

**РАСПРЕДЕЛЕННЫЕ БАЙЕСОВСКИЕ
ПРОЦЕДУРЫ РАСПОЗНАВАНИЯ
ТЕКСТОВОЙ ИНФОРМАЦИИ**

В настоящее время наблюдается быстрый рост объема информации, обрабатываемой компьютерными системами. Это приводит к необходимости разработки эффективных методов распознавания текстовой информации. Одним из таких методов являются распределенные байесовские процедуры. В данной статье рассматриваются алгоритмы таких процедур, а также их реализация на платформе MapReduce. Приведены результаты экспериментальных исследований, подтверждающие эффективность предложенных методов. В частности, показано, что предлагаемые алгоритмы позволяют значительно сократить время обработки данных по сравнению с традиционными методами. Это особенно важно для задач, требующих обработки больших объемов информации в реальном времени. Кроме того, рассмотрены вопросы оптимизации параметров моделей и выбора признаков. Приведены ссылки на научные публикации, посвященные данной теме.

MapReduce.

[3]).

MapReduce [4],

MapReduce.

MapReduce. $A, |A| < \infty$ $B, |B| < \infty$ -
 $f: A^l \mapsto B$ -
 $(a_1, \dots, a_l) \in A^l$
 $f(a_1, \dots, a_l) \in B.$ $f(a_1, \dots, a_l)$
MapReduce
 $: m(\cdot) \oplus \cdot,$
 $f(a_1, \dots, a_l) = m(a_1) \oplus m(a_2) \oplus \dots \oplus m(a_l),$ (1)
 $m: A \mapsto B$ -
 $a_i \in A, i = \overline{1, l}$
 $m(a_i) \in B.$ $m(a_i)$
 $\oplus: B \times B \mapsto B,$
:
 $b_1 \oplus b_2 = b_2 \oplus b_1, \forall b_1, b_2 \in B$ (2)
 $(b_1 \oplus b_2) \oplus b_3 = b_1 \oplus (b_2 \oplus b_3), \forall b_1, b_2, b_3 \in B$ (3)
 $\exists e \in B: b \oplus e = b, \forall b \in B$ (4)
(2) - (4),
 $m(a_i)$

ISSN 2616-938 . . 2018, 1 47

... ..

 MapReduce. -
 - Apache Spark, -
 : , , , -
 : -
 [5].
 , Y - D -
 , $E = D \times Y$.
 $\tau = (e_1, \dots, e_l) \in E^l$, l E -
 $g : X \mapsto Y$, $g(x(d)) \in Y$ $d \in D$ -
 $x(d) \in X$, X - $|X| < \infty$. -
 $g(x) = \arg \max_{y \in Y} P(x|y)P(y)$, (5)
 $P(x|y)$ - $x \in X$ $y \in Y$, $P(y)$ -
 $y \in Y$.
 $d \in D$,
 $x(d) = (x_1, \dots, x_n) \in X$.
 $y \in Y$
 $P(x_1, \dots, x_n | y) = \prod_{i=1}^n P_i(x_i | y)$, (6)
 $P_i(x_i | y)$ - i - x_i
 y . $P_i(x_i | y)$ $P(y)$ (6) -
 $\hat{P}_i(x_i | y)$ $\hat{P}(y)$, (3)
 $g(x) = \arg \max_{y \in Y} [\ln \hat{P}(y) + \sum_{i=1}^n \ln \hat{P}_i(x_i | y)]$. (7)
 $P_i(x_i | y)$, $P(y)$

Spark. Apache

$|V| < \infty$. D V , Y , (\cdot, \cdot, \cdot) .

$$x(d) = (x_1(d), \dots, x_n(d)) \in X, \quad x_i(d) \in \{0, 1\}, \quad i = \overline{1, n}, \quad x_i(d) = 1, \quad i \in V, \quad d \in D, \quad x_i(d) = 0. \quad (6)$$

$$P(x_1, \dots, x_n | y) = \prod_{i=1}^n [p_{iy} x_i + (1 - p_{iy})(1 - x_i)], \quad (8)$$

$$p_{iy} = \sum_{i \in V} p_{iy}, \quad y \in Y. \quad (5)$$

$$g_{\pi, \theta}(x_1, \dots, x_n) = \arg \max_{y \in Y} \pi_y \prod_{i=1}^n [\theta_{iy} x_i + (1 - \theta_{iy})(1 - x_i)], \quad (9)$$

$$\pi \in \mathbb{R}^{|Y|}, \quad P(i), \quad \theta \in \mathbb{R}^{n \times |Y|}, \quad P_{ij}, \quad \tau = (e_1, \dots, e_l).$$

MapReduce.

