で用いられている。時系列クラスタリングを行うには、時系列間の非類似度を決定する必要があるが、分類対象の時系列に応じた非類似度の選択はクラスタリングの精度に大きく影響を及ぼす。それゆえ、これまでにさまざまな非類似度が提案されてきた。

Liao (2005) では、時系列間の非類似度を、生データベース、特徴ベース、モデルベースという3つのカテゴリに分けてレビューを行った [7]。ここで、生データベースは二つの時系列の値の単純なユークリッド距離や、カルバックライブラー情報量<sup>2</sup>、ガウスノイズを仮定した時の誤差など、生データそのものから導出されるものである。また、特徴ベースは、時系列間で相互相関関数を計算したり、スペクトル変換、あるいはウェーブレット変換<sup>3</sup>を行った上でその距離を測るものである。そして、モデルベースは、後述する ARIMA モデルをはじめ、正規混合モデル、p 値など、時系列を生成したモデルを仮定し、その仮定のもとで非類似を測定するものである。

Liao (2005) の3つの分類に対し、Montero and Vilar (2015) は非類似度を、モデルフリー、モデルベース、複雑さベース、予測ベースの4つに分けてレビューを行った [9]。ここで、モデルフリーを構成する要素は Liao の生データベース・特徴ベースと似ており、モデルベースもほとんど同じ定義である。一方、複雑さベースには、コルモゴロフの複雑性<sup>4</sup>や、パーミュテーション分布に基づく複雑指標<sup>5</sup>などが位置付けられている。また、予測ベースは Montero and Vilar (2015) が提案した新たな非類似度であり、時系列間の予測値の距離を指標としたものである。

本論では、さまざまな非類似度の中で、ARIMA(自己回帰和文移動平均) モデルベースの手法に着目する。ARIMA モデルは、時系列の周期成分、トレンド、確率的変動、ランダム

<sup>2</sup>情報理論におけるデータ間の距離や、確率論における確立分布間の距離の指標としてよく用いられる。

<sup>3</sup>一言で言うと、時間スケールを考慮したフーリエ変換。現象のスケール変化を許容した上で周波数特性を抽出することができる。

<sup>4</sup>情報理論における文字列の複雑性の指標。プログラムにおいて、特定の文字列を記述可能な最小の長さ。

<sup>5</sup>二つの時系列の値を全てランダムに入れ替える作業を全通り繰り返した際に生成される分布。仮説検定をモチベーションとして考案された。

要素を捉えることが可能なモデルである。このモデルを利用する理由は以下の2点である。  
1つ目は、本論で扱う時系列がさまざまな周期特性を持っている点である。本論で扱う時系列は、ウミガメの輪郭座標の変化を波として定量化したものであり、大まかには前肢と後肢の動きを反映したものである。前肢は遊泳において動力を得る役割を果たすため、振幅が大きく周期も長い運動を行う一方で、後肢は振幅が小さく周期も短い。これらの性質に加え、遊泳パターンによって、さまざまなスペクトルが卓越する。

2つ目の理由は、クラスタリングの前準備にARモデル (ARIMAの単純化) を利用するためである。筆者の研究の目的は、ウミガメの遊泳パターンの分類であり、これを達成するにはクラスタリングを行う前に、数時間にわたる時系列を定常時系列<sup>6</sup>に分割する必要がある。この分割を行う際、各区間を異なるARモデルに従う時系列と考えると、ARモデルの当てはまりの良さ (AIC; 赤池情報量基準) を用いて時系列を分割することが可能になる。それゆえ、クラスタリングにおいてもAR(IMA)モデルを仮定することで、一貫した解析を行う必要がある。

本論では、まず2.2節においてARモデルを用いた時系列の分類手法を導入する。次に、2.3節においてARIMAモデルを仮定したモデルベースの指標のレビューを行う。ここでは、提案された年度順に、(1)Piccolo (1986)[10]、(2)Maharaj (1996)[8]、(3)Kalpakis and Putagunta (2001)[5]の指標を扱う。そして、実際にこれらの手法を用いて時系列を定常区間に分類し、レビューの結果選んだ非類似度を用いてクラスタリングを行った。

## 2 手法

### 2.1 材料

本論では、ウミガメの輪郭形状の変化に対して形状解析を行うことで抽出した2本の時系列pc1、pc2を用いる (動画1のtimeseriesを参照のこと)。ここで、pc1が正の時は両前肢を広げた姿勢を表し、負のときは丸まっている姿勢を表す。一方、pc2が正の時は右前肢だけ横に開いた姿勢を表し、負のときは左前肢だけ横に開いた姿勢を表す (動画1のpc plotを参

---

<sup>6</sup>一言で言うと、時間が経っても同じ確立分布に従う時系列のこと。厳密には強定常と言われる。

照のこと)。それゆえ筆者は pc1 と pc2 の変化をそれぞれ、体軸を中心に非対称的な漕ぎ運動 (rowing) と、対称的な羽ばたき運動 (flapping) の特徴量として捉えた。

## 2.2 局所定常 AR モデルによる分割

筆者は、これら pc1 と pc2 の時系列に対して AR モデルを適用し、分割前に比べて分割後の時系列のモデルの当てはまりがよくなる点を網羅的に探すことで、各個体で 1 時間の時系列を最小の長さが 10 秒以上となるような短い時系列に分割する。以下の分割法は Kitagawa (2020) [6] を参考に実装した。

