



Рис. 2. Пример схемы «звезды» [4]



Рис. 3. Пример схемы «снежинка» [4]

Достоинствами ROLAP-систем являются:

- возможность использования ROLAP с хранилищами данных и различными OLTP-системами;
- возможность манипулирования большими объемами данных;
- безопасность и администрирование обеспечивается реляционными СУБД.

Недостатки:

- получение агрегатов и листовых данных происходит медленнее, чем в MOLAP и HOLAP;
- функциональность систем ограничивается возможностями SQL, так как аналитические запросы пользователя транслируются в SQL-операторы выборки;
- сложно пересчитывать агрегированные значения при изменениях начальных данных;
- сложно поддерживать таблицы агрегатов.

Гибридный OLAP

В гибридных OLAP сочетаются черты ROLAP и MOLAP. В моделях HOLAP используются преимущества и минимизируются недостатки обеих архитектур. В HOLAP-системах структура куба и предварительно обработанные агрегаты хранятся

в многомерной базе данных. Это позволяет обеспечить быстрое извлечение агрегатов из структур MOLAP. Значения нижнего уровня иерархии в HOLAP остаются в реляционной витрине данных, которая служит источником данных для куба.

К достоинствам подхода можно отнести комбинирование технологии ROLAP для разреженных данных и MOLAP для плотных областей, а к недостаткам – необходимость поддерживания MOLAP и ROLAP.

Заключение

Можно сделать вывод, что выбор модели OLAP зависит от требований к скорости загрузки данных, требуемого дискового пространства и т.д. После выбора модели и настройки OLAP корпоративная информация предприятия будет представлена аналитику в необходимой форме, что позволит пересмотреть стратегию дальнейшего развития предприятия.

Литература

1. Способы аналитической обработки данных для поддержки принятия решений [Электронный ресурс]. – Режим доступа: <http://infovisor.ivanovo.ru/press/paper04.html>, свободный.
2. Кузьмин А. Н. Методы и модели обработки информации в хранилищах данных /А. Н. Кузьмин //Математическое моделирование, численные методы и комплексы программ. 2006. – С.193.
3. OLAP [Электронный ресурс]. Режим доступа: <http://ru.wikipedia.org/wiki/OLAP> свободный
4. Оперативная аналитическая обработка данных: концепции и технологии [Электронный ресурс]. – Режим доступа: http://www.olap.ru/basic/olap_and_ida.asp, свободный
5. OLAP и многомерные базы данных [Электронный ресурс]. – Режим доступа: <http://www.rae.ru/monographs/141-4638>, свободный

АВТОМАТИЧЕСКИЙ АНАЛИЗ ДОКУМЕНТОВ НА ЕСТЕСТВЕННЫХ ЯЗЫКАХ

Губин М.Ю.

Томский политехнический университет
634050, Россия, г. Томск, пр-т Ленина, 30

E-mail: gubin.m.u@gmail.com

Введение

Одной из фундаментальных проблем машинной обработки текстов является то, что естественные человеческие языки обладают большой выразительностью и сложностью, существенное влияние на смысл текста в них оказывает контекст и эмоциональная составляющая. Понимание есте-

ственного языка включает куда больше, чем разбор предложений на индивидуальные части речи и поиск значений слов в словаре. Оно базируется на обширном фоновом знании о предмете, идиомах, используемых в этой области, а также на способности применять общее контекстуальное знание для понимания недомолвок и неясностей, присущих

щих естественной человеческой речи. Поэтому системы, использующие натуральные языки с гибкостью и общностью, характерными для человеческой речи, лежат за пределами существующих методологий [1]. Однако, для определённых условий (когда документ имеет достаточно строгую грамматическую структуру, а следовательно – содержит достаточно информативную формальную составляющую) данная задача решаема с достаточно высоким качеством распознавания смысла [2]. В этой статье будут описаны условия, выполнение которых необходимо для успешного распознавания, и предлагаемый алгоритм, работающий с достаточно высокой степенью точности в описанном частном случае.

Постановка задачи

Данный алгоритм решает задачу создания метаописаний документов для последующего семантического поиска по ним на данном множестве документов D_i , относящихся к одной предметной области. Под документом D_i в рамках данного исследования будем понимать фрагмент текста на естественном языке.

