

POLITECHNIKA WARSZAWSKA
WYDZIAŁ MATEMATYKI I NAUK INFORMACYJNYCH

PRACA DYPLOMOWA MAGISTERSKA
NA KIERUNKU MATEMATYKA
SPECJALNOŚĆ STATYSTYKA MATEMATYCZNA I ANALIZA DANYCH

**ESTYMACJA W MODELU COXA
METODĄ STOCHASTYCZNEGO SPADKU GRADIENTU
NA PRZYKŁADZIE DANYCH
Z THE CANCER GENOME ATLAS**

AUTOR:
MARCIN PIOTR KOSIŃSKI

PROMOTOR:
PROF. NDZW. DR HAB. INŻ PRZEMYSŁAW BIECEK

WARSZAWA, LIPIEC 2015

.....
podpis promotora

.....
podpis autora

Spis treści

Wprowadzenie	5
1. Estymacja metodą największej wiarygodności	7
1.1. Estymacja	7
1.2. Metoda największej wiarygodności	9
1.3. Asymptotyczne własności estymatora największej wiarygodności	9
1.3.1. Zgodność estymatora największej wiarygodności	10
1.3.2. Asymptotyczna normalność estymatora największej wiarygodności	11
2. Model Coxa	15
2.1. Estymacja analityczna w oparciu o metodę największej wiarygodności dla funkcji pseudo/sub-wiarygodności	15
2.2. Estymacja numeryczna w oparciu o metodę stochastycznego spadku gradientu rzędu I dla funkcji pseudo/sub-wiarygodności	15
3. Numeryczne metody estymacji	17
3.0.1. Ogólne pojęcia związane ze zbieżnością algorytmu	17
3.1. Algorytmy spadku wzdłuż gradientu	17
3.1.1. Algorytm Cauchy’ego	17
3.1.2. Algorytm Raphsona-Newtona	17
3.2. Algorytmy stochastycznego spadku wzdłuż gradientu	17
3.2.1. Metoda estymacji stochastycznego spadku gradientu I	17
3.2.2. Metoda estymacji stochastycznego spadku gradientu II	17
4. Zaimplementowany algorytm	19
5. Analiza danych genomycznych - model Coxa z estymacją metodą stochastycznego spadku gradientu	21
5.1. Opis i pobranie danych	21
5.2. Analiza	21
A. Wykorzystane narzędzia	23
B. Kody w R	25
C. Dokumentacja pakietu RTCGA	27
Literatura	29

Wprowadzenie

Rozdział 1

Estymacja metodą największej wiarogodności

*The making of maximum likelihood was one of the most important developments in 20th century statistics. It was the work of one man but it was no simple process (...).
John Aldrich o R. A. Fisher'ze, 1997 [1]*

1.1. Estymacja

Estymacja to dział wnioskowania statystycznego będący zbiorem metod pozwalających na uogólnianie wyników badania próby losowej na nieznaną postać i parametry rozkładu zmiennej losowej całej populacji oraz szacowanie błędów wynikających z tego uogólnienia [18].

W statystyce matematycznej zakłada się, że rozkład prawdopodobieństwa opisujący doświadczenie należy do rodziny $\{\mathbb{P}_\theta : \theta \in \Theta\}$, ale nie zna się parametru θ .

Definicja 1.1. *Estymatorem parametru θ nazywamy dowolną statystykę $T = T(X)$ o wartościach w zbiorze Θ .*

Interpretuje się T jako przybliżenie θ i często estymator θ oznacza symbolem $\hat{\theta}$.

Pewne estymatory mające odpowiednie własności są preferowane nad inne ze względu na większą precyzję bądź ufność oszacowania danego estymatora. Poniżej przedstawione są 2 ważne definicje związane z jakością estymatorów [13], gdy rozmiar próbki X_1, \dots, X_n jest duży. Mówi się wtedy o własnościach asymptotycznych estymatorów, które z matematycznego punktu widzenia, są twierdzeniami granicznymi, w których n dąży do nieskończoności. Dzięki tym twierdzeniom możliwe jest opisanie w przybliżeniu zachowania estymatorów dla dostatecznie dużych próbek. Niestety, teoria asymptotyczna nie dostarcza informacji o tym, jak duża powinna być próbka, żeby przybliżenie było dostatecznie dobre.

Definicja 1.2. *Estymator $\hat{g}(X_1, \dots, X_n)$ wielkości $g(\theta)$ jest **nieobciążony**, jeśli dla każdego n*

$$\mathbb{E}\hat{g}(X_1, \dots, X_n) = g(\theta).$$

Definicja 1.3. Estymator $\hat{g}(X_1, \dots, X_n)$ wielkości $g(\theta)$ jest **zgodny**, jeśli dla każdego $\theta \in \Theta$

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta(|\hat{g}(X_1, \dots, X_n) - g(\theta)| \leq \varepsilon) = 1,$$

dla każdego $\varepsilon > 0$.

