

# Оглавление

<b>Предисловие .....</b>	<b>12</b>
От издательства .....	15
<b>Введение .....</b>	<b>16</b>
<b>Почему Data Science? .....</b>	<b>18</b>
<b>Глава 1. Об основах без лишних слов.....</b>	<b>21</b>
1.1. Подготовка данных.....	22
Формат данных.....	23
Типы переменных .....	24
Выбор переменных .....	25
Конструирование признаков .....	25
Неполные данные .....	26
1.2. Выбор алгоритма.....	27
Обучение без учителя.....	28
Обучение с учителем .....	29
Обучение с подкреплением .....	30
Другие факторы.....	31

---

1.3. Настройка параметров .....	31
1.4. Оценка результатов .....	33
Метрики классификации .....	34
Метрика регрессии .....	35
Валидация.....	36
1.5. Краткие итоги .....	38
<b>Глава 2. Кластеризация методом k-средних .....</b>	<b>39</b>
2.1. Поиск кластеров клиентов.....	40
2.2. Пример: профили кинозрителей .....	41
2.3. Определение кластеров.....	42
Сколько кластеров существует?.....	44
Что включают кластеры? .....	46
2.4. Ограничения .....	48
2.5. Краткие итоги.....	49
<b>Глава 3. Метод главных компонент .....</b>	<b>51</b>
3.1. Изучение пищевой ценности .....	52
3.2. Главные компоненты.....	53
3.3. Пример: анализ пищевых групп.....	56
3.4. Ограничения .....	61
3.5. Краткие итоги.....	64
<b>Глава 4. Ассоциативные правила .....</b>	<b>65</b>
4.1. Поиск покупательских шаблонов.....	66
4.2. Поддержка, достоверность и лифт .....	67

4.3. Пример: ведение продуктовых продаж .....	69
4.4. Принцип Apriori .....	72
Поиск товарных наборов с высокой поддержкой.....	73
Поиск товарных правил с высокой достоверностью или лифтом.....	74
4.5. Ограничения .....	75
4.6. Краткие итоги.....	76

**Глава 5. Анализ социальных сетей ..... 77**

5.1. Составление схемы отношений.....	78
5.2. Пример: геополитика в торговле оружием.....	80
5.3. Лувенский метод .....	84
5.4. Алгоритм PageRank .....	86
5.5. Ограничения .....	90
5.6. Краткие итоги .....	91

**Глава 6. Регрессионный анализ ..... 93**

6.1. Выведение линии тренда.....	94
6.2. Пример: предсказание цен на дома .....	95
6.3. Градиентный спуск .....	98
6.4. Коэффициенты регрессии.....	101
6.5. Коэффициенты корреляции.....	102
6.6. Ограничения .....	104
6.7. Краткие итоги.....	106

**Глава 7. Метод k-ближайших соседей и обнаружение  
аномалий ..... 107**

7.1. Пищевая экспертиза.....	108
------------------------------	-----

---

7.2. Яблоко от яблони недалеко падает .....	109
7.3. Пример: истинные различия в вине .....	111
7.4. Обнаружение аномалий.....	113
7.5. Ограничения .....	114
7.6. Краткие итоги.....	115
<b>Глава 8. Метод опорных векторов.....</b>	<b>117</b>
8.1 «Нет» или «о, нет!»?.....	118
8.2. Пример: обнаружение сердечно-сосудистых заболеваний .....	118
8.3. Построение оптимальной границы.....	120
8.4. Ограничения .....	124
8.5. Краткие итоги.....	125
<b>Глава 9. Дерево решений.....</b>	<b>127</b>
9.1. Прогноз выживания в катастрофе .....	128
9.2. Пример: спасение с тонущего «Титаника» .....	128
9.3. Создание дерева решений .....	131
9.4. Ограничения .....	133
9.5. Краткие итоги.....	135
<b>Глава 10. Случайные леса.....</b>	<b>137</b>
10.1. Мудрость толпы .....	138
10.2. Пример: предсказание криминальной активности.....	139
10.3. Ансамбли .....	144
10.4. Бэггинг.....	145

10.5. Ограничения.....	147
10.6. Краткие итоги .....	148

**Глава 11. Нейронные сети ..... 149**

11.1. Создание мозга .....	150
11.2. Пример: распознавание рукописных цифр.....	152
11.3. Компоненты нейронной сети.....	156
11.4. Правила активации .....	159
11.5. Ограничения.....	161
11.6. Краткие итоги .....	165

