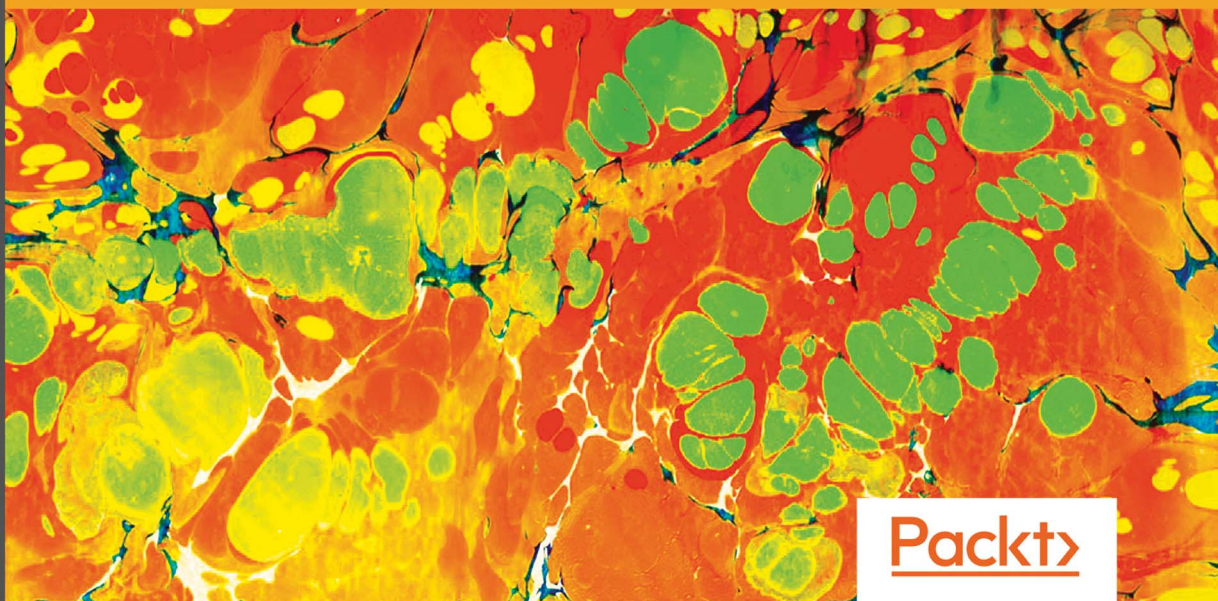


Машинное обучение для алгоритмической торговли на финансовых рынках. Практикум

Разработка инвестиционных стратегий на основе интеллектуальных,
обучаемых на данных алгоритмов и их реализация на языке Python



Стефан Янсен

Packt>



Hands-On Machine Learning for Algorithmic Trading

Design and implement investment strategies based on smart algorithms that learn from data using Python

Stefan Jansen

Packt>

BIRMINGHAM - MUMBAI

Стефан Янсен

Машинное обучение для алгоритмической торговли на финансовых рынках. Практикум

Санкт-Петербург

«БХВ-Петербург»

2020

УДК 004.8+336
ББК 32.973.26-018.1+65.26
Я65

Янсен С.

Я65 Машинное обучение для алгоритмической торговли на финансовых рынках. Практикум: Пер. с англ. — СПб.: БХВ-Петербург, 2020. — 560 с.: ил.
ISBN 978-5-9775-6595-0

Книга посвящена практике применения машинного обучения с целью создания мощных алгоритмических стратегий для успешной торговли на финансовых рынках. Изложены базовые принципы работы с данными: оценивание наборов данных, доступ к данным через API на языке Python, доступ к финансовым данным на платформе Quandl и управление ошибками предсказания. Рассмотрены построение и тренировка алгоритмических моделей с помощью Python-библиотек pandas, Seaborn, StatsModels и sklearn и построение, оценка и интерпретация моделей $AR(p)$, $MA(q)$ и $ARIMA(p, d, q)$ с использованием библиотеки StatsModels. Описано применение библиотеки PyMC3 для байесового машинного обучения, библиотек NLTK, sklearn (Scikit-learn) и spaCy для назначения меток финансовым новостям и классифицирования документов, библиотеки Keras для создания, настройки и оценки нейронных сетей прямого распространения, рекуррентных и сверточных сетей. Показано, как применять трансферное обучение к данным спутниковых снимков для предсказания экономической активности и как эффективно использовать подкрепляемое обучение для достижения оптимальных результатов торговли.

Для финансовых аналитиков и программистов на языке Python

УДК 004.8+336
ББК 32.973.26-018.1+65.26

Группа подготовки издания:

| | |
|-----------------------|--------------------------|
| Руководитель проекта | <i>Евгений Рыбаков</i> |
| Зав. редакцией | <i>Екатерина Сависте</i> |
| Перевод с английского | <i>Андрея Логунова</i> |
| Компьютерная верстка | <i>Ольги Сергиенко</i> |
| Оформление обложки | <i>Карины Соловьевой</i> |

© Packt Publishing 2019. First published in the English language under the title 'Hands-On Machine Learning for Algorithmic Trading – (9781789346411)'

Впервые опубликовано на английском языке под названием 'Hands-On Machine Learning for Algorithmic Trading – (9781789346411)'

"БХВ-Петербург", 191036, Санкт-Петербург, Гончарная ул., 20.

ISBN 978-1-78934-641-1 (англ.)
ISBN 978-5-9775-6595-0 (рус.)

© Packt Publishing 2019
© Перевод на русский язык, оформление.
ООО "БХВ-Петербург", ООО "БХВ", 2020

Оглавление

| | |
|--|-----------|
| Об авторе..... | 23 |
| О рецензентах..... | 24 |
| Комментарии переводчика..... | 25 |
| О терминологии | 26 |
| Предисловие | 28 |
| Для кого эта книга предназначена | 29 |
| Что эта книга охватывает..... | 29 |
| Получение максимальной отдачи от этой книги | 32 |
| Скачивание файлов с примерами исходного кода..... | 32 |
| Скачивание цветных изображений | 33 |
| Принятые в книге условные обозначения | 33 |
| Глава 1. Машинное обучение для торговли на финансовых рынках..... | 34 |
| Как читать эту книгу | 35 |
| Чего ожидать | 36 |
| Кто должен прочесть эту книгу | 37 |
| Как эта книга организована | 37 |
| Каркас — от данных к конструированию стратегии | 37 |
| Основы машинного обучения..... | 38 |
| Обработка естественного языка | 39 |
| Глубокое и подкрепляемое обучение | 40 |
| Что вам нужно для успеха | 40 |
| Источники данных..... | 41 |
| Репозиторий GitHub | 41 |
| Библиотеки Python..... | 42 |
| Рост популярности машинного обучения в инвестиционной индустрии | 42 |
| От торговли электронной к торговле высокочастотной | 43 |
| Факторное инвестирование и умные бета-фонды | 46 |
| Алгоритмические первопроходцы превосходят людей в широком масштабе | 51 |
| Фонды, ведомые машинным обучением, привлекают 1 трлн долларов AUM..... | 52 |

| | |
|---|------------|
| Возникновение квантоментальных фондов | 53 |
| Инвестиции в стратегический потенциал..... | 54 |
| Машинное обучение и альтернативные данные | 54 |
| Привлечение "толпы" в разработку алгоритмов торговли..... | 56 |
| Конструирование и исполнение стратегии торговли | 57 |
| Привлечение источников данных и управление данными | 58 |
| Исследование и оценивание альфа-факторов | 58 |
| Оптимизация портфеля и риск-менеджмент | 60 |
| Бэктестирование стратегии | 60 |
| Машинное обучение и стратегии алгоритмической торговли | 61 |
| Примеры использования машинного обучения для торговли на финансовых рынках | 62 |
| Глубинный анализ данных для извлечения признаков | 62 |
| Контролируемое обучение для создания и агрегирования альфа-факторов | 63 |
| Размещение финансовых средств среди портфельных активов | 64 |
| Тестирование торговых идей | 64 |
| Подкрепляемое обучение | 65 |
| Резюме | 65 |
| Глава 2. Рыночные и фундаментальные данные..... | 66 |
| Как работать с рыночными данными..... | 67 |
| Микроструктура рынка | 68 |
| Торговые площадки..... | 68 |
| Типы ордеров | 71 |
| Работа с данными ордерной книги..... | 72 |
| Протокол FIX | 72 |
| Данные ордерной книги TotalView-ITCH биржи Nasdaq..... | 73 |
| Регуляризация тиковых данных | 83 |
| Доступ к рыночным данным через API | 87 |
| Дистанционный доступ к данным с помощью библиотеки pandas | 87 |
| Платформа Quantopian | 90 |
| Библиотека zipline..... | 91 |
| Платформа Quandl | 93 |
| Другие поставщики рыночных данных | 93 |
| Как работать с фундаментальными данными | 94 |
| Данные финансовой отчетности..... | 95 |
| Автоматизированная обработка — XBRL..... | 95 |
| Построение временного ряда с фундаментальными данными | 96 |
| Другие источники фундаментальных данных | 101 |
| Макроэкономические и отраслевые данные с помощью pandas-datareader | 101 |
| Эффективное хранение данных с помощью библиотеки pandas..... | 102 |
| Резюме | 103 |
| Глава 3. Альтернативные данные для финансов..... | 104 |
| Революция альтернативных данных | 105 |
| Источники альтернативных данных | 107 |
| Физические лица..... | 107 |
| Производственные процессы предприятия | 108 |
| Датчики..... | 109 |

| | |
|--|------------|
| Оценивание альтернативных наборов данных | 110 |
| Критерии оценивания | 111 |
| Качество сигнального содержимого | 111 |
| Качество данных | 113 |
| Технические аспекты | 115 |
| Рынок альтернативных данных | 115 |
| Поставщики данных и примеры использования | 116 |
| Данные социальных настроений | 117 |
| Спутниковые данные | 118 |
| Геолокационные данные | 119 |
| Данные электронных почтовых квитанций | 119 |
| Работа с альтернативными данными | 119 |
| Выскабливание данных службы OpenTable | 119 |
| Извлечение данных из HTML с помощью библиотек requests и BeautifulSoup | 120 |
| Знакомство с библиотекой Selenium — использование браузерной автоматизации | 121 |
| Построение набора данных бронирований ресторанов | 122 |
| Следующий шаг — библиотеки Scrapy и splash | 124 |
| Стенограммы телеконференций о корпоративных заработках | 125 |
| Разбор HTML с помощью регулярных выражений | 126 |
| Резюме | 128 |
| Глава 4. Исследование альфа-факторов | 129 |
| Выработка альфа-факторов | 130 |
| Важные категории факторов | 131 |
| Импульсные и сентиментные факторы | 132 |
| Стоимостные факторы | 136 |
| Волатильностные и размерные факторы | 140 |
| Качественные факторы | 141 |
| Как преобразовывать данные в факторы | 143 |
| Полезные методы библиотек pandas и NumPy | 144 |
| Факторы, встроенные в платформу Quantopian | 147 |
| Библиотека TA-Lib | 147 |
| Поиски сигналов — как использовать библиотеку zipline | 148 |
| Архитектура — событийная симуляция торговли | 148 |
| Единственный альфа-фактор из рыночных данных | 149 |
| Сочетание факторов из разнообразных источников данных | 152 |
| Разделение сигнала и шума — как использовать библиотеку alphas | 154 |
| Создание форвардных финансовых возвратов и факторных квантилей | 155 |
| Предсказательная результативность по факторным квантилям | 156 |
| Информационный коэффициент | 159 |
| Оборачиваемость фактора | 161 |
| Ресурсы альфа-факторов | 162 |
| Альтернативные библиотеки алгоритмической торговли | 162 |
| Резюме | 163 |
| Глава 5. Оценивание стратегии | 164 |
| Как строить и тестировать инвестиционный портфель с помощью библиотеки zipline | 165 |
| Торговля по плану и ребалансировка портфеля | 165 |

| | |
|---|-----|
| Как измерять результативность с помощью библиотеки <code>pyfolio</code> | 167 |
| Коэффициент Шарпа | 167 |
| Фундаментальный закон активного менеджмента | 169 |
| Внутривыборочная и вневыборочная результативность с помощью библиотеки <code>pyfolio</code> | 170 |
| Получение входа в библиотеку <code>pyfolio</code> из библиотеки <code>alphalens</code> | 170 |
| Получение входа в библиотеку <code>pyfolio</code> из бэктеста <code>zipline</code> | 171 |
| Форвардное тестирование вневыборочных возвратов | 171 |
| Сводная статистика результативности | 172 |
| Периоды просадки и влияние факторов | 173 |
| Моделирование событийного риска | 174 |
| Как избегать ловушек бэктестирования | 175 |
| Сложности данных | 176 |
| Систематическое смещение из-за забега вперед | 176 |
| Систематическое смещение из-за выживших | 176 |
| Управление выбросами | 176 |
| Нерепрезентативный период | 177 |
| Вопросы реализации | 177 |
| Результативность на основе рыночной цены | 177 |
| Торговые издержки | 177 |
| Определение времени сделок | 178 |
| Прочесывание данных и переподгонка при бэктестировании | 178 |
| Минимальная длина бэктеста и дефлированный коэффициент Шарпа | 179 |
| Оптимальное прекращение бэктестирования | 179 |
| Как управлять портфельным риском и возвратностью | 180 |
| Среднедисперсная оптимизация | 182 |
| Как это работает | 182 |
| Эффективная граница на Python | 182 |
| Сложности и недостатки | 185 |
| Альтернативы среднедисперсной оптимизации | 186 |
| Портфель $1/n$ | 186 |
| Минимально-дисперсный портфель | 186 |
| Глобальная портфельная оптимизация — подход Блэка — Литтермана | 187 |
| Как определять размер ставок — правило Келли | 187 |
| Паритет риска | 190 |
| Риск-факторное инвестирование | 191 |
| Иерархический паритет риска | 191 |
| Резюме | 192 |

Глава 6. Процесс машинного обучения..... 194

| | |
|--|-----|
| Усвоение регулярностей из данных | 195 |
| Контролируемое обучение | 197 |
| Неконтролируемое обучение | 198 |
| Применения | 198 |
| Кластерные алгоритмы | 198 |
| Снижение размерности | 199 |
| Подкрепляемое обучение | 200 |

| | |
|---|------------|
| Рабочий поток машинного обучения | 201 |
| Простое пошаговое руководство — k ближайших соседей | 201 |
| Очертить рамки задачи — цели и метрики | 202 |
| Предсказание против статистического вывода | 202 |
| Регрессионные задачи | 204 |
| Классификационные задачи | 205 |
| Собрать и подготовить данные | 209 |
| Разведать, извлечь и выработать признаки | 209 |
| Использование теории информации для оценивания признаков | 210 |
| Отобрать автоматически обучающийся алгоритм | 211 |
| Сконструировать и настроить модели | 211 |
| Компромисс между смещением и дисперсией | 211 |
| Недоподгонка против переподгонки | 212 |
| Управление компромиссом | 213 |
| Кривые усвоения | 214 |
| Как использовать перекрестный контроль для отбора модели | 215 |
| Как реализовать перекрестный контроль на языке Python | 215 |
| Перекрестный контроль | 216 |
| Настройка параметров с помощью библиотеки <i>sklean</i> | 219 |
| Кривые перекрестного контроля с помощью библиотеки <i>yellowbricks</i> | 220 |
| Кривые усвоения | 220 |
| Настройка параметров с помощью интерфейса <i>GridSearchCV</i> и конвейера | 221 |
| Сложности перекрестного контроля в финансах | 221 |
| Перекрестный контроль временного ряда в библиотеке <i>sklearn</i> | 222 |
| Прочистка, наложение эмбарго и комбинаторный перекрестный контроль | 222 |
| Резюме | 223 |
| Глава 7. Линейные модели | 224 |
| Линейная регрессия для статистического вывода и предсказания | 225 |
| Множественная линейная регрессионная модель | 226 |
| Как формулировать модель | 227 |
| Как тренировать модель | 228 |
| Наименьшие квадраты | 228 |
| Оценивание максимального правдоподобия | 229 |
| Градиентный спуск | 230 |
| Теорема Гаусса — Маркова | 230 |
| Как проводить статистический вывод | 232 |
| Как диагностировать и устранять проблемы | 233 |
| Качество подгонки | 234 |
| Гетероскедастичность | 235 |
| Внутрирядовая корреляция | 236 |
| Мультиколлинеарность | 236 |
| Как выполнять линейную регрессию на практике | 237 |
| Обычные наименьшие квадраты из библиотеки <i>StatsModels</i> | 237 |
| Стохастический градиентный спуск с помощью библиотеки <i>sklearn</i> | 239 |
| Как строить линейную факторную модель | 240 |
| От модели <i>SARIM</i> к пятифакторной модели Фама — Френча | 241 |
| Получение рисков факторов | 242 |
| Регрессия Фама — Макбета | 244 |

| | |
|--|------------|
| Усадочные методы: регуляризация для линейной регрессии..... | 248 |
| Как хеджироваться от перепогонки..... | 248 |
| Как работает гребневая регрессия..... | 249 |
| Как работает регрессия лассо..... | 251 |
| Как применять линейную регрессию для предсказания финансовых возвратов..... | 251 |
| Подготовка данных..... | 251 |
| Создание универсума и временной горизонт..... | 251 |
| Расчет целевого финансового возврата..... | 252 |
| Отбор и преобразование альфа-факторов..... | 253 |
| Очистка данных — пропущенные данные..... | 253 |
| Разведывательный анализ данных..... | 254 |
| Кодирование категориальных переменных с помощью фиктивных значений..... | 254 |
| Создание форвардных возвратов..... | 255 |
| Линейная регрессия методом OLS с использованием библиотеки StatsModels..... | 256 |
| Диагностическая статистика..... | 256 |
| Линейная регрессия по методу OLS с использованием библиотеки sklearn..... | 257 |
| Кастомизированный перекрестный контроль временного ряда..... | 257 |
| Отбор признаков и цели..... | 258 |
| Перекрестный контроль модели..... | 258 |
| Тестовые результаты — информационный коэффициент и RMSE..... | 259 |
| Гребневая регрессия с использованием библиотеки sklearn..... | 261 |
| Настройка регуляризационных параметров с помощью перекрестного контроля..... | 261 |
| Результаты перекрестного контроля и траектории гребневых коэффициентов..... | 262 |
| Десятка ведущих коэффициентов..... | 262 |
| Лассо-регрессия с использованием библиотеки sklearn..... | 263 |
| Перекрестно-контрольный информационный коэффициент и траектория лассо-регрессии..... | 264 |
| Линейная классификация..... | 265 |
| Логистическая регрессионная модель..... | 265 |
| Целевая функция..... | 266 |
| Логистическая функция..... | 266 |
| Оценивание максимального правдоподобия..... | 267 |
| Как проводить статистический вывод с помощью библиотеки StatsModels..... | 268 |
| Как применять логистическую регрессию для предсказания..... | 270 |
| Как предсказывать ценовые движения с помощью библиотеки sklearn..... | 270 |
| Резюме..... | 272 |
| Глава 8. Модели временных рядов..... | 273 |
| Аналитические инструменты для диагностики и извлечения признаков..... | 274 |
| Как разложить временной ряд на регулярности..... | 275 |
| Как вычислять статистику скользящего окна..... | 276 |
| Скользящие средние и экспоненциальное сглаживание..... | 277 |
| Как измерять автокорреляцию..... | 277 |
| Как диагностировать и достигать стационарности..... | 278 |
| Преобразования временного ряда..... | 279 |
| Как диагностировать и обращаться с единичными корнями..... | 280 |
| Единично-корневые тесты..... | 282 |
| Как применять преобразования временных рядов..... | 283 |

| | |
|--|------------|
| Модели одномерных временных рядов | 285 |
| Как строить авторегрессионные модели | 285 |
| Как выявлять число сдвигов | 286 |
| Как диагностировать подгонку модели | 286 |
| Как строить модели скользящего среднего | 287 |
| Как выявлять число сдвигов | 287 |
| Связь между моделями AR и MA | 287 |
| Как строить модели ARIMA и их расширения | 287 |
| Как выявлять число членов моделей AR и MA | 288 |
| Добавление признаков — ARMAX | 289 |
| Добавление сезонного исчисления последовательных разностей — SARIMAX | 289 |
| Как прогнозировать макроэкономические фундаментальные показатели | 290 |
| Как использовать модели временных рядов для прогнозирования волатильности | 292 |
| Модель авторегрессионной условной гетероскедастичности (ARCH) | 293 |
| Обобщение модели ARCH — модель GARCH | 294 |
| Как строить модель прогнозирования волатильности | 294 |
| Модели многомерных временных рядов | 298 |
| Система уравнений | 298 |
| Векторная авторегрессионная модель (VAR) | 299 |
| Как использовать модель VAR для прогнозов макроэкономических фундаментальных показателей | 300 |
| Коинтеграция — временные ряды с общим трендом | 303 |
| Тестирование на коинтеграцию | 304 |
| Как использовать коинтеграцию для стратегии парной торговли | 305 |
| Резюме | 305 |
| Глава 9. Байесово машинное обучение | 307 |
| Как работает байесово машинное обучение | 308 |
| Как обновлять допущения на основе эмпирического наблюдения | 309 |
| Точный вывод: оценивание апостериорного максимума | 310 |
| Как отбирать априорные распределения | 311 |
| Как не усложнять вывод — сопряженные априорные распределения | 312 |
| Как динамически оценивать вероятности движения цены актива | 312 |
| Приближенный вывод: стохастический и детерминированный подходы | 314 |
| Стохастический вывод на основе выборок | 315 |
| Генерирование выборок методами Монте-Карло марковских цепей | 315 |
| Вариационный вывод | 318 |
| Вероятностное программирование с помощью библиотеки PyMC3 | 319 |
| Байесово машинное обучение с помощью библиотеки Theano | 319 |
| Рабочий поток библиотеки PyMC3 | 320 |
| Определение модели — байесова логистическая регрессия | 320 |
| Приближенный вывод — метод Монте-Карло марковской цепи | 323 |
| Приближенный вывод — вариационный Байес | 324 |
| Диагностика модели | 324 |
| Предсказание | 327 |
| Практические применения | 327 |
| Байесов коэффициент Шарпа и сравнение результативности | 328 |

| | |
|--|------------|
| Байесовы модели временных рядов | 330 |
| Стохастические модели волатильности | 330 |
| Резюме | 330 |
| Глава 10. Деревья решений и случайные леса | 331 |
| Деревья решений | 332 |
| Как деревья усваивают и применяют правила принятия решения | 332 |
| Как использовать деревья решений на практике | 334 |
| Как готовить данные | 334 |
| Как кодировать собственный класс перекрестного контроля | 335 |
| Как строить регрессионное дерево | 336 |
| Как строить классификационное дерево | 338 |
| Как визуализировать дерево решений | 340 |
| Как оценивать предсказания дерева решений | 340 |
| Свойство важности признаков | 342 |
| Переподгонка и регуляризация | 342 |
| Как регуляризовать дерево решений | 343 |
| Обрезка дерева решений | 345 |
| Как настраивать гиперпараметры | 345 |
| Класс <i>GridSearchCV</i> для деревьев решений | 346 |
| Как обследовать древесную структуру | 347 |
| Кривые усвоения | 348 |
| Сильные и слабые стороны деревьев решений | 349 |
| Случайные леса | 350 |
| Ансамблевые модели | 350 |
| Как бэггинг снижает дисперсию модели | 351 |
| Бутстрап-агрегированные деревья решений | 353 |
| Как строить случайный лес | 355 |
| Как тренировать и настраивать случайный лес | 356 |
| Свойство важности признаков для случайных лесов | 358 |
| Внепакетное тестирование | 359 |
| Сильные и слабые стороны случайных лесов | 360 |
| Резюме | 360 |
| Глава 11. Градиентно-бустинговые машины | 362 |
| Адаптивный бустинг | 363 |
| Алгоритм AdaBoost | 364 |
| AdaBoost в библиотеке <i>sklearn</i> | 366 |
| Градиентно-бустинговые машины | 368 |
| Как тренировать и настраивать модели на основе GBM | 370 |
| Размер ансамбля и досрочная остановка | 371 |
| Усадка и темп усвоения | 371 |
| Подвыборка и стохастический градиентный бустинг | 372 |
| Как использовать градиентный бустинг с помощью библиотеки <i>sklearn</i> | 372 |
| Как настраивать параметры с помощью класса <i>GridSearchCV</i> | 374 |
| Влияние параметров на тестовые отметки | 375 |
| Как тестировать на отложенном наборе данных | 377 |

| | |
|---|------------|
| Быстро масштабируемые реализации градиентно-бустинговых машин | 377 |
| Как алгоритмические инновации стимулируют результативность | 378 |
| Аппроксимация второпорядковой функции потерь | 378 |
| Упрощенные алгоритмы поиска разбиов | 379 |
| Поуровневый рост против полистового роста | 380 |
| Тренировка на основе GPU | 381 |
| Отсеивающая регуляризация DART — отсев для деревьев | 381 |
| Трактовка категориальных признаков | 382 |
| Дополнительные признаки и оптимизации | 382 |
| Как использовать библиотеки XGBoost, LightGBM и CatBoost | 383 |
| Как создавать двоичные форматы данных | 383 |
| Как оценивать результаты | 388 |
| Результаты перекрестного контроля между моделями | 388 |
| Как интерпретировать результаты GBM | 391 |
| Свойство важности признаков | 392 |
| Графики частичной зависимости | 393 |
| Аддитивные объяснения Шепли (SHapley) | 395 |
| Как резюмировать значения SHAP по признакам | 396 |
| Как использовать графики силы для объяснения предсказания | 397 |
| Как анализировать взаимодействие признаков | 398 |
| Резюме | 399 |
| Глава 12. Неконтролируемое обучение | 401 |
| Снижение размерности | 402 |
| Линейные и нелинейные алгоритмы | 403 |
| Проклятие размерности | 405 |
| Линейное снижение размерности | 407 |
| Анализ главных компонент | 407 |
| Анализ независимых компонент | 413 |
| Анализ главных компонент для алгоритмической торговли | 415 |
| Рисковые факторы, ведомые данными | 415 |
| Характеристические портфели | 418 |
| Усвоение проекций в топологическом многообразии | 421 |
| Алгоритм t-SNE | 423 |
| Алгоритм UMAP | 424 |
| Кластеризация | 425 |
| Кластеризация на основе k средних | 426 |
| Оценивание качества кластеров | 427 |
| Иерархическая кластеризация | 429 |
| Визуализация — дендрограммы | 430 |
| Плотностная кластеризация | 432 |
| Алгоритм DBSCAN | 432 |
| Иерархический алгоритм DBSCAN | 432 |
| Модели гауссовых смесей | 433 |
| Алгоритм максимизации ожиданий | 433 |
| Иерархический паритет риска | 434 |
| Резюме | 436 |

| | |
|---|------------|
| Глава 13. Работа с текстовыми данными | 438 |
| Как извлекать признаки из текстовых данных..... | 439 |
| Сложности обработки естественного языка..... | 439 |
| Рабочий поток обработки естественного языка..... | 440 |
| Разбор и лексемизация текстовых данных | 441 |
| Лингвистическая аннотация | 442 |
| Семантическая аннотация..... | 442 |
| Закрепление меток..... | 442 |
| Примеры использования | 443 |
| От текста к лексемам — конвейер обработки естественного языка | 443 |
| Конвейер по обработке естественного языка с помощью библиотек spaCy и textacy | 443 |
| Разбор, лексемизирование и аннотирование предложения | 444 |
| Пакетная обработка документов | 446 |
| Определение границ предложений | 446 |
| Распознавание именованных сущностей..... | 447 |
| <i>N</i> -граммы..... | 447 |
| Потоковый API библиотеки spaCy..... | 448 |
| Многоязычная обработка естественного языка | 448 |
| Обработка естественного языка с помощью библиотеки TextBlob | 450 |
| Выделение основ слов и лемматизация | 450 |
| Полярность и субъективность настроений..... | 451 |
| От лексем к числам — терм-документная матрица..... | 451 |
| Модель мешка слов | 452 |
| Измерение сходства документов..... | 452 |
| Терм-документная матрица с помощью библиотеки sklearn..... | 453 |
| Применение класса <i>CountVectorizer</i> | 454 |
| Классы <i>TfidfTransformer</i> и <i>TfidfVectorizer</i> | 456 |
| Предобработка текста — краткий обзор | 458 |
| Классификация текста и sentimentный анализ..... | 458 |
| Наивный байесов классификатор..... | 459 |
| Памятка по теореме Байеса..... | 459 |
| Допущение об условной независимости | 460 |
| Классифицирование новостных статей | 461 |
| Тренировка и оценивание мультиномиального наивного байесова классификатора..... | 461 |
| Sentimentный анализ | 462 |
| Данные социальной сети Twitter..... | 462 |
| Отзывы о деятельности предприятий — конкурсные данные веб-сайта Yelp..... | 463 |
| Резюме..... | 466 |
| Глава 14. Тематическое моделирование..... | 467 |
| Усвоение скрытых тем: цели и подходы | 468 |
| От линейной алгебры к иерархическим вероятностным моделям | 469 |
| Латентно-семантическое индексирование | 469 |
| Как реализовать LSI с помощью библиотеки sklearn..... | 471 |
| Сильные и слабые стороны..... | 472 |

| | |
|--|------------|
| Вероятностный латентный семантический анализ..... | 474 |
| Как реализовать pLSA с помощью библиотеки sklearn | 475 |
| Латентное размещение Дирихле | 476 |
| Как работает LDA..... | 476 |
| Распределение Дирихле | 477 |
| Генеративная модель..... | 477 |
| Реконструирование процесса..... | 478 |
| Как оценивать темы LDA..... | 479 |
| Перплексивность | 479 |
| Тематическая когерентность | 479 |
| Как реализовать LDA с помощью библиотеки sklearn..... | 480 |
| Как визуализировать результаты LDA с помощью библиотеки pyLDAvis..... | 481 |
| Как реализовать LDA с помощью библиотеки gensim..... | 482 |
| Тематическое моделирование применительно к телеконференциям о корпоративных зарплатах | 484 |
| Предобработка данных..... | 485 |
| Тренировка и оценивание модели..... | 485 |
| Проведение экспериментов | 487 |
| Тематическое моделирование применительно к отзывам о деятельности предприятий на веб-сайте Yelp | 488 |
| Резюме | 490 |
| Глава 15. Векторное вложение слов..... | 491 |
| Как векторные вложения слов кодируют семантику | 492 |
| Как нейронно-языковые модели усваивают словоупотребление в контексте | 492 |
| Модель Word2vec — усвоение векторных вложений в широком масштабе..... | 493 |
| Модельная цель — упрощенная активационная функция softmax | 494 |
| Автоматическое обнаружение фраз..... | 495 |
| Как оценивать векторные вложения — векторная арифметика и аналогии | 496 |
| Как использовать предварительно натренированные словарные векторы | 497 |
| GloVe — глобальные векторы для представления слов..... | 498 |
| Как тренировать собственные векторные вложения слов | 499 |
| Архитектура модели Skip-Gram в библиотеке Keras..... | 499 |
| Оценивание контрастивное к шуму | 499 |
| Компоненты модели | 500 |
| Словарные векторы из финансовой отчетности SEC с использованием библиотеки gensim..... | 500 |
| Предобработка | 500 |
| Автоматическое обнаружение фраз..... | 501 |
| Тренировка модели..... | 502 |
| Оценивание модели | 502 |
| Влияние параметрических настроек на результативность | 503 |
| Сентиментный анализ с помощью модели Doc2vec..... | 504 |
| Тренировка модели Doc2vec на сентиментных данных Yelp | 504 |
| Создание входных данных | 505 |
| Бонус — модель Word2vec для машинного перевода..... | 508 |
| Резюме | 508 |

| | |
|--|------------|
| Глава 16. Дальнейшие действия..... | 509 |
| Ключевые итоги и извлеченные уроки..... | 510 |
| Данные — единственный самый важный ингредиент | 510 |
| Контроль качества | 510 |
| Интеграция данных | 511 |
| Компетентное знание предметной области помогает разблокировать ценность данных | 511 |
| Выработка признаков и исследование альфа-факторов | 512 |
| Машинное обучение — это комплект инструментов для решения задач с помощью данных | 512 |
| Диагностика модели помогает ускорить оптимизацию | 513 |
| Принятие решений без бесплатного обеда..... | 513 |
| Управление компромиссом между смещением и дисперсией | 514 |
| Определение адресных модельных целей | 514 |
| Верификационная проверка оптимизации | 515 |
| Остерегайтесь переподгонки к историческим данным | 515 |
| Как проникать в сущность черно-ящичных моделей..... | 515 |
| Машинное обучение для торговли на финансовых рынках на практике | 516 |
| Технологии управления данными | 516 |
| Системы управления базами данных..... | 517 |
| Технологии больших данных — Hadoop и Spark | 517 |
| Инструменты машинного обучения..... | 518 |
| Онлайновые торговые платформы | 519 |
| Платформа Quantopian | 519 |
| Платформа QuantConnect..... | 519 |
| Платформа QuantRocket..... | 520 |
| Заключение..... | 521 |
| Глоссарий..... | 522 |
| Предметный указатель..... | 546 |

ГЛАВЫ, ПОМЕЩЕННЫЕ В ЭЛЕКТРОННЫЙ АРХИВ

| | |
|--|----------|
| Глава 17. Глубокое обучение..... | 1 |
| Глубокое обучение и искусственный интеллект | 2 |
| Сложности высокоразмерных данных..... | 3 |
| Глубокое обучение как усвоение представлений | 4 |
| Как глубокое обучение извлекает иерархические признаки из данных | 5 |
| Аппроксимация универсальных функций..... | 6 |
| Глубокое обучение и усвоение проекций в топологическом многообразии..... | 6 |
| Как глубокое обучение соотносится с машинным обучением и искусственным интеллектом | 7 |
| Как конструировать нейронную сеть..... | 8 |
| Как работают нейронные сети..... | 9 |
| Простая архитектура сети прямого распространения | 10 |

| | |
|---|-----------|
| Ключевые варианты выбора конструкции | 11 |
| Стоимостные функции | 12 |
| Выходные элементы | 12 |
| Скрытые элементы | 13 |
| Как регуляризовать глубокие нейронные сети | 14 |
| Штрафы по норме параметров | 14 |
| Досрочная остановка | 15 |
| Отсев | 15 |
| Оптимизация для глубокого обучения | 15 |
| Стохастический градиентный спуск | 16 |
| Импульс | 17 |
| Адаптивные темпы усвоения | 17 |
| Как строить нейронную сеть на языке Python | 18 |
| Входной слой | 18 |
| Скрытый слой | 19 |
| Выходной слой | 20 |
| Прямое распространение | 21 |
| Перекрестно-энтропийная стоимостная функция | 21 |
| Как тренировать нейронную сеть | 22 |
| Как реализовать обратное распространение с помощью Python | 22 |
| Как вычислять градиент | 23 |
| Градиент функции потери | 23 |
| Градиенты выходного слоя | 24 |
| Градиенты скрытого слоя | 24 |
| Собирая все вместе | 25 |
| Тренировка сети | 26 |
| Как применять библиотеки глубокого обучения | 28 |
| Как применять библиотеку Keras | 28 |
| Как применять инструмент визуализации TensorBoard | 30 |
| Как применять библиотеку PyTorch 1.0 | 31 |
| Как создать загрузчик данных в библиотеке PyTorch | 32 |
| Как формировать нейросетевую архитектуру | 33 |
| Как тренировать модель | 34 |
| Как оценивать модельные предсказания | 35 |
| Как применять библиотеку TensorFlow 2.0 | 35 |
| Как оптимизировать нейросетевые архитектуры | 35 |
| Создание временного ряда финансовых возвратов от акции для предсказания ценового движения актива | 36 |
| Определение нейросетевой архитектуры с заполнителями | 37 |
| Определение кастомизированной метрики потери для досрочной остановки | 38 |
| Выполнение GridSearchCV для настройки нейросетевой архитектуры | 38 |
| Как еще больше улучшить результаты | 40 |
| Резюме | 40 |
| Глава 18. Сверточные нейронные сети | 41 |
| Как работают сети ConvNet | 42 |
| Как работает сверточный слой | 43 |
| Сверточный этап — обнаружение локальных признаков | 45 |

| | |
|--|-----------|
| Детекторный этап — добавление нелинейности | 47 |
| Сводный этап — понижающая дискретизация признакововой карты | 47 |
| Источник вдохновения из нейробиологии | 48 |
| Опорные архитектуры сети ConvNet | 49 |
| Сеть LeNet5 — первая современная CNN-сеть (1998) | 49 |
| Сеть AlexNet — закрепление CNN-сетей на технологической карте (2012) | 50 |
| Сеть VGGNet — ориентация на фильтры меньших размеров | 51 |
| Сеть GoogLeNet — меньше параметров посредством модуля Inception | 52 |
| Сеть ResNet — современное состояние дел | 53 |
| Эталоны | 53 |
| Извлеченные уроки | 54 |
| Компьютерное зрение за пределами классификации — обнаружение и сегментация | 54 |
| Как конструировать и тренировать CNN-сеть с помощью языка Python | 55 |
| Сеть LeNet5 и набор данных MNIST с библиотекой Keras | 55 |
| Как готовить данные | 55 |
| Как определить архитектуру | 56 |
| Сеть AlexNet и набор данных CIFAR10 с библиотекой Keras | 58 |
| Как готовить данные с помощью аугментации снимка | 58 |
| Как формировать архитектуру модели | 59 |
| Как использовать CNN-сеть с данными временного ряда | 60 |
| Трансферное обучение — более быстрая тренировка с меньшим количеством данных | 62 |
| Как надстраивать над предварительно натренированной CNN-сети | 62 |
| Как извлекать бутылочные признаки | 63 |
| Как продолжить дальнейшую тренировку предварительно натренированной модели | 64 |
| Как обнаруживать объекты | 66 |
| Набор данных номеров домов Google Street View | 66 |
| Как определять CNN-сеть с многочисленными выходами | 67 |
| Новейшие разработки | 68 |
| Быстрое обнаружение объектов на спутниковых снимках | 68 |
| Как капсульные сети улавливают позу | 68 |
| Резюме | 68 |
| Глава 19. Рекуррентные нейронные сети | 70 |
| Как работают RNN-сети | 71 |
| Развертывание вычислительного графа с циклами | 73 |
| Обратное распространение во времени | 74 |
| Альтернативные архитектуры RRN-сетей | 74 |
| Выходная рекуррентия и вмешательство учителя | 74 |
| Двунаправленные RNN-сети | 75 |
| Кодировочно-декодировочные архитектуры и механизм внимания | 75 |
| Как конструировать глубокие RNN-сети | 76 |
| Сложности усвоения долгосрочных зависимостей | 77 |
| Элементы с долгой краткосрочной памятью | 77 |
| Вентильные рекуррентные элементы | 79 |
| Как строить и тренировать RNN-сети с помощью языка Python | 79 |
| Регрессия одномерного временного ряда | 80 |
| Как привести данные временного ряда в форму для RNN-сети | 80 |
| Как определять двухслойную RNN-сеть, используя единственный слой LSTM | 82 |

| | |
|---|-----------|
| Уложенные друг на друга элементы LSTM для классификации временного ряда | 84 |
| Как готовить данные | 84 |
| Как формировать архитектуру | 86 |
| Регрессия многомерного временного ряда..... | 88 |
| Загрузка данных..... | 88 |
| Подготовка данных..... | 88 |
| Определение и тренировка модели | 89 |
| Элементы LSTM и векторные вложения слов для классификации сентимента | 91 |
| Загрузка данных IMDB с отзывами о кинофильмах..... | 91 |
| Определение архитектур вложения и RNN-сети | 92 |
| Сентиментный анализ с предварительно натренированными словарными векторами | 93 |
| Предобработка текстовых данных | 94 |
| Загрузка предварительно натренированных вложений GloVe | 94 |
| Резюме | 95 |
| Глава 20. Автокодировщики и генеративные состязательные сети..... | 97 |
| Как работают автокодировщики | 98 |
| Нелинейное снижение размерности..... | 99 |
| Сверточные автокодировщики | 100 |
| Ограничения по разреженности с регуляризованными автокодировщиками..... | 100 |
| Исправление поврежденных данных с помощью шумоподавляющих автокодировщиков | 101 |
| Автокодировщики с отображением одной последовательности в другую | 101 |
| Вариационные автокодировщики | 102 |
| Конструирование и тренировка автокодировщиков с использованием языка Python | 102 |
| Подготовка данных..... | 103 |
| Однослойный автокодировщик прямого распространения | 104 |
| Определение кодировщика | 105 |
| Определение декодировщика | 105 |
| Тренировка модели..... | 106 |
| Оценивание результатов | 106 |
| Автокодировщик прямого распространения с ограничениями по разреженности | 107 |
| Автокодировщик глубокого прямого распространения..... | 107 |
| Визуализация кодирования..... | 108 |
| Сверточные автокодировщики | 109 |
| Шумоподавляющие автокодировщики..... | 110 |
| Как работают GAN-сети | 111 |
| В чем различие между генеративными и дискриминативными моделями | 111 |
| Как работает состязательная тренировка | 112 |
| Как эволюционируют архитектуры GAN-сетей | 113 |
| Глубокая сверточная GAN-сеть (DCGAN)..... | 113 |
| Условные GAN-сети..... | 114 |
| Успешные и быстро развивающиеся приложения GAN-сетей | 114 |
| Сеть CycleGAN — неспаренная трансляция из изображения в изображение | 114 |
| Сеть StackGAN — синтез изображений "текст в фото" | 114 |
| Сверхразрешающая способность фотореалистичных изображений..... | 115 |
| Синтетические временные ряды с помощью рекуррентных cGAN-сетей | 115 |

| | |
|---|------------|
| Построение GAN-сетей с использованием языка Python..... | 115 |
| Определение дискриминаторной сети..... | 116 |
| Определение генераторной сети..... | 116 |
| Совмещение обеих сетей для определения GAN-сети..... | 117 |
| Состязательная тренировка..... | 117 |
| Оценивание результатов..... | 118 |
| Резюме..... | 118 |
| Глава 21. Подкрепляемое обучение..... | 120 |
| Ключевые элементы подкрепляемого обучения..... | 121 |
| Компоненты интерактивной системы подкрепляемого обучения..... | 122 |
| Политика — от состояний к действиям..... | 122 |
| Вознаграждения — самообучение на основе действий..... | 123 |
| Ценностная функция — хорошие решения в долгосрочной перспективе..... | 123 |
| Модельные агенты против безмодельных..... | 124 |
| Как решать задачи подкрепляемого обучения..... | 124 |
| Ключевые сложности в решении задач подкрепляемого обучения..... | 125 |
| Учет заслуг..... | 125 |
| Разведывание против эксплуатации..... | 125 |
| Фундаментальные подходы к решению задач подкрепляемого обучения..... | 125 |
| Динамическое программирование — цикл по ценностным значениям и цикл по политике..... | 127 |
| Конечные задачи марковского процесса принятия решений..... | 127 |
| Последовательности состояний, действий и вознаграждений..... | 127 |
| Ценностные функции — как оценивать долгосрочное вознаграждение..... | 128 |
| Уравнение Беллмана..... | 129 |
| От ценностной функции к оптимальной политике..... | 129 |
| Цикл по политике..... | 130 |
| Цикл по ценностным значениям..... | 131 |
| Обобщенный цикл по политике..... | 132 |
| Динамическое программирование на языке Python..... | 132 |
| Настройка решетчатого мира GridWorld..... | 133 |
| Вычисление матрицы переходов..... | 135 |
| Цикл по ценностным значениям..... | 136 |
| Цикл по политике..... | 137 |
| Решение задач марковского процесса принятия решений с помощью библиотеки rumpdptoolbox..... | 137 |
| Выводы..... | 138 |
| Q -обучение..... | 138 |
| Компромисс между разведыванием и эксплуатацией — ε -жадная политика..... | 139 |
| Алгоритм Q -обучения..... | 139 |
| Тренировка агента Q -обучения с помощью языка Python..... | 139 |
| Глубокое подкрепляемое обучение..... | 140 |
| Аппроксимация ценностной функции с помощью нейронной сети..... | 141 |
| Алгоритм глубокого Q -обучения и его расширения..... | 141 |
| Воспроизведение опыта..... | 142 |
| Медленно изменяющаяся целевая сеть..... | 142 |
| Двойное глубокое Q -обучение..... | 143 |

| | |
|---|------------|
| Платформа OpenAI Gym — лунный посадочный модуль | 143 |
| Сеть DDQN с использованием библиотеки TensorFlow | 144 |
| Архитектура сети DQN | 144 |
| Настройка среды платформы OpenAI | 145 |
| Гиперпараметры | 145 |
| Вычислительный граф сети DDQN | 146 |
| Результативность | 148 |
| Подкрепляемое обучение для торговли на финансовых рынках | 148 |
| Конструирование торговой среды в платформе OpenAI | 148 |
| Элементарная торговая игра | 149 |
| Результативность глубокого Q -обучения на фондовом рынке | 150 |
| Резюме | 151 |
| Приложение 1. Настройка среды программирования на языке Python | 152 |
| Настройка дистрибутива Anaconda | 152 |
| Создание среды | 153 |
| Установка программных библиотек | 155 |
| Установка библиотеки TensorFlow | 158 |
| Блокноты Jupyter | 159 |
| Приложение 2. Скачивание и подготовка данных | 162 |
| Цены Quandl Wiki | 162 |
| Метаданные цен Wiki | 163 |
| Цены фондового индекса S&P 500 | 163 |
| Компонентные ценные бумаги фондового индекса S&P 500 | 164 |
| Метаданные по торгуемым компаниям США | 164 |
| Индексы вексельных цен | 165 |

Об авторе

Стефан Янсен (Stefan Jansen), дипломированный финансовый аналитик, является основателем и ведущим исследователем данных в компании Applied AI (<http://www.applied-ai.solutions/>), где он консультирует компании и стартапы списка Fortune 500 из разных отраслей промышленности по переводу деловых целей в стратегию данных и искусственного интеллекта, комплектует команды исследователей данных и разрабатывает решения на основе автоматизированного (машинного) обучения. До своей нынешней должности он был управляющим партнером и ведущим исследователем данных в международной инвестиционной фирме, в которой организовал и внедрил применение прогнозной аналитики и инвестиционных исследований в практику.

Он также был руководителем глобального финансово-технологического стартапа, работающего на 15 рынках, работал во Всемирном банке, консультировал центральные банки на развивающихся рынках и вел проекты на шести языках на четырех континентах. Стефан имеет степень магистра Гарвардского и Берлинского университетов и преподает науку о данных в частной школе General Assembly (<https://generalassemb.ly/>) и на платформе онлайн-обучения Datacamp (<https://www.datacamp.com/>).

"Благодарю Packt за эту возможность и команду, которая воплотила ее в конечный результат, в особенности хочу поблагодарить Снегил Колте (Snehal Kolte) за оказанное содействие в процессе редактирования. Многие мои сторонники заслуживают упоминания, но профессор Цекхаузер (Zeckhauser) из Гарварда выделяется вдохновляющим интересом к творческому использованию количественных методов для решения задач. Я в долгу перед родителями за то, что они поощряли мое любопытство и поддерживали меня. Но больше всего я благодарен Мариане за то, что она делает все это стоящим того".

Стефан Янсен

О рецензентах

Дуг Ортис (Doug Ortiz) является опытным архитектором корпоративного облака, больших данных, анализа данных и решений, который проектировал, конструировал, разрабатывал, реконструировал и интегрировал корпоративные решения. Его знания также охватывают облачные платформы Amazon Web Services, Azure, Google Cloud, деловую аналитику, базы данных Hadoop, Spark, NoSQL и SharePoint. Он является основателем проекта астрофизического моделирования Illustris.

"Огромное спасибо моей замечательной жене Милле, а также Марии, Николаю и нашим детям за их поддержку".

Дуг Ортис

Сандипан Дей (Sandipan Dey) является исследователем данных с широким спектром интересов, включая такие темы, как машинное обучение, глубокое обучение, обработка изображений и компьютерное зрение. Он работал во многих областях науки о данных, включая рекомендательные системы, предсказательные модели для индустрии событий, модели датчиковой локализации, сентиментный анализ и машинную прогностику. Он получил степень магистра в области компьютерных наук в Университете Мэриленда, округ Балтимор, его статьи были опубликованы в бюллетенях нескольких конференций и в журналах IEEE по глубинному анализу данных.

Он имеет сертификаты более 100 массовых открытых курсов дистанционного обучения (МООС) по науке о данных, машинному обучению, глубокому обучению, обработке изображений и родственным курсам/специализациям. Регулярно обновляет свой блог (sandipanweb) и является энтузиастом машинного обучения.

Комментарии переводчика

В эпоху всеобщей цифровизации и сетевого взаимодействия сфера биржевой торговли становится все более демократичной, позволяя любому попробовать в ней свои силы, используя многочисленные открытые источники данных и торговые платформы. Эта книга является одной из первых на русском языке, полностью посвященной алгоритмической торговле на финансовых рынках с использованием технологии машинного обучения, и она безусловно вызовет интерес как у профессионалов в области биржевой торговли, так и у специалистов по машинному обучению и, возможно, станет настольной у каждого кванта и разработчика автоматически обучающихся систем.

