

教師なし学習を用いたCOVID-19新規感染者数のクラスター分析

Cluster Analysis on the Number of COVID-19 Newly Infected People using Unsupervised Learning

小松 隆行*

Takayuki Komatsu

概要

本稿では、インターネット上に公開されている COVID-19 (新型コロナウイルス) に関するオープンデータに関して、相関分析とクラスター分析を行う。対象とするデータは、日本の全都道府県、日本を含むアジア諸国、およびヨーロッパ諸国の新規感染者数の日次データである。これらを各々時系列データと捉え、地域間の時系列データを統計学的に相関分析し、さらに機械学習の教師なし学習を用いてデータサイエンス的な視点から似たものを同一グループに分類するクラスター分析を試みる。

1. はじめに

2020年1月頃から2年を経た現在も継続している COVID-19 (新型コロナウイルス) のパンデミックは、人類の健康や社会的危機、経済的危機を生じさせている。ワクチン接種が進んでいるにもかかわらず、感染者の急拡大と収束を繰り返している状況にある。これに関連する情報は、数値データのオープンデータとして、全世界の多くの機関がインターネット上に公開しており、それらをダウンロードし分析や解析を行うことが可能となっている。これらのサイトから新規感染者や現在患者数、死者数、病床占有率など多種類のデータがオープンデータとして CSV 形式や Excel 形式でダウンロード可能である。その種類によっては、データ取得間隔の長さや欠損などがあるが、これらのデータの推移をグラフ化するなどして可視化したり、基本的な統計分析したりした結果や知見は、日本では行政やメディアのインターネットサイト⁽¹⁾などにおいて、世界の国と地域のデータは世界保健機関 (WHO) のサイト⁽²⁾で公開されている。北海道においても道内の大都市や振興局単位での感染状況のデータが公開されており、ダウンロードし解析することを可能にしている。

本報告の目的は、COVID-19 の新規感染者数のデータのある非定常の情報源から得られる実測値データと考え、これらの時系列データのみから統計学的手法とデータサイエンスの機械学習の手法を用

いることによって何らかの知見を得ることである。具体的には、相関分析や機械学習の教師なし学習と呼ばれている手法の Time Series k-means 法⁽³⁾によるクラスタリング (自動グループ分け) を試みて考察を行うことにする。

2. 全国47都道府県の新規感染者数に関する分析

まず、日本全国の47都道府県毎のデータ⁽¹⁾について分析してみる。ここでは、全国的に同時に急拡大した第5波よりも前の、都道府県毎の新規感染者数のデータの増減の傾向に差異が見られた第1波から第4波までの期間である2020年1月16日から2021年6月28日までのデータを使用することにする。

2.1 全国47都道府県間の相関分析

まず、相関分析を行う。第1波から第4波の期間、2020/1/16~2021/6/28のデータを用いて、都道府県間の相関係数を算出し、その値のヒートマップとして図1に示した。縦軸と横軸には、47都道府県名が配置され、座標の交差した場所に算出された相関係数が記載されている。都道府県名は、その所在地が概ね北から南の順に地方毎にまとめられて記載されている。この図1では、相関係数の大小に依存して該当するマス目の色が水色 (相関係数 0.0) から赤色 (相関係数 1.0) まで変化して表示されている。相関係数が0.9以上場合は、2つの変数間に