**Глава 12. А/В-тестирование и многорукие бандиты ..... 167**

12.1. Основы А/В-тестирования.....	168
12.2. Ограничения А/В-тестирования.....	169
12.3. Стратегия снижения эpsilon.....	169
12.4. Пример: многорукие бандиты .....	171
12.5. Забавный факт: ставка на победителя.....	174
12.6. Ограничения стратегии снижения эpsilon.....	175
12.7. Краткие итоги .....	176

**Приложения ..... 179**

Приложение А. Обзор алгоритмов обучения без учителя .....	180
Приложение В. Обзор алгоритмов обучения с учителем .....	181
Приложение С. Список параметров настройки.....	182

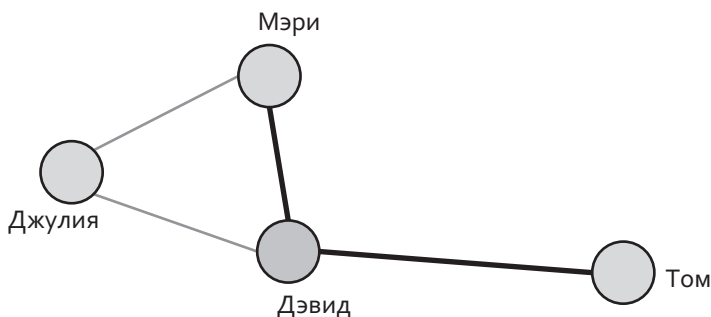
Приложение D. Другие метрики оценки.....	183
Метрики классификации .....	183
Метрики регрессии.....	186
<b>Глоссарий .....</b>	<b>188</b>
<b>Литература и ссылки на источники .....</b>	<b>199</b>
Источники на английском языке .....	199
Литература на русском языке.....	202
<b>Об авторах.....</b>	<b>204</b>

## 5.1. Составление схемы отношений

Большинство из нас имеет множество кругов общения, включающих такие категории людей, как родственники, коллеги или одноклассники. Чтобы выяснить, как устроены отношения всех этих людей, определив, например, активных персон и то, как они влияют на групповую динамику, мы можем воспользоваться методом под названием *анализ социальных сетей* (Social Network Analysis). Этот метод можно применять в вирусном маркетинге, моделировании эпидемий и даже для стратегий в командных играх. Тем не менее он больше известен своим использованием для анализа отношений в социальных сетях, что и дало ему название. На рис. 1 пример того, как анализ социальных сетей показывает отношения.

Рисунок 1 показывает сеть из четырех индивидов, также известную как *граф*, в котором каждый из этих персон представлен *узлом* (node). Отношения между узлами представлены линиями, называемыми *ребрами* (edges).

Каждое ребро может иметь *вес*, показывающий силу отношений.



**Рис. 1.** Простая сеть друзей. Более близкие отношения показаны утолщенными линиями

Из рис. 1 мы можем заключить:

- Дэвид имеет больше всех связей, будучи знакомым с остальными тремя персонами;
- Том не знает никого, кроме Дэвида, с которым они близкие друзья;
- Джулия знает Мэри и Дэвида, но не близка с ними.

Кроме отношений анализ социальных сетей может строить схемы и для других сущностей, при условии, что между ними есть связи. В этой главе мы воспользуемся им для анализа международной сети торговли оружием, чтобы выявить доминирующие силы и их сферы влияния.



## 5.2. Пример: геополитика в торговле оружием

Мы получили данные о двусторонних трансферах основных видов обычных вооружений из *Стокгольмского международного института по исследованию проблем мира*. Военные поставки были выбраны в качестве косвенного показателя двусторонних отношений, поскольку должны свидетельствовать о тесной связи стран на международной арене.

В этом анализе мы стандартизировали стоимость оружия на уровне цен 1990 года в долларах США, после чего приняли в расчет только сделки, сумма которых превысила 100 млн долларов. Чтобы учесть флуктуации в торговле оружием, обусловленные производственными циклами новых технологий, мы рассмотрели 10-летний период, с 2006 по 2015 год, построив сеть из 91 узла и 295 ребер.

Для визуализации сети использовался *силовой алгоритм* (force-directed algorithm): узлы без связей отталкиваются друг от друга, а связанные узлы, наоборот, притягиваются с той степенью близости, которая отражает силу их связи (рис. 2). Например, максимальный объем торговли зафиксирован между Россией и Индией (\$ 22,3 млрд), поэтому эти государства соединены толстой линией и близко расположены.