Книга написана высококлассным профессионалом и для профессионалов. Автор подробно разбирает все магистральные и новейшие технические решения машинного обучения применительно к торговле на финансовых рынках, предоставляя обширный справочный материал и примеры использования в своем репозитории на GitHub. Постоянно обновляемая кодовая база книги изначально была создана с использованием Ubuntu 18.04. Ее репозиторий содержит файлы операционных сред Ubuntu и Mac. Дальнейшие обновления репозитория будут содержать адаптированные версии для Windows. При переводе книги исходный код был частично протестирован в среде Windows 10. При тестировании исходного кода за основу взят Python версии 3.7.2. Время перевода — апрель-май 2019 г.

С целью расширить аудиторию книги за счет разработчиков обучающихся систем и, наоборот, облегчить биржевым торговцам работу с технологией машинного обучения настоящий перевод снабжен сносками в виде комментариев и определений терминов из финансов, машинного обучения, теории вероятностей и статистики, которые облегчат понимание материала и помогут "зацепиться" за тему. Все терминологические сноски в конце книги сведены в глоссарий основных терминов. Почти для всех терминов приведены ссылки на источник. Перевод книги содержит два приложения, посвященные настройке среды программирования на языке Python и подготовке данных, используемых на протяжении всей книги (см. электронный архив к книге по адресу <ftp://ftp.bhv.ru/9785977565950.zip> или на странице книги на сайте издательства www.bhv.ru).

О терминологии

Финансовый возврат и возвратность. Одним из ключевых понятий биржевой торговли является *доходность*, или, если точнее, *финансовый возврат* (return) на инвестированный капитал. С точки зрения биржевого торговца (трейдера), если представить ценовую информацию в виде барного временного графика, то финансовый возврат — это не что иное, как разница между ценой актива в момент времени $t-1$ и ценой в момент времени t ($\text{return}[t] - \text{return}[t-1]$). Если торговец сделал ставку в момент времени $t-1$, то он получает положительный, нулевой или отрицательный финансовый возврат в момент времени t . С точки зрения портфельного менеджера, это возврат на портфельные инвестиции в течение любого периода оценивания, включающего изменение рыночной стоимости портфеля. Процентное соотношение этого показателя называется *возвратностью* (rate of return, $(\text{return}[t] - \text{return}[t-1]) / \text{return}[t-1]$). Цель машинного обучения в финансах состоит в том, чтобы натренировать автоматически обучающуюся систему предсказывать состояние рынка и цену актива/стоимость портфеля в момент времени $t+1$ с целью достижения максимальной возвратности.

Результативность. В англоязычной литературе термин "результативность" (performance) широко применяется как в финансах, так и в машинном обучении, в особенности в финансах, где самое главное — это финансовый результат, получаемый от актива или инвестиционного портфеля. Данный термин ассоциируется не с созданием чего-либо за единицу времени (производительность, productivity) или эффекта (эффективность, effectiveness), а с созданием результата, который затем сравнивается с неким эталоном. Например, результативность актива (asset performance) означает способность предприятия брать операционные ресурсы, управлять ими и в результате получать прибыльную возвратность, превышающую индекс или кредитную ставку. Результативность в МО означает некое достижение по исполнении задачи, измеряемое общепринятыми стандартными мерами точности, полноты, стоимости и скорости. Хотя этот термин синонимичен термину "эффективность", однако эффективность имеет собственное слово и не раз встречается в этой книге отдельно, а термин "результативность" контрастнее подчеркивает нацеленность на результат и поэтому предпочтителен.

Терминология машинного обучения. В центре внимания машинного обучения и его подобласти, глубокого обучения, находится автоматически обучающаяся система, т. е. система, способная с течением времени приобретать новые знания и улучшать свою работу, используя поступающую информацию. В зарубежной специализированной литературе для *передачи* знаний ученику и *получения* знаний учеником существуют отдельные термины — train (*натренировать*) и learn (*выучить*, усвоить), где тренировка — это работа, которую выполняет исследователь-проектировщик для получения обучившейся модели, в основе которой лежит обучающийся алгоритм, по сути искатель минимумов (или максимумов) для надлежащим образом сформулированных функций, а усвоение, или самообучение, — это работа, которую выполняет алгоритм-ученик по заучиванию связей и регулярно-

стей в данных или изменению и закреплению своего поведения. Когда же для обоих английских терминов в русской спецлитературе используется один-единственный термин "обучение", то он несет в себе двусмысленность, потому что под ним может подразумеваться и передача знаний, и получение знаний одновременно, как, например, в случае с термином "машинное обучение", который может означать и тренировку алгоритмических машин, и способность таких машин автоматически обучаться, что нередко вносит путаницу и терминологический разброд в переводной литературе при решении дилеммы "training-learning" в то время, как появление в зарубежной технической литературе термина learning в любом виде подразумевает исключительно второе — *самообучение, заучивание, усвоение* алгоритмом регулярностей или параметров. Отсюда вытекает одно важное следствие: английский термин machine learning обозначает *усвоение знаний алгоритмической машиной*, а, следовательно, более соответствовать оригиналу будет термин "*машинное самообучение*" или "*автоматическое обучение*". Весомым аргументом в пользу этих вариантов термина является и то, что с начала 1960-х и до середины 1980-х годов в ходу был похожий термин — "обучающиеся машины" (см. работы А. Тьюринга, К. Шеннона, Н. Винера, Н. Нильсона, Я. З. Цыпкина и др.). В настоящем переводе, следуя принципам здравого смысла и бритвы Оккама, за основу принят зарубежный подход.

<https://ru.tradingview.com/u/capissimo/#published-scripts>

Предисловие

Наличие и доступность разнообразных данных повысила спрос на компетентные знания в области стратегий алгоритмической торговли. Благодаря этой книге вы выберете машинное обучение (МО), будете его применять к широкому спектру источников данных и создавать мощные алгоритмические стратегии.

Книга начинается с введения существенных элементов, таких как оценивание наборов данных, доступ к данным через API с помощью Python, использование платформы Quandl для доступа к финансовым данным и управление ошибками предсказания. Далее мы рассмотрим различные технические решения машинного обучения и автоматически обучающиеся алгоритмы, которые могут использоваться для построения и тренировки алгоритмических моделей с помощью программных Python-библиотек pandas, Seaborn, StatsModels и sklearn. Потом мы построим, оценим и дадим интерпретацию моделей $AR(p)$, $MA(q)$ и $ARIMA(p, d, q)$ с использованием библиотеки StatsModels. Вы примените байесовы понятия "априорное распределение", "наблюдение" и "апостериорное распределение" для того, чтобы различать понятие неопределенности с помощью библиотеки PyMC3. Затем мы задействуем библиотеки NLTK, sklearn (Scikit-learn) и spaCy для назначения sentimentных меток финансовым новостям и классифицирования документов для извлечения торговых сигналов. Мы научимся конструировать, строить, настраивать и оценивать нейронные сети прямого распространения, рекуррентные нейронные сети (RNN-сети) и сверточные нейронные сети (CNN-сети), используя библиотеку Keras для разработки изошренных алгоритмов. Вы примените трансферное обучение к данным спутниковых снимков для предсказания экономической активности. Наконец, мы применим подкрепляемое обучение для достижения оптимальных результатов торговли.

Прочитав книгу до конца, вы будете способны применять алгоритмическую торговлю, реализуя умные инвестиционные стратегии.

Для кого эта книга предназначена

Книга предназначена для аналитиков и исследователей данных и разработчиков на языке Python, а также инвестиционных аналитиков и портфельных менеджеров, работающих внутри финансово-инвестиционной индустрии. Если вы хотите реализовать эффективную алгоритмическую торговлю, разрабатывая интеллектуальные разведывающие стратегии с использованием автоматически обучающихся алгоритмов, то настоящая книга — именно то, что вам нужно! Некоторое понимание приемов языка Python и технических решений МО является обязательным.

Что эта книга охватывает

Глава 1 "Машинное обучение для торговли на финансовых рынках" отождествляет центральную тему данной книги, обрисовывая роль машинного обучения в формировании и оценивании сигналов для разработки и исполнения торговой стратегии. В ней описывается стратегический процесс от генерирования и моделирования гипотез, отбора данных и бэктестирования до оценивания и исполнения стратегий в портфельном контексте, включая риск-менеджмент.

Глава 2 "Рыночные и фундаментальные данные" охватывает источники данных и работу с исходными тиковыми данными, поступающими с бирж, и данными финансовой отчетности, а также способы доступа к многочисленным поставщикам общедоступных данных, на которые мы будем опираться на протяжении всего повествования в этой книге.

Глава 3 "Альтернативные данные для финансов" предоставляет категории и критерии для оценивания качества стремительно растущего числа источников и поставщиков. В ней также демонстрируются способы создания альтернативных наборов данных путем выскабливания веб-сайтов, например для сбора стенограмм телеконференций о корпоративных заработках для использования в алгоритмах естественно-языковой обработки и сентиментного анализа во второй условной части книги.

Глава 4 "Исследование альфа-факторов" обеспечивает каркас для понимания того, как работают факторы и как измерять их эффективность, например, используя информационный коэффициент (information coefficient, IC). В ней демонстрируется способ конструирования альфа-факторов из данных с помощью библиотек Python в режиме офлайн и на веб-платформе Quantopian. В ней также вводятся библиотека *zipline* для бэктестирования факторов и библиотека *alphalens* для оценивания их предсказательной мощности.

Глава 5 "Оценивание стратегии" знакомит с тем, как строить, тестировать и оценивать стратегии торговли, используя исторические данные с помощью библиотеки *zipline* в режиме офлайн и на веб-платформе Quantopian. Она представляет и демонстрирует способы вычисления результативности портфеля и метрических показателей риска с помощью библиотеки *ryfolio*. В ней также рассматриваются спосо-

бы управления методологическими сложностями бэктестирования стратегий и вводятся методы оптимизации стратегии с точки зрения портфельного риска.

Глава 6 "Процесс машинного обучения" закладывает фундамент, давая описания того, как формулировать, тренировать и настраивать автоматически обучающиеся модели и оценивать их предсказательную результативность как систематический рабочий процесс.

Глава 7 "Линейные модели" показывает, как использовать линейную и логистическую регрессию для выведения заключения и предсказания и применять регуляризацию для управления риском переподргонки. Она знакомит с торговой платформой Quantopian и демонстрирует способы построения факторных моделей и предсказания цен на финансовые активы.

Глава 8 "Модели временных рядов" посвящена одномерным и многомерным временным рядам, включая векторные авторегрессионные модели и коинтеграционные тесты, а также способы их применения к стратегиям парной торговли.

Глава 9 "Байесово машинное обучение" знакомит с тем, как формулировать вероятностные модели и как извлечение выборок из вероятностных распределений методами Монте-Карло марковских цепей (Markov Chain Monte Carlo, MCMC) и вариационный Байес обеспечивают аппроксимирование вывода. В ней также иллюстрируются способы применения библиотеки PyMC3 для вероятностного программирования с целью понимания сущности параметрической и модельной неопределенности.

Глава 10 "Деревья решений и случайные леса" показывает, как строить, тренировать и настраивать нелинейные древесные модели для проникновения в сущность данных и предсказания. В ней представлены древесные ансамблевые модели и иллюстрируется то, как случайные леса используют агрегирование бутстраповских выборок для преодоления некоторых недостатков деревьев решений.

Глава 11 "Градиентно-бустинговые машины" демонстрирует использование библиотек XGBoost, LightGBM и CatBoost для выполнения высокорезультативной тренировки моделей и предсказания. Кроме того, в ней подробно анализируется настройка многочисленных гиперпараметров.

Глава 12 "Неконтролируемое обучение" посвящена применению задач снижения размерности и кластеризации в алгоритмической торговле. В ней используется анализ главных и независимых компонент для извлечения ведомых данными рисков факторов. В ней также представлено несколько технических решений для кластеризации и показано использование иерархической кластеризации для размещения финансовых средств среди портфельных активов.

Глава 13 "Работа с текстовыми данными" демонстрирует способы конвертирования текстовых данных в числовой формат. В ней также показывается применение классификационных алгоритмов из второй условной части, связанных с сентиментным анализом, к крупным наборам данных.

Глава 14 "Тематическое моделирование" посвящена применению байесова неконтролируемого обучения с целью извлечения латентных тем, которые способны ре-

зюмировать большое число документов и предлагать более эффективные способы разведывания текстовых данных либо использования тем в качестве признаков в классификационной модели. В ней демонстрируется способ применения этого технического решения к стенограммам телеконференций о корпоративных заработках, источники которых привлечены в *главе 3*, и к годовым отчетам, предъявляемым в Комиссию по ценным бумагам и биржам (Securities and Exchange Commission, SEC).

Глава 15 "Векторное вложение слов" посвящена использованию нейронных сетей для усвоения передовых языковых признаков в форме словарных векторов, которые улавливают семантический контекст намного лучше, чем традиционные текстовые признаки, и представляют собой очень перспективное направление для извлечения торговых сигналов из текстовых данных.

Глава 16 "Последующие шаги" представляет собой резюме всех предыдущих глав.

Глава 17 "Глубокое обучение" знакомит с библиотеками Keras, TensorFlow и PyTorch, наиболее популярными каркасами глубокого обучения, которые мы будем использовать на протяжении всей четвертой условной части. В ней также представлены технические решения для тренировки и настройки, включая регуляризацию, и дан краткий обзор общепринятых архитектур.

Глава 18 "Сверточные нейронные сети" охватывает CNN-сети, очень мощные сетевые архитектуры, предназначенные для классификационных задач с широкомасштабными неструктурированными данными. Мы введем успешные архитектурные проекты, натренируем CNN-сеть на спутниковых данных, например для предсказания экономической активности, и применим трансферное обучение для ускорения тренировочного процесса.

Глава 19 "Рекуррентные нейронные сети" показывает, как RNN-сети могут быть полезными для моделирования отображений последовательности в последовательность, в том числе для временных рядов. В ней демонстрируется, как RNN-сеть улавливает нелинейные регуляризации за более продолжительные периоды.

Глава 20 "Автокодировщики и генеративные состязательные сети" обращается к вопросам неконтролируемого глубокого обучения, включая автокодировщики для нелинейного сжатия многомерных данных и генеративные состязательные сети (Generative Adversarial Networks, GAN) — одно из самых важных последних нововведений, которое применяется для генерирования синтетических данных.

Глава 21 "Подкрепляемое обучение" знакомит с автоматическим обучением на основе максимизации вознаграждения за принимаемое решение, позволяющее конструировать и тренировать агентов, которые учатся оптимизировать решения с течением времени в ответ на окружающую их среду. Вы увидите процесс построения агента, который откликается на рыночные сигналы, с использованием библиотеки Open AI Gym (<https://gym.openai.com>), так называемого "тренажерного зала искусственного интеллекта".



Главы 17–21 (а также приложения и цветные иллюстрации к главам 1–15) представлены в электронном архиве данной книги, размещенном на сайте издательства. Электронный архив можно скачать по ссылке <ftp://ftp.bhv.ru/9785977565950.zip> и со страницы книги на веб-сайте издательства по адресу www.bhv.ru.

Получение максимальной отдачи от этой книги

Для чтения данной книги требуется только базовое знание языка Python и элементарное понимание технических решений машинного обучения.

Скачивание файлов с примерами исходного кода¹

Файлы с примерами можно скачать с вашего аккаунта по адресу <http://www.packtpub.com/> для всех книг издательства Packt Publishing, которые вы приобрели. Если вы купили эту книгу в другом месте, то можно посетить <http://www.packtpub.com/support> и зарегистрироваться там, чтобы получить файлы прямо по электронной почте.

Скачать файлы с примерами можно, выполнив следующие шаги:

1. Войдите на наш веб-сайт или зарегистрируйтесь там, используя ваш адрес электронной почты и пароль.
2. Наведите указатель мыши на вкладку **SUPPORT** вверху страницы.
3. Щелкните по разделу **Code Downloads & Errata**, посвященному примерам программного кода и опечаткам.
4. Введите название книги в поле поиска.

Скачав файл, пожалуйста, убедитесь, что вы распаковали или извлекли папку, воспользовавшись последней версией указанных ниже архиваторов:

- ◆ WinRAR/7-Zip для Windows;
- ◆ Zipeg/iZip/UnRarX для Mac OS;
- ◆ 7-Zip/PeaZip для Linux.

Помимо этого, комплект примеров программного кода, прилагаемый к данной книге, размещен на GitHub в разделе **Packt** по адресу <https://github.com/PacktPublishing/Hands-On-Machine-Learning-for-Algorithmic-Trading>. В случае обновления программного кода он будет обновлен в существующем репозитории GitHub.

¹ Данный раздел относится к электронным архивам оригинального издания книги Stefan Jansen "Hands-On Machine Learning for Algorithmic Trading". Описание электронного архива к русскому изданию см. в конце книги. — *Прим. ред.*

Мы также располагаем другими комплектами примеров программного кода, которые можно выбрать из нашего богатого каталога книг и видеороликов, предлагаемого на странице <https://github.com/PacktPublishing/>. Можете убедиться сами!

Скачивание цветных изображений

Мы также предоставляем PDF-файл с цветными изображениями скриншотов/диаграмм, используемых в этой книге. Вы можете скачать его здесь: https://static.packt-cdn.com/downloads/9781789346411_ColorImages.pdf.²

Принятые в книге условные обозначения

В этой книге используется ряд текстовых условных обозначений.

- ◆ Код в тексте обозначает кодовые слова в тексте, имена таблиц базы данных, ввод данных пользователем и дескрипторы Twitter. Приведем пример: "Алгоритм продолжает исполняться после вызова функции `run_algorithm()` и возвращает тот же кадр данных с бэктекстовой результативностью".
- ◆ Блок кода выглядит следующим образом:


```
interesting_times = extract_interesting_date_ranges(returns=returns)
(interesting_times['Fall2015']).to_frame('pf')
.join(benchmark_rets)
.add(1).cumprod().sub(1)
.plot(lw=2, figsize=(14, 6),
      title='Паника после голосования по Брекситу'))
```
- ◆ **Полужирный шрифт** обозначает элемент интерфейса — команду меню, кнопку и др.
- ◆ *Курсивом* выделены новые термины, важные слова или слова, отображаемые на экране.
- ◆ Узким шрифтом выделены имена файлов, папок.



Данный элемент обозначает подсказку или совет.



Данный элемент обозначает предупреждение или предостережение.

² PDF-файл с цветными изображениями к русскому изданию книги представлен в электронном архиве, который доступен по ссылке <ftp://ftp.bhv.ru/9785977565950.zip> и со страницы книги на веб-сайте издательства по адресу www.bhv.ru. — Прим. ред.

1

Машинное обучение для торговли на финансовых рынках

Алгоритмическая торговля на финансовых рынках опирается на компьютерные программы, которые исполняют алгоритмы, автоматизирующие некоторые или все элементы торговой стратегии. Алгоритмы представляют собой последовательность шагов или правил по достижению цели и могут принимать множество форм. В случае *машинного обучения* (МО — machine learning, ML) эти алгоритмы ставят своей задачей усвоить другие алгоритмы, а именно правила по достижению цели на основе данных, такой как минимизация ошибки предсказания.

Указанные алгоритмы кодируют различные действия портфельного менеджера, который наблюдает за рыночными транзакциями и анализирует соответствующие данные, принимая решения о размещении ордеров на покупку или продажу. Последовательность ордеров определяет элементы содержимого инвестиционного портфеля, которые со временем призваны произвести финансовые возвраты на инвестированный капитал, привлекательные для поставщиков капитала, с учетом их аппетита к риску.

В конечном счете цель активного инвестиционного менеджмента заключается в достижении альфы, т. е. финансовых возвратов, превышающих эталонный показатель, используемый для оценивания. Фундаментальный закон активного инвестиционного менеджмента применяет *информационное соотношение* (information ratio, IR), которым выражается стоимость активного менеджмента как отношение портфельных финансовых возвратов сверх возвратов эталонного показателя, обычно индекса, к волатильности этих возвратов. Он аппроксимирует это информационное соотношение как произведение *информационного коэффициента* (information coefficient, IC), которым измеряется качество прогноза, как корреляция возвратов с результатами, и широты стратегии, выраженной как квадратный корень из числа ставок.

Следовательно, ключом к генерированию альфы является прогнозирование. Успешные предсказания, в свою очередь, требуют опережающей информации или превосходящей способности обрабатывать публичную информацию. Алгоритмы обеспечивают оптимизацию на протяжении всего инвестиционного процесса, от размещения финансовых средств среди портфельных активов до генерирования идей,

исполнения сделок и риск-менеджмента. Использование МО для алгоритмической торговли, в частности, нацелено на более эффективное использование обычных и альтернативных данных с целью порождения как более качественных, так и более действенных прогнозов, тем самым улучшая добавочную стоимость активного менеджмента.

Исторически алгоритмическая торговля на финансовых рынках обычно более узко определялась как автоматизация исполнения торговли с целью минимизации затрат, предлагаемых стороной продажи, но мы рассмотрим более полную перспективу, поскольку использование алгоритмов и, в частности, машинного обучения повлияло на более широкий спектр деятельности от генерирования идей и конструирования альфа-факторов до размещения финансовых средств среди портфельных активов, калибровки позиций и тестирования и оценивания стратегий.

В этой главе рассматривается более широкая картина того, как использование МО стало критически важным источником конкурентного преимущества в инвестиционной индустрии и где именно оно вписывается в инвестиционный процесс, обеспечивая условия для стратегий алгоритмической торговли.

В этой главе мы рассмотрим следующие темы:

- ◆ как организована эта книга и кто должен ее читать;
- ◆ как МО стало играть стратегическую роль в алгоритмической торговле;
- ◆ как конструировать и исполнять торговую стратегию;
- ◆ как МО наращивает добавочную стоимость в стратегии алгоритмической торговли.

Как читать эту книгу

Если вы читаете эти строки, то, вероятно, знаете, что МО стало стратегическим потенциалом во многих отраслях промышленности, включая инвестиционную индустрию. Взрывной рост цифровых данных, который стимулирует значительную долю роста МО, оказывает особенно сильное влияние на инвестиции, которые уже имеют долгую историю использования сложных моделей по обработке информации. Из размаха торговли по всем классам финансовых активов следует, что в дополнение к рыночным и фундаментальным данным, которые раньше были в центре внимания аналитических усилий, становится актуальным широкий спектр новых альтернативных данных.

Возможно, вы также столкнулись с пониманием того, что успешное применение МО или науки о данных требует интеграции статистических знаний, вычислительных навыков и компетентного знания предметной области на индивидуальном или командном уровне. Другими словами, важно задавать правильные вопросы, выявлять и понимать данные, которые способны обеспечивать ответы, разворачивать широкий спектр инструментов для получения результатов и интерпретировать их таким образом, который приводит к правильным решениям.

Следовательно, в этой книге рассматривается комплексный подход к применению МО в области инвестиций и торговли на финансовых рынках. В этом разделе мы расскажем о том, чего ожидать, как достигать своих целей и что вам нужно для достижения своих целей и получения удовлетворения от всего этого процесса.

Чего ожидать

Цель этой книги — ознакомить вас со стратегической перспективой, концептуальным пониманием и практическими инструментами, позволяющими наращивать добавочную стоимость от применения МО в торговом и инвестиционном процессах. В этой связи она охватывает МО не как отдельное мероприятие, а как важный элемент данного процесса.

Прежде всего, она охватывает широкий спектр алгоритмов контролируемого, неконтролируемого и подкрепляемого обучения, широко используемых для извлечения сигналов из разнообразных источников данных, относящихся к различным классам финансовых активов. Она знакомит с рабочим потоком МО и концентрирует внимание читателя на примерах практического его применения с соответствующими данными и многочисленными примерами исходного кода. Однако она также развивает математический и статистический фундамент с целью обеспечения настройки алгоритма или интерпретации результатов.

В этой книге признается, что инвесторы могут извлекать добавочную стоимость из сторонних данных больше, чем другие отрасли промышленности. Как следствие, в ней охвачены не только способы работы с рыночными и фундаментальными данными, но и способы привлечения, оценивания, обработки и моделирования альтернативных источников данных, таких как неструктурированные текстовые и графические данные.

Настоящая книга увязывает использование МО с исследованием и оцениванием альфа-факторов, с количественными и факторно обусловленными стратегиями и вводит портфельный менеджмент в качестве контекста для развертывания стратегий, объединяющих несколько альфа-факторов. В ней также подчеркивается, что МО может наращивать добавочную стоимость, выходя за пределы предсказаний, относящихся к ценам на индивидуальные финансовые активы, например к размещению финансовых средств среди портфельных активов, и обращается к проблеме рисков ложных обнаружений, возникающих в результате использования МО с крупными наборами данных для разработки торговой стратегии.

Не следует удивляться тому, что эта книга не содержит рекомендаций по инвестициям или готовых алгоритмов торговли. Вместо этого она знакомит со строительными блоками, необходимыми для выявления, оценивания и комбинирования наборов данных, подходящих для той или иной инвестиционной цели, отбирает и применяет автоматически обучающиеся алгоритмы к этим данным, а также разрабатывает и тестирует стратегии алгоритмической торговли, основываясь на полученных результатах.

Кто должен прочесть эту книгу

Вы найдете книгу информативной, если вы аналитик, исследователь данных или инженер в области МО с пониманием финансовых рынков и интересом к стратегиям торговли. Вы также должным образом оцените ее как инвестиционный профессионал, который стремится эффективно задействовать МО для принятия более оптимальных решений.

Если ваша квалификация — программное обеспечение и МО, вы можете бегло просмотреть или просто пропустить вводный материал по МО. Если ваши компетентные знания принадлежат области инвестиций, то вы, вероятно, будете знакомы с некоторым или всем финансовым контекстом. Вы, пожалуй, найдете книгу наиболее полезной в качестве краткого обзора ключевых алгоритмов, строительных блоков и примеров использования, чем в качестве источника специализированного анализа конкретного алгоритма или стратегии. Однако в этой книге мы исходим из того, что вы заинтересованы в продолжении изучения этой очень динамичной области. С этой целью в ней приводятся справочные материалы со ссылками на многочисленные ресурсы с целью поддержания вашего движения вперед к индивидуальным стратегиям торговли, в которых эффективно задействуются и развиваются охватываемые ею фундаментальные методы и инструменты.

Вы не должны испытывать трудностей в использовании языка Python 3 и различных научно-вычислительных библиотек, таких как NumPy, pandas или SciPy, и должны быть заинтересованы подхватить на этом пути целый ряд других. Некоторый опыт работы с МО и библиотекой sklearn был бы полезен, однако мы кратко рассмотрим базовый рабочий поток и будем давать различные справочные материалы со ссылками на всевозможные ресурсы с целью заполнения пробелов или более глубокого погружения в рассматриваемую тему.

Как эта книга организована

В книге представлено всестороннее введение в то, как МО может наращивать добавочную стоимость в конструировании и исполнении стратегий торговли. Книга состоит из четырех условных частей, которые охватывают различные аспекты процесса привлечения источников данных и разработки стратегии, а также разные решения всевозможных задач МО.

Каркас — от данных к конструированию стратегии

Первая часть обеспечивает каркас для разработки стратегий алгоритмической торговли. Она сосредоточена на данных, которые питают энергией обсуждаемые в этой книге автоматически обучающиеся (МО) алгоритмы и стратегии, описывает, как МО может использоваться для выведения торговых сигналов, и как развертывать и оценивать стратегии в качестве компонента инвестиционного портфеля.

В оставшейся части этой главы резюмируется то, как и почему МО заняло центральное место в инвестициях, описывается процесс торговли и обрисовывается то, как МО может наращивать добавочную стоимость.

Глава 2 "Рыночные и фундаментальные данные" охватывает источники и работу с исходными тиковыми данными, поступающими с бирж, и данными финансовой отчетности, а также способы доступа к многочисленным поставщикам общедоступных данных, на которые мы будем опираться на протяжении всей этой книги.

Глава 3 "Альтернативные данные для финансов" предоставляет категории и критерии для оценивания качества стремительно растущего числа источников и поставщиков. В ней также демонстрируются способы создания альтернативных наборов данных путем выскабливания веб-сайтов, например для сбора стенограмм телеконференций о корпоративных заработках для использования в алгоритмах естественно-языковой обработки и сентиментного анализа во второй части книги.

Глава 4 "Исследование альфа-факторов" обеспечивает каркас для понимания того, как работают факторы и как измерять их эффективность, например, используя информационный коэффициент (IC). В ней демонстрируется способ конструирования альфа-факторов из данных с помощью библиотек Python в режиме офлайн и на веб-платформе Quantopian. В ней также вводятся библиотека `zipline` для бэктестирования факторов и библиотека `alphalens` для оценивания их предсказательной мощности.

Глава 5 "Оценивание стратегии" знакомит с тем, как строить, тестировать и оценивать стратегии торговли, используя исторические данные с помощью библиотеки `zipline` в режиме офлайн и на веб-платформе Quantopian. Она представляет и демонстрирует способы вычисления результативности портфеля и метрических показателей риска с помощью библиотеки `pyfolio`. В ней также рассматриваются способы управления методологическими сложностями бэктестирования стратегий и вводятся методы оптимизации стратегии с точки зрения портфельного риска.

Основы машинного обучения

Вторая условная часть книги охватывает фундаментальные контролируемые и неконтролируемые обучающиеся алгоритмы и иллюстрирует их применение к стратегиям торговли на финансовых рынках. Она также знакомит с веб-платформой Quantopian, где можно эффективно задействовать и комбинировать данные и методы МО, разработанные в этой книге, для реализации алгоритмических стратегий, которые исполняют сделки на живых рынках.

Глава 6 "Процесс машинного обучения" закладывает фундамент, давая описания того, как формулировать, тренировать и настраивать автоматически обучающиеся модели и оценивать их предсказательную результативность в качестве систематического рабочего процесса.

Глава 7 "Линейные модели" показывает, как использовать линейную и логистическую регрессию для выведения заключения и предсказания и применять регуляризацию для управления риском перепогонки. Она знакомит с торговой платформой Quantopian и демонстрирует способы построения факторных моделей и предсказания цен на финансовые активы.

Глава 8 "Модели временных рядов" посвящена одномерным и многомерным временным рядам, включая векторные авторегрессионные модели и коинтеграционные тесты, а также способы их применения к стратегиям парной торговли.

Глава 9 "Байесово машинное обучение" знакомит с тем, как формулировать вероятностные модели и как извлечение выборок из вероятностных распределений методами Монте-Карло марковских цепей (Markov Chain Monte Carlo, MCMC) и вариационные байесовы методы обеспечивают аппроксимирование вывода. В ней также иллюстрируются способы применения библиотеки PyMC3 для вероятностного программирования с целью понимания сущности параметрической и модельной неопределенности.

Глава 10 "Деревья решений и случайные леса" показывает, как строить, тренировать и настраивать нелинейные древесные модели для проникновения в сущность данных и предсказания. В ней представлены древесные ансамблевые модели и иллюстрируется то, как случайные леса используют бутстраповское агрегирование для преодоления некоторых недостатков деревьев решений.

Глава 11 "Градиентно-бустинговые машины" демонстрирует, как использовать библиотеки xgboost, lightgbm и catboost для выполнения высокорезультативной тренировки и предсказания. Кроме того, в ней подробно анализируется настройка многочисленных гиперпараметров.

Глава 12 "Неконтролируемое обучение" посвящена применению задач снижения размерности и кластеризации в алгоритмической торговле. В ней используется анализ главных и независимых компонент для извлечения ведомых данными рисков факторов. В ней также представлено несколько технических решений для кластеризации и показано использование иерархической кластеризации для размещения финансовых средств среди портфельных активов.

Обработка естественного языка

В третьей условной части основное внимание уделяется текстовым данным и вводятся современные методы неконтролируемого обучения для извлечения высококачественных сигналов из этого ключевого источника альтернативных данных.

Глава 13 "Работа с текстовыми данными" демонстрирует способы конвертирования текстовых данных в числовой формат. В ней также показано применение классификационных алгоритмов из второй условной части, связанных с сентиментным анализом, к большим наборам данных.

Глава 14 "Тематическое моделирование" посвящена применению байесова неконтролируемого обучения с целью извлечения латентных тем, которое способно резюмировать большое число документов и предлагать более эффективные способы разведывания текстовых данных либо использования тем в качестве признаков в классификационной модели. В ней демонстрируется способ применения этого технического решения к стенограммам телеконференций о корпоративных зарплатах, источники которых будут привлечены в *главе 3*, и к годовым отчетам, по-

даваемым в Комиссию по ценным бумагам и биржам (Securities and Exchange Commission, SEC).

Глава 15 "Векторное вложение слов" посвящена использованию нейронных сетей для заучивания ультрасовременных языковых признаков в форме словарных векторов, которые улавливают семантический контекст намного лучше, чем традиционные текстовые признаки, и представляют собой очень перспективное направление для извлечения торговых сигналов из текстовых данных.

Глубокое и подкрепляемое обучение

Четвертая условная часть знакомит с глубоким обучением и подкрепляемым обучением.

Глава 17 "Глубокое обучение" знакомит с библиотеками Keras, TensorFlow и PyTorch — наиболее популярными каркасами глубокого обучения — и иллюстрирует способы тренировки и настройки различных архитектур.

Глава 18 "Сверточные нейронные сети" иллюстрирует применение CNN-сетей с графическими и текстовыми данными.

Глава 19 "Рекуррентные нейронные сети" знакомит с RNN-сетями для данных временных рядов.

Глава 20 "Автокодировщики и генеративные состязательные сети" показывает способы применения глубоких нейронных сетей для неконтролируемого автоматического обучения с автокодировщиками и знакомит с GAN-сетями, которые производят синтетические данные.

Глава 21 "Подкрепляемое обучение" демонстрирует использование автоматического обучения с максимизацией подкрепления для создания динамических агентов, которые заучивают функцию линии поведения, основываясь на вознаграждениях, с использованием программного инструмента Open AI gym.



Главы 17–21 (а также приложения и цветные иллюстрации к *главам 1–15*) представлены в электронном архиве данной книги, размещенном на сайте издательства. Электронный архив можно скачать по ссылке <ftp://ftp.bhv.ru/9785977565950.zip> и со страницы книги на веб-сайте издательства по адресу www.bhv.ru.

Что вам нужно для успеха

Изложение в книге вращается вокруг применения автоматически обучающихся алгоритмов к разным наборам данных. Значительный дополнительный материал размещен в репозитории GitHub, который обеспечивает повторение изученных тем и эксперименты с примерами, обсуждаемыми в книге. В нем представлены дополнительные сведения и инструкции, а также многочисленные справочные материалы.

Источники данных

Мы будем использовать свободно доступные исторические данные из рыночных, фундаментальных и альтернативных источников. *Главы 2 и 3* охватывают характеристики и обращаются к этим источникам данных и знакомят с ключевыми поставщиками, которых мы будем использовать на протяжении всей книги. Сопутствующий репозиторий GitHub (*см. далее*) содержит инструкции о том, как получать или создавать некоторые наборы данных, которые мы будем использовать повсюду, и включает некоторые меньшие наборы данных.

Несколько примеров источников данных, которые мы будем привлекать и с которыми будем работать, включают следующие, но не ограничиваются только ими:

- ◆ данные ордерной книги протокола передачи данныхITCH биржи NASDAQ;
- ◆ документы официальной отчетности в системе электронного сбора, анализа и извлечения данных (Electronic Data Gathering, Analysis and Retrieval, EDGAR) Комиссии по ценным бумагам и биржам США (U.S. Securities and Exchange Commission, SEC);
- ◆ стенограммы телеконференций о корпоративных заработках из портала Seeking Alpha (<https://seekingalpha.com/>);
- ◆ дневные цены рыночной площадки Quandl и другие точки данных более чем по 3000 акций США;
- ◆ различные макроэкономические фундаментальные данные Федеральной резервной системы и др.;
- ◆ крупный набор данных отзывов о деятельности предприятий веб-сайта Yelp и наборы данных социальной сети Twitter;
- ◆ снимковые данные по нефтяным танкерам.

Некоторые данные имеют размер в несколько гигабайт (например, данные биржи NASDAQ и документы финансовой отчетности, подаваемые в комиссию SEC). Такие моменты будут оговариваться в блокнотах Jupyter.

Репозиторий GitHub

Репозиторий GitHub содержит блокноты Jupyter, которые подробнее иллюстрируют многие понятия и модели. Отсылки к блокнотам рассеяны по всей книге, где они используются. Каждая глава имеет собственный каталог с отдельными инструкциями, где это необходимо, а также ссылку на содержимое конкретной главы.

Блокноты Jupyter — отличный инструмент для создания воспроизводимых вычислительных повествований. Он позволяет пользователям создавать и обмениваться документами, которые соединяют в себе "живой" код с повествовательным текстом, математическими уравнениями, визуализациями, интерактивными элементами управления и другим богатым выводимым материалом. Они также предостав-

ляют строительные блоки для интерактивных вычислений с данными, такие как файловый браузер, консоли и текстовый редактор¹.



Файлы исходного кода размещены по адресу
<https://github.com/PacktPublishing/Hands-On-Machine-Learning-for-Algorithmic-Trading>.

Библиотеки Python

В данной книге используется Python 3.7, и рекомендуется применение мини-версии дистрибутива Anaconda — miniconda — для установки менеджера пакетов conda и создания среды conda с целью установки реквизитных библиотек. В связи с этим репозиторий GitHub содержит манифестный файл `environment.yml`. Пожалуйста, обратитесь к инструкциям по установке, указанным в файле `README` в репозитории GitHub².

Рост популярности машинного обучения в инвестиционной индустрии

За последние несколько десятилетий инвестиционная индустрия значительно эволюционировала и продолжает это делать в условиях возросшей конкуренции, технического прогресса и сложных экономических условий. В этом разделе мы рассмотрим несколько ключевых трендов, которые сформировали инвестиционную среду в целом, а также контекст алгоритмической торговли на финансовых рынках, в частности, и связанные с этим темы, которые будут повторяться на протяжении всей книги.

Тренды, которые привели алгоритмическую торговлю и МО к нынешнему видному положению, включают:

- ◆ изменения в микроструктуре рынка, такие как распространение электронной торговли и интеграция рынков по всем классам финансовых активов и географиям;
- ◆ разработку инвестиционных стратегий с учетом влияния рисков факторов в отличие от классов финансовых активов;
- ◆ революции в вычислительной мощности, генерировании и управлении данными, а также аналитических методах;

¹ Если по каким-то причинам файл блокнота из репозитория книги открыть не удастся, то имеется обходной путь — перейти на веб-сайт обозревателя блокнотов nbviewer по адресу **<https://nbviewer.jupyter.org/>**, вставить в поисковое поле адрес нужного блокнота и нажать кнопку **Go!**. — Прим. перев.

² См. подробное описание организации работы с блокнотами Jupyter и настройку дистрибутива Anaconda в приложении 1. — Прим. перев.

- ♦ повышенную результативность первопроходцев в алгоритмической торговле по сравнению с человеческими, дискреционными инвесторами.

В дополнение к этому финансовые кризисы 2001 и 2008 гг. повлияли на подход инвесторов к диверсификации и риск-менеджменту и привели к появлению недорогих пассивных инвестиционных механизмов в форме *биржевых фондов*³ (ETF). На фоне низкой доходности и низкой волатильности после кризиса 2008 г. экономные инвесторы перевели 2 трлн долларов из активно управляемых взаимных фондов⁴ в пассивно управляемые биржевые фонды ETF. Конкурентное давление также отражается в более низких хедж-фондовых комиссионных, которые к 2017 г. упали с традиционных 2%-х комиссионных за годовой менеджмент и 20%-й доли в прибыли в среднем соответственно до 1,48 и 17,4%.

От торговли электронной к торговле высокочастотной

Электронная торговля радикально продвинулась в плане возможностей, объема, охвата классов финансовых активов и географии с тех пор, как в 1960-х годах сети начали маршрутизировать цены на компьютерные терминалы.

Фондовые рынки, т. е. рынки долевых ценных бумаг⁵ частных компаний, возглавили этот тренд во всем мире. Правила обработки ордеров комиссией SEC от 1997 г. внесли в биржи конкуренцию через *электронные коммуникационные сети*⁶ (ECN). Сети ECN представляют собой автоматизированные *альтернативные системы торговли* (Alternative Trading Systems, ATS), которые сочетают ордера на покупку и продажу по конкретным ценам, прежде всего на долевые ценные бумаги и валюты, и зарегистрированы как брокеры-дилеры. Она позволяет крупным брокерам и индивидуальным биржевым торговцам в различных географических точках торговать напрямую без посредников как на биржах, так и в нерабочее время. Темные пулы

³ Биржевой фонд (exchange-traded fund, ETF), или торгуемый на бирже фонд, — это ликвидная ценная бумага, которая отслеживает фондовый индекс, товар, облигации или корзину финансовых активов. См. <https://www.investopedia.com/terms/e/etf.asp>. — Прим. перев.

⁴ Взаимный фонд (mutual fund), или фонд взаимных инвестиций, — это финансовый механизм, состоящий из денежных средств, собранных от многочисленных инвесторов с целью инвестирования в ценные бумаги, такие как акции, облигации, инструменты денежного рынка и другие финансовые активы. Представляет собой портфель активов, тщательно отобранных и приобретенных профессиональными финансистами на вложения многих тысяч мелких вкладчиков. См. <https://www.investopedia.com/terms/m/mutualfund.asp>. — Прим. перев.

⁵ Долевая ценная бумага (equity) — это ценная бумага, обеспечивающая право собственности в компании через покупку обыкновенных либо привилегированных акций (долей). Термин equity также означает нетто-стоимость или чистую стоимость компании. См. http://www.morningstar.com/InvGlossary/equity_definition_what_is.aspx. — Прим. перев.

Рынок долевых ценных бумаг (equity market), или фондовый рынок, — это рынок, где выпускаются в обращение и торгуются акции, т. е. доли собственности в компании, через биржу или внебиржевые рынки. См. <https://www.investopedia.com/terms/s/shareholdersequity.asp>. — Прим. перев.

⁶ Электронные коммуникационные сети ECN (electronic communication networks) — это электронная система осуществления сделок купли-продажи биржевых товаров. — Прим. перев.

ликвидности — это еще один тип альтернативных систем торговли, которые позволяют инвесторам размещать ордера и торговать, не раскрывая свою информацию публично, как в ордерной книге, поддерживаемой биржей. Темные пулы, выросшие из постановления комиссии SEC 2007 г., часто размещаются в крупных банках и подлежат регламентированию со стороны комиссии SEC.

С ростом влияния электронной торговли алгоритмы затратно эффективного исполнения прошли ускоренную эволюцию, и их принятие быстро распространилось со стороны продажи на сторону покупки⁷ и по всем классам активов. Автоматизированная торговля появилась примерно в 2000 г. как инструмент на стороне продажи, направленный на затратно эффективное исполнение сделок, который распределяет ордера во времени, ограничивая их влияние на финансовый рынок⁸. Эти инструменты распространились на сторону покупки и становились все более изощренными, принимая в расчет, например, транзакционные издержки и ликвидность, а также краткосрочные ценовые и объемные прогнозы.

Прямой доступ к рынку (Direct Market Access, DMA) дает биржевому торговцу более высокий контроль над исполнением, позволяя ему отправлять ордера непосредственно на биржу, используя инфраструктуру и идентификацию участника рынка у брокера, который является членом биржи. Спонсируемый доступ снимает с брокеров обязанность осуществлять контроль предсделочных рисков и формирует основу для *высокочастотной торговли* (high-frequency trading, HFT).

Под термином "высокочастотная торговля" имеются в виду автоматизированные сделки с финансовыми инструментами, которые исполняются с чрезвычайно низкой временной задержкой в микросекундном интервале, и где участники занимают позиции в течение очень коротких периодов. Ее цель заключается в выявлении и использовании неэффективности микроструктуры рынка, институциональной инфраструктуры торговых площадок. За последние 10 лет высокочастотная торговля существенно выросла и по оценкам составляет примерно 55% объема торговли на фондовых рынках США и около 40% на фондовых рынках Европы. То же самое произошло на фьючерсных рынках, где высокочастотная торговля выросла примерно до 80% от объемов фьючерсов на иностранную валюту и двух третей от объемов фьючерсов на процентные ставки и 10-летние казначейские векселя (согласно ежегодному "Обзору доступа к финансовым службам FAS 2016").

Стратегии высокочастотной торговли ориентированы на получение малой прибыли в расчете на сделку с использованием пассивных или агрессивных стратегий. *Пас-*

⁷ Сторона покупки (buy side) и сторона продажи (sell side) относятся к фирмам, которые соответственно приобретают и торгуют товарами и услугами. Применительно к финансам, сторона покупки относится к покупке ценных бумаг и включает инвестиционных менеджеров, пенсионные фонды и хедж-фонды, а сторона продажи относится к фирмам, которые выпускают, продают или торгуют ценными бумагами, и включает инвестиционные банки, консультативные фирмы и корпорации. — *Прим. перев.*

⁸ Влияние на финансовый рынок (market impact) означает влияние, которое участник рынка оказывает на цену финансового актива, когда он покупает или продает его. — *Прим. перев.*

⁹ См. <http://data.imf.org/?sk=E5DCAB7E-A5CA-4892-A6EA-598B5463A34C>. — *Прим. перев.*

сивные стратегии включают арбитражную торговлю¹⁰, ориентированную на получение прибыли от очень небольших ценовых разниц для одного и того же финансового актива или его производных инструментов, торгуемых на разных площадках. Агрессивные стратегии включают предвосхищение заявок или инициацию импульса. Предвосхищение заявок, также именуемое *обнаружением ликвидности*, охватывает алгоритмы, которые размещают небольшие разведывательные ордера с целью обнаружения скрытой ликвидности у крупных институциональных инвесторов и торгуют перед крупным ордером с целью извлечь выгоду из последующих ценовых движений. Инициация импульса подразумевает алгоритм, исполняющий и аннулирующий серию ордеров с целью заманить другие алгоритмы высокочастотной торговли агрессивнее покупать (или продавать), и извлекает выгоду из результирующих ценовых изменений.

В связи с этим регуляторы выражали озабоченность по поводу потенциальной взаимосвязи между некоторыми агрессивными стратегиями высокочастотной торговли и повышенной хрупкостью и волатильностью рынка, например, во время молниеносного обвала в мае 2010 г., волатильности рынка казначейских обязательств в октябре 2014 г. и молниеносного обвала промышленного индекса Доу-Джонса в августе 2015 г. более чем на 1000 пунктов. Вместе с тем из-за присутствия высокочастотной торговли рыночная ликвидность увеличивалась вместе с объемами торгов, что снижало совокупные транзакционные издержки¹¹.

Сочетание сниженных торговых объемов на фоне снижения волатильности и растущей стоимости технологии и доступа как к данным, так и к торговым площадкам привело к финансовому давлению. По оценкам, совокупные поступления высокочастотной торговли от американских акций впервые с 2008 г. упали ниже 1 млрд долларов по сравнению с 7,9 млрд долларов в 2009 г.

Этот тренд привел к отраслевой консолидации за счет различных приобретений, например, крупнейшей зарегистрированной проприетарной торговой фирмой Virtu Financial и долевыми инфраструктурными инвестициями, такими как новый маршрут с ультранизкой временной задержкой Go West между Чикаго и Токио. Одновременно с этим такие стартапы, как Alpha Trading Lab, обеспечивают наличие торговой инфраструктуры и данных высокочастотной торговли, демократизируя высокочастотную торговлю за счет привлечения разработчиков алгоритмов через Интернет в обмен на долю прибыли.

¹⁰ Арбитраж (arbitrage) — это практика использования преимуществ разницы цен между двумя или более рынками и заключения комбинации совпадающих сделок, которые извлекают выгоду из таких дисбалансов. Например, под валютным арбитражем понимается одновременная покупка и продажа валюты на разных валютных рынках с целью получения прибыли от разницы обменных курсов в двух местах. См. <https://www.thefreelibrary.com/Currency+arbitrage+as+a+tool+of+corporate+financial+management-a0469641512>. — Прим. перев.

¹¹ Транзакционные издержки (transaction cost) — это стоимость осуществления любой экономической сделки, сопровождающая участие в рынке. В инвестировании это затраты, понесенные при покупке или продаже активов, такие как комиссии и спред.

См. https://en.wikipedia.org/wiki/Transaction_cost. — Прим. перев.

Факторное инвестирование и умные бета-фонды

Финансовый возврат¹², обеспечиваемый активом, является функцией неопределенности или риска, связанного с финансовой инвестицией. Инвестирование в долевые ценные бумаги подразумевает, например, принятие корпоративного делового риска, а инвестиция в облигации — принятие риска дефолта.

Поскольку конкретные характеристики рисков несут в себе предсказания финансовых возвратов, выявление и прогнозирование поведения этих рисковых факторов¹³ становится первостепенным направлением при конструировании инвестиционной стратегии. Это позволяет добиться ценных торговых сигналов и является ключом к превосходящим результатам активного менеджмента. Понимание рисковых факторов в финансовой индустрии со временем существенно эволюционировало и повлияло на то, как МО используется для алгоритмической торговли.

Теория современного инвестиционного портфеля (MPT)¹⁴ ввела разграничение между источниками идиосинкратического и систематического риска для определенного актива. Идиосинкратический риск¹⁵ может быть устранен путем диверсификации, но систематический риск — нет. В начале 1960-х годов *модель ценообразования капитальных активов* (CAPM)¹⁶ выявила один-единственный фактор, сти-

¹² Финансовый возврат (return), или возврат на инвестицию или финансовая отдача, — деньги, сделанные или потерянные на инвестициях. Возврат выражается номинально как изменение денежной стоимости инвестиции с течением времени, либо в процентах из соотношения прибыли к инвестициям. В этом случае он именуется возвратностью.

