

Лекция 6. АНАЛИЗ ТЕКСТОВОЙ И ГРАФИЧЕСКОЙ ИНФОРМАЦИИ

Интеллектуальный анализ текстов (ИАТ, англ. text mining) — направление в искусственном интеллекте, целью которого является получение информации из коллекций текстовых документов, основываясь на применении эффективных в практическом плане методов машинного обучения и обработки естественного языка. Название «интеллектуальный анализ текстов» перекликается с понятием «интеллектуальный анализ данных» (ИАД, англ. data mining), что выражает схожесть их целей, подходов к переработке информации и сфер применения; разница проявляется лишь в конечных методах, а также в том, что ИАД имеет дело с хранилищами и базами данных, а не электронными библиотеками и корпусами текстов.

Text Mining - это набор технологий и методов, предназначенных для извлечения информации из текстов. Основная цель - дать аналитику возможность работать с большими объемами исходных данных за счет автоматизации процесса извлечения нужной информации.

Задачи Text Mining. Ключевыми группами задач ИАТ являются: категоризация текстов, извлечение информации и информационный поиск, обработка изменений в коллекциях текстов, а также разработка средств представления информации для пользователя.

В соответствии с уже сформированной методологии к основным задачам Text Mining относятся:

- классификация (classification)

При классификации текстов используются статистические корреляции для построения правил размещения документов в определенные категории. Задача классификации - это классическая задача распознавания, где по некоторой контрольной выборке система относит новый объект к той или другой категории.

Особенность систем Text Mining заключается в том, что количество объектов и их атрибутов может быть очень большой, поэтому должны быть предусмотрены интеллектуальные механизмы оптимизации процесса классификации.

В существующих сегодня системах классификация применяется, например, в таких задачах: группировка документов в intranet-сетях и на Web-сайтах, размещение документов в определенные папки,

сортировка сообщений электронной почты, избирательное распространение новостей подписчикам.

- кластеризация (clustering)

Кластеризация базируется на признаках документов, которые использует лингвистические и математические методы без использования определенных категорий. Результат - таксономия или визуальная карта, которая обеспечивает эффективный охват больших объемов данных.

Кластеризация в Text Mining рассматривается как процесс выделения компактных подгрупп объектов с близкими свойствами. Система должна самостоятельно найти признаки и разделить объекты по подгруппам. Кластеризация, как правило, precedes классификации, поскольку разрешает определить группы объектов. Различают два основных типа кластеризации - иерархическую и бинарную.

Кластеризация применяется при реферировании больших документальных массивов, определение взаимосвязанных групп документов, упрощения процесса просмотра при поиске необходимой информации, нахождения уникальных документов из коллекции, выявления дубликатов или очень близких по содержанию документов.

- построение семантических сетей

Построение семантических сетей или анализ связей, которые определяют появление дескрипторов (ключевых фраз) в документе для обеспечения навигации.

- извлечение фактов, понятий (feature extraction),

Извлечение фактов, предназначенное для получения некоторых фактов из текста с целью улучшения классификации, поиска и кластеризации.

- суммаризация (summarization),
- ответ на запросы (question answering),
- тематическое индексирование (thematic indexing),
- поиск по ключевым словам (keyword searching).

Также в некоторых случаях набор дополняют средства поддержки и создание таксономии (oftaxonomies) и тезаурусов (thesauri)

Этапы Text Mining. Процесс анализа текста можно разделить на 5 этапов, как показано на рис. 6.1.

Первый этап Text Mining заключается в поиске информации для его дальнейшей обработки.

Вторым этапом Text Mining является предварительная обработка текста. Данный этап включает в себя следующие пункты:

-



- отбора документов из коллекции
- пометки определённых терминов в тексте

- Определение частых наборов слов и объединение их в ключевые понятия (Apriori)
- Идентификация фактов в текстах и извлечение их характеристик
- Применение шаблонов

Так, извлечение ключевых понятий с помощью шаблонов подразделяется на 2 этапа: локальный анализ и анализ понятий (рис. 6.2).

Допустим нам дан текст:

«Петр Сергеевич Иванов покинул должность вице-президента известной фабрики ООО «Анкор». Его заменил Иван Андреевич Сидоров»

Первым этапом будет произведен лексический анализ. Текст делится на предложения и лексемы. Словарь должен включать специальные термины, имена людей, названия городов, префиксы компаний («ООО», «ЗАО», «АО»). Лексемы: «Петр», «Иван» - имена, «ООО» - префикс фирмы.

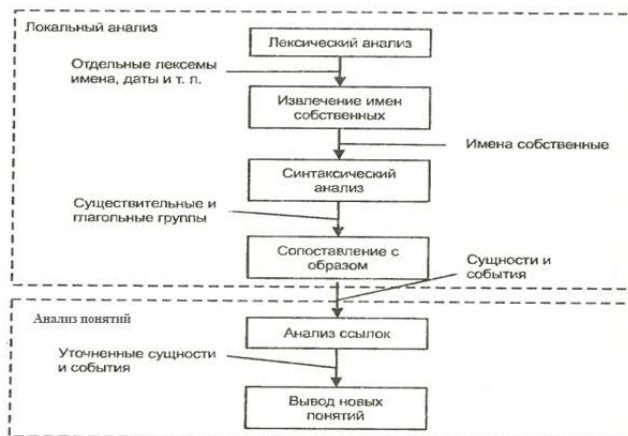


Рис. 6.2 Использование шаблонов

Далее идет извлечение имен. Имена идентифицируются с помощью образцов (регулярных выражений), которые строятся на основе частей речи, синтаксических и орфографических свойств. В результате получим следующую структуру:

- Петр Сергеевич Иванов с типом «человек»
- Иван Андреевич Сидоров с типом «человек»
- ООО «Анкор» с типом «фирма»

[имя собственное тип: человек Петр Сергеевич Иванов] должность вице-президента известной фабрики [имя собственное тип: фирма ООО «Анкор»]. Его заменил [имя собственное тип: человек Иван Андреевич Сидоров].

Синтаксический анализ заключается в построении структур для групп имён существительных (имя сущ. + его модификации) и глагольных групп (глагол+ вспомогательные части).

1. Помечаются все основные группы имён сущ. меткой «сущ.»

2. Помечаются глагольные группы меткой «гл.»

Для каждой группы имён существительных создаётся сущность. В данном примере их будет 6. Наборы образцов используют для укрупнения групп имён существительных. Образцы объединяют 2 группы имён существительных и промежуточные слова в большую группу (рис. 6.3).

[сущ. сущность: е1 Петр Сергеевич Иванов] [гл. покинул] [сущ. сущность: е2 должность вице-президента известной фирмы ООО "Анкор"]. [сущ. сущность: е5 Его] [гл. заменил] [сущ. сущность: е6 Иван Андреевич Сидоров].

Таким образом, список сущностей обновится следующим образом:

- ☐ е1 — тип: человек, имя: "Петр Сергеевич Иванов";
- ☐ е2 — тип: должность, значение: "вице-президент" фирмы: е3;
- ☐ е3 — тип: фирма, имя: "ООО "Анкор";
- ☐ е5 — тип: человек;
- ☐ е6 — тип: человек, имя: "Иван Андреевич Сидоров".

Рис. 6.3 Группировка структур

Конечным этапом идет сопоставление с образом. Для извлечения событий и отношений используются образцы, которые получаются за счёт расширения образцов, описанные ранее. Событие преемственности должности извлекается с помощью следующих образцов: человек покинул должность, человек заменяется человеком (рис. 6.4).

В анализе понятий первым происходит анализ ссылок, т.е. разрешение ссылок, представленных местоимениями и описываемыми группами имён сущ. Так, для разрешения ссылки е5 будет выполняться поиск первой предшествующей сущности с типом «человек». В нашем примере такой сущностью является е1. В результате ссылки на е5 должны быть заменены ссылками на е1. Таким образом, список сущностей и событий обновится следующим образом, как показано на рис. 6.5.

[событие: *e7* Петр Сергеевич Иванов покинул должность вице-президента известной фирмы ООО "Анкор"]. [событие: *e8* Его заменил Иван Андреевич Сидоров].

Список сущностей обновляется следующим образом:

- *e1* — тип: человек, имя: "Петр Сергеевич Иванов";
- *e2* — тип: должность, значение: "вице-президент" фирмы: *e3*;
- *e3* — тип: фирма, имя: "ООО "Анкор";
- *e5* — тип: человек;
- *e6* — тип: человек, имя: "Иван Андреевич Сидоров";
- *e7* — тип: покинул, человек: *e1*, должность: *e2*;
- *e8* — тип: заменил, человек: *e6*, человек: *e5*.

Рис. 6.4 Выделений структур событий

- *e1* — тип: человек, имя: "Петр Сергеевич Иванов";
- *e2* — тип: должность, значение: "вице-президент" фирмы: *e3*;
- *e3* — тип: фирма, имя: "ООО "Анкор";
- *e6* — тип: человек, имя: "Иван Андреевич Сидоров";
- *e7* — тип: покинул, человек: *e1*, должность: *e2*;
- *e8* — тип: заменил, человек: *e6*, человек: *e1*.

Рис. 6.5 Новое выделение структур

В результате последовательности действий можно получить следующие извлечённые ключевые понятия (таблица 6.1.)

Таблица 6.1.

	Событи е	Человек	Должность	Фирма
1	Покину л	Петр Сергеевич Иванов	Вице- президент	ООО «Анкор»
2	Вступил	Иван Андреевич Сидоров	Вице- президент	ООО «Анкор»

Методы классификации текстовых документов. Классификация — это отнесение каждого документа в определённый класс с заранее известными параметрами, полученными на этапе обучения. Число классов строго ограничено. Существует три подхода к задаче классификации текстов.

Во-первых, классификация не всегда осуществляется с помощью компьютера. Например, в обычной библиотеке тематические рубрики

присваиваются книгам вручную библиотекарем. Подобная ручная классификация дорога и неприменима в случаях, когда необходимо классифицировать большое количество документов с высокой скоростью.

Другой подход заключается в написании правил, по которым можно отнести текст к той или иной категории. Например, одно из таких правил может выглядеть следующим образом: "если текст содержит слова производная и уравнение, то отнести его к категории математика". Специалист, знакомый с предметной областью и обладающий навыком написания регулярных выражений, может составить ряд правил, которые затем автоматически применяются к поступающим документам для их классификации. Этот подход лучше предыдущего, поскольку процесс классификации автоматизируется и, следовательно, количество обрабатываемых документов практически не ограничено. Более того, построение правил вручную может дать лучшую точность классификации, чем при машинном обучении (см. ниже). Однако создание и поддержание правил в актуальном состоянии (например, если для классификации новостей используется имя действующего президента страны, соответствующее правило нужно время от времени изменять) требует постоянных усилий специалиста.

Наконец, третий подход основывается на машинном обучении. В этом подходе набор правил или, более обще, критерий принятия решения текстового классификатора, вычисляется автоматически из обучающих данных (другими словами, производится обучение классификатора). Обучающие данные — это некоторое количество хороших образцов документов из каждого класса. В машинном обучении сохраняется необходимость ручной разметки (термин разметка означает процесс приписывания класса документу). Но разметка является более простой задачей, чем написание правил. Кроме того, разметка может быть произведена в обычном режиме использования системы. Например, в программе электронной почты может существовать возможность помечать письма как спам, тем самым формируя обучающее множество для классификатора — фильтра нежелательных сообщений. Таким образом, классификация текстов, основанная на машинном обучении, является примером обучения с учителем, где в роли учителя выступает человек, задающий набор классов и размечающий обучающее множество.

Выделяют следующие методы классификации:

1. CI – Concept Indexing. Разбивает множество документов методом рекурсивной бисекции, т.е. разделяя множество документов на две части на каждом шаге рекурсии. Метод может использовать информацию, полученную на этапе обучения.

2. SOM – Self-Organizing Maps. Производит классификацию документов с использованием самонастраивающейся нейронной сети.

3. Метод опорных векторов. Метод позволяет найти оптимальную прямую или гиперплоскость для разделения данных на классы.

Методы кластеризации текстовых документов. Кластеризация – разбиение множества документов на кластеры – подмножества, параметры которых заранее неизвестны. Количество кластеров может быть произвольным или фиксированным. На данный момент существует множество методов, осуществляющих кластеризацию или классификацию документов. Некоторые методы могут использовать несколько альтернативных алгоритмов.

1. Custom Search Folders – позволяет сузить результаты поиска путём распределения их по «папкам» (folders). Эта система уже реализована в поисковом сервере и имеет название NorthernLight.

2. LSA/LSI – Latent Semantic Analysis/Indexing. Путём факторного анализа множества документов выявляются латентные (скрытые) факторы, которые в дальнейшем являются основой для образования кластеров документов;

3. STC – Suffix Tree Clustering. Кластеры образуются в узлах специального вида дерева – суффиксного дерева, которое строится из слов и фраз входных документов;

4. Single Link, Complete Link, Group Average – эти методы разбивают множество документов на кластеры, расположенные в древовидной структуре – dendrogramm, получаемой с помощью иерархической аггломеративной кластеризацией;

5. Scatter/Gather. Представляется как итеративный процесс, сначала разбивающий (scatter) множество документов на группы и представлении затем этих групп пользователю (gather) для дальнейшего анализа. Далее процесс повторяется снова над конкретными группами.

6. K-means. Относится к неиерархическим алгоритмам. Кластеры представлены в виде центроидов (см. ниже), являющихся “центром массы” всех документов, входящих в кластер.

К числовым методам кластеризации присущи два способа выполнения кластеризации множества документов: top-down, bottom-

ур. Top-down – весь имеющийся объем документов изначально рассматривается как единый кластер, и происходит его деление на более мелкие составляющие, до тех пор, пока не сработает остановочный критерий. Как правило, этим критерием является количество кластеров. Bottom-up – изначально рассматривает каждый документ, как отдельный кластер. В процессе работы наиболее близкие документы объединяются в кластеры, содержащие всё больше и больше документов. На остановку метода влияет остановочный критерий. К методам top-down относится CI в его «необучаемой» реализации. Bottom-up - это Single Linkage, Complete Linkage, Group Average и алгоритмы Buckshot и Fractionation применяемые в методе Scatter/Gather.

Аннотирование текстовых документов. Автоматическое реферирование, аннотирование (summarization) — построение краткого содержания документа по его полному тексту. Задача аннотирования документов является актуальной для любых хранилищ информации: от библиотек до интернет-порталов. Аннотирование требуется также и конкретному человеку, например, для быстрого ознакомления с интересующей его публикацией или с подборкой статей по одной тематики.

В настоящее время наиболее распространено ручное аннотирование, к достоинствам которого можно отнести, безусловно, высокое качество составления аннотации - ее "осмысленность". Типичные недостатки ручной системы аннотирования - высокие материальные затраты и присущая ей низкая скорость. Хорошее аннотирование предполагает содержание в аннотации предложений, представляющих максимальное количество тем, представленных в документе, при минимальной избыточности.

Процесс аннотирования распадается на три этапа:

1. Анализ исходного текста.
2. Определение его характерных фрагментов.
3. Формирование соответствующего вывода.

Большинство современных работ концентрируются вокруг разработанной технологии реферирования одного документа. Выделяют два основных подхода к автоматическому аннотированию текстовых документов:

- Извлечение — предполагает выделение наиболее важных фрагментов (чаще всего это предложения) из исходного текста и соединение их в аннотацию.

- **Обобщение** — предполагает использование предварительно разработанных грамматик естественных языков, тезаурусы, онтологические справочники и др., на основании которых выполняется переформулирование исходного текста и его обобщение.

В подходе, основанном на извлечении фрагментов методов сопоставлении шаблонов, выделяют наиболее лексически и статистически значимые части. В результате аннотация в данном случае создается простым соединением выбранных фрагментов.

В большинстве методов, основанных на данном подходе, используются весовые коэффициенты, вычисляемые для каждого фрагмента. Вычисления выполняются в соответствии с такими характеристиками, как расположение фрагмента в тексте, частота появления, частота использования в ключевых предложениях, а также показатели статистической значимости. Общий вид формулы вычисления веса фрагмента текста U выглядит следующим образом

$$Weight(U) = Location(U) + KeyPhrase(U) + StatTerm(U) + AddTerm(U).$$

Весовой коэффициент расположения ($Location$) в данной модели зависит от того, где во всем тексте или в отдельно взятом параграфе появляется данный фрагмент — в начале, в середине или в конце, а также используется ли он в ключевых разделах, например, во вводной части или в заключении. Ключевые фразы представляют собой лексические резюмирующие конструкции, такие как "в заключение", "в данной статье", "согласно результатам анализа" и т. д. Весовой коэффициент ключевой фразы ($KeyPhrase$) может зависеть также и от принятого в данной предметной области оценочного термина, например, "отличный" (наивысший коэффициент) или "малозначущий" (значительно меньший коэффициент).

Кроме того, при назначении весовых коэффициентов в этой модели учитывается показатель статистической важности ($StatTerm$). Статистическая важность вычисляется на основании данных, полученных в результате анализа автоматической индексации, при которой вычисляются весовые коэффициенты-ты лексем (например, методами TFIDF или TLTF).

И наконец, эта модель предполагает просмотр терминов в фрагменте текста и определение его весового коэффициента в соответствии с дополнительным наличием терминов ($AddTerm$) — появляются ли они также в заголовке, в колонтитуле, в первом параграфе и в пользовательском запросе. Выделение приоритетных терминов, наиболее точно отражающих интересы пользователя, — это

один из путей настроить аннотацию на конкретного человека или группу.

В подходе обобщения для подготовки аннотации требуются мощные вычислительные ресурсы для систем обработки естественных языков (NLP — Natural Language Processing), в том числе грамматики и словари для синтаксического разбора и генерации естественно-языковых конструкций. Кроме того, для реализации этого метода нужны некие онтологические справочники, отражающие соображения здравого смысла, и понятия, ориентированные на предметную область, для принятия решений во время анализа и определении наиболее важной информации. Данный подход предполагает использование двух основных типов методов. Первый тип опирается на традиционный лингвистический метод синтаксического разбора предложений. В этом методе применяется также семантическая информация для аннотирования деревьев разбора. Процедуры сравнения манипулируют непосредственно деревьями с целью удаления и перегруппировки частей, например, путем сокращения ветвей на основании некоторых структурных критериев, таких как скобки или встроенные условные или подчиненные предложения. После такой процедуры дерево разбора существенно упрощается, становясь, по существу, структурной "выжимкой исходного текста.

Второй тип методов аннотирования опирается на понимание естественного языка. Синтаксический разбор также входит составной частью в такие методы анализа, но деревья разбора в этом случае не порождаются. Напротив, формируются концептуальные структуры, отражающие всю исходную информацию, которая аккумулируется в текстовой базе знаний. В качестве структур могут быть использованы формулы логики предикатов или такие представления, как семантическая сеть или набор фреймов. Примером может служить шаблон банковских транзакций (заранее определенное событие), в котором перечисляются организации и лица, принимающие в нем участие, дата, объем перечисляемых средств, тип транзакции и т.д.

Подход, основанный на извлечении фрагментов, легко настраивается для обработки больших объемов информации. Из-за того, что работа таких методов основана на выборке отдельных фрагментов, предложений или фраз, текст аннотации, как правило, лишен связности. С другой стороны, такой подход выдает более сложные аннотации, которые нередко содержат информацию, дополняющую исходный текст. Так как он опирается на формальное представление информации в документе, то его можно настроить на

достаточно высокую степень сжатия, например, для рассылки сообщений на мобильные устройства.

Подход, основанный на обобщении и предполагающий опору на знания, как правило, требует полноценных источников знаний. Это является серьезным препятствием для его широкого распространения. Поэтому разработчики средств автоматического аннотирования все больше склоняются к гибридным системам, а исследователям все более успешно удается объединять статистические методы и методы, основанные на знаниях.

Visual Mining. VDM - процесс извлечения скрытых, не выраженных явным образом полезных знаний из больших наборов данных (data set); помогает специалистам обнаруживать новые тенденции, выявлять скрытые связи и закономерности в массивах данных, представленных в самых разных форматах. Состоит из четырёх направлений: визуализации данных (data visualization), визуализации результатов извлечения знаний из данных (mining result visualization), визуализации процесса извлечения знаний из данных (mining process visualization) и интерактивного визуального анализа данных (interactive visual mining) с помощью методов визуализации

За счёт того, что пользователь напрямую работает с данными, представленными в виде визуальных образов, которые он может рассматривать с разных сторон и под любыми углами зрения, в прямом смысле этого слова, он может получить дополнительную информацию, которая поможет ему более чётко сформулировать цели исследования.

Таким образом, визуальный анализ данных можно представить как процесс генерации гипотез. При этом сгенерированные гипотезы можно проверить или автоматическими средствами (методами статистического анализа или методами Data Mining), или средствами визуального анализа.

Визуальный анализ данных обычно выполняется в три этапа:

- беглый анализ - позволяет идентифицировать интересные шаблоны и сфокусироваться на одном или нескольких из них;
- увеличение и фильтрация - идентифицированные на предыдущем этапе шаблоны отфильтровываются и рассматриваются в большем масштабе;
- детализация по необходимости - если пользователю нужно получить дополнительную информацию, он может визуализировать более детальные данные

Методы геометрических преобразований. Методы геометрических преобразований визуальных образов направлены на трансформацию многомерных наборов данных с целью отображения их в декартовом и в недекартовом геометрических пространствах. Данный класс методов включает в себя математический аппарат статистики. К методам геометрических преобразований относятся:

- матрица диаграмм разброса;
- параллельные координаты;
- методы, ориентированные на пиксели
- рекурсивные шаблоны;
- циклические сегменты;
- иерархические образы
- наложение измерений.
- точки и матрицы;
- гипердоли;
- поверхностные и объемные графики, контуры;
- параллельные координаты;
- текстуры и растры.

Матрица диаграмм разброса (Scatterplot Matrix) является комбинацией отдельных диаграмм разброса, что позволяет отображать более одного атрибута. Значения атрибутов отображаются в диагональных ячейках матрицы, а остальные ячейки представляют собой отношения между ними. Например, на рис. 6 показана матрица 5×5 . Вдоль диагонали изображаются гистограммы пяти атрибутов, а, например, ячейка (2, 3) представляет отношение атрибута 2 с атрибутом 3. Соответственно, ячейка (3, 2) представляет отношение атрибута 3 с атрибутом 2 (рис.6.6.).

В данном методе визуализации могут быть использованы такие типы взаимодействия, как соприкосновение и связывание. Например, когда пользователь наводит курсор, или щелкает мышью на определенной точке, или выбирает несколько точек в одной из ячеек, представляющих отношение, то в остальных ячейках матрицы могут подсвечиваться эквивалентные точки.

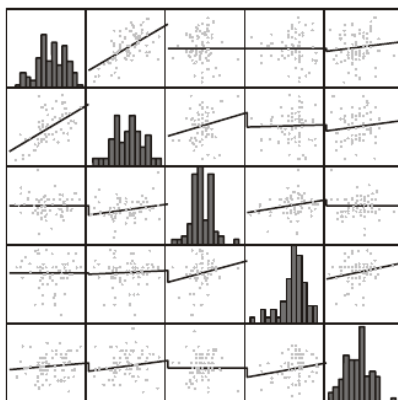


Рис. 6.6 Пример матрицы диаграмм разброса

Гипердоли являются модификацией матрицы диаграмм разброса. Основная концепция та же, за исключением того, что в ячейках матрицы отображаются скалярные функции. Таким образом, в диагональных ячейках матрицы отображается скалярная функция, представляющая отдельные атрибуты, а в остальных ячейках — скалярное отношение нескольких атрибутов.

Пользователь может взаимодействовать с данным представлением, описав визуальный фокус и диапазон значений (например, так, как в ячейке (2, 3) на рис. 6.7). При этом отображаться будут только данные в заданном диапазоне. Перемещая фокус, пользователь может быстро исследовать другие данные из близлежащих диапазонов.

