

# Пакет 'отримле'

29 мая 2021 г.

Тип Пакет

Название Надежная кластеризация на основе моделей

Описание Выполняет надежный кластерный анализ с учетом выбросов и шума, которые невозможно подобрать. тед любым кластером. Данные моделируются смесью гауссовых распределений. и шумовой составляющей, которая представляет собой неправильное равномерное распределение, охватывающее всю Европу. клидовое пространство. Параметры оцениваются по (псевдо) максимальной вероятности. Это подходит- с помощью алгоритма EM-типа. См. Коретто и Хен- ниг (2016) <doi: 10.1080 / 01621459.2015.1100996> и Коретто и Хен- ниг (2017) <<https://jmlr.org/papers/v18/16-382.html>>.

Версия 2.0

Дата 2021-05-28

Импортирует статистику, утилиты, графику, grDevices, mvtnorm, parallel, foreach, doParallel, robustbase, mclust

Лицензия GPL (> = 2)

LazyData ИСТИНА

Требуется компиляция нет

Автор Пьетро Коретто [aut, cre] (Домашняя страница: <<https://pietro-coretto.github.io/>>),  
Кристиан Хенниг [aut] (Домашняя страница: <<https://www.unibo.it/sitoweb/christian.hennig/en>>)

Сопровождающий Pietro Coretto <[pcoretto@unisa.it](mailto:pcoretto@unisa.it)>

Репозиторий CRAN

Дата / Публикация 2021-05-29 06:30:02 UTC

## Задокументированные темы R:

|                         |                   |
|-------------------------|-------------------|
| денежная купюра .....   | <a href="#">2</a> |
| генератор.отримле. .... | <a href="#">3</a> |
| InitClust. ....         | <a href="#">4</a> |
| kerndenscluster. ....   | <a href="#">6</a> |
| kerndensmeasure. ....   | <a href="#">7</a> |
| kerndensp. ....         | <a href="#">8</a> |

|                     |                        |
|---------------------|------------------------|
| 2                   | денежная купюра        |
| kmeanfun. ....      | <a href="#">10</a>     |
| отримле. ....       | <a href="#">11</a>     |
| отримлег. ....      | <a href="#">15</a>     |
| отримлесимг. ....   | <a href="#">17</a>     |
| сюжет.отримле. .... | <a href="#">21 год</a> |
| сюжет. ....         | <a href="#">24</a>     |
| ободок. ....        | <a href="#">25</a>     |
| Индекс              | <a href="#">30</a>     |

денежная купюра                      Данные о швейцарских банкнотах

Описание

Данные из таблиц 1.1 и 1.2 (стр. 5-8) Flury и Riedwyl (1988). Есть шесть измерений  
Сделано на 200 швейцарских банкнотах (старые швейцарские 1000 франков). Банкноты относятся к двум классам  
одинаковый размер: подлинный и поддельный.

Применение

данные (банкнота)

Формат

Data.frame размером 200x7 со следующими переменными:

Класс фактор с классами: подлинный, поддельный

Длина Длина банкноты (мм)

Левая ширина левого края (мм)

Правая ширина правого края (мм)

Нижняя ширина нижнего поля (мм)

Верхняя ширина верхнего поля (мм)

Диагональ Длина диагонали (мм)

Источник

Флури, Б. и Ридвил, Х. (1988). Многомерная статистика: практический подход. Лондон: Чепмен & Зал.

## Описание

При этом используются данные и выходные данные [otrimle](#) или [rimle](#) для создания нового искусственного набора данных размера исходных данных с использованием шума и пропорций кластера из выходных данных кластеризации. Кластеры затем генерируются из многомерных нормальных распределений с параметрами, оцененными с помощью [otrimle](#), шум генерируется передискретизацией из того, что оценивается как компонент шума с заданными весами апостериорной вероятностью того, что все наблюдения являются шумом. См. Хенниг и Коретто (2021 г.).

## Применение

generator.otrimle (данные, подходят)

## Аргументы

|                 |   |
|-----------------|---|
| данные          | что-то, что можно принудить к матрице. Набор данных.                |
| соответствовать | выходной объект <a href="#">otrimle</a> или <a href="#">обода</a> . |

## Значение

Список с данными компонентов, кластеризация.

|        |                                 |
|--------|---------------------------------|
| данные | матрица сгенерированных данных. |
|--------|---------------------------------|

|    |  |
|----|--|
| CS | вектор целых чисел. Индикатор кластеризации. |
|----|--|

## Авторы)

Кристиан Хенниг <[christian.hennig@unibo.it](mailto:christian.hennig@unibo.it)> [https://www.unibo.it/sitoweb/christian\\_hennig/en/](https://www.unibo.it/sitoweb/christian_hennig/en/)

## Рекомендации

Хенниг, К. и П. Коретто (2021 г.). Подход адекватности для определения количества кластеров для OTRIMLE кластеризация на основе робастной смеси Гаусса. Появиться в Австралии и Новой Зеландии Статистический журнал, <https://arxiv.org/abs/2009.00921>.

## Смотрите также

[kerndensp](#), [kerndensmeasure](#), [otrimle](#), [оболок](#)

## Примеры

```
данные (банкнота)
selectdata <- c (1: 30,101: 110,117: 136,160: 161)
set.seed (555566)
x <- банкнота [selectdata, 5: 7]
ox <- otrimle (x, G = 2, ncores = 1)
str (generator.otrimle (x, ox))
```

## Описание

Вычисляет начальное назначение кластера на основе комбинации шума на основе ближайшего соседа. и агломеративная иерархическая кластеризация на основе критериев максимального правдоподобия для Gaussianские модели смесей.

## Применение

InitClust (данные, G, k = 3, knnd.trim = 0,5, modelName = VVV)

## Аргументы

|                 |  |
|-----------------|--|
| данные          | Числовой вектор, матрица или фрейм данных наблюдений. Строки соответствуют наблюдения и столбцы соответствуют переменным. Категориальные переменные и NA значения не допускаются.  |
| грамм           | Целое число, определяющее количество кластеров.  |
| k               | Целое число, определяющее количество рассматриваемых ближайших соседей на каждую используемую точку. для шага шумоподавления (см. Подробности).  |
| knnd.trim       | Число в [0,1), которое определяет долю точек, инициализированных как шум. Обычно knnd.trim <= 0,5 (см. Подробности).   |
| название модели | Строка символов, указывающая модель ковариации, которая будет использоваться. Возможные модели находятся:<br>"E": равная дисперсия (одномерная)<br>"V": сферический, с переменной дисперсией (одномерный)<br>«EP»: сферический, равного объема<br>«VP»: сферическая, неодинакового объема.<br>"EEE": эллипсоид, равный объем, форма и ориентация.<br>«VVV»: эллипсоид, переменного объема, формы и ориентации (по умолчанию).<br>Смотрите подробности. |

## InitClust

5

## Подробности

Инициализация основана на Coretto and Hennig (2017). Сначала выполняются два шага:

Шаг 1 (шаг шумоподавления): для каждой точки данных вычислить расстояние до ближайших соседей (k-NN).

Все точки с k-NN больше, чем (1-knnd.trim)-квантиль k-NN, инициализируются как шум. Интерпретация k заключается в том, что (k-1), но не k, точки, расположенные близко друг к другу, все еще могут интерпретироваться как шум или выбросы

Шаг 2 (этап кластеризации): выполните иерархическую кластеризацию на основе модели (MBHC), предложенную в Фрейли (1998). Этот шаг выполняется с помощью [hc](#). Входной аргумент modelName передается [hc](#). См. Подробности о [hc](#) для получения более подробной информации.

Если на предыдущем шаге 2 не удалось предоставить кластеры G, каждый из которых содержит не менее 2 различных точек данных, выполняется замена классической иерархической кластеризацией, реализованной в [hclust](#). Наконец, если [hclust](#) не укажет допустимый раздел, пробуются до десяти случайных разделов.

## Значение

Целочисленный вектор, определяющий начальное назначение кластера, где 0 обозначает шум / выбросы.

## Рекомендации

Фрейли, К. (1998). Алгоритмы гауссовской иерархической кластеризации на основе моделей. SIAM Журнал на Научные вычисления 20: 270-281.