### 2.2.1 AR モデルと最小二乗法

$y_n$  を  $n$  番目の時系列データとすると、 $m$  次の AR(自己回帰) モデルは、

$$y_n = \sum_{i=1}^m a_i y_{n-i} + \epsilon, \quad (1)$$

と書くことができる。ここで  $\epsilon$  は平均 0、分散  $\sigma^2$  の正規分布に従う確率変数である。 $a_i$  は一般的な回帰における回帰係数に相当することから、AR モデルは自身の過去のデータの加重平均  $\sum_{i=1}^m a_i y_{n-i}$  に正規誤差を加えることで現在のデータが生成されることを仮定したモデルであるといえる。このとき時系列の  $M(> m)$  番目から  $N$  番目までに注目すると、回帰式は  $N - M$  次元ベクトル  $y$  と  $(N - M) \times m$  行列  $Z$

$$y = \begin{bmatrix} y_{M+1} \\ y_{M+2} \\ \vdots \\ y_N \end{bmatrix}, \quad Z = \begin{bmatrix} y_M & y_{M-1} & \cdots & y_{M-m} \\ y_{M+1} & y_M & \cdots & y_{M-m+1} \\ \vdots & \vdots & \ddots & \vdots \\ y_{N-1} & y_{N-2} & \cdots & y_{N-m} \end{bmatrix}, \quad (2)$$

を用いて、

$$y = Za + \epsilon, \quad (3)$$

と書くことができる。

78 この回帰モデルにおいて、回帰係数  $a$ 、分散  $\sigma^2$ 、説明変数  $x$  が与えられている時、対数尤  
79 度<sup>7</sup>は、

$$l(a_1, a_2, \dots, a_m, \sigma^2) = \sum_{n=M+1}^N \log p(y_n | a_1, a_2, \dots, a_m, \sigma^2, y_{n-1}, y_{n-2}, \dots, y_{n-m}), \quad (4)$$

80 と与えられる。ここで、 $\theta = (a_1, a_2, \dots, a_m, \sigma^2)$  とすると、データの正規性から、

$$l(\theta) = -\frac{N-M}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=M+1}^N \left( y_n - \sum_{i=1}^m a_i y_{n-i} \right)^2, \quad (5)$$

81 となる。この対数尤度が極大となる  $\sigma$  は、

$$\frac{\partial l(\theta)}{\partial \sigma^2} = -\frac{N-M}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{n=M+1}^N \left( y_n - \sum_{i=1}^m a_i y_{n-i} \right)^2, \quad (6)$$

82 を 0 とおくことにより、

$$\hat{\sigma}^2 = \frac{1}{N-M} \sum_{n=M+1}^N \left( y_n - \sum_{i=1}^m a_i y_{n-i} \right)^2 \quad (7)$$

83 となる。この分散の推定値を (5) 式に代入すると、対数尤度は

$$\begin{aligned} l(a_1, a_2, \dots, a_m) &= -\frac{N-M}{2} \log(2\pi\hat{\sigma}^2) - \frac{N-M}{2} \\ &= \log \hat{\sigma}^2 + \text{Const}, \end{aligned} \quad (8)$$

84 と与えられる。このことから、分散を最小化することで AR 係数の最尤推定値を得ることが  
85 できることがわかる。

## 86 2.2.2 ハウスホルダー法

87 分散は、(3) の行列表現により、

$$|\epsilon|^2 = |y - Za|^2 \quad (9)$$

---

<sup>7</sup>確率過程に基づく当てはまりの良さの指標

と書くことができる。ここで  $|A|$  はベクトル  $A$  のユークリッドノルムを表す。分散の最小化のためにこの式を微分して 0 とすることで、回帰係数  $a$  の推定値は、

$$\hat{a} = (Z^T Z)^{-1} Z^T y, \quad (10)$$

と得られる。これはいわゆる最小二乗法であるが、この逆行列の計算は時系列が長くなればなるほど計算量が膨大となる。

そこで、 $(N - M) \times (N - M)$  の直行行列  $U$  を用いた変換

$$|y - Za|^2 = |U(y - Za)|^2 = |Uy - UZa|^2 \quad (11)$$

を考える。 $a$  の最尤推定値には影響がないことに注意するとよい。ここで、 $UZ$  を扱いやすい形に変換してやることができれば、計算を簡単にすることができる。

以下では、 $U$  によるハウスホルダー変換を考える。まず、行列  $Z$  の右側にデータのベクトル  $y$  を付加した  $(N - M) \times (m + 1)$  行列

$$X = [Z \mid y], \quad (12)$$

を考える。この行列に対して適当な  $U$  をかけると、 $(N - M) \times (m + 1)$  の上三角行列に変換することができる。

$$UX = [UZ \mid Uy] = S = \begin{bmatrix} s_{1,1} & s_{1,2} & \cdots & s_{1,m} & s_{1,(m+1)} \\ 0 & s_{2,2} & \cdots & s_{2,m} & s_{2,(m+1)} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & s_{m,m} & s_{m,(m+1)} \\ 0 & 0 & \cdots & 0 & s_{(m+1),(m+1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & s_{(N-M),(m+1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & 0 \end{bmatrix} \quad (13)$$