Для реализации семантического поиска по документам, необходимо создать достаточно полные семантические метаописания документов T_i .

Семантическое метаописание документа строится согласно онтологии предметной области O , представляющей собой набор понятий C_i , связанных между собой отношениями R_i . Также в онтологию предметной области входят экземпляры объектов E_i . Понятия, отношения и экземпляры имеют одну или более текстовых меток T_i . Текстовая метка T_i элемента онтологии – слово либо словосочетание естественного русского языка, соответствующее некоторому элементу онтологии.

Для построения базового семантического метаописания на основе текста документа для каждого его предложения L_i формируется семантическая сеть, представляющая собой граф, состоящий из множества вершин W_i и соединяющих их рёбер L_i . Элементарная сеть представляет результат синтаксического анализа и дополнительных семантических трансформаций дерева синтаксических зависимостей между словами в отдельном предложении. Вершинами W_i семантической сети являются сущности, встречающиеся в предложении, а рёбра L_i представляют собой семантические отношения между сущностями. Семантические сети предполагается получать из результатов синтаксического разбора текстов на естественных языках. Задача синтаксического разбора текстов на данный момент в различной степени решена для русского [6, 7] и английского [3, 4, 5] языка. Также существуют работы по синтаксическому разбору текстов на французском, норвежском, корейском и греческом [4], а также испанском и японском [4, 5] языках. В данной работе рассматривается частный случай с русским языком.

Программный интерфейс большинства существующих семантических анализаторов позволяет получить для каждой сущности набор направленных связей, исходящих от нее к другим сущностям. Направление связи обычно соответствует направлению синтаксического подчинения (для равноправных однородных членов предложения пара одинаковых направленных связей идет в обе стороны). Семантические сети, соответствующие описанным выше критериям, могут быть использованы в разрабатываемом алгоритме с незначительными преобразованиями.

Семантическое метаописание – это набор извлечённых из предложений документа RDF-триплетов T_i , представляющих собой кортежи вида $\langle S_i, P_i, O_i \rangle$, где S_i включен в объединение C_i и E_i , P_i включен в R_i , а O_i включен в объединение C_i и E_i .

Также для ускорения актуализации метаданных алгоритмом генерируются частотные характеристики слов в документе – TF и IDF терминов [8].

Алгоритм формирования метаданных отдельного документа

На вход алгоритма поступает исходный текст файла, а также набор текстовых меток элементов онтологии.

Шаги алгоритма:

1. Производится семантический анализ текста. Выходом этого шага является программная структура, содержащая всю требуемую информацию о тексте – слова с номером их начальных символов, смысловые связи между словами, обнаруженные и преобразованные в RDF триплеты (части предложений, соответствующие одному из описанных выше фреймов). Эта программная структура приводится к семантической сети, пригодной для обработки алгоритмом.

2. Подсчитывается количество вхождений слов в текст. При этом не учитываются так называемые «стоп-слова». Стоп-словами являются предлоги, союзы и частицы. Остальные слова нормализуются и количество вхождений подсчитывается именно для нормы слова.

3. Составляется ранговое распределение слов в документе. Слова с одинаковым количеством вхождений объединяются в классы, которые затем нумеруются в порядке убывания количества вхождений слов-членов класса в тексте, начиная с 1 [8]

4. Производится поиск класса, слова в котором являются значимыми для текста, с наибольшим номером. Все классы, идущие после него, отсеиваются и в дальнейшей работе алгоритма не участвуют. [8]

5. Выставляется первичное значение «веса» слов в документе. Оно равняется N_{max}/N_i , где N_{max} – количество вхождений слов первого ранга, а N_i – количество вхождений слова t_i [8].

6. Производятся корректировки значений весов для упорядоченных пар слов, входящих в одни и те же триплеты либо предложения.

7. Из множества выделенных из текста RDF-триплетов выбираются:

Триплеты, каждая из позиций которых (субъект, предикат и объект) заняты в естественноязыковом представлении вхождением метки (соответственно, субъект и объект – метками понятия либо экземпляра, а предикат – меткой свойства).