* 北海道科学大学未来デザイン学部メディアデザイン学科

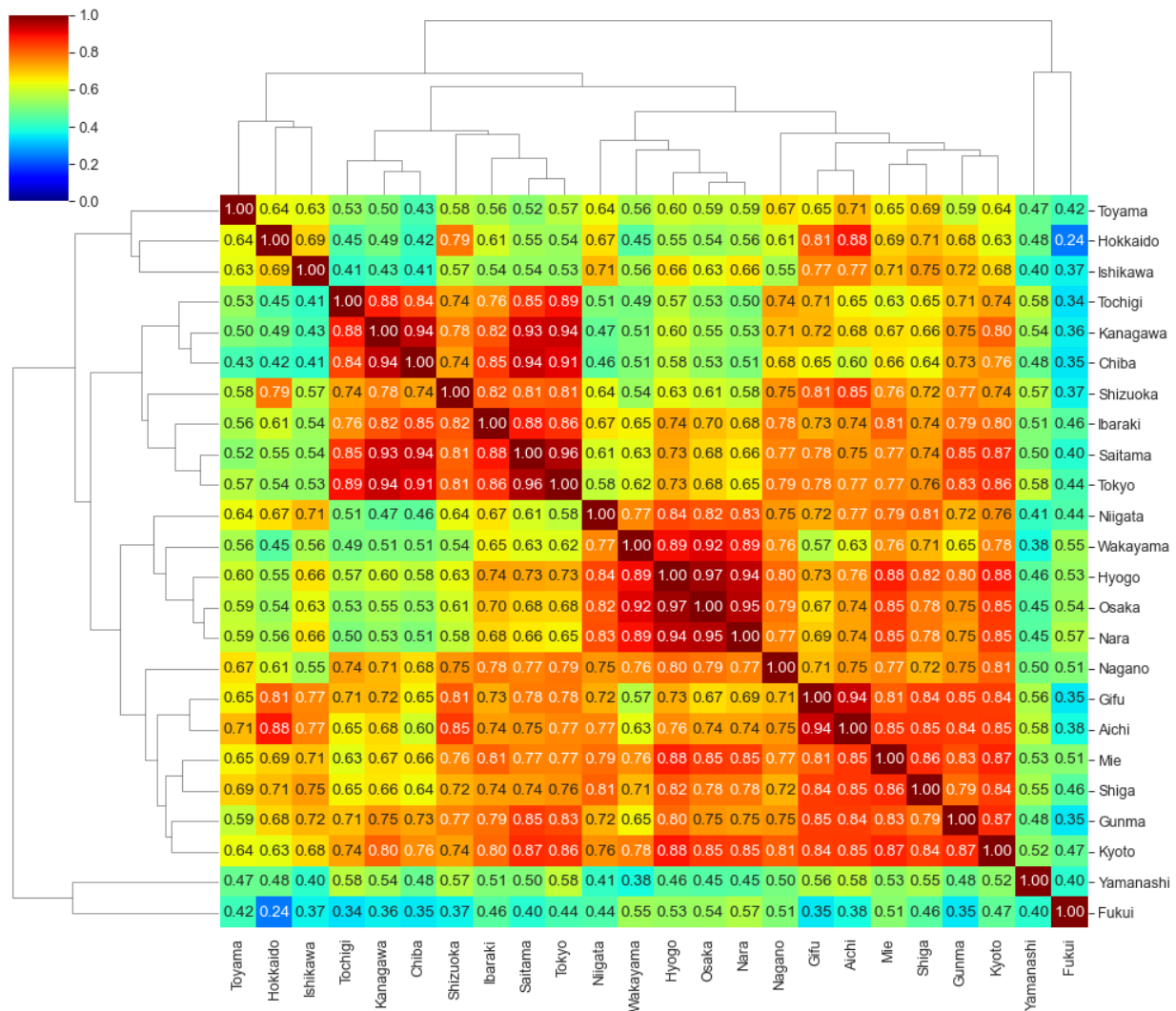


図2 関東/中部/関西の相関係数のヒートマップとデンドログラム (2020/1/16~2021/6/28)

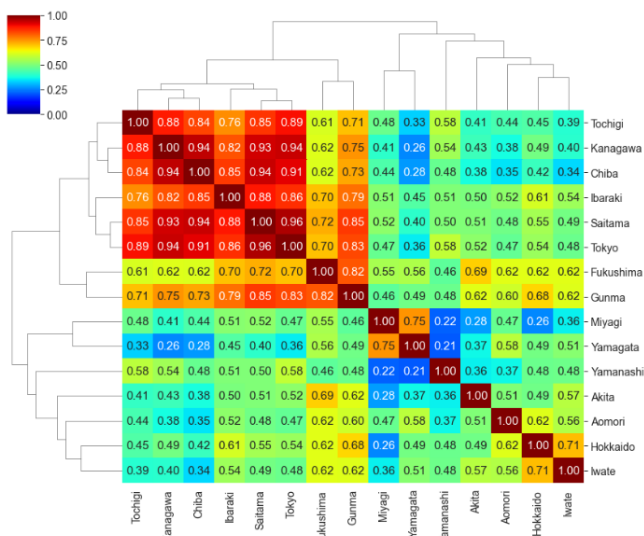


図3 北海道/東北/関東の相関係数のヒートマップとデンドログラム (2020/1/16~2021/6/28)

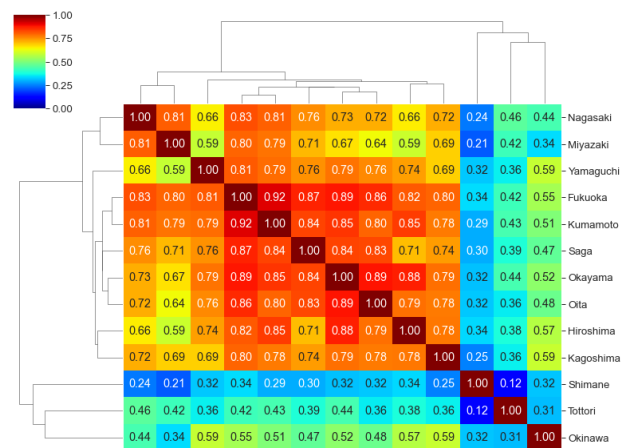


図4 中国/九州の相関係数のヒートマップとデンドログラム (2020/1/16~2021/6/28)