После анализа получившейся сети лувенским методом (Louvain Method, описан в следующем разделе) геополитические альянсы были сгруппированы в три кластера.

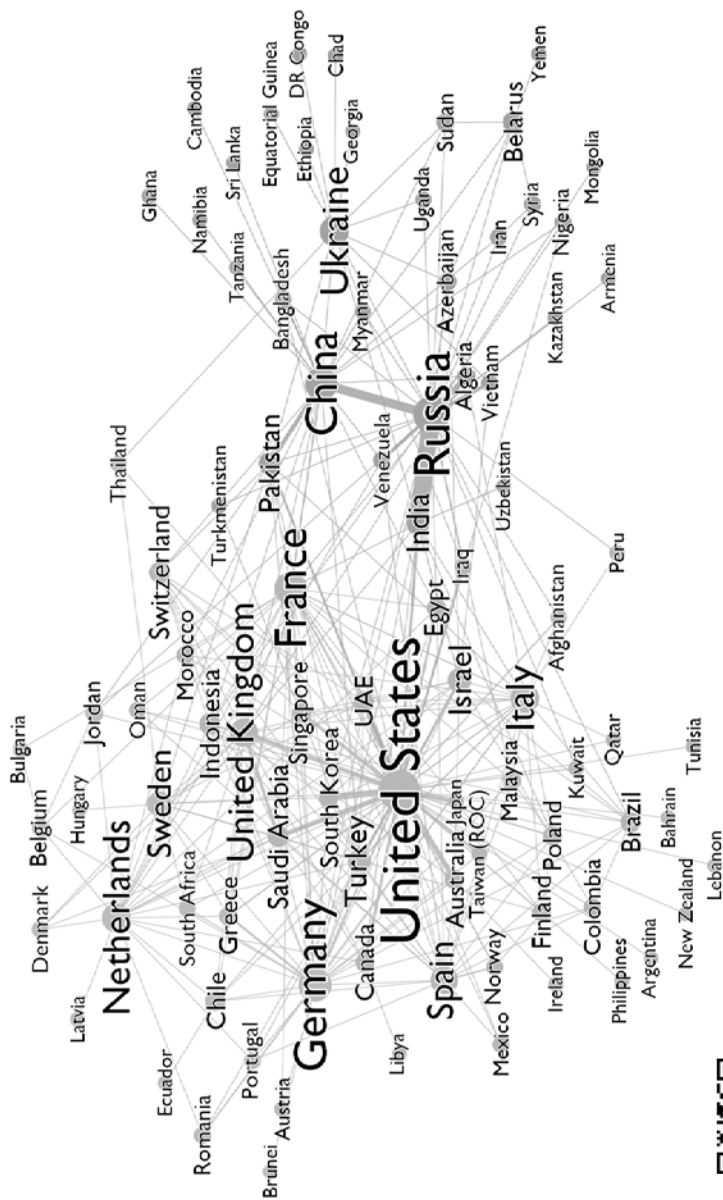


Рис. 2. Сеть стран, исходя из военных поставок

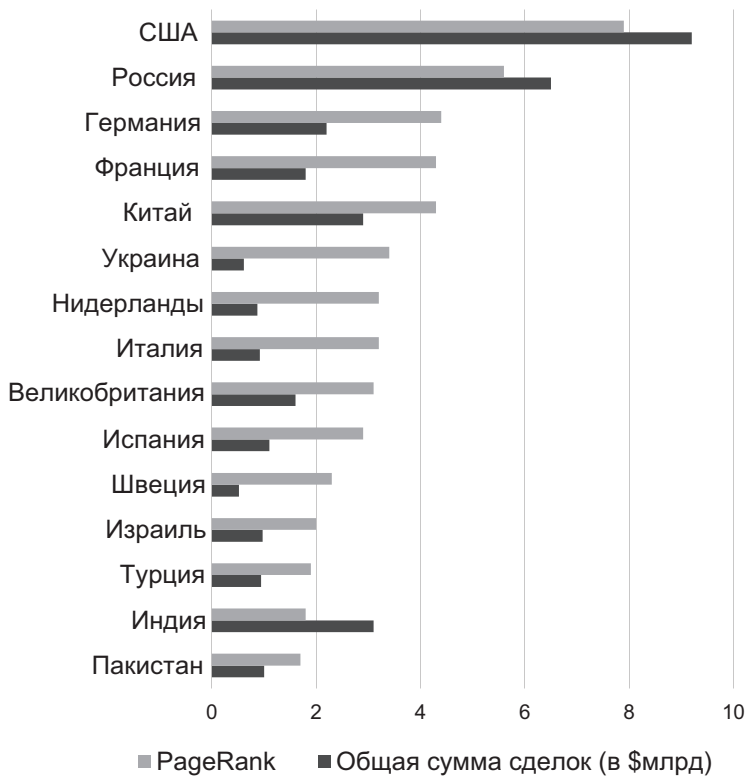


- **Светло-серый:** это крупнейший кластер, в котором доминируют США и который включает их союзников, таких как Великобритания и Израиль.
- **Светлый:** в нем лидирует Германия, и он включает в основном европейские страны, а также тесно связан со светло-серым кластером.
- **Темный:** в этом кластере доминируют Россия и Китай, он дистанцирован от двух других и включает в основном азиатские и африканские государства.

Кластеры отражают геополитические реалии XXI столетия, такие как долгосрочные альянсы между западными нациями, поляризацию между демократическими и коммунистическими странами и возрастающую роль противостояния между США и Китаем.

Кроме группировки в кластеры мы также проранжировали отдельные страны по уровню их влияния, воспользовавшись алгоритмом PageRank (описывается дальше). На рис. 3 представлены 15 самых влиятельных государств, которые также отмечены на рис. 2 более крупными узлами и подписями.

Согласно нашему анализу, в пятерку самых влиятельных стран входят США, Россия, Германия, Франция и Китай. Эти результаты подтверждаются тем обстоятельством, что четыре из пяти этих государств имеют влияние еще и как члены Совета Безопасности ООН.



**Рис. 3.** Самые влиятельные страны в торговле оружием, согласно алгоритму PageRank. Значение PageRank для каждой страны показано светлым, а торговый объем — темным

В следующих разделах мы рассмотрим методы, использованные для выделения кластеров и ранжирования стран.

### 5.3. Лувенский метод

Как видно на рис. 2, можно найти кластеры сети путем группировки узлов. Изучение этих кластеров поможет понять, чем различаются части сети и как они соединены.