См. <https://www.investopedia.com/terms/r/return.asp>. — *Прим. перев.*

¹³ Рисковый фактор (risk factor) — это измеримая характеристика или элемент, изменение которого может повлиять на стоимость актива, например, обменный курс, процентная ставка и рыночная цена. См. <http://www.businessdictionary.com/definition/risk-factor.html>. — *Прим. перев.*

¹⁴ Теория современного инвестиционного портфеля (modern portfolio theory, MPT), или портфельная теория Марковица, — это теория, которая описывает то, как инвесторы, не склонные к риску, могут создавать портфели для оптимизации или максимизации ожидаемой возвратности на основе заданного уровня рыночного риска, подчеркивая, что риск является неотъемлемой частью более высокого вознаграждения. См. <https://www.investopedia.com/terms/m/modernportfoliotheory.asp>. — *Прим. перев.*

Портфель (portfolio), как правило, представляет собой группу финансовых активов, таких как акции, облигации, товары, валюты и эквиваленты денежных средств, а также их фондовые аналоги, включая взаимные, биржевые и закрытые фонды. См. <https://www.investopedia.com/terms/p/portfolio.asp>. — *Прим. перев.*

¹⁵ Идиосинкратический риск (idiosyncratic risk) — это риск, который относится к факторам, влияющим на отдельный актив, по сравнению с систематическим риском, который относится к более широким трендам, сказывающимся на всем рынке в целом.

См. <https://www.investopedia.com/search?q=idiosyncratic+asset>. — *Прим. перев.*

¹⁶ Модель ценообразования капитальных активов (Capital Asset Pricing Model, CAPM) описывает связь между систематическим риском и ожидаемой возвратностью от активов, в особенности от акций. Указанная модель является однофакторной, и единственным фактором служит бета-фактор, который, согласно модели, показывает, насколько хорошо акция движется относительно рынка. Акции, которые двигались выше рынка, имели более высокий бета-фактор и, следовательно, более высокий риск и возврат.

См. <https://www.investopedia.com/terms/c/capm.asp>. — *Прим. перев.*

мулирующий финансовые возвраты от всех активов: финансовый возврат, приносимый рыночным портфелем¹⁷ сверх казначейских векселей. Рыночный портфель состоял из всех торгуемых ценных бумаг, взвешенных по их рыночной стоимости. Систематическое влияние на актив рынка измеряется бетой, которая представляет собой корреляцию между возвратами, приносимыми активом, и рыночным портфелем.

Признание, что риск того или иного актива зависит не от него в отдельности, а от того, как он движется по отношению к другим активам и рынку в целом, стало крупным концептуальным прорывом. Другими словами, активы получают рисковую премию не из-за их идиосинкратических особенностей, а из-за влияния на них базовых факторных рисков.

Вместе с тем крупный пласт академической литературы и многолетний опыт инвестирования опровергли предсказание CAPM-модели о том, что рисковые премии финансовых активов зависят от влияния на них лишь одного-единственного фактора, измеряемого бетой финансового актива. Напротив, с тех пор были обнаружены многочисленные дополнительные рисковые факторы. Фактор — это квантифицируемый сигнал, атрибут или любая величина, которая исторически коррелировала с будущими возвратами от акций и, как ожидается, останется коррелированной в будущем.

Эти рисковые факторы были объявлены аномалиями, поскольку они противоречили *гипотезе об эффективном рынке* (Efficient Market Hypothesis, ЕМН), которая доказывала, что рыночное равновесие всегда будет формировать цену ценных бумаг в соответствии с CAPM-моделью, вследствие чего никакие другие факторы не должны иметь предсказательной мощности. Экономическая теория в основе факторов может быть либо рациональной, когда премии за факторные риски компенсируют низкие финансовые возвраты в плохие времена, либо поведенческой, когда агенты не в состоянии задействовать избыточные возвраты за счет арбитража.

Хорошо известные аномалии включают стоимостные, размерные и импульсные эффекты, которые помогают предсказывать возвраты, учитывая при этом рыночный фактор CAPM-модели. Размерный эффект зиждется на том факте, что малые фирмы систематически превосходят по результативности крупные, как было обнаружено Банцем (Banz, 1981) и Рейнганумом (Reinganum, 1981). Стоимостной эффект (Basu, 1982) утверждает, что фирмы с низкими метриками стоимости демонстрируют повышенную результативность. Это говорит о том, что фирмы с низкими ценовыми мультипликаторами, такими как соотношение цены к корпоративным заработкам или цены к балансовой стоимости, показывают более высокую результативность, чем их более дорогие сверстники (как было предложено изобретателями инвестиций в стоимость Бенджамином Грэмом (Benjamin Graham) и Дэвидом Доддом (David Dodd) и популяризировано Уорреном Баффетом (Warren Buffett)).

¹⁷ Рыночный портфель (market portfolio) — это теоретическая диверсифицированная группа инвестиций, каждый актив которой взвешен пропорционально его суммарному присутствию на рынке. См. <https://www.investopedia.com/terms/m/market-portfolio.asp>. — Прим. перев.



Корпоративные заработки (earnings) обычно относятся к чистому доходу после уплаты налогов, иногда именуемому нижней строкой отчета о прибылях и убытках, или прибылями компании, согласно формуле:

Валовая выручка

минус Стоимость реализованной продукции/услуг

= Чистая выручка

плюс Другие доходы

минус Операционные расходы

= Валовая прибыль (EBITDA, корпоративный заработок)

минус Налоги, проценты, износ и амортизация

= Чистая прибыль (нераспределенная прибыль,
это тоже корпоративный заработок).

Корпоративные заработки являются основным определяющим фактором долевой цены компании, поскольку заработки и связанные с ними обстоятельства могут указывать на прибыльность и успешность предприятия в долгосрочной перспективе. Они являются, пожалуй, самым важным и наиболее изученным показателем в финансовой отчетности компании, поскольку показывают прибыльность по сравнению с оценками со стороны аналитиков и руководства компании (см. <https://www.investopedia.com/terms/e/earnings.asp>. — Прим. перев.)

Импульсный эффект, обнаруженный в конце 1980-х годов, в частности, Клиффордом Эснесом (Clifford Asness), партнером-основателем глобальной компании по инвестиционному менеджменту AQR, констатирует, что акции с хорошим импульсом, с точки зрения недавних 6–12-месячных возвратов, имеют более высокие возвраты в перспективе, чем слабо-импульсные акции с аналогичным рыночным риском. Исследователи еще обнаружили, что стоимостные и импульсные факторы объясняют возвраты от акций за пределами США, а также от других классов финансовых активов, таких как облигации, валюты и товары, и дополнительные рискованные факторы.

Применительно к ценным бумагам с фиксированной доходностью стоимостная стратегия называется *выездом на кривой отдачи*¹⁸ и является формой премии за продолжительность. Применительно к товарам она называется *скользящей возвратностью*¹⁹ с положительной возвратностью для фьючерсов с наклоном кривой

¹⁸ Выезд на кривой отдачи (riding the yield curve) — это торговая стратегия, включающая покупку долгосрочной облигации и ее продажу до ее созревания с целью получения прибыли от снижения доходности, которое происходит в течение срока действия облигации.

См. <https://www.investopedia.com/terms/r/ridingtheyieldcurve.asp>. — Прим. перев.

¹⁹ Скользящая возвратность (roll return, rolling return) представляет собой среднегодовую возвратность за период, заканчивающийся указанным годом.

См. <https://www.investopedia.com/terms/r/rollingreturns.asp>. — Прим. перев.

вверх и отрицательной возвратностью в противном случае. Применительно к иностранным валютам стоимостная стратегия называется *кэрри-трейдом*²⁰.

Существует также премия за неликвидность. Ценные бумаги, которые являются более неликвидными, торгуются по низким ценам и имеют высокие средние избыточные возвраты по сравнению с их более ликвидными аналогами. Облигации с более высоким риском дефолта, как правило, в среднем имеют более высокие возвраты, отражая премию за кредитный риск. Поскольку инвесторы готовы платить за страхование от высокой волатильности, когда возвраты тяготеют к обрушению, продавцы защиты от волатильности на опционных рынках чаще всего получают высокие возвраты.

Многофакторные модели определяют риски в более широких и разнообразных терминах, чем просто рыночный портфель. В 1976 г. Стивен Росс (Stephen Ross) предложил теорию арбитражного ценообразования, которая утверждала, что инвесторы получают компенсацию за многочисленные систематические источники риска, которые нельзя диверсифицировать. Три наиболее важными макроэкономическими факторами являются рост, инфляция и волатильность наряду с производительностью, демографическим и политическим рисками. В 1992 г. Юджин Фама (Eugene Fama) и Кеннет Френч (Kenneth French) объединили рискованные факторы долевых ценных бумаг²¹ — размер и стоимость — с рыночным фактором в единую модель, которая лучше объясняла срезовые (кросс-секционные) финансовые возвраты от акций. Позже они добавили модель, которая также включала рискованные факторы облигаций, одновременно объясняя возвраты для обоих классов активов.

Особо привлекательным аспектом рискованных факторов является их низкая или отрицательная корреляция. Например, стоимостные и импульсные рискованные факторы имеют отрицательную корреляцию, снижая риск и повышая скорректированные на риск возвраты за пределами выгоды, вытекающей из рискованных факторов. Более того, за счет кредитного плеча²² и длинно-коротких инвестиционных стратегий²³ фак-

²⁰ Кэрри-трейд (carry-trade, carry) — это стратегия, в соответствии с которой высокодоходная валюта финансирует сделку с низкодоходной валютой. Иными словами, это получение прибыли на валютном рынке за счет разной величины процентных ставок.

См. <https://www.investopedia.com/terms/c/currencycarrytrade.asp>. — *Прим. перев.*

²¹ Фактор (factor) — это характерная, поддающаяся квантификации особенность актива с существенной информацией о риске и возвратности. В случае долевых ценных бумаг наиболее известные и наиболее документированные факторы включают стоимость, размер, импульс, низкую волатильность и качество. — *Прим. перев.*

²² Кредитное плечо (leverage), или финансовый рычаг, — это дополнительная покупательная способность, создаваемая маржинальной торговлей, позволяющая эффективно платить меньше полной цены за актив, используя заемные средства. Кредитное плечо обычно представлено в виде соотношения: например, если на торговом счете имеется 10 тыс. долларов и вы занимаете еще 10 тыс. долларов, то ваше кредитное плечо равно 2:1. — *Прим. перев.*

²³ Длинно-короткая инвестиционная стратегия (long-short strategy) — это инвестиционная стратегия, обычно связанная с хедж-фондами и с некоторыми прогрессивными традиционными менеджерами активами, которая заключается в покупке длинных акций, т. е. с ожидаемым ростом цены, и продаже

торные стратегии могут быть объединены в рыночно-нейтральные подходы. Сочетание длинных позиций в ценных бумагах, находящихся под влиянием положительных рисков, с короткими позициями в ценных бумагах, находящихся под влиянием отрицательных рисков, позволяет собирать динамические рисковые премии.

Как следствие, факторы, объяснявшие возвраты выше и за пределами однофакторной CAPM-модели, встраивались в так называемые инвестиционные стили, которые накренивают портфели в пользу одного или нескольких факторов, и в результате активы начали мигрировать в факторно-ориентированные портфели. Финансовый кризис 2008 г. подчеркнул, насколько метки классов финансовых активов могут вводить в заблуждение и создавать ложное ощущение диверсифицированности, когда инвесторы не смотрят на лежащие в основе факторные риски, поскольку все классы активов обрушились вместе.

За последние несколько десятилетий квантитативное факторное инвестирование²⁴ эволюционировало от простого подхода, основанного на двух или трех стилях, до многофакторных умных или экзотических бета-продуктов. В 2017 г. умные бета-фонды²⁵ пересекли отметку 1 трлн долларов активов под управлением, что свидетельствует о популярности гибридной инвестиционной стратегии, сочетающей активный и пассивный менеджмент. Умные бета-фонды используют пассивную стратегию, но модифицируют ее в соответствии с одним или несколькими факторами, такими как более дешевые акции или их фильтрация в соответствии с выплатами дивидендов, генерируя более высокие финансовые возвраты. Этот рост совпал с усилением критики по поводу высоких комиссионных, взимаемых традиционными активными менеджерами, а также усилением контроля за их работой.

Продолжающееся обнаружение и успешное прогнозирование рисков факторов, которые индивидуально либо в сочетании с другими рисковыми факторами значительно влияют на будущие возвраты от активов по всем классам активов, является ключевой движущей силой всплеска интереса к МО в инвестиционной индустрии и будет ключевой темой на протяжении всей этой книги.

коротких акций, т. е. с ожидаемым снижением цены. Прилагательные "длинный" и "короткий" обозначают направление движения цены соответственно вверх и вниз, не обязательно обозначая продолжительность. — *Прим. перев.*

²⁴ Факторное инвестирование (factor investing) — это стратегия, которая выбирает ценные бумаги по атрибутам, связанным с более высокой возвратностью. Существует два основных типа факторов, определяющих возвратность активов: макроэкономические факторы и стилевые факторы. Первые охватывают широкие риски по всем классам активов, в то время как вторые призваны объяснить возвратности и риски в рамках классов активов. Некоторые общие макроэкономические факторы включают кредит, инфляцию и ликвидность, в то время как стилевые факторы среди прочих охватывают стиль, стоимость и импульс. См. <https://www.investopedia.com/terms/f/factor-investing.asp>. — *Прим. перев.*

²⁵ Умное бета (smart beta) инвестирование сочетает в себе преимущества пассивного инвестирования и преимущества активных инвестиционных стратегий. Цель умного бета-инвестирования — получить альфа, снизить риск или увеличить диверсификацию за стоимость ниже, чем традиционный активный менеджмент, и слегка выше, чем прямое индексное инвестирование. Оно стремится отыскать наилучшую структуру оптимально диверсифицированного портфеля.

См. <https://www.investopedia.com/terms/s/smart-beta.asp>. — *Прим. перев.*

Алгоритмические первопроходцы превосходят людей в широком масштабе

Репутация и рост *активов под управлением* (Assets Under Management, AUM) фирм²⁶, которые возглавили алгоритмическую торговлю, сыграли ключевую роль в формировании интереса инвесторов и последующих отраслевых усилий по воспроизведению их успеха. Систематические фонды отличаются от высокочастотной торговли тем, что в отличие от преимуществ от чистой скорости сделки могут проводиться значительно дольше, стремясь использовать возможности арбитража.

Систематические стратегии, которые в основном или исключительно полагаются на алгоритмическое принятие решений, приобрели большую известность благодаря внедрившему их математику Джеймсу Саймонсу (James Simons), который в 1982 г. основал хедж-фонд Renaissance Technologies и построил ее в первую квантовую фирму. Ее скрытный фонд Medallion, который закрыт для посторонних, по оценкам, с 1982 г. зарабатывал финансовый возврат в размере 35% в годовом исчислении.

Три самых известных квантитативных хедж-фонда, D.E. Shaw, Citadel и Two Sigma, которые используют систематические стратегии, основанные на алгоритмах, в 2017 г. впервые добрались до 20 лучших исполнителей на рынке с точки зрения суммарного долларового дохода для инвесторов после сборов и с момента создания.

Квантитативный хедж-фонд D.E. Shaw, основанный в 1988 г., с 47 млрд долларов AUM в 2018 г. присоединился к списку под номером 3. Квантитативный хедж-фонд Citadel, запущенный в 1990 г. Кеннетом Гриффином (Kenneth Griffin), управляет 29 млрд долларов и занимает 5 место, и квантитативный хедж-фонд Two Sigma, запущенный лишь в 2001 г. бывшими сотрудниками D.E. Shaw Джоном Овердеком (John Overdeck) и Дэвидом Сигелом (David Siegel), вырос с 8 млрд долларов AUM в 2011 г. до 52 млрд долларов в 2018 г. Хедж-фонд Bridgewater, запущенный в 1975 г., с более чем 150 млрд долларов AUM продолжает лидировать благодаря своему фонду Pure Alpha Fund, который также встраивает систематические стратегии.

Аналогичным образом, в списке 100 лучших хедж-фондов журнала Institutional Investor за 2017 г. (<https://www.institutionalinvestor.com>) пять из шести ведущих фирм при принятии инвестиционных решений в значительной степени или полностью опираются на компьютеры и торговые алгоритмы — и все они увеличивают свои финансовые активы в условиях, которые в иных отношениях были бы сложными. Несколько квантитативно-ориентированных фирм поднялись на несколько ступеней и в некоторых случаях увеличили свои активы на двузначные проценты. Номер 2 в списке, глобальная компания по управлению инвестициями *Applied*

²⁶ Активы под управлением (assets under management, AUM) служат мерой общей рыночной стоимости всех финансовых активов, которыми финансовое учреждение, такое как взаимный фонд, венчурная компания или брокерский дом, управляет от имени своих клиентов и себя. — Прим. перев.

Quantitative Research (AQR) в 2017 г. увеличила свои хедж-фондовые активы на 48% до 69,7 млрд долларов и управляла 187,6 млрд долларов в масштабах всей фирмы.

Среди всех хедж-фондов, ранжированных по совокупной результативности за последние 3 года, квантовые фонды, управляемые компанией Renaissance Technologies, достигли позиций 6 и 24, хедж-фонд Two Sigma достиг позиции 11, хедж-фонд D. E. Shaw — позиций 18 и 32, хедж-фонд Citadel — позиций 30 и 37. Помимо лучших исполнителей на рынке, за последние несколько лет хорошо зарекомендовали себя алгоритмические стратегии. За последние 5 лет квантовые хедж-фонды росли примерно на 5,1% в год, в то время как за тот же период средний хедж-фонд вырастал на 4,3% в год.

Фонды, ведомые машинным обучением, привлекают 1 трлн долларов AUM

Три общеизвестные революции в вычислительной мощности, данных и методах МО сделали принятие систематических, ведомых данными стратегий не только убедительнее и экономически эффективнее, но и ключевым источником конкурентного преимущества.

Вследствие этого алгоритмические подходы находят более широкое применение не только в хедж-фондовой индустрии, которая была первопроходцем этих стратегий, но и в более широком круге менеджеров финансовыми активами и даже в пассивно управляемых механизмах, таких как биржевые фонды ETF. В частности, предсказательная аналитика с использованием машинного обучения и алгоритмической автоматизации играет все более заметную роль на всех этапах инвестиционного процесса по всем классам финансовых активов, от генерирования идей и исследований до формулирования стратегии и конструирования портфеля, исполнения сделок и риск-менеджмента.

Оценочные размеры отрасли различаются по той причине, что нет объективного определения квантитативного или алгоритмического фонда, и многие традиционные хедж-фонды или даже взаимные фонды и биржевые фонды ETF внедряют компьютерные стратегии или интегрируют их в дискреционную среду в подходе "человек плюс машина".

По оценкам инвестиционного банка Morgan Stanley за 2017 г., за последние 6 лет алгоритмические стратегии росли на 15% в год и контролируют около 1,5 трлн долларов между хедж-фондами, взаимными фондами и умными биржевыми бета-фондами. Другие отчеты предполагают, что индустрия квантитативных хедж-фондов уже практически превысила 1 трлн долларов AUM, почти удвоив свой размер с 2010 г. на фоне оттока из традиционных хедж-фондов. В отличие от этого, согласно последнему глобальному докладу поставщика хедж-фондовых данных Hedge Fund Research (<https://www.hedgefundresearch.com/>), суммарный капитал хедж-фондовой индустрии достиг 3,21 трлн долларов.

По оценкам исследовательской фирмы Preqin, почти 1500 хедж-фондов совершают большинство сделок с помощью компьютерных моделей. Квантитативные хедж-фонды теперь ответственны за 27% всех сделок инвесторами с акциями США по сравнению с 14% в 2013 г. Но многие используют исследователей данных — или квантов²⁷ — которые, в свою очередь, используют машины для построения крупных статистических моделей²⁸.

Вместе с тем в последние годы фонды перешли к истинному МО, где искусственно-интеллектуальные системы могут быстро анализировать большие объемы данных и самосовершенствоваться посредством такого анализа. Недавние примеры включают фирмы по инвестиционному менеджменту Rebellion Research, Sentient и Aidyia, которые опираются на эволюционные алгоритмы и глубокое обучение при разработке полноавтоматических инвестиционных платформ, приводимых в действие *искусственным интеллектом* (ИИ).

Из стержневой хедж-фондовой индустрии принятие алгоритмических стратегий распространилось на взаимные фонды и даже на пассивно управляемые биржевые фонды в форме умных бета-фондов и на дискреционные фонды в форме квантоментальных²⁹ подходов.

Возникновение квантоментальных фондов

В активном инвестиционном менеджменте сформировалось два четких подхода: систематическое (или квантовое) инвестирование и дискреционное инвестирование. Систематические подходы опираются на алгоритмы повторяющегося и ведомого данными подхода по выявлению инвестиционных возможностей по многочисленным ценным бумагам; в отличие от этого дискреционный подход предусматривает углубленный анализ меньшего числа ценных бумаг. Эти два подхода становятся все более похожими, поскольку фундаментальные менеджеры используют все больше подходов, приводимых в действие наукой о данных.

Теперь даже фундаментальные торговцы вооружаются квантитативными методами, что, согласно инвестиционного банка Barclays, составляет 55 млрд долларов систематических финансовых активов. Нейтральные к конкретным компаниям, квантитативные фонды торгуют регулярностями и динамикой по широкому спектру ценных бумаг. Как показывают данные, собранные инвестиционным банком Barclays, на квантов теперь приходится около 17% от суммарного объема хедж-фондовых активов.

²⁷ Квант (quant) — это биржевой специалист по квантитативным методам. См. <https://habr.com/company/iticapital/blog/307854/>. — *Прим. перев.*

²⁸ См. статью "The Quants Run Wall Street Now" ("Теперь на Уолл-стрит заправляют кванты") в Wall Street Journal. — <https://www.wsj.com/articles/the-quants-run-wall-street-now-1495389108>.

²⁹ Квантоментальный (quantamental) — неологизм, который образован из двух терминов "квантитативный подход" и "фундаментальный подход" и означает сочетание новейших квантитативных методов, в том числе на основе машинного обучения и ИИ, с классическими методами на основе фундаментальных величин. — *Прим. перев.*

Американский хедж-фонд Point72 Asset Management с финансовыми активами в размере 12 млрд долларов переводит около половины своих портфельных менеджеров на подход "человек плюс машина". Хедж-фонд Point72 также инвестирует десятки миллионов долларов в группу, которая анализирует большие объемы альтернативных данных и передает результаты биржевым торговцам.

Инвестиции в стратегический потенциал

Рост инвестиций во взаимосвязанный потенциал — технологии, данные и, самое главное, квалифицированных людей — подчеркивает, насколько важной стала алгоритмическая торговля с использованием МО для конкурентного преимущества, в особенности в свете растущей популярности пассивных индексированных инвестиционных механизмов, таких как биржевые фонды ETF, после финансового кризиса 2008 г.

Инвестиционный банк Morgan Stanley отметил, что только 23% его квантовых клиентов говорят, что они не рассматривают возможность использования или уже не используют МО, по сравнению с 44% в 2016 г.

Например, компания Guggenheim Partners LLC построила в Национальной лаборатории Лоуренса Беркли в Калифорнии так называемый суперкомпьютерный кластер за 1 млн долларов с целью помочь ее инвестиционным фондам перемалывать числа. Электричество для компьютеров стоит еще 1 млн долларов в год.

Квантитативная инвестиционная группа AQR опирается на научные исследования с целью выявления и систематического учета торговых факторов, которые со временем доказали свое превосходство над более широким рынком. Указанная фирма привыкла избегать ведомых чисто компьютером стратегий квантовых коллег, таких как Renaissance Technologies или DE Shaw. Однако в последнее время AQR начала искать прибыльные регулярности на рынках, используя МО, занимаясь разбором и анализом новых наборов данных, таких как спутниковые снимки теней, отбрасываемых нефтяными скважинами и танкерами.

Ведущая фирма BlackRock, с более чем 5 трлн долларов AUM, также делает ставки на алгоритмы с целью победить менеджеров дискреционных фондов, обильно инвестируя в систематическую трейдинговую фирму SAE, которую она приобрела во время финансового кризиса. Американская компания по менеджменту финансовыми активами Franklin Templeton приобрела Random Forest Capital — инвестиционную компанию, ориентированную на долги и управляемую данными — за нераскрываемую сумму, надеясь, что ее технология сможет поддержать более широкого менеджера финансовыми активами.

Машинное обучение и альтернативные данные

Хедж-фонды давно занимались поисками альфа через информационное преимущество и способностей обнаруживать новые некоррелированные сигналы. Исторически сюда входили такие вещи, как проприетарные опросы покупателей или избира-

телей перед выборами или референдумами. Время от времени использование корпоративных инсайдеров, диагностов и экспертных сетей с целью расширения знаний об отраслевых трендах или компаниях пересекает правовые грани: после 2010 г. индустрию потрясла серия судебных преследований биржевых торговцев, портфельных менеджеров и аналитиков за использование инсайдерской информации.

В отличие от этого, информационное преимущество от использования традиционных и альтернативных источников данных с помощью МО связано не с экспертными и отраслевыми сетями или доступом к корпоративному управлению, а скорее с возможностью сбора крупных объемов данных и их реально-временного анализа.

Революция в использовании данных в стратегиях алгоритмической торговли была совершена тремя трендами, которые способны еще больше сдвинуть инвестиционную отрасль от дискреционных стилей в сторону квантитативных:

- ◆ экспоненциальное увеличение объема цифровых данных;
- ◆ увеличение вычислительной мощности и емкости хранения данных по более низкой цене;
- ◆ достижения в методах МО для анализа многосложных наборов данных.

Обычные данные охватывают экономическую статистику, торговые данные или корпоративные отчеты. Альтернативные данные намного шире и охватывают такие источники, как спутниковые снимки, кредитно-карточные продажи, сентиментный анализ, данные мобильной геолокации и выискивание веб-сайтов, а также преобразование данных, генерируемых в результате обычного протекания бизнеса, в ценную развединформацию. Они охватывают, в принципе, любой источник данных, содержащий торговые сигналы, которые можно извлечь с помощью МО.

Например, данные страховой компании о продажах новых полисов автострахования не только служат индикатором объемов продаж новых автомобилей, но и могут быть разбиты на бренды или географии. Многие поставщики выискивают веб-сайты в поисках ценных данных, начиная со скачивания приложений и отзывов пользователей до бронирования билетов авиакомпаний и номеров в гостиницах. Социально-медийные сайты также могут выискиваться в поисках подсказок о мнениях потребителей и трендах.

Как правило, наборы данных являются крупными и требуют хранения, доступа и анализа с использованием масштабируемых решений для параллельной обработки данных, таких как Hadoop и Spark; согласно финансовому конгломерату Германии Deutsche Bank, существует более 1 млрд веб-сайтов с более чем 10 трлн разных веб-страниц с 500 экзабайтами (или 500 млрд гигабайт) данных. И более 100 млн веб-сайтов добавляется в Интернет с каждым годом.

Проникновение в суть перспектив компании в режиме реального времени, задолго до того, как ее результаты будут опубликованы, может быть получено из снижения списков вакансий на ее веб-сайте, внутреннего рейтинга ее руководителя сотрудниками на сайте найма Glassdoor или падения средней цены одежды на ее веб-сайте.

Это может сочетаться со спутниковыми снимками парковок и данными геолокации с мобильных телефонов, которые показывают, сколько людей посещают магазины. С другой стороны, стратегические шаги можно извлечь из скачка в распределении должностей по конкретным функциональным областям или в определенных географических районах.

К числу наиболее ценных источников относятся данные, которые непосредственно отражают потребительские расходы, причем кредитно-карточная информация является первичным источником. Эти данные дают лишь частичное представление о трендах продаж, но в сочетании с другими данными могут давать важную информацию. Компания Point72, например, анализирует 80 млн операций по кредитным картам каждый день. Мы разведаем различные источники, примеры их использования и способы их детального оценивания в *главе 3*.

За последние 2 года по мере того, как индустрия менеджмента финансовыми активами пыталась оживить свое угасающее состояние, инвестиционные группы более чем удвоили свои расходы на наборы альтернативных данных и исследователей данных. В декабре 2018 г. было 375 альтернативных поставщиков данных, приводимых на портале **alternativedata.org** (спонсируемом провайдером YipitData).

В прошлом году³⁰ менеджеры финансовыми активами в общей сложности потратили 373 млн долларов на наборы данных и наем новых сотрудников для их анализа, что на 60% больше, чем в 2016 г., и, согласно опросу инвесторов на портале **alternativedata.org**, вероятно, потратят в общей сложности 616 млн долларов в этом году. По прогнозам портала, к 2020 г. общие расходы превысят 1 млрд долларов. Некоторые оценки еще выше: консалтинговая компания Optimus оценивает, что инвесторы тратят около 5 млрд долларов в год на альтернативные данные, и ожидает, что в ближайшие годы эта индустрия будет расти на 30% в год.

По мере того, как интенсифицируется конкуренция за ценные источники данных, соглашения об эксклюзивности становятся одной из ключевых характеристик контрактов с источниками данных, преследуя цель сохранения информационного преимущества. В то же время растет озабоченность по поводу конфиденциальности, и в настоящее время регулирующие органы начали рассматривать в основном нерегулируемую отрасль поставщиков данных.

Привлечение "толпы" в разработку алгоритмов торговли

В последнее время несколько фирм по алгоритмической торговле начали предлагать инвестиционные платформы, которые обеспечивают доступ к данным и среде программирования для сбора рискованных факторов за счет прямого вовлечения людей через Интернет (краудсорсинг), которые становятся частью инвестиционной стратегии или целых алгоритмов торговли. Главные примеры включают платформы WorldQuant, Quantopian и запущенную в 2018 г. Alpha Trading Labs.

Фирма по количественному инвестиционному менеджменту WorldQuant с 2007 г. управляла более чем 5 млрд долларов для глобальной инвестиционной фирмы

³⁰ Имеется в виду 2017 год. — *Прим. ред.*

Millennium Management с 34,6 млрд долларов AUM и в 2018 г. объявила, что она запустит свой первый публичный фонд. Она использует сотни исследователей данных и еще больше частично занятых работников по всему миру в своей альфа-фабрике, которая организует инвестиционный процесс как количественный сборочный конвейер. Как утверждается, эта фабрика произвела 4 млн успешно протестированных альфа-факторов для включения в более сложные стратегии торговли и нацелена на 100 млн. Каждый альфа-фактор представляет собой алгоритм, который стремится предсказать будущее изменение цены финансового актива. Другие команды затем объединяют эти альфа-факторы в стратегии и стратегии в портфели, размещают финансовые средства между портфелями и осуществляют риск-менеджмент, избегая стратегий, которые каннибализируют друг друга.

Конструирование и исполнение стратегии торговли

МО способно наращивать добавочную стоимость на нескольких этапах жизненного цикла торговой стратегии и опирается на ключевые ресурсы инфраструктуры и данных. Таким образом, эта книга ориентирована на рассмотрение того, каким образом технические решения МО вписываются в более широкий процесс конструирования, исполнения и оценивания стратегий.

Стратегия алгоритмической торговли приводится в действие комбинацией альфа-факторов, которые преобразовывают один или несколько источников данных в сигналы, а те, в свою очередь, предсказывают будущие возвраты от финансовых активов и запускают ордера на покупку или продажу. Главы 2 и 3 охватывают вопросы получения источников данных и управления ими, которые представляют собой сырую и единственную наиболее важную движущую силу успешной стратегии торговли.

Глава 4 описывает методологически обоснованный процесс управления риском ложных обнаружений, который увеличивается с количеством данных. Глава 5 обеспечивает контекст для исполнения и оценивания результативности стратегии торговли.

На рис. 1.1 схематично показаны шаги стратегии алгоритмической торговли.

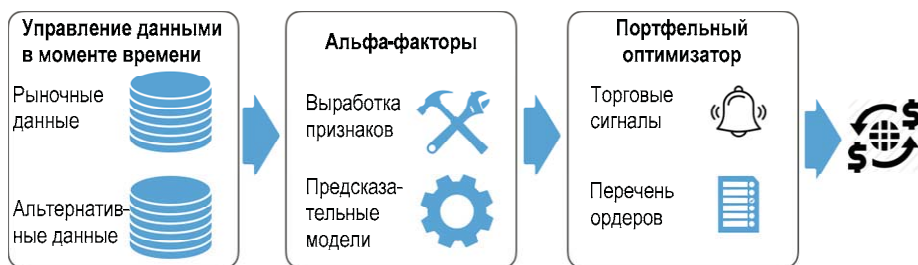


Рис. 1.1. Шаги стратегии алгоритмической торговли

Давайте кратко рассмотрим эти шаги, которые мы подробно обсудим в последующих главах.

Привлечение источников данных и управление данными

Показательная эволюция данных с точки зрения объема, разнообразия и скорости является необходимым условием и движущей силой применения МО к алгоритмической торговле. Расширяющееся предложение данных требует от активного менеджмента отождествлять потенциальную добавочную стоимость, включая следующие шаги:

1. Выявить и оценить источники рыночных, фундаментальных и альтернативных данных, содержащих альфа-сигналы, которые не затухают слишком быстро.
2. Развернуть новую либо обратиться к существующей облачной масштабируемой инфраструктуре данных и аналитическим инструментам, таким как каркасы распределенных вычислений на основе кластеров Hadoop или Spark, с целью обеспечения быстрого и гибкого доступа к данным.
3. Тщательно управлять и курировать данные во избежание систематического смещения из-за забегания вперед³¹ путем их настройки на нужную частоту на основе привязки к определенной точке во времени (point-in-time, PIT). Это означает, что данные могут отражать только имеющуюся и известную в тот момент информацию. Автоматически обучающиеся алгоритмы, натренированные на искаженных исторических данных, почти наверняка окажутся безуспешными во время живой торговли.

Исследование и оценивание альфа-факторов

Альфа-факторы призваны извлекать из данных сигналы, предсказывающие возвраты от финансовых активов для данного инвестиционного универсума на торговом горизонте. Фактор принимает одно-единственное значение для каждого финансового актива во время его оценивания, но может объединять в себе одну или несколько входных переменных. Этот процесс предусматривает шаги, описанные на рис. 1.2.

Исследовательская фаза рабочего потока торговой стратегии включает конструирование, оценивание и комбинирование альфа-факторов. МО играет большую роль в этом процессе, потому что сложность факторов увеличивалась, поскольку инвесторы реагируют и на затухание сигнала более простых факторов, и на гораздо более богатые данные, доступные сегодня.

³¹ Систематическое смещение из-за забегания вперед (look-ahead bias) — систематическая ошибка, возникающая из-за использования информации, которая не была публичной в момент принятия симулируемого решения. — *Прим. перев.*

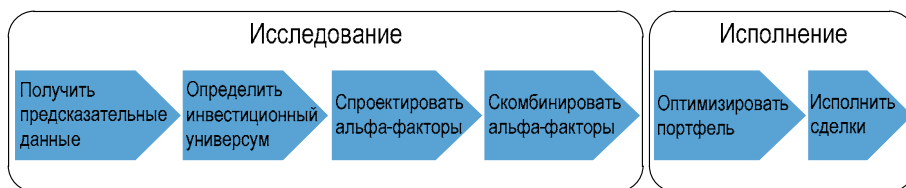


Рис. 1.2. Рабочий поток торговой стратегии

Разработка предсказательных альфа-факторов требует исследования связей между входными данными и целевыми возвратами, изобретательной выработки признаков, а также тестирования и тонкой настройки преобразований данных для оптимизации предсказательной мощности входных данных.

Преобразования данных варьируются от простых непараметрических ранжирований до сложных ансамблевых моделей или глубоких нейронных сетей в зависимости от количества сигнала на входах и сложности связей между входами и целью. Многие более простые факторы возникали в результате научных исследований, и в течение последних нескольких десятилетий их использование ширилось в финансовой индустрии.

Для минимизации рисков ложных обнаружений из-за глубинного анализа данных (data mining) и поскольку на протяжении десятилетий финансы были предметом исследований, приведших к нескольким Нобелевским премиям, инвесторы предпочитают полагаться на факторы, которые согласуются с теориями о финансовых рынках и поведении инвесторов. Изложение этих теорий выходит за рамки этой книги, но справочные материалы по ним покажут направления более глубокого погружения в этот важный рамочный аспект стратегий алгоритмической торговли.

С целью подтверждения сигнального содержимого кандидатного альфа-фактора необходимо получить робастную оценку его предсказательной мощности в средах, репрезентативных для рыночного режима³², в течение которого данный фактор будет использоваться в стратегии. Надежные оценки требуют избегания многочисленных методологических и практических ловушек, включая использование данных, которое индуцирует систематические смещения из-за выживших³³ или из-за забегания вперед, не отражающее реалистичную информацию с привязкой к определенной точке во времени, либо неспособность вносить поправки на систематическое смещение вследствие многочисленных тестов на одних и тех же данных.

Сигналы, выводимые из альфа-факторов, часто являются слабыми, но достаточно мощными в сочетании с другими факторами или источниками данных, например, модулируя сигнал в зависимости от рыночного или экономического контекста.

³² Рыночный режим (market regime) — это расширенный период, где специфическая комбинация активов доминирует над ценовой динамикой целевого инструмента. — Прим. перев.

³³ Систематическое смещение из-за выживших (survivorship bias) — ситуация, когда в качестве инвестиционного универсума используется текущий, и, следовательно, игнорируется тот факт, что некоторые компании могли обанкротиться, и ценные бумаги в результате были исключены из списка. — Прим. перев.

Оптимизация портфеля и риск-менеджмент

Альфа-факторы эмитируют входные и выходные сигналы, которые приводят к ордерам на покупку или продажу, а исполнение ордеров приводит к элементам содержимого портфеля. Рисковые профили отдельных позиций взаимодействуют, создавая рисковый профиль конкретного портфеля. Портфельный менеджмент предусматривает оптимизацию весов позиций с целью достижения желаемого портфельного риска и возвращает профиль, который согласуется с совокупными инвестиционными целями. Этот процесс является весьма динамичным, включая постоянно эволюционирующие рыночные данные.

Исполнение сделок в ходе этого процесса требует уравнивания дилеммы биржевого торговца: быстрое исполнение имеет тренд на подстегивание издержек из-за влияния на финансовый рынок, в то время как медленное исполнение может создавать дефицит реализации, когда реализованная цена отклоняется от цены, которая преобладала во время принятия решения. Риск-менеджмент осуществляется на протяжении всего процесса портфельного менеджмента, внося поправки в элементы содержимого портфеля или прибегая к хеджу³⁴ в зависимости от наблюдаемых или предсказанных изменений в рыночной среде, влияющих на рисковый профиль портфеля.

Бэктестирование стратегии

Встраивание инвестиционной идеи в алгоритмическую стратегию требует обширного тестирования с научным подходом, который пытается отклонить идею, основываясь на показателях ее результативности в альтернативных вневыборочных рыночных сценариях. Тестирование может предусматривать симулированные данные с целью улавливания сценариев, считающихся возможными, но не отраженных в исторических данных.

Механизм бэктестирования³⁵ стратегии должен реалистично симулировать исполнение стратегии, достигая несмещенных оценочных показателей результативности и риска. В дополнение к потенциальным систематическим смещениям, вносимым данными или небезупречным использованием статистики, бэктестовый механизм должен точно представлять практические аспекты оценивания торговых сигналов, размещения ордеров и исполнения в соответствии с рыночными условиями.

³⁴ Хеджирование (hedging) — это действие по снижению риска для портфеля или торговли. Например, если свечной график предполагает, что рынок биткойна очень нерешителен, и торговец считает, что цена собирается подняться, то он может купить еще несколько BTC, одновременно торгуя коррелированным активом (например, ETH). — *Прим. перев.*

³⁵ Бэктестирование (backtesting), или ретроактивное тестирование, ретротестирование, — это тестирование на исторических данных с целью получения результатов и анализа риска и возвратности инвестиций, прежде чем рисковать любым фактическим капиталом. — *Прим. перев.*

Машинное обучение и стратегии алгоритмической торговли

Квантитативные стратегии проходили процесс эволюции и становились изощреннее в три волны.

1. В 1980–90-х годах сигналы часто выявлялись в результате научных исследований и использовали единственный вход или очень мало входов, выводимых из рыночных и фундаментальных данных. Эти сигналы сейчас в основном стали коммодизированными, т. е. потеряли осмысленную дифференциацию, и доступны как биржевые фонды ETF. Это такие сигналы, как базовые стратегии разворота к среднему ценовому уровню.
2. В 2000-х годах процветало факторное инвестирование. Фонды использовали алгоритмы для выявления финансовых активов, находящихся под влиянием рисков факторов, таких как стоимость или импульс, в поисках арбитражных возможностей. Погашения в первые дни финансового кризиса вызвали квантовую встряску в августе 2007 г., которая каскадно пронеслась по индустрии факторно-ориентированных фондов. Эти стратегии теперь также доступны в виде только длинных умных бета-фондов, которые накренивают портфели в соответствии с заданным набором рисков факторов.
3. Третья эра диктуется инвестициями в потенциал МО и альтернативными данными с целью генерирования прибыльных сигналов для повторяемых стратегий торговли. Затухание факторов является серьезной проблемой: как показывает практика, к моменту публикации избыточные финансовые возвраты от новых аномалий падают на четверть с момента обнаружения и более чем на 50% после публикации, из-за конкуренции и скученности.

Существует несколько категорий стратегий торговли, в которых для исполнения торговых правил используются алгоритмы:

- ◆ краткосрочные сделки, которые ориентированы на получение прибыли от небольших ценовых движений, например из-за арбитража;
- ◆ поведенческие стратегии, ориентированные на извлечение прибыли от предвосхищения поведения других участников рынка;
- ◆ программы, ориентированные на оптимизацию исполнения сделок;
- ◆ большая группа торговой активности на основе предсказанной оценки стоимости активов.

Рассмотренные выше фонды со стратегиями высокочастотной торговли (HFT) весьма заметно опираются на короткие периоды владения финансовыми активами, извлекая выгоды из незначительных ценовых изменений, основываясь на арбитраже цен спроса-предложения (bid-ask) или статистическом арбитраже. Поведенческие алгоритмы обычно работают в условиях низкой ликвидности и призваны предвосхищать возможность существенного влияния на цену шагов более крупного

игрока. Ожидание воздействия на цену основано на алгоритмах-снифферах (вынюхивающих алгоритмах), которые генерируют проникновение в сущность стратегий других участников рынка, либо на рыночных регулярностях, таких как вынужденные сделки, осуществляемые биржевыми фондами ETF.

Программы исполнения сделок призваны ограничивать влияния сделок на финансовый рынок и варьируются от простого деления сделок на части до заключения сделок на определенную сумму или объем на регулярной основе методом *ценообразования, средневзвешенного по времени* (time-weighted average pricing, TWAP), и *ценообразования, средневзвешенного по объему* (volume-weighted average pricing, VWAP). Простые алгоритмы привлекают исторические регулярности, в то время как более изощренные алгоритмы учитывают транзакционные издержки, дефицит реализации или предсказанные ценовые движения. Эти алгоритмы могут работать на уровне ценных бумаг или портфеля, например, для реализации сделок с многими деривативами или со срезами (перекрестными сечениями) активов.

Примеры использования машинного обучения для торговли на финансовых рынках

МО извлекает сигналы из широкого спектра рыночных, фундаментальных и альтернативных данных и может применяться на всех этапах процесса стратегии алгоритмической торговли. Ключевые применения охватывают:

- ◆ глубинный анализ данных для выявления регулярности и извлечения признаков;
- ◆ контролируемое обучение для генерирования рисковых факторов или альфа и создания торговых идей;
- ◆ агрегирование отдельных сигналов в стратегию;
- ◆ размещение финансовых средств среди портфельных активов в соответствии с усвоенными алгоритмом рисковыми профилями;
- ◆ тестирование и оценивание стратегий, в том числе посредством использования синтетических данных;
- ◆ интерактивная, автоматизированная детализация стратегии с использованием подкрепляемого обучения.

Мы кратко остановимся на некоторых из этих приложений и определим точки, где продемонстрируем их использование в последующих главах.

Глубинный анализ данных для извлечения признаков

Затратно-эффективное оценивание крупных, многосложных наборов данных требует обнаружения сигналов в широком масштабе. В данной книге приведено несколько примеров.

- ◆ Теория информации является полезным инструментом извлечения признаков, который улавливает потенциальные сигналы и может использоваться в автома-

тически обучающихся моделях. В *главе 4* мы используем взаимную информацию и анализируем потенциальные значения индивидуальных признаков для алгоритма контролируемого обучения, который предсказывает возвраты от активов.

- ◆ В *главе 12* мы представляем различные технические решения для создания признаков из высокоразмерных наборов данных. В *главе 14* мы применяем эти технические решения к текстовым данным.
- ◆ Мы подчеркиваем модельно-специфичные способы получения сущностного понимания предсказательной мощности отдельных переменных. Мы используем новый теоретико-игровой подход, именуемый *аддитивными объяснениями Шепли* (SHapley Additive exPlanations, SHAP), закрепляя предсказательную результативность за индивидуальными признаками в сложных градиентно-бустинговых машинах с большим числом входных переменных.

Контролируемое обучение для создания и агрегирования альфа-факторов

Главным обоснованием применения МО к торговле на финансовых рынках является получение предсказаний фундаментальных показателей активов, ценовых движений или рыночных условий. Стратегия может привлекать несколько автоматически обучающихся алгоритмов, которые надстраиваются друг над другом. Нисходящие модели могут генерировать сигналы на уровне портфеля путем интегрирования предсказаний о перспективах отдельных активов, об ожиданиях рынка капитала и о корреляции между ценными бумагами. В качестве альтернативы предсказания автоматически обучающихся алгоритмов могут информировать дискреционные сделки, как в описанном выше квантоментальном подходе. Предсказания автоматически обучающихся алгоритмов могут также ориентироваться на конкретные рискованные факторы, такие как стоимость или волатильность, либо применять технические подходы, такие как следование за трендом или разворот к среднему ценовому уровню³⁶.

- ◆ В *главе 3* мы проиллюстрируем способы работы с фундаментальными данными по созданию входов в модели оценивания стоимости, приводимых в действие машинным обучением.
- ◆ В *главах 13–15* мы используем альтернативные данные отзывов о деятельности предприятий, которые могут применяться для проецирования выручки компании в будущее, в качестве входа в упражнение по оцениванию стоимости.

³⁶ Разворот к среднему ценовому уровню (mean reversion) — это концепция, используемая в финансах, которая предполагает, что цены финансовых активов и историческая возвратность в конечном итоге возвращаются к долгосрочному среднему значению или усредненному уровню всего набора данных (любая цена, которая отклоняется далеко от долгосрочной нормы, снова развернется, вернувшись в свое понятное состояние). Это среднее или усредненное значение может быть историческим средним значением цены или возвратности или другим соответствующим средним значением, таким как рост экономики или средняя возвратность отрасли.

См. <https://www.investopedia.com/terms/m/meanreversion.asp>. — Прим. перев.

- ♦ В *главе 8* мы демонстрируем способы прогнозирования макропеременных в качестве входов в рыночные ожидания и способы прогнозирования рисков факторов, таких как волатильность.
- ♦ В *главе 19* мы вводим *рекуррентные нейронные сети* (RNN-сети), которые достигают превосходящей результативности с данными нелинейных временных рядов.

Размещение финансовых средств среди портфельных активов

МО используется для размещения финансовых средств среди портфельных активов³⁷, применяя модели деревьев решений, которые вычисляют иерархическую форму паритета риска. Как результат, характеристики риска стимулируются не классами финансовых активов, а регулярностями в ценах финансовых активов, и достигают превосходящих показателей компромисса между риском и возвратностью³⁸.

В *главах 5 и 12* мы проиллюстрируем, как иерархическая кластеризация извлекает ведомые данными классы риска, которые лучше отражают корреляционные регулярности, чем обычное определение классов финансовых активов.

Тестирование торговых идей

Бэктестирование является критически важным шагом, позволяющим отбирать успешные стратегии алгоритмической торговли. Перекрестный контроль (cross validation) с использованием синтетических данных является ключевым техническим решением МО для генерирования надежных вневыборочных результатов в сочетании с надлежащими методами внесения поправок на множественное тестирование. Природа временных рядов финансовых данных требует внесения модификаций в стандартный подход во избежание систематического смещения из-за забегания вперед или иного загрязнения данных, используемых для тренировки, проверки и тестирования. Кроме того, ограниченная доступность исторических данных

³⁷ Размещение финансовых средств среди портфельных активов (asset allocation) — инвестиционная стратегия, призванная балансировать риск и вознаграждение путем выделения и перераспределения средств по портфельным активам в соответствии с целями, рискоустойчивостью и инвестиционным горизонтом человека. Три основных класса активов — долевые ценные бумаги (акции), активы с фиксированным доходом, денежные средства и их эквиваленты — имеют разные уровни риска и возвратности, поэтому каждый из них будет вести себя по-разному с течением времени.

См. <https://www.investopedia.com/terms/a/assetallocation.asp>. — Прим. перев.

³⁸ Компромисс между риском и возвратностью (risk-return trade-off) гласит, что потенциальная возвратность возрастает с увеличением риска. Используя этот принцип, люди связывают низкий уровень неопределенности с низкой потенциальной возвратностью, а высокий уровень неопределенности или риска с высокой потенциальной возвратностью. Согласно компромиссу "риск-возвратность", вложенные деньги могут давать более высокую возвратность только в том случае, если инвестор примет более высокую вероятность убытков.