$$\begin{aligned}
& \tau = (e_1, \dots, e_l) & \pi \in \mathbb{R}^{|Y|} & - \\
& \pi_i(\tau) = \frac{c_i(\tau)}{l}, & & (10) \\
& c_i(\tau) - & i \in Y & \tau \in E^l. \\
& (c_1(\tau), \dots, c_{|Y|}(\tau)) = c(e_1, \dots, e_l) & \text{MapReduce,} & \\
& & : E^l \mapsto \mathbb{Z}^{|Y|} & -
\end{aligned}$$

$$\begin{aligned}
& c(e_1, \dots, e_l) = m(e_1) \oplus m(e_2) \oplus \dots \oplus m(e_l). & (11) \\
& m : E \mapsto \{0, 1\}^{|Y|} \\
& e_k \in E & m(e_k) \in \{0, 1\}^{|Y|}, & , \\
& & \dots &
\end{aligned}$$

$$\begin{aligned}
& m_i(e_k) = u(j, y(e_k)), & (12) \\
& y(e_k) - & (&) & e_k, \quad \delta(\cdot, \cdot) -
\end{aligned}$$

$$\delta(x, y) = \begin{cases} 1, & x = y \\ 0, & x \neq y \end{cases}. \quad (13)$$

$$\begin{aligned}
& \oplus : \{0, 1\}^{|Y|} \times \{0, 1\}^{|Y|} \mapsto \mathbb{Z}^{|Y|} \\
& , & (2) - (4) & , \\
& \text{MapReduce.} \\
& \text{"} \in \mathbb{R}^{n \times |Y|} \\
& \ddagger = (e_1, \dots, e_l) \in E^l
\end{aligned}$$

$$\begin{aligned}
& \theta_{ij}(\tau) = \frac{c_{ij}(\tau)}{c_i(\tau)}, & (14) \\
& c_{ij}(\tau) - & j \in Y, \\
& i \in V \quad \tau \in E^l, \quad c_i(\tau) - & i \in Y \quad \tau \in E^l. \\
& c(\tau) \in \mathbb{Z}^{n \times |Y|} & \text{MapReduce,} & -
\end{aligned}$$

$$c(e_1, \dots, e_l) = m(e_1) \oplus m(e_2) \oplus \dots \oplus m(e_l), \quad (15)$$

$$\begin{aligned}
 & m : E \mapsto \mathbb{Z}^{n \times |Y|} && e_k, \\
 k = \overline{1, l} & && m(e_k) \in \mathbb{Z}^{n \times |Y|}, && y(e_k) - \\
 & x(d(e_k)) \in \mathbb{Z}^n && d(e_k) \in D && e_k \in D \times Y,
 \end{aligned}$$

$$m_{ij}(e_k) = x_i(d(e_k)) \cup (j, y(e_k)). \quad (16)$$

$$\begin{aligned}
 \oplus : \mathbb{Z}^{n \times |Y|} \times \mathbb{Z}^{n \times |Y|} &\mapsto \mathbb{Z}^{n \times |Y|} && - \\
 &&& (2) - (4), && -
 \end{aligned}$$

$$\begin{aligned}
 & d \in D && V && - \\
 & c && y \in Y. && d \in D && - \\
 n- & && , \quad n = |V| && x(d) \in \mathbb{Z}^{|V|}, && - \\
 x_i(d) & && i \in V && d. && - \\
 & && (6) && && -
 \end{aligned}$$

$$P(x_1, \dots, x_n | y) = K(x) \prod_{i=1}^n p_{iy}^{x_i}, \quad (17)$$

$$\begin{aligned}
 p_{iy} - & && i \in V && y \in Y (&& i \in V \\
 & && && && y \in Y), \quad K(x_1, \dots, x_n) = \frac{(\sum_{i=1}^n x_i)!}{\prod_{i=1}^n (x_i)!} - && -
 \end{aligned}$$

$$\begin{aligned}
 & && (5) && && - \\
 & && g_{\pi, \theta}(x_1, \dots, x_n) = \arg \max_{y \in Y} \pi_y \prod_{i=1}^n \theta_{iy}^{x_i}. && && (18) \\
 & && (9), && && -
 \end{aligned}$$

$$\begin{aligned}
 \pi & && \theta && \tau = (e_1, \dots, e_n). && - \\
 & && \pi && \text{MapReduce,} && - \\
 (10), & && && \theta && - \\
 (14) & && && && -
 \end{aligned}$$

$$\theta_{ij}(\tau) = \frac{c_{ij}(\tau)}{\sum_{k=1}^n c_{kj}(\tau)}. \quad (19)$$

(8), (17)

MapReduce

MapReduce.

, M.A.

MapReduce.

B.A. Biletskyy, M.A. Gupal

DISTRIBUTED BAYESIAN PROCEDURES OF THE TEXT INFORMATION RECOGNITION

In this paper, we consider Bernoulli and Multinomial variations of Bayesian Machine Learning procedures, as well as their distributed implementations based on MapReduce. We discuss their distributed implementation and use cases.

1. Hilbert M., López P. The World's Technological Capacity to Store, Communicate, and Compute Information. *Science*. 2011. Vol. 332. P. 60–65.
2. Moore G. Cramming more components onto integrated circuits. *Electronics*. 1965. Vol. 38. P. 114–117.
3. Gilbert S., Lynch N. Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services. *ACM SIGACT News*. 2002. Vol. 33. P. 51–59.
4. Dean J., Ghemawat S. MapReduce: simplified data processing on large clusters. Proceedings of the 6th conference on Symposium on Operating Systems Design & Implementation. 2004. Vol. 6. P. 10 – 17.
5. Manning C., Raghavan P., Schütze H. Introduction to information retrieval. *Cambridge University Press*. 2008. 482 p.

20.05.2018

Об авторах: