

Impact Factor:

ISRA (India) = 4.971
ISI (Dubai, UAE) = 0.829
GIF (Australia) = 0.564
JIF = 1.500

SIS (USA) = 0.912
PIHII (Russia) = 0.126
ESJI (KZ) = 8.716
SJIF (Morocco) = 5.667

ICV (Poland) = 6.630
PIF (India) = 1.940
IBI (India) = 4.260
OAJI (USA) = 0.350

SOI: [1.1/TAS](#) DOI: [10.15863/TAS](#)

International Scientific Journal Theoretical & Applied Science

p-ISSN: 2308-4944 (print) e-ISSN: 2409-0085 (online)

Year: 2020 Issue: 05 Volume: 85

Published: 30.05.2020 <http://T-Science.org>

QR – Issue



QR – Article



Olga Evgenievna Filatova

Peter the Great St. Petersburg Polytechnic University
Bachelor's Student
Institute of Computer Science and Technology

Oleg Yurievich Sabinin

Peter the Great St. Petersburg Polytechnic University
Candidate of Engineering Sciences, Docent
Institute of Computer Science and Technology

CLASSIFICATION AUTOMATION OF TEXT DOCUMENTS USING ORACLE MEANS

Abstract: The purpose of this article is to review methods for automating the classification of text documents and describe the features of solving this problem using Oracle tools (in particular, Oracle Data Mining and Oracle Text).

Key words: text mining, documents classification, oracle, oracle text.

Language: Russian

Citation: Filatova, O. E., & Sabinin, O. Y. (2020). Classification automation of text documents using Oracle means. *ISJ Theoretical & Applied Science*, 05 (85), 456-461.

Soi: <http://s-o-i.org/1.1/TAS-05-85-85> **Doi:**  <https://dx.doi.org/10.15863/TAS.2020.05.85.85>

Scopus ASCC: 1700.

АВТОМАТИЗАЦИЯ КЛАССИФИКАЦИИ ТЕКСТОВЫХ ДОКУМЕНТОВ С ПОМОЩЬЮ СРЕДСТВ ORACLE

Аннотация: Целью данной статьи является обзор методов автоматизации классификации текстовых документов и описание особенностей решения этой задачи с использованием средств Oracle (в частности Oracle Data Mining и Oracle Text).

Ключевые слова: анализ текста, классификация документов, oracle, oracle text.

Введение

С каждым годом растет спрос на цифровой документооборот. Уже сейчас проводятся попытки полностью уйти от бумажного документооборота. Так, например, 24 апреля 2020 года опубликован федеральный закон [1], согласно которому в период до 31 марта 2021 проводится эксперимент по ведению отдельными работодателями электронных документов, касающихся трудовых отношений с работниками, без дублирования на бумажном носителе. Поэтому проблемы, связанные с обработкой текстовых документов, в наше время очень актуальны.

Работа с большим количеством документов в цифровой форме приводит к тому, что необходимо организовать их хранение таким образом, чтобы облегчить поиск нужных документов в дальнейшем, поэтому одной из основных задач при обработке текстовой информации является классификация.

Задача классификации текстов возникла еще в начале 60-х годов, но только в начале 90-х для ее решения стали применяться информационные технологии благодаря возросшему прикладному интересу и доступности более мощного оборудования.

В настоящее время классификация текстов применяется для решения многих задач [2]:

Impact Factor:

ISRA (India) = 4.971
ISI (Dubai, UAE) = 0.829
GIF (Australia) = 0.564
JIF = 1.500

SIS (USA) = 0.912
РИИЦ (Russia) = 0.126
ESJI (KZ) = 8.716
SJIF (Morocco) = 5.667

ICV (Poland) = 6.630
PIF (India) = 1.940
IBI (India) = 4.260
OAJI (USA) = 0.350

определение темы или автора текста, фильтрация документов и писем, приходящих на электронную почту, модерация сообщений на форумах и в социальных сетях, индексирование документов, автоматическая генерация метаданных и т.д.

Формальное описание задачи классификации текстов

Классификация текстовых документов заключается в систематизации документов по категориям.

Множество документов представляется в виде:

$$D = \{d_1, \dots, d_i, \dots, d_n\} \quad (1)$$

Категории документов представляются множеством:

$$C = \{c_r\}, \text{ где } r = 1, \dots, m \quad (2)$$

В задаче классификации требуется на основе этих данных построить процедуру, которая заключается в нахождении наиболее вероятной категории из множества C для исследуемого документа d_i .

Большинство методов классификации текстов основаны на предположении, что документы, относящиеся к одной категории, содержат одинаковые признаки (слова или словосочетания), и наличие или отсутствие таких признаков в документе однозначно определяет его принадлежность или непринадлежность к той или иной категории.

Таким образом, для каждой категории должно быть определено множество признаков – словарь лексем, которые включают слова и/или словосочетания, характеризующие категорию:

$$F(C) = \bigcup c_r, \text{ где } F(c_r) = \langle f_1, \dots, f_k, \dots, f_z \rangle \quad (3)$$

Каждый документ также имеет признаки, по которым его можно отнести с некоторой степенью вероятности к одной или нескольким категориям:

$$F(d_i) = \langle f_1^i, \dots, f_l^i, \dots, f_y^i \rangle \quad (4)$$

Множество признаков всех документов должно совпадать с множеством признаков всех категорий, т. е.:

$$F(C) = F(D) = \bigcup F(d_i) \quad (5)$$

Отметим, что данные наборы признаков являются отличительной чертой классификации текстовых документов от классификации объектов в Data Mining, которые характеризуются набором атрибутов.

Решение об отнесении документа d_i к категории c_r принимается на основании пересечения признаков документа и признаков категории:

$$F(d_i) \cap F(c_r) \quad (6)$$

Задача методов классификации состоит в том, чтобы наилучшим образом выбрать такие признаки и сформулировать правила, на основе которых будет приниматься решение об отнесении документа к категории.

Методы классификации текстов

Существует два противоположных подхода [3, с. 226] к формированию множества признаков категорий $F(C)$ и построению правил: машинное обучение и экспертный метод.

До конца 80-х годов наиболее популярным подходом к классификации текстов был экспертный подход, состоящий в ручном определении набора правил, кодирующих экспертные знания о том, как классифицировать документы по данным категориям.

Достоинством данного подхода является отсутствие необходимости в построении обучающей выборки.

Недостатком данного подхода является необходимость в словарях, которые сложно построить для больших предметных областей.

В 90-е годы экспертный подход постепенно начал уступать место подходу с использованием машинного обучения, в соответствии с которым строится классификатор текста путем обучения на некотором наборе текстов.

Достоинствами данного подхода являются точность, сравнимая с той, которая достигается при экспертном методе, значительная экономия с точки зрения рабочей силы экспертов и возможность использования для различных предметных областей.

Недостатком данного подхода является то, что обучающая выборка должна включать достаточно документов для каждой из категорий, чтобы впоследствии избежать неправильной классификации.

Иногда применяется комбинация двух описанных подходов. В этом случае выделение ключевых слов для категорий и составление правил выполняются автоматически на основе обучающей выборки. В то же время составленные правила строятся в таком виде, чтобы в дальнейшем эксперт имел возможность корректировать эти правила вручную. В некоторых ситуациях такой подход очень полезен. Например, при возникновении проблемы переобучения классификатора. При использовании комбинации подходов после обучения классификатора у эксперта есть возможность проверки правил, удаления тех, которые вызывают сомнения и добавления некоторых новых правил для классификации к уже сгенерированным с помощью машинного обучения.

Технология Text Mining

Задачи обработки текстовой информации [3,4,5], и в том числе классификации текстов, относятся к задачам технологии Text Mining (интеллектуальный анализ текста).

Технология Text Mining представляет собой одну из разновидностей методов Data Mining (интеллектуального анализа данных) и

Impact Factor:

ISRA (India) = 4.971
ISI (Dubai, UAE) = 0.829
GIF (Australia) = 0.564
JIF = 1.500

SIS (USA) = 0.912
РИИЦ (Russia) = 0.126
ESJI (KZ) = 8.716
SJIF (Morocco) = 5.667

ICV (Poland) = 6.630
PIF (India) = 1.940
IBI (India) = 4.260
OAJI (USA) = 0.350

предназначена для решения задач обработки данных, представленных в текстовом виде. Основная проблема работы с текстами заключается в том, что текстовые документы практически невозможно преобразовать в количественное или табличное представление, поэтому информация хранится в исходном неструктурированном виде, что сильно затрудняет ее обработку и не позволяет использовать привычные алгоритмы Data Mining.

Как правило, Text Mining подразумевает процесс структурирования вводных текстовых данных, извлечение шаблонов из уже структурированных данных, и финальную оценку и интерпретацию полученных результатов [5].

Средства Oracle для реализации задач Text Mining

На данный момент существует большое множество различных программных продуктов [6] для реализации методов технологии Data Mining и Text Mining.

В частности, компания Oracle предоставляет два продукта для анализа и обработки текстовой информации.

1) Oracle Data Mining (Oracle Text Mining) – отдельно устанавливаемая опция программного обеспечения СУБД Oracle, доступная только для версии Enterprise Edition (EE). Чтобы использовать Oracle Data Mining, необходима лицензия для опции Data Mining. [7]

2) Oracle Text – это технология, включенная в базовую функциональность СУБД Oracle. Она использует внутренние компоненты Oracle Data Mining для предоставления некоторых возможностей интеллектуального анализа данных. Данный продукт существовал и в предыдущих версиях Oracle под названиями ConText и interMedia Text. [8]

Преимущества продуктов Oracle для анализа текстовой информации

- высокая производительность и широкие возможности поиска;
- ранжирование результатов;
- ведение статистики;
- возможность оптимизации поисковых индексов;
- встроенная поддержка русского языка (на уровне основ слов);
- возможность работы с наиболее популярными форматами документов, такими как ASCII, HTML, Microsoft Word, Adobe Acrobat (PDF) и WordPerfect.

Сравнение Oracle Data Mining и Oracle Text

Oracle Data Mining и Oracle Text обладают некоторыми различиями с точки зрения их применения для различных задач анализа данных.

Из табл.1 очевидно, что Oracle Text обладает меньшим набором возможностей, так как предназначен только для работы с текстом.

Но этот продукт имеет важное преимущество по сравнению с Data Mining – чтобы использовать Oracle Text и его возможности интеллектуального анализа данных, нет необходимости лицензировать опцию Data Mining. Начиная с 11g, все версии обеспечения Oracle поддерживают технологию Oracle Text, поэтому есть возможность использовать бесплатную версию Oracle Express Edition.

Исходя из рассмотренных особенностей продуктов Oracle, можно сделать вывод, что для автоматизации классификации текстовых документов лучше подходит Oracle Text.

Таблица 1. Возможности Oracle Data Mining и Oracle Text для решения задач анализа данных

Задачи анализа данных	Oracle Data Mining	Oracle Text
Обнаружение аномалий	SVM	Нет поддержки
Ассоциация	MDL	Нет поддержки
Важность атрибута	Apriori	Нет поддержки
Классификация	SVM, GLM или Naive Bayes	SVM, Decision Tree или с помощью правил записанных пользователем
Кластеризация	k -Means	k-Means
Извлечение функций	NMF	Нет поддержки
Регрессия	SVM или GLM	Нет поддержки

Алгоритмы классификации, доступные в Oracle Text

Для классификации текстовых документов успешно используются многие методы и алгоритмы классификации Data Mining: Naive Bayes, метод наименьших квадратов, C4.5, SVM и др. [9]

Очевидно, что требуется модификация этих методов для работы с текстовой информацией.

Oracle Text предлагает различные подходы к классификации документов:

- **rule-based classification** (классификация, основанная на правилах) – пользователь может написать правила классификации самостоятельно;

Impact Factor:

ISRA (India) = 4.971
ISI (Dubai, UAE) = 0.829
GIF (Australia) = 0.564
JIF = 1.500

SIS (USA) = 0.912
РИИЦ (Russia) = 0.126
ESJI (KZ) = 8.716
SJIF (Morocco) = 5.667

ICV (Poland) = 6.630
PIF (India) = 1.940
IBI (India) = 4.260
OAJI (USA) = 0.350

– **supervised classification** (контролируемая классификация) - Oracle Text создает правила классификации на основе набора обучающего набора документов (которые заранее классифицированы);

Классификация на основе правил является основным способом создания приложения классификации с помощью Oracle Text. Реализация этого подхода проста: создается таблица, в которую добавляются правила классификации, затем эти правила индексируются.

Но для реализации подхода с помощью машинного обучения необходимо использовать контролируемую классификацию, так как она позволяет автоматизировать процесс выявления правил для классификации, но при этом сохраняет возможность влиять на него.

Технология Oracle Text поддерживает два алгоритма контролируемой классификации текстов [8]: SVM и Decision Trees.

1) **Метод опорных векторов** (SVM – Support vector machine) – позволяет классифицировать объекты путем поиска гиперплоскости, которая разделяет два множества объектов. [9; 10]

Гиперплоскость – это функция. В самом простом случае гиперплоскость – это линия. Обычно ищут оптимальную гиперплоскость – такую, расстояние от которой до ближайшей точки множества максимально.

Идея для поиска проста - построить две параллельные гиперплоскости, между которыми нет никаких точек множества, и максимизировать расстояние между ними. Ближайшие к параллельным гиперплоскостям точки называются опорными векторами (support vectors), которые и дали название методу.

Как только гиперплоскость найдена, классификация становится делом техники: класс нового вектора определяется тем, с какой стороны он оказался от гиперплоскости.

Модификация данного алгоритма при классификации текстовых документов состоит в том, каждый документ представляется в виде вектора, в котором i -тый элемент - это мера вхождения i -того слова словаря (списка всех слов, которые встречаются в документах) в документ. В итоге получаются вектора большой размерности, где большинство элементов — нули.

При каждом проходе классификатор определяет, какие документы относятся к одной из категорий, поэтому его необходимо запускать столько раз, сколько категорий документов участвуют в классификации.

Преимущества алгоритма:

- высокая точность классификации;
- возможность работы с небольшим набором данных для обучения.

Недостатки алгоритма:

- сложная интерпретируемость результатов алгоритма;

- неустойчивость по отношению к выбросам в исходных данных;

- большое количество признаков сильно сказывается на длительности работы.

2) **Метод деревьев решений** (DT – Decision Trees) – позволяет сформировать правила классификации в виде иерархической древовидной структуры. [9; 11]

Деревом решений называют ациклический граф, по которому производится классификация объектов (в нашем случае текстовых документов), описанных набором признаков. Каждый узел дерева содержит условие ветвления по одному из признаков. Количество ветвлений зависит от того, сколько значений имеет выбранный признак. В процессе классификации осуществляются последовательные переходы от одного узла к другому в соответствии со значениями признаков объекта. Классификация считается завершенной, когда достигнут один из листьев (конечных узлов) дерева. Значение этого листа определяет класс, к которому относится рассматриваемый объект.

На практике обычно используют бинарные деревья решений, в которых условия ветвлений представлены простой проверкой наличия признака в документе.

Преимущества алгоритма:

- относительно простая программная реализация алгоритма;

- легкая интерпретируемость результатов работы алгоритма.

Недостатки алгоритма:

- неустойчивость алгоритма по отношению к выбросам в исходных данных;

- большой объем данных для получения точных результатов.

Особенности реализации алгоритмов в Oracle Text

Oracle Text включает в себя несколько пакетов, которые позволяют выполнять различные действия, начиная от синхронизации индекса Oracle Text до выделения разделов в документах [8].

Для решения задачи классификации необходимо иметь привилегии на запуск двух пакетов:

- Пакет CTX_CLS позволяет выполнять классификацию документов.

- Пакет CTX_DDL содержит процедуры и функции для создания и управления настройками необходимыми для текстовых индексов.

Этапы реализации классификации в Oracle Text

1. Выбор классификатора

Impact Factor:

ISRA (India) = 4.971
ISI (Dubai, UAE) = 0.829
GIF (Australia) = 0.564
JIF = 1.500

SIS (USA) = 0.912
РИИЦ (Russia) = 0.126
ESJI (KZ) = 8.716
SJIF (Morocco) = 5.667

ICV (Poland) = 6.630
PIF (India) = 1.940
IBI (India) = 4.260
OAJI (USA) = 0.350

Для каждого из алгоритмов в Oracle Text предусмотрен свой тип классификатора:

- SVM_CLASSIFIER – для реализации алгоритма SVM;
- RULE_CLASSIFIER – для реализации алгоритма Decision Tree.

При использовании типа SVM_CLASSIFIER правила генерирует в двоичном формате, а при использовании типа RULE_CLASSIFIER – в виде логических правил, которые легко проверяются и изменяются человеком.

2. Индексирование текстов в Oracle Text

Текстовые возможности СУБД Oracle основаны на использовании специального вида полнотекстовых индексов. Полнотекстовый индекс — это индекс, в котором перечислены все слова, встречающиеся в тексте, и указаны позиции, на которых эти слова встречаются.

Oracle Text имеет в готовом виде полнотекстовые индексы четырех типов. Но при решении задачи классификации текстовых документов нужны только два из них:

1) CTXSYS.CONTEXT – основной тип индекса – применяется для индексирования содержимого одной или нескольких колонок таблицы.

2) CTXSYS.CTXRULE - применяется при реализации классификации документов.

3. Обучение классификатора

При контролируемой классификации для обучения классификатора и автоматизации этапа написания правил применяется процедура CTX_CLS.TRAIN, которая использует обучающий набор документов для выведения правил классификации.

В результате обучения генерируются правила классификации, которые записываются в специальную таблицу. Сгенерированные правила необходимо проиндексировать путем создания индекса CTXSYS.CTXRULE.

4. Классификация новых документов

Классификация новых документов производится путем простого SQL запроса с использованием оператора MATCHES.

Если используется тип классификатора SVM_CLASSIFIER, то оператор MATCHES возвращает оценку в диапазоне от 0 до 100; более высокое число указывает на большую уверенность в результате классификации.

Если используется тип классификатора RULE_CLASSIFIER, то оператор MATCHES

возвращает либо 100 (документ соответствует критериям), либо 0 (не соответствует).

Описание возможной программы для классификации текстовых документов

Для решения задачи классификации текстовых документов предлагается создать набор процедур на языке PL/SQL, которые будут выполнять следующие действия:

1) Создание объектов базы данных, необходимых для хранения информации о текстовых документах, категориях классификации, результатах классификации.

2) Подготовка данных для обучения классификатора. Этот этап включает загрузку информации о документах (в том числе текстов документов), категориях и вспомогательных данных в базу данных, а также предварительную обработку текстовых данных.

3) Обучение классификатора на подготовленных ранее обучающих данных.

4) Классификация текстовых документов из тестового набора с помощью полученных правил.

Заключение

В статье рассмотрена задача классификации текстовых документов.

Для автоматизации процесса решения этой задачи необходимо определить правила классификации. Они могут быть получены двумя методами – экспертным или с помощью машинного обучения. Оба подхода применимы при реализации решения задачи средствами Oracle.

