

周期時系列の統計解析

(17) コサイン関数の変数クラスター分析

nino

2021年 2月 1日

複数のコサイン関数に変数クラスター分析を適用した場合、どのようなクラスターが形成されるのかについて検討した。その結果、コサイン関数は周期別にクラスターが形成され、また、それらクラスター内では位相差により分類されると考えられた。

周期時系列

周期時系列として、次のコサイン関数を用いた。

$$x_{T,\alpha} = \cos[2\pi(t+\alpha)/T] \quad (1)$$

ここで、 T : 周期, α : 位相, そして, t : 時間 ($t=0,1,\dots,m$) である。

具体例として、1hr間隔で、期間が2日間 ($n=24 \times 2=48\text{hr}$, $t=0,1,\dots,47\text{hr}$) の時系列について考える。時間と角度の関係は、 $1\text{hr}=360^\circ / 24\text{hr}=15^\circ$ に相当する。

例えば、周期 $T=24\text{hr}$ で位相 $\alpha=0\text{hr}$ の時系列 $x_{24,0}$ は、

$$x_{24,0} = \cos[2\pi(t+0)/24] \quad (2)$$

で表される。同様に、時系列 $x_{24,4}$, $x_{12,0}$, $x_{12,4}$ はそれぞれ次式で表される。

$$x_{24,4} = \cos[2\pi(t+4)/24] \quad (3)$$

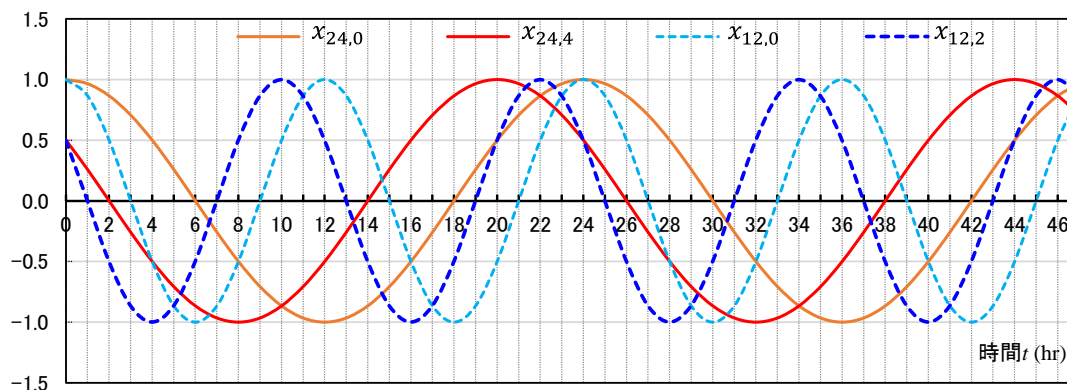
$$x_{12,0} = \cos[2\pi(t+0)/12] \quad (4)$$

$$x_{12,2} = \cos[2\pi(t+2)/12] \quad (5)$$

$x_{24,4}$ は周期が $x_{24,0}$ と同じ24hrだが位相はそれより4hr ($360^\circ \times 4\text{hr} / 24\text{hr} = 60^\circ$) 進み, $x_{12,0}$ は周期が12hrで位相が0hr, そして, $x_{12,2}$ は周期が $x_{12,0}$ と同じだが位相はそれより2hr ($360^\circ \times 2\text{hr} / 12\text{hr} = 60^\circ$) 進みである。

これら4つの時系列 $x_{24,0}$, $x_{24,4}$, $x_{12,0}$, $x_{12,2}$ (「Case 1」とする) を図1に示した。

図1 (Case 1)



Case 1 の変数クラスター分析

最初に, Case1について変数クラスター分析を適用し検討した. 距離関数は相関係数 r を含む $\sqrt{2(1-r)}$ を, また, クラスタリング手法としてウォード法を用いた. 距離関数 $\sqrt{2(1-r)}$ の意味については後述する.

表 1 にCase 1 の相関係数 r を示した. 表中では, $x_{24,0}$ と $x_{24,4}$ および $x_{12,0}$ と $x_{12,1}$ についてそれぞれ□で囲んである. また, 図 2 にクラスター数を 2 とした場合における樹形図を示した.

