

О разработке рекомендательной системы, предлагающей книги по предпочтениям пользователя

Волкова Л.Л., НИУ ВШЭ, МИЭМ; МГТУ им. Н.Э. Баумана

liliyavolkova@itas.miem.edu.ru

Токарева М.М., НИУ ВШЭ, МИЭМ

rit1336@yandex.ru

Ланко А.А., МГТУ им. Н.Э. Баумана

aleksanchezz@mail.ru

Аннотация

Данная работа посвящена разработке рекомендательной системы, принимающей в расчёт предпочтения пользователя в книгах: параметрические, жанровые, основанные на прочитанных произведениях. Проведён обзор применения стратегий рекомендательных систем в данной задаче. Описаны ключевые положения для создания рекомендательной системы, предлагающей читателю книги на основе введённых им предпочтений, а также ключевые методы, которые следует положить в основу подсистему извлечения данных из текстов произведений. Рассмотрена перспектива расширения системы данными, извлечёнными из рецензий на книги.

Введение

С распространением систем каталогизации и поиска совершение выбора стало за широтой его и проще, и сложнее одновременно. В связи с этим популярность получили методы автоматизации выбора и принятия решения и, наконец, рекомендательные системы. Математические методы позволили механизировать оценку альтернатив (к примеру, проектов аэропортов, расположенных в различных местах и имеющих различные критерии), исключив из системы принятия решения человеческий фактор, заменённый алгоритмом, опирающимся на занесённые в цифровой вид знания экспертов. Затем методы принятия решения и рекомендательные системы повернулись лицом к более широкой аудитории, которой они представляют свои гипотезы на различных сайтах и сервисах по подбору товаров, услуг, впечатлений. Существуют порталы, рекомендуемые кинофильмы и путешествия на основании тем или иным способом интерпретированных вкусов пользователей, и одно из перспективных направлений расширения области применения таких реко-

mendательных систем – книжный сектор. Помимо очевидных факторов новизны, стоимости, заявленного жанра и, что реже, тематики книги, отдельный интерес представляет оценка стиля книги – стиля, которым она написана. Спрос на услугу выбора книг именно по предпочтениям пользователя такой системы будет иметь место, и добиться приемлемого качества – разрешимая задача, поскольку существуют работы по стилометрии текстов, а книги имеют, в отличие от кратких новостных статей, достаточный объём, чтобы составить их портреты, которые затем можно сравнивать.

1 Постановка задачи

Книги – одно из самых важных богатств человечества, поскольку они передают опыт, устремления, достижения мысли прошлых поколений. Книги могут увлекать, пусть даже и научные, могут вдохновлять, могут быть друзьями и учителями. В цифровую эру книги не утратили своей значимости, но изменилась форма книжных клубов. Теперь мнениями о прочитанном делятся в сети на специализированных форумах, а также пишут рецензии и отзывы в книжных отделах электронных магазинов и в электронных отделах книжных настоящих. Задача разработки рекомендательной системы, учитывающей не рейтинг (или не только его), но жанровые и иные предпочтения пользователя, является актуальной.

2 Разработка рекомендательной системы

В данном разделе будут рассмотрены стратегии рекомендательных систем. На основании анализа будет предложен подход, который может служить основой для РС, и будет изучен ряд прикладных задач, которые нужно решить в процессе разработки.

Многие рекомендательные системы (РС) используют подход, основанный на обобщённых рейтингах, но при этом снижается точность соответствия поисковой выдачи пользовательскому запросу (хотя такая практика широко распространена). О вкусах не спорят, поэтому интерес представляет разработка именно параметрического поискового механизма, который позволит искать пусть не наиболее популярные, но тематические произведения по предпочтениям пользователя, а также в перспективе подбирать похожие объекты в базе данных по метрикам схожести, применяющимся к таксономической модели.

2.1 Стратегии рекомендательных систем

РС, позволяющие учесть различные аспекты и детали пользовательских предпочтений, можно классифицировать по стратегиям следующим образом.

1) Системы с коллаборативной фильтрацией группируют пользователей со схожими предпочтениями для дальнейшей рекомендации. Подход заключается в выявлении «единомышленников» и в составлении рекомендаций с учётом их мнений, т.е. по похожим профилям (сюда часто включается история выбора пользователей и их оценок). Существуют два подхода: кластеризация по пользователям и по предпочтениям [Resnick, Varian, 1997]. Недостатком таких систем является т.н. проблема холодного старта, когда данные о новом пользователе ещё не накоплены.