99 このとき、1 から  $m$  列までが  $UZ$  に、 $m + 1$  列が  $Uy$  に対応するため、

$$|Uy - UZa|^2 = \left\| \begin{bmatrix} s_{1,(m+1)} \\ s_{2,(m+1)} \\ \vdots \\ s_{m,(m+1)} \\ s_{(m+1),(m+1)} \\ \vdots \\ s_{(N-M),(m+1)} \\ \vdots \\ 0 \end{bmatrix} - \begin{bmatrix} s_{1,1} & s_{1,2} & \cdots & s_{1,m} \\ 0 & s_{2,2} & \cdots & s_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_{m,m} \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} \right\|^2 \quad (14)$$

100 となる。これは、

$$|Uy - UZa|^2 = \left\| \begin{bmatrix} s_{1,(m+1)} \\ s_{2,(m+1)} \\ \vdots \\ s_{m,(m+1)} \end{bmatrix} - \begin{bmatrix} s_{1,1} & s_{1,2} & \cdots & s_{1,m} \\ 0 & s_{2,2} & \cdots & s_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_{m,m} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} \right\|^2 + s_{(m+1),(m+1)}^2 + \cdots + s_{(N-M),(m+1)}^2, \quad (15)$$

101 と書き直すことができる。ここでユークリッドノルムの外は  $a$  に関係しないため、ユークリッ  
102 ドノルムの中を 0 にすればよい。つまり、 $a$  の最尤推定値は

$$\begin{bmatrix} s_{1,1} & s_{1,2} & \cdots & s_{1,m} \\ 0 & s_{2,2} & \cdots & s_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_{m,m} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} = \begin{bmatrix} s_{1,(m+1)} \\ s_{2,(m+1)} \\ \vdots \\ s_{m,(m+1)} \end{bmatrix} \quad (16)$$

103 の解として求めることができる。例えば、下から 2 行は

$$s_{(m-1),(m-1)}a_{(m-1)} + s_{(m-1),m}a_m = s_{(m-1),(m+1)}, \quad s_{m,m}a_m = s_{m,(m+1)}, \quad (17)$$

104 となるため、 $\hat{a}_m = s_{m,(m+1)}/s_{m,m}$  であり、これを代入することで、1つ前を

$$\hat{a}_{(m-1)} = \frac{s_{(m-1),(m-1)} - s_{(m-1),m}\hat{a}_m}{s_{(m-1),(m-1)}}, \quad (18)$$

105 と求めることができる。これを一般化すると、

$$\hat{a}_i = \frac{s_{i,m+1} - s_{i,m}\hat{a}_m - \cdots - s_{i,(i+1)}\hat{a}_m}{s_{i,i}}, \quad (19)$$

106 と書くことができる。このとき、(15)は

$$\hat{\sigma}_m^2 = \frac{1}{N-M} \sum_{j=m+1}^{N-M} s_{j,m+1}^2 \quad (20)$$

107 となり、モデルの AIC <sup>8</sup>は、

$$AIC_m = (N-M)(\log 2\pi\hat{\sigma}^2 + 1) + 2(m+1), \quad (21)$$

108 と計算することができる。

109 ハウスホルダー変換の優れた点は2つ存在する。1つ目は、ひとたび  $m$  次の AIC を求める  
110 際に行列  $S$  を計算しておけば、

$$\hat{\sigma}_j^2 = \frac{1}{N-M} \sum_{i=j+1}^{N-M} s_{j,m+1}^2 \quad (22)$$

111 と  $j$  次のモデルの分散を求めることで、 $j$  次のモデルの AIC が計算できる点である。2 点目  
112 は、時系列を加えた際の AIC の変化は、新たなデータベクトル  $y'$  と行列  $Z'$  により作成した  
113  $X'$  を  $S$  に対して列方向に結合し、新たに得た  $S'$  に対してハウスホルダー変換を行うことで、  
114 逐次的に計算できる点である。これらの点から、ハウスホルダー変換を用いた最小二乗法は、  
115 他の手法に比べて計算が非常に軽い。

116 ハウスホルダー変換を用いた最小二乗法を用いて、pc1 と pc2 時系列を客観的に分割する。  
117 具体的には、pc1 や pc2 の区間時系列  $\mathbf{y}_k = y_{k,1}, y_{k,2}, \dots, y_{k,N}$  に対し、候補点  $y_{k,c}, \dots, y_{k,C}$   
118 で時系列を分割した際の AIC を計算する。次に、候補点の中から最も AIC の低い、すなわち

---

<sup>8</sup>AIC は、 $-2 \times$  対数尤度  $+ 2 \times$  パラメータ数によって計算される。AIC の導出は [1] を参照のこと



119 AR モデルの当てはまりの良い点  $y_{k,s}$  をこの区間の分割点として決定する。そして、次に分割  
 120 を行う区間時系列は、前の区間で選択された分割点が  $y_{k,s}$  であることから、 $y_{n+1,1} = y_{k,s} + 1$   
 121 と設定する。これを繰り返すことで、時系列全体を区間時系列に分割する。本論では、ウミ  
 122 ガメの最小の遊泳行動継続時間を 10 秒以上と考え、 $N$  には 40 秒を設定し、分割候補点は 40  
 123 秒を (10,20,10) に分けた中間の 20 秒の中から選択した<sup>9</sup>。

## 124 2.3 AR モデルを仮定した時系列クラスタリング

125 2.2 節の操作によって分割した区間時系列  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k$  を、各区間の AR 係数を特徴量と  
 126 して分類することを考える。