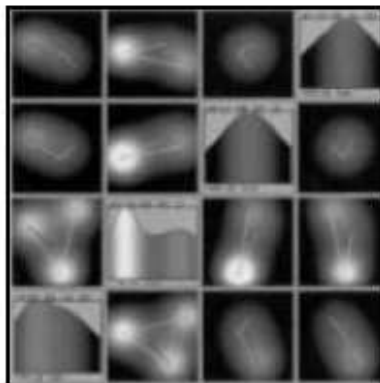


Рис. 6.7 Пример гипердолей

На ранних фазах визуального анализа большие величины непрерывных данных могут отображаться с помощью объема. Объемный рендеринг позволяет пользователю видеть внутреннюю часть объемных графиков. Цвета, яркость и полупрозрачность используются, чтобы изобразить различия распределений и значения атрибутов. Подвижность объемных графиков используется, чтобы визуализировать различные их слои.

Объемные графики (рис. 6.8) представляют собой 3D-плоскость, на которой отображается отношение между данными. Контурные линии используются для соединения точек, соответствующих данным с одинаковыми атрибутами. Однако представление большого количества данных с помощью данного метода может быть затруднено из-за густоты точек и, как следствие, затемненности и неясности изображения.

Еще одним распространенным методом геометрических преобразований является метод параллельных координат. Данный метод предполагает представление атрибутов параллельными линиями на недеклартовой плоскости. Данные представляются кривыми линиями, которые пересекают линии атрибутов. Точки пересечений соответствуют значениям соответствующих атрибутов, отображаемых данных. На рис. 6.9 приведен пример для данных, характеризующихся 10-ю измерениями.

Это достаточно простой способ представления многомерных данных, но при большом количестве линий получается большая зашумленность изображения, что приводит к неинформативности визуализации.

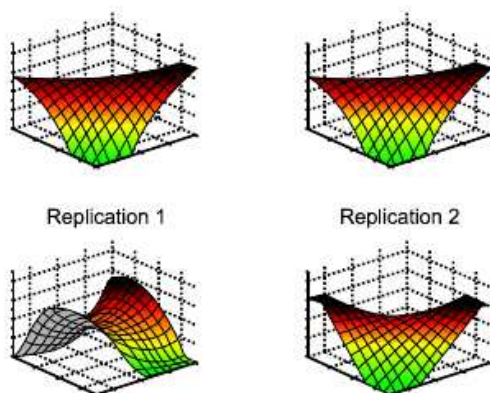


Рис. 6.8 Пример объемных графиков

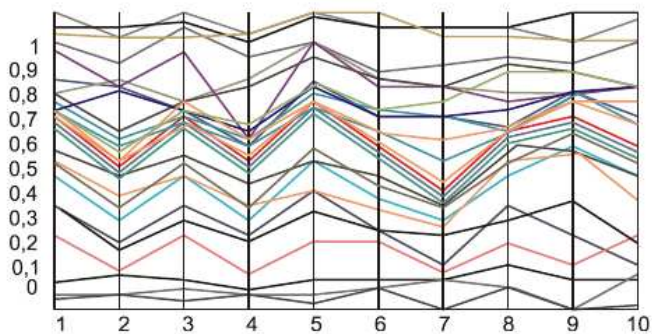


Рис. 6.9 Пример параллельных координат

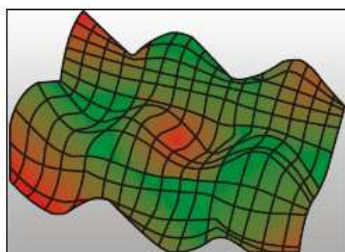


Рис. 6.10 Пример отображения текстур

Текстурная и растровая визуализации используют способность человека к преаттентивному (подсознательному) восприятию информации. Такой метод в совокупности с различными визуальными свойствами (такими как подсветка и интенсивность) позволяет отобразить большое количество атрибутов. Например, на рис. 6.10 с помощью текстуры представляется векторная и контурные диаграммы на плоскости.