См. <https://www.investopedia.com/terms/r/riskreturntradeoff.asp>. — Прим. перев.

привела к появлению альтернативных подходов, использующих синтетические данные:

- ◆ мы продемонстрируем различные методы тестирования автоматически обучающихся моделей с использованием рыночных, фундаментальных и альтернативных данных, которые позволяют получать обоснованные оценки вневыборочных ошибок.
- ◆ в *главе 20* мы представим GAN-сети, которые способны производить высококачественные синтетические данные.

Подкрепляемое обучение

Финансовая торговля происходит на конкурентном, интерактивном рынке. Подкрепляемое обучение, или автоматическое обучение на основе максимизации получаемого подкрепления, призвано тренировать агентов, которые усваивают функцию политики, основываясь на получаемых вознаграждениях.

В *главе 21* мы представим ключевые алгоритмы подкрепляющего обучения, такие как Q-обучение и архитектура Дупа, и продемонстрируем тренировку алгоритмов подкрепляемого обучения для торговли на финансовых рынках с использованием среды OpenAI Gym.

Резюме

В этой главе мы познакомили вас со стратегиями алгоритмической торговли и проследили за тем, как МО стало ключевым ингредиентом в конструировании и комбинировании альфа-факторов, которые, в свою очередь, являются ключевыми движущими силами результативности инвестиционного портфеля. Мы рассмотрели различные отраслевые тренды вокруг стратегий алгоритмической торговли, проследили появление альтернативных данных и использование МО для эксплуатации этих новых источников информационных преимуществ.

Кроме того, мы познакомились с процессом конструирования стратегии алгоритмической торговли, важными типами альфа-факторов и тем, как мы будем использовать МО для конструирования и исполнения наших стратегий. В последующих двух главах мы более подробно рассмотрим нефть, которая питает любую стратегию алгоритмической торговли — рыночные, фундаментальные и альтернативные источники данных — с использованием МО.

2

Рыночные и фундаментальные данные

Данные всегда были существенной движущей силой торговли на финансовых рынках, и биржевые торговцы давно прилагали усилия к тому, чтобы получать преимущество, имея доступ к опережающей информации. Эти усилия восходят, по крайней мере, к слухам о том, что дом Ротшильдов получил значительную выгоду от покупки облигаций после заблаговременного извещения о победе британцев при Ватерлоо, которое через пролив Ла-Манш принесли голуби.

Сегодня инвестиции в более быстрый доступ к данным принимают форму консорциума Go West ведущих фирм *высокочастотной торговли* (HFT), который соединяет *Чикагскую товарную биржу* (Chicago Mercantile Exchange, CME) с Токио. Временная задержка полного цикла между CME и фондовой биржей BATS в Нью-Йорке сократилась до теоретического предела 8 миллисекунд вследствие того, что биржевые торговцы конкурируют между собой за эксплуатацию арбитражных возможностей.

Традиционно инвестиционные стратегии в основном опираются на публично доступные данные при ограниченных усилиях по созданию или приобретению частных наборов данных. В случае долевого ценного бумага фундаментальные стратегии использовали финансовые модели, строившиеся на отчетных финансовых данных, возможно, в сочетании с отраслевыми данными либо макроданными. Стратегии, мотивируемые техническим анализом, извлекают сигналы из рыночных данных, таких как цены и объемы.

Автоматически обучающиеся алгоритмы могут эффективнее эксплуатировать рыночные и фундаментальные данные, в частности в сочетании с альтернативными данными, что является темой следующей главы. В последующих главах мы рассмотрим несколько технических решений, которые фокусируются на рыночных и фундаментальных данных, таких как классические и современные технические решения по работе с временными рядами, включая *рекуррентные нейронные сети* (RNN-сети).

В этой главе представлены источники рыночных и фундаментальных данных, а также среда, в которой они создаются. Осведомленность о различных типах orde-

ров и торговой инфраструктуре важна, поскольку они влияют на бэкстестовое симулирование торговой стратегии. Мы также проиллюстрируем способы применения языка Python для доступа и работы с торговыми и финансово-отчетными данными.

В частности, в этой главе будут рассмотрены следующие темы:

- ◆ как рыночная микроструктура формирует рыночные данные;
- ◆ как реконструировать ордерную книгу из тиковых данных с помощью протоколаITCH биржи NASDAQ;
- ◆ как резюмировать тиковые данные, используя различные типы баров;
- ◆ как работать с электронными документами официальной отчетности, кодированными на *расширяемом языке деловой отчетности* XBRL (eXtensible Business Reporting Language);
- ◆ как разбирать, анализировать и комбинировать рыночные и фундаментальные данные для создания временного ряда "цена/корпоративные заработки" (Price/Earnings, P/E);
- ◆ как получать доступ к различным рынкам и источникам фундаментальных данных с помощью Python.

Как работать с рыночными данными

Рыночные данные являются результатом размещения и обработки ордеров на покупку и продажу в ходе торговли финансовыми инструментами на многочисленных рынках. Эти данные отражают институциональную среду торговых площадок, в том числе правила и регламенты, которые регулируют ордера, исполнение сделок и ценообразование.

Алгоритмические торговцы используют автоматически обучающиеся алгоритмы для анализа потока ордеров на покупку и продажу и результирующей объемной и ценовой статистики для извлечения торговых сигналов или признаков, которые улавливают сущность, например, динамики спроса и предложения или поведения определенных участников рынка.

Сначала мы проведем краткий обзор институциональных признаков, которые влияют на симулирование торговой стратегии во время бэктеста. Затем мы посмотрим на то, как тиковые данные могут быть реконструированы из источника ордерной книги¹. Далее мы выделим несколько методов, которые регуляризируют тиковые данные и ориентированы на максимизацию информационного содержимого.

¹ Ордерная книга (order book) — это электронный список заявок на покупку и продажу конкретной ценной бумаги или финансового инструмента, организованный по уровню цен. Она имеет и другие названия — *стакан заявок*, книга приказов, или глубина рынка.

<https://www.investopedia.com/terms/o/order-book.asp>. — Прим. перев.

Наконец, мы проиллюстрируем получение доступа к различным интерфейсам поставщиков рыночных данных и отметим несколько поставщиков.

Микроструктура рынка

Микроструктура рынка — это отрасль финансовой экономики, которая исследует торговый процесс и организацию смежных рынков. Ее институциональные детали довольно сложны и разнообразны по всем классам финансовых активов и их производным, торговым площадкам и географиям. Мы дадим лишь краткий обзор ключевых понятий и потом погрузимся в данные, генерируемые торговлей. Справочные материалы в репозитории GitHub отсылают к нескольким источникам, которые рассматривают эту тему во всех подробностях.

Торговые площадки

Торговля финансовыми инструментами осуществляется на организованных главным образом электронных биржах и вне бирж². Биржа является центральным рынком, где встречаются покупатели и продавцы, и покупатели конкурируют друг с другом за самую высокую цену спроса (bid), в то время как продавцы конкурируют за самую низкую цену предложения (ask)³.

В США и за рубежом существует большое число бирж и альтернативных торговых площадок. В табл. 2.1 перечислены несколько крупных глобальных бирж и объемы торгов за 12 месяцев, завершившихся в марте 2018 г. в различных классах финансовых активов, включая деривативы. Как правило, на долю меньшинства финансовых инструментов приходится подавляющая часть торговли.

Биржи могут опираться на билатеральные системы торговли, или системы, приводимые в действие ордерами, в которых ордера на покупку и продажу сопоставляются в соответствии с определенными правилами. Ценообразование может происходить на аукционах, например на *Нью-Йоркской фондовой бирже* (NYSE), где сопоставляются самые высокие и самые низкие цены, либо через дилеров, которые покупают у продавцов и продают покупателям.

² Термин "вне биржи" (over-the-counter, OTC) относится к процессу торговли ценными бумагами для компаний, которые не котируются на официальной бирже, такой как Нью-Йоркская фондовая биржа (NYSE). Ценные бумаги, которые торгуются вне биржи, торгуются не на централизованной бирже, а через дилерскую сеть. См. <https://www.investopedia.com/terms/o/otc.asp>. — Прим. перев.

³ Цена предложения (ask) — цена, которую продавец готов принять за ценную бумагу или другой финансовый инструмент. Также упоминается как предложение (offer). Цена спроса (bid), или заявка — самая высокая цена, которую любой покупатель готов заплатить за данную ценную бумагу в данный момент времени. Как правило, заявка (bid) ниже, чем предложения (ask), и разница между ними называется спредом между ценой спроса и ценой предложения.

См. <http://www.businessdictionary.com/definition/ask.html>

и <http://www.investorwords.com/469/bid.htm>. — Прим. перев.

Таблица 2.1. Перечень глобальных бирж и объемы торгов за 12 месяцев (март 2018 г.)

| Биржа | Акции | | | | | Описание |
|--|--------------------------------------|---------------------------|----------------------------|------------------------------|---------------------------------|---|
| | рыночная капитализация, млн долларов | число котируемых компаний | объем в день, млн долларов | число долей в день, тыс. шт. | число опционов в день, тыс. шт. | |
| NYSE | 23 138 626 | 2 294 | 78 410 | 6 122 | 1 546 | Нью-Йоркская фондовая биржа, расположенная на Уолл-Стрит, главная фондовая биржа США, крупнейшая в мире по обороту |
| NASDAQ — US (National Association of Securities Dealers Automated Quotation) | 10 375 718 | 2 968 | 65 026 | 7 131 | 2 609 | Автоматизированные котировки Национальной ассоциации дилеров по ценным бумагам) — американская биржа, специализирующаяся на акциях высокотехнологичных компаний (производство электроники, программного обеспечения и т. п.) |
| Japan Exchange Group Inc. | 6 287 739 | 3 618 | 28 397 | 3 361 | 1 | Японская корпорация финансовых услуг, которая управляет несколькими фондовыми биржами, включая Токийскую фондовую биржу и фондовую биржу Осаки |
| Шанхайская фондовая биржа | 5 022 691 | 1 421 | 34 736 | 9 801 | | Китайская фондовая биржа |
| Euronext | 4 649 073 | 1 240 | 9 410 | 836 | 304 | Панъевропейская фондовая биржа, имеющая филиалы в Бельгии, Франции, Нидерландах и Португалии |
| Hong Kong Exchanges and Clearing | 4 443 082 | 2 186 | 12 031 | 1 174 | 516 | Крупный финансовый холдинг, образованный в 2000 г. путем объединения компаний Гонконгской фондовой биржи, Гонконгской фьючерсной биржи (Hong Kong Futures Exchange) и Гонконгской расчетно-клиринговой компании по сделкам с ценными бумагами (Hong Kong Securities Clearing Company) |

Таблица 2.1 (окончание)

| Биржа | Акции | | | | | Описание |
|---|---|--------------------------------------|-------------------------------------|---------------------------------------|--|--|
| | рыночная капитализа- ция, млн долларов | число котируе- мых компаний | объем в день, млн долларов | число долей в день, тыс. шт. | число опционов в день, тыс. шт. | |
| LSE Group (London Stock Exchange Group) | 3 986 413 | 2622 | 10 398 | 1011 | | Британская фондовая биржа и финансовая информационная компания |
| Шэньчжэньская фондовая биржа | 3 547 312 | 2110 | 40 244 | 14 443 | | Китайская фондовая биржа |
| Deutsche Boerse AG | 2 339 092 | 506 | 7825 | 475 | | Немецкая биржа, созданная в форме акционерного общества и осуществляющая организацию работы различных рынков ценных бумаг |
| BSE India Limited | 2 298 179 | 5439 | 602 | 1105 | | Бомбейская фондовая биржа |
| National Stock Exchange of India Limited | 2 273 286 | 1952 | 5092 | 10 355 | | Национальная фондовая биржа Индии |
| BATS Global Markets — US (Better Alternative Trading System) | | | | | 1243 | Компания, управляющая тремя частными фондовыми биржами, расположенная в пригороде Канзас-Сити, штат Миссури, США |
| CBOE (Chicago Board Options Exchange) | | | | | 1811 | Чикагская биржа опционов |
| ISE (International Securities Exchange) | | | | | 1204 | Международная биржа ценных бумаг, дочерняя компания американской много- национальной корпорации финансовых услуг NASDAQ, Inc. Является электронной биржей опционов |

Многие биржи используют специальных посредников — маркетмейкеров⁴, которые обеспечивают ликвидность, т. е. возможность торговать, создавая рынки определенных ценных бумаг. Нью-Йоркская фондовая биржа (NYSE), например, обычно имеет одного назначенного маркетмейкера, который обеспечивает упорядоченную торговлю каждой ценной бумагой, в то время как биржа *NASDAQ* имеет их несколько. Посредники могут выступать в качестве дилеров, которые ведут торговлю от своего имени, или брокеров, которые ведут торговлю от имени других.

Биржи раньше принадлежали ее членам, но часто переходили в корпоративную собственность по мере того, как рыночные реформы усиливали конкуренцию. Создание биржи NYSE датируется 1792 годом, в то время как биржа *NASDAQ* начала свою деятельность в 1971 г. как первый в мире электронный фондовый рынок и взяла на себя большинство биржевых сделок, которые исполнялись вне биржи. Только на одних фондовых рынках США торговля фрагментирована по 13 биржам и более чем 40 альтернативным торговым площадкам, каждая из которых сообщает о сделках на консолидированной ленте, но с разными временными задержками.

Типы ордеров

Торговцы могут размещать различные типы ордеров на покупку или продажу. Некоторые ордера гарантируют немедленное исполнение, в то время как другие могут указывать ценовой порог или иные условия, которые запускают исполнение. Ордера, как правило, действительны в течение одного торгового дня, если не указано иное.

Рыночный ордер (*market order*) гарантирует немедленное исполнение ордера по прибытии на торговую площадку по цене, которая преобладает в данный момент. Напротив, лимитный ордер (*limit order*) исполняется только в том случае, если рыночная цена находится выше (ниже) предела лимитного ордера на продажу (покупку). Стоп-ордер (*stop order*), в свою очередь, становится активным только тогда, когда рыночная цена поднимается выше (опускается ниже) указанной в стоп-ордере цены на покупку (продажу). Стоп-ордер на покупку может использоваться для ограничения убытков коротких продаж. Стоп-ордера тоже могут иметь пределы.

За ордерами могут закрепляться многочисленные другие условия — ордера "все или ничего" (*all or none order*) не препятствуют частичному исполнению, заполня-

⁴ Маркетмейкер (*market maker*) — брокер, дилер или инвестиционная компания, которая принимает на себя рыночный риск (системный риск), вступая во владение финансовым активом и торгуя им в качестве принципала. Проще говоря, это организации, которые предоставляют ликвидность бирже путем размещения лимитных ордеров в ее книге заказов для того, чтобы сделки могли совершаться в диапазоне цен. Маркетмейкеры обязаны постоянно выдавать котировки цен спроса и предложения, а также гарантировать полную продажу или поглощение финансового актива по определенной цене. — *Прим. перев.*

ются⁵ только при наличии указанного числа акций и могут быть действительны в течение дня или дольше. Они требуют специальной обработки и не видны участникам рынка. Ордера "заполнить или уничтожить" (fill or kill order) тоже предотвращают частичное исполнение, но досрочно отменяют, в случае если они не исполняются немедленно. Ордера "исполнить немедленно или аннулировать" (immediate or cancel order) немедленно покупают или продают число акций, которое имеется в наличии, и отменяют остаток. Ордера по неустановившемуся рынку (not-held order)⁶ позволяют брокеру принимать решение о времени и цене исполнения. Наконец, ордера на открытии/закрытии рынка (market on open/close order) исполняются на открытии или закрытии рынка или вблизи этого времени. Допускаются частичные исполнения.

Работа с данными ордерной книги

Первичным источником рыночных данных является ордерная книга, или стакан заявок, которая постоянно обновляется в режиме реального времени в течение дня, отражая всю торговую деятельность. Биржи обычно предлагают эти данные в качестве реально-временной службы и могут предоставлять некоторые исторические данные бесплатно.

Торговая активность отражается в многочисленных сообщениях участников рынка о торговых ордерах. Эти сообщения, как правило, соответствуют протоколу *электронного обмена финансовой информацией* FIX (Financial Information eXchange), служащему для обмена транзакциями с ценными бумагами и рыночными данными в режиме реального времени, либо протоколу нативного обмена.

Протокол FIX

Так же, как и коммуникационный протокол SWIFT, служащий для обмена сообщениями внутри операционного офиса (например, для расчетов по сделкам), протокол FIX де-факто является стандартом обмена сообщениями для связи до и во время исполнения сделок между биржами, банками, брокерами, клиринговыми фирмами и другими участниками рынка. Холдинговая компания Fidelity Investments и инвестиционный банк Salomon Brothers ввели протокол FIX в 1992 г. с целью обеспечения электронной связи между брокерами-дилерами и институциональными клиентами, которые к тому времени обменивались информацией по телефону.

⁵ Заполнение (fill) — это действие по исполнению или удовлетворению ордера на ценную бумагу или биржевой товар. Это основной ордер при сделках с акциями, облигациями или любым другим видом ценных бумаг. Досрочное аннулирование (kill) — это запрос на отмену сделки между ее размещением и ее исполнением. — *Прим. перев.*

⁶ Ордер по неустановившемуся рынку (not-held order) — это рыночный или лимитный ордер, в котором клиент не желает автоматически совершать сделки на внутреннем рынке (установившемся рынке, market held) и дает брокеру или торговцу время и свободу выбора цены с целью получить наилучшую возможную цену. — *Прим. перев.*

Указанный протокол стал популярным на глобальных фондовых рынках до того, как распространился на рынки иностранной валюты, рынки ценных бумаг с фиксированной доходностью и деривативов, и далее на постторговый рынок для поддержки сквозной обработки. Биржи предоставляют доступ к сообщениям FIX как реально-временному каналу данных, который анализируется алгоритмическими торговцами для отслеживания рыночной активности и, например, выявления следа участников рынка и предвосхищения их следующего шага.

Последовательность сообщений позволяет реконструировать ордерную книгу. Масштаб транзакций на многочисленных биржах создает большой объем (~10 Тбайт) неструктурированных данных, которые крайне сложно обрабатывать и, следовательно, способными выступить источником конкурентного преимущества.

Протокол FIX, в настоящее время в версии 5.0, является бесплатным и открытым стандартом с большим сообществом аффилированных отраслевых профессионалов. Он самоописуем, как более поздний формат XML, и сеанс FIX поддерживается базовым управляющим протоколом передачи TCP (Transmission Control Protocol). Указанное сообщество постоянно добавляет новый функционал.

Указанный протокол поддерживает пары "ключ — значение", разделенные символом вертикальной черты, а также теговый синтаксис FIXML. Образец сообщения, запрашивающего серверный логин, выглядит следующим образом:

```
8=FIX.5.0|9=127|35=A|59=theBroker.123456|56=CSERVER|34=1|32=20180117-08:03:04|57=TRADE|50=any_string|98=2|108=34|141=Y|553=12345|554=passwd|10=131|
```

Существует несколько общедоступных реализаций протокола FIX на языке Python, которые можно использовать для формулирования и разбора сообщений FIX. Брокерская фирма Interactive Brokers предлагает *интерфейс межмашинной связи* "компьютер — компьютер" (computer-to-computer Interface, CTCI) на основе FIX для автоматической торговли (см. раздел ресурсов для этой главы в репозитории GitHub).

Данные ордерной книги TotalView-ITCH биржи Nasdaq

Хотя протокол FIX занимает подавляюще большую долю рынка, вдобавок к нему биржи предлагают нативные протоколы. Биржа NASDAQ предлагает протокол прямой передачи данных TotalView-ITCH, который позволяет абонентам отслеживать индивидуальные ордера на долевые финансовые инструменты от размещения до исполнения или аннулирования.

В результате можно реконструировать ордерную книгу, которая отслеживает список активных лимитных ордеров на покупку и продажу конкретной ценной бумаги или финансового инструмента. Ордерная книга показывает глубину рынка в течение дня, перечисляя число акций, которое запрашивается или предлагается в каждой ценовой точке. Она также может выявлять участника рынка, ответственного за конкретные ордера на покупку и продажу, если только он не размещен анонимно.

Глубина рынка⁷ — это ключевой индикатор ликвидности и потенциального влияния на цену значительных рыночных ордеров.

Помимо сопоставления рыночных и лимитных ордеров, биржа NASDAQ также проводит аукционы или кроссы, которые исполняют большое число сделок на открытии и закрытии рынка⁸. Кроссы приобретают все большую важность, поскольку пассивное инвестирование продолжает расти, и торговцы ищут возможности для исполнения более крупных блоков акций. Протокол TotalView также распространяет *индикатор чистого дисбаланса ордеров (NOI)*⁹ для кроссов на открытия и закрытия и кросса IPO/Halt биржи NASDAQ.

Разбор двоичных сообщенийITCH

Спецификация ITCH v5.0 объявляет более 20 типов сообщений, связанных с системными событиями, характеристиками акций, размещением и модификацией лимитных ордеров и исполнением сделок. Она также содержит информацию о чистом дисбалансе ордеров перед кроссом на открытии и закрытии.

Биржа NASDAQ предлагает образцы дневных двоичных файлов за несколько месяцев. Репозиторий GitHub для этой главы содержит блокнот Jupyter `build_order_book.ipynb`, который иллюстрирует, как разбирать и анализировать образец файла сообщений ITCH и реконструировать исполненные сделки и ордерную книгу для любого заданного тика.

В табл. 2.2 показана частота наиболее распространенных типов сообщений для образца файла, датированного 27 марта 2019 г.

Для каждого сообщения в спецификации указаны компоненты, их соответствующая длина и типы данных (табл. 2.3).

⁷ Глубина рынка (market depth) — это способность рынка конкретного актива поддерживать крупные ордера на данный актив без существенного движения цены актива. Чем больше открытых лимитных ордеров на обеих сторонах книги ордеров биржи, тем больше глубина этой книги. Глубина также тесно связана с ликвидностью: чем глубже книга ордеров для актива, тем больше ликвидности книга ордеров предоставляет этому активу. — *Прим. перев.*

⁸ Кросс (cross) — это данные, которые биржа NASDAQ собирает и публикует по всем заинтересованностям на покупку и продажу за две минуты до своего открытия или закрытия. Они называются кроссом на открытии (opening cross) или кроссом на закрытии (closing cross). Такое распространение информации о заинтересованности в цене помогает ограничить сбой в ликвидности на открытии биржи. Кросс на закрытии биржи сопоставляет запросы и предложения по данной акции с целью создать окончательную цену дня.

См. <https://www.investopedia.com/terms/n/net-order-imbalance-indicator-noii.asp>. — *Прим. перев.*

⁹ Индикатор чистого дисбаланса ордеров (Net Order Imbalance Indicator, NOI) — это информация о кроссах на открытие и закрытие на фондовом рынке NASDAQ, предоставляемая пользователям рынка перед исполнением кроссов. Индикатор NOI показывает истинное предложение и спрос на акции на основе фактических ордеров на покупку и продажу за 10 минут до закрытия рынка и за 5 минут до открытия рынка.

См. <https://www.investopedia.com/terms/n/net-order-imbalance-indicator-noii.asp>. — *Прим. перев.*

Таблица 2.2. Частота наиболее распространенных типов сообщений ITCH

| Тип сообщения | Влияние ордерной книги | Число сообщений |
|---------------|---|-----------------|
| A | Новый неатрибутированный лимитный ордер | 186 296 811 |
| D | Ордер аннулирован | 181 953 144 |
| U | Ордер аннулирован и заменен | 34 555 656 |
| E | Полное или частичное исполнение; возможно, многочисленные сообщения для одного и того же исходного ордера | 7 331 047 |
| X | Модифицирован после частичного аннулирования | 5 495 709 |
| F | Добавлен атрибутированный ордер | 1 477 447 |
| P | Сообщение о сделке (не кросс-сделке) | 1 105 314 |
| C | Исполнен полностью или частично по цене, отличной от первоначальной отображаемой цены | 137 556 |
| Q | Сообщение о кросс-сделке | 17 445 |

Таблица 2.3. Компоненты сообщения ITCH и их соответствующая длина и типы данных

| Имя | Сдвиг | Длина | Значение | Примечания |
|---|-------|-------|---------------|---|
| Тип сообщения (message type) | 0 | 1 | F | Добавить сообщение об атрибуции MPID-ордера |
| Локализация акции (stock locate) | 1 | 2 | Целочисленное | Код локализации, идентифицирующий ценную бумагу |
| Номер для отслеживания (tracking number) | 3 | 2 | Целочисленное | Внутренний для NASDAQ трек-номер |
| Временная метка (timestamp) | 5 | 6 | Целочисленное | Наносекунды с полуночи |
| Ссылочный номер ордера (order reference number) | 11 | 8 | Целочисленное | Уникальный ссылочный номер нового ордера |
| Индикатор покупки/продажи (buy/sell indicator) | 19 | 1 | Альфа | Тип ордера: B — ордер на покупку, S — ордера на продажу |
| Доли (паи) собственности (shares) | 20 | 4 | Целочисленное | Число долей (долевых ценных бумаг) для ордеров, добавляемых в книгу |
| Акция (stock) | 24 | 8 | Альфа | Символ акции, дополняемый пробелами справа |

Таблица 2.3 (окончание)

| Имя | Сдвиг | Длина | Значение | Примечания |
|-------------------------|-------|-------|----------|---|
| Цена (price) | 32 | 4 | Цена (4) | Отображаемая (дисплейная) цена нового ордера |
| Атрибуция (attribution) | 32 | 4 | Альфа | Идентификатор участника рынка биржи NASDAQ, связанный с ордером |

Язык Python предоставляет модуль `struct` для разбора двоичных данных с помощью форматных строковых значений, которые отождествляют элементы сообщения, указывая длину и тип различных компонентов байтовой строки, как указано в спецификации.

Давайте пройдемся по критически важным шагам разбора и анализа сообщений о сделках и реконструкции ордерной книги.

1. Анализатор протоколаITCHN опирается на спецификации сообщений, предоставляемые в виде CSV-файла (созданного сценарием `create_message_spec.py`) и собирает форматные строки в соответствии со словарем `formats`:

```
formats = {
    ('integer', 2): 'H', # целое длиной 2 => форматная строка 'H'
    ('integer', 4): 'I',
    ('integer', 6): '6s', # целое длиной 6 => разобрать как строку,
                        # конвертировать позже
    ('integer', 8): 'Q',
    ('alpha', 1) : 's',
    ('alpha', 2) : '2s',
    ('alpha', 4) : '4s',
    ('alpha', 8) : '8s',
    ('price_4', 4): 'I',
    ('price_8', 8): 'Q',
}
```

2. Анализатор транслирует спецификации сообщений в форматные строки и именованные кортежи `namedtuple`, которые захватывают содержимое сообщения:

```
# Получить спецификации ITCH и создать кортежи в формате (тип, длина)
message_types = pd.read_csv('message_types.csv')
message_types.loc[:, 'formats'] = (message_types[['value', 'length']]
                                   .apply(tuple, axis=1).map(formats))

# Форматирование полей альфа
alpha_fields = message_types[message_types.value == 'alpha'].set_index('name')
alpha_msgs = alpha_fields.groupby('message_type')
alpha_formats = {k: v.to_dict() for k, v in alpha_msgs.formats}
alpha_length = {k: v.add(5).to_dict() for k, v in alpha_msgs.length}
```

```
# Генерирование классов сообщений в виде именованных кортежей и форматных строк
message_fields, fstring = {}, {}
for t, message in message_types.groupby('message_type'):
    message_fields[t] = namedtuple(typename=t, field_names=message.name.tolist())
    fstring[t] = '>' + ''.join(message.formats.tolist())
```

3. Поля с типом `alpha` требуют постобработки, как определено в функции `format_alpha`:

```
def format_alpha(mtype, data):
    for col in alpha_formats.get(mtype).keys():
        if mtype != 'R' and col == 'stock': # имя акции только
            # в сводном сообщении 'R'
            data = data.drop(col, axis=1)
            continue
        data.loc[:, col] = data.loc[:, col].str.decode("utf-8").str.strip()

    if encoding.get(col): # целочисленное кодирование
        data.loc[:, col] = data.loc[:, col].map(encoding.get(col))
    return data
```

4. Двоичный файл за один день содержит более 300 млн сообщений в объеме более 9 Гбайт. Этот сценарий итеративно добавляет разобранный результат в файл быстрого формата HDF5 во избежание ограничений по памяти (см. последний раздел этой главы, где даны дополнительные сведения об этом формате). Следующий ниже (упрощенный) фрагмент кода обрабатывает двоичный файл и производит разобранные ордера, сохраняемые по типу сообщения:

```
with file_name.open('rb') as data:
    while True:
        # определить размер сообщения в байтах
        message_size = int.from_bytes(data.read(2), byteorder='big', signed=False)
        # получить тип сообщения, прочитав первый байт
        message_type = data.read(1).decode('ascii')
        # создать структуру данных для улавливания результата
        if not messages.get(message_type):
            messages[message_type] = []

        message_type_counter.update([message_type])

        # прочитать и сохранить сообщение
        record = data.read(message_size - 1)
        message = message_fields[message_type]._make(unpack(
            fstring[message_type], record))
        messages[message_type].append(message)

        # обработать системные события, такие как открытие/закрытие рынка
        if message_type == 'S':
            timestamp = int.from_bytes(message.timestamp, byteorder='big')
```

```

print('\n',
      event_codes.get(message.event_code.decode('ascii'), 'Error'))
print('\t{0}\t{1:,.0f}'.format(timedelta(seconds=timestamp * 1e-9),
                                message_count))
if message.event_code.decode('ascii') == 'C':
    store_messages(messages)
    break

message_count += 1
if message_count % 2.5e7 == 0:
    timestamp = int.from_bytes(message.timestamp, byteorder='big')
    print('\t{0}\t{1:,.0f}\t{2}'.format(timedelta(seconds=timestamp * 1e-9),
                                        message_count,
                                        timedelta(seconds=time() - start)))

    store_messages(messages)
    messages = {}

```

5. Как и ожидалось, на малое число из более чем 8500 долевых ценных бумаг, торговавшихся в этот день, приходится большая часть сделок:

```

with pd.HDFStore(itch_store) as store:
    stocks = store['R'].loc[:, ['stock_locate', 'stock']]
    trades = store['P'].append(store['Q'].rename(columns={'cross_price': 'price'}),
                               sort=False).merge(stocks)

trades['value'] = trades.shares.mul(trades.price)
trades['value_share'] = trades.value.div(trades.value.sum())
trade_summary = trades.groupby('stock').value_share.sum().sort_values(ascending=False)
trade_summary.iloc[:50].plot.bar(figsize=(14, 6), color='darkblue',
                                  title='Доля торговавшейся стоимости')

plt.gca().yaxis.set_major_formatter(FuncFormatter(lambda y, _: '{:,.0f}'.format(y)))

```

Мы получаем следующий график (рис. 2.1).

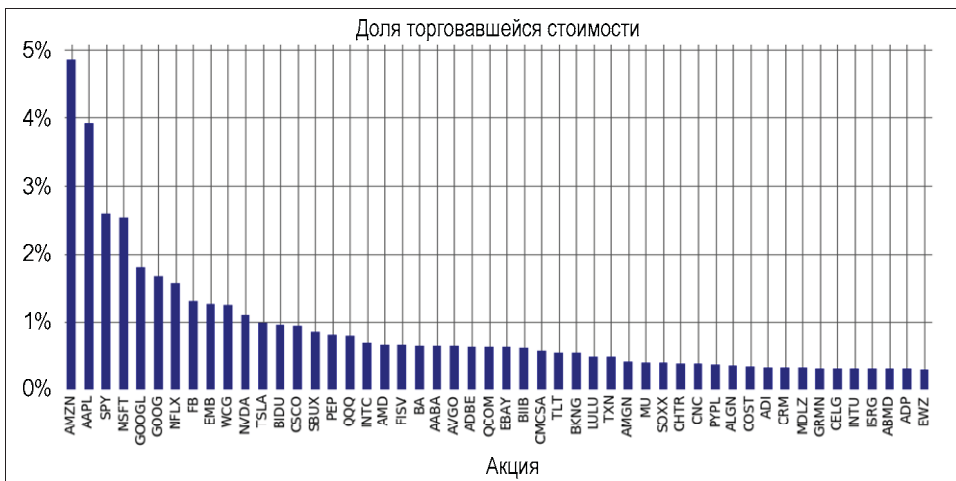


Рис. 2.1. Гистограмма торговавшейся стоимости по видам акций

Реконструирование сделок и ордерной книги

Разобранные сообщения позволяют воссоздать поток ордеров на тот или иной день. Тип сообщения 'R' содержит список всех акций, торговавшихся в течение того или иного дня, включая информацию о *первичных публичных размещениях*¹⁰ (IPO) и торговых ограничениях.

В течение дня в ордерную книгу добавляются новые ордера, а исполненные и аннулированные ордера из нее удаляются. Правильный учет сообщений, которые ссылаются на ордера, размещенные на предыдущую дату, требует отслеживания ордерной книги за несколько дней, но мы здесь этот аспект проигнорируем.

Функция `get_messages()` иллюстрирует, как собирать ордера для одиночной акции, которая влияет на торговлю (обратитесь к спецификации ИТСН за подробностями по каждому сообщению; немного в упрощенном виде см. блокнот):

```
def get_messages(date, stock=stock):
    """Собрать сообщения о сделках для заданной акции"""

    with pd.HDFStore(itch_store) as store:
        stock_locate = store.select('R',
                                    where='stock = stock').stock_locate.iloc[0]
        target = 'stock_locate = stock_locate'

        data = {}

        # Типы торговых сообщений
        messages = ['A', 'F', 'E', 'C', 'X', 'D', 'U', 'P', 'Q']
        for m in messages:
            data[m] = store.select(m,
                                   where=target).drop('stock_locate',
                                                       axis=1).assign(type=m)

        order_cols = ['order_reference_number',
                      'buy_sell_indicator', 'shares', 'price']
        orders = pd.concat([data['A'], data['F']], sort=False,
                           ignore_index=True).loc[:, order_cols]

        for m in messages[2: -3]:
            data[m] = data[m].merge(orders, how='left')

        data['U'] = data['U'].merge(orders, how='left',
                                    right_on='order_reference_number',
                                    left_on='original_order_reference_number',
                                    suffixes=('', '_replaced'))
```

¹⁰ Первичное публичное размещение (initial public offering, IPO) — самое первое предложение акций частной корпорации для общественности. — *Прим. перев.*


```

data['Q'].rename(columns={'cross_price': 'price'}, inplace=True)
data['X']['shares'] = data['X']['cancelled_shares']
data['X'] = data['X'].dropna(subset=['price'])

data = pd.concat([data[m] for m in messages],
                  ignore_index=True, sort=False)
data['date'] = pd.to_datetime(date, format='%m%d%Y')
data.timestamp = data['date'].add(data.timestamp)
data = data[data.printable != 0]

drop_cols = ['tracking_number', 'order_reference_number',
             'original_order_reference_number', 'cross_type',
             'new_order_reference_number', 'attribution',
             'match_number', 'printable', 'date', 'cancelled_shares']

return data.drop(drop_cols,
                  axis=1).sort_values('timestamp').reset_index(drop=True)

```

Реконструировать успешные сделки, т. е. ордера, которые исполняются в отличие от тех, которые были аннулированы, из связанных со сделкой типов сообщений C, E, P и Q, относительно просто:

```

def get_trades(m):
    """Скомбинировать сообщения C, E, P и Q в торговые записи"""

    trade_dict = {'executed_shares': 'shares', 'execution_price': 'price'}
    cols = ['timestamp', 'executed_shares']
    trades = pd.concat([m.loc[m.type == 'E',
                             cols + ['price']].rename(columns=trade_dict),
                        m.loc[m.type == 'C',
                             cols +
                             ['execution_price']].rename(columns=trade_dict),
                        m.loc[m.type == 'P', ['timestamp', 'price', 'shares']],
                        m.loc[m.type == 'Q', ['timestamp', 'price',
                                                'shares']].assign(cross=1),
                        ], sort=False).dropna(subset=['price']).fillna(0)

    return trades.set_index('timestamp').sort_index().astype(int)

```

Ордерная книга отслеживает лимитные ордера, при этом различные ценовые уровни для ордеров на покупку и продажу составляют глубину ордерной книги. Реконструирование ордерной книги для заданного уровня глубины требует следующих шагов:

1. Функция `add_orders()` накапливает ордера на продажу по возрастанию и ордера на покупку по убыванию для заданной временной метки до желаемого уровня глубины:

```
def add_orders(orders, buysell, nlevels):
    """Добавлять ордера вплоть до нужной глубины, заданной аргументом
       nlevels.
       Продавать в возрастающем, покупать в убывающем порядке"""

    new_order = []
    items = sorted(orders.copy().items())
    if buysell == 1:
        items = reversed(items)
    for i, (p, s) in enumerate(items, 1):
        new_order.append((p, s))
        if i == nlevels:
            break

    return orders, new_order
```

2. Мы перебираем все сообщения ИТСН и обрабатываем ордера и их замещения в соответствии с требованиями спецификации:

```
order_book = {-1: {}, 1: {}}
current_orders = {-1: Counter(), 1: Counter()}
message_counter = Counter()
nlevels = 100

for message in messages.itertuples():
    i = message[0]
    if i % 1e5 == 0 and i > 0:
        print('{:,.0f}\t\t{}'.format(i, timedelta(seconds=time() - start)))
        save_orders(order_book, append=True)
        order_book = {-1: {}, 1: {}}
        start = time()
    if np.isnan(message.buy_sell_indicator):
        continue
    message_counter.update(message.type)

    buysell = message.buy_sell_indicator
    price, shares = None, None

    if message.type in ['A', 'F', 'U']:
        price = int(message.price)
        shares = int(message.shares)

        current_orders[buysell].update({price: shares})
        current_orders[buysell], new_order =
            add_orders(current_orders[buysell],
                       buysell, nlevels)
        order_book[buysell][message.timestamp] = new_order
```

```

if message.type in ['E', 'C', 'X', 'D', 'U']:
    if message.type == 'U':
        if not np.isnan(message.shares_replaced):
            price = int(message.price_replaced)
            shares = -int(message.shares_replaced)
    else:
        if not np.isnan(message.price):
            price = int(message.price)
            shares = -int(message.shares)

if price is not None:
    current_orders[buysell].update({price: shares})
    if current_orders[buysell][price] <= 0:
        current_orders[buysell].pop(price)
    current_orders[buysell],
    new_order = add_orders(current_orders[buysell],
                           buysell, nlevels)
    order_book[buysell][message.timestamp] = new_order

```

Число ордеров на различных ценовых уровнях, отмеченных на рис. 2.2 с использованием разной интенсивности для ордеров на покупку и продажу, визуализирует глубину ликвидности в любой момент времени. На графике слева показано, как распределение цен лимитных ордеров было взвешено в сторону ордеров на покупку по более высоким ценам. График справа показывает эволюцию лимитных ордеров и цен в течение всего торгового дня: темная линия отслеживает цены исполненных сделок в часы работы рынка, в то время как красные (в нижней части графика) и синие точки (в верхней части графика) обозначают индивидуальные лимитные ордера на поминутной основе (подробнее см. блокнот).

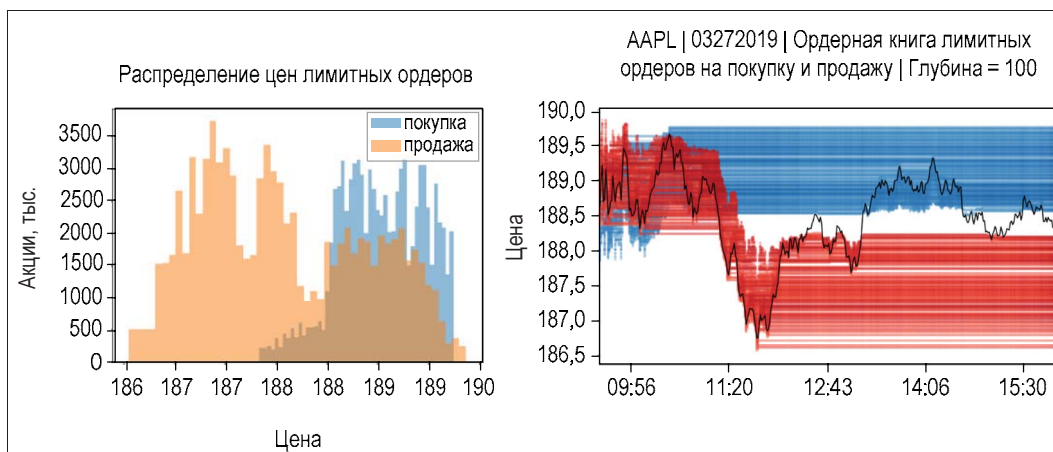


Рис. 2.2. Число ордеров на разных ценовых уровнях по акциям компании Apple (AAPL)

Регуляризация тиковых данных

Торговые данные индексируются наносекундами и очень зашумлены. Например, скачки цен спроса-предложения (bid-ask) вынуждают цену осциллировать между ценами спроса и предложения, когда инициализация сделки чередуется между рыночными ордерами на покупку и продажу. Для улучшения соотношения "шум/сигнал" и улучшения статистических свойств нам необходимо передискретизировать и регуляризовать тиковые данные путем агрегирования торговой активности.

Обычно мы собираем цену открытия (первую цену), цену минимума, цену максимума и цену закрытия (последнюю цену) за агрегируемый период, а также цену, *средневзвешенную по объему* (volume-weighted average price, VWAP), число торгуемых долей (акций) и временную метку, ассоциированную с этими данными.

Дополнительные подробности см. в блокноте `normalize_tick_data.ipynb` в папке этой главы в репозитории GitHub.

Тиковые бары

График сырых тиковых ценовых и объемных данных для акций AAPL выглядит следующим образом:

```
data_path = Path('data')
itch_store = str(data_path / 'itch.h5')      # 'data\\itch.h5'
order_book_store = str(data_path / 'order_book.h5')
stock, date = 'AAPL', '20190327'
title = '{} | {}'.format(stock, pd.to_datetime(date).date())

with pd.HDFStore(itch_store) as store:
    s = store['S'].set_index('event_code').drop_duplicates() # системные события
    s.timestamp = s.timestamp.add(pd.to_datetime(date)).dt.time
    market_open = s.loc['Q', 'timestamp']
    market_close = s.loc['M', 'timestamp']

with pd.HDFStore(order_book_store) as store:
    trades = store['{}/trades'.format(stock)]

trades.price = trades.price.mul(1e-4) # отформатировать цену
trades = trades[trades.cross == 0]    # исключая данные из кроссов открытия/закрытия

trades = trades.between_time(market_open, # только часы работы рынка
                             market_close).drop('cross', axis=1)

tickBars = trades.copy()
tickBars.index = tickBars.index.time
tickBars.price.plot(figsize=(10, 5),
                    title='{} | {}'.format(stock, pd.to_datetime(date).date()),
                    lw=1)
```

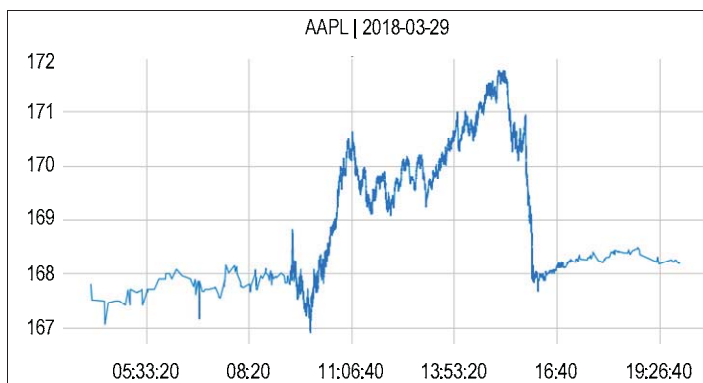


Рис. 2.3. График сырых тиковых данных с ценами и объемами акций AAPL

Приведенный выше фрагмент кода строит график, представленный на рис. 2.3.

Тиковые возвраты далеки от нормального распределения, о чем свидетельствует низкое p -значение функции `scipy.stats.normaltest`:

```
from scipy.stats import normaltest

normaltest(tickBars.price.pct_change().dropna())

NormaltestResult(statistic=11417.148036373566, pvalue=0.0)
```

Временные бары

Временные бары предусматривают агрегирование сделок по периоду:

```
def get_bar_stats(agg_trades):
    vwap = agg_trades.apply(lambda x: np.average(x.price,
                                                  weights=x.shares)).to_frame('vwap')

    ohlc = agg_trades.price.ohlc()
    vol = agg_trades.shares.sum().to_frame('vol')
    txn = agg_trades.shares.size().to_frame('txn')

    return pd.concat([ohlc, vwap, vol, txn], axis=1)

resampled = trades.resample('1Min')
time_bars = get_bar_stats(resampled)
```

Этот результат можно отобразить в виде графика цены-объема:

```
def price_volume(df, price='vwap', vol='vol', supitle=title):
    fig, axes = plt.subplots(nrows=2, sharex=True, figsize=(15,8))
    axes[0].plot(df.index, df[price])
    axes[1].bar(df.index, df[vol], width=1/(len(df.index)), color='r')

    # форматирование
    xfmt = mpl.dates.DateFormatter('%H:%M')
```

```

axes[1].xaxis.set_major_locator(mpl.dates.HourLocator(interval=3))
axes[1].xaxis.set_major_formatter(xfmt)
axes[1].get_xaxis().set_tick_params(which='major', pad=25)
axes[0].set_title('Цена', fontsize=14)
axes[1].set_title('Объем', fontsize=14)
fig.autofmt_xdate()
fig.suptitle(suptitle)
fig.tight_layout()
plt.subplots_adjust(top=0.9)

```

```
price_volume(time_bars)
```

Приведенный выше фрагмент кода строит следующий график (рис. 2.4).

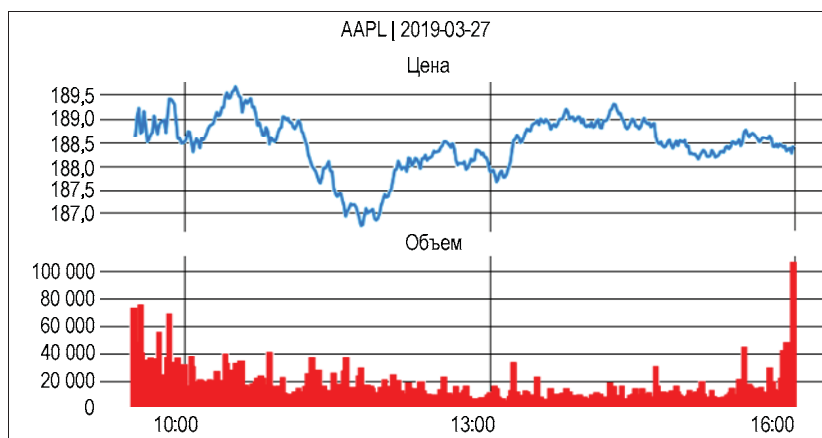


Рис. 2.4. График цены-объема с временными барами

Либо в виде свечного графика с использованием графопостроительной библиотеки Bokeh:

```

resampled = trades.resample('5Min') # 5-минутные интервалы
                                     # для более четкой печати
df = get_bar_stats(resampled)

increase = df.close > df.open
decrease = df.open > df.close
w = 2.5 * 60 * 1000 # 2.5 мин в миллисекундах

WIDGETS = "pan, wheel_zoom, box_zoom, reset, save"

p = figure(x_axis_type='datetime', tools=WIDGETS,
           plot_width=1500, title="Свечной график AAPL")
p.xaxis.major_label_orientation = pi/4
p.grid.grid_line_alpha=0.4

```

```

p.segment(df.index, df.high, df.index, df.low, color="black")
p.vbar(df.index[increase], w, df.open[increase], df.close[increase],
       fill_color="#D5E1DD", line_color="black")
p.vbar(df.index[decrease], w, df.open[decrease], df.close[decrease],
       fill_color="#F2583E", line_color="black")
show(p)

```

В результате получим следующий свечной график (рис. 2.5).

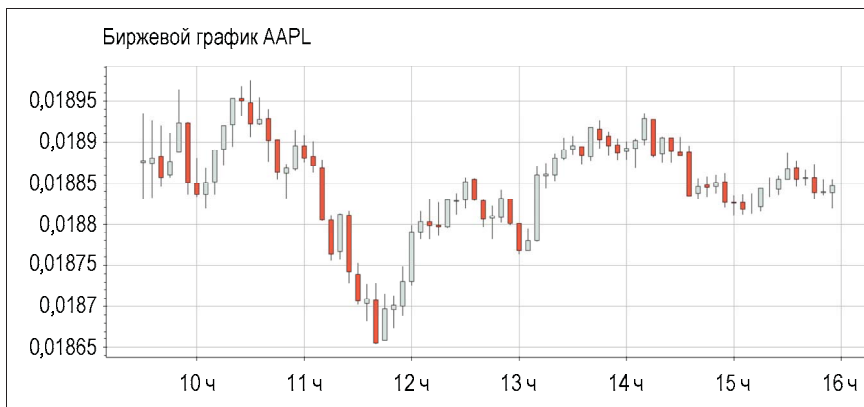


Рис. 2.5. Свечной график движения цен на акции AAPL

Объемные бары

Временные бары сглаживают некоторые шумы, содержащиеся в сырых тиковых данных, но могут не учитывать фрагментацию ордеров. Алгоритмическая торговля, сконцентрированная на исполнении, может ориентироваться на соответствие *средневзвешенной по объему цены* (VWAP) за тот или иной период и будет делить одиночный ордер на многочисленные сделки и размещать ордера в соответствии с историческими регулярностями. Временные бары трактуют один и тот же ордер по-разному, даже если никакой новой информации на рынок не поступало.

