



МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ
ИМЕНИ М.В. ЛОМОНОСОВА

Факультет вычислительной математики и кибернетики
Кафедра автоматизации систем вычислительных комплексов

Треско Константин Игоревич

Исследование и разработка средств многотемной классификации веб-страниц

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА

Научные руководители:

к.ф.-м.н. М.И. Петровский

Д.В. Царев

Аннотация

Целью данной работы является исследование и разработка экспериментального прототипа системы сбора и многотемной классификации веб-страниц, с которыми работали пользователи, находящиеся внутри одной сети. Отличие многотемной классификации заключается в том, что классы могут пересекаться или даже быть вложенными.

Разрабатываемая система сбора должна удовлетворять требованиям масштабируемости (линейного роста расхода ресурсов при увеличении числа подключений) и защищенности (пользователь не должен иметь возможности фальсифицировать собранные данные). Модуль классификации данных должен обеспечивать многотемную классификацию на основе машинного обучения с возможностью добавления и удаления тематик.

В ходе работы был проведен обзор существующих средств многотемной классификации и средств реализации системы сбора. Также был реализован экспериментальный прототип системы сбора и многотемной классификации, удовлетворяющий поставленным выше требованиям, проведено тестирование, показавшее масштабируемость и защищенность разработанного средства.

1 Введение

1.1 Задача классификации

Настоящая работа посвящена исследованию и разработке программных средств сбора и многотемной классификации текстовых данных веб-страниц. Задача классификации многотемных документов (multi-label classification), заключается в определении принадлежности документа к одному или нескольким классам (из predetermined набора классов) на основании анализа совокупности признаков, характеризующих данный документ. В отличие от традиционной задачи классификации, классы могут пересекаться или быть вложенными, то есть документ может принадлежать нескольким классам. Классы, к которым принадлежит документ называются релевантными.

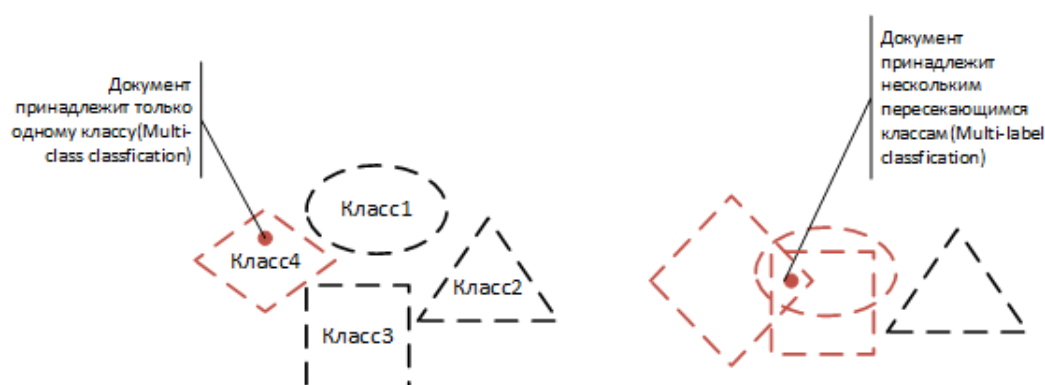


Рисунок 1 – Многотемная и многоклассовая классификация

1.2 Области применения классификации веб-данных

На сегодняшний день существует ряд задач, для решения которых требуются системы сбора и классификации контента, частным случаем которого является веб-контент:

- **Анализ работы сотрудников организации.** Определение тематик документов, с которыми работает пользователь. На основе полученной информации применение политик безопасности.

- **Фильтрации контента, для обнаружения информации, относящейся к определенному делу (eDiscovery[3]).** Документы и электронная почта могут фильтроваться в зависимости от присвоенных им классов, чтобы гарантировать, что только материалы с запрошенной бизнес-информацией в рамках расследуемого дела были найдены и сохранены.
- **Определения конфиденциальности документа, для предотвращения утечек информации** По различным экспертным оценкам в настоящее время наибольшие риски для информационной безопасности представляют не внешние, а внутренние угрозы [1]. Поэтому существует необходимость в средствах определения конфиденциальности документа.

Далее будут рассмотрены классы индустриальных систем, функционал которых включает в себя классификацию текстовых данных (в том числе веб-данных), с которыми работают пользователи.

- **Системы управления корпоративным контентом (англ. Enterprise Content Management, ECM)** - программные решения для управления информационными ресурсами предприятия предоставляют программные средства сбора, анализа, управления, накопления, хранения и доставки документов в масштабах организации. В настоящее время большинство ECM систем включают в себя системы обнаружения данным, связанных с определенным судебным делом (англ. eDiscovery). Средства eDiscovery обеспечивают процесс, с помощью которого организации находят, получают, сохраняют и анализируют документы, связанные с делом.
- **Системы предотвращения утечки данных (англ. Data Loss Prevention, DLP)** - программные решения для предотвращения утечек конфиденциальной информации и минимизации других рисков, связанных с внутренними угрозами.

1.3 Основные методы сбора и классификации данных в корпоративных системах

Поиск и анализ информации в ECM системах происходит следующим образом:

- **Сбор данных.** Системы сбора содержат компоненту, установленную на пользовательском компьютере и отвечающую за мониторинг деятельности пользователя. В терминологии IBM данные компоненты называются Искателями [6].
- **Применение методов анализа документов.** В терминологии IBM данный этап называется аналитическим конвейером [7].

- **Индексация данных.** Компоненты индексации добавляют в индекс информацию о новых и измененных документах [8, 9, 11].

Существуют три основных подхода к классификации данных в ЕСМ системах

- **Классификация на основе обучающей выборки** (IBM Content Classification [2], Symantec eDiscovery [4]).
- **Классификация на основе заданных правил.** Принадлежность документа к тому или иному классу определяется на основе эвристических правил (сигнатур). Правила могут формироваться на основе контекста – имя отправителя, директория создания и т.п., а также на основе контента – шаблоны текста, ключевые слова и т.п.
- **Определение категорий в неизвестных документах.** В задаче кластеризации нет предопределенного набора классов. Исходное множество документов разбивается на подмножество таким образом, чтобы документы в различных подмножествах существенно отличались.

В DLP системах выделяют следующие технологии классификации данных:

- **Цифровые отпечатки (англ. digital fingerprint)** - технология предназначена для защиты больших по объему документов, содержание которых не изменяется или меняется незначительно. Детектор цифровых отпечатков позволяет автоматически обнаруживать в анализируемом тексте цитаты из документов-образцов, содержащих конфиденциальную информацию.
- **Анализ шаблонов** - технология предназначена для детектирования алфавитно-цифровых объектов по шаблону данных (маске) и позволяет наиболее эффективно выявлять факты пересылки персональных данных или финансовой информации. Кроме того, данная технология может использоваться как вспомогательный метод для обнаружения фактов несанкционированной пересылки внутренних документов, содержащих формализованные данные, образованные по определенному шаблону (например, договоров или счетов в случае детектирования банковских реквизитов, кодов классификаторов и т.д.).
- **Машинное обучение** (InfoWatch [5], Symantec VML [12]). Примерная схема работы классификатора такова: на вход подается обучающий набор документов, состоящий как из конфиденциальных так и не конфиденциальных документов. По этим наборам производится обучение и строится статистическая модель. Далее полученная модель используется для классификации неизвестным документов. При обнаружении конфиденциального документа производятся действия, предписанные политиками безопасности.

Основное преимущество машинного обучения заключается в том, что в отличие от других описанных технологий, оно предназначено для работы не со статическими, а с постоянно меняющимися документами.

Таблица 1 – Подходы к классификации

	Достоинства	Недостатки
Цифровые отпечатки - обнаружение в тексте цитат из документов-образцов	Высокая точность детектирования статичных документов	Чувствительность к текстовым изменениям
Анализ шаблонов - анализ текстов на основе словарей и регулярных выражений	Эффективность в детектировании формализованных данных	Не применим для детектирования неформализованных данных
Машинное обучение – построение на основе обучающего набора статистической модели	Работают напрямую с содержимым документов Обучаемость Возможность обобщения	Требуется формирование обучающего набора Возможность ложноотрицательных и ложноположительные срабатываний Возможность обобщения

1.4 Специфика задачи классификации текстовых веб-данных

Веб – страницы имеют следующие особенности:

- Содержимое веб-страниц постоянно меняется.
- Содержимое имеет многотемную природу, то есть каждая из страниц может одновременно принадлежать к нескольким тематикам.



Рисунок 2 – Многотемная природа документа

1.5 Выводы

На сегодняшний день существует ряд прикладных задач, требующих классификации текстовых данных. Обзор индустриальных систем показал, что наиболее актуальными технологиями классификации текстовых данных являются:

- Метод шаблонов.
- Метод цифровых отпечатков.
- Метод машинного обучения.

Первые две технологии применимы только к статическим данным, в то время как метод машинного обучения позволяет адаптироваться к тому, что содержимое и состав анализируемых данных постоянно меняется.

Для задачи классификации текстовых веб-данных из рассмотренных технологий подходит только машинное обучение с возможностью дообучения, так как веб – данные и набор классов постоянно меняются.

Рассмотренные системы решают задачу многоклассовой классификации, так как документ может принадлежать только к одному из предопределенного набора

классу. Так как большая часть веб – данных имеет многотемную природу, то существует актуальность разработки системы сбора и многотемной классификации текстовых веб - данных пользователей. С учетом большого объема анализируемой информации к модулю классификации предъявляется требование масштабируемости.

Так как в рассмотренных системах присутствуют компоненты сбора информации, расположенные на пользовательских машинах, то при разработке системы сбора нужно учитывать то, что деятельность компонент никак не должна сказываться на работе пользователя. Также в рассмотренных системах пользователь не имеет доступа к собираемой информации, поэтому система сбора должна удовлетворять требованию защищенности.

2 Постановка задачи

Разработать архитектуру и реализовать прототип системы сбора и многотемной классификации текстовых веб-данных пользователя в соответствии с требованиями:

- Модуль сбора должен обеспечивать:
 - Масштабируемость (линейный рост расхода ресурсов при росте числа подключений).
 - Защищенность (пользователь не должен иметь возможности фальсифицировать данные).
 - Функционирование под ОС Windows и браузером IE.
- Модуль классификации должен обеспечивать:
 - Многотемную классификацию на основе машинного обучения с возможностью дообучения.
 - Возможность добавления и удаления тематик.

3 Обзор

Для реализации прототипа системы многотемной классификации веб-страниц, с которыми работал пользователь необходимо осуществить:

- **Сбор** просмотренных пользователем веб-страниц.
- **Модуль классификации**, который будет определять принадлежность документа предопределенным классам.

Компонента модуля сбора будет установлена на каждой из пользовательских машин и будет состоять из:

- Расширения для браузера, сохраняющего html код просматриваемой пользователем веб – страницы в локальную базу данных, расположенную на пользовательском компьютере.
- Модуля сбора данных.

3.1 Расширение для браузера

Для решения задачи необходимо расширение для браузера, имеющее возможность сохранить html код просматриваемой пользователем веб – страницы в локальную базу данных, без подтверждения пользователем.

- **ВНО** (Browser Helper Object) - DLL-модуль, разработанный как плагин для Internet Explorer для обеспечения дополнительной функциональности. В ВНО API существует возможность получения доступа к DOM (Document Object Model[15]) текущей страницы. Также существует возможность обработки событий и управления навигацией. Для задачи перехвата контента данное решение подходит. Минусом является то, что данное решение подходит только для браузера IE. Но ВНО имеет возможность записи и чтения данных из файловой системы пользователя, без его участия.
- **Kynext** позволяет перехватывать содержимое веб – страницы. Поддерживает IE, Firefox, Safari, and Chrome, но для написания необходимо использовать проприетарный язык Kynetx Rules Language. Еще одним минусом является то, что работа написанного расширения зависит от работоспособности самого расширения Kynext.
- **WebMynd** поддерживает IE, Firefox, Safari, and Chrome. Пока доступна бета версия, и не известно, будет ли дальнейшая поддержка продукта. Поддерживает JavaScript API. Не является бесплатным программным обеспечением.

Таблица 2 – Средства расширения для браузера

	ВНО	Crossrider	Kynext	WebMynd
Кроссплатформенность	-(только IE)	+	+	+
Бесплатность	+	+	+	-
Поддерживаемые языки	C++,C#	JavaScript	Kynext Rules Language	JavaScript
Возможность записи в локальную ФС	+	-	-	-

- **Crossrider** позволяет быстро создавать кросс-браузерные расширения. Использует один API и поддерживает JavaScript и jQuery, так что разработчик с базовыми знаниями JavaScript может писать и поддерживать свой код. Имеется возможность писать под IE, Firefox, Safari, and Chrome. Бесплатен в использовании. Существует документация и демо – видео для некоторых задач. Доступ к файловой системе пользователя только с его разрешения.

3.2 Методы многотемной классификации

Можно выделить три основных подхода к многотемной классификации: [27]

- «Оптимизационный» подход [20, 21, 28].
- Подход на основе декомпозиции в набор независимых бинарных проблем.
- Подход на основе ранжирования [29, 30].

3.2.1 Методы, основанные на "оптимизационном" подходе

К числу алгоритмов, использующих данный подход, относятся алгоритмы, решающие оптимизационную задачу. Рассмотрим некоторые из данных алгоритмов:

- Основанные на AdaBoost алгоритме (AdaBoost.MH [20], ADTBoost.MH [21])
- минимизируется функция Hamming Loss), где

$$HammingLoss = \frac{1}{k} \sum_{i=1}^{k-1} \frac{1}{|Y|} |(f_Z(x_i)) \nabla(y_i)|$$

где $a \nabla(b) = (a \cup b) \setminus (a \cap b)$, $a \subseteq Y, b \subseteq Y$, k - размер тестового набора. Другими словами, критерий оценивает, сколько раз пара «пример-метка» классифицируется не правильно.

- Multi-Label-kNN - [22] максимизируются апостериорные вероятности принадлежности классам

. **Методы, основанные на алгоритме AdaBoost.** К числу данных методов можно отнести AdaBoost.MH [20] и ADTBoost.MH [21]. AdaBoost вызывает слабые классификаторы в цикле $t = 1, \dots, T$. После каждого вызова обновляется распределение весов D_t , которые отвечают важности каждого из объектов обучающего множества для классификации. На каждой итерации веса каждого неверно классифицированного объекта возрастают, таким образом новый комитет классификаторов «фокусирует своё внимание» на этих объектах.

Алгоритм работы AdaBoost [10]:

Дано: $(x_1, y_1), \dots, (x_m, y_m)$ $x_i \in X, y_i \in Y = \{-1, +1\}$

Изначально присваиваем всем классам одинаковый вес $D_t(i) = \frac{1}{m}, i = 1, \dots, m$.
for $t = 1, \dots, T$

- Находим классификатор, минимизирующий взвешенную ошибку классификации $h_t = \arg \min_{h_j \in \mathcal{H}} \epsilon_j$, где $\epsilon_j = \sum_{i=1}^m D_t(i) [y_i \neq h_j(x_i)]$.
- Если величина $\epsilon_t \geq 0.5$, то останавливаемся.
- Выбираем $\alpha_t \in \mathbf{R}$, обычно $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$, где ϵ_t взвешенная ошибка классификатора h_t .
- $D_{t+1}(i) = \frac{D_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$, где Z_t является нормализующим параметром (выбранным так, чтобы D_{t+1} являлось распределением вероятностей, то есть $\sum_{i=1}^m D_{t+1}(i) = 1$).

Итоговый классификатор имеет вид: $H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$

Таким образом, на каждом шаге выбирается оптимальный классификатор, и веса тех объектов x_i , которые он предсказал правильно, уменьшаются, а тех, которые он предсказал неверно - увеличиваются, чтобы на следующем шаге был выбран классификатор, который лучше распознает объекты, неверно распознанные предыдущим.

С точки зрения сформулированных требований метод AdaBoost имеет следующие недостатки:

- Алгоритм не имеет возможности пошагового дообучения, поскольку принцип работы метода рассчитан на то, что обучающий набор задан заранее целиком.

- Для достижения качества классификации, удовлетворяющего потребностям актуальных прикладных задач, этот метод требует большое число шагов на этапе обучения, а следовательно, обучение выполняется достаточно долго.

В случае многотемного AdaBoost заключительная гипотеза для ранжирования классов получается как взвешенное голосование слабых гипотез (гипотез, полученных при решении задачи многоклассовой классификации), причём вес слабой гипотезы вычисляются с учётом качества классификации этой гипотезой примеров из обучающего набора. Заключительная гипотеза multi-label классификации получается на основе гипотезы ранжирования путём отсечения по нулевому порогу релевантности.

Multi-Label-kNN. Данный метод представляет собой подход многотемного обучения, основанный на методе k ближайших соседей. Основным принципом метода ближайших соседей является то, что объект присваивается тому классу, который является наиболее распространённым среди соседей данного элемента. В многотемном варианте kNN расстояние выражается через косинус угла между векторами признаков элементов: $\rho = 1 - \cos(\vec{a}_1, \vec{a}_2)$.

Пусть H_1^l - событие, что тестовый пример \vec{a} имеет метку l , H_0^l - иначе, $E_k^l (j \in 0, \dots, k)$ - событие, что j соседей имеют метку l , $\vec{C}_{\vec{a}}$ - вектор подсчета членства. Вектор категорий $\vec{y}_{\vec{a}}$ определяется как:

$$\vec{y}_{\vec{a}} = \arg \max_{b \in 0,1} P(H_b^l | E_{\vec{C}_{\vec{a}}}^l), l \in Y$$

, где Y - набор классов.

Данный метод имеет следующие недостатки:

- Большие вычислительные затраты.
- Большой объём памяти, необходимый для хранения и работы сразу со всем набором обучающих данных.

3.2.2 Методы, основанные на декомпозиции в набор независимых бинарных проблем

Работу алгоритмов можно описать следующей последовательностью действий:

- Декомпозиция исходной проблемы в набор независимых бинарных проблем («каждый-против-остальных»). Для каждого из N классов создается одна бинарная проблема.
- В бинарной проблеме для класса l все обучающие примеры, помеченные этим классом, считаются положительными, а все остальные обучающие примеры считаются отрицательными.
- Далее к каждой из N полученных бинарных проблем применяется бинарный алгоритм обучения (например, двухклассовая SVM [13]).

- В результате получаются бинарных классификаторов (гипотез), каждый из которых независимо от остальных оценивает релевантность каждого из классов.
- Для каждого из N классов будет построена решающая функция бинарной классификации.
- Знак каждой решающей функции дает предсказание релевантности темы l для данного документа.

Достоинства и недостатки данного подхода:

- Так как решение о принадлежности документа принимается независимо от остальных классов, то имеется возможность добавления и удаления классов без необходимости обучения "с нуля".
- С использованием метода опорных векторов достигается высокая точность классификации, отсутствует возможность пошагового обучения.
- При применении алгоритма персептрона для бинарной классификации обеспечивается возможность дообучения, но в результате, метод, как правило, имеет низкую точность.
- Из-за того, что количество бинарных подзадач равно числу классов, подход имеет высокую вычислительную сложность.
- Строятся независимые классификаторы, которые не учитывают корреляции между классами.

3.2.3 Методы, основанные на подходе ранжирования с последующим отсечением нерелевантных классов

Работу алгоритмов можно разделить на два этапа

- Первый этап состоит в обучении алгоритма ранжирования, который упорядочивает все классы по степени их релевантности для заданного классифицируемого объекта (ММР [23, 24], k-NN [25]).
- Второй этап заключается в построении функции многотемной классификации, отделяющей релевантные классы от нерелевантных.

Метод Multiclass-Multilabel Perceptron. Алгоритм ранжирования ММР поддерживает набор из n прототипов тем - $\vec{w}_1, \dots, \vec{w}_n$. При получении документа алгоритм модифицирует предложенную гипотезу, то есть изменяет набор прототипов (весов) $\vec{w}_1, \dots, \vec{w}_n$. Это осуществляется для всех документов из обучающей выборки. Окончательной гипотезой является набор прототипов после одного прохода.

Пусть \vec{x} - документ, $\vec{w}_1, \dots, \vec{w}_n$ - набор прототипов. Ранжирование осуществляется следующим на основе скалярного произведения, т.е. тема r ранжируется выше, чем тема s , если $(w_r, x) > (w_s, x)$.

Пусть $y \in Y$ - набор релевантных тем, где Y - множество тем. Совершенным называется такое ранжирование, при котором, для $\forall s \in y$ и для $\forall k (k \in Y \setminus y)$ s ранжируется выше, чем k , т.е. $(s, y) > (k, y)$.

Качество совершенного ранжирования определяется размером промежутка между самой низкой оценкой среди релевантных тем и самой высокой оценкой среди нерелевантных тем: $\min \vec{w}_y - \max \vec{w}_k$.

Алгоритм ранжирования ММР получает и анализирует обучающие данные последовательно, пример за примером. Данный алгоритм удобен при работе с очень большими наборами обучающих данных, так как он эффективен по памяти и имеет возможность дообучения. Однако данный алгоритм не имеет возможность добавления тем без необходимости заново обучать модель ранжирования.

3.3 Метод классификации многотемных документов на основе подхода попарных сравнений

В лаборатории Технологий Программирования было предложено новое решение, которое включает:

- Новый алгоритм ранжирования, основанный на модифицированном для случая существенно пересекающихся классов методе попарных сравнений с помощью набора бинарных классификаторов и вычислении степеней принадлежности классам с использованием обобщенной модели Брэдли-Терри с «ничьей» [26].
- Новый алгоритм построения пороговой функции отсека нерелевантных классов, строящий пороговую функцию не в исходном пространстве характеристик (как это делает большинство традиционных алгоритмов), а в пространстве релевантностей классов, что позволяет упростить вид пороговой функции, значительно сократить вычисления и в большинстве случаев увеличить точность

3.4 Выводы из обзора

К компоненте сбора предъявлялись следующие требования:

- Работа компоненты не должна сказываться на деятельности пользователя.

– Пользователь не должен иметь доступа к собранным данным.

Из рассмотренных средств написания расширения для браузера обоим требованиям удовлетворяет ВНО. Первое требование удовлетворяется возможностью записи в локальную файловую систему без подтверждения. Второе требования удовлетворяется установкой прав доступа на папку, в которую сохраняются собранные данные.

Таблица 3 – Средства расширения для браузера

Подход	Недостатки
Оптимизационный подход	Нет возможности дообучения Большие затраты памяти Нет возможности удаления и добавления классов
Подход на основе декомпозиции в набор независимых бинарных проблем	Нет учета корреляции между классами Высокая вычислительная сложность
Подход на основе ранжирования	Высокая вычислительная сложность Нет возможности дообучения

Обзор существующих методов классификации многотемных документов показал, что обозначенным в постановке задачи требованиям удовлетворяют только методы, основанные на подходе «каждый-против-остальных», но качество классификации не удовлетворяет потребностям прикладных задач.

Поэтому был выбран алгоритм, разработанный в лаборатории Технологий Программирования, основанный на подходе попарных сравнений и удовлетворяющий всем поставленным требованиям.

4 Исследование и построение решения

Проведенный анализ систем ЕСМ и DLP в контексте задач классификации показывает, что каждый из пользователей является источником анализируемой информации, а для анализа собираемая информация должна храниться в единой базе данных. Для управления потоком данных из различных источников используется мультиагентная система [19]. Под агентом подразумевается автономный процесс, способный реагировать на среду исполнения и вызывать изменения в среде, возможно, в кооперации с пользователями или с другими агентами.

4.1 Разбиение на подзадачи

Для решения задачи сбора и многотемной классификации текстовых данных веб – страниц пользователя был выбран мультиагентный подход [19], поэтому система состоит из:

- **Агента мониторинга**, собирающего информацию и отправляющего ее агенту консолидации
- **Агента консолидации**, сохраняющего информацию со всех источников в единую базу данных
- **Модуля классификации**, предоставляющего возможность определения принадлежности документа к одному из предопределенного набора классов

4.1.1 Агент мониторинга

Агент мониторинга состоит из нескольких компонент:

- Расширение для браузера.
- Модуль передачи данных агенту консолидации.

Расширение для браузера, написанное с помощью ВНО, считывает html код просматриваемой пользователем веб-страницы и сохраняет ее в локальную базу данных. Сохранение в локальную базу данных осуществляется для контроля нагрузки на агент консолидации и возможности отправки данных по расписанию, а также на случай потери связи с агентом консолидации.

Модуль передачи данных агенту подключается к локальной базе данных и отправляет хранящуюся в ней информацию агенту консолидации, при получении

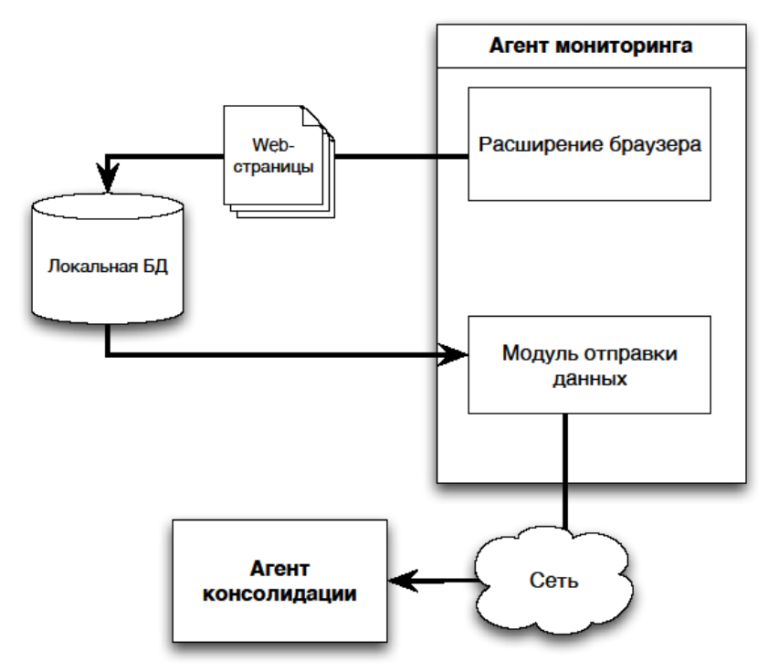


Рисунок 3 – Архитектура агента мониторинга

ответа от агента консолидации отправленные данные удаляются из локальной базы данных, чтобы размер локальной базы данных не увеличивался постоянно.

Каждый из пользователей является источником собираемой информации. При этом задача требует того, чтобы пользователь не имел доступа к собираемым данным. В случае реализации агента мониторинга с помощью ВНО защищенность может осуществляться с помощью установки прав доступа на папку, в которой хранятся собираемые данные.

4.1.2 Агент консолидации

Агент консолидации сохраняет данные, полученные от всех агентов мониторинга, в единую базу данных. Также в базу данных заносится дополнительная информация о страницах:

- Логин пользователя, посетившего веб-страницу.
- Имя компьютера, на котором находится пользователь.
- Дата посещения.
- ID веб-страницы.

При реализации необходимо учитывать то, что количество пользователей может быть достаточно велико, поэтому необходимо обеспечить масштабируемость агента сбора.

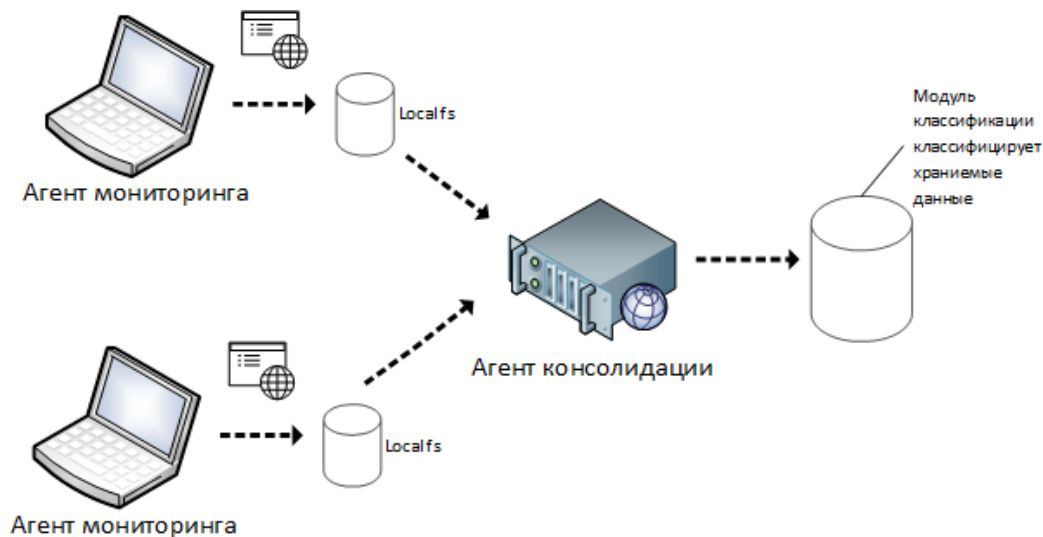


Рисунок 4 – Архитектура агента сбора

4.1.3 Модуль классификации

В ходе обзора было принято решение использовать модуль классификации, реализованный в лаборатории Технологий Программирования, основанный на подходе попарных сравнений (декомпозиция типа «каждый-против-каждого»).

Предложенное решение включает новый алгоритм ранжирования, основанный на модифицированном для случая существенно пересекающихся классов, и новый алгоритм построения пороговой функции отсеечения нерелевантных классов. Он

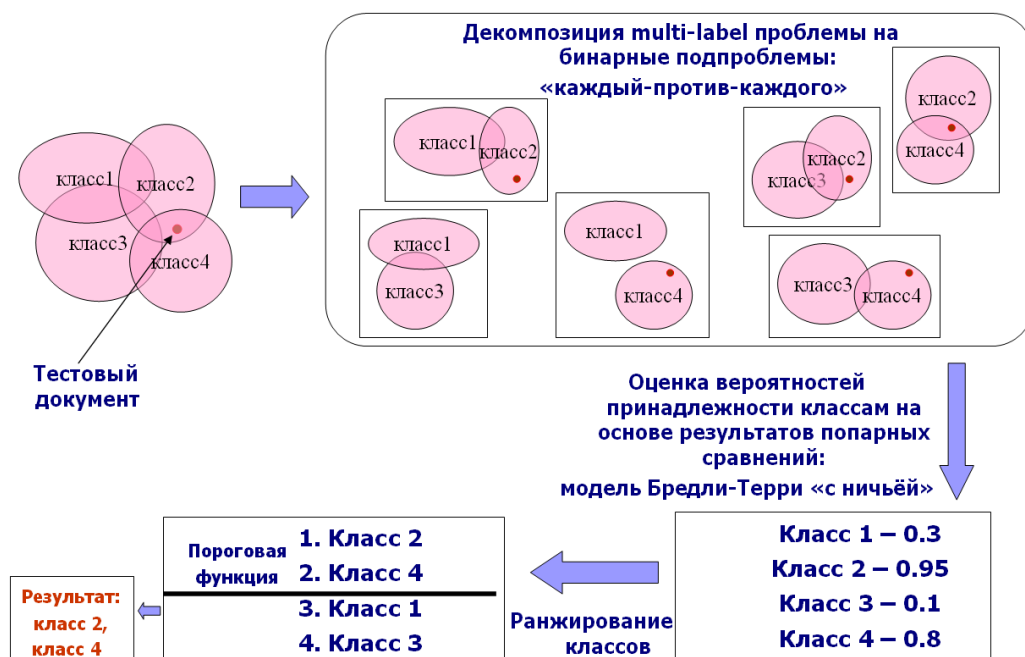


Рисунок 5 – Работа модуля классификации

предоставляет следующие сценарии работы:

- **Обучение.** Построение модели классификации на основе совокупности заранее рубрицированных гипертекстовых документов.
- **Классификация.** Применение построенной модели к новому классифицируемому документу.
- **Дообучение.** Модификация модели классификации на основе дообучения на новых документах с релевантными для них тематиками.
- **Удаление темы.** Удаление тематики классификации из модели без необходимости последующего обучения "с нуля".

5 Описание практической части

В данном разделе будет рассмотрена архитектура разработанного программного средства, представлен основной сценарий работы с ним, а также пояснены детали реализованных механизмов. Также будут приведены некоторые характеристики функционирования разработанного средства.

5.1 Общая архитектура разработанного средства

При разработке прототипа системы использовался мультиагентный подход, прототип состоит из:

- Агента мониторинга.
- Агента консолидации.
- Модуля многотемной классификации.

Архитектура реализованного прототипа представлена ниже.



Рисунок 6 – Архитектура прототипа

5.1.1 Агент мониторинга

Задача сбора просмотренных пользователем веб-страниц осуществляется с помощью расширения для браузера IE.

При открытии окна браузера расширение создает подключение к существующей базе данных, реализованной с помощью SQLite. Если базы данных нет, то создается новая, имеющая таблицу со следующими полями:

- Имя пользователя
- Имя компьютера
- URL просмотренной страницы
- Дата
- HTML код страницы

Каждый раз, когда пользователь загружает новую страницу, происходит событие, по которому html код просмотренной веб-страницы сохраняется в локальную базу данных.

Расширение написано с помощью ВНО (Browser Helper Object) на языке программирования C#, объем кода - 350 строк

5.1.2 Агент сбора

Каждый агент, расположенный на пользовательском компьютере, подключается к базе данных, в которую были записаны данные расширением для браузера, и отправляет данные агенту сбора.

При настройке системы можно указать временной интервал, по которому будут отправляться данные. При отправке сохраняется время, когда была отправлена последняя просмотренная пользователем веб-страница.

При получении ответа от агента сбора, все записи, просмотренные до сохраненной временной метки, удаляются, так как агент сбора успешно их сохранил, если же ответ не получен, то посылаются все записи, которые хранятся в базе.

Соединение агента мониторинга и агента сбора осуществляется с помощью TCP/IP сокетов. Каждый новый агент сбора обрабатывается в агенте консолидации асинхронно в отдельном потоке, что обеспечивает высокую скорость взаимодействия.

Детали реализации модуля передачи данных агенту консолидации:

- Количество строк кода - 250.
- Язык реализации - C#.
- Доступ к локальной файловой системе без подтверждения пользователя осуществляется с помощью утилиты `icacls`. [16], которая позволяет менять Integrity Levels [16] файлов. Так как Internet Explorer имеет право записи только в папки с Low Integrity уровнем, то для записи в нужное место необходимо создать папку и указать Integrity Level.

Агент сбора при получении посылает ответ агенту мониторинга, закрывает соединение и записывает полученные данные в единую базу данных, при этом к каждой записи добавляется ID.

Детали реализации агента консолидации:

- Агент сбора был написан на примере асинхронного сокет сервера [14]
- Язык написания - C#
- Количество строк кода - 600

5.1.3 Модуль многотемной классификации

С учетом требований, сформулированных в постановке задачи, для многотемной классификации был выбран модуль, разработанный в лаборатории Технологий Программирования. Модуль многотемной классификации состоит из:

- компонент лексического анализа (парсер) – осуществляет разбор, выделение признаков и преобразование гипертекстовых документов во внутреннее представление;
- компонент вычисления меры сходства – определяет значения близости между документами (значения функции ядра) на основе выданного парсером представления и осуществляет кэширование этих значений;
- классификатор – строит дообучаемую модель классификации и на её основе осуществляет классификацию многотемных гипертекстовых документов.

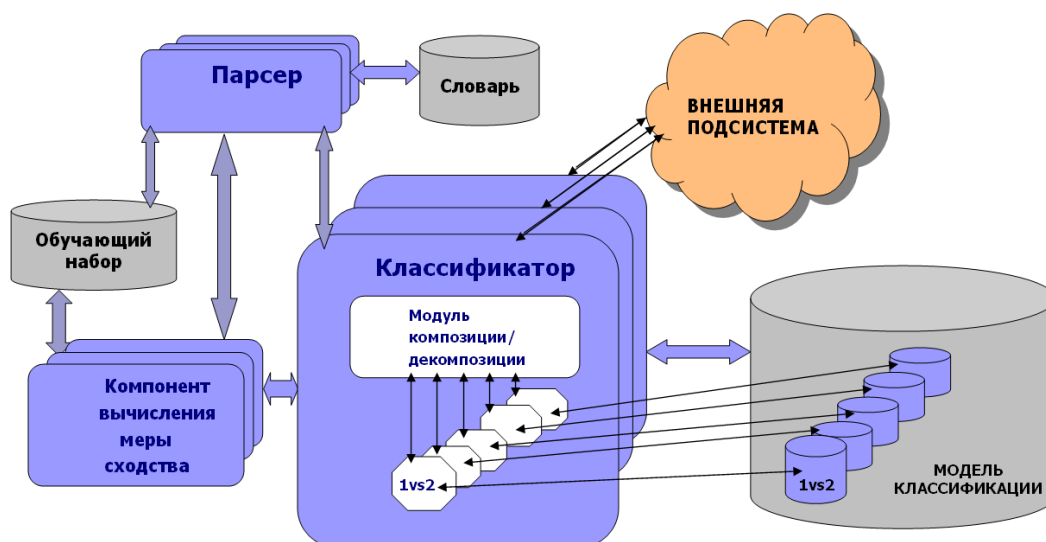


Рисунок 7 – Архитектура модуля классификации

Для интеграции системы сбора с модулем многотемной классификации был написан модуль, который позволяет:

- Производить парсинг html кода страницы.
- Обучать классификатор, подавая ему на вход тренировочный набор, содержащий веб-страницу и тематику, к которой данный документ принадлежит.
- Выгружать веб-страницы из заданной базы данных и подавать их на вход классификатору. На выходе создается .csv файл с весами тематик для каждого документа.

Модуль написан на языке программирования Python, объем - 200 строк

5.2 Экспериментальные исследования

В ходе обзора существующих решений было выдвинуто требование масштабируемости агента консолидации.

Так как разработка велась на языке программирования C# в среде разработки Visual Studio, то для проведения нагрузочного тестирования было решено воспользоваться встроенными средствами Visual Studio [18].

Ход тестирования:

- Для нагрузочного тестирования были созданы Unit тесты [17], содержащие код компоненты агента мониторинга, который взаимодействует с агентом консолидации
- Visual Studio позволяет выбирать такие параметры сценария тестирования как: начальное количество клиентов, максимально количество клиентов, и скорость увеличения количества клиентов.
- Далее, задавая параметры сценариев, было проведено исследование зависимости нагрузки ЦП и оперативной памяти от количества подключенных клиентов.

5.2.1 Исследование зависимости нагрузки ЦП от количества подключенных клиентов

Сценарий тестирования:

- Начальное количество клиентов - 1. Каждый подключенный клиент с заданной частотой посылает агенту консолидации просмотренные пользователем веб-страниц.
- С установленной скоростью увеличивается количество подключенных агентов мониторинга.

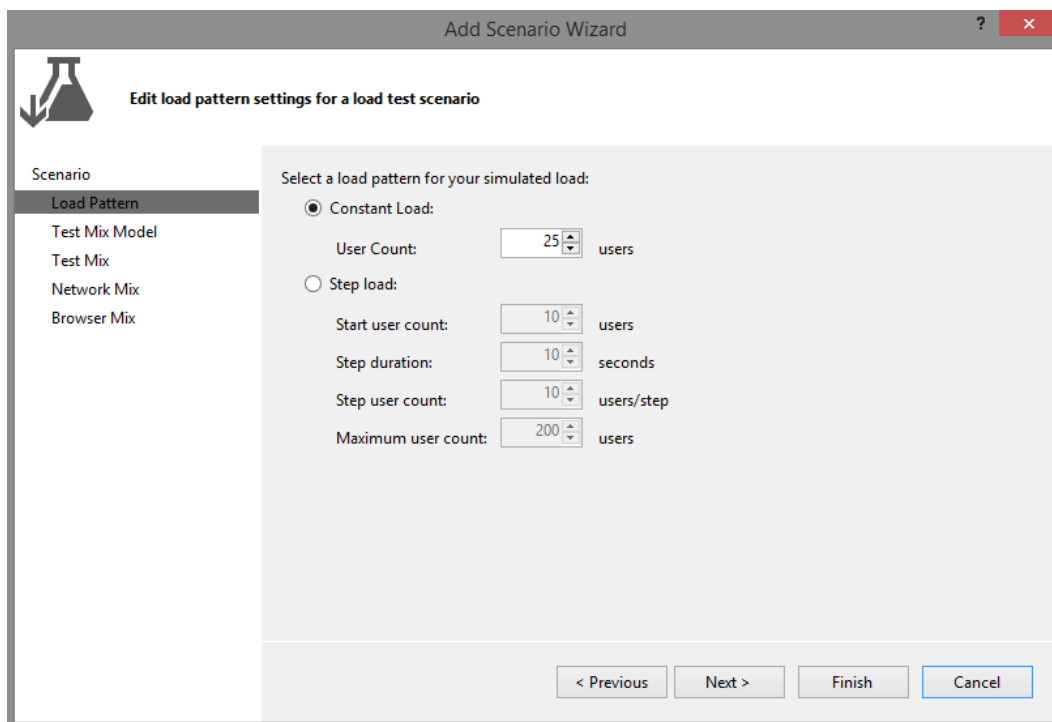


Рисунок 8 – Параметры сценария тестирования

- С помощью средств Visual Studio отрисовывается график, отражающий зависимость нагрузки ЦП от количества подключенных клиентов.

Результаты тестирования проиллюстрированы ниже

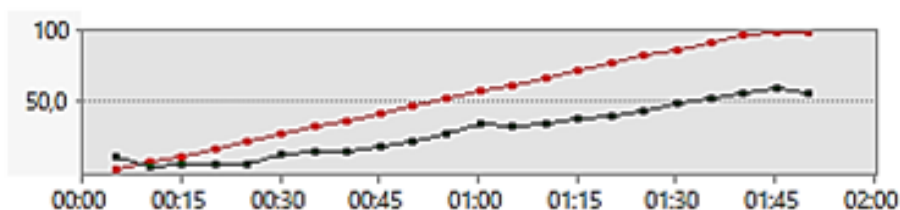


Рисунок 9 – Зависимость нагрузки ЦП от количества клиентов

5.2.2 Исследование использования оперативной памяти от количества клиентов

Сценарий тестирования аналогичен сценарию из предыдущего пункта. Результаты тестирования приведены ниже:

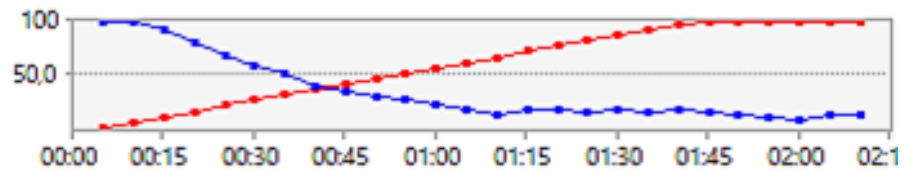


Рисунок 10 – Зависимость использования оперативной памяти от количества клиентов

5.3 Результаты тестирования

Зависимость нагрузки центрального процессора от количества подключенных клиентов линейна, что говорит о масштабируемости разработанного средства.

Увеличение использования оперативной памяти при увеличении количества подключенных клиентов линейно, что также говорит о масштабируемости разработанного средства.

6 Заключение

В настоящей выпускной квалификационной работе были проведены исследования методов многотемной классификации веб-страниц. В ходе работ обозревали современные индустриальные системы, которые решают ряд задач, требующих сбор и классификации веб-страниц, с которыми работали пользователи. Были рассмотрены различные подходы к классификации текстовых данных: методов шаблонов, метод цифровых, методы на основе машинного обучения. Было принято решение, что требованиям поставленной задачи удовлетворяет только методы, основанные на машинном обучении.

Был проведен анализ основных подходов к решению задачи многотемной классификации: "оптимизационный" подход, подход на основе декомпозиции в набор независимых бинарных проблем, подход на основе ранжирования, выявлены недостатки, не позволяющие применять их для решения поставленной задачи. Также был рассмотрен и выбран метод многотемной классификации на основе подхода попарных сравнений, учитывающий недостатки традиционных подходов.

Была осуществлена программная реализация прототипа системы сбора и многотемной классификации веб-страниц, состоящей из:

- Агента мониторинга.
- Агента консолидации.
- Модуля многотемной классификации.

С помощью средств модульного и нагрузочного тестирования, предоставленного средой Microsoft Visual Studio были реализованы сценарии нагрузочного тестирования и измерены следующие параметры функционирования разработанного прототипа системы:

- Зависимость нагрузки центрального процессора от количества подключенных клиентов.
- Зависимость количества свободной оперативной памяти от количества подключенных клиентов.

Результаты тестирования показали линейный рост нагрузки центрального процесса при увеличении количества подключенных клиентов, что говорит о масштабируемости разработанного прототипа.

Также результаты тестирования показали линейное убывание свободной оперативной памяти при увеличении количества клиентов, что говорит о масштабируемости.

Реализация сбора данных с помощью ВНО позволила сохранять данные о просмотренных пользователем веб-страницах в локальную файловую систему в директорию с ограниченным правом доступа, что обеспечивает защищенность (обычный пользователь не имеет доступа к собранным данным) разработанного прототипа.

Список литературы

1. Аналитический Центр InfoWatch, Безопасность информации в корпоративных информационных системах. Внутренние угрозы (<http://www.infowatch.ru/analytics/reports/4609>)
2. *Component overview (Content Classification 8.8.0)*
3. *Electronic discovery*. (wikipedia.org/wiki/Electronic_discovery)
4. *Предиктивное обучение* (<http://www.symantec.com/ru/ru/predictive-coding/>)
5. *Infowatch БКФ* (<http://www.infowatch.ru/technologies>)
6. *text122. Component overview (IBM Watson Content Analytics 3.5.0)* (http://www-01.ibm.com/support/knowledgecenter/SS5RWK_3.5.0/com.ibm.discovery.es.nav.doc/i1ysaov)
7. *Официальная документация IBM в Интернет* (http://www-01.ibm.com/support/knowledgecenter/SS5RWK_3.5.0/com.ibm.discovery.es.nav.doc/i1ysaov)
8. *Content Analytics. Официальный сайт OpenText в Интернет* (<http://www.opentext.com/what-we-do/products/discovery/content-analytics>)
9. *EMC Kazeon File Intelligence* (<http://www.emc.com/content-management/emc-kazeon-file-intelligence.htm>)
10. *Работа AdaBoost алгоритма* (<https://ru.wikipedia.org/wiki/AdaBoost>)
11. *Using the Taxonomy Proposer to discover new categories* (<http://www-01.ibm.com/support/knowledgecenter/>)
12. *Symantec Machine Learning* (<http://eval.symantec.com/mktginfo/enterprise/>)
13. *Boutell M. R., Luo J., Shen X., Brown C.M. Learning multi-label scene classification* Pattern Recognition. 2004. №37. pp. 1757-1771.
14. *Асинхронный сокет сервер* ([https://msdn.microsoft.com/ru-ru/library/fx6588te\(v=vs.110\).aspx](https://msdn.microsoft.com/ru-ru/library/fx6588te(v=vs.110).aspx))
15. *Document Object Model* https://ru.wikipedia.org/wiki/Document_Object_Model
16. *Утилита icacls* (<https://msdn.microsoft.com/en-us/library/bb625965.aspx>)
17. *Create and run unit tests.* (<https://www.visualstudio.com/en-us/get-started/code/create-and-run-unit-tests-vs>)
18. *Walkthrough: Creating and Running a Load Test Containing Unit Tests* ([https://msdn.microsoft.com/en-us/library/vstudio/ff355993\(v=vs.110\).aspx](https://msdn.microsoft.com/en-us/library/vstudio/ff355993(v=vs.110).aspx))
19. *Таненбаум Э., Ван Стеен М. Распределенные системы. Принципы и парадигмы*
20. *Schapire R. E., Singer Y. Boostexter A boosting-based system for text categorization*
21. *Comite F. D., Gilleron R., Tommasi M. Learning multi-label alternating decision tree from texts and data* Machine Learning and Data Mining in Pattern Recognition, MLDM 2003 Proceedings, Lecture Notes in Computer Science 2734. Berlin, 2003. pp. 35–49s

22. Zhang M.-L., Zhou Z.-H. *A k-nearest neighbor based algorithm for multi-label classification* Proceedings of the 1st IEEE International Conference on Granular Computing (GrC'05). Beijing, China, 2005. pp. 718-721
23. C. Crammer, Y. Singer. *A family of additive online algorithms for category ranking* Machine Learning Research. №3. 2003. pp. 1025–1058
24. Crammer C., Singer Y. *A new family of online algorithms for category ranking* Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval. Tampere, Finland, 2002. pp. 151 – 158
25. Minh Duc Cao, Xiaoying Gao. *Combining Content and Citation for Scientific Document Classification* AI2005, LNAI 3809, 2005. pp. 143-152
26. P.V. Rao and L.L. Kupper. Ties in paired-comparison experiments A generalization of the Bradley–Terry model, Amer. Statist. Assoc, 62, 1967. pp. 194–204.
27. Глазкова В.В. Исследование и разработка методов построения программных средств классификации многотемных гипертекстовых документов
28. Zhang M.-L., Zhou Z.-H. *A k-nearest neighbor based algorithm for multi-label classification* Proceedings of the 1st IEEE International Conference on Granular Computing (GrC'05). Beijing, China, 2005. pp. 718-721
29. C. Crammer, Y. Singer. *A family of additive online algorithms for category ranking* Machine Learning Research. №3. 2003. pp. 1025–1058
30. Minh Duc Cao, Xiaoying Gao. *Combining Content and Citation for Scientific Document Classification* AI2005, LNAI 3809, 2005. pp. 143-152.