表 1

相関係数	$x_{24,0}$	$x_{24,4}$	$x_{12,0}$	$x_{12,2}$
$x_{24,0}$	1			
$x_{24,4}$	0.500	1		
$x_{12,0}$	0.000	0.000	1	
$x_{12,2}$	0.000	0.000	0.500	1

図 2

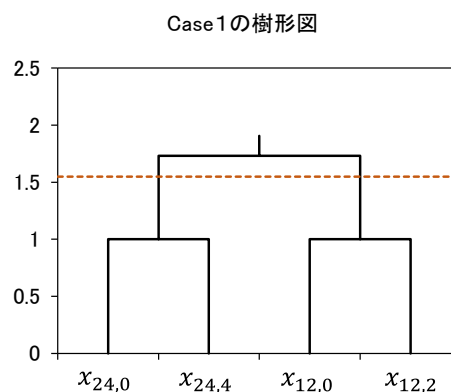


表 1 をみると, $x_{24,0}$ と $x_{24,4}$ および $x_{12,0}$ と $x_{12,2}$ はそれぞれ同じグループを形成し, 相関係数はともに0.5 ($=\cos 60^\circ$) を示した. 相関係数は周期時系列の位相差のコサインで表されるからである (参考文献 1). また, $x_{24,\alpha}$ と $x_{12,4\alpha}$ は周期が異なるため, 両者の相関係数はすべてゼロを示した (参考文献 2).

樹形図 (図 2) では, $x_{24,0}$ と $x_{24,4}$ および $x_{12,0}$ と $x_{12,1}$ はそれぞれクラスターを形成し, 両クラスター内での距離はともに1 ($=\sqrt{2(1-0.5)}$) で同じ値を示した. また, 両クラスター間における合併後の距離は1.73となった.

距離関数の意味について

今回用いた距離関数 $\sqrt{2(1-r)}$ はどのように解釈すれば良いのだろうか? 距離関数の意味について検討してみた.

クラスター分析では, ベクトル間の類似性を表す指標として, コサイン類似度を用いる場合がある. そして, このコサイン類似度を用いてユークリッド距離を表すことができる (詳細は参考文献 3, 4 を参照). ここでは, 2次元ベクトルの場合について概要を記した.

2点 $P(p_1, p_2)$, $Q(q_1, q_2)$ に対してベクトル $\vec{p}(p_1, p_2)$, $\vec{q}(q_1, q_2)$ とすると, この2点のユークリッド距離 $d(P, Q)$ は,

$$\begin{aligned}
 d(P, Q) &= |\vec{p} - \vec{q}| \\
 &= \sqrt{|\vec{p} - \vec{q}|^2} \\
 &= \sqrt{|\vec{p}|^2 - 2\vec{p} \cdot \vec{q} + |\vec{q}|^2}
 \end{aligned} \tag{6}$$

となる.

ここで、両ベクトルの内積 $\vec{p}, \vec{q} = |\vec{p}||\vec{q}|\cos(\vec{p}, \vec{q})$ を式(6)に代入すると、

$$d(P, Q) = \sqrt{|\vec{p}|^2 - 2|\vec{p}||\vec{q}|\cos(\vec{p}, \vec{q}) + |\vec{q}|^2} \quad (7)$$

となり、さらに、 \vec{p}, \vec{q} を正規化すると、 $|\vec{p}| = 1, |\vec{q}| = 1$ であるから、

$$d(P, Q) = \sqrt{2(1 - \cos(\vec{p}, \vec{q}))} \quad (8)$$

が得られた。式(8)は、正規化した場合におけるユークリッド距離 $d(P, Q)$ とコサイン類似度 $\cos(\vec{p}, \vec{q})$ の関係を示している。

一方、コサイン類似度 $\cos(\vec{p}, \vec{q})$ は相関係数 r に等しいから(参考文献5)、

$$d(P, Q) = \sqrt{2(1 - r)} \quad (9)$$

が成り立つ。式(9)は、今回用いた距離関数に等しい。このように、ベクトルを用いてクラスター分析におけるユークリッド距離が導出された。では、周期時系列を用いた場合は、どうだろうか？