П. Коретто и К. Хенниг (2017). Последовательность, отказоустойчивость и алгоритмы для надежной неправильной кластеризации максимального правдоподобия. Журнал исследований в области машинного обучения, Vol. 18 (142), С. 1-39. <https://jmlr.org/papers/v18/16-382.html>

## Авторы

Пьетро Коретто <[pcoretto@unisa.it](mailto:pcoretto@unisa.it)> <https://pietro-coretto.github.io>

## Смотрите также

[hc](#)

## Примеры

```
## Загрузить данные о швейцарских банкнотах
данные (банкнота)
x <- банкнота [, - 1]

## Начальные кластеры с аргументами по умолчанию
init <- InitClust (данные = x, G = 2)
печать (инициализация)

## Выполнить отримле
a <- otrimle (data = x, G = 2, initial = init,
             loglcd = c (-Inf, -50, -10), ncores = 1)
plot (a, what = "кластеризация", data = x)
```

## Kerndenscluster

Агрегированное расстояние до эллиптической унимодальной плотности по кластерам

## Описание

Это вызывает [kerndensp](#) для вычисления и агрегирования распределения на основе [плотности](#) и главных компонент. соотношения между многомерными данными и одномерным эллиптическим распределением относительно среднего значения данных для всех кластеры в кластеризации на основе смеси, генерируемые [otrimle](#) или [rimle](#). Для использования в [отримлеге](#).

## Применение

```
kerndenscluster (x, подгонка, maxq = qnorm (0,9995), kernn = 100)
```

## Аргументы

|                 |   |
|-----------------|---|
| Икс             | что-то, что можно принудить к матрице. Набор данных.  |
| соответствовать | выходной объект <a href="#">отримле</a> или <a href="#">обода</a> .   |
| maxq            | положительное числовое. Одномерные плотности оцениваются между средним (x) -maxq и mean (x) + maxq.   |
| Kernn           | целое число. Количество точек, в которых оценивается одномерная плотность, входной параметр <a href="#">плотности</a> n. Это должно быть равно. |

## Подробности

См. Hennig and Coretto (2021), Sec. 4.2. kerndenscluster вызывает [kerndensp](#) для всех кластеров и

агрегирует полученные меры в виде корня из суммы квадратов.

#### Значение

Список с компонентами `ddpi`, `ddpm`, `measure`.

`ddpi` список выходов [kerndensp](#) для всех кластеров.

`ddpm` вектор мер-компонентов `ddpi`.

мера Окончательный результат агрегирования.

#### Авторы)

Кристиан Хенниг <[christian.hennig@unibo.it](mailto:christian.hennig@unibo.it)> [https://www.unibo.it/sitoweb/christian\\_hennig/en/](https://www.unibo.it/sitoweb/christian_hennig/en/)

#### Рекомендации

Хенниг, К. и П. Коретто (2021 г.). Подход адекватности для определения количества кластеров для OTRIMLE кластеризация на основе робастной смеси Гаусса. Появиться в Австралии и Новой Зеландии Статистический журнал, <https://arxiv.org/abs/2009.00921>.

#### Kerndensmeasure

7

#### Смотрите также

[kerndensp](#), [kerndensmeasure](#), [otrimle](#), [ободок](#)

#### Примеры

```
данные (банкнота)
selectdata <- c (1: 30,101: 110,117: 136,160: 161)
set.seed (555566)
x <- банкнота [selectdata, 5: 7]
ox <- otrimle (x, G = 2, ncores = 1)
kerndenscluster (x, ox) $ мера
```

#### Kerndensmeasure

Статистическое измерение близости к симметричному унимодальному распределению

#### Описание

Расстояние на основе плотности между одномерными данными и унимодальным симметричным распределением около среднего значения данных основано на Pons (2013, стр.79), адаптировано Хеннигом и Коретто (2021), см. подробности.

#### Применение

```
kerndensmeasure (x, weights = rep (1, nrow (as.matrix (x))), maxq = qnorm (0,9995),
kermn = 100)
```

#### Аргументы

`Икс` вектор. Одномерный набор данных.

`веса` неотрицательный вектор. Относительные веса наблюдений (будут стандартизованы до поддерживать до одного внутри).

`maxq` плотности оцениваются между средним (x) -maxq и средним (x) + maxq.

`Kernn` целое число. Количество точек, в которых оценивается плотность, входной параметр n

от [плотности](#). Это должно быть равно.

#### Подробности

Функция [плотности](#) используется для вычисления оценки плотности ядра на основе данных. Керна значения плотности затем упорядочиваются от наибольшего до наименьшего. Начиная от самого большого до наименьшего, формируются пары из двух значений (наибольшее и наибольшее наибольшее, третье и четвертое наибольшее и скоро). Каждая пара заменяется двумя копиями среднего двух значений. Затем с каждой стороны среднее значение каждой копии размещается от самого большого до самого маленького, и это дает симметричный плотность около среднего. Среднеквадратичная разница между этой плотностью и исходной плотностью вычисляется.

Стр. 8

8

Kerndensp

#### Значение

Список с компонентами `sr`, `srh`, `measure`.

`sr` вектор сгенерированных значений симметричной плотности от наибольшего к наименьшему (всего один копия, `sr kernn / 2` значения).

`srh` у-компонента [плотности](#)-выход.

мера среднеквадратичная разница между плотностями.

#### Авторы

Кристиан Хенниг <[christian.hennig@unibo.it](mailto:christian.hennig@unibo.it)> [https://www.unibo.it/sitoweb/christian\\_hennig/en/](https://www.unibo.it/sitoweb/christian_hennig/en/).

#### Рекомендации

Хенниг, К. и П. Коретто (2021 г.). Подход адекватности для определения количества кластеров для OTRIMLE кластеризация на основе робастной смеси Гаусса. Появиться в Австралии и Новой Зеландии Статистический журнал, <https://arxiv.org/abs/2009.00921>.

Понс, О. (2013). Статистическая проверка непараметрических гипотез: асимптотическая теория. Мировая наука tific, Сингапур.

#### Смотрите также

[Kerndensp](#)

#### Примеры

```
set.seed(124578)
x <- runif(20)
str(kerndensmeasure(x))
```

Kerndensp

Близость многомерного распределения к эллиптическому одномодальному распределению  
тион

#### Описание

Расстояние на основе плотности и главных компонент между многомерными данными и одномодальным эллиптическим распределением о средних данных, см. Hennig and Coretto (2021). Для использования в [kerndenscluster](#).

#### Применение

```
kerndensp (x, weights = rep (1, nrow (as.matrix (x))), siglist, maxq = qnorm (0,9995),
kernn = 100)
```

## Kerndensp

9

### Аргументы

|         |   |
|---------|---|
| Икс     | что-то, что можно принудить к матрице. Набор данных.  |
| веса    | неотрицательный вектор. Относительные веса наблюдений (будут стандартизированы до поддерживать до одного внутри).                               |
| сиглист | список с компонентами cov (ковариационная матрица), center (среднее) и n.obs (num-количество наблюдений).                                       |
| maxq    | положительное числовое. Одномерные плотности оцениваются между средним (x) -maxq и mean (x) + maxq.   |
| Kernn   | целое число. Количество точек, в которых оценивается одномерная плотность, входной параметр <a href="#">плотности</a> n. Это должно быть равно. |

### Подробности

См. Hennig and Coretto (2021), Sec. 4.2. [kerndensmeasure](#) выполняется на основных компонентах Икс. Полученные показатели стандартизируются с помощью [kmeanfun](#) и [ksdfun](#), а затем агрегируются как среднее значение. квадрат положительных значений, см. Hennig and Coretto (2021). PCS вычисляется [princomp](#) и всегда будет использовать siglist, а не статистику, вычисленную по x.

### Значение

Список с компонентами cml, cm, pca, stanmeasure, measure.

|             |   |
|-------------|---|
| cml         | Список выходных данных <a href="#">kerndensmeasure</a> для основных компонентов.                          |
| cm          | вектор компонент меры <a href="#">кернаденсера</a> для главного компонента ненцы.                         |
| stanmeasure | вектор компонент стандартизированной меры <a href="#">кернаденсера</a> для принципа ципальные компоненты. |
| pca         | вывод <a href="#">princomp</a> .  |
| мера        | Окончательный результат агрегирования.  |

### Авторы)

Кристиан Хенниг <christian.hennig@unibo.it> [https://www.unibo.it/sitoweb/christian\\_hennig/en/](https://www.unibo.it/sitoweb/christian_hennig/en/)

### Рекомендации

Хенниг, К. и П. Коретто (2021 г.). Подход адекватности для определения количества кластеров для OTRIMLE кластеризация на основе робастной смеси Гаусса. Появиться в Австралии и Новой Зеландии Статистический журнал, <https://arxiv.org/abs/2009.00921>.

### Смотрите также

[kerndensmeasure](#), [kerndenscluster](#)



10

Kmeanfun

## Примеры

```
set.seed(124578)
x <- cbind(runif(20), runif(20))
siglist <- список(cov = cov(x), center = colMeans(x), n.obs = 20)
kerndensp(x, siglist = siglist) $ мера
```

Kmeanfun

Среднее и стандартное отклонение статистики унимодальности

## Описание

Эти функции аппроксимируют среднее значение и стандартное отклонение вычисленной статистики унимодальности. by [kerndensmeasure](#), предполагающий стандартные гауссовские данные, зависящие от количества наблюдений. Они были выбраны на основе моделирования с 74 различными значениями n. Используется для стандартных [isation](#) в [kerndensp](#).