Definicja 1.4. Estymator $\hat{g}(X_1, \dots, X_n)$ wielkości $g(\theta)$ jest **mocno zgodny**, jeśli

$$\mathbb{P}_\theta\left(\lim_{n \rightarrow \infty} \hat{g}(X_1, \dots, X_n) = g(\theta)\right) = 1.$$

Zgodność (mocna zgodność) znaczy tyle, że

$$\hat{g}(X_1, \dots, X_n) \rightarrow g(\theta), \quad (n \rightarrow \infty)$$

według prawdopodobieństwa (prawie na pewno). Interpretacja jest taka: estymator jest uznany za zgodny, jeśli zmierza do estymowanej wielkości przy nieograniczonym powiększaniu badanej próbki.

Jednak zgodność (nawet w mocnym sensie) nie jest specjalnie satysfakcjonującą własnością estymatora, a zaledwie minimalnym żądaniem, które powinien spełniać każdy przyzwoity estymator. Dlatego od niektórych estymatorów żąda się silniejszych właściwości, takich jak asymptotyczna normalność.

Definicja 1.5. Estymator $\hat{g}(X_1, \dots, X_n)$ wielkości $g(\theta)$ jest **asymptotycznie normalny**, jeśli dla każdego $\theta \in \Theta$ istnieje funkcja σ^2 , zwana asymptotyczną wariancją, taka że

$$\sqrt{n}(\hat{g}(X_1, \dots, X_n) - g(\theta)) \xrightarrow{D} \mathcal{N}(0, \sigma^2(\theta)), \quad (n \rightarrow \infty).$$

Oznacza to, że rozkład prawdopodobieństwa statystyki $\hat{g}(X_1, \dots, X_n)$ jest dla dużych n zbliżony do rozkładu

$$\mathcal{N}\left(g(\theta), \frac{\sigma^2(\theta)}{n}\right).$$

Inaczej mówiąc, estymator jest asymptotycznie normalny, gdy:

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta\left(\frac{\sqrt{n}}{\sigma(\theta)}(\hat{g}(X_1, \dots, X_n) - g(\theta)) \leq a\right) = \Phi(a).$$

Asymptotyczna normalność mówi, że estymator nie tylko zbiega do nieznanego parametru, ale również że zbiega wystarczająco szybko, jak $\frac{1}{\sqrt{n}}$.

Jeśli estymator jest asymptotycznie normalny, to jest zgodny, choć nie musi być *mocno zgodny*.

W dalszej części tego rozdziału zostanie wprowadzone pojęcie estymatora największej wiarygodności oraz zostaną udowodnione dla niego jego właściwości, co utwierdzi w przekonaniu, że metoda największej wiarygodności, przy odpowiednich założeniach, jest metodą konstrukcji rozsądnych estymatorów.

1.2. Metoda największej wiarygodności

Metodę największej wiarygodności wprowadził R. A. Fisher w 1922 r. [6], dla której po raz pierwszy procedurę numeryczną zaproponował już w 1912 r. [5]. O burzliwym procesie powstawania metody, o zmianach w jej uzasadnieniu, o koncepcjach, które powstały w obrębie tej metody takich jak parametr, statystyka, wiarygodność, dostateczność czy efektywność oraz o podejściach, które Fisher odrzucił tworząc podstawy pod nową teorię można przeczytać w obszernej pracy dokumentalnej [1].

Metoda ta, jako alternatywa dla metody najmniejszych kwadratów [11], [7], była rozwijana i szeroko stosowana później przez wielu statystyków i wciąż znajduje obszerne zastosowania w wielu obszarach estymacji statystycznej, np. [9], [10], [12].

Aby zdefiniować estymator oparty o metodę największej wiarygodności, należy najpierw wprowadzić pojęcie funkcji wiarygodności.

Definicja 1.6. *Funkcją wiarygodności nazywamy funkcję $L : \Theta \rightarrow \mathbb{R}$ daną wzorem*

$$L(\theta; x_1, \dots, x_n) = f(\theta; x_1, \dots, x_n),$$

którą rozważamy jako funkcję parametru θ przy ustalonych wartościach obserwacji x_1, \dots, x_n , gdzie

$$f(\theta; x_1, \dots, x_n) = \begin{cases} \mathbb{P}_\theta(X_1 = x_1, \dots, X_n = x_n), & \text{dla rozkładów dyskretnych,} \\ f_\theta(x_1, \dots, x_n), & \text{dla rozkładów absolutnie ciągłych.} \end{cases}$$

Oznacza to, że wiarygodność jest właściwie tym samym, co gęstość prawdopodobieństwa, ale rozważana jako funkcja parametru θ , przy ustalonych wartościach obserwacji $x = X(\omega)$.

