

# Верификация кластеризации временных рядов

Кривонос Анна

Московский государственный университет им. М. В. Ломоносова

*Курс:* Моя первая научная статья

*Научный руководитель:* О. В. Сенько

2023

Важнейшими задачами эпидемиологии являются исследования влияния различных факторов на ход эпидемического процесса, а также прогнозирование развития эпидемии. Например, заболеваемость коронавирусной инфекцией (COVID-19), поразившей практически весь мир в 2019-2021 годах, протекала в разных регионах и странах мира по-разному в зависимости от состояния систем здравоохранения, климатических, социально-экономических, демографических условий, других характеристик регионов. Для решения обеих задач могут быть применены современные методы машинного обучения и анализа данных. Целью настоящей работы является поиск оптимальной схемы использования кластерного анализа, являющегося популярным и эффективным инструментом современного анализа данных, для изучения эпидемического процесса.

## Постановка задачи

В качестве данных в данной задаче используются кривые темпа роста Covid-19 для различных стран мира, а также для отдельных регионов России.

Необходимо провести кластеризацию временных рядов, которая позволяла бы оценить ход эпидемиологического процесса и влияние на него различных факторов. Для этого проводится верификация полученной кластеризации. Для верификации кластеризации рассматривается подход, основанный на проверке нулевой гипотезы о равновероятности различных соответствий мер сходства между двумя временными рядами.

## Описание метода кластеризации

На первом этапе вычислялась мера сходства между всевозможными парами эпид-кривых  $S_i$  и  $S_j$  через подбор лага  $l$  из отрезка  $[0, 20]$ .

Для каждого  $l$  вычислялся коэффициент корреляции  $p_l^+$  между рядами  $S_i(0), \dots, S_i(n-l)$  и  $S_j(l), \dots, S_j(n)$  и коэффициент корреляции  $p_l^-$  между рядами  $S_i(l), \dots, S_i(n)$  и  $S_j(0), \dots, S_j(n-l)$ .

$$\rho(S_i, S_j) = \max\{p_0^+, p_0^-, \dots, p_{20}^+, p_{20}^-\}$$

Пусть  $P_{m \times m}$  является матрицей сходства  $m$  эпид-кривых. На диагонали симметричной матрицы  $P_{m \times m}$  находятся единицы, а внедиагональными элементами являются максимальные коэффициенты корреляции  $\rho(S_i, S_j)$ , рассчитанные согласно приведённой выше процедуре.

## Описание метода кластеризации

После подсчёта мер близости между кривыми использовался метод иерархической агломеративной кластеризации. В качестве меры сходства двух групп эпидкривых  $G'$  и  $G''$  использовалось среднее значение меры сходства между эпидкривыми из разных групп:

$$P(G', G'') = \frac{1}{m' m''} \sum_{i=1}^{m'} \sum_{j=1}^{m''} p(S_i, S_j)$$

Процесс слияния кластеров прекращался, если мера сходства  $P$  между любыми двумя кластерами в текущей кластеризацией не окажется ниже 0.5.

## Метод верификации

Пронумеруем элементы матрицы  $P_{m \times m}$ , находящиеся выше диагонали.

Пусть  $I$  взаимно однозначное отображение множества  $\{(i, j) | i, j = 1, \dots, m, i < j\}$  в  $\{1, \dots, M\}$ , где  $M = \frac{m(m-1)}{2}$

Пусть  $f$  – перестановка элементов множества  $\{1, \dots, M\}$ .

По перестановке  $f$  строится матрица  $P_{m \times m}^f$ . Пусть  $k = I(i, j)$ .

Тогда элементу матрицы  $P_{m \times m}^f$  в позиции  $(i, j)$  присваивается элемент  $P_{m \times m}$ , который находится в позиции  $(i^*, j^*) = I^{-1}[f(k)]$ .

Элементы ниже главной диагонали заполняются симметричным образом.