## Применение

```
kmeanfun(сущ)
ksdfun(сущ.)
```

## Аргументы

n                      целое число. Количество наблюдений.

## Значение

Полученное среднее (kmeanfun) или стандартное отклонение (ksdfun).

## Авторы)

Кристиан Хенниг <christian.hennig@unibo.it> [https://www.unibo.it/sitoweb/christian\\_hennig/en/](https://www.unibo.it/sitoweb/christian_hennig/en/)

## Смотрите также

[Kerndensp](#)

## Примеры

```
кмеанфун(50)
ксдфун(50)
```

## Описание

otrimle ищет G кластеры приблизительно гауссовой формы с / без шума / выбросов. В настройка метода, контролирующего уровень шума, выбирается адаптивно на основе данных.

## Применение

otrimle (данные, G, начальное = NULL, logicd = NULL, npr.max = 0,5, erc = 20, beta = 0, iter.max = 500, tol = 1e-06, ncores = NULL, monitor = TRUE)

## Аргументы

|          |   |
|----------|---|
| данные   | Числовой вектор, матрица или фрейм данных наблюдений. Строки соответствуют наблюдения и столбцы соответствуют переменным. Категориальные переменные и NA значения не допускаются.   |
| грамм    | Целое число, определяющее количество кластеров.   |
| исходный | Целочисленный вектор, определяющий начальное назначение кластера, где 0 обозначает шум / выбросы. Если NULL (по умолчанию) инициализация выполняется с помощью <a href="#">InitClust</a> .  |
| логичный | Вектор, определяющий сетку логарифмических (icd) значений, где icd обозначает неправильное совпадение. постоянная плотность. Если logicd = NULL, рассматривается сетка по умолчанию. Чистый гауссовский Подгонка модели смеси (полученная, когда $\log(\text{icd}) = -\text{Inf}$ ) включена в значение по умолчанию. путь поиска.  |
| npr.max  | Число в [0,1), определяющее максимальную долю шума / выбросов. Этот определяет ограничение доли шума. Если npr.max = 0, одно решение без вычисляется шумовая составляющая (соответствующая logicd = -Inf.   |
| эрг      | Число >= 1, указывающее максимально допустимое соотношение между внутри кластера собственные значения ковариационной матрицы. Это определяет ограничение на собственное отношение. erc = 1 обеспечивает сферические кластеры с равными ковариационными матрицами. Большой erc позволяет для больших межкластерных расхождений ковариации. Чтобы облегчить При настройке erc предлагается масштабировать столбцы данных (см. <a href="#">масштаб</a> ) в любое время единицы измерения различных переменных совершенно несовместимы. |
| бета     | Неотрицательная константа. Это коэффициент бета-штрафа, введенный в Коретто и Хенниг (2016).  |
| iter.max | Целочисленное значение, определяющее максимальное количество итераций, разрешенных в лежащий в основе ECM-алгоритм.   |
| тол      | Критерий остановки для лежащего в основе ECM-алгоритма. Итерация ECM останавливается если два последовательных неправильных значения логарифма правдоподобия находятся в пределах допуска.  |
| ncores   | целочисленное значение, определяющее количество ядер, используемых для параллельных вычислений. Когда ncores = NULL (по умолчанию), количество доступных ядер г определяется, и (г-1) из них используются (см. подробности). Если ncores = 1, параллельный бэкэнд не запускается.   |
| монитор  | логично. Если ИСТИНА, сообщения о ходе выполнения печатаются на экране.   |

## Подробности

Функция otrimle вычисляет решение OTRIMLE на основе ECM-алгоритма (ожидание алгоритм условной максимизации), предложенный в Coretto and Hennig (2017).

Критерий otrimle минимизируется по логической сетке логарифмических (icd) значений с использованием параллельных вычислений. ing на основе [foreach](#). Обратите внимание, что, в зависимости от настройки BLAS / LAPACK, увеличение ncores может не привести к желаемому сокращению времени вычислений. Последнее происходит при оптимизации линейного используются процедуры алгебры (например, OpenBLAS, Intel Math Kernel Library (MKL) и т. д.). Эти оптимальные

В миниатюрных разделяемых библиотеках уже реализована многопоточность. Следовательно, в этом случае увеличение pscores может лишь незначительно сократить время вычисления.

Иногда могут быть наборы данных, для которых функция не дает решения, основанного на де-аргументы вины. Это соответствует `code = 0` и `flag = 1` или `flag = 2` на выходе (см. Раздел «Значение»). ниже). Обычно это происходит, когда возникают некоторые (или все) из следующих обстоятельств: (i) `erg is` слишком большой; (ii) `prg.max` слишком велик; (iii) выбор начального разбиения. Что касается (i) и (ii), это невозможно дать числовые ссылки, потому что, являются ли эти числа слишком большими / маленькими, сильно зависит от размера выборки и размерности данных. Ссылки, приведенные ниже, объясняют соотношение между этими величинами.

Предлагается попробовать следующее всякий раз, когда возникает ошибка `code = 0`. Установите логический диапазон достаточно широкий (например, `logisd = seq (-500, -5, length = 50)`), выберите `erg = 1` и низкий выбор `prg.max` (например, `prg.max = 0,02`). Следите за решением с помощью графика профилирования критериев ([ссылка](#)). Ас-в соответствии с логикой изменения графика профилирования критерия и увеличения `erg` и `prg.max` до точки когда получен "четкий" минимум в графике профилирования критериев. Если эта стратегия не сработает, Предлагается поэкспериментировать с другими начальными разделами (см. начальный выше).

TBA: Кристиан может добавить сюда кое-что о бета-версии.

Объект `ri`, возвращаемый функцией `rimle` (см. `Value`), соответствует вектору параметров `ri` в базовой псевдомодели (1), определенной в Coretto and Hennig (2017). С `logisd = -Inf` функция Римла аппроксимирует MLE для модели простой гауссовой смеси с собственным соотношением ковариационная регуляризация, в этом случае первый элемент вектора `ri` устанавливается в ноль, потому что шумовая составляющая не учитывается. В общем, для выборки `iid` из моделей конечной смеси В контексте эти параметры `ri` определяют ожидаемые пропорции кластеров. Из-за доли шума ограничение в RIMLE, бывают ситуации, когда это соединение может не произойти на практике. В последнее, вероятно, произойдет, когда как `logisd`, так и `prg.max` велики. Таким образом, расчетная ожидаемая пропорции кластеров указываются в объекте `exp` пропорции на выходе `oboda`, и это вычисляется на основе неправильных апостериорных вероятностей, указанных в `tau`. Увидеть Коретто и Хеннига (2017) для более подробного обсуждения этого вопроса.

Более ранняя приближенная версия алгоритма была первоначально предложена Коретто и Хеннигом. (2016). Программное обеспечение для исходной версии алгоритма можно найти в дополнительном материале. риалы Коретто и Хеннига (2016).

#### Значение

Объект `S3` класса `otrimle`, обеспечивающий оптимальное (по критерию OTRIMLE) кластеризация. Компоненты вывода следующие:

|     |   |
|-----|---|
| код | Целочисленный индикатор сходимости. <code>code = 0</code> , если решение не найдено (см. Подробности); <code>code = 1</code> , если при оптимальном значении <code>icd</code> соответствующий EM-алгоритм не сходились в пределах <code>em.iter.max</code> ; <code>code = 2</code> сходимость полностью достигнута. |
|-----|---|

|              |   |
|--------------|---|
| флаг         | Строка символов, содержащая один или несколько флагов, связанных с итерацией EM в оптимальный <code>icd</code> . <code>flag = 1</code> , если не удалось предотвратить численное вырождение несобственных апостериорных вероятностей (значение <code>tau</code> ниже). <code>flag = 2</code> , если принудительное исполнение ограничение доли шума не удалось по числовым причинам. <code>flag = 3</code> , если шум ограничение пропорции было успешно применено хотя бы один раз. <code>flag = 4</code> , если ограничение на собственное отношение было успешно применено хотя бы один раз. |
| iter         | Количество итераций, выполненных в базовом EM-алгоритме при оптимальном <code>icd</code> .  |
| логичный     | Результирующее значение оптимального журнала ( <code>icd</code> ).  |
| илоглик      | Результирующее значение неправдоподобной вероятности.   |
| критерий     | Итоговое значение критерия OTRIMLE.   |
| Пи           | Расчетный вектор параметров <code>ri</code> базовой псевдомодели (см. De-хвосты).   |
| иметь в виду | Матрица размерности <code>ncol</code> (данные) x <code>G</code> , содержащая средние параметры каждого кластер (по столбцам).   |

|              |   |
|--------------|---|
| cov          | Массив размером $ncol(\text{данные}) \times ncol(\text{данные}) \times G$ , содержащий ковариационную матрицу каждого кластера.   |
| tau          | Матрица размерности $nrow(\text{data}) \times \{1 + G\}$ , где $\tau[i, 1 + j]$ - оценка (несобственная) апостериорная вероятность того, что $i$ -е наблюдение принадлежит $j$ -му кластеру - тер. $\tau[i, 1]$ - это оценочная (несобственная) апостериорная вероятность того, что $i$ -е наблюдение относится к шумовой составляющей. |
| smd          | Матрица размерности $nrow(\text{данные}) \times G$ , где $smd[i, j]$ - квадрат Махаланобиса данных $[i, j]$ от среднего $[, j]$ согласно $cov[, j]$ .   |
| кластер      | Вектор целых чисел, обозначающий назначения кластера для каждого наблюдения. Это 0 для наблюдения, относящиеся к шуму / выбросам.   |
| размер       | Вектор целых чисел с размерами (количеством) каждого кластера.  |
| экспропорция | Вектор предполагаемых ожидаемых пропорций кластеров (см. Подробности).  |
| оптимизация  | Data.frame с профилированием оптимизации OTRIMLE. Для каждого значения $\log(\text{icd})$ исследуемый алгоритмом data.frame хранит логику, критерий, илоглик, код, флаг (определено выше), а enrg обозначает ожидаемую долю шума.   |