Definicja 1.7. *Estymatorem największej wiarygodności parametru θ , oznaczanym $ENW(\theta)$, nazywamy wartość parametru, w której funkcja wiarygodności przyjmuje supremum*

$$L(\hat{\theta}) = \sup_{\theta \in \Theta} L(\theta).$$

Niektóre pozycje w literaturze, w definicji estymatora największej wiarygodności, supremum zastępują wartością największą [17] str 14, [19] str 1 bądź [14] str 11.

1.3. Asymptotyczne własności estymatora największej wiarygodności

W tym podrozdziale zostanie wykazane, że estymator największej wiarygodności jest

1. asymptotycznie nieobciążony, (na pewno? - ostatecznie tego nie wykazuję a może warto?)
2. zgodny,
3. asymptotycznie normalny.

Dowody w tym rozdziale są znane w literaturze i opierają się o [15] i [19].

1.3.1. Zgodność estymatora największej wiarygodności

Chcąc wykazać zgodność estymatora największej wiarygodności (**przy pewnych warunkach regularności?**) przydatna będzie poniższa definicja i następujący Lemat.

Definicja 1.8. *Funkcja log-wiarygodności to funkcja spełniająca równanie*

$$l(\theta) = \log(L(\theta)).$$

Lemat 1.1. *Gdy θ_0 to prawdziwe maksimum funkcji wiarygodności, to dla każdego $\theta \in \Theta$*

$$\mathbb{E}_{\theta_0} l(\theta) \leq \mathbb{E}_{\theta_0} l(\theta_0).$$

Dowód. Rozważając różnicę:

$$\begin{aligned} \mathbb{E}_{\theta_0} l(\theta) - \mathbb{E}_{\theta_0} l(\theta_0) &= \mathbb{E}_{\theta_0} (l(\theta) - l(\theta_0)) = \mathbb{E}_{\theta_0} (\log f(\theta; X) - \log f(\theta_0; X)) \\ &= \mathbb{E}_{\theta_0} \log \frac{f(\theta; X)}{f(\theta_0; X)}, \end{aligned}$$

i pamiętając o tym, że $\log t \leq t - 1$, można dojść do

$$\begin{aligned} \mathbb{E}_{\theta_0} \log \frac{f(\theta; X)}{f(\theta_0; X)} &\leq \mathbb{E}_{\theta_0} \left(\frac{f(\theta; X)}{f(\theta_0; X)} - 1 \right) = \int \left(\frac{f(\theta; x)}{f(\theta_0; x)} - 1 \right) f(\theta_0; x) dx \\ &= \int f(\theta; x) dx - \int f(\theta_0; x) dx = 1 - 1 = 0. \end{aligned}$$

Obie całki równają się 1 jako, że są całkami z funkcji gęstości, zaś równość w nierówności zachodzi tylko wtedy, gdy $\mathbb{P}_\theta = \mathbb{P}_{\theta_0}$. ■

Dzięki temu wynikowi możliwe jest udowodnienie poniższego Twierdzenia.

Twierdzenie 1.2. *Pod pewnymi warunkami regularności nałożonymi na rodzinę rozkładów prawdopodobieństwa, estymator największej wiarygodności $ENW(\theta)$ jest zgodny, tzn.*

$$ENW(\theta) \rightarrow \theta \quad \text{dla} \quad n \rightarrow \infty.$$

Dowód.

1) Z definicji w $ENW(\theta)$ przyjmowana jest wartość największa funkcji $L(\theta)$, a więc tym bardziej funkcji $l(\theta) = \log L(\theta)$ oraz funkcji $l_n(\theta) = \frac{1}{n} l(\theta) = \frac{1}{n} \log L(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(\theta; X_i)$, gdyż ekstremum jest niezmiennicze ze względu na monotoniczną transformację i liniowe przekształcenie jakim jest podzielenie przez n .

2) Z Lematu 1.1 wynika, że θ_0 maksymalizuje $\mathbb{E}_{\theta_0} l(\theta)$.

3) Z Prawa Wielkich Liczb wynika, że $l_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(\theta; X_i) \rightarrow \mathbb{E}_{\theta_0} l(\theta)$, co ostatecznie oznacza, że $ENW(\theta)$ jest zgodny. ■

1.3.2. Asymptotyczna normalność estymatora największej wiarygodności

Fisher w swojej karierze wprowadził wiele pożytecznych pojęć stosowanych do dziś. Jednym z nich jest Informacja Fishera, która zostanie wykorzystana w dowodzie asymptotycznej normalności estymatora największej wiarygodności.