2.2 全国47都道府県のクラスタリング

ここではまず、全国47都道府県毎に平均と標準偏差で標準化された新規感染者数の14日平均の推移をTime series k-means法⁽³⁾(Euclid計量)を用いてクラスタリングを行った。クラスター数はエルボー法により10とした。図5は、同じクラスターに分類された都道府県毎の推移を、そのクラスター毎にまとめてプロットしたグラフである。横軸は日付ですべて同一であり、縦軸は標準化による標準偏差の何倍かを表している。それに対応してスケールと最大値がクラスター毎に異なることに留意する。

図5から、関東、東北と北陸、関西と隣接県、などがそれぞれ大きなクラスターを形成し、隣接県同士の鳥取と島根、宮城と山形がそれぞれ同じクラスターとなっていることが分かる。

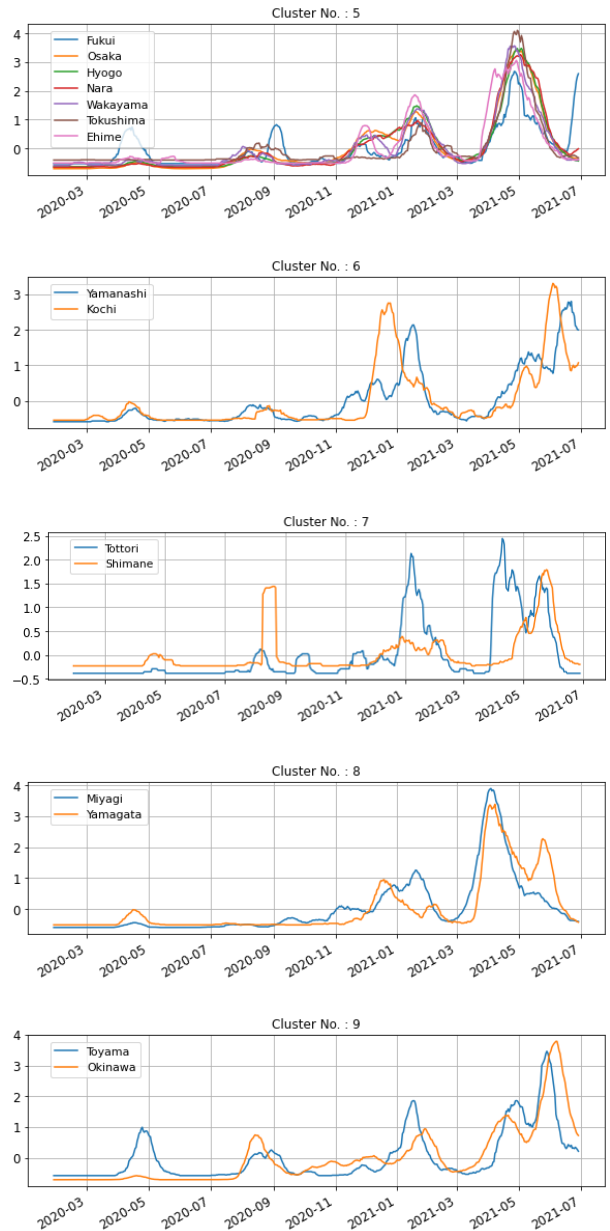
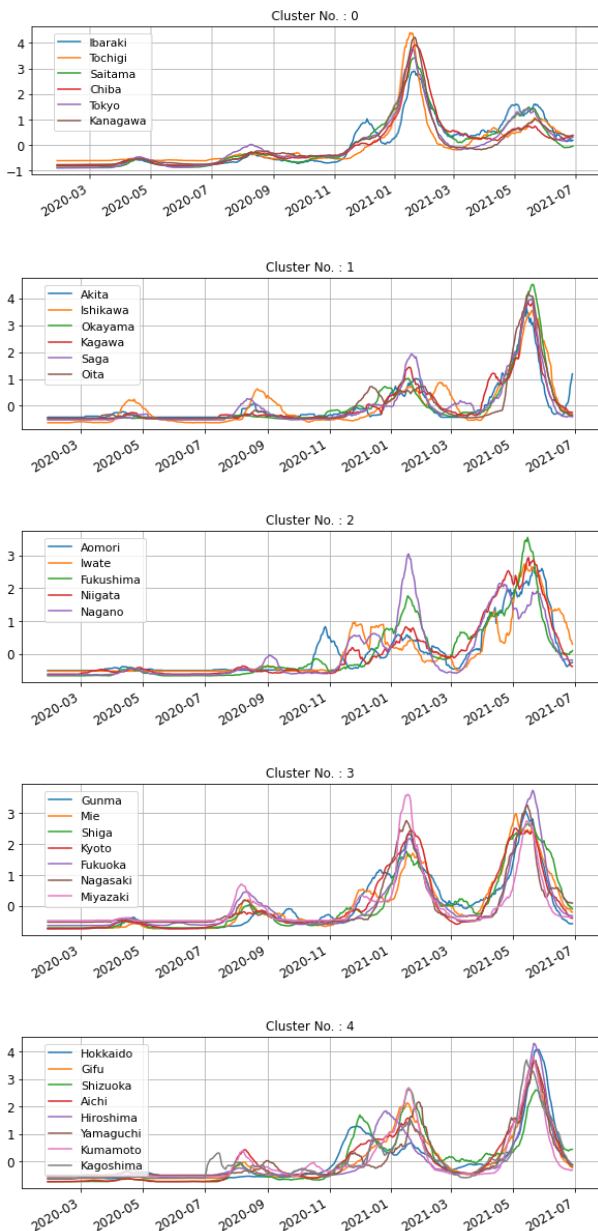