### 127 2.3.1 Piccolo (1990) の指標

128 Piccolo は、時系列  $\mathbf{y}_p, \mathbf{y}_q$  から求められた AR 係数  $\mathbf{a}_p = a_{p,1}, a_{p,2}, \dots$  と  $\mathbf{a}_q = a_{q,1}, a_{q,2}, \dots$   
 129 に関して、

$$d(\mathbf{y}_p, \mathbf{y}_q)_{pic} = \left\{ \sum_{j=1}^{\infty} (a_{p,j} - a_{q,j}) \right\}^2, \quad (23)$$

130 を時系列間の非類似度として提案した。この指標は以下に述べる特徴を持つ。

131 1-pic. 古典的な距離の定義 (非負性、対称性、三角不等式が成立する) を満たす点

132 2-pic. 定常時系列に対しては必ず導出することができる点

133 3-pic. AR 係数が 0 であっても適用可能である点

134 4-pic. 時系列のばらつきを考慮しないため、振幅のスケールに依存しない点

135 5-pic. 片方の時系列の AR 係数が全て 0 のとき (すなわちホワイトノイズであるとき)、距離は  
 136 もう一方の時系列の係数の絶対値となる点

137 6-pic. AR 係数に基づくため、現象の時間スケールに依存しない点

138 7-pic. 任意の  $AR(j)$  モデルに従う時系列間の距離は有界である点

<sup>9</sup>実際には 1 秒 5 フレームに分割している

### 2.3.2 Maharaj (1996) の指標

Maharaj は、時系列  $\mathbf{y}_p, \mathbf{y}_q$  から求められた  $m$  次の AR 係数  $\mathbf{a}_p$  と  $\mathbf{a}_q$  に関して、自由度  $m$  のカイ二乗分布に従う検定統計量を時系列間の非類似度として提案した。以下ではこの指標を導出し、特性を述べる。

2つの時系列  $\mathbf{y}_p, \mathbf{y}_q$  を生成したプロセスの違いについての検定を行う場合に、以下の帰無仮説、対立仮説を考える。

$H_0$ : 2つの定常時系列を生成したプロセスが同じで  $\mathbf{a}_p = \mathbf{a}_q$

$H_1$ : 2つの定常時系列を生成したプロセスが異なり  $\mathbf{a}_p \neq \mathbf{a}_q$

Bhansall (1978) [4] は、長さが  $T$  の時系列に対し、AR 次数  $m$  が

$$\frac{m^3}{T} \rightarrow 0, \quad (24)$$

を満たすとき、

$$\frac{1}{\sqrt{T}} \sum_{j=m+1}^{\infty} |a_{p,j}| \rightarrow 0, \quad (25)$$

となるという Berk (1974) [3] の結果を利用して、推定された AR 係数  $\hat{\mathbf{a}}_p$  と真の AR 係数  $\mathbf{a}_p$  の差が、

$$\sqrt{T}(\hat{\mathbf{a}}_p - \mathbf{a}_p) \sim N(0, \sigma_p^2 R_p^{-1}(m)), \quad (26)$$

と漸近的に正規分布に従うことを導いた。 $\sigma_{a_p}^2$  は時系列が含むホワイトノイズの分散であり、

$R_p(m)$  は定常時系列の自己共分散行列における  $m \times m$  部分行列である。

$$R_p = \begin{bmatrix} C_0 & C_1 & \cdots & C_m & \cdots \\ C_1 & C_0 & \cdots & C_{m-1} & \cdots \\ \vdots & \vdots & \ddots & \vdots & \cdot \\ C_m & C_{m-1} & \cdots & C_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \quad (27)$$

153 ここで、

$$C_k = \frac{1}{T} \sum_{n=k+1}^T (a_{p,n} - \mu_p)(a_{p,n-k} - \mu_p), \quad (28)$$

154 はラグ  $k$  の自己共分散関数と呼ばれる。

155 この漸近的な性質から、2本の時系列から推定された AR 係数は、

$$\hat{\mathbf{a}}_p \sim N\left(\mathbf{a}_p, \frac{1}{T} \frac{\sigma_p^2}{R_p(m)}\right), \quad \hat{\mathbf{a}}_q \sim N\left(\mathbf{a}_q, \frac{1}{T} \frac{\sigma_q^2}{R_q(m)}\right), \quad (29)$$

156 に従う。これより、

$$\hat{\mathbf{a}}_p - \hat{\mathbf{a}}_q \sim N\left(\mathbf{a}_p - \mathbf{a}_q, \frac{1}{T} [\sigma_p^2 R_p^{-1}(m) + \sigma_q^2 R_q^{-1}(m)]\right) \quad (30)$$

157 となるため、

$$d(\mathbf{y}_p, \mathbf{y}_q)_{mah} = (\hat{\mathbf{a}}_p - \hat{\mathbf{a}}_q)^T [\sigma_p^2 R_p^{-1}(m) + \sigma_q^2 R_q^{-1}(m)]^{-1} (\hat{\mathbf{a}}_p - \hat{\mathbf{a}}_q) \sim \chi^2(m), \quad (31)$$

158 が成り立つ。ここで、 $\chi^2(m)$  は自由度  $m$  のカイ二乗分布である。すなわち  $d(\mathbf{y}_p, \mathbf{y}_q)$  は、本  
159 節の初めに設定した仮設検定における検定統計量であり、Maharaj はこれを非類似度の指標  
160 として用いることを提案した。

161 Maharaj はいくつかの ARIMA モデルから生成した時系列に対してこの非類似度の指標を  
162 適用し、以下の特徴を強調した。

163 1-mah. 非負性と対称性を持つ点

164 2-mah. 全ての定常時系列に対し計算することが可能である点

165 3-mah. Piccolo の指標と異なり、時系列の振幅のスケールに影響を受ける点

166 4-mah. Piccolo の指標に比べ、クラスタリングの早い段階で類似したクラスタを形成すること  
167 ができる点

### 2.3.3 Kalpakis and Putagunta (2001) の指標

Kalpakis and Putagunta は、時系列のケプストラム (cepstral coefficients) を非類似度の指標として提案した。ケプストラムの導出は、Atal (1974) [2] を参考にした。

AR モデル、

$$y_n = \sum_{i=1}^m a_i y_{n-i} + \epsilon, \quad (32)$$

は  $y_{n-1} = B y_n$  を満たす時間シフトオペレータ  $B$  によって

$$y_n = \sum_{i=1}^m a_i B^i y_n + \epsilon \quad (33)$$

と表現することができる。これは、

$$y_n = \frac{1}{1 - \sum_{i=1}^m a_i B^i} \epsilon = \sum_{i=1}^{\infty} g_i B^i = g(B) \epsilon \quad (34)$$

と書くことができる。定常 AR モデルにおいてこの式は、過去から現在までのホワイトノイズの足し合わせが現在の時系列を生じていることを表し、 $g_i$  をインパルス関数、 $g(B)$  を伝達関数と呼ぶ。ここで、対数を取った伝達関数のべき展開を考える。

$$\log g(B) = \sum_{j=1}^{\infty} \phi_j \exp(-j f n), \quad (35)$$

この  $\phi_j$  がケプストラムである。この式において  $f$  が周波数の時、ケプストラムは  $\log g(B)$  のフーリエ係数であると解釈することができる。ケプストラムはもともと音声の解析のために考案された。これは、声が周波数の大きな声帯成分と周波数の小さい声道成分からなり、その分離は伝達関数の  $f$  に対する周期性をもとに行われるためである。この点から、さまざまな周期運動が混在した時系列において、ケプストラムは有用な情報をもたらす。

ケプストラムと AR 係数の関係式を導出する。ケプストラムのべき式において、両辺を  $B$  で微分することを考える。

$$\frac{d}{dB} \log \left[ \frac{1}{1 - \sum_{i=1}^m a_i B^i} \right] = \frac{d}{dB} \sum_{i=1}^{\infty} \phi_i B^i, \quad (36)$$

184 これは、

$$\frac{\sum_{i=1}^m i a_i B^{i-1}}{1 - \sum_{i=1}^m a_i B^i} = \sum_{j=1}^{\infty} j \phi_j B^{j-1} \quad (37)$$

185 となり、

$$\sum_{i=1}^m i a_i B^{i-1} = (1 - \sum_{i=1}^m a_i B^i) \left( \sum_{j=1}^{\infty} j \phi_j B^{j-1} \right) \quad (38)$$

186 を得ることができる。この式を展開して整理すると、 $m$  次の AR モデルのケプストラムは、

$$\begin{aligned} \phi_1 &= a_1, \\ \phi_n &= a_n + \sum_{j=1}^{n-1} \left( 1 - \frac{j}{n} \right) a_j \phi_{n-j}, \quad \text{if } 1 < n \leq m, \\ \phi_n &= \sum_{j=1}^m \left( 1 - \frac{j}{n} \right) a_j \phi_{n-j}, \quad \text{if } m < n \end{aligned} \quad (39)$$

187 と AR 係数によって与えられる。

188 Kalpakis and Puttagunta (2001) [5] は、二つの時系列間の非類似度を、ケプストラム間の  
189 ユークリッド距離

$$d(\mathbf{y}_p, \mathbf{y}_q)_{cep} = \left\{ \sum_{j=1}^{\infty} (\phi_{p,j} - \phi_{q,j}) \right\}^2, \quad (40)$$

190 と定義した。また、シミュレーションした ARMIA 時系列に対して、この非類似度に基づい  
191 たクラスタリングを行い、以下の特徴を挙げた。

192 1-cep. 次数の小さい係数ほど重みが大きく設定されており、非類似度に大きく影響を及ぼす点

193 2-cep. 次数が大きくなればなるほど、影響が急速に 0 に近づくため、はじめの数個の係数が  
194 うまく減少を捉えている必要がある点

195 3-cep. 他の時系列指標 (主に特徴ベースの手法) よりも AR モデルの分類精度が高い点

196 4-cep. 振幅の平均やスケールに対して不変である点

197 5-cep. 時系列の並行移動に対して不変である点

198 6-cep. 周期性に対して敏感である点

199 7-cep. インパルス関数の積算はケプストラムでは加法となるため、時系列の変化を減法によっ  
200 て検出できる可能性がある点

## 201 2.4 非類似度の利用とクラスタリング

202 これら3つの非類似度のレビューから、本研究では  $d_{pic}$  と  $d_{cep}$  を用いることにする。これ  
203 には2つの理由がある。1つ目は、 $d_{mah}$  は時系列間の距離を計算することはできるが、相対  
204 座標が与えられない点である。距離だけが与えられた場合においてもクラスタリングを実行  
205 することはできるが、それにはクラスタ間の距離を近似しつつ更新する必要がある。さらに、  
206 クラスタの解釈を行う際には、後述するクラスタ重心が大きな役割を果たすが、相対座標が  
207 与えられないため、 $d_{mah}$  によって構築されたクラスタの特徴を可視化することができない。  
208 2つ目は、 $d_{mah}$  が仮設検定に基づき構築された指標である点である。Maharaj の仮設検定に  
209 おいては、2つの時系列間の距離が等しいという帰無仮説を利用するが、これは多次元時系  
210 列間の比較やクラスタという概念と整合的ではない。

211 非類似度  $d_{pic}$  と  $d_{cep}$  を計算する際には、まず区間時系列  $y_1, y_2, \dots, y_k$  の AR 係数を *arfit*  
212 関数 [6] を用いて求めた。 $d_{cep}$  の導出においては、式 (39) を用いてこれらの AR 係数からケ  
213 プストラムを計算した。求めたこれらの相対座標に基づき、*dist* 関数 [11] によって各時系列  
214 間の距離を求めた。そして、計算した距離を用いて *hclust* [11] によって階層クラスタリング  
215 を行った。階層クラスタリングの手法として重心を用いるものには、異なるクラスタの重心  
216 間の距離を用いてクラスタを統合する重心法と、クラスタを統合した際のクラスタ内の重心  
217 からの距離の分散が最も小さくなるクラスタを逐次的に形成する Ward 法が存在する。この  
218 うち、重心法でクラスタリングを行った場合、全ての時系列が1つのクラスタに統合されて  
219 しまったため、結果では Ward 法の結果のみを示す (図2のデンドログラムも参照のこと)。

220 作成されたクラスタの解釈には、各クラスタの重心を用いた。各クラスタの重心から再計  
221 算された AR 係数を用いて pc1 と pc2 のシミュレーションを行い、シミュレーション結果か  
222 らウミガメの輪郭形状として復元することで、クラスタを特徴付ける遊泳パターンを調べた  
223 (動画2、3を参照のこと)。

### 3 結果

幼体のアオウミガメ 10 匹の輪郭形状の変化に関する時系列  $pc1$  と  $pc2$  を AR モデルの AIC によって分割した。このとき、幼体のウミガメの運動は 10 フレーム (2 秒) 以下の周期運動によって構成されていたことから、AR モデルの最大次数は 20 に設定した。ここで、動画 1 の timeseries は分割された  $pc1$  と  $pc2$  の時系列を表し、プロットした frame(表記上は time) に対応する水槽での相対座標 (Position)、輪郭形状 (Shape)、 $pc$  平面上での値 ( $pc$  plot) を示す。<sup>10</sup> この動画から、AR モデルによる分割は、漕ぎ運動から羽ばたき運動、その逆、あるいはどちらとも取れない中間的な泳ぎへの変化点を抽出できたと言える。

得られた区間時系列どうしの距離を非類似度  $d_{pic}, d_{cep}$  を用いて計算し、Ward 法によるクラスタリングを行った。クラスタリングの結果、各非類似度に対してデンドログラム (図 1、2) を得た。非類似度間で生成されたデンドログラムには、3 つの違いが見られた。

1.  $d_{pic}$  を用いた方ではクラスタが急増する高さ 4 の前まで 4 つの均等なクラスタが存在するのに対し、 $d_{cep}$  を用いた方は 3 つの均等なクラスタが存在する点
2.  $d_{pic}$  を用いた方に比べ、 $d_{cep}$  を用いた方が高さの低いところでクラスタが急増している点
3.  $d_{pic}$  を用いた方に比べ、 $d_{cep}$  を用いた方が最も大きな 2 つのクラスタが統合されるまでの高さが高くなっている。

1 点目から、最適クラスタ数を  $d_{pic}, d_{cep}$  に対しそれぞれ 3、4 と設定した。これは、 $d_{cep}$  の方は早い段階でまとまりを作っているといえる。2 点目と 3 点目からは、 $d_{pic}$  の方が早い段階でクラスタが形成され、クラスタ間の非類似度が大きく計算されていることがわかる。

得られたクラスタ 1-3 の重心を、AR 係数 (piccolo)、あるいはケプストラムとして求めた。ケプストラムによって求めた重心は、式 (39) を用いることで AR 係数に変換した。得られた AR 係数を用いて、クラスタ 1-4 を特徴づける  $pc1$  と  $pc2$  をシミュレーションによって生成

---

<sup>10</sup> クラスタリングの結果も示されているが、動画が重くなるので 1000 ステップにトリミングしている。今回は分割にのみ注目されたい。

した。そして、得られた  $pc1$  と  $pc2$  の値から輪郭形状を復元することで、動画 2-1,2-2,2-3,2-4,3-1,3-2,3-3 を作成した。

$d_{pic}$  を用いたクラスタリングによって作成された動画において、クラスタ 1-4 の大きな違いは見受けられなかった。強いてあげれば、クラスタ 1 と 2 に比べ、3 と 4 はばたついており、flapping が少ない印象を受けるが、明確な違いとはいえない。 $d_{cep}$  でも、クラスタ 1 に比べて 2 と 3 がばたついており、flapping が少ない印象を受けるが、これも明確な違いとはいえない。この原因については、4. 議論で述べる。