2) Основанные на контенте системы используют метрику вкусов пользователя, которая извлекается из предложенного ему при регистрации опросника, а в дальнейшем пополняется его фактическим выбором [Candillier et al., 2009]. В зависимости от встроеного алгоритма машинного обучения такие системы могут быть чувствительны к обучению. Как отмечается в обзоре [Сергеев, 2008], пользователю не стоит оценивать подряд все книги из школьной программы, иначе система может сосредоточиться в своих рекомендациях на детской литературе и русской классике. Также рекомендуется не ограничиваться бестселлерами и добавить в список предпочтений некоторую специфику, дабы уточнить и разнообразить выдачу рекомендательной системы [Сергеев, 2008].

3) Прецедентные системы (частный случай контентных) борются с проблемой изменения предпочтений пользователя во времени

[Smyth, 2007; Bridge et al., 2006]: выделяются типажии пользователей либо сценарии поездок, описывающие данный случай (от англ. *case*). Так, можно выделить литературу для детей, для субкультур эмо (Г. Лавкрафт, С. Кинг, Ф. Кафка), ролевиков (жанры стимпанк, фэнтези и др.), инженеров («Хроники лаборатории», «Понедельник начинается в субботу»), научных сотрудников (Д. Лодж – «Академический обмен», Т. Пратчетт – «Наука плоского мира»). Впрочем, в случае книг отличия по профессиональному признаку либо субкультурам будут похожи на результаты системы с коллаборативной фильтрацией: так, аудитория порталов habrahabr.ru и geektimes.ru такова, что предпочтения, которые отследит коллаборативная фильтрация, совпадут с прецедентами.

4) Гибридные системы [Adomavicius, Tuzhilin, 2005; Falk, 2017].

В случае рекомендательных систем, ориентированных на книги, как и в других рекомендательных системах, возникает проблема холодного старта [Poirier et al., 2010]. Данные о новом пользователе следует собрать. Частное решение возможно, если доступны данные социальных сетей (так, профиль социальной сети Вконтакте включает графу «любимые книги»). Если у пользователя заполнен перечень прочитанных либо любимых книг, можно автоматически извлечь такой перечень соответствующим модулем взаимодействия с социальной сетью, а также, если оценка не проставлена, пользователю можно сразу предложить оценить те самые книги.

Одним из частных решений проблемы холодного старта является определение некоторого набора книжных «маршрутов» для типажа предпочтений (экскурсия по классической фантастике, по «попаданцам», по «романам взросления», по научно-популярным книгам о физике). Такие фиксированные выборки следует разрабатывать вручную либо выявлять на материале кластеризации пользователей существующего книжного сервиса.

Рекомендательную систему можно разработать как гибридную: анализировать контент и использовать прецеденты. Последние возможно моделировать набором тегов либо группами тегов (что подводит к переходу к нечёткой модели). Если реализовать открытое множество меток, присваиваемых объектам, это может повлечь появление шума, поэтому желательна модерация меток (возможно, на началах краудсорсинга). Если предлагать

пользователю выбрать интересующую его тематику из облака тегов, прошедших модерацию (что решает проблему холодного старта, но не даёт исчерпывающего представления о предпочтениях пользователя), задача сводится к отбору книг с заданными метками.

Одним из направлений перспективного развития можно считать автоматическое извлечение меток на основе аннотаций и рецензий, собранных на существующих ресурсах (в том числе на сайтах электронных книжных магазинов), по аналогии с [Pronoza et al., 2014; Lande et al., 2014; Масленникова, Ягунова, 2016]. Следующим шагом может стать разработка онтологии предметной области либо легковесной онтологии [Токарева и др., 2016], что позволит определять правила схожести книг на этом основании.

2.2 Анализ книг и рецензий

В части выборки данных следует рассмотреть отдельно два класса: полные тексты произведений и рецензии на книги.