Триплеты, одна из позиций которых занята вхождением ключевого слова, а две других – вхождением метки, так называемые триплеты «кандидаты».

Выход алгоритма – метаописание документа, в которое входит набор записей вида $\langle E_i, S_i \rangle$, где E_i – идентификатор элемента онтологии (так называемый URI – Universal Resource Identifier), а S_i – индекс значимости этого элемента для документа. При этом S_i имеет вид $S_i = \langle S_{iT}, S_{iDF}, S_{iC} \rangle$, где S_{iT} – коэффициент значимости элемента с точки зрения документа (модифицированный коэффициент TF), S_{iDF} – коэффициент значимости элемента с точки зрения набора документов (коэффициент IDF), S_{iC} – итоговое значение коэффициента значимости термина. В метаописание также входят все обнаруженные в тексте триплеты, все позиции которых заняты вхождениями меток элементов онтологии.

Кроме того, по завершении работы алгоритм генерирует набор вспомогательных записей, уменьшающих время возможной последующей повторной обработки документа.

Результаты работы алгоритма – семантические метаописания, которые позволяют реализовать семантический поиск и семантическую навигацию по обработанному множеству текстов. Качество распознавания находится на уровне примерно 60 % от распознавания человеком, в зависимости

от полноты онтологии предметной области и глубины анализа текста.

Литература

1. Люгер Д.Ф. Искусственный интеллект: стратегии и методы решения сложных проблем. 4-е издание. – М.: Вильямс, 2003. – 864 с.

2. Хорошилов А.А. Белоногов Г.Г. Калинин Ю.П. Компьютерная лингвистика и перспективные информационные технологии: теория и практика. // НТИ. Сер. 2. Информ. процессы и системы / ВИНИТИ. – 2004. – N 8. – С.30-43.

3. Poon H., Domingos P. Unsupervised semantic parsing. ACL Anthology. A Digital Archive of Research Papers in Computational Linguistics / [Электронный ресурс]. Режим доступа: www.aclweb.org/anthology/D/D09/D09-1001.pdf, свободный (дата обращения: 02.10.2010).

4. Deep linguistic processing with hpsg. [Электронный ресурс]. – Режим доступа: <http://www.delph-in.net>, свободный (дата обращения: 02.10.2010).

5. Сайт лаборатории speech technology копорации microsoft. [Электронный ресурс]. – Режим доступа: <http://research.microsoft.com/en-us/groups/srg/default.aspx>, свободный (дата обращения: 02.10.2010).

6. Сайт рабочей группы «Автоматическая обработка текстов». [Электронный ресурс] / Режим доступа: <http://aot.ru/>, свободный (дата обращения: 02.10.2010).

7. Сайт компании RCO [Электронный ресурс] / Режим доступа: <http://www.rco.ru>, свободный (дата обращения: 02.10.2010).

8. Thomas Roelleke , Jun Wang, TF-IDF uncovered: a study of theories and probabilities // Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval, July 20-24, 2008, Singapore, Singapore

ДЕТЕКТИРОВАНИЕ ЛИЦ С ПОМОЩЬЮ СВЕРТОЧНОЙ НЕЙРОННОЙ СЕТИ

Калиновский И.А.

Научный руководитель: д.т.н., профессор Спицын В.Г.

Томский политехнический университет

634050, Россия, г. Томск, пр-т Ленина, 30

E-mail: kua_21@mail.ru

Введение

Задача выделения лиц на изображениях или в видеопотоке является одной из классических в области обработки изображений и компьютерного зрения. Потребность в таких алгоритмах обусловлена необходимостью в автоматизации различных процессов, связанных с обеспечением безопасности, учетом и контролем доступа.

Исследования в этой области ведутся уже более 15 лет. Предложено множество алгоритмов, начиная от простых статистических моделей и заканчивая методами машинного обучения и 3D моделированием лица. Однако нельзя ска-

зать, что эта задача решена полностью, т.к. не разработан алгоритм, позволяющий надежно детектировать лица при любом распределении освещенности, различных поворотах, наклонах и масштабах лица, перекрытии лица объектами, а также при низком разрешении изображения и наличии шумов. В таблице 1 перечислены некоторые существующие методы выделения лиц, а так же их преимущества и недостатки [1].