Объемные бары предлагают альтернативу за счет агрегирования торговых данных в соответствии с объемом. Это делается следующим образом:

```

trades_per_min = trades.shares.sum()/(60*7.5) # минимум за торговый день
trades['cumul_vol'] = trades.shares.cumsum()
df = trades.reset_index()
by_vol = df.groupby(df.cumul_vol.div(trades_per_min).round().astype(int))
volBars = pd.concat([by_vol.timestamp.last().to_frame('timestamp'),
                    get_bar_stats(by_vol)], axis=1)
price_volume(volBars.set_index('timestamp'))

```

Приведенный выше фрагмент кода строит график, представленный на рис. 2.6.

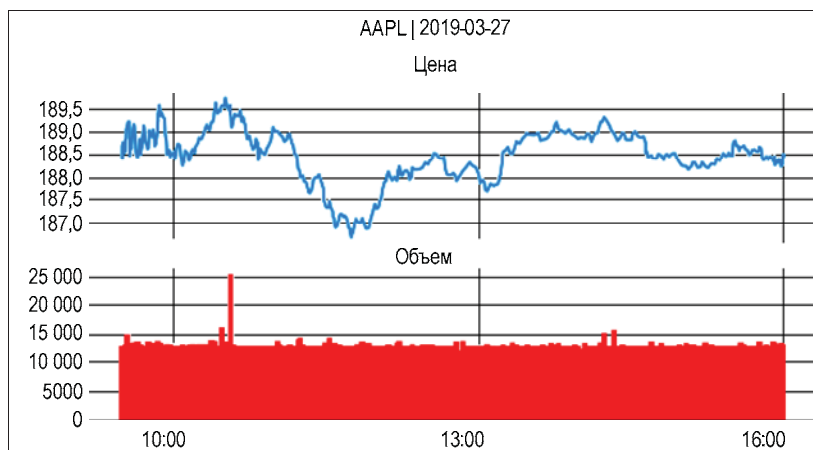


Рис. 2.6. График цены-объема с объемными барами

Долларовые бары

При значительном изменении цен на активы или после дробления акций стоимость определенного количества акций изменяется. Объемные бары отражают этот факт неправильно и могут затруднить сравнение торгового поведения за разные периоды времени, которые отражают такие изменения. В этих случаях метод объемного бара должен быть скорректирован за счет произведения акций и цены, в результате чего получаются долларовые бары.

Доступ к рыночным данным через API

Существует несколько вариантов доступа к рыночным данным через API с помощью языка Python. Сначала мы представим несколько источников, встроенных в библиотеку `pandas`. Затем мы кратко представим торговую платформу Quantopian, поставщика данных Quandl и библиотеку бэктестирования, которую будем использовать позже в этой книге, а также перечислим несколько дополнительных вариантов доступа к различным типам рыночных данных. В папке `data_providers` репозитория GitHub содержится несколько блокнотов, иллюстрирующих использование этих вариантов.

Дистанционный доступ к данным с помощью библиотеки `pandas`

Библиотека `pandas` обеспечивает доступ к данным, размещаемым на веб-сайтах, посредством функции `read_html` и доступ к конечным точкам API различных поставщиков данных посредством родственной библиотеки `pandas-datareader`.

Чтение `html`-таблиц

Скачивание содержимого одной или нескольких `html`-таблиц, например для компонентов индекса S&P 500 из Википедии, работает следующим образом:

```
sp_url = 'https://en.wikipedia.org/wiki/List_of_S%26P_500_companies'
sp = pd.read_html(sp_url, header=0)[0] # возвращает список
                                         # для каждой таблицы

sp.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 505 entries, 0 to 504
Data columns (total 9 columns):
Symbol                505 non-null object
Security              505 non-null object
SEC filings           505 non-null object
GICS Sector           505 non-null object
GICS Sub Industry     505 non-null object
Headquarters Location 505 non-null object
Date first added      402 non-null object
CIK                   505 non-null int64
Founded               172 non-null object
dtypes: int64(1), object(8)
memory usage: 35.6+ KB
```

Библиотека pandas-datareader для рыночных данных

Библиотека pandas используется с целью обеспечения прямого доступа к API поставщиков данных, однако этот функционал переместился в соответствующую библиотеку pandas-datareader. Стабильность API варьируется в зависимости от политик поставщика, и по состоянию на июнь 2018 г. в версии 0.7 этой библиотеки доступны источники¹¹, перечисленные в табл. 2.4.

Таблица 2.4. Перечень поставщиков данных, поддерживаемых библиотекой pandas-datareader

| Источник | Область действия | Комментарий | Описание |
|------------------------------|---|----------------------------------|---|
| Yahoo! Finance | Цены на конец дня (EOD), дивиденды, данные по дроблению акций и валютным парам | Нестабилен | Поставщик финансовой информации, https://finance.yahoo.com/ |
| Tiingo | Цены на долевые ценные бумаги, взаимные фонды и биржевые фонды на конец дня (EOD) | Требуется бесплатная регистрация | Платформа финансового анализа, https://www.tiingo.com |
| The Investors Exchange (IEX) | Исторические цены на акции, данные ордерных книг | Ограничен пятью годами | Биржа, специализирующаяся на защите инвесторов и эмитентов, https://www.iextrading.com/ |

¹¹ Полный перечень из 15 источников данных смотрите на странице официальной документации по адресу https://pydata.github.io/pandas-datareader/devel/remote_data.html. — Прим. перев.

Таблица 2.4 (окончание)

| Источник | Область действия | Комментарий | Описание |
|---------------|--|---|--|
| Robinhood | Цены на долевые ценные бумаги на конец дня (EOD) | Ограничен одним годом | Финансовые услуги, https://robinhood.com/ |
| Quandl | Рыночная площадка с широким диапазоном цен на финансовые активы | Для премиальных данных требуется подписка | Финансово-экономические и альтернативные данные, https://www.quandl.com/ |
| NASDAQ | Свежие тикерные символы, торговавшиеся на бирже NASDAQ, плюс некоторая дополнительная информация | | Веб-сайт биржи NASDAQ, https://www.nasdaq.com/ |
| Stooq | Данные нескольких индексов фондового рынка | | Польский веб-сайт, https://stooq.com/ |
| MOEX | Данные Московской фондовой биржи | | https://www.moex.com/ |
| Alpha Vantage | Цены на конец дня (EOD) на акции и валютные пары | | Сообщество исследователей, инженеров и менеджеров-профессионалов, https://www.alphavantage.co/ |
| Fama/French | Факторные возвраты и исследовательские портфели из библиотеки данных FF Data Library | | Веб-страница библиотеки данных Кеннета Френча, http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html |

Доступ и извлечение данных следует аналогичному API для всех источников, как это проиллюстрировано для Yahoo! Finance:

```
import pandas_datareader.data as web
from datetime import datetime

start = '2016' # принимает строковые значения
end = datetime(2019, 5, 24) # или объекты datetime

yahoo= web.DataReader('FB', 'yahoo', start=start, end=end)
yahoo.info()

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 845 entries, 2015-12-31 to 2019-05-10
Data columns (total 6 columns):
High      845 non-null float64
Low       845 non-null float64
Open      845 non-null float64
Close     845 non-null float64
```

```
Volume      845 non-null int64
Adj Close    845 non-null float64
dtypes: float64(5), int64(1)
memory usage: 46.2 KB
```

Биржа инвесторов Investors Exchange

Биржа IEX (Investors Exchange)¹² представляет собой альтернативную торговую площадку, запущенную в ответ на разногласия по поводу высокочастотной торговли (HFT), описанные в спорной книге Майкла Льюиса (Michael Lewis) "Flash Boys: Высокочастотная революция на Уолл-стрит" ("Flash Boys: A Wall Street Revolt"). Она призвана замедлять скорости торговли с целью создания более ровного игрового поля и ускоренно набирает обороты с момента своего запуска в 2016 г., довольствуясь на июнь 2018 г. небольшой долей рынка около 2,5%.

В дополнение к историческим данным о цене и объеме на конец дня биржа IEX обеспечивает реально-временную глубину котировок книги, которые предлагают агрегированный размер ордеров по цене и стороне сделки. Эта служба также включает в себя информацию о последней цене и размере сделки. Ниже приведен стандартный вызов:

```
start = datetime(2017, 2, 9)
iex = web.DataReader('FB', 'iex', start)
iex.info()

<class 'pandas.core.frame.DataFrame'>
Index: 566 entries, 2017-02-09 to 2019-05-10
Data columns (total 5 columns):
open      566 non-null float64
high      566 non-null float64
low       566 non-null float64
close     566 non-null float64
volume    566 non-null int64
dtypes: float64(4), int64(1)
memory usage: 26.5+ KB
```

Дополнительные примеры см. в блокноте `datareader.ipynb`.

Платформа Quantopian

Quantopian — это инвестиционная фирма, которая предлагает исследовательскую платформу для прямого вовлечения людей через Интернет (краудсорсинга) с целью отбора алгоритмов торговли. После бесплатной регистрации она позволяет членам

¹² Investors Exchange (IEX, биржа инвесторов) — фондовая биржа, базирующаяся в США, основана в 2012 г. и запущена в качестве национальной биржи ценных бумаг в сентябре 2016 г. 24 октября 2017 г. IEX получила нормативное одобрение SEC на включение компаний в список. IEX зарегистрировала свою первую публичную компанию Interactive Brokers 5 октября 2018 г. — *Прим. перев.*

исследовать торговые идеи, используя широкий спектр источников данных. Она также предлагает среду для бэктестирования алгоритма на исторических данных, а также для его форвардного тестирования вне выборки^{13, 14} с живыми данными. В качестве награды она размещает инвестиционные средства среди высокорезультативных алгоритмов, авторы которых (на момент написания этих строк) имеют право на 10%-ю долю прибыли.

Исследовательская платформа Quantopian состоит из среды на основе блокнота Jupyter, служащей для изучения и анализа результативности альфа-факторов. Кроме того, имеется *интерактивная среда разработки* (IDE) для программирования алгоритмических стратегий и бэктестирования результатов с использованием исторических данных, начиная с 2002 г. с частотой минутного бара.

Пользователи также могут симулировать работу алгоритмов, т. е. торговать виртуально, занимаясь так называемой бумажной торговлей с живыми данными. Платформа Quantopian предоставляет различные наборы рыночных данных, включая ценовые и объемные данные по долевым ценным бумагам и фьючерсам США, а также корпоративные фундаментальные показатели долевого ценного бумага компаний США, и интегрирует многочисленные альтернативные наборы данных.

Мы рассмотрим платформу Quantopian подробнее в *главе 4* и будем опираться на ее функционал на протяжении всей книги, поэтому без колебаний открывайте там аккаунт (подробности см. в репозитории GitHub).

Библиотека zipline

Библиотека zipline содержит функционал алгоритмической торговли, которая приводит в действие платформу бэктестирования и живой торговли Quantopian. Она также доступна в автономном режиме для разработки стратегии с использованием лимитированного числа бесплатных комплектов данных, которые могут поглощаться и использоваться для тестирования результативности торговых идей перед переносом результата в онлайн-платформу Quantopian для бумажной и живой торговли.

¹³ Форвардное тестирование (forward testing) — это симулирование реальных рыночных данных виртуально, только на бумаге, т. е. хотя вы и движетесь по рынкам вживую, на самом деле вы не вкладываете реальные деньги и занимаетесь виртуальной торговлей с целью лучше понять движения рынков. См. https://en.wikipedia.org/wiki/Walk_forward_optimization. — *Прим. перев.*

¹⁴ Внутри выборки и вне выборки (in-sample и out-of-sample) — тестирование обычно проводится путем разбиения набора данных на внутривыборочный (in-sample) период, используемый для первоначального оценивания параметров и отбора модели, и вневыборочный (out-of-sample) период, используемый для оценивания результативности предсказания. В частности, если имеются данные, скажем, за 3 года, необходимые для расчета волатильности, то модель, используемая в течение этого периода, будет "внутри выборки". Но если использовать исторические данные для предсказания на перспективу, то оценивание будет выполняться за период времени, для которых нет данных (вне выборки). Таким образом, обычно понятие "вне выборки" относится к предсказанию там, где у нас нет данных. Технически, даже использование модели для оценивания сегодняшней волатильности на основе исторической выборки является предсказанием вне выборки, потому что у нас нет сиюминутной волатильности. См. https://ec.europa.eu/eurostat/statistics-explained/index.php/Glossary:In-sample_vs._out-of-sample_forecasts. — *Прим. перев.*

Следующий ниже фрагмент кода иллюстрирует, как библиотека `zipline` позволяет получать доступ к ежедневным данным акций ряда компаний. Сценарии `zipline` можно запускать в блокноте `Jupyter`, используя магическую функцию с тем же именем.

Сначала необходимо инициализировать контекст нужными символами ценных бумаг. Мы также будем использовать переменную `context`. Затем `zipline` вызывает функцию `handle_data`, где мы используем метод `data.history()` для заглядывания назад на один период и добавления данных за последний день в CSV-файл:

```
%load_ext zipline
%%zipline --start 2010-1-1 --end 2018-1-1 --data-frequency daily
from zipline.api import order_target, record, symbol

def initialize(context):
    context.i = 0
    context.assets = [symbol('FB'), symbol('GOOG'), symbol('AMZN')]

def handle_data(context, data):
    df = data.history(context.assets, fields=['price', 'volume'],
                      bar_count=1, frequency="1d")
    df = df.to_frame().reset_index()

    if context.i == 0:
        df.columns = ['date', 'asset', 'price', 'volume']
        df.to_csv('stock_data.csv', index=False)
    else:
        df.to_csv('stock_data.csv', index=False, mode='a', header=None)
        context.i += 1

df = pd.read_csv('stock_data.csv')
df.date = pd.to_datetime(df.date)
df.set_index('date').groupby('asset').price.plot(lw=2, legend=True, figsize=(14, 6));
```

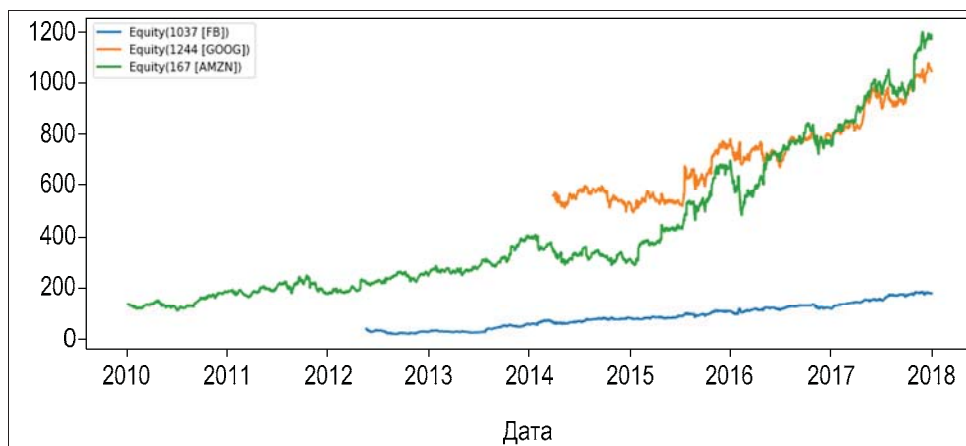


Рис. 2.7. Построение ценового графика нескольких финансовых активов с помощью библиотеки `zipline`

Приведенный выше фрагмент кода строит график, представленный на рис. 2.7.

В следующих главах мы подробнее рассмотрим возможности библиотеки `zipline` и, в особенности, онлайн-платформы `Quantopian`.

Платформа Quandl

Платформа `Quandl` предоставляет широкий спектр источников данных, как бесплатных, так и по подписке, используя `API Python`. Зарегистрируйтесь и получите бесплатный ключ `API`, который позволит совершать более 50 вызовов `API` в день. Данные `Quandl` охватывают несколько классов финансовых активов помимо доле-вых ценных бумаг и содержат валютные пары рынка `Forex`, ценные бумаги с фиксированной доходностью, индексы, фьючерсы и опционы, а также товары.

Использовать ее интерфейс `API` очень просто. Он документирован и гибок и имеет многочисленные методы, выходящие за пределы скачивания данных одного цено-вого ряда, например, включая массовые скачивания или поиск метаданных. Следующий ниже вызов `API` возвращает цены на нефть с 1986 г. согласно котировкам Министерства энергетики США:

```
import quandl
```

```
oil = quandl.get('EIA/PET_RWTC_D').squeeze()
```

```
oil.plot(lw=2, title='Цена сырой нефти WTI')
```

Приведенный выше фрагмент кода построит график, представленный на рис. 2.8.



Рис. 2.8. Ценовой график финансового актива, полученный с помощью библиотеки `quandl`

Другие поставщики рыночных данных

Целый ряд поставщиков предлагает рыночные данные для различных классов финансовых активов. Приведем несколько примеров в соответствующих категориях.

- ◆ Биржи, которые получают все бóльшую долю своей выручки от все более широкого спектра служб данных, как правило, используя подписку.
- ◆ Агентства Bloomberg и Thomson Reuters, которые уже давно являются ведущими агрегаторами данных с совокупной долей более 55% на рынке финансовых данных с его суммарным объемом 28,5 млрд долларов. Помимо них, растут более мелкие конкуренты, такие как агрегатор FactSet, или появляются платформы, такие как money.net и Quandl, а также платформы Trading Economics (<https://tradingeconomics.com>) или Barchart (<https://www.barchart.com/>).
- ◆ Изобилие специализированных поставщиков данных. Одним из примеров является поставщик LOBSTER (Limit Order Book SysTem, <https://lobsterdata.com/>), который агрегирует данные ордерной книги биржи NASDAQ в режиме реального времени.
- ◆ Бесплатные поставщики данных, в том числе поставщик Alpha Vantage (<https://www.alphavantage.co/>), который предлагает API на языке Python для реально-временных данных рынка долевых ценных бумаг, рынка Forex и рынка криптовалют, а также технические индикаторы.
- ◆ Краудсорсинговые инвестиционные фирмы, питаемые за счет прямого вовлечения людей через Интернет, которые предоставляют исследовательские платформы с доступом к данным, помимо Quantopian включают фирму Alpha Trading Labs (<https://www.alphatradinglabs.com/>), стартовавшую в марте 2018 г., которая обеспечивает инфраструктуру высокочастотной торговли и данные.

Как работать с фундаментальными данными

Фундаментальные данные относятся к экономическим движущим силам, определяющим стоимость ценных бумаг. Характер данных зависит от класса финансовых активов:

- ◆ применительно к долевым ценным бумагам и корпоративным кредитам они охватывают корпоративные финансовые показатели, а также отраслевые и общеэкономические данные;
- ◆ применительно к государственным облигациям они охватывают международные макроданные и иностранную валюту;
- ◆ применительно к товарам они охватывают факторы, определяющие спрос и предложение, специфичные для конкретного финансового актива, такие как погодные данные для сельскохозяйственных культур.

Мы сосредоточимся на фундаментальных финансовых показателях долевых ценных бумаг компаний США¹⁵, где легче получить доступ к данным. По всему миру

¹⁵ Фундаментальные финансовые показатели акционерного капитала (equity fundamentals) включают в себя денежный поток, доходность активов, соотношение собственного/заемного капитала (консервативное лeverажное соотношение), историю удержания прибыли для финансирования будущего роста, обоснованность управления капиталом для максимизации заработков и финансовых возвратов акционеров. См. <https://www.investopedia.com/articles/fundamental/03/022603.asp>. — Прим. перев.

существует более 13 000 публичных компаний, которые генерируют 2 млн страниц годовых отчетов и более 30 000 часов телеконференций о корпоративных заработках¹⁶. В алгоритмической торговле фундаментальные данные и признаки, вырабатываемые из этих данных, могут использоваться для непосредственного получения торговых сигналов, например в качестве стоимостных индикаторов, и являются важным входом в предсказательные модели, включая автоматически обучающиеся модели.

Данные финансовой отчетности

Комиссия по ценным бумагам и биржам (SEC) требует от эмитентов США, т. е. зарегистрированных на бирже компаний и ценных бумаг, включая взаимные фонды, подавать три квартальных финансовых отчета (по форме 10-Q) и один годовой отчет (по форме 10-K) в дополнение к другим всевозможным требованиям регламента к предъявляемой официальной отчетности.

С начала 1990-х годов комиссия SEC сделала эти документы финансовой отчетности доступными через свою систему *электронного сбора, анализа и извлечения данных* (Electronic Data Gathering, Analysis, and Retrieval, EDGAR). Они являются первичными источниками данных для фундаментального анализа долевых и других ценных бумаг, таких как корпоративные кредиты, где стоимость зависит от перспектив бизнеса и финансового состояния эмитента.

Автоматизированная обработка — XBRL

Автоматизированный анализ регламентных документов официальной отчетности стал намного проще с тех пор, как комиссия SEC ввела расширяемый язык деловой отчетности XBRL (Extensible Business Reporting Language) — свободный, открытый и глобальный стандарт для электронного представления и обмена производственными отчетами. Язык XBRL основан на XML; он опирается на *таксономии*, которые определяют смысл элементов отчета и связывают их с тегами, которые подчеркивают соответствующую информацию в электронной версии отчета. Одна из таких таксономий представляет *общепринятые принципы бухгалтерского учета США* (Generally Accepted Accounting Principles, GAAP).

Комиссия SEC ввела добровольную подачу документов официальной отчетности в формате XBRL в 2005 г. в ответ на бухгалтерские скандалы, после чего с 2009 г. начала требовать этот формат со всех заявителей отчетности и продолжает расши-

¹⁶ Телеконференция о корпоративных заработках (earnings call) — это встреча или веб-трансляция, в которой публичная компания обсуждает финансовые результаты отчетного периода. Данный термин происходит от понятия "заработки в расчете на долю в акционерном капитале" (earnings per share, EPS), т. е. число в нижней строке отчета о прибылях и убытках, разделенное на число долей в обращении. Базирующийся в США Национальный институт по связям с инвесторами (NIRI) сообщает, что 92% компаний, представленных своими членами, проводят такие телеконференции и что практически все они транслируются через Интернет. См. https://en.wikipedia.org/wiki/Earnings_call. — Прим. перев.

рять обязательный охват других регламентных документов официальной отчетности. Комиссия SEC имеет свой веб-сайт, на котором приводится текущая таксономия, формирующая содержимое разных документов официальной отчетности и которые могут использоваться для извлечения специфических элементов.

Следующие далее наборы данных содержат информацию, извлеченную из вложений EX-101, предъявляемых в Комиссию в плоском формате данных, тем самым помогая пользователям в потреблении этих данных с целью анализа. Данные отражают информацию, отобранную из финансовых отчетов, помеченных XBRL-тегами. В настоящее время она включает числовые данные из квартальной и годовой финансовой отчетности, а также некоторые дополнительные поля (например, поле стандартной промышленной классификации (Standard Industrial Classification, SIC)).

Существует несколько способов отслеживания и доступа к фундаментальным финансовым данным, подаваемым в комиссию SEC:

- ◆ электронные каналы принятых документов официальной отчетности в рамках *публичной службы распространения информации* (Public Dissemination Service, PDS) системы EDGAR доступны за отдельную плату;
- ◆ комиссия SEC обновляет *RSS-каналы* каждые 10 минут, в которых перечисляются структурированные материалы раскрываемых сведений;
- ◆ существуют публичные индексные файлы для извлечения через FTP всех поданных документов официальной отчетности с целью автоматической обработки;
- ◆ наборы данных финансовых отчетов (и примечаний) содержат разобранные XBRL-данные из всех финансовых отчетов и сопроводительных примечаний.

Комиссия SEC также публикует журнальные файлы через **SEC.gov**, содержащие поисковый интернет-трафик относительно поданных в систему EDGAR документов официальной отчетности. Правда, это делается с шестимесячной задержкой.

Построение временного ряда с фундаментальными данными

Область действия наборов данных финансовых отчетов и примечаний ограничивается числовыми данными, извлеченными из первичных финансовых отчетов (балансового отчета, отчетов о прибылях и убытках, о потоке денежных средств, об изменении в собственном капитале и о совокупном доходе) и сносков к этим отчетам. Данные имеются, начиная с 2009 г.

Извлечение набора данных финансовых отчетов и примечаний

Следующий ниже фрагмент кода скачивает и извлекает все исторические документы официальной отчетности, содержащиеся в наборах данных *финансовых отчетов и примечаний* (Financial Statement and Notes, FSN) для заданного диапазона кварталов (дальнейшие подробности см. в блокноте `edgar_xbrl.ipynb`):

```

SEC_URL = 'https://www.sec.gov/files/dera/data/financial-statement-and-notes-data-sets/'

first_year, this_year, this_quarter = 2014, 2018, 3
past_years = range(2014, this_year)
filing_periods = [(y, q) for y in past_years for q in range(1, 5)]
filing_periods.extend([(this_year, q) for q in range(1, this_quarter + 1)])

for i, (yr, qtr) in enumerate(filing_periods, 1):
    filing = f'{yr}q{qtr}_notes.zip'
    path = data_path / f'{yr}_{qtr}' / 'source'
    response = requests.get(SEC_URL + filing).content

    with ZipFile(BytesIO(response)) as zip_file:
        for file in zip_file.namelist():
            local_file = path / file

            with local_file.open('wb') as output:
                for line in zip_file.open(file).readlines():
                    output.write(line)

```

Эти данные являются довольно крупными, и для обеспечения более быстрого доступа, чем это позволяют делать исходные текстовые файлы, лучше всего конвертировать текстовые файлы в двоичный, столбчатый формат `parquet` (см. разд. "Эффективное хранение данных с помощью библиотеки `pandas`" далее в этой главе о сопоставлении производительности различных вариантов хранения данных, совместимых с кадрами данных `DataFrame` библиотеки `pandas`):

```

for f in data_path.glob('**/*.tsv'):
    file_name = f.stem + '.parquet'
    path = Path(f.parents[1]) / 'parquet'
    df = pd.read_csv(f, sep='\t', encoding='latin1', low_memory=False)
    df.to_parquet(path / file_name)

```

Данные финансовых отчетов и примечаний FSN организованы по каждому кварталу в восемь наборов файлов, которые содержат информацию о подаче отчетности, числах, таксономических тегах, презентации и многое другое. Каждый набор данных состоит из строк и полей и предоставляется в виде текстового файла с разделением полей данных символом табуляции (табл. 2.5).

Таблица 2.5. Поля текстового файла FSN

| Формат файла | Набор данных | Описание |
|--------------|-------------------|---|
| SUB | Подача отчетности | Идентифицирует каждую подачу отчетности в формате XBRL по компании, форме, дате и т. д. |
| TAG | Тег | Определяет и объясняет каждый таксономический тег |

Таблица 2.5 (окончание)

| Формат файла | Набор данных | Описание |
|--------------|---------------|---|
| DIM | Размерность | Добавляет детали к числовым и текстовым данным |
| NUM | Числовой | Содержит одну строку для каждой несовпадающей точки данных в документе официальной отчетности |
| TXT | Простой текст | Содержит все нечисловые поля XBRL |
| REN | Отрисовка | Содержит информацию для визуализации на сайте SEC |
| PRE | Презентация | Содержит детализацию о представлении тегов и чисел в первичных отчетах |
| CAL | Калькуляция | Показывает арифметические взаимосвязи между тегами |

Извлечение всей ежеквартальной официальной отчетности компании Apple

Набор данных с отчетностью содержит уникальные идентификаторы, необходимые для извлечения документов официальной отчетности: *центральный индексный ключ* (Central Index Key, CIK) и учетный номер (adsh). Ниже приведена информация о квартальной официальной отчетности компании Apple по форме 10-Q за первый квартал 2018 г. (2018Q1):

```
apple = sub[sub.name == 'APPLE INC'].T.dropna().squeeze()
key_cols = ['name', 'adsh', 'cik', 'name', 'sic', 'countryba',
            'strba', 'cityba', 'zipba', 'basl', 'form', 'period',
            'fy', 'fp', 'filed']
apple.loc[key_cols]
```

| | |
|-----------|----------------------|
| name | APPLE INC |
| adsh | 0000320193-18-000070 |
| cik | 320193 |
| name | APPLE INC |
| sic | 3571 |
| countryba | US |
| strba | CA |
| cityba | CUPERTINO |
| zipba | 95014 |
| basl | ONE APPLE PARK WAY |
| form | 10-Q |
| period | 20180331 |
| fy | 2018 |
| fp | Q2 |
| filed | 20180502 |

Используя центральный индексный ключ, мы можем выявлять всю историческую квартальную официальную отчетность по компании Apple и объединять эту информацию для получения 26 форм 10-Q и девяти годовых форм 10-K:

```
aapl_subs = pd.DataFrame()
for sub in data_path.glob('**/sub.parquet'):
    sub = pd.read_parquet(sub)
    aapl_sub = sub[(sub.cik.astype(int) == apple.cik) &
                  (sub.form.isin(['10-Q', '10-K']))]
    aapl_subs = pd.concat([aapl_subs, aapl_sub])
```

```
aapl_subs.form.value_counts()
```

```
10-Q 15
```

```
10-K 4
```

С помощью учетного номера adsh для каждого документа официальной отчетности теперь мы можем опираться на таксономии при отборе соответствующих XBRL-тегов (перечисленных в TAG-файле) из файлов NUM и TXT для получения вызывающих интерес числовых или текстовых/сносковых данных.

Прежде всего, давайте извлечем все имеющиеся числовые данные из 19 официальных отчетов компании Apple:

```
aapl_nums = pd.DataFrame()
for num in data_path.glob('**/num.parquet'):
    num = pd.read_parquet(num)
    aapl_num = num[num.adsh.isin(aapl_subs.adsh)]
    aapl_nums = pd.concat([aapl_nums, aapl_num])

aapl_nums.ddate = pd.to_datetime(aapl_nums.ddate, format='%Y%m%d')

aapl_nums.shape

(28281, 16)
```

Построение временного ряда "цена/заработки"

Девятилетняя история подачи документов официальной отчетности предоставляет нам в общей сложности более 28 тыс. числовых значений. Мы можем выбрать полезное поле, такое как *заработки на долю собственности в акционерном капитале разбавленные* (Diluted Earnings per Share, Diluted EPS), которое можно объединить с рыночными данными для расчета популярной метрики оценивания рыночной стоимости — *соотношения цены к заработкам* (P/E).

Однако нам все-таки нужно учесть дробление акций в соотношении 7:1, которое компания Apple предприняла 4 июня 2014 г., и скорректированные заработки

в расчете на акцию перед дроблением, чтобы сделать заработки сопоставимыми, как показано в следующем фрагменте кода:

```
field = 'EarningsPerShareDiluted'
stock_split = 7
split_date = pd.to_datetime('20140604')

# Отфильтровать по тегу; оставить только те значения,
# которые измеряются I кварталом
eps = aapl_nums[(aapl_nums.tag == 'EarningsPerShareDiluted') & \
                (aapl_nums.qtrs == 1)].drop('tag', axis=1)

# Оставить только наиболее недавнюю точку данных
# из каждого документа официальной отчетности
eps = eps.groupby('adsh').apply(lambda x: x.nlargest(n=1, columns=['ddate']))

# Скорректировать заработки вниз перед дроблением акций
eps.loc[eps.ddate < split_date, 'value'] = eps.loc[eps.ddate < split_date,
'value'].div(7)
eps = eps[['ddate', 'value']].set_index('ddate').squeeze()

# Создать заработки на долю (eps) за предшествующие
# 12 месяцев из квартальных данных
eps = eps.rolling(4, min_periods=4).sum().dropna()
```

Для получения данных по ценам акций Apple с 2009 г. можно использовать поставщика Quandl:

```
import pandas_datareader.data as web
symbol = 'AAPL.US'
aapl_stock = web.DataReader(symbol, 'quandl', start=eps.index.min())

# Обеспечить, чтобы даты были выровнены
# с данными заработка на акцию (eps)
aapl_stock = aapl_stock.resample('D').last()
```

Теперь у нас есть данные для вычисления соотношения цены к заработкам (P/E) с учетом предшествующих 12 месяцев за весь период:

```
pe = aapl_stock.AdjClose.to_frame('price').join(eps.to_frame('eps'))
pe = pe.fillna(method='ffill').dropna()
pe['P/E Ratio'] = pe.price.div(pe.eps)
axes = pe.plot(subplots=True, figsize=(16,8), legend=False, lw=2);
```

Приведенный фрагмент кода строит график, представленный на рис. 2.9.

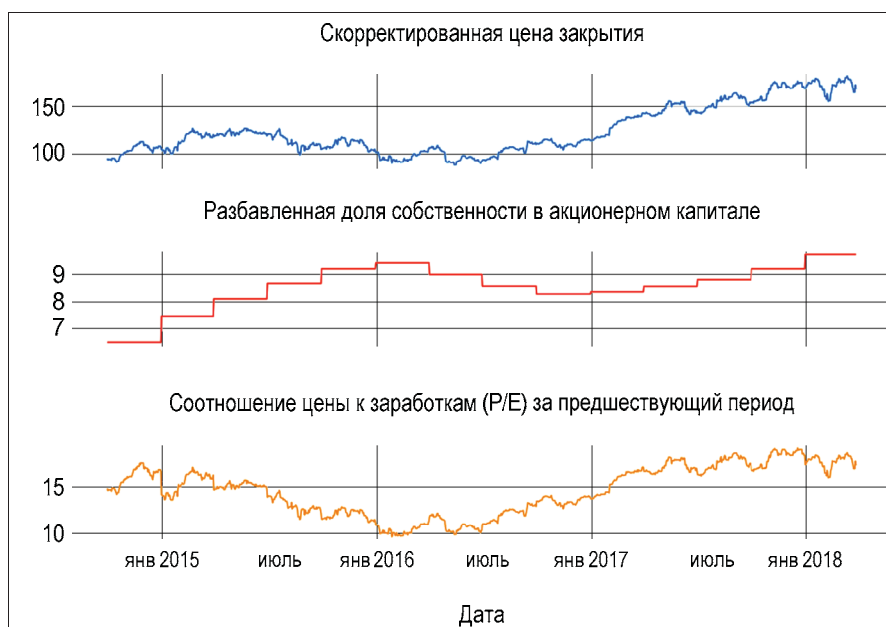


Рис. 2.9. График соотношения цены к заработкам (P/E) с учетом предшествующих 12 месяцев за весь период по акциям Apple

Другие источники фундаментальных данных

Существует ряд других источников фундаментальных данных. Многие из них доступны с помощью представленного ранее модуля `pandas-datareader`. Дополнительные данные можно получить непосредственно от некоторых организаций, таких как МВФ, Всемирный банк или крупные национальные статистические учреждения во всем мире (справочные материалы см. в репозитории GitHub).

Макроэкономические и отраслевые данные с помощью `pandas-datareader`

Библиотека `pandas-datareader` обеспечивает доступ в соответствии с условиями, представленными в конце предыдущего раздела о рыночных данных. Она охватывает API для многочисленных источников глобальных фундаментальных макроэкономических и отраслевых данных, включая следующие:

- ◆ *библиотека данных Кеннета Френча* — рыночные данные о портфелях, охватывающие размер, стоимость и импульсные факторы, дезагрегированные по отраслям (https://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html);
- ◆ *сент-луисские данные ФРС* (St. Louis FED, FRED) — данные Федеральной резервной системы по экономике и финансовым рынкам США (<https://fred.stlouisfed.org/categories>);

- ♦ *Всемирный банк* — глобальная база данных о долговременном низкочастотном экономическом и социальном развитии и демографии (<https://data.worldbank.org/>);
- ♦ *ОЭСР* — аналогичные данные для стран ОЭСР (<https://data.oecd.org/>);
- ♦ *Enigma* — различные наборы данных, включая альтернативные источники (<https://www.enigma.com/public-data>);
- ♦ *Евростат* — ориентированные на ЕС экономические и социально-демографические данные (<https://ec.europa.eu/eurostat/data/database>).

Эффективное хранение данных с помощью библиотеки pandas

В этой книге мы будем использовать самые разные наборы данных, и поэтому стоит сравнить основные форматы на эффективность и производительность. В частности, мы сравним следующие форматы:

- ♦ CSV — стандартный формат плоского текстового файла с разделением полей данных запятыми;
- ♦ HDF5 — иерархический формат данных, первоначально разработанный в Национальном центре суперкомпьютерных вычислений, представляет собой быстрый и масштабируемый формат хранения числовых данных, доступный в библиотеке pandas с использованием библиотеки PyTables;
- ♦ parquet — двоичный, столбчатый формат хранения в составе экосистемы Apache Hadoop, который обеспечивает эффективное сжатие и кодирование данных; был разработан в Cloudera и Twitter. Он доступен для библиотеки pandas через библиотеку pyarrow, возглавляемую Уэсом Маккинни, первоначальным автором библиотеки pandas.

В блокноте `storage_benchmark.ipynb` приводится сравнение производительности указанных библиотек с использованием тестового кадра данных, который может быть сконфигурирован так, чтобы содержать числовые или текстовые данные. Для библиотеки HDF5 мы тестируем как фиксированный, так и табличный формат. Табличный формат позволяет выполнять запросы и может пополняться.

На диаграммах рис. 2.10 показана производительность чтения и записи 100 тыс. строк с 1000 столбцами случайных вещественных чисел и 1000 столбцами случайного 10-символьного строкового значения либо просто со столбцами из 2000 вещественных чисел.

- ♦ В случае чисто числовых данных формат HDF5 работает лучше всего, при этом табличный формат HDF5 также имеет совместно с CSV наименьший объем памяти на уровне 1,6 Гбайт. Фиксированный формат HDF5 использует в 2 раза больше места, при этом формат parquet использует 2 Гбайт.
- ♦ В случае смеси числовых и текстовых данных формат parquet работает значительно быстрее, при этом формат HDF5 использует свое преимущество при чте-

нии по сравнению с форматом CSV (который имеет очень низкую производительность записи в обоих случаях).

В упомянутом выше блокноте показано, как конфигурировать, тестировать и хронометрировать работу с помощью волшебной команды `%%timeit`, и в то же время показано применение соответствующих команд библиотеки `pandas`, необходимых для использования этих форматов хранения.

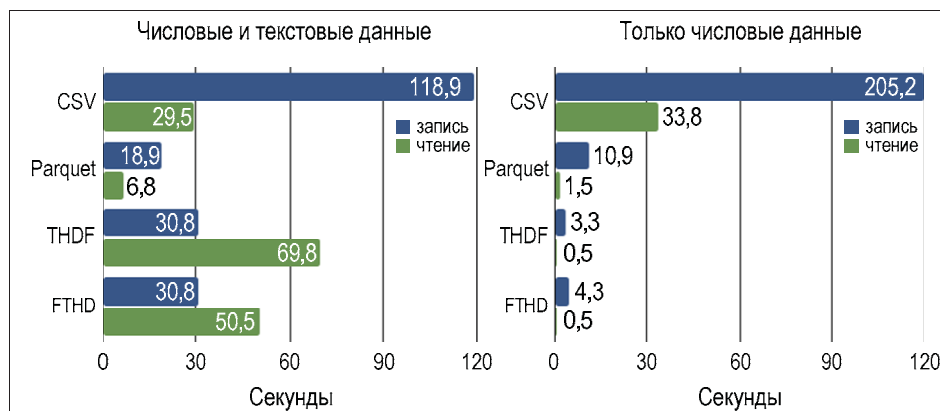


Рис. 2.10. Сравнение производительности чтения и записи разных форматов хранения данных: THDF — табличный HDF; FHDF — фиксированный HDF

Резюме

В этой главе были представлены источники рыночных и фундаментальных данных, которые составляют стержень большинства стратегий торговли на финансовых рынках. Вы узнали о многочисленных способах доступа к этим данным и о том, как выполнять предобработку сырой информации, что позволяет приступить к извлечению торговых сигналов с использованием методов машинного обучения, которые мы вскоре представим.

Прежде чем мы перейдем к разработке и оцениванию стратегий торговли и использованию автоматически обучающихся моделей, необходимо рассмотреть альтернативные наборы данных, которые появились в последние годы и были значительной движущей силой роста популярности МО для алгоритмической торговли на финансовых рынках.

3

Альтернативные данные для финансов

Подстегиваемые взрывным ростом Глобальной паутины и мобильных сетей цифровые данные продолжают расти экспоненциально на фоне прогресса в технологии обработки, хранения и анализа новых источников данных. Экспоненциальный рост доступности и способности управлять более разнообразными цифровыми данными, в свою очередь, был решающей силой, стоящей за резким улучшением результативности МО, которое стимулирует инновации в разных отраслях, включая инвестиционную индустрию.

Масштабы революции данных необычайны: только за последние два года было создано 90% всех данных, которые существуют в мире сегодня, и к 2020 г. каждый из 7,7 млрд людей во всем мире, как ожидается, будет производить 1,7 Мбайт новой информации каждую секунду каждый день. С другой стороны, еще в 2012 г. анализировалось и использовалось только 0,5% всех данных, в то время как по оценкам к 2020 г. 33% данных будет иметь добавочную стоимость. Разрыв между доступностью данных и их использованием, вероятно, быстро сократится, поскольку глобальные инвестиции в аналитику к 2020 г. превысят 210 млрд долларов, в то время как потенциал создания добавочной стоимости будет многократно выше.

В этой главе объясняется, как люди, производственные процессы и датчики генерируют альтернативные данные, а также каким образом обеспечивается фундамент для ориентации в расширяющемся предложении альтернативных данных для инвестиционных целей и их оценивания. В ней демонстрируется рабочий поток, от получения до предобработки и хранения данных, полученных с использованием Python посредством выскабливания веб-страниц с целью подготовки почвы для применения МО. В заключение приводятся примеры источников, поставщиков и приложений.

В этой главе будут рассмотрены следующие темы:

- ◆ как революция альтернативных данных породила новые источники информации;
- ◆ как физические лица, производственные процессы и датчики генерируют альтернативные данные;

- ◆ как оценивать расширяющееся предложение альтернативных данных, используемых для алгоритмической торговли;
- ◆ как работать с альтернативными данными на языке Python, в частности, получаемыми в результате выскабливания Интернета;
- ◆ важные категории и поставщики альтернативных данных.

Революция альтернативных данных

Лавина данных, стимулируемая цифровизацией, сетевым взаимодействием и резким падением расходов на хранение, привела к глубоким качественным изменениям в природе информации, доступной для предсказательной аналитики, часто резюмируемым пятью латинскими буквами V.

- ◆ *Объем (Volume)* — объем данных, генерируемых, собираемых и хранимых, на порядки больше, поскольку побочный продукт онлайн- и офлайн-активности, транзакций, записей и других источников и объемов продолжает расти с потенциалом для анализа и хранения.
- ◆ *Скорость (Velocity)*: данные генерируются, передаются и обрабатываются, становясь доступными с реально-временной скоростью или близкой к ней.
- ◆ *Разнообразие (Variety)* — данные организованы в форматах, которые больше не ограничиваются структурированными табличными формами, такими как CSV-файлы или таблицы реляционных баз данных. Вместо этого новые источники создают полуструктурированные форматы, такие как JSON или HTML, и неструктурированное информационное содержимое, включая сырые текстовые, снимковые и аудио- или видеоданные, ставя новые задачи по оформлению данных, подходящих для автоматически обучающихся алгоритмов.
- ◆ *Достоверность (Veracity)* — разнообразие источников и форматов значительно затрудняет подтверждение достоверности информационного содержимого данных.
- ◆ *Ценность (Value)* — определение добавочной стоимости новых наборов данных может быть гораздо более времязатратным и ресурсоемким, а также более неопределенным, чем раньше.

Для алгоритмической торговли новые источники данных предлагают информационное преимущество, если они предоставляют доступ к информации, недоступной из традиционных источников, или предоставляют доступ скорее. Следуя глобальным трендам, инвестиционная индустрия стремительно расширяется за пределы рыночных и фундаментальных данных в сторону альтернативных источников, пожиная альфа через информационное преимущество. Ежегодные расходы на данные, технологические возможности и связанный с ними потенциал, как ожидается, к 2020 г. будет увеличиваться с нынешних 3 млрд долларов на 12,8% ежегодно.

Сегодня инвесторы могут получать реально-временной доступ к макроэкономическим или корпоративным данным, которые исторически были доступны только с гораздо более низкой частотой. Вот примеры использования новых источников данных:

- ◆ онлайн-ценовые данные на репрезентативный набор товаров и услуг могут быть использованы для измерения инфляции;
- ◆ число посещений магазинов или покупок позволяет делать реально-временные оценки продаж или экономической активности отдельной компании или отрасли;
- ◆ спутниковые снимки могут выявлять сельскохозяйственные урожаи либо активность на шахтах или нефтяных вышках до того, как эта информация будет доступна в других местах.

По мере стандартизации и внедрения больших наборов данных информация, содержащаяся в обычных данных, вероятно, утратит большую часть своей предсказательной ценности.

Более того, потенциал обработки и интеграции разнообразных наборов данных и применения МО позволяет извлекать сущностную информацию. В прошлом количественные подходы опирались на простую эвристику, создавая рейтинги компаний с использованием исторических данных по таким метрическим показателям, как отношение рыночной цены к балансовой учетной стоимости, в то время как автоматически обучающиеся алгоритмы синтезируют новые метрики и усваивают и адаптируют такие правила с учетом эволюционирующих рыночных данных. Эти идеи создают новые возможности для улавливания классических инвестиционных тематик, таких как стоимость, импульс, качество или сентимент:

- ◆ *импульс* — МО способно выявлять влияние на финансовые активы движений рыночных цен, отраслевого сентимента или экономических факторов;
- ◆ *стоимость* — алгоритмы способны анализировать крупные объемы экономических и отраслевых структурированных и неструктурированных данных, помимо финансовой отчетности, с целью предсказания внутренней стоимости компании;
- ◆ *качество* — изощренный анализ интегрированных данных позволяет оценивать отзывы клиентов или сотрудников, электронную коммерцию или трафик приложений с целью идентификации прироста в доле рынка или других качественных движущих сил, лежащих в основе корпоративных заработков.

Вместе с тем на практике полезные данные зачастую не находятся в свободном доступе, а альтернативные наборы данных требуют тщательного оценивания, дорогостоящего их приобретения, тщательного менеджмента и изощренного анализа для извлечения торгуемых сигналов.

Источники альтернативных данных

Наборы альтернативных данных генерируются многими источниками, но могут быть классифицированы на высоком уровне как преимущественно создаваемые:

- ♦ *физическими лицами*, которые публикуют посты в социальных медиа, отзываются о продуктах или используют поисковые системы;
- ♦ *предприятиями*, которые регистрируют коммерческие транзакции, в частности кредитно-карточные платежи, или осуществляют деятельность в рамках снабженческой сети в качестве посредников;
- ♦ *датчиками*, которые, среди прочего, улавливают экономическую активность посредством снимков, в частности, со спутников или камер видеонаблюдения, или благодаря регулярностям продвижения, таким как вышки-ретрансляторы сотовой связи.

Природа альтернативных данных продолжает ускоренно эволюционировать по мере появления новых источников данных, в то время как источники, ранее именовавшиеся альтернативными, становятся частью магистрального потока. Например, *Балтийский фрахтовый индекс сухогрузного тоннажа* (BDI)¹ собирает данные от нескольких сотен судоходных компаний для аппроксимирования спроса/предложения на сухогрузы, и теперь доступен в терминале агентства Bloomberg.

Альтернативные данные включают сырые данные, а также данные, которые агрегированы или обработаны в той или иной форме с целью наращивания добавочной стоимости. Например, некоторые поставщики стремятся извлекать торгуемые сигналы, такие как сентиментные оценки, или отметки (score). Мы рассмотрим различные типы поставщиков в *главе 4*.

Альтернативные источники данных различаются в решающих аспектах, которые определяют их ценность или сигнальное содержимое для стратегий алгоритмической торговли. Мы рассмотрим эти аспекты в следующем разделе, посвященном оцениванию альтернативных наборов данных.

Физические лица

Физические лица автоматически создают электронные данные благодаря своей онлайн-активности, а также посредством своей автономной деятельности, поскольку последняя фиксируется электронным способом и часто связана с онлайн-выми удостоверениями личности. Данные, генерируемые физическими лицами, часто неструктурированные, находятся в текстовых, снимковых или видеоформатах и распространяются через несколько платформ.

¹ Балтийский фрахтовый индекс сухогрузного тоннажа (Baltic Dry Index, BDI, БФИ) — торговый индекс, ежедневно рассчитываемый Балтийской биржей (Baltic Exchange). Индекс отражает стоимость перевозок сухого груза (уголь, руда, зерно и т. п.) морем по двадцати основным торговым маршрутам. — *Прим. перев.*

Сюда входят:

- ◆ социально-медийные сообщения, такие как мнения или реакции на универсальных веб-сайтах, в том числе Twitter, Facebook или LinkedIn, либо на веб-сайтах отзывов о деятельности предприятий, таких как Glassdoor или Yelp;
- ◆ электронно-коммерческая активность, которая отражает интерес или восприятие продуктов на таких сайтах, как Amazon или Wayfair;
- ◆ поисковая активность с использованием платформ, таких как Google или Bing;
- ◆ использование мобильных приложений, скачиваний и отзывов;
- ◆ личные данные, такие как трафик обмена мгновенными сообщениями.