#### Рекомендации

Коретто, П. и К. Хеннинг (2016). Надежная неправильная максимальная вероятность: настройка, вычисление и сравнение с другими методами робастной гауссовой кластеризации. Журнал американской статистики Ассоциация, Vol. 111 (516), стр. 1648–1659. DOI: [10.1080/01621459.2015.1100996](https://doi.org/10.1080/01621459.2015.1100996)

П. Коретто и К. Хеннинг (2017). Последовательность, отказоустойчивость и алгоритмы для надежной неправильной кластеризации максимального правдоподобия. Журнал исследований в области машинного обучения, Vol. 18 (142), С. 1-39. <https://jmlr.org/papers/v18/16-382.html>

#### Авторы)

Пьетро Коретто <[pcoretto@unisa.it](mailto:pcoretto@unisa.it)> <https://pietro-coretto.github.io>

#### Смотрите также

[plot.otrimle](#), [InitClust](#), [оболоч.](#)

#### Примеры

```
## Загрузить данные о швейцарских банкнотах
данные (банкнота)
x <- банкнота[, -1]

## Выполните отримле кластеризацию с аргументами по умолчанию
set.seed(1)
a <- ottrimle(data = x, G = 2, logicd = c(-Inf, -50, -10), ncores = 1)

## Кластеризация графиков
plot(a, data = x, what = "кластеризация")

## Постройте профилирование критерия OTRIMLE
сюзет(a, what = "критерий")

## Постройте неправильное профилирование логарифмической вероятности
сюзет(a, what = "iloglik")

## PP график кластерного эмпирического взвешенного квадрата Махаланобиса
## расстояния относительно целевого распределения pchisq(, df = ncol(data))
сюзет(a, what = "fit")
plot(a, what = "fit", cluster = 1)
```

```
## Выполните то же отримле, что и раньше, с ненулевым штрафом
set.seed(1)
b <- otrimle (data = x, G = 2, beta = 0.5, logictd = c (-Inf, -50, -10), ncores = 1)

## Кластеризация графиков
plot (b, data = x, what = "кластеризация")

## Постройте профилирование критерия OTRIMLE
слюжет (b, what = "критерий")

## Постройте неправильное профилирование логарифмической вероятности
слюжет (b, what = "iloglik")

## PP график кластерного эмпирического взвешенного квадрата Махаланобиса
## расстояния относительно целевого распределения rchisq (, df = ncol (data))
слюжет (b, what = "fit")
plot (b, what = "fit", cluster = 1)
```

отримлег

15

```
## Не работать:
## Выполните тот же пример, используя более тонкую сетку по умолчанию logictd
## значения с использованием нескольких ядер
##
a <- otrimle (data = x, G = 2)

## Проверьте критерий otrimle-vs-logictd
слюжет (a, что = критерий)

## Минимум достигается при $ logictd = -9, и исследуя $ оптимизацию, это
## видно, что интервал [-12,5, -4] ограничивает оптимальную
## решение. Мы ищем с более мелкой сеткой, расположенной около минимума
##
b <- otrimle (данные = x, G = 2, logictd = seq (-12,5, -4, length.out = 25))

## Проверьте критерий otrimle-vs-logictd
слюжет (b, что = критерий)

## Проверьте разницу между двумя кластерами
таблица (A = a $ cluster, B = b $ cluster)

## Проверить различия в оценочных параметрах
##
colSums (abs (a $ mean - b $ mean)) ## Расстояние L1 для средних векторов
применить ({a $ cov-b $ cov}, 3, попп, type = "F") ## Расстояние Фробениуса для ковариаций
с (Шум = abs (a $ prp-b $ prp), abs (a $ spr-b $ spr)) ## Абсолютная разница в пропорциях

## Конец (не запускается)
```

отримлег

OTRIMLE для диапазона количества кластеров с кластерами на основе плотности-статистика качества

Описание

Вычисляет оптимально настроенную робастную неправильную кластеризацию максимального правдоподобия (OTRIMLE), см. [отримле](#) вместе со статистикой качества кластера Q на основе плотности (Hennig and Coretto 2021) для

диапазон значений количества кластеров.

#### Применение

otrimleg (набор данных, G = 1: 6, многоядерный = ИСТИНА, ncores = detectCores (логический = ЛОЖЬ) -1, erno = 20, beta0 = 0, fixlogid = NULL, monitor = 1, dmaxq = qnorm (0,9995))

#### Аргументы

|              |  |
|--------------|--|
| набор данных | что-то, что может быть преобразовано в матрицу переменных, умноженных на наблюдения. В набор данных. |
| грамм        | вектор целых чисел (обычно начиная с 1). Количество рассматриваемых кластеров эред.                  |

## Стр.16

16

отримлег

|              |  |
|--------------|--|
| многоядерный | логично. Если TRUE, параллельные вычисления используются через функцию <a href="#">mclapply</a> из пакет параллельный; прочтите там предупреждения, если вы собираетесь использовать это; это не работает в Windows.   |
| ncores       | целое число. Количество ядер для распараллеливания.  |
| эрк          | Число больше или равно единице, указывающее максимально допустимое соотношение между собственными значениями внутрикластерной ковариационной матрицы. См. <a href="#">Отримле</a> .  |
| beta0        | <a href="#">Неотрицательная</a> константа, штрафной член за шум, который передается в <a href="#">отримле</a> как бета., см. документацию там.   |
| fixlogid     | числовое значение NULL. Значение логарифма неправильной логической постоянной плотности, см. <a href="#">оболок</a> , который запускается вместо <a href="#">отримле</a> , если он не равен NULL. NULL означает, что <a href="#">отримле</a> определяет это по данным. |
| монитор      | 0 или 1. Если 1, на экран выводятся сообщения о ходе выполнения.   |
| dmaxq        | числовой. Передано как maxq в <a href="#">kerndensmeasure</a> . Интервал, рассматриваемый для одномерная оценка плотности (-maxq, maxq).   |

#### Подробности

Для оценки количества кластеров это должно быть [вызвано otrimlesim](#). Выход [otrimleg](#) не предназначен для непосредственного использования для оценки количества кластеров, см. Hennig and Коретто (2021 г.).

#### Значение

otrimleg возвращает список, содержащий компоненты solution, iloglik, ibic, criterion, logid, noiseprob, denscrit, ddpm. Все это списки или векторы, номер компонента которых является количеством кластеров.

|           |   |
|-----------|---|
| решение   | список объектов вывода <a href="#">отримле</a> или <a href="#">обода</a> .  |
| илоглик   | вектор неправильных значений правдоподобия из <a href="#">отримле</a> .   |
| ибик      | вектор неправильных значений BIC (маленький - хорошо), вычисленный из iloglik и количество параметров. Обратите внимание, что поведение неправильной вероятности несовместимо со стандартным использованием BIC, поэтому это экспериментальный и не стоит доверять при выборе количества кластеров. |
| критерий  | вектор значений критерия OTRIMLE, см. <a href="#">отримле</a> .   |
| noiseprob | вектор расчетных пропорций шума, экспропорция [1] из <a href="#">отримле</a> .  |
| denscrit  | вектор статистики качества кластера Q на основе плотности (Hennig and Coretto 2021) как обеспечивается компонентом меры <a href="#">ядра kerndensmeasure</a> .  |
| ddpm      | список вектора кластерных мер качества кластеров на основе плотности как про-с помощью компонента <a href="#">ddpm</a> программы <a href="#">kerndensmeasure</a> .  |

Авторы)

Кристиан Хенниг <christian.hennig@unibo.it> [https://www.unibo.it/sitoweb/christian\\_hennig/en/](https://www.unibo.it/sitoweb/christian_hennig/en/)

Стр.17

отримлесимг

17

Рекомендации

Коретто, П. и К. Хенниг (2016). Надежная неправильная максимальная вероятность: настройка, вычисление и сравнение с другими методами робастной гауссовой кластеризации. Журнал американской статистики Ассоциация, Vol. 111 (516), стр. 1648–1659. DOI: 10.1080 / 01621459.2015.1100996

П. Коретто и К. Хенниг (2017). Последовательность, отказоустойчивость и алгоритмы для надежной неправильной кластеризации максимального правдоподобия. Журнал исследований в области машинного обучения, Vol. 18 (142), С. 1-39. <https://jmlr.org/papers/v18/16-382.html>

Хенниг, К. и П. Коретто (2021 г.). Подход адекватности для определения количества кластеров для OTRIMLE кластеризация на основе робастной смеси Гаусса. Появиться в Австралии и Новой Зеландии Статистический журнал, <https://arxiv.org/abs/2009.00921>.