## 4 議論

AR モデルの当てはまりから、時系列を直感に沿う形で分割することができた。本研究で用いたデータは小さい水槽で録られたものであり、行動のレパートリーに休息や潜水はみられず、視覚的にはほとんどが漕ぎ運動と羽ばたき運動から構成されていた。この点、視覚的な最適クラスタ数は 2 であり、大きなクラスタでは漕ぎ運動と羽ばたき運動が卓越することが予想できるが、 $d_{pic}$  と  $d_{cep}$  によって作成されたクラスタに明確な違いは見られなかった。この原因は、 $pc1$  と  $pc2$  を独立な時系列と考え、AIC の計算、クラスタリング、シミュレーションを行ったためだと考えられる。 $pc1$  と  $pc2$  は主成分分析によって導出されるが、形状に対する主成分分析であるため、時間方向には相関が存在する。例えば、羽ばたき運動が卓越する際は、 $pc2$  が大きく変動し、 $pc1$  はほとんど変化しないはずであるが、この関連性が捉えられていないために、不明瞭で不自然な輪郭形状の変化が得られたと考えられる。

この点を検証するため、時刻 10100-10400(目視で定常そうなところをピックアップした)でそれぞれ独立と仮定して推定した AR 係数 ( $arfit$  によって推定 [6]) と、相関を考慮して推定した MAR 係数 ( $marfit$  によって推定 [6]) を用いてシミュレーションを行った結果を、それぞれ動画 4-1 と 4-2 に示す。また、AR、MAR 係数の推定に用いた元動画から計算した  $pc1$  と  $pc2$  を用いて復元した形状変化を動画 4-3 に示す。 $pc1$  と  $pc2$  を独立と仮定してシミュレーションした動画 4-1 では、例えば右手を出した状態からいきなり左手が伸びる、両手を広げた状態からいきなり縮こまるなど、不自然な形状変化がたびたびみられ、全体的にばたつい



272 ている印象を受ける。一方、pc1 と pc2 の相関を考慮してシミュレーションした動画 4-2 で  
273 は、形状が滑らかに変化しており、動画 4-3 と見分けがつかない程度に元の運動が再現され  
274 ているといえる。また、時系列を見ても、pc1 の変動が大きい場合には pc2 の変動は小さい、  
275 あるいは pc2 の変動が大きい場合には pc1 の変動は小さい、といった漕ぎ運動と羽ばたき運  
276 動が再現されている。

277 これらのことから、直感に沿ったクラスタを形成するには、pc1 と pc2 の相関を考慮した  
278 MAR モデルの距離を用いてクラスタリングを行う必要があると考えられる。このためには、  
279 元の時系列の分割も MAR モデルによって行い、クラスタ間距離には  $d_{pic}$  を MAR 係数に拡  
280 張したものをを用いる必要がある。このとき  $d_{cep}$  の拡張の拡張も可能であれば比較対象とすべ  
281 きである。

282 本論では、視覚的に解釈が容易なことから、分割、クラスタリング、シミュレーションの  
283 全てにおいて pc1 と pc2 のみを用いた。しかし、抽出した形状変化そのものである動画 1 と、  
284 pc1 と pc2 によって復元された動画 4-3 を比較すると、重要な運動は捉えることができてい  
285 るが、十分に捉えているとも言い難い。それゆえ、MAR モデルを用いてもなお精度が良く  
286 ない場合には、pc3 以降の軸を考慮することも視野に入れる必要があるだろう。本論で用い  
287 た手法はほとんどが線形計算からなるため、pc3 以降の軸を考慮しても、計算時間が膨大に  
288 なることはないと考えられる。

289 本論によるクラスタリングによって、遊泳パターンの継続時間やその変化を定量化するこ  
290 とが可能である。図 3 に、 $d_{pic}$  によって作成したクラスタに代表される遊泳の継続時間 (上)  
291 とその変化の割合 (下) を示す。図に記載されている 1-4 はそれぞれクラスタ 1-4 に対応する。  
292 行動の継続時間や変化は、個体のストレス評価において重要となる (詳しくは次回のゼミ資  
293 料を参照のこと)。それゆえ、本研究の手法がブラッシュアップされ、直感に合うクラスタを  
294 形成することができれば、ストレス評価に応用できると考えられる。

## 295 References

- 296 [1] 竹村彰通. (2020). 現代数理統計学.

- [2] Atal, B. S. (1974). Effectiveness of linear prediction characteristics of the speech wave for automatic speaker identification and verification. the Journal of the Acoustical Society of America, 55(6), 1304-1312.
- [3] Berk, K. N. (1974). Consistent autoregressive spectral estimates. The Annals of Statistics, 489-502.
- [4] Bhansali, R. J. (1978). Linear prediction by autoregressive model fitting in the time domain. The Annals of Statistics, 224-231.
- [5] Kalpakis, K., Gada, D., & Puttagunta, V. (2001). Distance measures for effective clustering of ARIMA time-series. In Proceedings 2001 IEEE international conference on data mining (pp. 273-280). IEEE.
- [6] Kitagawa, G. (2020) Introduction to Time Series Modeling with Applications in R. Chapman & Hall/CRC.
- [7] Liao, T. W. (2005). Clustering of time series data—a survey. Pattern recognition, 38(11), 1857-1874.
- [8] Maharaj, E. A. (1996). A significance test for classifying ARMA models. Journal of Statistical Computation and Simulation, 54(4), 305-331.
- [9] Montero, P., & Vilar, J. A. (2015). TSclust: An R package for time series clustering. Journal of Statistical Software, 62, 1-43.
- [10] Piccolo, D. (1990). A distance measure for classifying ARIMA models. Journal of time series analysis, 11(2), 153-164.
- [11] R Core Team. (2024). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.