Полные тексты книг, в том числе классических, доступны в электронных библиотеках. Каждая книга может быть подвергнута анализу с целью выделения тех маркеров или параметров, которые могут быть использованы для сужения поискового запроса с последующей выдачей рекомендации. Так, используя инструменты анализа текстов на естественном языке [Большакова и др., 2011], можно выделять сущности и факты из произведения (страна и/или локация, в которой происходит действие, ключевых персонажей и др. [Alexandrov et al., 2005; Клышинский, Кочеткова, 2014; Шелманов, 2012; Starostin et al., 2016]) и взаимосвязи между ними. Отношения между сущностями могут храниться в виде набора правил, графа, семантической сети, онтологии или легковесной онтологии [Cimiano, 2006; Giunchiglia et al., 2006; Токарева и др., 2016; Бородин, Строганов, 2016]. Считается, что «Именно текст... является релевантной единицей стилистических исследований» (М.А.К. Хэллидей). И при наличии корпуса книг можно натренировать методы автоматического извлечения данных из книг. Далее по результатам опроса фокус-группы (на материале апробации рекомендательной системы) и гипотетического успеха сервиса анализ вновь издающихся книг может проводиться в закрытом режиме на стороне издательства, если оно сочтёт целесообразным после выхода книги либо через некоторый

срок после него занести профиль книги в рекомендательную систему. Положительная сторона такого подхода заключается в рекомендации книг, не вошедших в верхние строки рейтинга, тем читателям, которым такие книги будут интересны по описанию их содержания [Gelbukh et al., 2006]: мира, типа сюжета и др.

Отдельный интерес представляет стилометрия [Волкова, 2015], поскольку оценка стиля и жанра текста и стиля автора тоже является перспективным направлением. В частности, оценке подлежат сложность текста [Мизернов, Гращенко, 2015], самоповторы автора [Гращенко, Романишин, 2015], богатство используемого словаря [Stamatatos et al., 2000], особенности лексики [Фоменко, Фоменко, 1996; Иомдин и др., 2013; Волкова, 2015], показатель употребления синонимов и антонимов [Ефремова и др., 2010], частеречный спектр текстов [Браславский, 2003], нарушения проективности [Гладкий, 2007], особенности строения фраз [Wang, Zhang, 2008]. Если учитывать такие параметры, возможно дать частное решение задаче поиска похожих книг: рекомендации тогда будут основываться на схожести жанра и почерка автора.

Тексты рецензий, если их выгрузить с профильных сайтов, могут быть подвергнуты процедуре выделения аспектов и заполнения фрейма о данном произведении, по аналогии с [Pronoza et al., 2014; Lande et al., 2014; Масленникова, Ягунова, 2016] (ручное выделение признаков и разработка метода для извлечения данных из текстов отзывов). Это весьма перспективный подход, тем более интересный, что рецензии доступны и на только что выпущенные книги, что позволит быстро извлечь характеристики произведения в случае полноты анализируемого обзора или обзоров.

Оба подхода могут использоваться по отдельности или комбинироваться в зависимости от того, насколько давно вышла из печати книга, от пользовательской настройки поиска.

Дополнительно данные могут извлекаться из интернета: существуют различные порталы, предоставляющую информацию о книгах. Такие сайты можно разделить на магазины и на системы с отзывами, в том числе профильные социальные сети. Ввиду отсутствия гарантии полноты книжной графы профиля в социальной сети, модуль выгрузки данных из последней может быть разве что дополнительным в системе: полные данные позволяют

решить проблему холодного старта, их отсутствие акцентирует важность решения проблемы холодного старта другими средствами.

2.3 Параметры книг

Следует по возможности выделять следующие свойства книг:

- автор и название;
- жанр;
- цикл книг (если есть);
- время написания;
- страна происхождения автора;
- страна (в случае, если события привязаны к реальной географии) либо мир, где происходит действие;
- время действия (книги также могут быть сгруппированы по описываемым эпохам);
- популярность книги;
- стиль автора.

Кроме того, существует мнение, что следует разделять литературу на «женскую» и «мужскую», однако если вводить такое членение, то делать это следует только факультативно, с возможностью отключения критерия пользователем в интерфейсе.

Многопараметрический поиск потребует ранжирования критериев. Одним из способов ввести ранжирование в сам запрос является проставление важности критерия (блокирующий, не важен, важен), как это сделано в интерфейсе рекомендательной системы Triplantica [Triplantica, 2016]. Полученная выдача может быть ранжирована при помощи весовой функции, которая требует параметризации на стадии разработки системы. Результирующую выборку можно выдать целиком в порядке убывания релевантности. В части визуализации можно использовать средства кластеризации, а также выводить рядом с названием каждой книги перечень меток, которые поставлены ей в соответствие.

При переходе от многопараметрического поиска к поиску похожих книг потребуется определить схожесть. С одной стороны, есть возможность создания семантической сети или онтологии (возможно, легковесной), включающей определённые экспертом отношения между книгами и их параметрами, либо набора правил над классификацией книг [Токарева и др., 2016]. С другой стороны, пользователю рекомендательной системы можно рекомендовать уточнить, хочет ли он искать книги, схожие по одной или несколь-

ким выбранным категориям, описанным выше, либо запросить сравнение по текстам книг. В последнем случае интерес представляет выделение фактов из текста произведения, чтобы создать машинное представление канвы сюжета [Бодрова, Бочаров, 2014] (которая может быть отнесена к некоторому жанру либо использована в качестве образца при поиске) и некоторые детали повествования, а также определение функционального стиля текста (художественный, публицистический, научный, смешение стилей), стиля автора [Волкова, Ланко, 2016].

2.4 К прототипу

Разнообразие настроек работы рекомендательной системы сводится воедино в метод гибкого параметрического поиска, допускающий поиск по образцу. Адаптивный поиск может основываться на комбинаторной оптимизации [Токарева и др., 2016], на деревьях принятия решений и иных методах классификации в случае реализации рекомендательной системы, основанной на контенте либо прецедентах. В частности, для решения проблемы холодного старта на начальном этапе можно основываться на наборе предварительно сформированных выборок произведений, которые можно извлечь с существующих статических ресурсов, посвящённым книгам. Отличительной особенностью описываемого проекта является поиск схожих произведений по образцу с учётом содержания, функционального и авторского стилей. Интерес представляет также лексический анализ произведений, что позволит сделать поиск более узконаправленным. Следует провести тестирование прототипа системы на фокус-группе. Предполагается, что широта возможностей поиска привлечёт разносторонний круг пользователей и удовлетворит вкусам искушённых читателей.

Заключение

В результате проведенной работы изучены проблемы, возникающие при создании рекомендательной системы, посвящённой книгам. Рассмотрены стратегии рекомендательных систем применительно к задаче рекомендации книг по пользовательским параметрическим предпочтениям. Предложены пути решения проблемы холодного старта. Приведены ключевые группы методов, которые следует встроить в подсистему извлечения дан-

ных из книг и рецензий, включая анализ языковых средств и особенностей построения текста. Выделены ключевые параметры книг, которые следует выделять и закладывать в основу разрабатываемой системы и, возможно, базы знаний о понятиях предметной области. Предложены основные положения гибкой рекомендательной системы, допускающей различные стратегии поиска книг.