*Лувенский метод* — один из способов определения кластеров сети. Он подбирает различные кластерные конфигурации, чтобы: 1) максимизировать число и силу связей между узлами в одном кластере; 2) минимизировать при этом связи между узлами различных кластеров. Степень удовлетворения этим двум условиям известна как *модулярность* (modularity), и более высокая модулярность — признак более оптимальных кластеров.

Чтобы получить оптимальную конфигурацию кластеров, лувенский метод итеративно проходит следующие стадии.

**Стадия 0:** рассматривает каждый узел в качестве кластера, то есть начинает с числа кластеров, равного числу узлов.

**Стадия 1:** меняет кластерное членство узла, если это приводит к улучшению модулярности. Если модулярность больше нельзя улучшить, узел остается на месте. Это повторяется для каждого узла до тех пор, пока изменения кластерного членства не будут исчерпаны.

**Стадия 2:** строит грубую версию сети, в которой каждый кластер, найденный на стадии 1, представлен отдельным

узлом, и объединяет бывшие межкластерные соединения в утолщенные ребра этих новых узлов в соответствии с их весом.

**Стадия 3:** повторяет стадии 1 и 2 до тех пор, пока не закончатся дальнейшие изменения членства и размера связей.

Таким образом, лувенский метод помогает нам выявить более значимые кластеры, начав с обнаружения малых из них, а затем при необходимости соединяя их.

Простота и эффективность делают лувенский метод популярным решением для кластеризации сети. Однако он имеет свои ограничения.

**Важные, но малые кластеры могут быть поглощены.** Итеративный процесс слияния кластеров может привести к тому, что значимые, но небольшие кластеры будут обойдены вниманием. Чтобы избежать этого, мы можем при необходимости проверять идентифицированные кластеры на промежуточных фазах итераций.

**Множество возможных кластерных конфигураций.** Для сетей, содержащих перекрывающиеся или вложенные кластеры, определить оптимальное кластерное решение может оказаться трудным. Тем не менее, когда имеются несколько решений с высокой модулярностью, мы можем сверить кластеры с другими информационными источниками, что мы и проделали на рис. 2, приняв во внимание географическое местоположение и политическую идеологию.