Definicja 1.9. Niech X będzie zmienną losową o gęstości f_θ , zależnej od jednowymiarowego parametru $\theta \in \Theta \subset \mathbb{R}$. **Informację Fishera** zawartą w obserwacji X nazywa się funkcję

$$I(\theta) = \mathbb{E}_\theta(l'(\theta; X))^2 = \mathbb{E}_\theta\left(\frac{\partial}{\partial\theta} \log f_\theta(X)\right)^2, \quad (1.1)$$

gdzie odpowiednio

$$I(\theta) = \int \left(\frac{\partial}{\partial\theta} \log f_\theta(x)\right)^2 f_\theta(x) dx \text{ dla zmiennej ciągłej;}$$

$$I(\theta) = \sum_x \left(\frac{\partial}{\partial\theta} \log f_\theta(x)\right)^2 f_\theta(x) \text{ dla zmiennej dyskretnej.}$$

W dowodzie asymptotycznej normalności estymatora największej wiarygodności kluczowymi założeniami są poniższe warunki regularności. Rodzina gęstości musi być dostatecznie regularna aby pewne kroki rachunkowe w dalszych rozumowaniach były poprawne.

Definicja 1.10. Warunki regularności.

- (i) Informacja Fishera jest dobrze określona. Zakłada się, że Θ jest przedziałem otwartym, istnieje pochodna $\partial/\partial\theta \log f_\theta$, całka/suma we wzorze (1.1) jest bezwzględnie zbieżna i $0 < I(\theta) < \infty$.
- (ii) Wszystkie gęstości f_θ mają ten sam nośnik, to znaczy zbiór $x \in X : f_\theta(x) > 0$ nie zależy od θ .
- (iii) Można przenosić pochodną przed znak całki, czyli zamienić kolejność operacji różniczkowania $(\partial/\partial\theta)$ i całkowania $\int \dots dx$.

Dzięki wprowadzeniu takich założeń otrzymano bardzo przydatne właściwości Informacji Fishera.

Stwierdzenie 1.3. Jeśli spełnione są warunki regularności (1.10) to:

- (i) $\mathbb{E}_\theta \frac{\partial}{\partial\theta} \log f_\theta(X) = 0$,
- (ii) $I(\theta) = \text{Var}_\theta\left(\frac{\partial}{\partial\theta} \log f_\theta(X)\right)$,
- (iii) $I(\theta) = -\mathbb{E}_\theta\left(\frac{\partial^2}{\partial\theta^2} \log f_\theta(X)\right)$.

Dowód tego stwierdzenia można znaleźć w [13].

Patrząc na postać pochodnej funkcji log-wiarygodności

$$l'(\theta_0; X) = (\log f(\theta_0; X))' = \frac{f'(\theta_0; X)}{f(\theta_0; X)},$$

można wywnioskować, że nieformalnie interpretacja Informacji Fishera jest miarą tego jak szybko zmieni się funkcja gęstości jeśli delikatnie zmieni się parametr θ w okolicach θ_0 . Biorąc kwadrat i wartość oczekiwaną, innymi słowy uśredniając po X , otrzymuje się uśrednioną wersję tej miary. Jeżeli Informacja Fishera jest duża, oznacza to, że gęstość zmieni się szybko gdyby poruszyć parametr θ_0 , co innymi słowy oznacza, że gęstość z parametrem θ_0 jest ‘znacząco inna’ i ‘może zostać łatwo odróżniona’ od gęstości z parametrami nie tak bliskimi θ_0 . Oznacza to, że możliwa estymacja θ_0 oparta o takie dane jest dobra. Z drugiej strony, jeżeli Informacja Fishera jest mała, oznacza to, że gęstość dla θ_0 jest bardzo podobna do gęstości z parametrami nie tak bliskimi do θ_0 , a co za tym idzie, dużo ciężiej będzie odróżnić tę gęstość, czyli estymacja będzie słabsza.

Dzięki pojęciu Informacji Fishera i warunkom regularności, możliwe jest udowodnienie poniższego Twierdzenia.

Twierdzenie 1.4. *Pod pewnymi warunkami regularności nałożonymi na rodzinę rozkładów prawdopodobieństwa, estymator największej wiarygodności jest asymptotycznie normalny,*

$$\sqrt{n}(ENW(\theta) - \theta_0) \rightarrow \mathcal{N}\left(0, \frac{1}{I(\theta_0)}\right).$$

Z Twierdzenia widać, że im większa Informacja Fishera tym mniejsza asymptotyczna wariancja estymatora prawdziwego parametru θ_0 .

Dowód. Ponieważ $ENW(\theta)$ maksymalizuje $l_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(\theta; X)$, to $l'_n(\theta) = 0$.