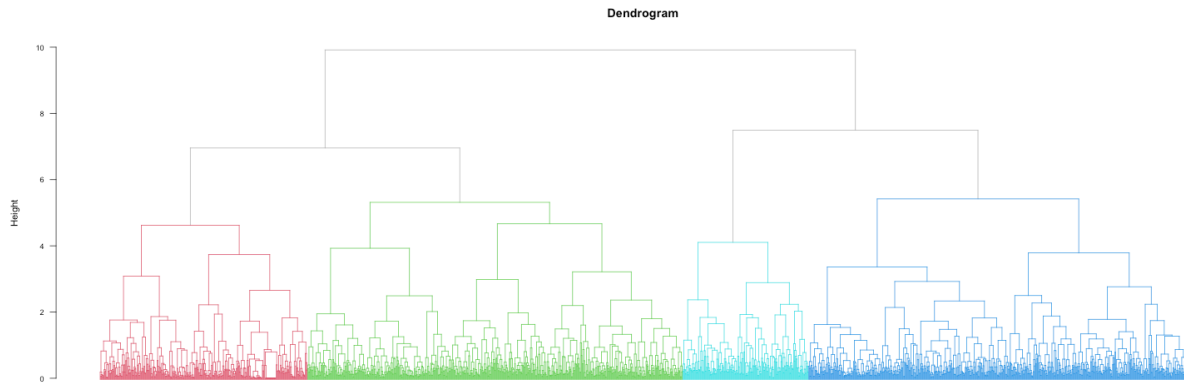


Figure 1:  $d_{pic}$  を用いて作成されたデンドログラム

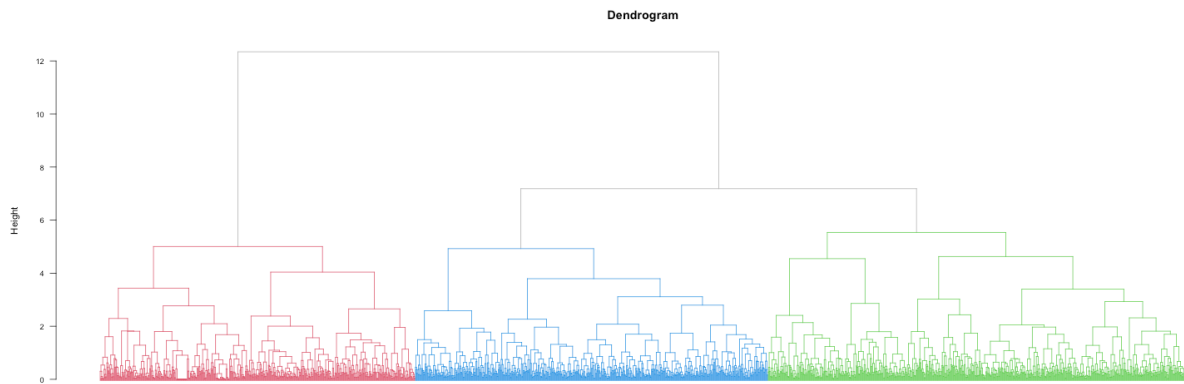


Figure 2:  $d_{cep}$  を用いて作成されたデンドログラム

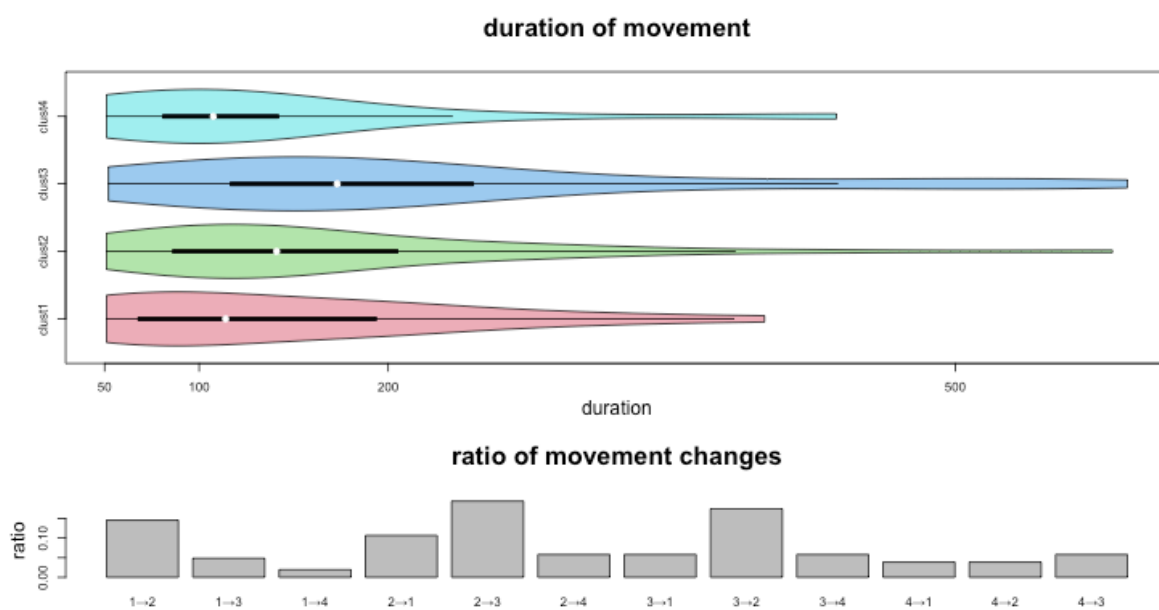


Figure 3:  $d_{pic}$  によって分類されたクラスターに代表される遊泳パターンの継続時間 (上) とその変化 (下)。図は Ekahi のもの。