図5 全国の標準化新規感染者数14日平均のTime Series k-means法によるクラスタリング

次に、全国47都道府県毎の10万人当たりの標準化された新規感染者数の14日平均の推移を、Time Series k-means法⁽³⁾(Euclid計量)を用いてクラスタリングを行った結果を図6に示す。クラスター数はエルボー法により9とした。最大値が10を超えるクラスターでは、東京、大阪、北海道、沖縄が単独でクラスタリングされている。最大値が6以上10以下のクラスターでは、東京の隣接県の埼玉と千葉と神奈川、愛知と岐阜、広島と岡山という隣接県同士が、それぞれ同じクラスターに分類されている。また、関東と九州の大都市の周辺の県のクラスター、東北の隣接県と関西と四国の隣接県、大都市から比

較的遠い県のクラスターも形成されている。

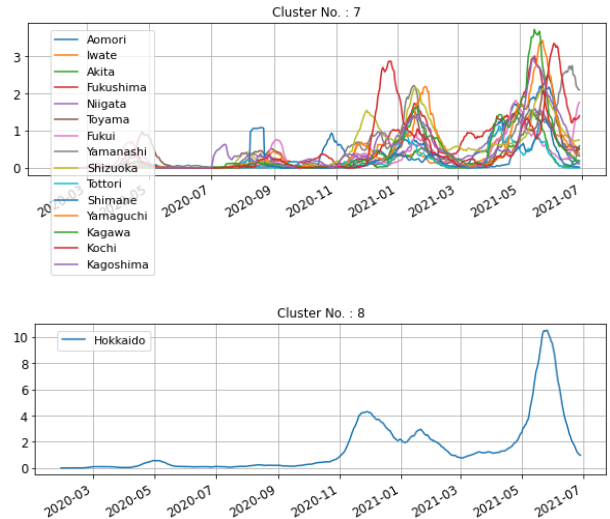
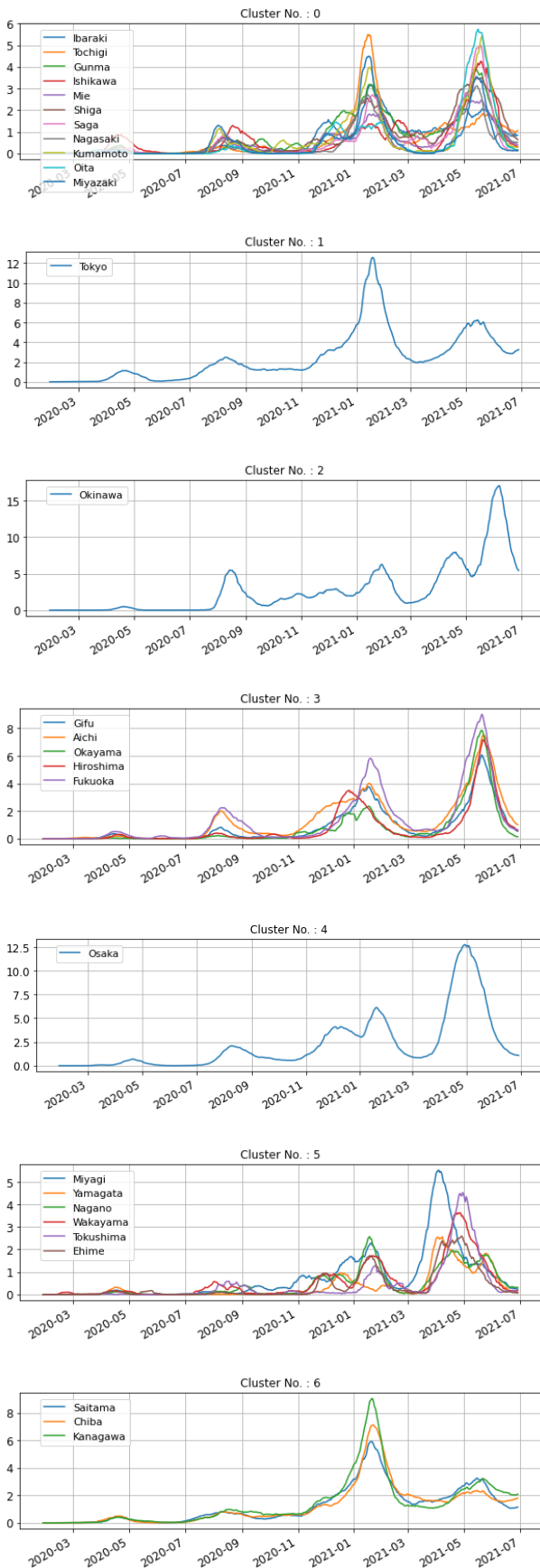