Список литературы

- А. А. Бодрова, В. В. Бочаров. Извлечение фактов об отношениях между персонажами из художественных текстов. Статьи международной конференции Диалог-2014, публикуемые на сайте. URL: www.dialog-21.ru/digests/dialog2014/materials/pdf/BodrovaAABocharovVV.pdf (дата обращения: 31.06.2014).
- Е. И. Большакова, Э. С. Клыпинский, Д. В. Ландэ, А. А. Носков, О. В. Пескова, Е. В. Ягунова. 2011. *Автоматическая обработка текстов на естественном языке и компьютерная лингвистика. Учебное пособие*. М.: МИЭМ.
- Д. С. Бородин, Ю. В. Строганов. 2016. *К задаче составления запросов к базам данных на естественном языке*. Новые информационные технологии в автоматизированных системах: материалы девятнадцатого научно-практического семинара, с. 119–126. М.: ИПМ им. М.В. Келдыша.
- П. Браславский. 2003. *Опыт автоматической классификации текстов по стилям (на материале документов Internet)*. Русский язык в Интернете. Сб. статей. Казань, 2003. С. 6–15.
- Л. Л. Волкова. 2015. *К задаче определения функционального стиля документа на естественном языке*. Новые информационные технологии в автоматизированных системах: материалы восемнадцатого научно-практического семинара, с. 615–626. М.: ИПМ им. М.В. Келдыша.
- Л. Л. Волкова, А. А. Ланко. 2016. *О кластеризации текстов по функциональным стилям с использованием триграмм из частей речи слов*. Инновационные, информационные и коммуникационные технологии: сборник трудов XIII Международной научно-практической конференции. / под ред. С.У. Увайсов. С. 193–196. М.: Ассоциация выпускников и сотрудников ВВИА им. проф. Жуковского.
- А. В. Гладкий. 2007. *Синтаксические структуры естественного языка*. Изд. 2. М.: ЛКИ.
- Л. А. Гращенко, Г. В. Романишин. 2015. *Опыт автоматизированного анализа повторов в научных текстах*. Новые информационные технологии в автоматизированных системах: материалы восемнадцатого научно-практического семинара, с. 582–590. М.: ИПМ им. М.В. Келдыша.
- Н. Э. Ефремова, Е. И. Большакова, А. А. Носков, В. Ю. Антонов. 2010. *Терминологический анализ текста на основе лексико-синтаксических шаблонов*. Компьютерная лингвистика и интеллектуальные технологии: По материалам ежегодной Международной конференции «Диалог» (Бекасово, 26–30 мая 2010 г.). Вып. 9 (16). С. 124–129. М.: Изд-во РГГУ.
- Б. Л. Йомдин, Б. Л., Лопухина А. А., Панина М. Ф., Носырев Г. В., Вилл М. В., Зайдельман Л. Я., Матиссен-Рожкова В. И., Винокуров Ф. Г., Выборнова А. Н. 2013. *Mag vel mot: изменения в языке на материале бытовой терминологии*. Компьютерная лингвистика и интеллектуальные технологии: По материалам ежегодной Международной конференции «Диалог» (Бекасово, 29 мая – 2 июня 2013 г.). Вып. 12 (19). Т. 1: Основная программа конференции. С. 311–324. М.: Изд-во РГГУ.
- Э. С. Клыпинский, Н. А. Кочеткова. 2014. *Метод извлечения технических терминов с использованием меры странности*. Новые информационные технологии в автоматизированных системах: материалы семнадцатого научно-практического семинара. С. 365–370. М.: ИПМ им. М.В. Келдыша.
- А. Масленникова, Е. В. Ягунова. 2016. *Извлечение информации и мнений на материале русскоязычных и англоязычных рецензий о ведущих музеях мира. Методика и предварительные результаты*. Новые информационные технологии в автоматизированных системах: материалы девятнадцатого научно-практического семинара, с. 68–74. М.: ИПМ им. М.В. Келдыша.
- И. Ю. Мизернов, Л. А. Гращенко. 2015. *Анализ методов оценки сложности текста*. Новые информационные технологии в автоматизированных системах: материалы восемнадцатого научно-практического семинара, с. 572–581. М.: ИПМ им. М.В. Келдыша.
- Л. В. Найханова. 2005. *Основные аспекты построения онтологий верхнего уровня и предметной области*. Интернет-порталы: содержание и технологии, вып. 3. Редкол.: А.Н. Тихонов (пред.) и др.; ФГУ ГНИИ ИТТ "Информика", с. 452–479. М.: Просвещение.
- А. Сергеев. 2008. *Рекомендательная революция*. «Вокруг света», апрель 2008. URL: <http://www.vokrugsveta.ru/vs/article/6226/> (дата обращения: 1.2.2017).
- М. М. Токарева, Л. Л. Волкова, А. П. о. Абдуллаев. 2016. *О рекомендательной маршрутной системе, основанной на оценке предпочтений*