Dalej, korzystając z Twierdzenia o Wartości Średniej:

$$\frac{g(a) - g(b)}{a - b} = g'(c) \text{ albo } g(a) = g(b) + g'(c)(a - b), \text{ dla } c \in [a, b],$$

gdzie $g(\theta) = l'_n(\theta)$, $a = ENW(\theta)$, $b = \theta_0$, można zapisać równość

$$0 = l'_n(ENW(\theta)) = l'_n(\theta_0) + l''_n(\theta_1)(ENW(\theta) - \theta_0), \text{ dla } \theta_1 \in [ENW(\theta), \theta_0],$$

a z niej przejść do postaci

$$\sqrt{n}(ENW(\theta) - \theta_0) = -\frac{\sqrt{n}l'_n(\theta_0)}{l''_n(\theta_1)}. \quad (1.2)$$

Z Lematu (1.1) wynika, że θ_0 maksymalizuje $\mathbb{E}_{\theta_0} l(\theta_0)$ czyli

$$\mathbb{E}_{\theta_0} l'(\theta_0) = 0, \quad (1.3)$$

a to można wstawić do licznika w równaniu (1.2)

$$\begin{aligned} \sqrt{n}l'_n(\theta_0) &= \sqrt{n}\left(\frac{1}{n} \sum_{i=1}^n l'(\theta_0) - 0\right) \\ &= \sqrt{n}\left(\frac{1}{n} \sum_{i=1}^n l'(\theta_0) - \mathbb{E}_{\theta_0} l'(\theta_0)\right) \rightarrow \mathcal{N}\left(0, \text{Var}_{\theta_0}(l'(\theta_0))\right), \end{aligned} \quad (1.4)$$

gdzie zbieżność wynika z Centralnego Twierdzenia Granicznego.

Następnie można rozważyć mianownik w równaniu (1.2). Dla wszystkich θ wynika

$$l''(\theta) = \frac{1}{n} \sum_{i=1}^n l''(\theta) \rightarrow \mathbb{E}_{\theta_0} l''(\theta)$$

z Prawa Wielkich Liczb.

Dodatkowo, ponieważ $\theta_1 \in [ENW(\theta), \theta_0]$ a $ENW(\theta)$ jest zgodny (poprzedni podrozdział), to ponieważ $ENW(\theta) \rightarrow \theta_0$, to też $\theta_1 \rightarrow \theta_0$, a wtedy

$$l''_n(\theta_1) \rightarrow \mathbb{E}_{\theta_0} l''(\theta_0) = -I(\theta_0)$$

z punktu (iii) ze Stwierdzenia (1.3).

Wtedy prawa strona równania (1.2), dzięki (1.4)

$$-\frac{\sqrt{n}l'_n(\theta_0)}{l''_n(\theta_1)} \xrightarrow{D} \mathcal{N}\left(0, \frac{\mathbb{V}ar_{\theta_0}(l'(\theta_0))}{(I(\theta_0))^2}\right).$$

Ostatecznie wariancja

$$\mathbb{V}ar_{\theta_0}(l'(\theta_0)) \stackrel{zdef}{=} \mathbb{E}_{\theta_0}(l'(\theta_0))^2 - (\mathbb{E}_{\theta_0} l'(\theta_0))^2 = I(\theta_0) - 0,$$

co wynika z definicji Informacji Fishera i (1.3).

■

Rozdział 2

Model Coxa

- 2.1. Estymacja analityczna w oparciu o metodę największej wiarygodności dla funkcji pseudo/sub-wiarygodności
- 2.2. Estymacja numeryczna w oparciu o metodę stochastycznego spadku gradientu rzędu I dla funkcji pseudo/sub-wiarygodności

Poszukujemy rozwiązań równości

$$\delta \ln L_n / \delta \theta = 0.$$

Tym razem w ogólnym przypadku zwykle nie znajdziemy analitycznego rozwiązania. W związku z tym jesteśmy zdani na metody iteracyjne. Poza tym, być może rozwiązanie problemu nie istnieje albo istnieje ich wiele. Zwykle używa się do tego celu, tj. znalezienia rozwiązania, metody Newtona, zwykle w literaturze statystycznej w zastosowaniu do tego problemu, nazywanej metodą Newtona-Raphsona. W efekcie w zasadzie dla każdego modelu z osobna należy badać własności asymptotyczne estymatora największej wiarygodności.

Rozdział 3

Numeryczne metody estymacji

3.0.1. Ogólne pojęcia związane ze zbieżnością algorytmu

Warunki stopu itp

3.1. Algorytmy spadku wzdłuż gradientu

3.1.1. Algorytm Cauchy’ego

3.1.2. Algorytm Raphsona-Newtona

3.2. Algorytmy stochastycznego spadku wzdłuż gradientu

3.2.1. Metoda estymacji stochastycznego spadku gradientu I

Algorytm SGD

3.2.2. Metoda estymacji stochastycznego spadku gradientu II

Rozdział 4

Zaimplementowany algorytm

Rozdział 5

Analiza danych genomicznych - model Coxa z estymacją metodą stochastycznego spadku gradientu