Сентиментный анализ, или анализ настроений, в социальных сетях стал очень популярным, потому что его можно применять к отдельным биржевым акциям, отраслевым корзинам или рыночным индексам. Наиболее распространенным источником является социальная сеть Twitter, вслед за которой идут различные поставщики новостей и блог-сайты. Предложение является конкурентоспособным, а цены становятся ниже, потому что это часто получается за счет неуклонно растущего коммодизирующего, т. е. устраняющего осмысленную товарную дифференциацию, веб-выскабливания. Надежные социально-медийные наборы данных, включающие блоги, твиты или видеоролики, как правило, имеют менее 5 лет истории, учитывая то, что потребители лишь недавно приняли эти инструменты в широком масштабе. Поисковая история, напротив, доступна с 2004 года.

Производственные процессы предприятия

Частные предприятия и государственные структуры производят и собирают многочисленные ценные источники альтернативных данных. Данные, получаемые в результате производственных процессов, часто имеют более высокую структурированность, чем данные, генерируемые физическими лицами. Они являются очень эффективными в качестве ведущего индикатора активности, который в противном случае доступен с гораздо более низкой частотой.

Данные, генерируемые производственными процессами, охватывают:

- ◆ данные платежно-карточных транзакций, предоставляемые процессорами и финансовыми учреждениями;
- ◆ выпускные корпоративные данные, производимые в результате обычной цифровизованной деятельности или ведения учета, такие как банковские выписки, данные сканера расчетно-кассового узла или заказы снабженческой сети;
- ◆ данные торговых потоков и микроструктуры рынка (например, данные ордерной книги L-2 и L-3, проиллюстрированные в *главе 2*);
- ◆ платежи компании, которые мониторятся кредитно-рейтинговыми агентствами или финансовыми учреждениями с целью оценивания ликвидности и кредитоспособности.

Кредитно-карточные транзакции и выпускные корпоративные данные, такие как данные о точках продаж, относятся к числу наиболее надежных и предсказуемых наборов данных. Кредитно-карточные данные доступны с историей примерно до 10 лет и, в разных временных лагах, вплоть до реального времени, в то время как отчеты о корпоративных заработках сообщаются ежеквартально с 2,5-недельным временным лагом. Временной горизонт и временной лаг отчетности по выпускным корпоративным данным широко варьируются в зависимости от источника. Рыночно-микроструктурные наборы данных имеют более чем 15-летнюю историю по сравнению с данными о потоках на стороне продажи, непрерывная история продаж которых, как правило, менее 5 лет.

Датчики

Данные, генерируемые сетевыми датчиками, встроенными в широкий спектр устройств, относятся к числу наиболее ускоренно растущих источников данных, стимулируемых распространением смартфонов и снижением стоимости спутниковых технологий.

Эта категория альтернативных данных, как правило, является очень неструктурированной и зачастую значительно большей по объему, чем данные, генерируемые физическими лицами или производственными процессами, и ставит задачи гораздо более высокой сложности в области обработки. Ключевые альтернативные источники данных в этой категории включают:

- ◆ получение спутниковых снимков с целью мониторинга экономической активности, такой как строительство, судоходство или снабжение товарами;
- ◆ геолокационные данные для отслеживания проходимости в розничных магазинах, например, с помощью добровольно предоставляемых смартфонных данных, или на транспортных маршрутах, например на судах или грузовиках;
- ◆ камеры, расположенные в интересующем месте;
- ◆ датчики погоды и загрязнения окружающей среды.

Интернет вещей (Internet of Things, IoT) еще больше ускорит масштабный сбор этого типа альтернативных данных путем внедрения сетевых микропроцессоров в персональные и коммерческие электронные устройства, такие как бытовая техника, общественные места и промышленные производственные процессы.

Датчиковые альтернативные данные, содержащие спутниковые снимки, данные потребления мобильных приложений или отслеживания местоположения сотовых телефонов, как правило, доступны с 3–4-летней историей.

Спутники

Ресурсы и сроки, необходимые для запуска спутника геопространственной визуализации резко сократились; вместо десятков миллионов долларов и лет подготовки стоимость упала до 100 тыс. долларов на то, чтобы разместить малый спутник в качестве вторичной полезной нагрузки для вывода на околоземную орбиту. Таким образом, компании могут получать гораздо более высокочастотный охват

(в настоящее время почти ежедневный) конкретных мест, используя целые флотилии спутников.

Примеры использования включают мониторинг экономической и коммерческой деятельности, которая может улавливаться с использованием аэрофотосъемочного покрытия, в частности сельскохозяйственного и минерального производства и поставок, строительства недвижимости или судов, промышленных инцидентов, таких как пожар или автомобильное и пешеходное движение в местах, представляющих интерес. Соответствующие датчиковые данные предоставляются дронами, которые используются в сельском хозяйстве для мониторинга сельскохозяйственных культур с использованием инфракрасного света.

Прежде чем данные спутниковых снимков можно будет надежно использовать в автоматически обучающихся моделях, возможно, потребуется решить ряд трудностей. Они включают учет погодных условий и, в частности, облачного покрова и сезонных эффектов, учет условий, например, приуроченных к праздникам, и нерегулярного покрытия конкретных мест — все они могут повлиять на качество предсказательных сигналов.

Геолокационные данные

Геолокационные данные — это еще одна ускоренно растущая категория альтернативных данных, генерируемых датчиками. Хорошо знакомым источником являются смартфоны, с помощью которых люди добровольно делятся своим географическим положением через приложение или с помощью беспроводных сигналов, таких как GPS, CDMA или Wi-Fi, измеряют пешеходный трафик вокруг интересующих их мест, таких как магазины, рестораны или места проведения мероприятий.

Более того, все большее число аэропортов, почтовых отделений и розничных магазинов устанавливают датчики, отслеживающие число и перемещения клиентов. Хотя первоначальная мотивировка к развертыванию этих датчиков часто заключалась в измерении воздействия маркетинговой деятельности, полученные данные могут также использоваться для оценивания проходимости или продаж. Датчики улавливания геолокации включают формирование трехмерного стерео видео- и тепловизионного изображения, что снижает озабоченности по поводу конфиденциальности, но хорошо работает с движущимися объектами. Кроме того, существуют датчики с крепежом к потолку, а также коврики, чувствительные к давлению. Некоторые поставщики используют несколько датчиков в сочетании, включая видео-, аудио- и мобильную локализацию для всестороннего учета передвижений покупателей, что включает не только число и продолжительность посещений, но и распространяется на конверсию и измерение повторных посещений.

Оценивание альтернативных наборов данных

Принципиальной целью альтернативных данных является предоставление информационного преимущества в конкурентном поиске торговых сигналов, которые производят альфа, а именно положительные, некоррелированные финансовые воз-

враты от инвестиций. На практике сигналы, извлеченные из альтернативных наборов данных, могут использоваться отдельно или комбинироваться с другими сигналами в рамках количественной стратегии. Независимое использование является жизнеспособным, если коэффициент Шарпа, генерируемый стратегией, основанной на одиночном наборе данных, достаточно высок, но это редко встречается на практике (подробности об измерении и оценивании сигналов см. в главе 4).

Квантовые фирмы создают библиотеки альфа-факторов, которые, являясь слабыми сигналами по отдельности, могут приносить привлекательные возмездия в сочетании. Как подчеркивалось в главе 1, инвестиционные факторы должны основываться на фундаментальном и экономическом обосновании, в противном случае они скорее всего являются результатом перепогонки к историческим данным, чем будут сохранять постоянство и генерировать альфа на новых данных.

Затухание сигнала вследствие конкуренции вызывает серьезную озабоченность, и по мере развития экосистемы альтернативных данных маловероятно, что многие наборы данных сохранят сигналы со значимым коэффициентом Шарпа. Эффективные стратегии продления периода полураспада² сигнального содержимого альтернативного набора данных включают соглашения об эксклюзивности или сосредоточение внимания на наборах данных, которые создают трудности с обработкой с целью повышения барьеров для входа.

Критерии оценивания

Альтернативный набор данных может быть оценен на основе качества его сигнального содержимого, качественных аспектов данных и различных технических аспектов.

Качество сигнального содержимого

Содержимое сигнала может быть оценено по целевому классу финансовых активов, стилю инвестирования, отношению к обычным рисковым премиям и, самое главное, по содержанию его альфа.

Классы финансовых активов

Большинство альтернативных наборов данных содержат информацию, непосредственно касающуюся долевого ценного бумага и товаров. Кроме того, интересные наборы данных, ориентированные на инвестиции в недвижимость, также умножились после того, как в 2006 г. компания Zillow, занимающаяся онлайн-базами данных недвижимости, успешно положила начало прогнозированию цен.

² Полураспад в физике — это время, затрачиваемое определенным количеством вещества на распад до половины его массы. Указанное понятие связано с измерением скорости этого процесса. Применительно к данной теме полураспад показывает медленность процесса либо время достижения ожидаемого значения. — *Прим. перев.*

По мере развития альтернативных источников мониторинга корпоративных платежей, в том числе для малого бизнеса, растут альтернативные данные корпоративного кредита. Прогнозные данные по ценным бумагам с фиксированной доходностью и процентным ставкам являются более свежим явлением, но продолжают увеличиваться по мере того, как все больше аккумулируется информации о продажах и ценах на продукцию.

Инвестиционный стиль

Большинство наборов данных сосредоточены на конкретных секторах и акциях и, как таковые, естественным образом, привлекательны для инвесторов в долевыми ценными бумагами с длинно-короткой инвестиционной стратегией. Поскольку масштабы и размах сбора альтернативных данных продолжают расти, альтернативные данные, вероятно, также станут актуальными для инвесторов в таких макроэкономических областях, как потребительский кредит, активность на развивающихся рынках и тренды в товарном секторе.

Некоторые альтернативные наборы данных могут использоваться в качестве индикаторов традиционных мер рыночного риска, в то время как другие сигналы актуальнее для высокочастотных торговцев, использующих количественные стратегии в течение короткого периода времени.

Рисковые премии

Как показывает практика, некоторые альтернативные наборы данных, такие как кредитно-карточные платежи или социально-медийный сентимент, дают сигналы, которые имеют низкую корреляцию (менее 5%) с традиционными рисковыми премиями³ на рынках долевыми ценными бумагами, такими как стоимость, импульс и качество волатильности. Как результат, комбинация сигналов, получаемых из таких альтернативных данных, со стратегией алгоритмической торговли, основанной на традиционных рисковых факторах, может стать важным строительным блоком для более диверсифицированного портфеля рисковых премий.

Содержимое альфа и его качество

Присущая сигналу сила, необходимая для обоснования инвестиции в альтернативный набор данных, естественным образом зависит от стоимости данных, при этом цены на альтернативные данные сильно варьируют. Данные, которые выставляют оценки социальным настроениям, можно получить за несколько тысяч долларов или меньше, в то время как стоимость набора данных о всеобъемлющих и регулярных кредитно-карточных платежах может стоить несколько миллионов в год.

Мы подробно разведем способы оценивания стратегий торговли, приводимых в действие альтернативными данными с использованием исторических данных, так

³ Рисковая премия (risk premium) — это финансовый возврат, превышающий безрисковый финансовый возврат, ожидаемый от инвестиций.

См. <https://www.investopedia.com/search?q=Risk+premium>. — *Прим. перев.*

называемые *бэктесты*, для оценивания величины альфа, содержащегося в наборе данных. В отдельных случаях набор данных может содержать альфа-сигнал, достаточный для того, чтобы приводить стратегию в действие автономно, но более типичной является ситуация, когда различные альтернативные и другие источники данных используются в сочетании. В таких случаях набор данных позволяет извлекать слабые сигналы, производящие малый положительный коэффициент Шарпа, который не смог бы самостоятельно получать капиталовложения, но вполне может обеспечивать стратегию портфельного уровня при его интегрировании с другими аналогичными сигналами. Однако это не гарантируется, поскольку существуют многочисленные альтернативные наборы данных, которые не имеют никакого альфа-содержимого.

Помимо оценивания содержимого альфа того или иного набора данных, также важно оценить, в какой степени сигнал является инкрементным или ортогональным, т. е. уникальным для этого набора данных или уже захваченным другими данными, и в последнем случае сравнивать затраты на этот тип сигнала.

Наконец, важно оценить потенциал стратегии, которая опирается на данный капитал, т. е. объем капитала, который может быть размещен без ущерба для ее успеха, поскольку лимит потенциала затруднит возмещение стоимости данных.

Качество данных

Качество набора данных является еще одним важным критерием, поскольку оно влияет на усилия, необходимые для их анализа и монетизации, а также на надежность содержащегося в них предсказательного сигнала. К аспектам качества относятся частота данных и продолжительность имеющейся в них истории, надежность или точность содержащейся в них информации, степень их соответствия действующему или потенциальному будущему регламенту и эксклюзивность их использования.

Правовые и репутационные риски

Использование альтернативных наборов данных может нести юридический или репутационный риск, в частности, если они включают следующие элементы:

- ◆ *материальную непубличную информацию* (Material Non-Public Information, MNPI), поскольку она влечет за собой нарушение правил инсайдерской торговли;
- ◆ *персонально идентифицируемую информацию* (Personally Identifiable Information, PII), в первую очередь с тех пор, как Европейский союз принял *общий регламент по защите данных* (General Data Protection Regulation, GDPR).

Соответственно, правовые требования и требования к соблюдению регламента нуждаются в тщательном рассмотрении. Кроме того, могут возникать конфликты интересов, когда поставщик данных одновременно является участником рынка, который активно торгует на основе этого набора данных.

Эксклюзивность

Вероятность того, что альтернативный набор данных содержит сигнал, достаточно предсказательный для того, чтобы приводить стратегию в действие на автономной основе с высоким коэффициентом Шарпа в течение значимого периода времени, обратно пропорциональна его доступности и простоте обработки. Другими словами, чем эксклюзивнее данные и чем сложнее их обрабатывать, тем больше шансов, что набор данных с альфа-содержимым сможет привести стратегию в действие без ускоренного затухания сигнала.

Публичные фундаментальные данные, которые снабжают стандартными финансовыми коэффициентами, содержат мало альфа и непривлекательны для автономной стратегии, но могут помочь диверсифицировать портфель рисков факторов. Для усвоения рынком больших и сложных наборов данных потребуется больше времени, и новые наборы данных будут появляться все чаще. Поэтому важно оценить, насколько другие инвесторы уже знакомы с тем или иным набором данных, и является ли этот поставщик наилучшим источником такого типа информации.

Дополнительные преимущества от эксклюзивности или раннего принятия нового набора данных могут возникать, когда предприятие только начинает продавать исчерпывающие данные, которые оно создало для других целей, поскольку в этом случае существует возможность повлиять на то, как данные собираются или курируются, или согласовать условия, которые лимитируют доступ конкурентов по крайней мере в течение определенного периода времени.

Временной горизонт

Более протяженная история крайне желательна для проверки предсказательной мощности набора данных в условиях разных сценариев. Ее наличие значительно варьирует от нескольких месяцев до нескольких десятилетий и имеет важные последствия для диапазона торговой стратегии, которая может быть построена и протестирована на основе таких данных. Мы упомянули некоторые диапазоны временных горизонтов для разных наборов данных, когда знакомили вас с главными типами источников.

Частота

Частота данных определяет, с какой периодичностью новая информация становится доступной и насколько дифференцированным может быть предсказательный сигнал за данный промежуток. Она также влияет на временной горизонт инвестиционной стратегии и варьируется от внутрисуточного, ежедневного, недельного или даже более низкой частоты.

Надежность

Естественно, степень, с которой данные точно отражают то, что они намереваются измерить, или насколько хорошо это может быть верифицировано, вызывает серьезную озабоченность, и она должна быть проверена посредством тщательного аудита. Это относится как к сырым, так и к обработанным данным, в которых необ-

ходимо проанализировать методологию, используемую для извлечения или агрегирования информации, с учетом отношения затрат к выгоде относительно предлагаемых к приобретению данных.

Технические аспекты

Технические аспекты касаются временной задержки или запаздывания отчетности и формата, в котором обеспечивается наличие этих данных.

Задержка

Поставщики данных часто предоставляют ресурсы пакетами, и временная задержка может быть вызвана процедурой сбора данных, последующей обработкой и передачей, а также регламентными или правовыми ограничениями.

Формат

Данные предоставляются в широком диапазоне форматов в зависимости от источника. Обработанные данные будут находиться в удобных для пользователя форматах и легко интегрироваться в существующие системы или запросы через надежный API. На другом конце спектра находятся объемистые источники данных, такие как видео, аудио или снимки, или проприетарный формат, которые требуют повышенных навыков их подготовки к анализу, но также обеспечивают более высокие барьеры для входа потенциальных конкурентов.

Рынок альтернативных данных

В 2018 г. инвестиционная индустрия собирается потратить на услуги передачи данных примерно 2–3 млрд долларов, и эта сумма, как ожидается, будет расти двузначными числами в год в соответствии с другими отраслями. Эти расходы включают получение альтернативных данных, инвестиции в соответствующую технологию и наем квалифицированных кадров.

Опрос аудиторско-консалтинговой компании Ernst and Young показывает, что в 2017 г. происходил значительный рост принятия альтернативных данных; например, 43% фондов используют данные, полученные в результате выскребывания Интернета, и почти 30% экспериментируют со спутниковыми данными. Основываясь на накопленном опыте, менеджеры фондов считают, что данные выскребывания Интернета и кредитно-карточные данные являются наиболее проникательными в отличие от геолокации и спутниковых данных, которые около 25% фондов считают менее информативными (рис. 3.1).

Рынок альтернативных поставщиков данных довольно фрагментирован, что как раз отражает ускоренный рост этой новой отрасли. Финансовый холдинг J. P. Morgan приводит список более 500 специализированных фирм данных, в то время как портал **AlternativeData.org** насчитывает более 300. Поставщики играют многочисленные роли, в том числе посреднические, в частности, в качестве консультантов,

агрегаторов и разработчиков технических решений; поддержка стороны продаж обеспечивает доставку данных в различных форматах, от сырых до полуобработанных данных или с сигналом, извлеченным в той или иной форме из одного или нескольких источников.

Для того чтобы проиллюстрировать их разнообразие, мы заострим внимание на размере главных категорий и рассмотрим несколько ярких примеров.

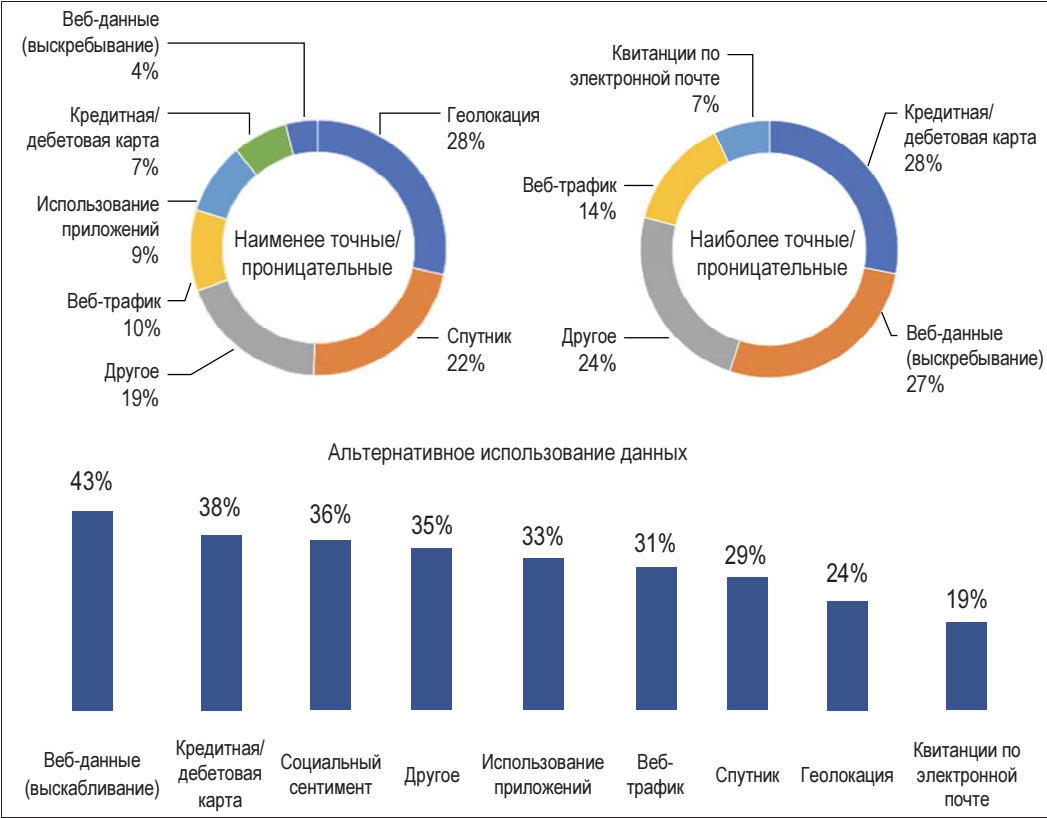


Рис. 3.1. Данные опроса, проведенного компанией Ernst and Young в 2017 г.

Поставщики данных и примеры использования

Портал **AlternativeData.org** (поддерживаемый провайдером YipitData) приводит несколько категорий, которые могут служить грубым индикатором деятельности в различных сегментах поставки данных (табл. 3.1). На сегодняшний день анализ социальных настроений является самой большой категорией, в то время как спутниковые и геолокационные данные ускоренно растут в последние годы.

Далее приведены небольшие примеры, которые призваны проиллюстрировать широкий круг поставщиков служб и потенциальных случаев их использования.

Таблица 3.1. Перечень категорий данных, согласно порталу *AlternativeData.org*

| Категория продукта | Число поставщиков | Цели |
|--|-------------------|--|
| Социальный сентимент | 48 | Сырые или обработанные социально-медийные данные, краткосрочные тренды |
| Спутники | 26 | Аэрофотосъемочный мониторинг среднесрочной экономической активности |
| Геолокация | 22 | Отслеживание розничного трафика, трафика коммерческой недвижимости или даже пешеходного движения |
| Веб-данные и трафик | 22 | Мониторинг поисковой заинтересованности, брендовой популярности и событий |
| Использование кредитных/дебетовых карт | 14 | Отслеживание ближнесрочных потребительских расходов и выручки предприятий |
| Использование приложений | 7 | Мониторинг продаж приложений или сбор вторичных данных |
| Электронные почтовые и потребительские квитанции | 6 | Отслеживание потребительских расходов по торговой сети, бренду, сектору или географии |
| Погода | 4 | Долгосрочные тренды, связанные с урожаем и сырьевыми товарами |
| Другие | 87 | — |

Данные социальных настроений

Анализ социальных настроений наиболее тесно связан с данными социальной сети Twitter. Компания Gnip была самым первым социально-медийным агрегатором, который предоставлял данные из многочисленных веб-сайтов с использованием API и был приобретен Twitter в 2014 г. за 134 млн долларов. Поисковые системы — это еще один источник, который стал заметным, когда исследователи опубликовали в журнале *Nature*, что инвестиционные стратегии, основанные на службе Google Тренды для таких терминов, как "долг", могут использоваться для прибыльной торговой стратегии в течение продолжительного периода (справочные материалы см. в репозитории GitHub по адресу <https://github.com/PacktPublishing/Hands-on-Machine-Learning-for-Algorithmic-Trading>).

Компания Datamir

Компания по поиску информации Datamir была основана в 2009 г. Она предоставляет анализ социальных настроений и новостей на основе эксклюзивного соглашения с Twitter. Данная компания является одним из крупнейших альтернативных

поставщиков, и в июне 2018 г. она привлекла дополнительные 392 млн долларов финансовых средств во главе с американской холдинговой компанией Fidelity, чья оцениваемая стоимость составляет в 1,6 млрд долларов, которая довела общее финансирование до 569 млрд долларов. Она делает акцент на реально-временных сигналах, извлекаемых из социально-медийных потоков с использованием машинного обучения, и обслуживает широкий круг клиентов, включая не только инвестиционные фирмы, находящиеся на стороне покупки и продажи, но и новостные организации и государственный сектор.

Социальная сеть StockTwits

StockTwits — это социальная сеть и микроблоговая платформа, где несколько сотен тысяч инвестиционных профессионалов обмениваются информацией и торговыми идеями в виде твитов StockTwit, которые просматриваются большой аудиторией по всей финансовой сети и социально-медийным платформам. Эти данные могут быть задействованы ввиду того, что они могут отражать сентимент инвесторов или сами по себе стимулировать сделки, которые, в свою очередь, воздействуют на цены. Справочные материалы в репозитории GitHub содержат ссылку на исследовательскую работу, в которой торговая стратегия строится на отобранных признаках.

Компания RavenPack

Поставщик аналитики больших данных компания RavenPack анализирует большое число разнообразных, неструктурированных, текстовых данных для получения структурированных индикаторов, включая сентиментные отметки, которые призваны содержать информацию, имеющую отношение к инвесторам. Лежащие в основе источники данных варьируются от премиальных новостных каналов и регламентной информации до пресс-релизов и более 19 тыс. веб-публикаций. Финансовый холдинг J. P. Morgan протестировал так называемые длинно-короткие стратегии инвестирования в суверенные облигации и долевые ценные бумаги на основе сентиментных отметок и достиг положительных результатов с низкой корреляцией с устоявшимися рисковыми премиями (см. справочные материалы).

Спутниковые данные

Поставщик приложений больших данных на основе ИИ компания RS Metrics, основанная в 2010 г., триангулирует геопространственные данные со спутников, дронов и самолетов, делая акцент на приложения, связанные с металлами и товарами, а также недвижимостью и промышленностью. Указанная компания предлагает сигналы, предсказательную аналитику, оповещения и приложения для конечных пользователей на основе собственных высокоразрешающих спутников. Примеры использования включают оценку розничного трафика с прицелом на определенные торговые сети или коммерческую недвижимость, а также производство и хранение определенных распространенных металлов или занятость на соответствующих производственных площадках.

Геолокационные данные

Компания Advan, основанная в 2015 г., обслуживает хедж-фондовых клиентов сигналами, выводимыми из данных мобильного трафика, ориентируясь на 1600 тикеров в различных секторах в США и ЕС. Эта компания собирает данные с помощью приложений, которые устанавливают геолокационные коды на смартфоны с прямым согласием пользователя и для повышения точности отслеживают местоположение с помощью нескольких каналов (таких как Wi-Fi, Bluetooth и сотовый сигнал). Примеры использования включают оценки клиентского трафика в физических местах хранения, которые, в свою очередь, могут использоваться в качестве входа в модели, предсказывающие валовую выручку торгуемых компаний.

Данные электронных почтовых квитанций

Поставщик решений на основе альтернативных данных компания Eagle Alpha, среди прочих услуг, предоставляет данные о крупном наборе онлайн-транзакций с использованием электронных почтовых квитанций, охватывающих более 5000 розничных торговцев, включая транзакционные данные уровня номенклатуры и единиц складского учета, классифицированные по 53 продуктовым группам. Финансовый холдинг J. P. Morgan проанализировал набор данных временных рядов, начиная с 2013 г., которые охватывали постоянную группу пользователей, активных на протяжении всего периода выборки. Набор данных содержал совокупные расходы, число заказов и число уникальных покупателей за определенный период.

Работа с альтернативными данными

Мы проиллюстрируем получение альтернативных данных с помощью веб-выскабливания, ориентируясь сначала на данные о рестораторах службы OpenTable, а затем перейдем к стенограммам телеконференций о корпоративных заработках, размещенных на портале Seeking Alpha.

Выскабливание данных службы OpenTable

Типичными источниками альтернативных данных являются веб-сайты с отзывами о деятельности предприятий, такие как Glassdoor или Yelp, которые передают инсайдерскую информацию, используя комментарии сотрудников или отзывы гостей. Эти данные служат ценным материалом для автоматически обучающихся моделей, которые ориентированы на предсказание перспектив предприятия или непосредственно его рыночной стоимости с целью получения торговых сигналов.

Данные должны быть извлечены из источника HTML, не принимая во внимание любые юридические препятствия. В качестве иллюстрации инструментов веб-выскабливания, которые предлагает язык Python, мы извлечем информацию о бронировании ресторанов из службы OpenTable. Данные такого рода могут использо-

ваться для прогнозирования экономической активности по географии, ценам на недвижимость или выручки ресторанной сети.

Извлечение данных из HTML с помощью библиотек requests и BeautifulSoup

В этом разделе мы запросим и разберем исходный код HTML. Мы будем использовать библиотеку requests для выполнения запросов по протоколу передачи гипертекста (HTTP) и извлечения исходного кода HTML и библиотеку BeautifulSoup для разбора и извлечения текстового содержимого.

Однако мы столкнемся с общим препятствием: веб-сайты могут запрашивать определенную информацию с сервера только после начальной загрузки страницы с использованием языка JavaScript. В результате прямой HTTP-запрос не будет успешным. Для того чтобы обойти этот тип защиты, мы будем использовать безголовый (headless) браузер, который извлекает содержимое веб-сайта, как это делает обычный браузер:

```
from bs4 import BeautifulSoup
import requests

# Задать и запросить URL; извлечь исходный код
url = "https://www.opentable.com/new-york-restaurant-listings"
html = requests.get(url)
html.text[:500]

'
    <!DOCTYPE html><html lang="en"><head><meta charset="utf-8"/>
<meta http-equiv="X-UA-Compatible" content="IE=9; IE=8; IE=7; IE=EDGE"/>
<title>Restaurant Reservation Availability</title>    <meta name="robots"
content="noindex" > </meta><link rel="canonical" href="https://www.opentable.com/
new-york-restaurant-listings" > </link>    <link rel="shortcut icon"
href="//components.otstatic.com/components/favicon/1.0.5/favicon/favicon.ico"
type="image/x-icon"/><link rel="icon" href="//comp'
```

Теперь можно применить библиотеку BeautifulSoup для разбора содержимого HTML, а затем разыскать все теги span с классом, ассоциированным с именами ресторанов, которые мы получаем путем обследования исходного кода, rest-row-name-text (см. репозиторий GitHub, где приводятся ссылки на инструкции по обследованию исходного кода веб-сайта):

```
# Разобрать сырой html => объект BeautifulSoup
soup = BeautifulSoup(html.text, 'html.parser')

# Для каждого тега span распечатать текст => имя ресторана
for entry in soup.find_all(name='span', attrs={'class': 'rest-row-name-text'}):
    print(entry.text)

Clarks
Kessler
```

```
Hane
Hansen
At Cremin
...
```

После того как вы определили интересующие элементы страницы, библиотека BeautifulSoup позволяет легко получить содержащийся текст. Если вы хотите получить ценовую категорию для каждого ресторана, то можно применить следующий фрагмент кода:

```
# Получить число долларовых знаков для каждого ресторана
for entry in soup.find_all('div', {'class': 'rest-row-pricing'}):
    price = entry.find('i').text
```

Однако при попытке получить число заказов вы получите пустой список, потому что для запроса этой информации веб-сайт использует код на языке JavaScript после завершения начальной загрузки страницы:

```
soup.find_all('div', {'class': 'booking'})
[]
```

Знакомство с библиотекой Selenium — использование браузерной автоматизации

Мы воспользуемся инструментом браузерной автоматизации Selenium для управления безголовым браузером Firefox, который будет анализировать содержимое HTML за нас.

Следующий фрагмент кода открывает браузер Firefox:

```
from selenium import webdriver
from selenium.webdriver.common.desired_capabilities import DesiredCapabilities

cap = DesiredCapabilities().FIREFOX

# Создать драйвер с именем Firefox
driver = webdriver.Firefox(capabilities=cap)
```

Давайте закроем браузер:

```
# Закроем его
driver.close()
```

Для извлечения исходного кода HTML с помощью библиотеки Selenium и браузера Firefox выполните следующие действия:

```
import time, re

# Посетить страницу opentable со списком
driver = webdriver.Firefox(capabilities=cap)
driver.get(url)
time.sleep(1) # Подождать 1 секунду
```



```
# Извлечь исходный код HTML
html = driver.page_source
html = BeautifulSoup(html, "lxml")

for booking in html.find_all('div', {'class': 'booking'}):
    match = re.search(r'\d+', booking.text)
    if match:
        print(match.group())
```



В среде ОС Windows 10 для работы браузера selenium надо скачать запакованный драйвер geckodriver с указанного адреса: <https://github.com/mozilla/geckodriver/releases/tag/v0.24.0>, распаковать исполнимый файл geckodriver.exe и разместить его в подкаталоге языка Python, задать строковую переменную с путем к драйверу, например, как показано ниже:

```
from selenium import webdriver
from selenium.webdriver.common.desired_capabilities import
DesiredCapabilities

cap = DesiredCapabilities().FIREFOX
path = "C:\\Python37\\Drivers\\geckodriver\\geckodriver.exe"
```

и использовать следующую инструкцию, возможно, еще добавив путь к установленному браузеру Firefox, например:

```
driver = webdriver.Firefox(
    capabilities=cap,
    executable_path=path,
    firefox_binary='C:\\Program Files\\Mozilla Firefox'
)
```

Построение набора данных бронирований ресторанов

Теперь для создания признака, который можно было бы использовать в модели для предсказания экономической активности в географических регионах или проходимости клиентов в определенных окрестностях, вам нужно только объединить все интересующие элементы веб-сайта.

С помощью библиотеки Selenium можно переходить по ссылкам на следующие страницы и быстро создавать набор данных из более чем 10 тыс. ресторанов в Нью-Йорке, который затем можно периодически обновлять для отслеживания временного ряда. Сначала мы формируем функцию разбора содержимого страниц, которые планируем обойти:

```
def parse_html(html):
    """Разобрать содержимое различных тегов
    из списка ресторанов OpenTable"""
```

```

data, item = pd.DataFrame(), {}
soup = BeautifulSoup(html, 'lxml')
for i, resto in enumerate(soup.find_all('div',
                                         class_='rest-row-info')):
    item['name'] = resto.find('span',
                             class_='rest-row-name-text').text

    booking = resto.find('div', class_='booking')
    item['bookings'] = re.search('\d+', booking.text).group() \
        if booking else 'NA'

    rating = resto.select('div.all-stars.filled')
    item['rating'] = int(re.search('\d+',
                                   rating[0].get('style')).group()) \
        if rating else 'NA'

    reviews = resto.find('span',
                          class_='star-rating-text--review-text')
    item['reviews'] = int(re.search('\d+', reviews.text).group()) \
        if reviews else 'NA'

    item['price'] = int(resto.find('div',
                                   class_='rest-row-pricing').find('i').text.count('$'))
    item['cuisine'] = resto.find('span',
                                 class_='rest-row-meta--cuisine').text
    item['location'] = resto.find('span',
                                  class_='rest-row-meta--location').text
    data[i] = pd.Series(item)
return data.T

```

Затем мы запускаем безголовый браузер, который продолжает нажимать кнопку **Next (Далее)** за нас и выхватывать результаты, отображаемые на каждой странице:

```

restaurants = pd.DataFrame()
driver = webdriver.Firefox(capabilities=cap)
url = "https://www.opentable.com/new-york-restaurant-listings"
driver.get(url)
page = collected = 0
while True:
    sleep(1)
    new_data = parse_html(driver.page_source)
    if new_data.empty:
        break
    if page == 0:
        new_data.to_csv('results.csv', index=False)
    elif page > 0:
        new_data.to_csv('results.csv', index=False, header=None, mode='a')

```

```
page += 1
collected += len(new_data)
print(f'Страница: {page} | Скачано: {collected}')
driver.find_element_by_link_text('Next').click()

driver.close()
restaurants = pd.read_csv('results.csv')
print(restaurants)
```

Веб-сайты продолжают меняться, поэтому в какой-то момент этот исходный код может перестать работать и потребует обновления для отслеживания свежей версии навигации по сайту и обнаружения программным роботом.

Следующий шаг — библиотеки Scrapy и splash

Библиотека Scrapy — это мощная библиотека, которая предназначена для построения программных роботов. Такие роботы следуют по ссылкам, извлекают содержимое и сохраняют разобранный результат в структурированном виде. В сочетании с безголовым браузером splash она также может интерпретировать исходный код на JavaScript и становиться эффективной альтернативой библиотеке Selenium. С помощью команды `scrapy crawl opentable` в каталоге `01_opentable` можно запустить веб-обходчика. В указанном каталоге результаты регистрируются в журнале `spider.log`:

```
from opentable.items import OpentableItem
from scrapy import Spider
from scrapy_splash import SplashRequest

class OpenTableSpider(Spider):
    name = 'opentable'
    start_urls =
        ['https://www.opentable.com/new-york-restaurant-listings']

    def start_requests(self):
        for url in self.start_urls:
            yield SplashRequest(url=url,
                                callback=self.parse,
                                endpoint='render.html',
                                args={'wait': 1},
                                )

    def parse(self, response):
        item = OpentableItem()
        for resto in response.css('div.rest-row-info'):
            item['name'] = resto.css('span.rest-row-name-text::text').extract()
            item['bookings'] = resto.css('div.booking::text').re(r'\d+')
            item['rating'] = resto.css('div.all-stars::attr(style)').re_first('\d+')
            item['reviews'] = \
                resto.css('span.star-rating-text--review-text::text').re_first(r'\d+')
```

```

item['price'] = len(resto.css('div.rest-row-pricing > i::text')).re('\$'))
item['cuisine'] = resto.css('span.rest-row-meta--cuisine::text').extract()
item['location'] = resto.css('span.rest-row-meta--location::text').extract()
yield item

```

Помимо извлечения отзывов об отдельных ресторанах или сетях и заказов в них, существует много других способов извлечения информации из этих данных.

К примеру, можно было бы также собрать и геокодировать адреса ресторанов для того, чтобы связать физическое местоположение ресторанов с другими интересующими районами, такими как популярные розничные торговые точки или окрестности с целью проникновения в суть конкретных аспектов экономической активности. Как упоминалось ранее, такие данные будут иметь особую ценность в сочетании с другой информацией.

Стенограммы телеконференций о корпоративных заработках

Текстовые данные являются важным альтернативным источником данных. Одним из примеров текстовой информации являются стенограммы телеконференций о корпоративных заработках, в которых руководители не только представляют последние финансовые результаты, но и отвечают на вопросы финансовых аналитиков. Инвесторы задействуют эти стенограммы для оценивания изменений в сентименте, акцентов на конкретных темах или стиля общения.

Мы проиллюстрируем выскабливание и стенограммы телеконференций о корпоративных заработках из популярного торгового портала **www.seekingalpha.com**:

```

import re
from pathlib import Path
from random import random
from time import sleep
from urllib.parse import urljoin

import pandas as pd
from bs4 import BeautifulSoup
from furl import furl
from selenium import webdriver
from selenium.webdriver.common.desired_capabilities import DesiredCapabilities

cap = DesiredCapabilities().FIREFOX

transcript_path = Path('transcripts')

SA_URL = 'https://seekingalpha.com/'
TRANSCRIPT = re.compile('Earnings Call Transcript')
# path = "C:\\Python37\\Drivers\\geckodriver\\geckodriver.exe"

```

```

next_page = True
page = 1
driver = webdriver.Firefox(capabilities=cap) # , executable_path=path
meta = {}
while next_page:
    print(f'Page: {page}')
    url = f'{SA_URL}/earnings/earnings-call-transcripts/{page}'
    driver.get(urljoin(SA_URL, url))
    response = driver.page_source
    page += 1
    soup = BeautifulSoup(response, 'lxml')
    links = soup.find_all(name='a', string=TRANSCRIPT)
    if len(links) == 0:
        next_page = False
    else:
        for link in links:
            transcript_url = link.attrs.get('href')
            article_url = furl(urljoin(SA_URL,
                                       transcript_url)).add({'part': 'single'})
            driver.get(article_url.url)
            html = driver.page_source
            result = parse_html(html)
            if result is not None:
                meta, participants, content = result
                meta['link'] = link
                store_result(meta, participants, content)
                sleep(5 + (random() - .5) * 2)

driver.close()

```

Разбор HTML с помощью регулярных выражений

Для сбора структурированных данных из неструктурированных стенограмм в дополнение к библиотеке BeautifulSoup можно использовать регулярные выражения.

Они позволяют собирать подробную информацию не только о компании и хронометраже телеконференции о корпоративных заработках, но и фиксировать присутствовавших и приписывать высказывания аналитикам и представителям компании:

```

def parse_html(html):
    # Главный анализатор HTML
    date_pattern = re.compile(r'(\d{2})-(\d{2})-(\d{2})')
    quarter_pattern = re.compile(r'(\bQ\d\b)')
    soup = BeautifulSoup(html, 'lxml')

    meta, participants, content = {}, [], []
    h1 = soup.find('h1', itemprop='headline')

```

```

if h1 is None:
    return
h1 = h1.text
meta['company'] = h1[:h1.find('(')].strip()
meta['symbol'] = h1[h1.find('(') + 1:h1.find(')')]

title = soup.find('div', class_='title')
if title is None:
    return
title = title.text
print(title)
match = date_pattern.search(title)
if match:
    m, d, y = match.groups()
    meta['month'] = int(m)
    meta['day'] = int(d)
    meta['year'] = int(y)

match = quarter_pattern.search(title)
if match:
    meta['quarter'] = match.group(0)

qa = 0
speaker_types = ['Executives', 'Analysts']
for header in [p.parent for p in soup.find_all('strong')]:
    text = header.text.strip()
    if text.lower().startswith('copyright'):
        continue
    elif text.lower().startswith('question-and'):
        qa = 1
        continue
    elif any([type in text for type in speaker_types]):
        for participant in header.find_next_siblings('p'):
            if participant.find('strong'):
                break
        else:
            participants.append([text, participant.text])
    else:
        p = []
        for participant in header.find_next_siblings('p'):
            if participant.find('strong'):
                break
        else:
            p.append(participant.text)
        content.append([header.text, qa, '\n'.join(p)])
return meta, participants, content

```

Мы сохраняем результат в нескольких CSV-файлах с целью легкого доступа во время применения МО для обработки естественного языка:

```
def store_result(meta, participants, content):
    path = transcript_path / 'parsed' / meta['symbol']
    if not path.exists():
        path.mkdir(parents=True, exist_ok=True)

    pd.DataFrame(content, columns=['speaker', 'q&a',
                                   'content']).to_csv(path / 'content.csv',
                                                       index=False)
    pd.DataFrame(participants, columns=['type',
                                         'name']).to_csv(path / 'participants.csv',
                                                           index=False)
    pd.Series(meta).to_csv(path / 'earnings.csv')
```

Подробности и справочные материалы о дополнительных ресурсах с расширенными возможностями для разработки приложений веб-выскабливания см. в файле README репозитория GitHub.

Резюме

В этой главе мы представили новые источники альтернативных данных, ставшие доступными в результате революции больших данных, включая физических лиц, производственные процессы и датчики, такие как спутники или устройства определения местоположения GPS. Мы представили каркас для оценивания альтернативных наборов данных с инвестиционной точки зрения и изложили ключевые категории и поставщиков с целью помочь вам ориентироваться в этой обширной и быстро расширяющейся области, которая обеспечивает входные данные, критически важные для стратегий алгоритмической торговли, использующих МО.

Мы развели мощные инструменты Python для сбора собственных наборов данных в широком масштабе, что позволяет в перспективе работать над получением своего частного информационного преимущества в качестве алгоритмического торговца, используя веб-выскабливание.

В следующей главе мы перейдем к разработке и оцениванию альфа-факторов, которые генерируют торговые сигналы, и рассмотрим способы их комбинирования в контексте портфеля.

4

Исследование альфа-факторов

Стратегии алгоритмической торговли приводятся в действие сигналами, которые указывают на моменты, когда следует покупать или продавать финансовые активы для получения положительных финансовых возвратов относительно эталона. Часть возврата от финансового актива, которая не объясняется влиянием эталона, называется *альфой*¹, и поэтому эти сигналы называются *альфа-факторами*.

Альфа-факторы призваны предсказывать ценовые движения финансовых активов в инвестиционном универсуме на основе имеющихся рыночных, фундаментальных или альтернативных данных. Фактор может объединять в себе одну или несколько входных переменных, но принимает единственное значение для каждого актива всякий раз, когда фактор оценивается стратегией. Торговые решения, как правило, опираются на относительные стоимости по всем финансовым активам. Стратегии торговли часто основываются на сигналах, эмитируемых многочисленными факторами, и мы увидим, что автоматически обучающиеся модели особенно хорошо подходят для эффективной интеграции различных сигналов, позволяя делать более точные предсказания.

Как показано на рис. 4.1, конструирование, оценивание и комбинирование альфа-факторов являются критически важными шагами во время исследовательской фазы рабочего потока алгоритмической торговой стратегии. В этой главе мы сосредоточимся на исследовательской фазе, а в следующей *главе 5* — на исполнительной фазе. Остальная часть этой книги будет посвящена использованию МО для обнаружения и комбинирования альфа-факторов.

В этой главе будет использован простой фактор разворота к среднему ценовому уровню (mean-reversal) с целью введения библиотеки симуляции алгоритмической

¹ Альфа (alpha) — это показатель, рассчитываемый для ценной бумаги или портфеля ценных бумаг, связывающий возвратность ценной бумаги (портфеля) с возвратностью близкого фондового индекса. Высокое значение альфа означает, что ценная бумага (портфель) работает лучше, чем ожидалось при его заданном бета (волатильности). Термин заимствован из статистики, где альфа (уровень значимости) — это вероятностный порог "необычности", который случайные результаты должны превзойти, чтобы фактические результаты считались статистически значимыми.

См. <https://ru.wikipedia.org/wiki/Альфа-коэффициент>. — Прим. перев.

торговли `zipline`, написанной на Python и обеспечивающей тестирование альфа-факторов для заданного инвестиционного универсума. Мы также воспользуемся библиотекой `zipline` в следующей главе при бэктестировании торговых стратегий в контексте портфеля. Далее мы обсудим ключевые метрики, служащие для оценивания предсказательной результативности альфа-факторов, включая информационный коэффициент и информационное соотношение, которые приводят к фундаментальному закону активного менеджмента.

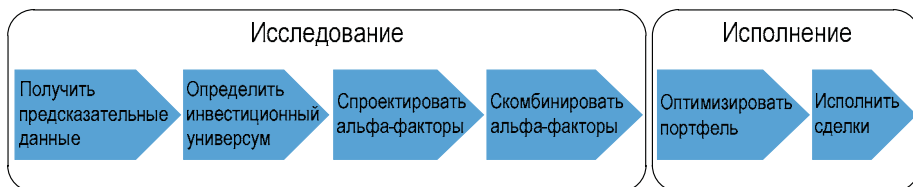


Рис. 4.1. Фазы рабочего потока алгоритмической торговой стратегии

В частности, в этой главе будут рассмотрены следующие темы:

- ◆ как характеризовать, обосновывать и измерять ключевые типы альфа-факторов;
- ◆ как создавать альфа-факторы с помощью выработки финансовых признаков;
- ◆ как использовать `zipline` в режиме офлайн для тестирования индивидуальных альфа-факторов;
- ◆ как применять библиотеку `zipline` на платформе Quantopian для комбинирования альфа-факторов и выявления более сложных сигналов;
- ◆ как информационный коэффициент измеряет предсказательную результативность альфа-фактора;
- ◆ как использовать библиотеку `alphalens` для оценивания предсказательной результативности и оборачиваемости.

Выработка альфа-факторов

Альфа-факторы — это преобразования рыночных, фундаментальных и альтернативных данных, содержащих предсказательные сигналы. Они призваны улавливать риски, стимулирующие возвраты от финансовых активов. Один набор факторов описывает фундаментальные, общеэкономические переменные, такие как рост, инфляцию, волатильность, производительность и демографический риск. Еще один набор состоит из торгуемых инвестиционных стилей, таких как инвестирование в рыночный портфель, инвестирование в рост стоимости и импульсное инвестирование.

Существуют также факторы, которые объясняют ценовые движения на основе экономических или институциональных условий финансовых рынков или поведения инвесторов, включая известные тенденциозности (смещения) такого поведения.

Экономическая теория, стоящая за факторами, может быть рациональной, когда факторы имеют высокие возвраты в течение продолжительного времени, компенсируя их низкие возвраты в плохие времена, либо поведенческой, где премии за факторные риски являются результатом, возможно, тенденциозного (смещенного) или не совсем рационального поведения агентов, которые не задействует арбитраж.

Постоянно ведется поиск и выявление новых факторов, которые могли бы лучше улавливать известные либо новые движущие силы финансовых возвратов. Джейсон Хсу (Jason Hsu), соучредитель компании по выявлению умных бета и распределению портфельных инвестиций Research Affiliates, управляющей около 200 млрд долларов, выявил около 250 факторов, которые к 2015 г. были опубликованы с эмпирическими наблюдениями в авторитетных журналах, и оценивает, что это число, вероятно, будет увеличиваться на 40 факторов в год. Для того чтобы избежать ложных обнаружений и гарантировать, что фактор дает непротиворечивые результаты, он должен обладать значимой экономической интуицией, дающей основания полагать, что указанный фактор отражает риски, которые будут компенсированы рынком.

Преобразования данных включают простую арифметику, такую как абсолютные или относительные изменения переменной во времени, соотношения между рядами данных или агрегации во временном окне, такие как простое или экспоненциальное скользящее среднее. Они также охватывают расчеты, получаемые в результате технического анализа ценовых регулярностей, таких как индекс относительной силы спроса в сопоставлении с предложением и многочисленные метрики, знакомые по фундаментальному анализу ценных бумаг.

Важные категории факторов

В идеализированном мире категории рисковых факторов должны быть независимыми друг от друга (ортогональными), давать положительные рисковые премии и формировать полное множество, охватывающее все размерности риска и объяснять систематические риски для финансовых активов того или иного класса. На практике эти требования соблюдаются лишь приближенно. В *главе 12* мы рассмотрим способы выведения синтетических ведомых данными рисковых факторов с использованием неконтролируемого обучения, в частности, анализа главных и независимых компонент.