- пользователя. Новые информационные технологии в автоматизированных системах: материалы девятнадцатого научно-практического семинара, с. 75–80. М.:ИПМ им. М.В. Келдыша.
- V. П. Фоменко, Т. Г. Фоменко. 1996. *Авторский инвариант русских литературных текстов. Предисловие А.Т. Фоменко*. Фоменко А.Т. Новая хронология Греции: Античность в средневековье. Т. 2, с. 768–820. М.: Изд-во МГУ.
- A. О. Шелманов 2012. *Метод автоматического выделения многословных терминов из текстов научных публикаций*. Труды тринадцатой национальной конференции по искусственному интеллекту с международным участием КИИ-2012, т. 1, с. 268–274. Белгород: БГТУ.
- G. Adomavicius, A. Tuzhilin. 2005. *Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions*. IEEE Trans. on Knowl. and Data Eng., vol. 17, pp. 734–749, June 2005. IEEE Computer Society, Washington, USA.
- M. Alexandrov, A. Gelbukh, P. Rosso. 2005. *An Approach to Clustering Abstracts*. LNCS 3513, pp. 275–285. Springer, Berlin.
- D. Bridge, M. Goker, L. McGinty, B. Smyth. 2006. *Case-based recommender systems*. Knowledge Engineering Review, vol. 20 (3), pp. 315–320. Cambridge University Press, New York, USA.
- L. Candillier, K. Jack, F. Fessant, F. Meyer. 2009. *State-of-the-art recommender systems*. Collaborative and Social Information Retrieval and Access-Techniques for Improved User Modeling, pp. 1–22. IGI Global, Hershey, USA.
- P. Cimiano. 2006. *Ontology Learning and Population from Text: Algorithms, Evaluation and Applications*. Springer, Berlin.
- K. Falk. 2017. Practical Recommender Systems. URL: <https://www.manning.com/books/practical-recommender-systems> (date of retrieval: 28.03.2017).
- A. Gelbukh, G. Sidorov, A. Guzmán-Arenas. 1999. *A method of describing document contents through topic selection*. Proc. of the String Processing and Information Retrieval Symposium and International Workshop on Groupware, pp. 73–80. IEEE, Los Alamitos.
- F. Giunchiglia, M. Marchese, I. Zaihrayeu. 2006. *Encoding classifications into lightweight ontologies*. Technical Report DIT-06-016, University of Trento, Italy, March 2006.
- D. Lande, A. Snarskii, E. Yagunova, E. Pronoza, S. Volskaya. 2014. *Network of Natural Terms Hierarchy as a Lightweight Ontology*. Thirteenth Mexican International Conference on Artificial Intelligence MICAI 2014, Tuxtla Gutiérrez, Mexico, 16–22 November 2014. Special session. Revised papers. Gelbukh, A., Espinoza, F. C., Galicia-Haro, S. N. (Eds.), pp. 16–23. IEEE, Los Alamitos.
- D. Poirier, F. Fessant, I. Tellier. 2010. *Reducing the cold-start problem in content recommender through opinion classification*. Proc. IEEE/WIC/ACM Int. Conf. WI-IAT, pp. 204–207. IEEE Computer Society, Washington, USA.
- E. Pronoza, S. Volskaya, E. Yagunova. 2014. *Corpus-based Information Extraction and Opinion Mining for the Restaurant Recommendation System*. Proceedings of the 2nd Statistical Language and Speech Processing. L. Besacier et al. (Eds.): SLSP 2014, LNAI, vol. 8791, pp. 272–284. Springer, Berlin.
- P. Resnick, H. R. Varian. 1997. *Recommender systems*. Commun. ACM, vol. 40 (3), pp. 56–58, March 1997. ACM, New York, USA.
- B. Smyth. 2007. *Case-based recommender*. The adaptive web, pp. 342–376. Springer-Verlag, Berlin, Heidelberg, Germany.
- E. Stamatatos, N. Fakotakis, G. Kokkinakis. 2000. *Automatic text categorization in terms of genre and author*. Comput. Linguist. 26, 2000, pp. 471–495.
- A. S. Starostin, V. V. Bocharov, S. V. Alexeeva, A. A. Bodrova, A. S. Chuchunkov, S. S. Dzhumaev, I. V. Efimenko, D. V. Granovskiy, V. F. Khoroshevsky, I. V. Krylova, M. A. Nikolaeva, I. M. Smurov, S. Y. Toldova. 2016. *FactRuEval 2016: Evaluation of Named Entity Recognition and Fact Extraction Systems for Russian*. Computational Linguistics and Intellectual Technologies: Proceedings of the Annual International Conference “Dialogue” (Moscow, June 1–4, 2016), issue 15 (22), pp. 702–720. RSUH, Moscow.
- Triplantica. 2016. URL: <http://www.triplantica.com> (дата обращения: 21.04.2016).
- M. Uschold, M. Gruninger. 2004. *Ontologies and semantics for seamless connectivity*. SIGMOD Rec., 33(4), pp. 58–64. ACM, New York, USA.
- L. Wang, K. Zhang. 2008. *Space efficient algorithms for ordered tree comparison*. Algorithmica, July 2008, Vol. 51, Issue 3, pp. 283–297. Springer, Berlin.