5.1. Opis i pobranie danych

5.2. Analiza

Dodatek A

Wykorzystane narzędzia

Dodatek B

Kody w R

Dodatek C

Dokumentacja pakietu RTCGA

R documentation

of all in ‘RTCGA/man’

May 5, 2015

R topics documented:

RTCGA-package	1
availableDataSets	2
availableDates	3
checkDataSetsAvailability	3
checkGenesNamesAvailability	4
downloadTCGA	5
infoTCGA	6
mergeTCGA	6
read.clinical	7
Index	8

RTCGA-package	<i>The Cancer Genome Atlas data integration</i>
---------------	---

Description

The Cancer Genome Atlas (TCGA) Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA. It contains clinical information, genomic characterization data, and high level sequence analysis of the tumor genomes. The key is to understand genomics to improve cancer care. RTCGA package offers download and integration of the variety and volume of TCGA data using patient barcode key, what enables easier data possession. This may have an beneficial influence on impact on development of science and improvement of patients’ treatment. Furthermore, RTCGA package transforms TCGA data to form which is convenient to use in R statistical package. Those data transformations can be a part of statistical analysis pipeline which can be more reproducible with RTCGA

Details

For more detailed information visit **RTCGA** wiki on [Github](#).

Author(s)

Marcin Kosinski [aut, cre] < m.p.kosinski@gmail.com >
Przemysław Biecek [aut] < przemyslaw.biecek@gmail.com >

See Also

Other RTCGA: [availableDataSets](#); [availableDates](#); [availableGenesNames](#), [checkGenesNamesAvailability](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [infoTCGA](#); [mergeTCGA](#); [read.clinical](#)

availableDataSets	<i>TCGA datasets' names</i>
-------------------	-----------------------------

Description

Enables to check TCGA datasets' names for current release date and cohort.

Usage

```
availableDataSets(cancerType, date = NULL)
```

Arguments

cancerType	A character of length 1 containing abbreviation (Cohort code) of types of cancers to check for available datasets' names on http://gdac.broadinstitute.org/ .
date	A NULL or character specifying from which date datasets' names should be checked. By default (date = NULL) the newest available date is used. All available dates can be checked on http://gdac.broadinstitute.org/runs/ or by using availableDates function. Required format "YYYY-MM-DD".

Value

A vector of available datasets' names to pass to the [downloadTCGA](#) function.

See Also

Other RTCGA: [RTCGA-package](#); [availableDates](#); [availableGenesNames](#), [checkGenesNamesAvailability](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [infoTCGA](#); [mergeTCGA](#); [read.clinical](#)

Examples

```
## Not run:
availableDataSets( "BRCA" )
availableDataSets( "OV", availableDates()[5] ) # error

## End(Not run)
```

availableDates	<i>TCGA datasets' releases dates</i>
----------------	--------------------------------------

Description

Enables to check dates of TCGA datasets' releases.

Usage

```
availableDates()
```

Value

A vector of available dates to pass to the [downloadTCGA](#) function.

See Also

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableGenesNames](#), [checkGenesNamesAvailability](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [infoTCGA](#); [mergeTCGA](#); [read.clinical](#)

Examples

```
## Not run:
availableDates()

## End(Not run)
```

checkDataSetsAvailability	<i>TCGA datasets' names availability</i>
---------------------------	--

Description

Enables to check TCGA datasets' names availability for current release date and cancer type.

Usage

```
checkDataSetsAvailability(cancerTypes, pattern = "Merge_Clinical.Level_1",
  date = NULL)
```

Arguments

cancerTypes	A character vector containing abbreviation (Cohort code) of types of cancers to check for availability of datasets' name on http://gdac.broadinstitute.org/ .
pattern	A character vector of length 1 containing a part of a dataset's name to be checked for availability for current date parameter. By default phrase "Merge_Clinical.Level_1" is checked.
date	A NULL or character specifying from which date datasets' names should be checked for availability. By default (date = NULL) the newest available date is used. All available dates can be checked on http://gdac.broadinstitute.org/runs/ or by using availableDates function. Required format "YYYY-MM-DD".

Value

A vector of available datasets names to pass to the [downloadTCGA](#) function.

See Also

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableDates](#); [availableGenesNames](#), [checkGenesNamesAvailability](#); [downloadTCGA](#); [infoTCGA](#); [mergeTCGA](#); [read.clinical](#)

Examples

```
## Not run:
checkDataSetsAvailability( "BRCA" )
checkDataSetsAvailability( c("BRCA", "OV") )
checkDataSetsAvailability( "BRCA", "Mutation_Packager_Calls.Level" )

## End(Not run)
```

checkGenesNamesAvailability

TCGA genes' names and availability in Merge_rnaseqv2__... dataset.

Description

availableGenesNames returns all available genes' names from genes' expressions dataset, where checkGenesNamesAvailability checks whether genes specified in genes are available in Merge_rnaseqv2__illumina (genes' expressions) dataset.