同一周期 T の2つの時系列 x_{T, α_i} と x_{T, α_j} の距離 $d(x_{T, \alpha_i}, x_{T, \alpha_j})$ について考える。両者の相関係数 r と位相差 $\alpha_i - \alpha_j = \varphi$ の関係は $r = \cos\varphi$ であるから(参考文献1)、式(9)より次式が得られる。

$$d(x_{T, \alpha_i}, x_{T, \alpha_j}) = \sqrt{2(1 - \cos\varphi)} \quad (10)$$

式(10)を式(8)と対比すると、 $\cos(\vec{p}, \vec{q}) = \cos\varphi$ となる。すなわち、2つのベクトルの角度と2つの周期時系列の位相差は等しい(参考文献6)。

さらに、式(10)を三角関数の半角公式を用いて変形すると、

$$\begin{aligned} d(x_{T, \alpha_i}, x_{T, \alpha_j}) &= \sqrt{2(1 - \cos\varphi)} \\ &= \sqrt{2(1 - (1 - 2\sin^2(\varphi/2)))} \\ &= 2\sin(\varphi/2) \end{aligned} \quad (11)$$

となり、同一周期の時系列間の距離は位相差 φ を含むサイン関数で表された。コサイン類似度がコサイン関数で表されているのと対比している。

このサイン関数で表された距離とコサイン類似度はどのような関係にあるのかを調べるため、2次元ベクトル図を用いて検討した。

例として、図3に半径1の円弧を含むベクトル図を示した。

点Pと点Qは半径1の円弧上にあり、それぞれを原点からのベクトル \vec{p}, \vec{q} とすると、 $|\vec{p}| = |\vec{q}| = 1$ である。これは先述の正規化したベクトルの場合に等しい。また、 \vec{p} と \vec{q} のなす角度を $\angle POQ = \varphi$ とする。点Sは原点Oから線分 \overline{PQ} に垂線を下した交点、点Rは点Qから線分 \overline{OP} に垂線を下した交点である。したがって、 $\angle POS = \angle QOS = \varphi/2$ であり、 $\triangle OPQ$ は線分 \overline{PQ} を底辺とする二等辺三角形となる。

図 3

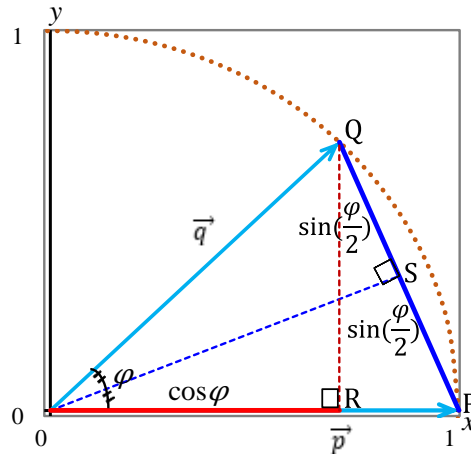


図 3 において、 \overline{OR} がコサイン類似度 $\cos\phi$ に相当し、 \vec{q} の x 軸 (\overline{OP}) への射影で表される。コサイン類似度は \vec{q} の \vec{p} との重なり程度と解釈できる。両ベクトルが完全に重なった時、すなわち、 $\phi=0^\circ$ の時にコサイン類似度は最大値 1 を示す。また、 $\phi=90^\circ$ の時にゼロを、 $\phi=180^\circ$ の時に最小値 -1 をとる。

一方、 $\overline{PS}=\overline{QS}=\sin(\phi/2)$ であるから、 $\overline{PQ}=2\sin(\phi/2)$ が成り立つ。これが距離に相当し、角度 ϕ と連動している。距離は \vec{q} の \vec{p} との乖離の程度を表しているといえる。距離は $\phi=0^\circ$ の時にゼロを、 $\phi=90^\circ$ の時に $2 \times \sin 45^\circ \approx 1.414$ を示す。また、 $\phi=180^\circ$ の時に最大値 2 を、 $\phi=-180^\circ$ の時に最小値 -2 を示す。距離は「コサイン類似度」と対比して「サイン乖離度」と命名できる。

なお、図 3 では、便宜上 \vec{p} を x 軸上に置いて説明したが、 \vec{p} と \vec{q} の関係はあくまで相対的なものである。

参考までに、表 2、表 3 に Case 1 の位相差 ϕ と距離 $2\sin(\phi/2)$ を示した。なお、表 2 の位相差 ϕ は式(2)～式(4)より明らかだが、 $\arccos(r)$ から求めることができる。