## Метод верификации

Значения индикатора качества  $I(C)$  для кластеризации  $C$ , полученной по исходной матрице сходства  $P_{m \times m}$  сравнивается со значениями индикатора качества  $I(C^f)$  для кластеризаций  $C^f$ , по матрицам сходства  $P_{m \times m}^f$ , генерируемым по  $\tilde{f}$  - случайному подмножеству перестановок элементов множества  $\{1, \dots, M\}$ .

В качестве  $\rho$ -значения используется доля перестановок, при которых значение индикатора качества  $I(C^f)$  для кластеризации  $C^f$  достигает или превосходит значения индикатора качества  $I(C)$  для кластеризации  $C$ ,

$$\rho = \frac{|\{f \in \tilde{f} | I(C^f) \geq I(C)\}|}{|\tilde{f}|}$$

## Метод верификации

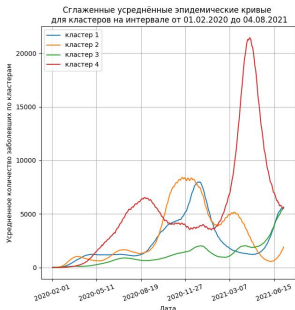
В качестве индикатора качества предлагается воспользоваться коэффициентом силуэта.

$$I(n) = \frac{b(n) - a(n)}{\max\{a(n), b(n)\}},$$

где  $a(n)$  - среднее расстояние от  $n$  до другого объекта внутри кластера,  $b(n)$  - среднее расстояние от  $n$  до объекта в другом кластере. Среднее значение всех силуэтов называют коэффициентом силуэта. В нашем случае в качестве метрики для оценки силуэта используется максимальный коэффициент корреляции Пирсона с лагом:

$$d_{i,j} = \rho_{\max}(S_i, S_j)$$





- ▶ кластер 1: Великобритания, Россия, США, Португалия и др.
- ▶ кластер 2: Австрия, Германия, Италия, Сербия, Болгария, Румыния и др. Географически локализованы главным образом в Европе.
- ▶ кластер 3: Алжир, Марокко, Бангладеш, Вьетнам и др. Страны с тропическим климатом.
- ▶ кластер 4: Аргентина, Колумбия, Парагвай, Уругвай, Камбоджа и др. Географически локализована в Южной Америки.

Для данных о динамике заболеваемости за временной промежуток от 22-01-2020 до 05-08-2021 среднее значение всех вычисленных силуэтов при разделении стран на 4 обозначенных кластера составило 0.36.

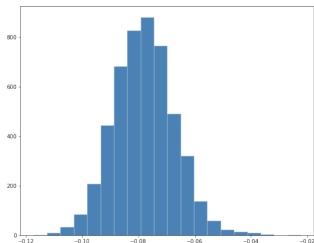


Рис. 5: Коэффициенты силуэта для 5000 различных экспериментов.

Для верификации результата предлагается вычислить коэффициент с помощью использования перестановок. Осуществим 5000 различных перестановок  $f$  и посчитаем коэффициенты силуэта. Максимальное значение коэффициента силуэта равно 0.02.

Проведем кластеризацию для случайных сгенерированных временных рядов из нормального распределения. Полученный коэффициент силуэта на основе алгоритма 3 для нашей кластеризации составил  $-0.004$ .

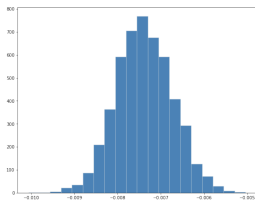


Рис. 6: Коэффициенты силуэта для 5000 различных экспериментов для случайных сгенерированных рядов.

Не смотря на отрицательное значение коэффициента силуэта он все же остается выше, чем коэффициенты для полученных перестановок - максимальный коэффициент силуэта для 5000 перестановок составил  $-0.005$ .

В работе описан метод кластеризации временных рядов, приведен способ оценки значимости кластеризации. В ходе экспериментов показана хорошая работа на эпидемиологических кривых Covid-19. В дальнейшем планируется усовершенствовать метод кластеризации с помощью введения остановки слияния кластеров в зависимости от вычисляемых  $\rho$ -значений для каждого кластера на каждом этапе слияния.