Usage

```
checkGenesNamesAvailability(rnaseqDir, genes)

availableGenesNames(rnaseqDir)
```

Arguments

rnaseqDir	A directory to a cancerType.rnaseqv2__illumina.rnaseqv2__unc_edu__Level_3__RSEM file
genes	A character - which genes to check for availability in a dataset.

Value

A vector containing genes' names that matched existing names.

See Also

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableDates](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [infoTCGA](#); [mergeTCGA](#); [read.clinical](#)

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableDates](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [infoTCGA](#); [mergeTCGA](#); [read.clinical](#)

Examples

```
## Not run:
  checkGenesNamesAvailability( rnaseqDir, "TP53" )

## End(Not run)
```

downloadTCGA	<i>Download TCGA data</i>
--------------	---------------------------

Description

Enables to download TCGA data from specified dates of releases of concrete Cohorts of cancer types. Pass a name of required dataset to the `dataSet` parameter. By default the Merged Clinical `dataSet` is downloaded (value `dataSet = "Merge_Clinical.Level_1"`) from the newest available date of release.

Usage

```
downloadTCGA(cancerTypes, dataSet = "Merge_Clinical.Level_1", destDir,
  date = NULL)
```

Arguments

<code>cancerTypes</code>	A character vector containing abbreviations (Cohort code) of types of cancers to download from http://gdac.broadinstitute.org/ .
<code>dataSet</code>	A part of the name of <code>dataSet</code> to be downloaded from http://gdac.broadinstitute.org/runs/ . By default the Merged Clinical <code>dataSet</code> is downloaded (value <code>dataSet = "Merge_Clinical.Level_1"</code>). Available datasets' names can be checked using availableDataSets function.
<code>destDir</code>	A character specifying a directory into which <code>dataSets</code> will be downloaded.
<code>date</code>	A NULL or character specifying from which date <code>dataSets</code> should be downloaded. By default (<code>date = NULL</code>) the newest available date is used. All available dates can be checked on http://gdac.broadinstitute.org/runs/ or by using availableDates function. Required format "YYYY-MM-DD".

See Also

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableDates](#); [availableGenesNames](#), [checkGenesNamesAvailability](#); [checkDataSetsAvailability](#); [infoTCGA](#); [mergeTCGA](#); [read.clinical](#)

Examples

```
## Not run:

dir.create( "hre")

downloadTCGA( cancerTypes = "BRCA", dataSet = "miR_gene_expression",
  destDir = "hre/" )

downloadTCGA( cancerTypes = c("BRCA", "OV"), destDir = "hre/" )

## End(Not run)
```

infoTCGA

Information about cohorts from TCGA project

Description

Function restores codes and counts for each cohort from TCGA project.

Usage

```
infoTCGA()
```

Value

A list with a tabular information from <http://gdac.broadinstitute.org/>.

See Also

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableDates](#); [availableGenesNames](#), [checkGenesNamesAvailability](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [mergeTCGA](#); [read.clinical](#)

mergeTCGA

Merge Clinical data with genes' Mutations and Expressions data

Description

mergeTCGA enables to

Usage

```
mergeTCGA(clinicalDir, rnaseqDir = NULL, mutationDir = NULL, genes)
```

Arguments

clinicalDir A directory to a cancerType.clin.merged.txt file. cancerType might be BRCA, OV etc. Can be checked using [infoTCGA](#) function.

rnaseqDir A directory to a cancerType.rnaseqv2__illuminahisec_rnaseqv2__unc_edu__Level_3__RSEM file, which is a set with gene's Expressions.

mutationDir A directory to a Mutation_Packager_Calls.Level1 folder where are genes' Mutations files.

genes For rnaseqDir - which genes' expressions to merge with clinical data in clinicalDir. For mutationDir which gene's mutations to merge with clinical data in clinicalDir.

Value

A cancerType.clin.merged.txt file is updated with newline containing informations about genes passed to genes argument.

Note

Original cancerType.clin.merged.txt file will be changed after performing merge operation.
Only one of rnaseqDir and mutationDir can be used at a time.

See Also

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableDates](#); [availableGenesNames](#),
[checkGenesNamesAvailability](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [infoTCGA](#); [read.clinical](#)

read.clinical	<i>Read from txt fo,e</i>
---------------	---------------------------

Description

TODO

Usage

```
read.clinical(clinicalDir, ...)
```

Arguments

clinicalDir	A directory to a cancerType.clin.merged.txt file. cancerType might be BRCA, OV etc.
-------------	---

Value

A data.frame with clinical data.

See Also

Other RTCGA: [RTCGA-package](#); [availableDataSets](#); [availableDates](#); [availableGenesNames](#),
[checkGenesNamesAvailability](#); [checkDataSetsAvailability](#); [downloadTCGA](#); [infoTCGA](#); [mergeTCGA](#)

Index

availableDataSets, [2](#), [2](#), [3–7](#)
availableDates, [2](#), [3](#), [3](#), [4–7](#)
availableGenesNames, [2–7](#)
availableGenesNames
 (checkGenesNamesAvailability),
 [4](#)

checkDataSetsAvailability, [2](#), [3](#), [3](#), [4–7](#)
checkGenesNamesAvailability, [2–4](#), [4](#), [5–7](#)

downloadTCGA, [2–4](#), [5](#), [6](#), [7](#)

infoTCGA, [2–6](#), [6](#), [7](#)

mergeTCGA, [2–6](#), [6](#), [7](#)

read.clinical, [2–7](#), [7](#)
RTCGA-package, [1](#)

Literatura

- [1] Aldrich J., (1997) *R. A. Fisher and the Making of Maximum Likelihood 1912 – 1922*, Statistical Science 1997, Vol. 12, No. 3, 162-176.
- [2] Biecek P., (2011) *Przewodnik po pakiecie R*, Rozprawa doktorska, Oficyna Wydawnicza GiS, wydanie II.
- [3] Bottou L., (2010) *Large-Scale Machine Learning with Stochastic Gradient Descent*.
- [4] Bottou L., (2012) *Stochastic Gradient Descent Tricks*.
- [5] Fisher R. A., (1912) *An absolute criterion for fitting frequency curves*.
- [6] Fisher R. A., (1922) *On the mathematical foundations of theoretical statistics*, Philos. Trans. Roy. Soc. London Ser. A 222 309-368.
- [7] Gauss C. F., (1809) *Theoria Motus Corporum Coelestium*.
- [8] Gągolewski M., (2014) *Programowanie w języku R*, Wydawnictwo Naukowe PWN.
- [9] Hutchinson J. B., (1928) *The Application of the "Method of Maximum Likelihood" to the Estimation of Linkage*, Genetics. 1929 Nov; 14(6): 519–537.
- [10] Kenward M. G., Lesaffre E. and Molenberghs G., (1994) *An Application of Maximum Likelihood and Generalized Estimating Equations to the Analysis of Ordinal Data from a Longitudinal Study with Cases Missing at Random*, Biometrics Vol. 50, No. 4 (Dec., 1994), pp. 945-953.
- [11] Legendre A. M., (1804) *Nouvelles methods pour la determination des orbites des comètes*.
- [12] Millar R. B., (2011) *Maximum Likelihood Estimation and Inference: With Examples in R, SAS and ADMB, chapter 6. Some Widely Used Applications of Maximum Likelihood*, John Wiley & Sons, Ltd.
- [13] Niemiro W., (2011) Skrypt do przedmiotu *Statystyka*,
<http://www-users.mat.umk.pl/~wniem/Statystyka/Statystyka.pdf>
- [14] Panchenko D., (2006), Notatki do otwartego kursu MIT *Statistics for Applications, Lecture 2: Maximum Likelihood Estimators.*,
<http://ocw.mit.edu/courses/mathematics/18-443-statistics-for-applications-fall-2006/>
- [15] Panchenko D., (2006), Notatki do otwartego kursu MIT *Statistics for Applications, Lecture 3: Properties of MLE: consistency, asymptotic normality. Fisher information.*,
<http://ocw.mit.edu/courses/mathematics/18-443-statistics-for-applications-fall-2006/>

-
- [16] R Core Team, (2013) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Wiedeń , ISBN 3-900051-07-0, <http://www.R-project.org/>.
 - [17] Rydlewski J., (2009) *Estymatory Największej Wiarogodności w Uogólnionych Modelach Regresji Nieliniowej*, Rozprawa doktorska.
 - [18] Wikipedia, encyklopedia wolnego dostępu wikipedia.pl
 - [19] Woodcock S, (2014), Notatki do otwartego kursu Uniwersytetu Simona Frasera *ECON 837, Lecture 11 Asymptotic Properties of Maximum Likelihood Estimators*, <http://www.sfu.ca/~swoodcoc/teaching/sp2014/econ837/11.mle.pdf>
 - [20] Zieliński R., (1990) *Siedem wykładów wprowadzających do statystyki matematycznej*, Warszawa, Wydawnictwo Naukowe PWN.

Marcin Piotr Kosiński
Nr albumu 265361

Warszawa, 19 maja 2015

Oświadczenie

Oświadczam, że pracę magisterską pod tytułem „Estymacja w modelu Coxa metodą stochastycznego spadku gradientu na przykładzie danych z The Cancer Genome Atlas”, której promotorem jest prof. ndzw. dr hab. inż Przemysław Biecek wykonałem samodzielnie, co poświadczam własnoręcznym podpisem.

.....

Marcin Piotr Kosiński