Смотрите также

[otrimle](#), [rimle](#), [otrimlesimg](#), [kerndensmeasure](#)

Примеры

```
данные (банкнота)
selectdata <- c (1: 30,101: 110,117: 136,160: 161)
x <- банкнота [selectdata, 5: 7]
obanknote <- otrimleg (x, G = 1: 2, multicore = FALSE)
```

отримлесимг

Подход адекватности количества кластеров для OTRIMLE

Описание

otrimlesimg вычисляет оптимально настроенную робастную неправильную кластеризацию максимального правдоподобия (OTRIMLE), см. [otrimle](#) для диапазона значений количества кластеров, а также для искусственных наборов данных simu- рассчитано по параметрам модели, оцененным на исходных данных. Резюме-методы присутствуют и оцените результаты так, чтобы наименьшее адекватное количество кластеров можно было найти как наименьшее. для которых значение статистики качества кластера Q на основе плотности исходных данных совместимо с его распределением на искусственных наборах данных с тем же количеством кластеров, см. Hennig и Coretto 2021 для подробностей.

Применение

```
otrimlesimg (набор данных, G = 1: 6, многоядерный = ИСТИНА,
ncores = detectCores (логический = FALSE) -1, erc = 20, beta0 = 0, simruns = 20,
sim.est.logicd = ЛОЖЬ,
монитор = 1)
```

```
## Метод S3 для класса otrimlesimgdens
сводка (объект, noisepenalty = 0,05, sdcutoff = 2
, ...)
```

```
## Метод S3 для класса summary.otrimlesimgdens
печать (x, ...)

## Метод S3 для класса summary.otrimlesimgdens
plot (x, plot = "критерий", penx = NULL,
peny = NULL, pencex = 1, cutoff = TRUE, ylim = NULL, ...)
```

## Аргументы

|                |   |
|----------------|---|
| набор данных   | что-то, что может быть преобразовано в матрицу переменных, умноженных на наблюдения. В набор данных.  |
| грамм          | вектор целых чисел (обычно начиная с 1). Количество рассматриваемых кластеров эред.   |
| многоядерный   | логично. Если TRUE, параллельные вычисления используются через функцию <a href="#">mclapply</a> из пакет параллельный; прочтите там предупреждения, если вы собираетесь использовать это; это не сработает в Windows.   |
| ncores         | целое число. Количество ядер для распараллеливания.   |
| эрк            | Число больше или равно единице, указывающее максимально допустимое соотношение между собственными значениями внутрикластерной ковариационной матрицы. См. <a href="#">Отримле</a> .   |
| beta0          | <a href="#">Неотрицательная</a> константа, штрафной член за шум, который передается в <a href="#">otrimle</a> как бета., см. документацию там.  |
| симруны        | целое число. Количество реплик искусственных наборов данных, взятых из каждой модели.   |
| sim.est.logicd | логический. Если ИСТИНА, логарифм неправильной логики постоянной плотности, см. <a href="#">отримле</a> , переоценивается при запуске <a href="#">otrimle</a> на искусственных наборах данных. Oth-<br>В противном случае значение, оцененное на исходных данных, принимается фиксированным. ИСТИНА требует гораздо больше времени вычислений, но можно рассматривать как генерирующие более реалистичные вариация. |
| монитор        | 0 или 1. Если 1, на экран выводятся сообщения о ходе выполнения.  |
| шум            | число от 0 до 1. p_0 в Хенниге и Коретто (2021 г.); обычно маленький.<br>Метод предпочитает рассматривать долю <= noisepenalty точек как выбросы.<br>для добавления кластера.   |
| sd cutoff      | числовой. c в формуле (7) Хеннига и Коретто (2021 г.). Кластеризация лечится как адекватный, если его значение измерения качества кластера на основе плотности Q откалибровано (т. е. среднее / стандартное отклонение) по значениям в искусственных наборах данных, сгенерированных из расчетной модели составляет <= sd cutoff.   |
| участок        | «критерий» или «шум», см. подробности.  |
| ручка          | FALSE, NULL или числовой. x-координата, откуда простота упорядочивания дана кластеризация (как тест на графике). Если ЛОЖЬ, это не добавляется к графику.<br>Если NULL, по умолчанию делается предположение о хорошей позиции (что не всегда работает. хорошо).   |
| копейка        | NULL или числовой. x-координата, откуда простота упорядочивания кластеров-<br>Теринг дан (как тест на графике). Если NULL, по умолчанию делается предположение для хорошего положение (что не всегда работает).   |



|          |   |
|----------|---|
| пенекс   | числовой. Коэффициент увеличения (параметр сех передается в <a href="#">легенду</a> ) для простота заказа, см. параметр repnх.  |
| отрезать | логично. Если ИСТИНА, "критерий" -пункт показывает значение отсечения, ниже которого количество кластеров адекватное, см. подробности.  |
| Илим     | вектор двух числовых значений, диапазон оси Y, который будет передан для построения <a href="#">графика</a> . Если NULL, диапазон выбирается автоматически (но может отличаться от <a href="#">графика по</a> умолчанию). |
| объект   | объект класса otrimlesimgdens, полученный при вызове <a href="#">otrimlesimg</a>  |
| Икс      | объект класса summary.otrimlesimgdens, полученный в результате вызова <a href="#">summary</a> функция над объектом класса otrimlesimgdens, полученным в результате вызова <a href="#">отримлесимг</a> .                   |
| ...      | необязательные параметры для передачи на <a href="#">график</a> .   |

#### Подробности

Метод полностью описан у Хеннига и Коретто (2021). Необходимые константы настройки для выбор оптимального количества кластеров, наименьшего процента дополнительного шума, который пользователь готов обменять на добавление еще одного кластера ( $p_0$  в документе,  $poisepenalty$  здесь) и критическое значение ( $c$  в документе,  $sd cutoff$  здесь) для соответствия стандартизированной плотности на основе качества мера  $Q$  передаются в итоговую функцию, которая требуется для выбора наилучшего (простейшего адекватное) количество кластеров.

Функция построения графика `plot.summary.otrimlesimgdens` может создавать два графика. Если сюжет = "критерий", стандартизованная мера качества кластера  $Q$  на основе плотности нанесена на график в зависимости от количества кластеров. Значения для смоделированных искусственных наборов данных являются точками, значения для исходного набора данных даны, как тип линии. Если `cutoff = "TRUE"`, критические значения (см. Выше) добавляются в виде красных крестиков; число кластеров является адекватным, если значение исходных данных ниже критического значения, т. е.  $Q$  не значительно больше, чем для искусственных наборов данных, созданных на основе подобранной модели. Используя `repnх`, упорядоченные номера кластеров от простейших до наименее простых также могут быть указаны на графике, где просто определяется как количество кластеров плюс расчетная доля шума, разделенная шумовым штрафом, см. выше. Выбранное количество кластеров является наиболее простым и адекватным, т. Е. что предпочтительны небольшое количество кластеров и низкая доля шума.

Если `plot = "noise"`, доля шума (черный) и простота (красный) наносятся на график против числа-бер кластеров.