表 2

位相差(°)	$x_{24,0}$	$x_{24,4}$	$x_{12,0}$	$x_{12,2}$
$x_{24,0}$	0			
$x_{24,4}$	60.0	0		
$x_{12,0}$	90.0	90.0	0	
$x_{12,2}$	90.0	90.0	60.0	0

表 3

距離	$x_{24,0}$	$x_{24,4}$	$x_{12,0}$	$x_{12,2}$
$x_{24,0}$	0			
$x_{24,4}$	1.000	0		
$x_{12,0}$	1.414	1.414	0	
$x_{12,2}$	1.414	1.414	1.000	0

表 2 で留意すべき点は、24hr 周期と 12hr 周期の時系列間における位相差が $\phi=90^\circ$ で示されていることである。 $\phi=90^\circ$ は、同一周期の時系列間における位相差であるが、この場合は周期が異なるため、相関係数がゼロになる(表 1)ことに起因している(参考文献 6)。そのため、表 3 に示すように、異周期の時系列間の距離は、同一周期で位相差 $\phi=90^\circ$ の時系列間の距離 ($2\sin 45^\circ \approx 1.414$) と同じ値になった。また、 $x_{24,0}$ と $x_{24,4}$ および $x_{12,0}$ と $x_{12,2}$ の距離はともに $2\sin 30^\circ = 1$ を示した。

Case 1 に周期時系列 $x_{24,2}$ を加えた 5 つの周期時系列の変数クラスター分析

Case 1 に時系列 $x_{24,2}$ を加えた 5 つの時系列 ($x_{24,0}$, $x_{24,2}$, $x_{24,4}$, $x_{12,0}$, $x_{12,2}$: Case 2 とする) の変数クラスター分析について検討した.

$$x_{24,2} = \cos[2\pi(t+2)/24] \tag{12}$$

$x_{24,2}$ は周期が24hrで位相が2hr ($360^\circ \times 2\text{hr} / 24\text{hr} = 30^\circ$) の時系列である.

図 4 にCase 2 の時系列を示した.

図 4 (Case2)

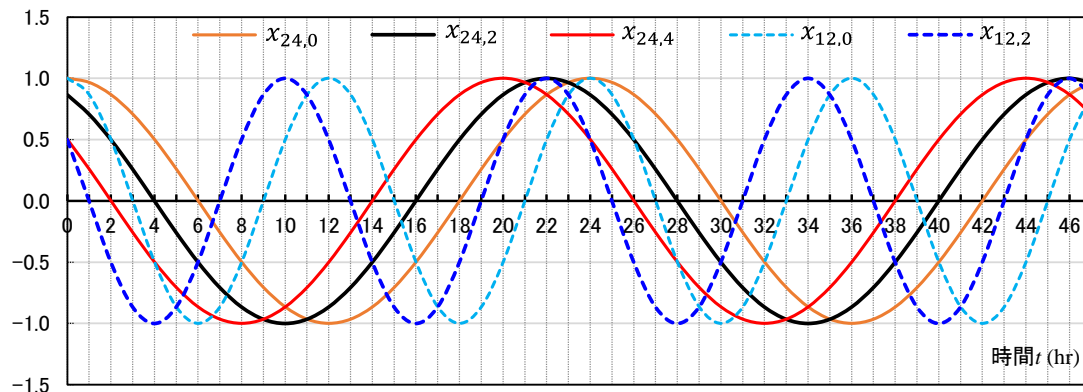


表 4, 表 5, 表 6 にCase 2 の相関係数, 位相差, 距離を示した. また, 図 5 にクラスター数を 2 とした場合の樹形図を示した.

表 4

相関係数	$x_{24,0}$	$x_{24,2}$	$x_{24,4}$	$x_{12,0}$	$x_{12,2}$
$x_{24,0}$	1				
$x_{24,2}$	0.866	1			
$x_{24,4}$	0.500	0.866	1		
$x_{12,0}$	0.000	0.000	0.000	1	
$x_{12,2}$	0.000	0.000	0.000	0.500	1