Рассмотрены возможности двух продуктов Oracle – Oracle Data Mining и Oracle Text. На основе анализа их преимуществ и недостатков сделан вывод, что для решения задачи классификации наилучшим выбором является Oracle Text. При изучении данного продукта определено, что для обучения классификатора может быть использован один из двух алгоритмов – SVM и Decision Tree. SVM следует применять в случаях, когда нужна высокая точность, а Decision Tree – когда нужно, чтобы основная часть правил была сгенерирована автоматически, но при этом оставалась возможность отредактировать наборы правил при необходимости.

Также в статье были рассмотрены особенности этапов реализации классификации текстов с использованием Oracle Text и предложен порядок разработки программы на PL/SQL для решения этой задачи.

Impact Factor:	ISRA (India) = 4.971	SIS (USA) = 0.912	ICV (Poland) = 6.630
	ISI (Dubai, UAE) = 0.829	PIHHI (Russia) = 0.126	PIF (India) = 1.940
	GIF (Australia) = 0.564	ESJI (KZ) = 8.716	IBI (India) = 4.260
	JIF = 1.500	SJIF (Morocco) = 5.667	OAJI (USA) = 0.350

References:

- (2020). *Federal'nyi zakon ot 24 aprelya g. N 122-FZ «O provedenii eksperimenta po ispol'zovaniyu elektronnykh dokumentov, svyazannykh s rabotoi»*. Opublikovan 24.04. na oficial'nom internet-portale pravovoj informacii Retrieved from <http://www.pravo.gov.ru>.
- Sebastiani, F. (2002). Machine Learning in Automated Text Categorization. *ACM Computing Surveys*, Vol.34, No.1, March, pp. 1-47.
- Barsegyan, A.A. (2009). *Analiz dannykh i protsessov: ucheb. posobie / A.A. Barsegyan, M.S. Kupriyanov, I.I. Kholod, M. D. Tess, S.I. Elizarov. — 3-e izd., pererab. i dop. (p.512)*. SPb.: BKhV-Peterburg.
- Shkurina, M.V., & Sabinin, O. Y. (2018). An overview and analysis of automatic text summarization methods. *ISJ Theoretical & Applied Science*, 12 (68), 282-286.
- Kutukova, E.S. (2013). Tekhnologiya Text mining. *Sbornik nauchny'x trudov SWorld*, 30 (4), 33-36.
- Kotel'nikova, Y.E. (2010). Obrabotka tekstovykh dokumentov i evolyutsiya avtomatizirovannykh sistem proektirovaniya. *Izvestiya vysshikh uchebnykh zavedenii. Priborostroenie*, 53 (6), 21-25.
- (n.d.). *Oracle Data Mining Concepts, 11g Release 2 (11.2)*, URL: Retrieved from https://docs.oracle.com/cd/E11882_01/datamine.112/e16808/title.htm (Date of access: 09.05.20).
- (n.d.). *Oracle Text Application Developer's Guide, 11g Release 2 (11.2)*, URL: Retrieved from https://docs.oracle.com/cd/E18283_01/text.112/e16594/title.htm (Date of access: 09.05.20).
- Batura, T.V. (2017). Metody` avtomaticheskoy klassifikacii tekstov. *Programmny'e produkty` i sistemy`*, 30 (1), 85-89. doi: 10.15827/0236-235X.030.1.085-099.
- Demidova, L.A., Nikul'chev, E.V., & Sokolova, Yu.S. (2016). Klassifikatsiya bol'shikh dannykh: ispol'zovanie SVM-ansamblei i SVM-klassifikatorov s modifitsirovannym roevym algoritmom. *Cloud of science*, 3 (1), 5-42.
- Kaftannikov, I.L., & Parasich, A.V. (2015). Osobennosti primeneniya derev'ev reshenii v zadachakh klassifikatsii. *Vestnik Yuzhno-Ural'skogo gosudarstvennogo universiteta*. Seriya: Komp'yuternye tekhnologii, upravlenie, radioelektronika, 15 (3), 26-32.