#### Значение

`otrimlesimg` возвращает список типа «`otrimlesimgdens`», содержащий компоненты `result`, `simresult`, `simruns`.

|           |  |
|-----------|--|
| результат | выходной объект <a href="#">otrimleg</a> (список результатов по исходным данным) запускается с параметром <a href="#">ters</a> , предоставленные <a href="#">otrimlesimg</a> . |
| simresult | список симуляций длины выходных объектов <a href="#">отримлега</a> для всех имитируемых артефактов. сиальные наборы данных.  |
| симруны   | входной параметр <code>simruns</code> .  |

`summary.otrimlesimgdens` возвращает список типа «`summary.otrimlesimgdens`» с компонентами. `G`, `simeval`, `ssimruns`, `npr`, `nprdiff`, `logicd`, `denscrit`, `peng`, `penorder`, `bestG`, `sd cutoff`, `bestresult`, `cluster`. симруны

|                                |  |
|--------------------------------|--|
| грамм                          | <code>otrimlesimg</code> входной параметр $G$ (количество кластеров).  |
| Симевал                        | список с компонентами <code>denscrit</code> , <code>meandens</code> , <code>sddens</code> , <code>standens</code> , <code>errors</code> , определены ниже. |
| ssimruns                       | <code>otrimlesimg</code> входной параметр <code>simruns</code> .   |
| энергетический ядерный реактор | вектор расчетных долей шума на исходных данных для всех чисел кластеры, экспорпорции [1] от <a href="#">отримле</a> .                                      |

|           |   |
|-----------|---|
| nprdiff   | вектор для всех количеств кластеров различий между оцененными наименьшими кластерами-пропорция <code>ter</code> и доля шума на исходных данных.   |
| логичный  | вектор журналов неправильных значений постоянной плотности на исходных данных для всех количество кластеров.  |
| denscrit  | вектор по всем числам кластеров статистики качества кластера на основе плотности <code>Q</code> на исходных данных, предоставленных компонентом меры <a href="#">ядро <code>kerdensmeasure</code></a> . |
| Пэн       | вектор значений простоты (см. Подробности) по всем номерам кластеров.   |
| наказание | простота порядка количества кластеров.  |
| bestG     | наилучшее (т.е. наиболее простое адекватное) количество кластеров.  |
| sdcutoff  | входной параметр <code>sdcutoff</code> .  |
| результат | вывод <a href="#">отримле</a> для лучшего числа кластеров <code>bestG</code> .  |
| кластер   | вектор кластеризации для наилучшего числа кластеров <code>bestG</code> . 0 соответствует шуму / выбросам.   |

Компоненты `simeval` компонента вывода `summary.otrimlesimgdens`:

|                             |   |
|-----------------------------|---|
| <code>denscritmatrix</code> | максимальное количество кластеров, умноженное на <code>simruns</code> матрицу <code>denscrit</code> -векторов для всех кластеризация на смоделированных данных.   |
| извилистый                  | вектор по количеству кластеров робастной оценки среднего значения плотности над смоделированными наборами данных, вычисленными с помощью <a href="#">scaleTau2</a> .  |
| <code>sddens</code>         | вектор по количеству кластеров робастной оценки стандартного отклонения <code>denscrit</code> по смоделированным наборам данных, вычисленный с помощью <a href="#">scaleTau2</a> .  |
| <code>Standens</code>       | вектор по количеству кластеров плотности исходных данных, стандартизированных по меанденс и <code>sddens</code> .   |
| ошибки                      | вектор по количеству кластеров, количество раз, которое <a href="#">отримле</a> привело к ошибке.<br><br><code>plot.summary.otrimlesimgdens</code> вернет результат <a href="#">номинальной()</a> перед любым-вещь была изменена функцией сюжета. |

Авторы)

Кристиан Хенниг <[christian.hennig@unibo.it](mailto:christian.hennig@unibo.it)> [https://www.unibo.it/sitoweb/christian\\_hennig/en/](https://www.unibo.it/sitoweb/christian_hennig/en/)

сюжет.

21 год

Рекомендации

Коретто, П. и К. Хенниг (2016). Надежная неправильная максимальная вероятность: настройка, вычисление и сравнение с другими методами робастной гауссовой кластеризации. Журнал американской статистики Ассоциация, Vol. 111 (516), стр. 1648–1659. DOI: [10.1080/01621459.2015.1100996](https://doi.org/10.1080/01621459.2015.1100996)

П. Коретто и К. Хенниг (2017). Последовательность, отказоустойчивость и алгоритмы для надежной неправильная кластеризация максимального правдоподобия. Журнал исследований в области машинного обучения, Vol. 18 (142), С. 1-39. <https://jmlr.org/papers/v18/16-382.html>

Хенниг, К. и П. Коретто (2021 г.). Подход адекватности для определения количества кластеров для OTRIMLE кластеризация на основе робастной смеси Гаусса. Появиться в Австралии и Новой Зеландии Статистический журнал, <https://arxiv.org/abs/2009.00921>.

Смотрите также

[otrimle](#), [rimle](#), [otrimleg](#), [kerdensmeasure](#)

Примеры

```
## otrimlesimg требует больших ресурсов компьютера, поэтому только небольшая часть данных
## используется для скорости.
данные (банкнота)
selectdata <- c (1: 30,101: 110,117: 136,160: 161)
set.seed (555566)
x <- банкнота [selectdata, 5: 7]

## simruns = 2 выбрано для скорости. На практике это не рекомендуется.
obanknote <- otrimlesimg (x, G = 1: 2, multicore = FALSE, simruns = 2, monitor = 0)
sobanknote <- резюме (obanknote)
печать (банкнота)
сюжет (банкнота, сюжет = "критерий", penx = 1,4)
сюжет (sobanknote, plot = "noise", penx = 1.4)
plot (x, col = sobanknote $ cluster + 1, pch = c ("N", "1", "2") [sobanknote $ cluster + 1])
```

сюжет. Методы построения для объектов OTRIMLE

Описание

Постройте надежные результаты кластеризации на основе модели: диаграмма рассеяния с информацией о кластеризации, оптимизация профилирование и кластерная подгонка.

Применение

```
## S3 метод для класса otrimle
plot (x, what = c («критерий», «илоглик», «соответствие», «кластеризация»),
      данные = NULL, поля = NULL, кластер = NULL, ...)
```

22 сюжет.

Аргументы

|         |   |
|---------|---|
| Икс     | Вывод из <a href="#">отримле</a>  |
| какие   | Тип графика. Это может быть одно из следующих значений: «критерий» (по умолчанию), «илоглик», «подгонка», «кластеризация». Смотрите подробности.  |
| данные  | Вектор данных, матрица или data.frame (или какое-то их преобразование), используемый для получение отримле объекта. Это актуально, только если what = "clustering".   |
| поля    | Вектор целых чисел, обозначающий переменные (количество столбцов данных), которые должны быть используется для парного графика, если what = "кластеризация". Когда поля = NULL, устанавливается значение 1: ncol (данные) (по умолчанию). |
| кластер | Целое число, обозначающее кластер, для которого возвращается подходящий график. Это только актуально, если что = "подходит".  |
| ...     | дальнейшие аргументы, переданные другим методам или от них.   |

Значение

Если что = "критерий" График с профилированием критерия оптимизации OTRIMLE. Критерий at log (icd) = - всегда отображается Inf.

If what = "iloglik" График с профилированием неправильной функции логарифма правдоподобия вдоль путь поиска для оптимизации OTRIMLE.

If what = "fit" График PP (график вероятности-вероятности) взвешенного эмпирического распределения функций расстояний Махаланобиса наблюдений от центров скоплений относительно цели получить распространение. Целевое распределение - это распределение хи-квадрат со степенями свободы

dom равно вес (данные). Веса задаются неправильными апостериорными вероятностями. Если cluster = NULL PP-графики создаются для всех кластеров, в противном случае кластер выбирает один PP сюжет временами.

If what = "clustering" Парная диаграмма рассеяния с членством в кластерах. Баллы, присвоенные Компонент шум / выбросы обозначен знаком +.

#### Рекомендации

Коретто, П. и К. Хенниг (2016). Надежная неправильная максимальная вероятность: настройка, вычисление и сравнение с другими методами робастной гауссовой кластеризации. Журнал американской статистики Ассоциация, Vol. 111 (516), стр. 1648–1659. DOI: [10.1080 / 01621459.2015.1100996](https://doi.org/10.1080/01621459.2015.1100996)

П. Коретто и К. Хенниг (2017). Последовательность, отказоустойчивость и алгоритмы для надежной неправильная кластеризация максимального правдоподобия. Журнал исследований в области машинного обучения, Vol. 18 (142), С. 1-39. <https://jmlr.org/papers/v18/16-382.html>

#### Авторы)

Пьетро Коретто <pcoretto@unisa.it> <https://pietro-coretto.github.io>

#### Смотрите также

[сюжет.](#)

сюжет.

23

#### Примеры

```
## Загрузить данные о швейцарских банкнотах
данные (банкнота)
x <- банкнота [, - 1]

## Выполните отрицательную кластеризацию на небольшой сетке логических значений
a <- otrimle (data = x, G = 2, loglcl = c (-Inf, -50, -10), ncores = 1)
печать (a)

## Кластеризация графиков
plot (a, data = x, what = "кластеризация")

## Постройте кластеризацию по выбранным полям
plot (a, data = x, what = "clustering", margins = 4: 6)

## Постройте кластеризацию по первым двум основным компонентам
z <- масштаб (x)% *% собственные (кор (x), симметричный = ИСТИНА) $ векторов
colnames (z) <- paste ("ПК", 1: ncol (z), sep = "")
plot (a, data = z, what = "кластеризация", поля = 1: 2)

## Постройте профилирование критерия OTRIMLE
сюжет (a, what = "критерий")

## Постройте неправильное профилирование логарифмической вероятности
сюжет (a, what = "iloglik")

## Подобрать график для всех кластеров
сюжет (a, what = "fit")

## Подобрать график для кластера 1
plot (a, what = "fit", cluster = 1)

## Не работать:
## Выполните тот же пример, используя более тонкую сетку по умолчанию loglcl
## значения с использованием нескольких ядер
```

```
## a<- otrimle (data = x, G = 2)

## Проверьте критерий otrimle-vs-logicd
сюжет (a, что = критерий)

## Минимум достигается при $ logicd = -9, и исследуя $ оптимизацию, это
## видно, что интервал [-12,5, -4] ограничивает оптимальную
## решение. Мы ищем с более мелкой сеткой, расположенной около минимума
##
b <- otrimle (данные = x, G = 2, logicd = seq (-12,5, -4, length.out = 25))

## Проверьте критерий otrimle-vs-logicd
сюжет (b, что = критерий)
```

24

plot.rimle

```
## Проверьте разницу между двумя кластерами
таблица (A = a $ cluster, B = b $ cluster)

## Проверить различия в оценочных параметрах
##
colSums (abs (a $ mean - b $ mean))          ## Расстояние L1 для средних векторов
применить ({a $ cov-b $ cov}, 3, norm, type = "F")      ## Расстояние Фробениуса для ковариаций
с (Шум = abs (a $ npr-b $ npr), abs (a $ cpr-b $ cpr)) ## Абсолютная разница в пропорциях

## Конец (не запускается)
```

plot.rimle

Методы сюжета для объектов RIMLE

## Описание

Постройте надежные результаты кластеризации на основе модели: диаграмму разброса с информацией о кластеризации и подбором кластеров.

## Применение

```
## Метод S3 для обода класса
plot (x, what = c ("соответствие", "кластеризация"),
      данные = NULL, поля = NULL, кластер = NULL, ...)
```

## Аргументы

|         |   |
|---------|---|
| Икс     | Выход из <a href="#">обода</a>  |
| какие   | Тип графика. Это может быть одно из следующих значений: «подходит» (по умолчанию), «кластеризация». Смотрите подробности.   |
| данные  | Вектор данных, матрица или data.frame (или какое-то их преобразование), используемый для получение обода объекта. Это актуально, только если what = "clustering".   |
| поля    | Вектор целых чисел, обозначающий переменные (количество столбцов данных), которые должны быть используется для парного графика, если what = "кластеризация". Когда поля = NULL, устанавливается значение 1: ncol (данные) (по умолчанию). |
| кластер | Целое число, обозначающее кластер, для которого возвращается подходящий график. Это только актуально, если что = "подходит".  |
| ...     | дальнейшие аргументы, переданные другим методам или от них.   |

## Значение

If what = "fit" График PP (график вероятности-вероятности) взвешенного эмпирического распределения функция расстояний Махаланобиса наблюдений от центров скоплений относительно цели получить распространение. Целевое распределение - это распределение хи-квадрат со степенями свободы dom равно ncol (данные). Веса задаются неправильными апостериорными вероятностями. Если

cluster = NULL PP-графики создаются для всех кластеров, в противном случае кластер выбирает один PP сюжет временами.

If what = "clustering" Парная диаграмма рассеяния с членством в кластерах. Баллы, присвоенные Компонент шум / выбросы обозначен знаком +.

ободок

25

#### Рекомендации

Коретто, П. и К. Хенниг (2016). Надежная неправильная максимальная вероятность: настройка, вычисление и сравнение с другими методами робастной гауссовой кластеризации. Журнал американской статистики Ассоциация, Vol. 111 (516), стр. 1648–1659. DOI: [10.1080 / 01621459.2015.1100996](https://doi.org/10.1080/01621459.2015.1100996)

П. Коретто и К. Хенниг (2017). Последовательность, отказоустойчивость и алгоритмы для надежной неправильная кластеризация максимального правдоподобия. Журнал исследований в области машинного обучения, Vol. 18 (142), С. 1-39. <https://jmlr.org/papers/v18/16-382.html>

#### Авторы)

Пьетро Коретто <pcoretto@unisa.it> <https://pietro-coretto.github.io>

#### Смотрите также

[отримле](#)

#### Примеры

```
## Загрузить данные о швейцарских банкнотах
данные (банкнота)
x <- банкнота [, - 1]

## Выполните кластеризацию римле с аргументами по умолчанию
set.seed (1)
a <- rimle (данные = x, G = 2)
печать (a)

## Кластеризация графиков
plot (a, data = x, what = "кластеризация")

## Постройте кластеризацию по выбранным полям
plot (a, data = x, what = "clustering", margins = 4: 6)

## Постройте кластеризацию по первым двум основным компонентам
z <- масштаб (x)% **% собственные (кор (x), симметричный = ИСТИНА) $ векторов
colnames (z) <- paste ("ПК", 1: ncol (z), sep = "")
plot (a, data = z, what = "кластеризация", поля = 1: 2)

## Подобрать график для всех кластеров
сюжет (a, what = "fit")

## Подобрать график для кластера 1
plot (a, what = "fit", cluster = 1)
```

ободок

Устойчивая неправильная кластеризация максимального правдоподобия

26 год

ободок

## Описание

`rimle` ищет  $G$  кластеры приблизительно гауссовой формы с / без шума / выбросов. В настройка метода, контролирующая уровень шума, является фиксированной и предоставляется пользователем или будет угадываться функцией довольно быстро и грязно ([otrimle](#) выполняет более сложные выбор на основе данных).

## Применение

`rimle` (данные,  $G$ , начальный = NULL, `loglcd` = NULL, `npr.max` = 0,5, `erc` = 20, `iter.max` = 500, `tol` = 1e-6)

## Аргументы

|                       |  |
|-----------------------|--|
| данные                | Числовой вектор, матрица или фрейм данных наблюдений. Строки соответствуют наблюдения и столбцы соответствуют переменным. Категориальные переменные и NA значения не допускаются.  |
| грамм                 | Целое число, определяющее количество кластеров.  |
| исходный              | Целочисленный вектор, определяющий начальное назначение кластера, где 0 обозначает шум / выбросы. Если NULL (по умолчанию) инициализация выполняется с помощью <a href="#">InitClust</a> .   |
| логичный              | Число <code>log (icd)</code> , где $0 \leq \text{icd} < \text{Inf}$ - значение неправильной константы. плотность ( <code>icd</code> ). Это настройка RIMLE для управления размером шума. Если <code>loglcd</code> = NULL (по умолчанию), значение <code>icd</code> определяется на основе данных. Чистый Подгонка модели гауссовой смеси получается с помощью <code>loglcd</code> = -Inf.  |
| <code>npr.max</code>  | Число в $[0,1)$ , определяющее максимальную долю шума / выбросов. Этот определяет ограничение доли шума. Если <code>npr.max</code> = 0 решение без шума компонент вычисляется (соответствует <code>logiked</code> = -Inf.  |
| эрк                   | Число $> 1$ , указывающее максимально допустимое соотношение между внутри кластера собственные значения ковариационной матрицы. Это определяет ограничение на собственное отношение. <code>erc</code> = 1 обеспечивает сферические кластеры с равными ковариационными матрицами. Большой <code>erc</code> позволяет для больших межкластерных расхождений ковариации. Чтобы облегчить При настройке <code>erc</code> предлагается масштабировать столбцы данных (см. <a href="#">масштаб</a> ) в любое время единицы измерения различных переменных совершенно несовместимы. |
| <code>iter.max</code> | Целочисленное значение, определяющее максимальное количество итераций, разрешенных в ECM-алгоритм (см. Подробности).   |
| тол                   | Критерий остановки для лежащего в основе ECM-алгоритма. Итерация ECM останавливается если два последовательных неправильных значения логарифма правдоподобия находятся в пределах допуска.   |

## Подробности

Функция `rimle` вычисляет решение RIMLE с использованием алгоритма ECM, предложенного в Coretto и Хенниг (2017).

Могут быть наборы данных, для которых функция не предоставляет решения на основе аргументов по умолчанию. `ments`. Это соответствует `code` = 0 и `flag` = 1 или `flag` = 2 на выходе (см. Раздел «Значение» ниже).

Обычно это происходит, когда возникают некоторые (или все) из следующих обстоятельств: (i) `log (icd)` слишком большой; (ii) `erc` слишком велик; (iii) `npr.max` слишком велик; (iv) выбор начального разбиения. В этих случаях предлагается найти подходящий интервал значений `icd` с помощью функции [otrimle](#). Детали

Раздел [otrimle](#) предлагает несколько действий, которые необходимо предпринять, когда возникает [ошибка](#) `code` = 0.

Объект  $\rho_i$ , возвращаемый функцией `gimle` (см. Value), соответствует вектору параметров  $\rho_i$  в базовой псевдомодели (1), определенной в Coretto and Hennig (2017). С `logicd = -Inf` функция Римла аппроксимирует MLE для модели простой гауссовой смеси с собственным соотношением ковариационная регуляризация, в этом случае первый элемент вектора  $\rho_i$  устанавливается в ноль, потому что шумовая составляющая не учитывается. В общем, для выборки iid из моделей конечной смеси. В контексте эти параметры  $\rho_i$  определяют ожидаемые пропорции кластеров. Из-за доли шума ограничение в RIMLE, бывают ситуации, когда это соединение может не произойти на практике. В последнее, вероятно, произойдет, когда как `logicd`, так и `prg.max` велики. Таким образом, расчетная ожидаемая пропорции кластеров указываются в объекте экспорпорции на выходе обода, и это вычисляется на основе неправильных апостериорных вероятностей, указанных в `tau`. Увидеть Коретто и Хеннига (2017) для более подробного обсуждения этого вопроса.

Более ранняя приближенная версия алгоритма была первоначально предложена Коретто и Хеннигом. (2016). Программное обеспечение для исходной версии алгоритма можно найти в дополнительном материале. риалы Коретто и Хеннига (2016).

#### Значение

Объект S3 класса `gimle`. Компоненты вывода следующие:

|                   |  |
|-------------------|--|
| код               | Целочисленный индикатор сходимости. <code>code = 0</code> , если решение не найдено (см. Подробности); <code>code = 1</code> , если EM-алгоритм не сходится в пределах <code>em.iter.max</code> ; <code>code = 2</code> сходимость полностью достигнута.   |
| флаг              | Строка символов, содержащая один или несколько флагов, связанных с итерацией EM в оптимальный <code>icd</code> . <code>flag = 1</code> , если не удалось предотвратить численное вырождение несобственных апостериорных вероятностей (значение <code>tau</code> ниже). <code>flag = 2</code> , если принудительное исполнение ограничение доли шума не удалось по числовым причинам. <code>flag = 3</code> , если принудительно Ограничение на собственное отношение не удалось по числовым причинам. <code>flag = 4</code> , если ограничение доли шума было успешно применено хотя бы один раз. <code>flag = 5</code> если ограничение на собственное отношение было успешно применено хотя бы один раз. |
| <code>iter</code> | Количество итераций, выполненных в базовом EM-алгоритме.   |
| логичный          | Значение журнала ( <code>icd</code> ).   |
| илогик            | Значение неправдоподобной вероятности.   |
| критерий          | Значение критерия OTRIMLE.   |
| Пи                | Расчетный вектор параметров $\rho_i$ базовой псевдомодели (см. De-хвосты).   |
| иметь в виду      | Матрица размерности <code>ncol</code> (данные) x <code>G</code> , содержащая средние параметры каждого кластер (по столбцам).  |
| <code>cov</code>  | Массив размером <code>ncol</code> (данные) x <code>ncol</code> (данные) x <code>G</code> , содержащий ковариационную матрицу каждого кластера.   |
| <code>tau</code>  | Матрица размерности <code>ngrow</code> (данные) x <code>{1 + G}</code> , где <code>tau[i, 1 + j]</code> - оценка (несобственная) апостериорная вероятность того, что $i$ -е наблюдение принадлежит $j$ -му кластеру - тер. <code>tau[i, 1]</code> - это оценочная (несобственная) апостериорная вероятность того, что $i$ -е наблюдение ция относится к шумовой составляющей.  |
| <code>smd</code>  | Матрица размерности <code>ngrow</code> (данные) x <code>G</code> , где <code>smd[i, j]</code> - квадрат Маха-расстояние <code>lanobis</code> данных <code>[i, ]</code> от среднего <code>[, j]</code> согласно <code>cov[, j]</code> .   |

|              |   |
|--------------|---|
| кластер      | Вектор целых чисел, обозначающий назначения кластера для каждого наблюдения. Это 0 для наблюдения, относящиеся к шуму / выбросам. |
| размер       | Вектор целых чисел с размерами (количеством) каждого кластера.  |
| экспорпорция | Вектор предполагаемых ожидаемых пропорций кластеров (см. Подробности).  |



Авторы)

Пьетро Коретто <pcoretto@unisa.it> <https://pietro-coretto.github.io>

#### Рекомендации

Коретто, П. и К. Хенниг (2016). Надежная неправильная максимальная вероятность: настройка, вычисление и сравнение с другими методами робастной гауссовой кластеризации. Журнал американской статистики Ассоциация, Vol. 111 (516), стр. 1648–1659. DOI: [10.1080 / 01621459.2015.1100996](https://doi.org/10.1080/01621459.2015.1100996)

П. Коретто и К. Хенниг (2017). Последовательность, отказоустойчивость и алгоритмы для надежной неправильной кластеризация максимального правдоподобия. Журнал исследований в области машинного обучения, Vol. 18 (142), С. 1-39. <https://jmlr.org/papers/v18/16-382.html>

#### Смотрите также

[plot.rimle](#), [InitClust](#), [отримле](#).

#### Примеры

```
## Загрузить данные о швейцарских банкнотах
данные (банкнота)
x <- банкнота [, - 1]

## -----
## ПРИМЕР 1:
## Выполните RIMLE с входами по умолчанию
## -----

set.seed (1)
a <- rimle (данные = x, G = 2)
печать (a)

## Кластеризация графиков
plot (a, data = x, what = "кластеризация")

## РР график кластерного эмпирического взвешенного квадрата Махаланобиса
## расстояния относительно целевого распределения rchisq (, df = ncol (data))
сюжет (a, what = "fit")
plot (a, what = "fit", cluster = 1)

## -----
## ПРИМЕР 2:
## Сравните решения для различных вариантов логики.
## -----

set.seed (1)
```

```
## Случай 1: бесшумное решение, которое соответствует чистой модели гауссовой смеси
b1 <- ободок (данные = x, G = 2, logicc = -Inf)
plot (b1, data = x, what = "кластеризация")
сюжет (b1, what = "fit")

## Случай 2: низкий уровень шума
b2 <- ободок (данные = x, G = 2, logicc = -100)
plot (b2, data = x, what = "кластеризация")
сюжет (b2, what = "fit")

## Случай 3: средний уровень шума
b3 <- ободок (данные = x, G = 2, logicc = -10)
plot (b3, data = x, what = "кластеризация")
сюжет (b3, what = "fit")

## Случай 3: большой уровень шума
```

```
b3 <- ободок (данные = x, G = 2, logict = 5)
plot (b3, data = x, what = "кластеризация")
сюжет (b3, what = "fit")
```

## Индекс

- \* Кластер
  - генератор.отримле, [3](#)
  - kerndenscluster, [6](#)
  - отримле, [15](#)
  - отримлесимг, [17](#)
- \* Датаген
  - генератор.отримле, [3](#)
- \* Наборы данных
  - банкнота, [2](#)
- \* htest
  - kerndenscluster, [6](#)
  - kerndensmeasure, [7](#)
  - kerndensp, [8](#)
  - kmeanfun, [10](#)
- \* Робастный
  - отримле, [15](#)
- отримле, [3](#), [6](#), [7](#), [11](#), [15](#)- [18](#), [20](#)- [22](#), [25](#), [26](#), [28](#)
- отримле, [6](#), [15](#), [16](#), [19](#), [21](#)
- отримлесимг, [16](#), [17](#), [17](#), [19](#)
- rag, [20](#)
- участок, [19](#)
- участок.отримле, [12](#), [14](#), [21](#), [22](#)
- участок.римле, [24](#), [28](#)
- plot.summary.otrimlesimgdens
  - (отримлесимг), [17](#)
- princomp, [9](#)
- принт.отримле (отримле), [11](#)
- print.rimle (ободок), [25](#)
- print.summary.otrimlesimgdens
  - (отримлесимг), [17](#)

отримлесимг, [17](#)

банкнота, [2](#)

плотность, [6-9](#)

foreach, [12](#)

генератор.отримле, [3](#)

hc, [5](#)

hclust, [5](#)

InitClust, [4, 11, 14, 26, 28](#)

kerndenscluster, [6, 8, 9](#)

kerndensmeasure, [3, 7, 7, 9, 10, 16, 17, 20, 21](#)

kerdensp, [3, 6-8, 8, 10](#)

kmeanfun, [9, 10](#)

ksdfun, [9](#)

ksdfun (kmeanfun), [10](#)

легенда, [19](#)

mclapply, [16, 18](#)

обод, [3, 6, 7, 14, 16, 17, 21, 24, 25](#)

масштаб, [11, 26](#)

scaleTau2, [20](#)

резюме, [19](#)

summary.otrimlesimgdens (отримлесимг),  
[17](#)