表 5

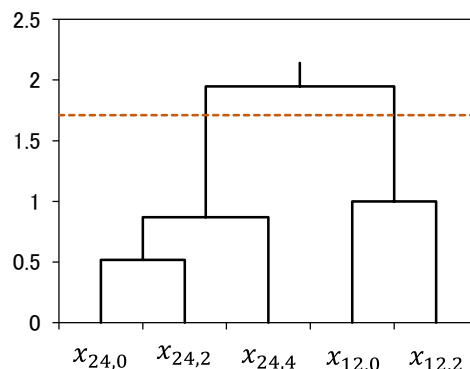
位相差($^\circ$)	$x_{24,0}$	$x_{24,2}$	$x_{24,4}$	$x_{12,0}$	$x_{12,2}$
$x_{24,0}$	0				
$x_{24,2}$	30.0	0			
$x_{24,4}$	60.0	30.0	0		
$x_{12,0}$	90.0	90.0	90.0	0	
$x_{12,2}$	90.0	90.0	90.0	60.0	0

表 6

距離	$x_{24,0}$	$x_{24,2}$	$x_{24,4}$	$x_{12,0}$	$x_{12,2}$
$x_{24,0}$	0				
$x_{24,2}$	0.518	0			
$x_{24,4}$	1.000	0.518	0		
$x_{12,0}$	1.414	1.414	1.414	0	
$x_{12,2}$	1.414	1.414	1.414	1.000	0

図 5

Case2の樹形図



相関係数 (表 4) によると, Case 2 は, Case 1 の場合と同様に, 24hr周期と12hr周期の 2 つのグループ間で無相関を示した. 位相差 (表 5) や距離 (表 6) も同様であった. 表 6 についてみると, $x_{24,0}$ と $x_{24,2}$ の距離および $x_{24,2}$ と $x_{24,4}$ の距離は, いずれの位相差も 30° であ

るので、 $0.518 (=2\sin 15^\circ)$ を示した。

樹形図(図5)では、Case 1の場合と同様に、周期別にクラスターが形成された。12hr周期の $x_{12,0}$ と $x_{12,2}$ の距離はCase 1と同様に1を示した。しかし、24hr周期のクラスターでは、まず、 $x_{24,0}$ と $x_{24,2}$ のクラスターを形成し、そのクラスターと $x_{24,4}$ とで新たにクラスターを形成している。 $x_{24,0}$ と $x_{24,2}$ の距離は0.518、そのクラスターと $x_{24,4}$ の合併後の距離は0.869を示した。したがって、24周期のクラスター内では3つの周期時系列は相互の位相差によって分類されると考えられた。さらに、周期別のクラスター間距離は1.95となった。

なお、このワード法による合併後の距離を試算したところ(参考文献7)、当然のことだが、一致した。試算した理由は、複数の市販ソフトやフリーソフトを使用して、距離を求めたところ異なる結果が得られることもあったためである。一般に統計ソフトによって得られた計算結果をそのまま使用する場合が多いと思うが、モデルデータなどを使用して、結果を確認しておくことが望ましい。

これまでの検討結果から、複数の周期時系列に変数クラスター分析を適用すると、時系列は周期別にクラスターを形成するとともに、それら同一周期の各クラスター内では位相差によりさらに分類されると考えられた。主成分分析もクラスター分析と同様な傾向を示しており(参考文献6)、両者には共通するものがある。クラスター分析と主成分分析はともに相関係数に基づいた統計手法であるためと考えられる。

参考文献

1. 物理のかぎしっぽ, 周期時系列の統計解析(1)相関係数と位相差
<http://hooktail.sub.jp/contributions/shuki01.pdf>
2. 物理のかぎしっぽ, 周期時系列の統計解析(13)周期時系列における主成分分析の意味
<http://hooktail.sub.jp/contributions/shuki13.pdf>
3. ユークリッド距離さえあればコサイン類似度が計算できる。
https://qiita.com/obake_kaiware/items/36104a479582063308f0
4. Cosine Distance as Similarity Measure in KMeans [duplicate]
<https://stats.stackexchange.com/questions/299013/cosine-distance-as-similarity-measure-in-kmeans>
5. 相関係数とコサイン類似度
<https://qiita.com/wsuzume/items/59ef9db4b3fb9d750dc1>
6. 物理のかぎしっぽ, 周期時系列の統計解析(12) 横断面データの主成分分析
<http://hooktail.sub.jp/contributions/shuki12.pdf>
7. 階層型クラスタリング(ワード法, 群平均法など)の計算過程
<https://analysis-navi.com/?p=1805>