Мы проведем обзор ключевых категорий факторов, выводимых из рыночных, фундаментальных и альтернативных данных, и типичные метрики, используемые для их улавливания. Мы также продемонстрируем способы реализации этих факторов для алгоритмов, протестированных на платформе Quantopian с использованием встроенных факторов, собственных вычислений с помощью библиотек NumPy и pandas или библиотеки TA-Lib, специально предназначенной для технического анализа.

Импульсные и сентиментные факторы

Импульсное инвестирование следует поговорке: тренд — ваш друг, или пусть победители работают. Импульсные рискованные факторы предназначены для того, чтобы лонговать на активах, которые показывали хорошую результативность, и шортить на активах с низкой результативностью за определенный период².

Предпосылка стратегий, опирающихся на этот фактор, заключается в том, что цены финансовых активов демонстрируют тренд, отражающийся в положительных внутридневных корреляциях. Такой ценовой импульс бросает вызов гипотезе об эффективных рынках, которая утверждает, что прошлые ценовые возвраты сами по себе не в состоянии предсказывать будущую результативность. Несмотря на теоретические аргументы об обратном, стратегии ценового импульса производили и производят положительные возвраты по всем классам финансовых активов и являются важной частью многих стратегий торговли на финансовых рынках.

Обоснование

Причины импульсного эффекта указывают на поведение инвесторов, устойчивые дисбалансы спроса и предложения, цикл положительной обратной связи между рискованными активами и экономикой или рыночной микроструктурой.

Поведенческие причины отражают тенденциозности (смещения) из-за недостаточного и чрезмерного реагирования на рыночные новости, поскольку инвесторы обрабатывают новую информацию с разной скоростью. После первоначальной недостаточной реакции на новости инвесторы часто экстраполируют прошлое поведение и создают ценовой импульс. Ралли технологических акций во время рыночного пузыря конца 1990-х годов было экстремальным примером. Психология страха и жадности также мотивирует инвесторов увеличивать влияние выигрышных активов и продолжать продавать убыточные активы.

Импульс также может иметь фундаментальные движущие силы, такие как цикл положительной обратной связи между рискованными активами и экономикой. Экономический рост подстегивает рост стоимости долевых ценных бумаг, и эффект результирующего богатства возвращается по каналу обратной связи в экономику через более высокие расходы, снова подпитывая рост. Положительная обратная связь между ценами и экономикой часто расширяет импульс в долевых ценных бумагах и кредитах до более длинных горизонтов, чем для облигаций, валют и товаров, где отрицательная обратная связь создает развороты, требующие гораздо более короткого инвестиционного горизонта. Еще одной причиной импульса может быть устойчивый дисбаланс спроса и предложения из-за рыночных трений, например,

² Лонговать (go long) — покупать активы, открывать длинную позицию на покупку в надежде на рост цены актива, причем это вовсе не означает, что позиция будет удерживаться продолжительное время. Торговец "лонгует", когда цена растет, давая ему возможность получать прибыль от покупки. Шортить (go short) — продавать активы, открывать короткую позицию на продажу в надежде на снижение цены актива. Торговец "шортит", когда цена падает, давая ему возможность получать прибыль от продажи. — *Прим. перев.*

когда товарное производство занимает значительные количества времени на то, чтобы скорректироваться под тренды спроса. Добыча нефти может годами задерживать повышенный спрос со стороны быстро развивающейся экономики, а устойчивый дефицит предложения может провоцировать и поддерживать ценовой импульс.

Эффекты рыночной микроструктуры также могут создавать ценовой импульс, связанный с поведенческими регулярностями, мотивирующими инвесторов покупать продукты и реализовывать стратегии, которые подражают ихтенденциозностям (смещениям). Например, торговая мудрость урезать убытки и давать прибыли работать побуждает инвесторов использовать такие стратегии торговли, как стоп-лосс (остановка убытка)³, *портфельное страхование с постоянной пропорцией*⁴, динамическое дельта-хеджирование или опционные стратегии, такие как защитные опционы put⁵. Эти стратегии создают импульс, поскольку они влекут за собой предварительное обязательство продавать, когда актив показывает себя менее успешно, и покупать, когда он показывает себя более успешно. Схожим образом, риск-паритетные стратегии (см. следующую главу) тяготеют к покупке низковолатильных активов, которые часто демонстрируют положительную результативность, и к продаже высоковолатильных активов, которые часто имели отрицательную результативность. Автоматическая перебалансировка портфелей с использованием этих стратегий чаще всего усиливает ценовой импульс.

³ Стратегия "стоп-лосс" — это стратегия торговли с преобладанием ордеров стоп-лосс, т. е. размещением сделок на бирже для немедленного их исполнения, в случае если актив достигает определенной ценовой точки. Как следует из названия, этот вид ордера предназначен для ограничения убытков. — *Прим. перев.*

⁴ Портфельное страхование с постоянной пропорцией (Constant Proportion Portfolio Insurance, CPPI) — это тип портфельного страхования, в котором инвестор устанавливает нижний уровень долларовой стоимости своего портфеля, а затем структурирует размещение активов вокруг этого решения. В CPPI используется два класса активов: рискованный актив (обычно акционерные капиталы или взаимные фонды) и консервативный актив, состоящий из денежных средств, их эквивалентов или казначейских облигаций. Процент, выделенный каждому из них, зависит от значения "подушки", определенной как разность текущей стоимости портфеля и значения нижнего уровня, и мультипликативного коэффициента, где большее число обозначает более агрессивную стратегию.

См. <https://www.investopedia.com/terms/c/cppi.asp>. — *Прим. перев.*

⁵ Защитный опцион put (protective put) — это стратегия риск-менеджмента, которую инвесторы используют для защиты от потери нереализованной прибыли в акциях или других активах. Опцион put действует как страховой полис — он стоит денег, которые уменьшает потенциальную прибыль инвестора от владения ценной бумагой, но и сокращает риск потери денег, если стоимость ценной бумаги снижается. См. <https://www.investopedia.com/terms/p/protective-put.asp>. — *Прим. перев.*

Опцион put (put option) — опционный контракт, дающий владельцу право, но не обязательство, продать определенное количество базового актива по оговоренной цене в течение определенного периода времени. Поскольку базовый актив обесценивается по стоимости, сам опцион put растет в цене; поэтому покупка опционов put является методом его шортирования (занятия короткой позиции). Является противоположностью опциона call, который дает держателю право купить базовую ценную бумагу по указанной цене до истечения срока действия опциона.

См. <https://www.investopedia.com/terms/p/putoption.asp>. — *Прим. перев.*

Ключевые метрические показатели

Импульсные факторы, как правило, выводятся из изменений в ценовом временном ряде путем выявления трендов и регулярностей. Они могут конструироваться на основе абсолютного или относительного финансового возврата, путем сравнения среза (перекрестного сечения) активов или анализа временного ряда актива внутри или между традиционными классами активов и на разных временных горизонтах.

Несколько популярных наглядных показателей приведено в табл. 4.1.

Таблица 4.1. Импульсные индикаторы

| Фактор | Описание |
|--|---|
| Индикатор относительной силы (relative strength indicator, RSI) | Индикатор RSI сравнивает недавние ценовые изменения по акциям с целью идентификации перекупленности или перепроданности акций. Высокий RSI (например, выше 70) указывает на перекупленность, а низкий RSI (например, ниже 30) — на перепроданность. Он использует среднее ценовое изменение за заданное число предыдущих торговых дней с положительными ценовыми изменениями Δ_p^{-up} и отрицательными ценовыми изменениями Δ_p^{-down} , вычисляя: $RSI = 100 - \frac{100}{1 + \Delta_p^{-up} / \Delta_p^{-down}}$ |
| Ценовой импульс (price momentum) | Указанный фактор вычисляет общий возврат за заданное число предыдущих торговых дней. В академической литературе принято использовать последние 12 месяцев, но исключать самый недавний месяц из-за эффекта краткосрочного разворота, часто наблюдаемого в самых недавних ценовых движениях, но более короткие периоды также широко используются |
| 12-месячный ценовой импульс с поправкой на волатильность (12-month price momentum vol adj) | 12-месячный ценовой импульс с поправкой на фактор волатильности нормализует суммарный возврат за предыдущие 12 месяцев, деля его на стандартное отклонение этих возвратов |
| Ценовое ускорение (price acceleration) | Ценовое ускорение вычисляет градиент тренда (с поправкой на волатильность), используя линейную регрессию на дневных ценах в течение более длительного и более короткого периода, например, один год и три месяца торговых дней, и сравнивает изменение наклона прямой в качестве меры ценового ускорения |
| Процентное отклонение от 52-недельной максимальной цены (percent off 52 week high) | Указанный фактор использует процентную разницу между самой недавней и самой высокой ценой за последние 52 недели |

Дополнительные сентиментные индикаторы представлены в табл. 4.2.

Таблица 4.2. Сентиментные индикаторы

| Фактор | Описание |
|---|--|
| Количество оценок заработков (earnings estimates count) | Указанный фактор ранжирует акции по числу консенсусных оценок в качестве индикатора аналитического покрытия и информационной неопределенности. Чем выше значение, тем лучше |
| <i>N</i> -месячное изменение в рекомендации (N month change in recommendation) | Указанный фактор ранжирует акции по изменению консенсусной рекомендации за предыдущий <i>N</i> -месячный период, где желательно, чтобы присутствовало улучшение (независимо от того, перешли ли они от активной продажи к продаже или от покупки к активной покупке и т. д.) |
| 12-месячное изменение в долях (паях) в обращении (12-month change in shares outstanding) | Указанный фактор служит мерой изменения числа долей в предприятии с поправкой на дробление за последние 12 месяцев, где отрицательное изменение подразумевает выкуп акций и является желательным, поскольку сигнализирует о том, что руководство рассматривает акции как дешевые по отношению к их внутренне присущей и, следовательно, будущей стоимости |
| 6-месячное изменение в целевой цене (6-month change in target price) | Указанный фактор отслеживает 6-месячное изменение в средней целевой аналитической цене, и более высокое положительное изменение, естественно, является более желательным |
| Пересмотры чистых заработков (net earnings revisions) | Указанный фактор отражает разницу между пересмотрами в сторону повышения и понижения оценок заработков в процентах от общего числа пересмотров |
| Уровень коротких продаж ⁶ по отношению к акциям в обращении (short interest to shares outstanding) | Указанный фактор представляет собой процент акций (долевых ценных бумаг) в обращении, которые в настоящее время продаются, т. е. продаются инвестором, который занял акцию (долевою ценную бумагу) и должен выкупить ее в более поздний день, спекулируя на том, что ее цена упадет. Следовательно, высокий уровень коротких продаж указывает на негативные настроения и, как ожидается, сигнализирует о плохой результативности в перспективе |

⁶ Уровень коротких продаж (short interest), т. е. продаж в надежде на падение цены, — это индикатор рыночного сентимента, представляющий собой число акций (долей собственности), которые инвесторы продали, но еще не покрыли. Он используется для определения того, надеются ли инвесторы на понижение (медвежий рынок) или повышение цены (бычий рынок) и иногда применяется в качестве встречного индикатора. См. <https://www.investopedia.com/terms/s/shortinterest.asp>. — Прим. перев.

Стоимостные факторы

Акции с низкими ценами относительно их фундаментальной стоимости тяготеют к обеспечению финансовых возвратов, превышающих эталон, взвешенный по капитализации. Стоимостные факторы отражают эту корреляцию и предназначены для подачи сигналов на покупку недооцененных активов, т. е. относительно дешевых, и продажу переоцененных и дорогих. По этой причине в основе любой стоимостной стратегии лежит модель оценивания стоимости, которая оценивает или служит индикатором справедливой или фундаментальной стоимости актива. *Справедливая стоимость* может быть определена как абсолютный ценовой уровень, спред по отношению к другим активам или диапазон, в котором актив должен торговаться (например, два стандартных отклонения).

Стоимостные стратегии опираются на разворот цен к справедливой среднеуровневой стоимости актива. Они исходят из того, что цены отклоняются от справедливой стоимости лишь временно из-за поведенческих эффектов, таких как чрезмерная реакция или стадный эффект, или эффектов ликвидности, таких как временное влияние на финансовый рынок или долговременные трения между спросом и предложением. Поскольку стоимостные факторы опираются на разворот к среднему уровню, они часто проявляют свойства, противоположные свойствам импульсных факторов. Применительно к долевым ценным бумагам противоположными стоимостным акциям являются ростовые акции с высокой оценкой рыночной стоимости из-за ожиданий роста.

Стоимостные факторы позволяют осуществлять широкий массив систематических стратегий, включая фундаментальное и рыночное оценивание стоимости, статистический арбитраж и стоимость относительно среза активов. Они часто реализуются в виде длинно-коротких портфелей без влияния других традиционных или альтернативных рисков факторов.

Фундаментально-стоимостные стратегии выводят справедливую стоимость активов из экономических и фундаментальных индикаторов, которые зависят от целевого класса актива. В случае ценных бумаг с фиксированной доходностью, валют и товаров указанные индикаторы включают, например, уровни и баланс счета движения капитала, экономическую активность, инфляцию или движения фондов. В случае долевого ценного бумага и корпоративного кредита стоимостные факторы восходят к анализу ценных бумаг Грэма (Graham) и Додда (Dodd) 1930-х годов, с тех пор ставшему известным благодаря Уоррену Баффету. Подходы на основе стоимости долевой ценной бумаги сравнивают цену акции с фундаментальными показателями, такими как балансовая стоимость, валовая выручка (верхняя строка отчета о прибылях и убытках), заработки в нижней строке (чистый доход после уплаты налогов) или различные показатели движения денежных средств.

Рыночно-стоимостные стратегии используют статистические модели или автоматически обучающиеся модели, выявляя ошибочную оценку стоимости активов из-за неэффективностей в предоставлении ликвидности. Статистический и индексный арбитраж являются яркими примерами, которые улавливают разворот временных

влияний на фондовый рынок в течение коротких временных горизонтов (мы рассмотрим парную торговлю в следующей главе). На более длительных временных горизонтах сделки по рыночной стоимости также задействуют сезонные эффекты в долевыми ценными бумагами и товарах.

Относительные стоимостные стратегии на срезе активов сосредоточены на относительной ошибочной оценке стоимости различных активов. Например, арбитраж по конвертируемым облигациям предусматривает сделки по относительной стоимости между акцией, кредитом и волатильностью отдельной компании. Относительная стоимость также охватывает сделки между волатильностями кредита и долевыми ценными бумагами, используя кредитные сигналы для торговли долевыми ценными бумагами, либо сделки между товарами и долевыми ценными бумагами.

Обоснование

Существуют как рациональные, так и поведенческие объяснения существования стоимостного эффекта. Мы приведем несколько ярких примеров из целого ряда исследований с дополнительными справочными материалами, перечисленными в репозитории GitHub.

С точки зрения рациональных и эффективных рынков стоимостная премия компенсирует более высокие реальные или субъективно воспринимаемые риски. Исследователи представили наблюдения о том, что стоимостные фирмы обладают меньшей гибкостью, необходимой для адаптации к неблагоприятным экономическим условиям, чем более экономные (поджарые) и гибкие ростовые компании, или что стоимостные риски связаны с высоким кредитным плечом и более неопределенными будущими заработками. Было также показано, что стоимостные и малокапитализированные портфели были также более чувствительны к макрошокам, чем ростовые и высококапитализированные портфели.

С поведенческой точки зрения стоимостная премия может быть объяснена систематическим смещением из-за неприятия убытка и ментального учета⁷. Инвесторы вполне могут быть менее обеспокоены убытками по активам с сильной недавней результативностью из-за подушек безопасности, предлагаемых предыдущим приростом. Такое неприятие убытка побуждает инвесторов воспринимать акции как менее рискованные, чем раньше, и дисконтировать их будущие денежные потоки по более низкой ставке. И наоборот, слабая недавняя результативность может побудить инвесторов повысить дисконтную ставку актива. Дифференцированные ожидания возвратности приводят к стоимостной премии, т. к. ростовые акции с высокоценовым мультипликатором относительно фундаментальных финансовых показателей хорошо зарекомендовали себя в прошлом, но в перспективе инвесто-

⁷ Ментальный учет (mental accounting) — это экономическая концепция, созданная экономистом Ричардом Талером (Richard H. Thaler, теория рациональных расчетов), которая утверждает, что люди классифицируют личные финансовые средства по-разному и поэтому склонны к иррациональному принятию решений в своем поведении по расходованию и инвестированию финансовых средств.

См. <https://www.investopedia.com/terms/m/mentalaccounting.asp>. — Прим. перев.

рам потребуется более низкий средний возврат из-за их систематически смещенного субъективного восприятия более низких рисков, в то время как обратное верно для стоимостных акций.

Ключевые метрические показатели

Существует большое число индикаторов оценивания стоимости, рассчитываемых на основе фундаментальных данных. Эти факторы могут комбинироваться в качестве входов в автоматически обучающуюся модель оценивания стоимости для предсказания цен. В последующих главах мы увидим примеры того, как некоторые из этих факторов используются на практике (табл. 4.3).

Таблица 4.3. Индикаторы оценивания стоимости

| Фактор | Описание |
|---|---|
| Отдача от потока денежных средств (cash flow yield) | Указанный коэффициент делит операционный денежный поток в расчете на акцию на цену акции. Более высокий коэффициент подразумевает более высокий денежный возврат для акционеров (если он выплачивается с использованием дивидендов или выкупа акций или с прибылью реинвестируется в предприятие) |
| Отдача от потока свободных денежных средств (free cash flow yield) | Указанный коэффициент делит поток свободных денежных средств в расчете на акцию, что отражает сумму денежных средств, доступных для распределения после необходимых расчетов по расходам и инвестициям, на акцию. Более высокая и растущая отдача от потока свободных денежных средств обычно рассматривается как сигнал о повышенной результативности |
| Возвратность денежного потока на инвестированный капитал (cash flow return on invested capital, CFROIC) | Указанный коэффициент служит мерой прибыльности денежного потока предприятия. Он делит операционный денежный поток на инвестированный капитал, определяемый как суммарный долг плюс чистые активы. Более высокая возвратность означает, что предприятие имеет больше денежных средств на заданную сумму инвестированного капитала, генерируя более высокую стоимость для акционеров |
| Денежный поток к суммарным активам (cash flow to total assets) | Указанный коэффициент делит операционный денежный поток на суммарные активы и показывает, сколько денежных средств предприятие может генерировать относительно своих активов, где подобно коэффициенту CFROIC более высокий коэффициент является желательным |
| Поток свободных денежных средств к стоимости предприятия (free cash flow to enterprise value) | Указанный коэффициент служит мерой потока свободных денежных средств, который предприятие генерирует по отношению к стоимости предприятия, измеряемой как совокупная стоимость собственного капитала и долга |

Таблица 4.3 (окончание)

| Фактор | Описание |
|--|---|
| EBITDA к стоимости предприятия (EBITDA to enterprise value) | Указанный коэффициент служит мерой заработков предприятия до отчислений на уплату процентов, налогов и амортизацию (Earnings before interest, taxes, depreciation and amortization, EBITDA), который является индикатором денежного потока относительно стоимости предприятия |
| Отдача от заработков (за 1 предшествующий год) (earnings yield (1 year trailing)) | Указанный коэффициент делит сумму заработков за последние 12 месяцев на последнюю рыночную цену (цену закрытия) |
| Отдача от заработков (за 1 последующий год) (earnings yield (1 year forward)) | Вместо фактических исторических заработков данный коэффициент делит скользящую 12-месячную прогнозную консенсусную аналитическую оценку заработков на последнюю цену, где консенсус состоит из (возможно, взвешенного) среднего значения прогнозов |
| Отношение "цена/зárботки к росту" (PEG ratio) | Отношение "цена/зárботки к росту" (Price/Earnings to Growth, PEG) делит отношение цены акции к зárботкам (P/E) на уровень роста заработков за указанный период. Данный коэффициент корректирует цену, уплаченную за доллар заработанных денег (измеряемым коэффициентом P/E) на рост заработков предприятия |
| Прогнозное соотношение "цена/заработки на 1 год вперед относительно сектора" (P/E 1 year forward (FY1) relative to sector) | Прогнозный коэффициент "цена/заработки" (P/E) по отношению к соответствующему секторальному P/E. Он призван устранять секторальное смещение обобщенного коэффициента P/E за счет секторальных различий в оценке стоимости |
| Отдача от продаж (sales yield) | Указанный коэффициент служит мерой оценки стоимости акций относительно их способности генерировать выручку. При прочих равных, акции с более высокими историческими соотношениями продаж к цене демонстрируют ожидаемо повышенную результативность |
| Отдача от продаж на 1 год вперед (sales yield FY1) | Коэффициент форвардных продаж к цене использует аналитический прогноз продаж в сочетании со (взвешенным) средним значением |
| Отдача от балансовой стоимости (book value yield) | Указанный коэффициент делит историческую балансовую стоимость на цену акции |
| Дивидендная отдача (dividend yield) | Текущий годовой дивиденд, деленный на последнюю цену закрытия. Оценка стоимости дисконтированного потока денежных средств исходит из того, что рыночная стоимость предприятия равна текущей стоимости ее будущих потоков денежных средств |

Волатильностные и размерные факторы

Низковолатильностный фактор улавливает избыточные финансовые возвраты от акций, чья волатильность, бета или идиосинкратический риск ниже среднего. Акции с более крупной рыночной капитализацией тяготеют к обладанию более низкой волатильностью, вследствие чего традиционный размерный фактор часто сочетается с более недавним волатильностным фактором.

Низковолатильностная аномалия представляет собой эмпирическую головоломку, которая противоречит основным принципам финансов. *Модель ценообразования капитальных активов* (Capital Asset Pricing Model, CAPM) и другие модели оценивания стоимости активов утверждают, что более высокий риск должен приносить более высокие возвраты, однако на многих рынках и в течение продолжительных периодов наблюдалось обратное, когда менее рискованные активы превосходили по результативности своих более рискованных сверстников.

Обоснование

Низковолатильностная аномалия противоречит гипотезе об эффективных рынках и допущениям CAPM-модели. Вместо этого было предложено несколько поведенческих объяснений.

Лотерейный эффект строится на эмпирическом наблюдении о том, что люди принимают ставки, которые напоминают лотерейные билеты с малым ожидаемым проигрышем, но крупным потенциальным выигрышем, даже если этот крупный выигрыш имеет довольно низкую вероятность. Если инвесторы воспринимают профиль компромисса между риском и возвратностью низкоценовой волатильной акции как похожий на лотерейный билет, то эта ставка может быть привлекательной. В результате инвесторы могут переплачивать за акции с высокой волатильностью и недоплачивать за акции с низкой волатильностью из-за своих систематически смещенных предпочтений. Систематическое смещение из-за репрезентативности⁸ предполагает, что инвесторы экстраполируют успех нескольких хорошо разрекламированных волатильных акций на все волатильные акции, игнорируя спекулятивный характер таких акций.

Инвесторы также могут быть слишком уверены в своей способности прогнозировать будущее, и различия в их мнениях увеличиваются для волатильных акций с более неопределенными результатами. Поскольку легче выражать позитивный взгляд путем занятия длинной позиции, лонгуя, т. е. владея активом, чем негативный взгляд путем занятия короткой позиции, шорта, оптимисты могут превзойти

⁸ Систематическое смещение из-за репрезентативности (representativeness bias), или систематическая ошибка из-за стереотипизации, означает, что менеджеры оценивают вероятность события на основе его близости с другими событиями. Горизонтальное смещение означает, что люди склонны классифицировать предмет по его аналогам и прогнозировать его в будущем в соответствии с его сходством. Вертикальное смещение подразумевает, что на финансовых рынках люди легко склонны судить или прогнозировать акции в соответствии с их собственной предысторией.

См. <https://www.hindawi.com/journals/mpe/2014/686201/>. — Прим. перев.

пессимистов по численности и продолжать подстегивать цену волатильных акций, что в итоге приводит к более низким возвратам.

Более того, инвесторы ведут себя по-разному на бычьих рынках и во время кризисов. Во время бычьих рынков дисперсия бета-факторов намного ниже, вследствие чего низковолатильностные акции не сильно уступают, если вообще это происходит, тогда как во время кризисов инвесторы стремятся владеть или придерживаются низковолатильностные акции, и дисперсия бета-факторов увеличивается. В результате низковолатильностные активы и портфели показывают более высокую результативность в долгосрочной перспективе.

Ключевые метрические показатели

Метрики, используемые для идентификации низковолатильностных акций, охватывают широкий спектр с реализованной волатильностью (ее стандартным отклонением) на одном конце и прогнозной (вмененной) волатильностью и корреляциями — на другом. Некоторые операционализируют низкую волатильность как низкую бета. Подтверждающие данные в пользу волатильностной аномалии выглядят робастными для разных метрик.

Качественные факторы

Качественный фактор призван улавливать избыточный финансовый возврат от компаний, которые являются высокоприбыльными, операционно эффективными, безопасными, стабильными и хорошо управляемыми, короче говоря, высококачественными, в сопоставлении с рынком. Рынки также, судя по всему, вознаграждают относительную определенность заработков и наказывают акции с высокой волатильностью заработков. Крен портфеля в сторону предприятий с высоким качеством давно пропагандируется отборщиками акций, которые опираются на фундаментальный анализ, но являются относительно новым явлением в количественных инвестициях. Главная трудность заключается в том, как определять качественный фактор непротиворечиво и объективно, используя количественные индикаторы в условиях субъективной природы качества.

Стратегии, основанные на автономных качественных факторах, тяготеют к тому, чтобы показывать результативность противочиклически, поскольку инвесторы платят премию с целью минимизации понижательных рисков и подстегивания оценки рыночной стоимости. По этой причине качественные факторы часто сочетаются с другими рисковыми факторами в многофакторной стратегии, чаще всего со стоимостью для порождения качества в разумной ценовой стратегии. Длинно-короткие качественные факторы тяготеют к обладанию отрицательной рыночной бетой, потому что они представляют собой длинные качественные акции, которые также имеют низкую волатильность, и короткие более волатильные низкокачественные акции. Отсюда, качественные факторы часто положительно коррелируют с низковолатильностными и импульсными факторами и отрицательно коррелируют с влиянием стоимости и широкого рынка.

Обоснование

Качественные факторы могут сигнализировать о превосходящей результативности, поскольку фундаментальные финансовые показатели, такие как устойчивая прибыльность, устойчивый рост в денежном потоке, разумное использование кредитного плеча, низкая потребность в финансировании на рынке капитала или низкий финансовый риск, лежат в основе спроса на доли собственности в частных компаниях и поддерживают цену таких компаний в долгосрочной перспективе. С точки зрения корпоративных финансов, качественная компания часто аккуратно управляет своим капиталом и снижает риск чрезмерного использования кредитного плеча или чрезмерной капитализации.

Поведенческое объяснение говорит о том, что инвесторы недостаточно реагируют на информацию о качестве, аналогично обоснованию импульса, когда инвесторы гонятся за победителями и продают проигравших. Еще один аргумент в пользу премий за качество — это стадный аргумент, подобный аргументу ростовых акций. Фондовым менеджерам вполне бывает легче оправдывать покупку компании ее сильными фундаментальными финансовыми показателями, даже когда она становится дорогой акцией, а не более волатильной (рискованной) стоимостной акцией.

Ключевые метрические показатели

Качественные факторы опираются на метрики, вычисляемые из балансового отчета и отчета о прибылях и убытках, которые указывают на прибыльность, отраженную в высокой прибыли или остатках потока денежных средств, операционной эффективности, финансовой прочности и конкурентоспособности в более широком смысле, поскольку она подразумевает способность поддерживать позицию прибыльности с течением времени.

Таким образом, качество было измерено с использованием валовой прибыльности (которая недавно была добавлена в факторную модель Фама — Френча (*см. главу 7*), возврата на инвестированный капитал, низкой волатильности корпоративных заработков или комбинации различных метрик прибыльности, качества заработков и кредитного плеча, с некоторыми вариантами, перечисленными в табл. 4.4.

Менеджмент корпоративных заработков осуществляется главным образом путем маневрирования накоплений⁹. Таким образом, размер накоплений часто используется в качестве индикатора качества корпоративных заработков: более высокие общие накопления по отношению к активам делают более вероятным низкое качество заработков. Однако это не является однозначным, поскольку накопления могут отражать маневрирование заработка, а также бухгалтерские оценки будущего роста предприятия.

⁹ Накопления (accruals) — это заработанные доходы и понесенные расходы, которые оказывают общее влияние на отчет о прибылях и убытках. — *Прим. перев.*

Таблица 4.4. Качественные индикаторы

| Фактор | Описание |
|--|---|
| Оборачиваемость активов (asset turnover) | Указанный фактор служит мерой эффективности использования предприятием своих активов, которые требуют капитала для получения выручки, и рассчитывается путем деления продаж на суммарные активы; более высокая оборачиваемость лучше |
| 12-месячное изменение оборачиваемости активов (asset turnover 12 month change) | Указанный фактор служит мерой изменения эффективности менеджмента в использовании активов для получения выручки за последний год. Акции с высоким уровнем эффективности, как правило, показывают опережающую результативность |
| Текущий коэффициент (current ratio) | Текущий коэффициент является показателем ликвидности, который служит мерой способности предприятия выплачивать краткосрочные обязательства. Он сравнивает текущие активы предприятия с его текущими обязательствами, и более высокий текущий коэффициент лучше с точки зрения качества |
| Процентное покрытие (interest coverage) | Указанный фактор служит мерой того, насколько легко предприятие сможет выплачивать проценты по своему долгу. Он рассчитывается путем деления EBIT предприятия (заруботков до уплаты процентов и налогов) на его процентные расходы. Желательнее более высокий коэффициент |
| Кредитное плечо (leverage) | Предприятие со значительно большим долгом, чем собственный (акционерный) капитал, рассматривается как имеющее завышенное кредитное плечо. Отношение долга к собственному капиталу обычно является обратно пропорциональным перспективам, при этом чем меньше кредитное плечо, тем лучше |
| Коэффициент выплаты (payout ratio) | Сумма корпоративных заработков, выплачиваемых в виде дивидендов акционерам. Акции с более высокими коэффициентами выплат размещаются в верхнем дециле, а акции с более низкими коэффициентами выплат — в нижнем дециле |
| Возвратность собственного капитала (return on equity, ROE) | Ранжирует акции на основе их исторического показателя возвратности (или рентабельности) собственного капитала и размещая в верхнем дециле те, у которых самая высокая возвратность ROE |

Как преобразовывать данные в факторы

Исходя из концептуального понимания категорий ключевых факторов, их обоснования и популярных метрических показателей, ключевой задачей является выявление новых факторов, которые могут лучше отражать риски, воплощенные ранее изложенными движущими силами возвратности, или же поиск новых. В любом случае для определения инкрементного прироста в сигнале важно сравнивать эффективность инновационных факторов с эффективностью известных факторов.

Полезные методы библиотек pandas и NumPy

Библиотеки NumPy и pandas являются ключевыми инструментами в вычислениях кастомизированных факторов. Блокнот `feature_engineering.ipynb` в каталоге `00-data` содержит примеры создания различных факторов. В указанном блокноте используются данные, сгенерированные блокнотом `create_datasets.ipynb` из папки `data` в корневом каталоге репозитория GitHub. Они хранятся в формате HDF5 для более быстрого доступа. Сравнение форматов хранения `parquet`, `HDF5` и `CSV` для кадров данных библиотеки pandas см. в блокноте `storage_benchmarks.ipynb` в каталоге для главы 2 репозитория GitHub.

Ниже приведены некоторые ключевые шаги в вычислении отобранных факторов из сырых данных акций. Дополнительные подробности и визуализации, которые мы здесь опустили, чтобы сэкономить пространство, см. в блокноте `feature_engineering.ipynb`.

Загрузка данных

Мы загружаем наборы данных Quandl с ценами акций, охватывающие рынки доле-вых ценных бумаг компаний США за 2000–2018 гг., используя `pd.IndexSlice` для выполнения операции среза на мультииндексе `pd.MultiIndex`, отбираем скорректированную цену закрытия и снимаем наложение столбцов по вертикали (`unstack`) на наложение по горизонтали для конвертации кадра данных в широкий формат с тикерами в столбцах и временными метками в строках:

```
idx = pd.IndexSlice
DATA_STORE = '../..data/assets.h5'
with pd.HDFStore(DATA_STORE) as store:
    prices = store['quandl/wiki/prices'].loc[idx['2000':'2018', :],
                                              'adj_close'].unstack('ticker')
prices.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 4706 entries, 2000-01-03 to 2018-03-27
Columns: 2286 entries, A to ZUMZ
dtypes: float64(2286)
memory usage: 82.1 MB
```

Передискретизация с дневной частоты на месячную

В целях сокращения времени на тренировку и для экспериментирования со стратегиями на более длинных временных горизонтах мы конвертируем частоту данных из операционной ежедневной частоты в частоту конца месяца, используя имеющуюся скорректированную цену закрытия:

```
monthly_prices = prices.resample('M').last()
```

Вычисление импульсных факторов

Для того чтобы уловить динамику временного ряда, которая захватывает, например, импульсные регулярности, мы вычисляем исторические финансовые возвраты с помощью метода `pct_change(n_periods)`¹⁰, т. е. возвраты за различные месячные периоды, как указано в списке временных сдвигов `lags`. Затем мы конвертируем широкоформатный результат обратно в длинный формат с помощью метода `stack()`, используем конвейер `pipe()` для применения метода `clip()` к результирующему кадру данных и винсоризируем¹¹ возвраты на уровнях [1%; 99%]; т. е. подрезаем выбросы в этих процентилях.

Наконец, мы нормализуем возвраты с помощью геометрического среднего. После использования метода `swaplevel()` для изменения порядка следования уровней `MultiIndex` мы получаем компаундизированные ежемесячные возвраты за шесть периодов от 1 до 12 месяцев:

```
outlier_cutoff = 0.01
data = pd.DataFrame()
lags = [1, 2, 3, 6, 9, 12]
for lag in lags:
    data[f'return_{lag}m'] = (monthly_prices
                             .pct_change(lag)
                             .stack()
                             .pipe(lambda x:
                                   x.clip(lower=x.quantile(outlier_cutoff),
                                       upper=x.quantile(1-outlier_cutoff)))
                             .add(1)
                             .pow(1/lag)
                             .sub(1)
                             )
data = data.swaplevel().dropna()
data.info()

<class 'pandas.core.frame.DataFrame'>
MultiIndex: 381216 entries, (A, 2001-01-31 00:00:00) to (ZUMZ, 2018-03-31 00:00:00)
Data columns (total 6 columns):
return_1m      381216 non-null float64
return_2m      381216 non-null float64
return_3m      381216 non-null float64
return_6m      381216 non-null float64
```

¹⁰ Функция `pct_change` вычисляется по формуле $(r_{\text{текущ}} - r_{\text{пред}})/r_{\text{пред}}$, где r — это финансовый возврат. — Прим. перев.

¹¹ Винсоризация (winsorizing) — это преобразование статистического показателя путем лимитирования экстремальных значений в статистических данных с целью уменьшения эффекта возможно многих выбросов. Данный метод назван в честь инженера-биоистатистика Чарльза П. Уинзора (1895–1951). Эффект тот же, что и при амплитудном ограничении в обработке сигналов. — Прим. перев.

```
return_9m      381216 non-null float64
return_12m     381216 non-null float64
dtypes: float64(6)
memory usage: 18.9+ MB
```

Эти результаты можно использовать для вычисления импульсных факторов, основываясь на разнице между возвратами за более длительные периоды и самым недавним месячным возвратом, а также для разности между 3- и 12-месячными возвратами следующим образом:

```
for lag in [2,3,6,9,12]:
    data[f'momentum_{lag}'] = data[f'return_{lag}m'].sub(data.return_1m)
data[f'momentum_3_12'] = data[f'return_12m'].sub(data.return_3m)
```

Использование сдвинутых во времени финансовых возвратов и разных периодов владения

Для использования сдвинутых во времени значений в качестве входных переменных или признаков¹², связанных с текущими наблюдениями, мы применим метод `shift()`, который сдвигает исторические возвраты вплоть до текущего периода:

```
for t in range(1, 7):
    data[f'return_1m_{t}'] = data.groupby(level='ticker').return_1m.shift(t)
```

Для вычисления финансовых возвратов на разных периодах владения мы схожим образом берем возвраты, вычисленные ранее за нормализованный период, и сдвигаем их назад, выравнивая их с текущими финансовыми признаками:

```
for t in [1,2,3,6,12]:
    data[f'target_{t}m'] = data.groupby(level='ticker')[f'return_{t}m'].shift(-t)
```

Вычисление факторных бета

Мы представим данные Фама — Френча для оценивания влияния на активы общих рисков факторов с использованием линейной регрессии в *главе 8*. Эмпирически было показано, что пять факторов Фама — Френча, а именно рыночный риск, размер, стоимость, операционная прибыльность и инвестиция, объясняют возвраты от финансовых активов, причем они обычно используются для анализа риск-возвратного профиля портфелей. Поэтому вполне естественно включать прошлые влияния факторов в качестве финансовых признаков в модели, призванные предсказывать будущие возвраты.

Мы можем получить доступ к историческим факторным возвратам с помощью библиотеки `pandas-datareader` и оценить исторические влияния факторов с по-

¹² Признак (feature) — индивидуальное измеряемое свойство или характеристика наблюдаемого явления. В матричной форме признак обычно представляется столбцом матрицы. Строки матрицы образуют наблюдения (также именуемые образцами, экземплярами или примерами). В данном контексте выборками, как правило, называется подмножество строк. — *Прим. перев.*

мощью линейно-регрессионного функционала `PandasRollingOLS` в библиотеке `rufinance` следующим образом:

```
factors = ['Mkt-RF', 'SMB', 'HML', 'RMW', 'CMA']
factor_data = web.DataReader('F-F_Research_Data_5_Factors_2x3',
                             'famafrench', start='2000')[0].drop('RF', axis=1)
factor_data.index = factor_data.index.to_timestamp()
factor_data = factor_data.resample('M').last().div(100)
factor_data.index.name = 'date'
factor_data.info()

T = 24
betas = (factor_data
         .groupby(level='ticker', group_keys=False)
         .apply(lambda x: PandasRollingOLS(window=min(T, x.shape[0]-1),
                                           y=x.return_1m,
                                           x=x.drop('return_1m',
axis=1)).beta))
```

Факторная модель Фама — Френча и линейная регрессия будут рассмотрены более подробно в *главе 7*. Дополнительные примеры см. в блокноте.

Факторы, встроенные в платформу Quantopian

Прилагаемый блокнот `factor_library.ipynb` содержит многочисленные примеры факторов, которые предоставляются платформой `Quantopian` или вычисляются из источников данных, доступных с помощью исследовательского API из блокнота.

Встроенные в платформу факторы можно использовать в сочетании с количественными библиотеками Python, в частности NumPy и pandas, для получения более сложных факторов из широкого спектра соответствующих источников данных, таких как цены долевых ценных бумаг компаний США, фундаментальные финансовые показатели от инвестиционно-аналитической компании Morningstar (<http://www.morningstar.com/>) и настроения инвесторов.

Например, коэффициент цены к продажам, обратный приведенной выше отдаче от продаж, имеется в рамках набора данных о фундаментальных финансовых показателях от компании Morningstar. Его можно использовать как часть конвейера, который будет описан подробнее, когда мы введем библиотеку `zipline`.

Библиотека TA-Lib

Библиотека `TA-Lib` охватывает многочисленные технические факторы. Ее реализация на языке Python доступна для локального использования, например, совместно с библиотеками `zipline` и `alphalens`, а также на платформе `Quantopian`. Прилагаемый к книге блокнот также иллюстрирует несколько технических индикаторов, доступных с помощью `TA-Lib`.

Поиски сигналов — как использовать библиотеку zipline

Исторически альфа-факторы для выявления сигналов на покупку или продажу использовали один вход и простые эвристики, пороги или квантильные отсечки. МО доказало свою эффективность в извлечении сигналов из более разнообразного и значительно более крупного множества входных данных, включая другие альфа-факторы, основанные на анализе исторических регулярностей. По этой причине стратегии алгоритмической торговли сегодня используют большое число альфа-сигналов, многие из которых могут быть слабыми индивидуально, но могут давать надежные предсказания при их сочетании автоматически обучающимся алгоритмом с другими модельными или традиционными факторами.

Библиотека zipline с открытым исходным кодом представляет собой событийную систему бэкестирования, сопровождаемую и используемую в производстве краудсорсинговым квантитативным инвестиционным фондом Quantopian (<https://www.quantopian.com/>) с целью обеспечения разработки алгоритмов и ведения живой торговли. Она автоматизирует реакцию алгоритма на торговые события и предоставляет ему текущие и исторические данные в определенной точке во времени (point-in-time, PIT), что позволяет избегать систематического смещения из-за забегания вперед.

Ее можно использовать в режиме офлайн в сочетании с комплектами данных для исследования и оценивания альфа-факторов. При ее применении на платформе Quantopian вы получаете доступ к более широкому набору фундаментальных и альтернативных данных. В этой главе мы также продемонстрируем исследовательскую среду Quantopian, а в следующей главе — бэктестовую интегрированную среду разработки (IDE). Исходный код этого раздела находится в подпапке 01_factor_research_evaluation папки этой главы в репозитории GitHub.

После установки и перед исполнением первого алгоритма вам необходимо принять комплект данных, который по умолчанию состоит из сопровождаемых сообществом Quandl данных о ценах на акции, дивидендах и дроблениях для 3 тыс. публично торгуемых компаний США. Для выполнения следующей ниже строки кода, которая сохранит данные в вашей домашней папке по маршруту `~/zipline/data/<комплект>`, вам понадобится ключ API платформы Quandl:

```
$ QUANDL_API_KEY=<ваш_ключ> zipline ingest [-b <комплект>]
```

Архитектура — событийная симуляция торговли

Алгоритм библиотеки zipline будет работать в течение определенного периода после начальной настройки и исполнять свою торговую логику при возникновении конкретных событий. Эти события приводятся в действие частотой торговли, могут быть запланированы алгоритмом и сводятся к вызову определенных методов библиотеки zipline. Алгоритм поддерживает состояние с помощью контекстного сло-

варя `context` и получает полезную информацию с помощью переменной `data`, содержащей текущие и исторические данные с привязкой к *определенной точке во времени* (point-in-time, PIT). Алгоритм возвращает кадр данных, содержащий метрики результативности портфеля при наличии сделок, а также пользовательские метрики, которые могут оказаться полезными для записи, например, значений факторов.

Алгоритм можно исполнить из командной строки, в блокноте Jupyter и с помощью функции `run_algorithm()`.

Алгоритм требует метода `initialize()`, который вызывается один раз при запуске симуляции. Этот метод можно использовать для добавления свойств в словарь `context`, доступный всем другим методам алгоритма, или регистрации конвейеров `pipeline`, которые выполняют более сложную обработку данных, например фильтрацию ценных бумаг на основе, скажем, логики альфа-факторов.

Исполнение алгоритма происходит с помощью необязательных методов, которые запланированы библиотекой `zipline` к исполнению автоматически либо через определенные пользователем интервалы. Метод `before_trading_start()` вызывается ежедневно перед открытием рынка и в первую очередь служит для выявления набора ценных бумаг, которыми алгоритм может торговать в течение дня. Метод `handle_data()` вызывается каждую минуту.

Конвейерный API `Pipeline` обеспечивает определение и вычисление альфа-факторов для среза ценных бумаг из исторических данных. Конвейер задает вычисления, которые генерируют столбцы в таблице со значениями в определенной точке во времени (PIT) для набора ценных бумаг. Он должен быть зарегистрирован с помощью метода `initialize()` и затем может исполняться по автоматическому или пользовательскому расписанию. Данная библиотека содержит многочисленные встроенные вычисления, такие как скользящие средние или полосы Боллинджера, которые можно использовать для быстрого вычисления стандартных факторов, а также позволяет создавать собственные факторы, как мы продемонстрируем далее.

Самое главное, что конвейерный API позволяет проводить исследование альфа-фактора модульно, потому что он отделяет вычисление альфа-фактора от остальной части алгоритма, в том числе включая размещение и выполнение торговых ордеров и ведение служебных учетных операций, касающихся элементов содержимого портфеля, значений и т. д.

Единственный альфа-фактор из рыночных данных

Сначала мы проиллюстрируем рабочий поток библиотеки `zipline` по исследованию альфа-факторов в офлайновой среде. В частности, мы разработаем и испытаем простой фактор разворота к среднему уровню. Данный фактор служит мерой степени, с которой недавняя результативность отклонилась от исторического среднего. Краткосрочный разворот — это общепринятая стратегия, которая пользуется преимуществом слабой предсказательной регулярности, состоящей в том, что рост цены на акцию, скорее всего, развернется назад к среднему уровню на горизонтах

от менее минуты до одного месяца. Подробности см. в блокноте `single_factor_zipline.ipynb`.

С этой целью указанный фактор вычисляет z-оценку для последнего месячного финансового возврата относительно скользящих ежемесячных возвратов за последний год. На данный момент мы не будем размещать какие-либо ордера, а просто проиллюстрируем реализацию собственного фактора `CustomFactor` и запишем результаты во время симуляции.

После небольших базовых настроек класс `MeanReversion` подклассирует класс `CustomFactor` и определяет метод `compute()`. Он по умолчанию создает входы, `inputs` в виде ежемесячных возвратов в годовом окне, также заданном по умолчанию, вследствие чего переменная `monthly_return` будет иметь 252 строки и один столбец для каждой ценной бумаги в наборе данных `Quandl` в определенный день.

Метод `compute_factors()` создает экземпляр фактора `MeanReversion` и конвейерные столбцы `long`, `short` и `ranking`. Первые два содержат булевы значения, которые могут использоваться для размещения ордеров, а последний отражает совокупный рейтинг для оценивания совокупной результативности фактора. Более того, он использует встроенный фактор `AverageDollarVolume` с целью ограничить вычисление более ликвидными акциями:

```
from zipline.api import (attach_pipeline,
                        date_rules,
                        time_rules,
                        order_target_percent,
                        pipeline_output,
                        record,
                        schedule_function,
                        get_open_orders,
                        calendars
                        )

from zipline.finance import commission, slippage
from zipline.pipeline import Pipeline, CustomFactor
from zipline.pipeline.factors import Returns, AverageDollarVolume
import numpy as np
import pandas as pd

MONTH = 21
YEAR = 12 * MONTH
N_LONGS = N_SHORTS = 25
VOL_SCREEN = 1000

class MeanReversion(CustomFactor):
    """Вычислить коэффициент последнего месячного возврата
    относительно 12-месячного среднего, нормализованный
    стандартным отклонением ежемесячных возвратов"""
```

```

inputs = [Returns(window_length=MONTH)]
window_length = YEAR

def compute(self, today, assets, out, monthly_returns):
    df = pd.DataFrame(monthly_returns)
    # вычислить z-оценку:
    # из последнего периода (-1) вычесть среднее mean() и
    # разделить результат на стандартное отклонение std()
    out[:] = df.iloc[-1].sub(df.mean()).div(df.std())

def compute_factors():
    """Создать факторный конвейер, включая разворот к среднему уровню,
    отфильтрованный по долларовому объему за 30 дней;
    уловить факторные ранги"""

    mean_reversion = MeanReversion()
    dollar_volume = AverageDollarVolume(window_length=30)
    return Pipeline(columns={'longs' : mean_reversion.bottom(N_LONGS),
                              'shorts' : mean_reversion.top(N_SHORTS),
                              'ranking': mean_reversion.rank(ascending=False)},
                    screen=dollar_volume.top(VOL_SCREEN))

```

Полученный результат позволил бы нам размещать длинные и короткие ордера. В следующей главе мы рассмотрим способ построения портфеля путем выбора периода перебалансировки и внесения поправок в элементы содержимого портфеля по мере поступления новых сигналов.

Метод `initialize()` регистрирует конвейер `compute_factors()`, а метод `before_trading_start()` обеспечивает работу конвейера на ежедневной основе. Функция `record()` добавляет столбец `ranking` конвейера, а также текущие цены активов в кадр данных `performance`, возвращаемый функцией `run_algorithm()`:

```

def initialize(context):
    """Настройка: зарегистрировать конвейер,
    запланировать перебалансировку и задать торговые параметры"""

    attach_pipeline(compute_factors(), 'factor_pipeline')

def before_trading_start(context, data):
    """Выполнить факторный конвейер"""

    context.factor_data = pipeline_output('factor_pipeline')
    record(factor_data=context.factor_data.ranking)
    assets = context.factor_data.index
    record(prices=data.current(assets, 'price'))

```


Наконец, зададим объекты `Timestamp` с временными метками начала и конца в терминах времени UTC, установим первичный капитал и запустим `run_algorithm()` со ссылками на ключевые исполняющие методы. Кадр данных `performance` содержит вложенные данные. Например, столбец `prices` состоит из объектов `pd.Series` в каждой ячейке. Следовательно, последующий доступ к данным будет легче, если их сохранить в формате консервации данных `pickle`:

```
start, end = pd.Timestamp('2015-01-01', tz='UTC'), pd.Timestamp('2018-01-01', tz='UTC')
capital_base = 1e7
```

```
performance = run_algorithm(start=start,
                             end=end,
                             initialize=initialize,
                             before_trading_start=before_trading_start,
                             capital_base=capital_base)
```

```
performance.to_pickle('single_factor.pickle')
```

В следующем далее разделе мы будем использовать факторные и ценовые данные, хранящиеся в кадре данных `performance`, для оценивания результативности фактора в различные периоды владения, но сначала рассмотрим способ создания более сложных сигналов, скомбинировав несколько альфа-факторов из разнообразного набора источников данных на платформе Quantopian.