図 6 10万人当たりの標準化された新規感染者数の14日平均のTime Series k-means法によるクラスタリング (2020/1/16~2021/6/28)

3. アジア諸国とヨーロッパ諸国についての分析

ここでは、2章と同様にアジア諸国とヨーロッパ諸国について分析を試みる。データは世界保健機関 (WHO) のサイト⁽²⁾から取得した。使用データの期間は、2020年1月16日から2021年6月28日までである。図7と図10は、それぞれアジア諸国間、及びヨーロッパ諸国間の相関係数のヒートマップとデンドログラムである。また図8と図9、及び図11は、それぞれアジア諸国とヨーロッパ諸国の標準化された新規感染者数7日平均のTime Series k-Means法⁽³⁾によるクラスタリングの結果である。クラスター数は、エルボー法により8及び10とした。

アジアでは、単独国のクラスターを除いて、ヒートマップと共に考察すると、おおよそインドその周辺国、タイとカンボジアとマレーシアとその周辺国、中国の周辺国 (モンゴルやベトナムなど)、日本とインドネシアと韓国などの少し離れた島国のような国々で、それぞれクラスターが形成されている。

ヨーロッパでは、ドイツとオランダとデンマークとスウェーデン、イタリアとスイスとオーストリア、ブルガリアとセルビアとハンガリー、フランスとベルギー、スペインとポルトガル、フィンランドとノルウェー、などの隣接国が各々同じクラスターに分類されている。島国のアイルランドやアイスランドは、単独クラスターとなっている。ヒートマップと共に考察すると、ドイツとその北の隣接国、イタリアとオーストリア周辺、旧東欧諸国周辺などで、それぞれ強い相関を示している。

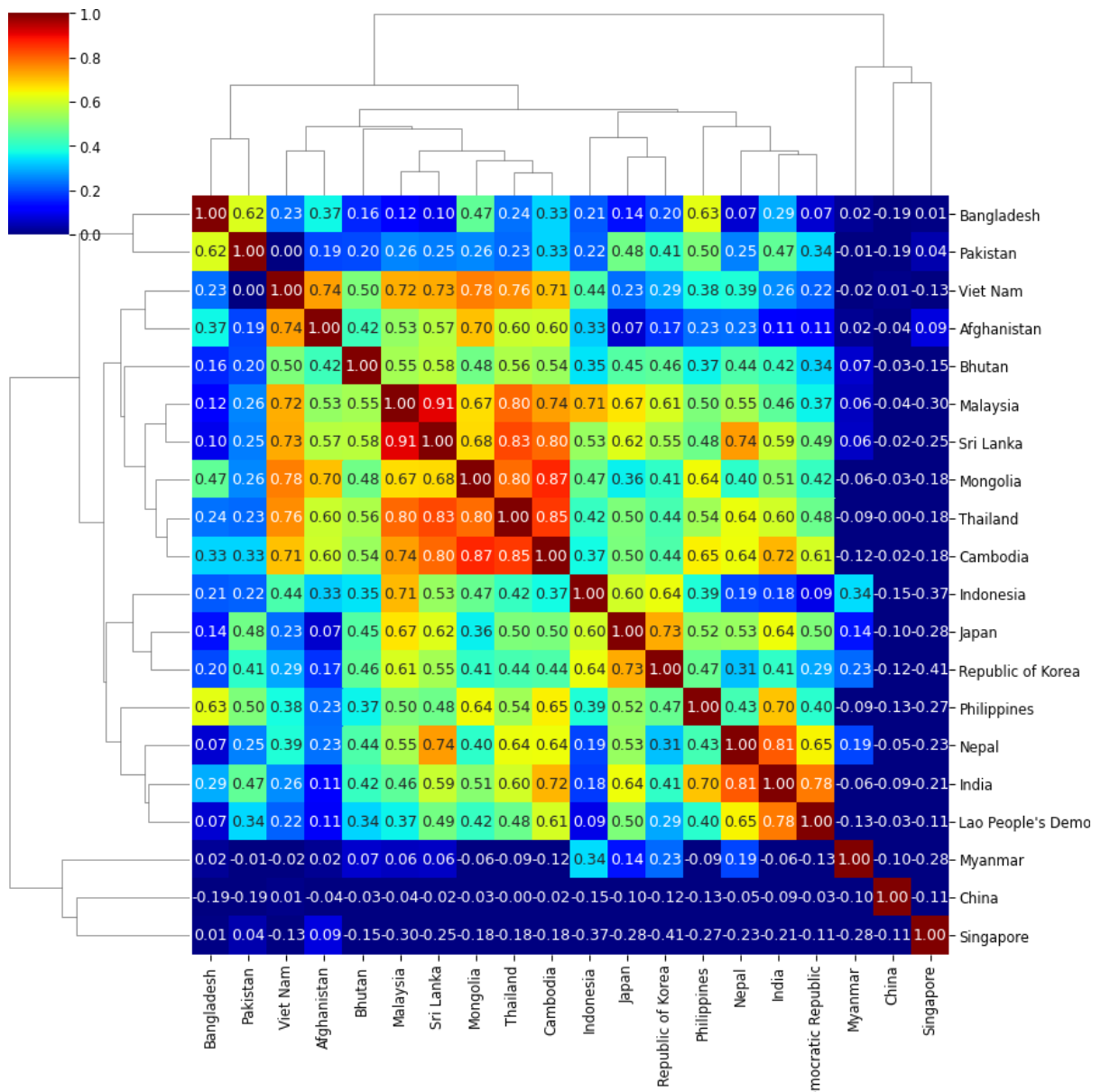


図7 アジア諸国間の相関係数のヒートマップとデンドログラム (2020/1/16~2021/6/28)

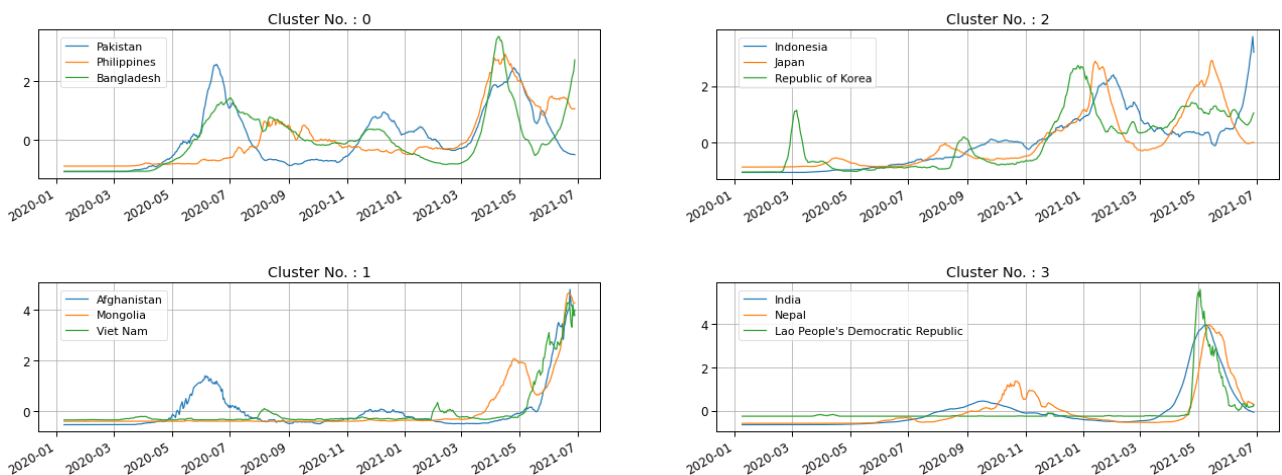


図8 アジア諸国の標準化新規感染者数7日平均のTime Series k-Means法によるクラスタリング(続く)

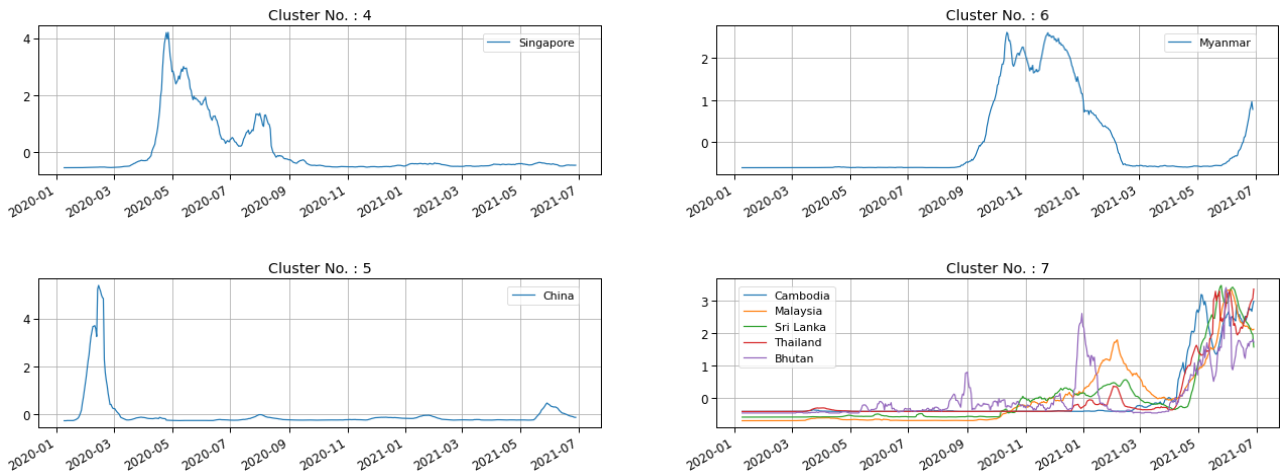


図9 アジアの標準化された新規感染者数7日平均のTime Series k-Means法によるクラスタリング

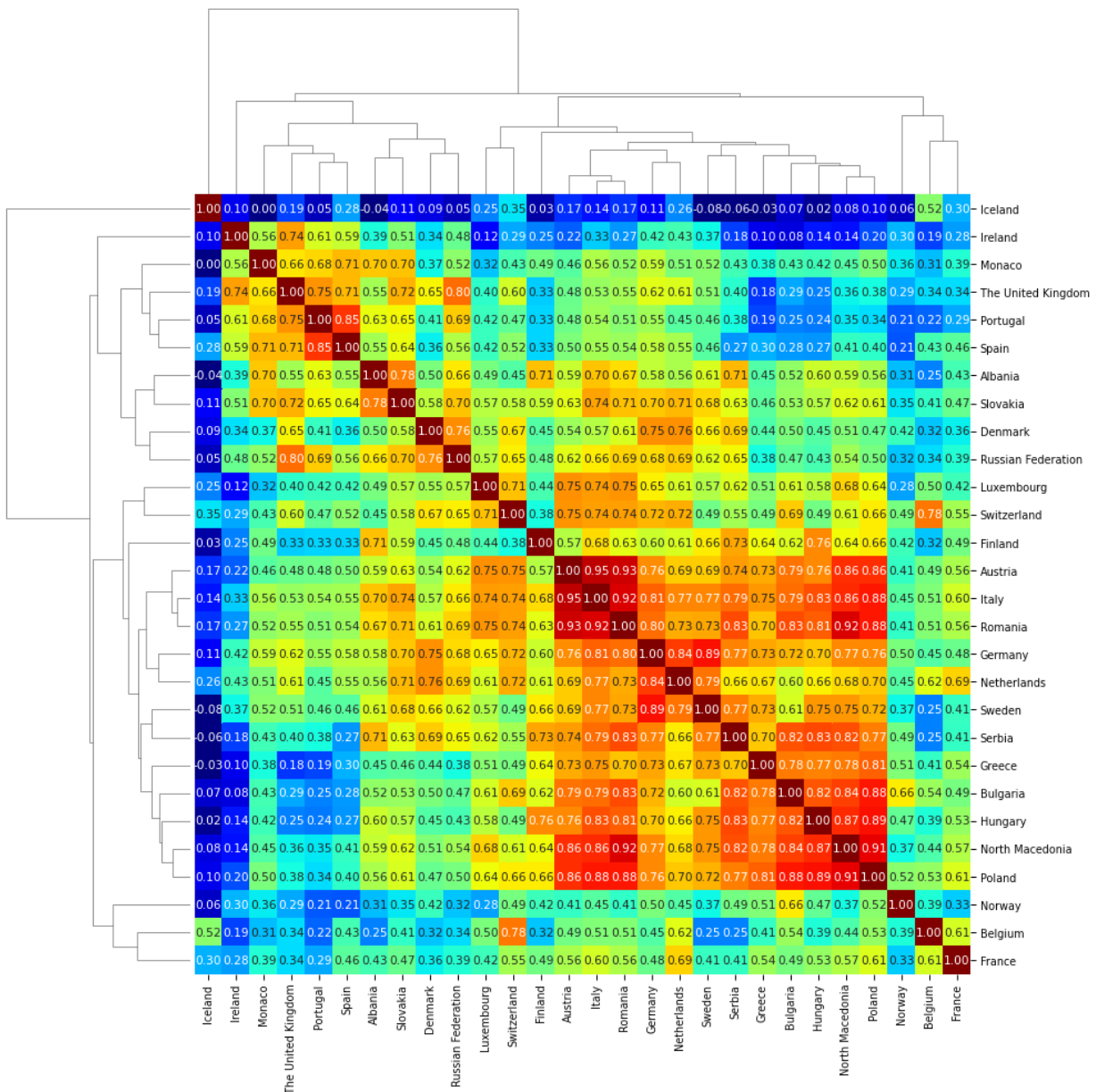


図10 ヨーロッパ諸国間の相関係数のヒートマップとデンドログラム (2020/1/16~2021/6/28)

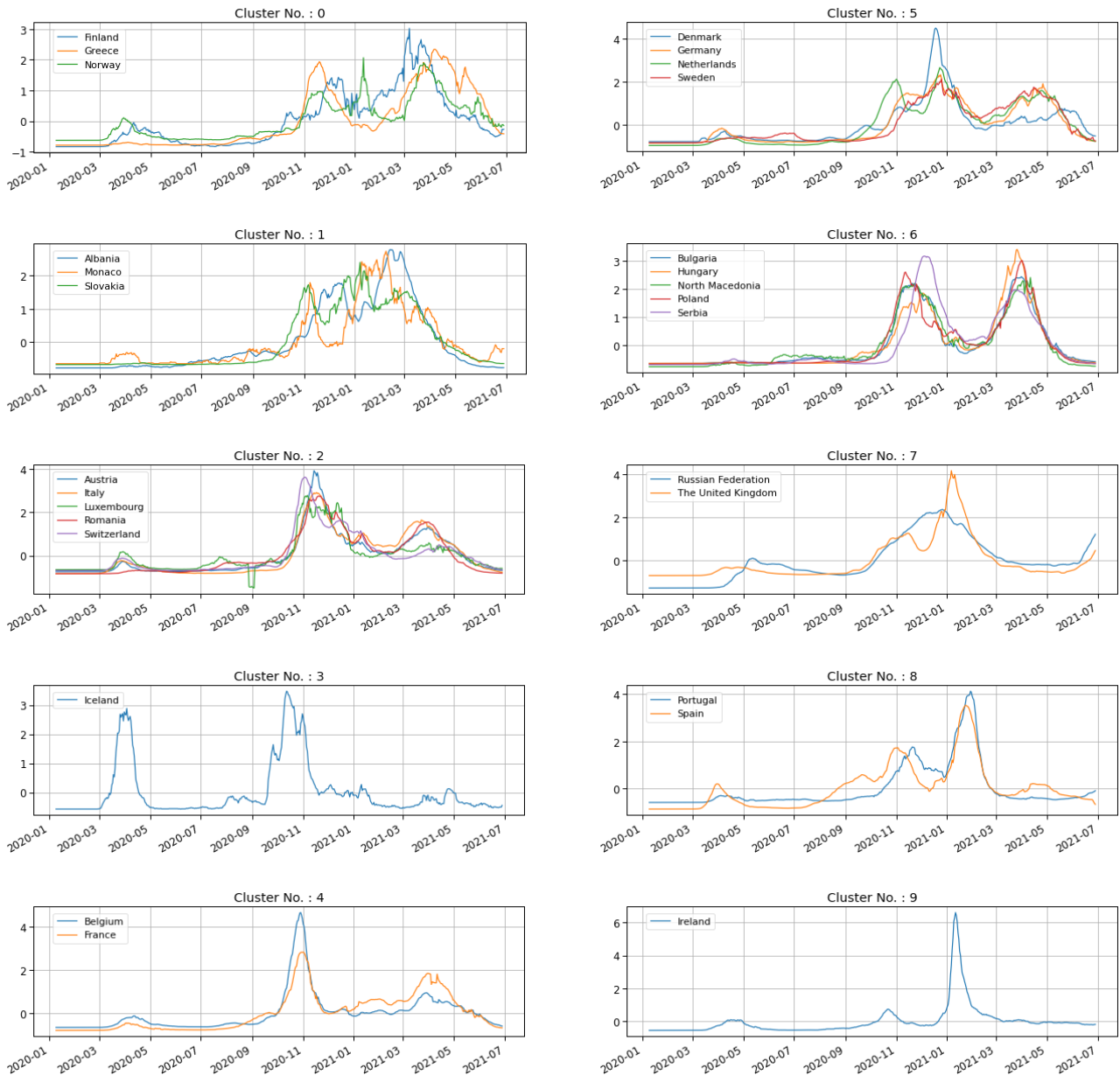


図 11 ヨーロッパの標準化された新規感染者数 7 日平均の Time Series k-Means 法によるクラスタリング

4. まとめ

COVID-19 の新規感染者数のオープンデータを時系列データと捉えて、相関分析と機械学習の教師なし学習によるクラスタリングを試みた。これらの分析だけでは、因果関係の有無を論じることはできないが、データサイエンスの視点からの分析におけるいくつかの知見は得られたと考えられる。今後は、新たなデータを加えた分析や、他のデータサイエンスや機械学習の分析手法との比較、それらを組み合わせた分析や応用が必要であると考えられる。

参考文献

(1) NHK特設サイト新型コロナウイルス: 新型コ

ロナデータ一覧：都道府県ごとの感染状況, 2021年6月30日, <https://www3.nhk.or.jp/news/special/coronavirus/data-widget/>.

(2) World Health Organization: WHO Coronavirus (COVID-19) Dashboard, 2022年2月27日, <https://covid19.who.int/info/>.

(3) Xiaohui Huang, Yunming Ye, Liyan Xiong, Raymond Y. K. Lau, Nan Jiang, Shaokai Wang: Time series k-means: A new k-means type smooth subspace clustering for time series data, Information Sciences, Volumes 367-368, 1 November, pp. 1-13, 2016.