Сочетание факторов из разнообразных источников данных

Исследовательская среда Quantopian приспособлена к ускоренному тестированию предсказательных альфа-факторов. Данный процесс очень похож, потому что он строится поверх библиотеки `zipline`, но предлагает гораздо более богатый доступ к источникам данных. В следующем ниже фрагменте кода показан способ вычисления альфа-факторов не только из рыночных данных, как ранее, но и из фундаментальных и альтернативных данных. Подробности см. в блокноте `multiple_factors_quantopian_research.ipynb`.

Платформа Quantopian предоставляет несколько сотен фундаментальных переменных Morningstar бесплатно, а также включает сигналы `stocktwits` в качестве примера альтернативного источника данных. Кроме того, существуют кастомизированные определения универсумов, такие как универсум торгуемых акций США `QTradableStocksUS`, применяющий несколько фильтров для ограничения бэктестового универсума акциями, которые, скорее всего, будут торговаться в реальных рыночных условиях:

```
from quantopian.research import run_pipeline
from quantopian.pipeline import Pipeline
from quantopian.pipeline.data.builtin import USEquityPricing
```

```

from quantopian.pipeline.data.morningstar import (income_statement,
                                                    operation_ratios,
                                                    balance_sheet)
from quantopian.pipeline.data.psychsignal import stocktwits
from quantopian.pipeline.factors import (CustomFactor,
                                          SimpleMovingAverage,
                                          Returns)
from quantopian.pipeline.filters import QTradableStocksUS

# Периоды в торговых днях
MONTH = 21
QTR = 4 * MONTH
YEAR = 12 * MONTH

```

Мы воспользуемся кастомизированным классом `AggregateFundamentals` для использования последней отчетной точки фундаментальных данных. Это позволит учесть тот факт, что данные фундаментальных финансовых показателей сообщаются ежеквартально, и что платформа Quantopian в настоящее время не обеспечивает простой способ агрегирования исторических данных, скажем, для получения суммы последних четырех кварталов, на скользящей основе:

```

class AggregateFundamentals(CustomFactor):
    def compute(self, today, assets, out, inputs):
        out[:] = inputs[0]

```

Мы снова задействуем собственный фактор `MeanReversion` из предыдущего фрагмента кода. Мы также вычислим несколько других факторов для заданного определения универсума, используя параметр `mask` метода `rank()`:

```

def compute_factors():
    universe = QTradableStocksUS()

    profitability = (AggregateFundamentals(inputs =
                                          [income_statement.gross_profit],
                                          window_length = YEAR) /
                    balance_sheet.total_assets.latest).rank(mask=universe)

    roic = operation_ratios.roic.latest.rank(mask=universe)

    ebitda_yield = (AggregateFundamentals(inputs =
                                          [income_statement.ebitda],
                                          window_length = YEAR) /
                    USEquityPricing.close.latest).rank(mask=universe)

    mean_reversion = MeanReversion().rank(mask=universe)

    price_momentum = Returns(window_length=QTR).rank(mask=universe)

```

```
sentiment = SimpleMovingAverage(inputs=[stocktwits.bull_minus_bear],
                                window_length=5).rank(mask=universe)

factor = profitability + roic + ebitda_yield + mean_reversion +
        price_momentum + sentiment

return Pipeline(
    columns={'Прибыльность': profitability,
            'ROIC': roic,
            'Отдача EBITDA': ebitda_yield,
            'Разворот к среднему (1M)': mean_reversion,
            'Сентимент': sentiment,
            'Ценовой импульс (3M)': price_momentum,
            'Альфа-фактор': factor})
```

Указанный алгоритм использует наивный метод объединения шести отдельных факторов, просто добавляя ранги активов по каждому из этих факторов. При предсказании будущих возвратов вместо равного веса мы хотели бы учитывать относительную важность и дополнительную информацию. Автоматически обучающиеся алгоритмы последующих глав дадут нам возможность делать именно это, используя тот же самый каркас бэктестирования.

Исполнение также опирается на функцию `run_algorithm()`, но возвратный кадр данных на платформе Quantopian содержит только значения факторов, созданные конвейером Pipeline. Это удобно тем, что указанный формат данных может использоваться в качестве входа для библиотеки `alphalens`, служащей для оценивания предсказательной результативности альфа-факторов.

Разделение сигнала и шума — как использовать библиотеку `alphalens`

Платформа Quantopian имеет Python-библиотеку с открытым исходным кодом `alphalens`, созданную для анализа результативности предсказательных акционных факторов, которая хорошо интегрируется с бэктестовой библиотекой `zipline` и библиотекой анализа портфельной результативности и рисков `pyfolio`, которую мы рассмотрим в следующей главе.

Библиотека `alphalens` обеспечивает анализ предсказательной мощности альфа-факторов и касается:

- ◆ корреляции сигналов с последующими финансовыми возвратами;
- ◆ прибыльности равного или факторно-взвешенного портфеля на основе (подмножества) сигналов;
- ◆ оборачиваемости факторов, указывающих на потенциальные торговые издержки;
- ◆ факторной результативности во время специфических событий;
- ◆ разбивки всего предыдущего по секторам.

Анализ может быть проведен с использованием свободных таблиц или отдельных вычислений и графиков. Сводные таблицы `tear sheet` проиллюстрированы в онлайн-репозитории с целью экономии пространства.

Создание форвардных финансовых возвратов и факторных квантилей

Для того чтобы задействовать библиотеку `alphalens`, нам нужно предоставить сигналы для универсума активов, подобные тем, которые возвращаются рангами фактора `MeanReversion`, и форвардные финансовые возвраты, зарабатываемые при инвестировании в актив в течение определенного периода владения. Подробности см. в блокноте `03_performance_eval_alphalens.ipynb`.

Мы восстановим цены из файла консервации `single_factor.pickle` (факторные данные `factor_data`) следующим образом:

```
performance = pd.read_pickle('single_factor.pickle')

prices = pd.concat([df.to_frame(d) for d, df in
                    performance.prices.items()], axis=1).T
prices.columns = [re.findall(r"\[(.+)\]", str(col))[0]
                  for col in prices.columns]
prices.index = prices.index.normalize()
prices.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 755 entries, 2015-01-02 to 2017-12-29
Columns: 1661 entries, A to ZTS
dtypes: float64(1661)
```

Блокнот `alpha_factor_evaluation.ipynb` в репозитории `GitHub` содержит более подробную информацию о том, как проводить посекторное оценивание.

Входные данные для библиотеки `alphalens` можно создать в необходимом формате, используя вспомогательную функцию `get_clean_factor_and_forward_returns`, которая также возвращает сигнальные квантили и форвардные финансовые возвраты за заданные периоды владения:

```
HOLDING_PERIODS = (5, 10, 21, 42) # Периоды владения
QUANTILES = 5
alphalens_data = get_clean_factor_and_forward_returns(factor=factor_data,
                                                       prices=prices,
                                                       periods=HOLDING_PERIODS,
                                                       quantiles=QUANTILES)
```

```
Dropped 14.5% entries from factor data: 14.5% in forward returns
computation and 0.0% in binning phase (set max_loss=0 to see potentially
suppressed Exceptions). max_loss is 35.0%, not exceeded: OK!
```

Кадр данных `alphalens_data` содержит финансовые возвраты на инвестиции в определенный актив на заданную дату для указанного периода владения, а также значение фактора, т. е. рейтинг разворота к среднему уровню `MeanReversion` актива на эту дату, и соответствующее квантильное значение (табл. 4.5).

Таблица 4.5. Финансовые возвраты на инвестиции в определенный актив на заданную дату для указанного периода владения

| date | asset | 5D | 10D | 21D | 42D | factor | factor_quantile |
|----------|-------|--------|--------|---------|---------|--------|-----------------|
| 01/02/15 | A | 0,07% | −5,70% | −2,32% | 4,09% | 2618 | 4 |
| | AAL | −3,51% | −7,61% | −11,89% | −10,23% | 1088 | 2 |
| | AAP | 1,10% | −5,40% | −0,94% | −3,81% | 791 | 1 |
| | AAPL | 2,45% | −3,05% | 8,52% | 15,62% | 2917 | 5 |
| | ABBV | −0,17% | −2,05% | −6,43% | −13,70% | 2952 | 5 |

Форвардные возвраты и сигнальные квантили являются основой для оценивания предсказательной мощности сигнала. Как правило, фактор должен давать заметно разные возвраты для несовпадающих квантилей, таких как отрицательные возвраты для нижнего квинтиля и положительные возвраты для верхнего квантиля факторных значений.

Предсказательная результативность по факторным квантилям

В качестве первого шага мы хотели бы визуализировать среднепериодный финансовый возврат по факторному квантилю. Мы можем применить встроенную функцию `mean_return_by_quantile` из модуля `performance` и функцию `plot_quantile_returns_bar` из модуля `plotting`:

```
from alphalens.performance import mean_return_by_quantile
from alphalens.plotting import plot_quantile_returns_bar

mean_return_by_q, std_err = mean_return_by_quantile(alphalens_data)
plot_quantile_returns_bar(mean_return_by_q);
```

Результатом является гистограмма, которая разбивает среднее значение форвардных финансовых возвратов для четырех разных периодов владения, основываясь на квинтиле факторного сигнала. Как вы видите, нижние квинтили дали заметно больше отрицательных результатов, чем верхние квинтили, за исключением самого длительного периода владения (рис. 4.2).

10-дневный период владения (10D) дает несколько более высокие результаты для первого и четвертого квантилей. Мы также хотели бы посмотреть на результативность за время инвестиций, стимулируемую каждым сигнальным квинтилем.

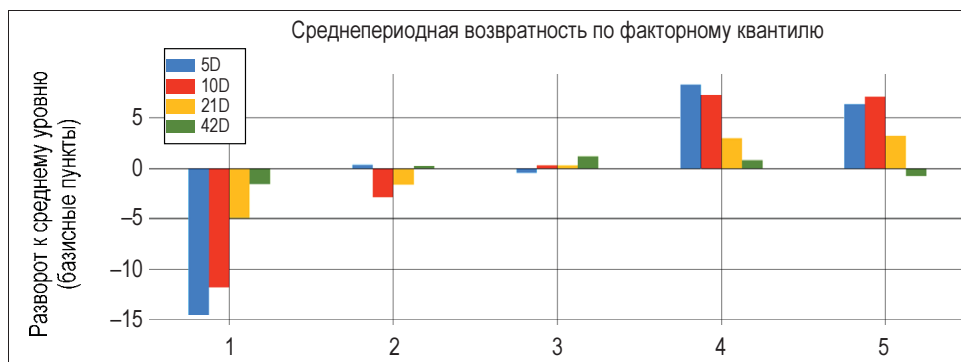


Рис. 4.2. Гистограмма среднего значения форвардных финансовых возвратов для четырех разных периодов владения

Мы вычислим ежедневные финансовые возвраты, в отличие от средних возвратов за 5-дневный период владения (5D), и библиотека `alphalens` скорректирует периодные возвраты на несоответствие между ежедневными сигналами и более длительным периодом владения (подробности см. в документации):

```
from alphalens.plotting import plot_cumulative_returns_by_quantile
```

```
mean_return_by_q_daily, std_err = mean_return_by_quantile(alphalens_data, by_date=True)
plot_cumulative_returns_by_quantile(mean_return_by_q_daily['5D'], period='5D');
```

Результирующий линейный график показывает, что в течение большей части этого трехлетнего периода два верхних квантиля значительно превосходили два нижних квантиля. Однако, как было предложено предыдущим графиком, сигналы четвертого квантиля произвели более высокую результативность, чем сигналы верхнего квантиля (рис. 4.3).

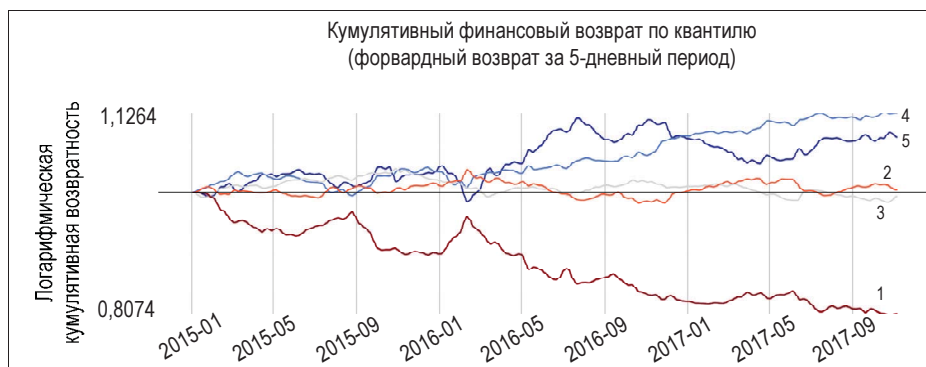


Рис. 4.3. Кумулятивный финансовый возврат по квантилю

Полезный для торговой стратегии фактор показывает приведенную выше регулярность, где кумулятивные финансовые возвраты развиваются по четко различающимся траекториям, поскольку он допускает длинно-короткую стратегию с более

низкими требованиями к капиталу и, соответственно, меньшим влиянием совокупного рынка.

Однако нам также необходимо учитывать не только средние значения, но и дисперсию периодных возвратов. С этой целью мы можем опереться на встроенную функцию `plot_quantile_returns_violin`:

```
from alphas.plotting import plot_quantile_returns_violin
```

```
plot_quantile_returns_violin(mean_return_by_q_daily);
```

Скрипичный график концентрации распределения подчеркивает, что диапазон ежедневных возвратов довольно широк и, несмотря на разные средние значения, выделить распределения можно очень ограниченно, вследствие чего в любой день различия в результативности между разными квантилями могут быть довольно ограниченными (рис. 4.4).

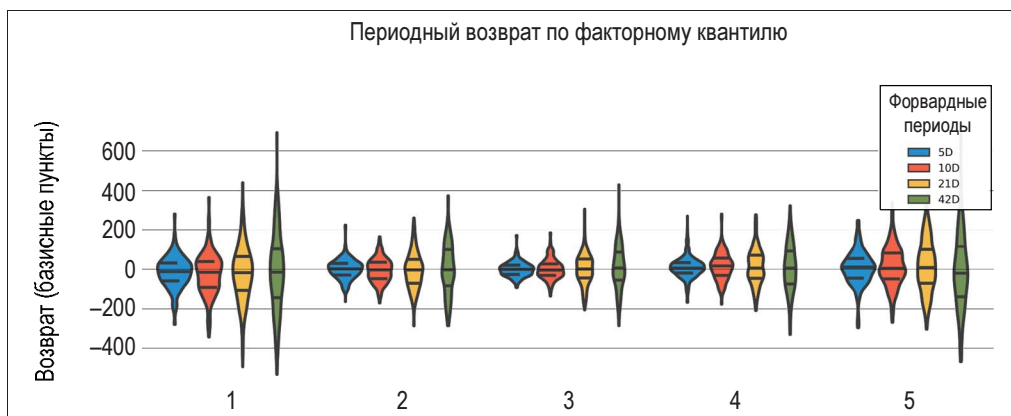


Рис. 4.4. Периодный возврат по факторному квантилю

Когда мы концентрируемся на оценивании одного-единственного альфа-фактора, мы упрощаем вещи, игнорируя практические вопросы, связанные с исполнением сделок, которые мы ослабим, когда обратимся к надлежащему бэктестированию в следующей главе. Некоторые из них включают:

- ◆ транзакционные издержки торговли;
- ◆ проскальзывание¹³, разницу между ценой при принятии решения и исполнением сделки, например, из-за влияния рынка.

¹³ Проскальзывание (slippage) — это разница в цене, по которой брокер получает указание от принципала выполнить ордер, и цене, по которой ордер фактически исполнен.

См. <http://www.businessdictionary.com/definition/slippage.html>. — Прим. перев.

Информационный коэффициент

Подавляющая часть этой книги посвящена разработке альфа-факторов с использованием автоматически обучающихся моделей. МО занимается оптимизацией некоторой предсказательной целевой задачи, и в этом разделе мы представим ключевые метрические показатели, используемые для измерения результативности альфа-фактора. Мы определим альфа как средний финансовый возврат, превышающий эталон.

Такое определение приводит к *информационному соотношению* (information ratio, IR), которым измеряется средний избыточный финансовый возврат на единицу риска, взятому путем деления альфы на ошибку прослеживания. Когда эталоном является безрисковая ставка, информационное соотношение соответствует хорошо известному коэффициенту Шарпа, и мы выделим важнейшие вопросы статистического измерения, которые возникают в типичном случае, когда возвраты не являются нормально распределенными. Мы также дадим объяснение фундаментальному закону активного менеджмента, который разбивает информационное соотношение на сочетание умения прогнозировать и способность стратегии эффективно задействовать умения прогнозировать.

Целью альфа-факторов является точное направленное предсказание будущих финансовых возвратов. Следовательно, естественной результативностной мерой является корреляция между предсказаниями со стороны альфа-фактора и форвардными возвратами от целевых активов, т. е. информационный коэффициент (подробнее см. разд. "Фундаментальный закон активного менеджмента" главы 5).

В этом случае лучше использовать непараметрический коэффициент ранговой корреляции Спирмена, которым измеряется то, насколько хорошо может быть описана связь между двумя переменными, при помощи монотонной функции, в отличие от корреляции Пирсона, которым измеряется сила линейной связи.

Мы получаем информационный коэффициент¹⁴ с помощью библиотеки `alphalens`, которая под капотом опирается на `scipy.stats.spearmanr` (пример того, как использовать библиотеку `SciPy` напрямую для получения *p*-значений, см. в репозитории `GitHub`). Функция `factor_information_coefficient` вычисляет периодическую корреляцию, а функция `plot_ic_ts` создает график временных рядов с одномесячным скользящим средним:

```
from alphalens.performance import factor_information_coefficient
from alphalens.plotting import plot_ic_ts

ic = factor_information_coefficient(alphalens_data)
plot_ic_ts(ic[['5D']])
```

¹⁴ Ошибка прослеживания (tracking error), иногда риск отклонения, — это стандартное отклонение возвратности портфеля от эталонной возвратности за определенный период времени. См. <https://financetrain.com/tracking-error-tracking-risk/>. — Прим. перев.

Этот график временного ряда показывает продолжительные периоды со значительно положительным скользящим средним информационного коэффициента (IC)¹⁵, рис. 4.5. IC, равный 0,05 или даже 0,1, допускает значительно повышенную результативность, если есть достаточные возможности для применения этого умения прогнозировать, как будет иллюстрировать фундаментальный закон активного менеджмента.

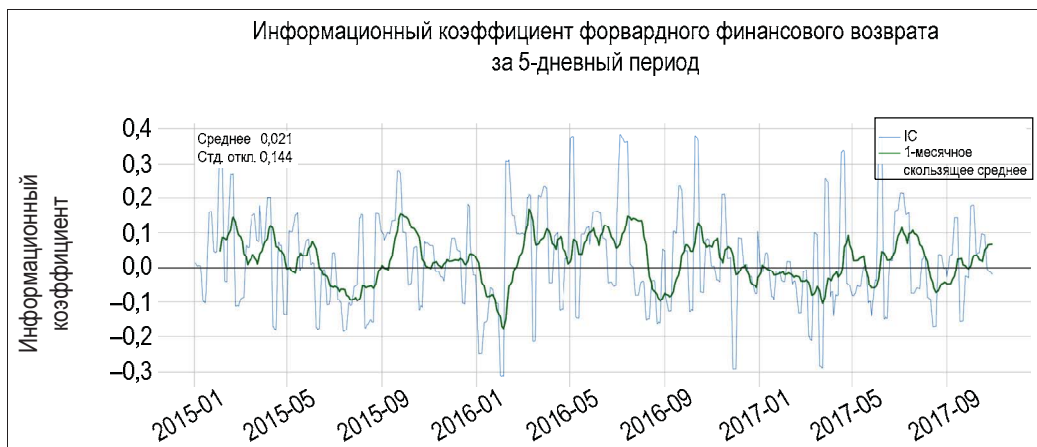


Рис. 4.5. Информационный коэффициент форвардного финансового возврата за 5-дневный период

График среднегодового значения IC показывает, как результативность фактора была исторически неравномерной:

```
ic = factor_information_coefficient(alphalens_data)
ic_by_year = ic.resample('A').mean()
ic_by_year.index = ic_by_year.index.year
ic_by_year.plot.bar(figsize=(14, 6))
```

В результате получится график, представленный на рис. 4.6.

Информационный коэффициент (IC) ниже 0,05, как и в этом случае, является низким, но значительным и может давать положительные остаточные возвраты относительно эталона, как мы увидим в следующем разделе.

Вызов функции `create_summary_tear_sheet(alphalens_data)` создает сводку статистических показателей IC, где скорректированный на риск IC получен путем деления среднего IC на стандартное отклонение IC, который также подвергается двусторонней статистической проверке на основе статистического показателя t (t -теста)

¹⁵ Информационный коэффициент (IC) не следует путать с информационным соотношением (IR). Последнее является мерой мастерства инвестиционного менеджера, сравнивая избыточную доходность менеджера с суммой принятого риска. Информационный коэффициент описывает корреляцию между предсказанным и фактическим возвратом, варьируясь в интервале от -1 до 1 ; иногда используется для измерения вклада финансового аналитика.

См. <https://www.investopedia.com/terms/i/information-coefficient.asp>. — Прим. перев.

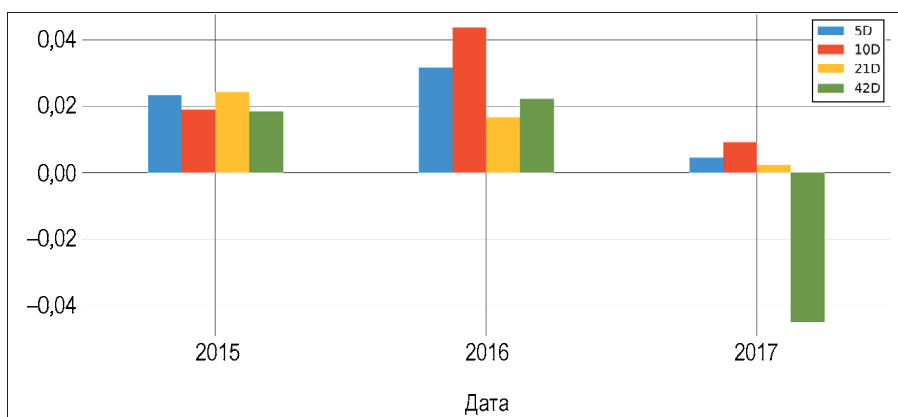


Рис. 4.6. Среднегодовое значение информационного коэффициента

с нулевой гипотезой¹⁶ $IC = 0$, используя для этого функцию `scipy.stats.ttest_1samp` (табл. 4.6).

Таблица 4.6. Сводка статистических показателей информационного коэффициента

| | 5D | 10D | 21D | 42D |
|------------------------------------|------|-------|-------|-------|
| Среднее значение IC | 0,01 | 0,02 | 0,01 | 0,00 |
| Стандартное отклонение IC | 0,14 | 0,13 | 0,12 | 0,12 |
| IC с поправкой на риск | 0,10 | 0,13 | 0,10 | 0,01 |
| Статистический показатель t (IC) | 2,68 | 3,53 | 2,53 | 0,14 |
| p -значение (IC) | 0,01 | 0,00 | 0,01 | 0,89 |
| Асимметрия IC | 0,41 | 0,22 | 0,19 | 0,21 |
| Экссесс IC | 0,18 | -0,33 | -0,42 | -0,27 |

Оборачиваемость фактора

Оборачиваемость фактора служит мерой частоты изменения активов, ассоциированных с определенным квантилем, т. е. сколько сделок требуется для корректировки портфеля в соответствии с последовательностью сигналов. Более конкретно, она служит мерой доли активов, находящихся в факторном квантиле в настоящее время, которая не была в этом квантиле в последний период. Приведенная далее табл. 4.7 создается этой командой:

```
create_turnover_tear_sheet(alphalens_data)
```

¹⁶ Нулевая гипотеза (null hypothesis) — это общее утверждение или принимаемая по умолчанию позиция относительно того, что между двумя наблюдаемыми событиями взаимосвязи не существует, или что причиной всему является случайность. — *Прим. перев.*

Доля активов, которые должны были присоединиться к квантильно-ориентированному портфелю, довольно высока, что говорит о том, что торговые издержки создают проблему для получения выгод от предсказательной результативности.

Таблица 4.7. Средняя оборачиваемость фактора

| Средняя оборачиваемость | 5D | 10D | 21D | 42D |
|-------------------------|-----|-----|-----|-----|
| Квантиль 1 | 59% | 83% | 83% | 41% |
| Квантиль 2 | 74% | 80% | 81% | 65% |
| Квантиль 3 | 76% | 80% | 81% | 68% |
| Квантиль 4 | 74% | 81% | 81% | 64% |
| Квантиль 5 | 57% | 81% | 81% | 39% |

Альтернативным взглядом на оборачиваемость фактора является корреляция ранга актива вследствие фактора за различные периоды владения. Приведем часть сводной таблицы (табл. 4.8).

Таблица 4.8. Корреляция ранга актива

| | 5D | 10D | 21D | 42D |
|--|-------|-------|--------|--------|
| Автокорреляция среднего факторного ранга | 0,711 | 0,452 | −0,031 | −0,013 |

Как правило, для удержания торговых издержек на приемлемом уровне предпочтительна более высокая стабильность.

Ресурсы альфа-факторов

Процесс исследования требует конструирования и отбора альфа-факторов по отношению к предсказательной мощности их сигналов. Стратегия алгоритмической торговли обычно строится на нескольких альфа-факторах, которые посылают сигналы для каждого актива. Как мы увидим в последующих главах, автоматически обучающаяся модель может агрегировать эти факторы с целью оптимизации того, как разнообразные сигналы транслируются в решения о выборе правильного времени и размеров отдельных позиций.

Альтернативные библиотеки алгоритмической торговли

Дополнительные Python-библиотеки с открытым исходным кодом для алгоритмической торговли и сбора данных включают перечисленные далее средства (справочные материалы см. в репозитории GitHub):

- ◆ бэктестовый инструмент QuantConnect (<https://www.quantconnect.com/>) является конкурентом платформы Quantopian;
- ◆ фирма по квантитативному инвестиционному менеджменту WorldQuant (<https://www.worldquant.com/>) предлагает онлайнovou конкуренцию и набирает соавторов из сообщества в хедж-фонд, питаемый за счет прямого вовлечения людей через Интернет;
- ◆ фирма Alpha Trading Labs (<https://www.alphatradinglabs.com/>) предлагает высокочастотную инфраструктуру тестирования с производственной моделью, аналогичной Quantopian;
- ◆ Python-библиотека для алгоритмической торговли PyAlgoTrade (см. на GitHub) фокусируется на бэктестировании и предлагает поддержку бумажной и живой торговли. Она позволяет оценивать идею торговой стратегии с использованием исторических данных и стремится делать это с минимальными усилиями;
- ◆ векторизованный каркас бэктестирования Pybacktest (см. на GitHub) использует библиотеку pandas и стремится быть компактным, простым и быстрым (данный проект в настоящее время приостановлен);
- ◆ более старый проект ultrafinance (см. на GitHub) объединяет реально-временной сбор финансовых данных, анализ и тестирование стратегий торговли;
- ◆ веб-сайт "Trading with Python" (Торговля с помощью Python) (<http://www.tradingwithpython.com/>) предлагает курсы и коллекцию функций и классов для квантитативной торговли;
- ◆ фирма Interactive Brokers (<https://www.interactivebrokers.co.uk>) предлагает API Python для живой торговли на своей платформе.

Резюме

В этой главе мы рассмотрели использование библиотеки `zipline` для событийного симулирования алгоритма торговли как в режиме офлайн, так и на онлайн-платформе Quantopian. Мы проиллюстрировали конструирование и оценивание индивидуальных альфа-факторов с целью извлечения сигналов для стратегии алгоритмической торговли из рыночных, фундаментальных и альтернативных данных и продемонстрировали наивный способ комбинирования нескольких факторов. Мы также представили библиотеку `alphalens`, которая позволяет всесторонне оценивать предсказательную результативность и торговую оборачиваемость сигналов.

Процесс конструирования портфеля, в свою очередь, имеет более широкую перспективу и направлен на оптимальное определение размеров позиций с точки зрения риска и возвратности. Теперь мы обратимся к всевозможным стратегиям балансировки риска и возвратности в портфельном процессе. Мы также подробнее рассмотрим сложности бэктестирования стратегий торговли на ограниченном наборе исторических данных и способы их преодоления.

5

Оценивание стратегии

Альфа-факторы приводят в действие алгоритмическую стратегию, транслируемую в сделки, которые, в свою очередь, порождают инвестиционный портфель. Финансовые возвраты и риск результирующего портфеля определяют успех стратегии. Тестирование стратегии предусматривает симулирование портфелей, генерируемых алгоритмом, в целях проверки результативности портфеля в рыночных условиях. Оценивание стратегии включает бэктестирование на исторических данных с целью оптимизации параметров стратегии и форвардное тестирование с целью проверки внутривыборочной результативности на новых вневыборочных данных и предотвращения ложных обнаружений, возникающих из-за приспособления стратегии к конкретным прошлым обстоятельствам.

В контексте портфеля положительные возвраты от активов могут компенсировать отрицательные ценовые движения нелинейным образом, вследствие чего совокупная вариация портфельных возвратов будет меньше, чем средневзвешенная вариация портфельных позиций, если только их возвраты идеально и положительно не коррелируют. В 1952 г. Гарри Марковиц разработал теорию управления современным инвестиционным портфелем на основе диверсификации, которая привела к среднedisперсной оптимизации (между средним значением и дисперсией): для заданного набора активов портфельные веса могут быть оптимизированы с целью снижения риска, измеряемого как стандартное отклонение возвратов для заданного ожидаемого уровня возвратов.

Модель ценообразования капитальных активов (CAPM) ввела рисковую премию как равновесное вознаграждение за владение активом, компенсирующее влияние одного-единственного рискового фактора — рынка — который невозможно диверсифицировать. Риск-менеджмент постепенно эволюционировал, становясь гораздо изощреннее по мере появления дополнительных рисковых факторов и более гранулярных вариантов влияния. Правило Келли является популярным подходом к динамической оптимизации портфеля, которое представляет собой выбор последовательности позиций с течением времени; в 1968 г. оно было прекрасно адаптировано на фондовом рынке Эдвардом Торпом (Edward Thorp) после его первоначального применения в азартных играх.

Как результат, существует несколько подходов к портфельной оптимизации, которые включают применение *машинного обучения* (МО) с целью усвоения иерархических взаимосвязей между активами и рассмотрения элементов их содержимого как дополнений или заменителей по отношению к рисковому профилю портфеля.

В этой главе мы рассмотрим следующие темы:

- ◆ как строить и тестировать инвестиционный портфель на основе альфа-факторов с помощью библиотеки `zipline`;
- ◆ как измерять портфельный риск и возвратность;
- ◆ как оценивать портфельную результативность с помощью библиотеки `pyfolio`;
- ◆ как управлять портфельными весами с помощью среднедисперсной оптимизации и ее альтернатив;
- ◆ как использовать машинное обучение для оптимизации размещения финансовых средств среди активов в контексте портфеля.



Примеры исходного кода данной главы расположены в каталоге `05_strategy_evaluation_and_portfolio_management` сопутствующего репозитория GitHub.

Как строить и тестировать инвестиционный портфель с помощью библиотеки `zipline`

В предыдущей главе мы ввели библиотеку `zipline` для симулирования вычисления альфа-факторов из предшествующих срезовых (кросс-секционных) рыночных, фундаментальных и альтернативных данных. Теперь мы будем использовать альфа-факторы для выведения сигналов и действия по сигналам на покупку и продажу. Мы отложим оптимизацию портфельных весов до конца этой главы, а пока просто назначим каждому элементу содержимого портфеля равнозначные позиции. Исходный код этого раздела находится в подпапке `01_trading_zipline`.

Торговля по плану и перебалансировка портфеля

Мы будем использовать собственной фактор разворота к среднему `MeanReversion`, разработанный в предыдущей главе (его реализацию см. в сценарии `alpha_factor_zipline_with_trades.py`).

Конвейер `Pipeline`, созданный функцией `compute_factors()`, возвращает таблицу со столбцами `long` и `short` для 25 акций с наибольшими отрицательными и положительными отклонениями их последнего месячного возврата от его среднегодового значения, нормализованного стандартным отклонением. Он также ограничил универсум 500 акциями с самым высоким средним объемом торговли за последние

30 торговых дней. Метод `before_trading_start()` обеспечивает ежедневное исполнение конвейера и запись результатов, включая текущие цены.

Новый метод `rebalance()` предъявляет торговые ордера методу `exec_trades()` для активов с длинными и короткими позициями, помеченными конвейером положительным и отрицательным весами. Он также дивестирует (снимает) любые текущие элементы содержимого портфеля, которые больше не включаются в факторные сигналы:

```
def exec_trades(data, assets, target_percent):
    """Разместить ордера на активы, используя
       целевое портфельное процентное соотношение"""

    for asset in assets:
        if data.can_trade(asset) and not get_open_orders(asset):
            order_target_percent(asset, target_percent)

def rebalance(context, data):
    """Вычислить длинные, короткие и устаревшие
       элементы содержимого портфеля;
       разместить торговые ордера"""

    factor_data = context.factor_data
    assets = factor_data.index

    longs = assets[factor_data.longs]
    shorts = assets[factor_data.shorts]
    divest = context.portfolio.positions.keys() - longs.union(shorts)

    exec_trades(data, assets=divest, target_percent=0)
    exec_trades(data, assets=longs,
                 target_percent=1 / N_LONGS if N_LONGS else 0)
    exec_trades(data, assets=shorts,
                 target_percent=-1 / N_SHORTS if N_SHORTS else 0)
```

Метод `rebalance()` работает в соответствии с правилами `date_rules` и `time_rules`, заданными вспомогательной функцией `schedule_function()` в начале недели, сразу после открытия рынка `market_open`, как это предусмотрено встроенным календарем `US_EQUITIES` (подробнее о правилах см. в документации). Также можно указать торговую комиссию как в относительном выражении, так и в виде минимальной суммы. Кроме того, существует возможность определить проскальзывание, т. е. стоимость неблагоприятного изменения цены между решением о сделке и ее исполнением:

```
def initialize(context):
    """Настройка: зарегистрировать конвейер,
       спланировать перебалансировку и задать торговые параметры"""
```

```

attach_pipeline(compute_factors(), 'factor_pipeline')
schedule_function(rebalance,
                  date_rules.week_start(),
                  time_rules.market_open(),
                  calendar=calendars.US_EQUITIES)

set_commission(us_equities=commission.PerShare(cost=0.00075,
                                                min_trade_cost=.01))
set_slippage(us_equities=slippage.VolumeShareSlippage(volume_limit=0.0025,
                                                       price_impact=0.01))

```

Алгоритм продолжает исполняться после вызова функции `run_algorithm()` и возвращает тот же кадр данных бэктестовой результативности, `performance`. Теперь мы обратимся к общим мерам портфельной возвратности и риска, а также тому, как их вычислять с помощью библиотеки `pyfolio`.

Как измерять результативность с помощью библиотеки `pyfolio`

Задачей МО является оптимизация целевых функций. В алгоритмической торговле целевыми функциями являются возврат и риск совокупного инвестиционного портфеля, как правило, по отношению к эталону (который может быть денежным или безрисковой процентной ставкой).

Для оценивания этих целевых функций существует несколько метрических показателей. Мы кратко рассмотрим наиболее часто применяемые метрики и способы их вычисления с помощью библиотеки `pyfolio`, которая также используется библиотекой `zipline` и платформой `Quantopian`. Мы рассмотрим способы применения этих показателей на платформе `Quantopian` при тестировании стратегии алгоритмической торговли.

Мы будем использовать несколько простых обозначений: пусть R равно временному ряду простых однопериодных портфельных возвратов, $R = (r_1, \dots, r_T)$, с даты 1 до даты T , и $R^f = (r_1^f, \dots, r_T^f)$ равно соответствующему временному ряду безрисковых процентных ставок такому, что $R^e = R - R^f = (r_1 - r_1^f, \dots, r_T - r_T^f)$ является избыточным возвратом.

Коэффициент Шарпа

Априорный *коэффициент Шарпа* (Sharpe ratio, SR) сравнивает математическое ожидание избыточного портфельного возврата с волатильностью этого избыточного возврата, измеряемой его стандартным отклонением. Он служит мерой компенсации как среднего избыточного возврата на единицу принятого риска:

$$\mu \equiv E(R_t);$$

$$\sigma_{R^e}^2 \equiv \text{Var}(R - R^f);$$

$$\text{SR} \equiv \frac{\mu - R_f}{\sigma_{R^e}}.$$

Математические ожидания возвратов и волатильностей не наблюдаемы, но могут оцениваться следующим образом с использованием исторических данных:

$$\hat{\mu}_{R^e} = \frac{1}{T} \sum_{t=1}^T r_t^e;$$

$$\hat{\sigma}_{R^e}^2 = \frac{1}{T} \sum_{t=1}^T (r_t^e - \hat{\mu}_{R^e})^2;$$

$$\widehat{\text{SR}} \equiv \frac{\hat{\mu}_{R^e}}{\hat{\sigma}_{R^e}}.$$

Если только безрисковая ставка не является волатильной (как на развивающихся рынках), стандартное отклонение избыточных и сырых возвратов будет похожим. Когда коэффициент Шарпа (SR) используется с эталоном, отличным от безрисковой ставки, например с фондовым индексом S&P 500¹, он называется *информационным соотношением* (information ratio, IR). В этом случае он служит мерой избыточного возврата портфеля, также именуемой *альфой*, относительно ошибки прослеживания, т. е. отклонения портфельных возвратов от эталонных возвратов².

Для *одинаково распределенных и взаимно независимых* (iid) финансовых возвратов математическое выведение распределения правила оценки коэффициента Шарпа для проверок на статистическую значимость следует из применения центральной предельной теоремы к $\hat{\mu}$ и $\hat{\sigma}^2$ в соответствии со статистической теорией больших выборок.

Однако финансовые возвраты нередко нарушают допущения об одинаковой распределенности и взаимной независимости. Эндрю Ло (Andrew Lo) вывел необходимые поправки в это распределение и временную агрегацию для возвратов, кото-

¹ S&P 500 — это фондовый индекс, в корзину которого включено 500 избранных акционерных компаний США, имеющих наибольшую капитализацию. Список принадлежит компании Standard & Poor's и ею же составляется. Индекс публикуется с 4 марта 1957 г.

² Информационное соотношение (information ratio) — это показатель возвратности портфеля, выходящий за пределы возвратности эталонного показателя по сравнению с волатильностью этой возвратности; вычисляется по формуле $(R_p - R_b)/\text{ошибка_прослеживания}$, где R_p — это возвратность инвестиционного портфеля, R_b — эталонная возвратность, *ошибка_прослеживания* — стандартное отклонение избыточной возвратности по отношению к эталонной возвратности. IR пытается определить стабильность результативности портфеля путем включения в расчет риска отклонения, или компоненты стандартного отклонения. См. <https://www.investopedia.com/terms/i/informationratio.asp>. — Прим. перев.

рые являются стационарными, но автокоррелированными. Это важно, потому что свойства временных рядов инвестиционных стратегий (например, среднелазворотных, импульсных и других форм внутрирядовой корреляции) могут оказывать нетривиальное влияние на само правило оценки коэффициента Шарпа, в особенности при пересчете коэффициента Шарпа в годовом исчислении из более высокочастотных данных (Lo 2002).

Фундаментальный закон активного менеджмента

Из высокого *информационного соотношения* (IR) следует привлекательная повышенная результативность относительно дополнительного принятого риска. Фундаментальный закон активного менеджмента разбивает IR на *информационный коэффициент* (information coefficient, IC), как меру умения прогнозировать, и способность применять это умение через независимые ставки³. Он резюмирует важность играть часто (высокая широта) и играть хорошо (высокий IC):

$$IR = IC \cdot \sqrt{\text{широта}} .$$

Информационный коэффициент IC служит мерой корреляции между альфа-фактором и форвардными возвратами, вытекающими из его сигналов, и улавливает точность умения менеджера прогнозировать. Широта стратегии измеряется независимым числом ставок, которые инвестор делает в определенный период времени, а произведение обоих значений пропорционально IR, которое также именуется *аттестационным риском* (appraisal risk, Treynor и Black).

Этот каркас был расширен за счет *передаточного коэффициента* (transfer coefficient, TC) для отражения портфельных ограничений (например, в отношении коротких продаж), которые могут лимитировать информационное соотношение ниже того уровня, который в иных случаях достигим при наличии IC или широты стратегии. TC служит индикатором эффективности, с которой менеджер транслирует идеи в портфельные ставки (Clarke et al., 2002).

Указанный фундаментальный закон важен, потому что он выделяет ключевые движущие силы повышенной результативности: важны как точные предсказания, так и способность делать независимые прогнозы и действовать в соответствии с этими прогнозами. На практике менеджеры с широким набором инвестиционных решений могут достигать значительных скорректированных на риск повышенных возвратов с информационными коэффициентами между 0,05 и 0,15 (по причинам экономики пространства симуляционная диаграмма не включена).

На практике оценивание широты стратегии затруднено с учетом срезовой (крессионной) и рядовой корреляции между прогнозами.

³ Информационный коэффициент (information coefficient) описывает корреляцию между предсказываемой и фактической возвратностью активов, варьируясь в диапазоне -1 до 1; иногда используется для измерения вклада финансового аналитика. См. <https://www.investopedia.com/terms/i/information-coefficient.asp>. — Прим. перев.

Внутривыборочная и вневыборочная результативность с помощью библиотеки `pyfolio`

Библиотека `pyfolio` обеспечивает анализ портфельной результативности и риска в выборке и вне выборки с использованием многочисленных стандартных метрических показателей. Она производит сводные таблицы, охватывающие анализ возвратов, позиций и транзакций, а также событийный риск в периоды рыночного стресса с использованием нескольких встроенных сценариев и, кроме того, включает байесов анализ внутривыборочной результативности.

Она опирается на портфельные возвраты и данные о позициях, а также может учитывать транзакционные издержки и потери от проскальзывания торговой деятельности. Метрики вычисляются с помощью библиотеки `empirical`, которая также может использоваться самостоятельно.

Кадр данных результативности, создаваемый бэктестовым механизмом библиотеки `zipline`, может быть транслирован в требуемый для `pyfolio` вход.

Получение входа в библиотеку `pyfolio` из библиотеки `alphalens`

Вместе с тем библиотека `pyfolio` также интегрируется с библиотекой `alphalens` напрямую и позволяет создавать входные данные для `pyfolio` с помощью функции `create_pyfolio_input`:

```
from alphalens.performance import create_pyfolio_input

qmin, qmax = factor_data.factor_quantile.min(),
             factor_data.factor_quantile.max()
input_data = create_pyfolio_input(alphalens_data,
                                  period='1D',
                                  capital=100000,
                                  long_short=False,
                                  equal_weight=False,
                                  quantiles=[1, 5],
                                  benchmark_period='1D')

returns, positions, benchmark = input_data
```

Существует два варианта определения того, как будут генерироваться портфельные веса:

- ◆ `long_short` — если равно `False`, то веса будут соответствовать факторным значениям, разделенным на их абсолютное значение, поэтому отрицательные факторные значения генерируют короткие позиции. Если равно `True`, то из факторных значений сначала вычитается среднее, поэтому длинные и короткие позиции нейтрализуют друг друга, и портфель является рыночно нейтральным;
- ◆ `equal_weight` — если равно `True` и `long_short` равно `True`, то активы будут разделены на две равные группы, верхняя/нижняя половина которых будет составлять длинные/короткие позиции.

Длинно-короткие портфели также могут быть созданы для групп, если данные `factor_data` включают, например, информацию о секторе по каждому активу.

Получение входа в библиотеку `pyfolio` из бэктеста `zipline`

Результат бэктеста библиотеки `zipline` можно конвертировать в необходимый для `pyfolio` вход, используя функцию `extract_rets_pos_txn_from_zipline`:

```
returns, positions, transactions = extract_rets_pos_txn_from_zipline(backtest)
```

Форвардное тестирование вневыборочных возвратов

Тестирование торговой стратегии предусматривает бэктестирование на исторических данных с целью точной настройки альфа-факторных параметров, а также форвардное тестирование на новых рыночных данных с целью подтверждения того, что стратегия хорошо работает вне выборки или что параметры переподогнаны к конкретным историческим обстоятельствам.

Библиотека `pyfolio` допускает назначение вневыборочного периода для симулирования форвардного тестирования. При тестировании стратегии для получения статистически надежных результатов необходимо учитывать многочисленные аспекты, которые мы рассмотрим далее.

Функция `plot_rolling_returns` показывает кумулятивные внутривыборочные и вневыборочные возвраты в сопоставлении с задаваемым пользователем эталоном (мы используем фондовый индекс S&P 500):

```
from pyfolio.plotting import plot_rolling_returns
```

```
oos_date = '2017-01-01'
plot_rolling_returns(returns=returns,
                    factor_returns=benchmark_rets,
                    live_start_date=oos_date,
                    cone_std=(1.0, 1.5, 2.0))
```

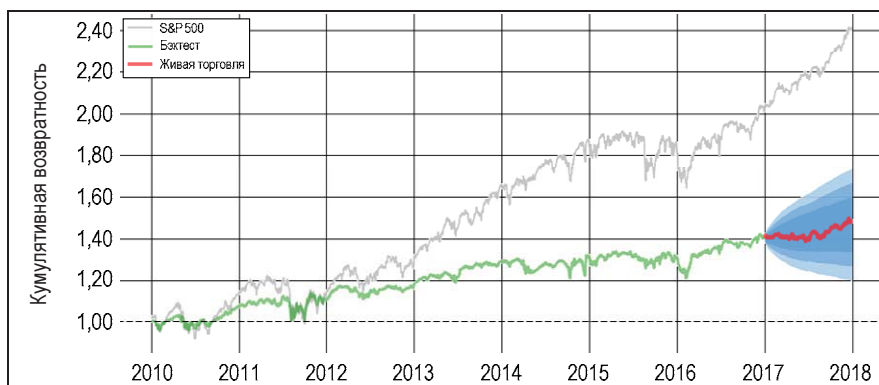


Рис. 5.1. Кумулятивные внутривыборочные и вневыборочные финансовые возвраты в сопоставлении с эталоном

Результирующий график включает конус, который показывает расширяющиеся доверительные интервалы, обозначающие, когда вневыборочный возврат выглядит маловероятным при наличии допущений о случайном блуждании (рис. 5.1). Здесь наша стратегия не показала хорошую результативность в сопоставлении с эталоном в течение симулированного вневыборочного периода за 2017 г.

Сводная статистика результативности

Библиотека `pyfolio` предлагает несколько аналитических функций и графиков. Статистическая сводка `perf_stats` показывает годовую и кумулятивные возвраты, волатильность, асимметрию и эксцесс возвратов, а также коэффициент Шарпа. Следующие дополнительные метрики (которые также могут быть рассчитаны индивидуально) являются наиболее важными:

- ◆ *максимальная просадка* (max drawdown) — самый высокий процент убытка от предыдущего пика;
- ◆ *коэффициент Калмар* (коэффициент просадки; calmar ratio — по названию фирмы CALifornia Managed Accounts Reports) — годовой портфельный возврат относительно максимальной просадки;
- ◆ *коэффициент омега* (omega ratio) — взвешенное по вероятности соотношение приростов и убытков для цели по возвратности, ноль по умолчанию;
- ◆ *коэффициент Сортино* (Sortino ratio, по фамилии создателя) — избыточный возврат по отношению к понижающему стандартному отклонению;
- ◆ *хвостовой коэффициент* (tail ratio) — размер правого хвоста (прироста, абсолютного значения 95-го перцентиля) относительно размера левого хвоста (убытков, абсолютного значения 5-го перцентиля);
- ◆ *дневная стоимость под риском* (daily value at risk, VaR)⁴ — убыток, соответствующий возврату на два стандартных отклонения ниже среднедневного значения;
- ◆ *альфа* (alpha) — портфельный возврат, необъясненный эталонным возвратом;
- ◆ *бета* (beta) — влияние эталонного показателя.

```
from pyfolio.timeseries import perf_stats

perf_stats(returns=returns,
            factor_returns=benchmark_rets,
            positions=positions,
            transactions=transactions)
```

⁴ Стоимость под риском (value at risk, VaR) — это статистический метод измерения портфельного риска: максимальная величина стоимости, потеря которой может ожидаться за данный временной горизонт. Например, если портфель имеет 95% дневной VaR в размере 100 тыс. долларов, то это означает, что, статистически говоря, существует 95% уверенность в том, что портфель не потеряет более 100 тыс. долларов стоимости в течение следующего дня. — Прим. перев.

Применительно к симулированному длинно-короткому портфелю, выведенному из фактора `MeanReversion`, мы получаем статистические показатели результативности, которые приведены в табл. 5.1.

Таблица 5.1. Сводные статистические показатели результативности

| Метрика | Все | Внутри выборки | Вне выборки |
|-----------------------------|---------|----------------|-------------|
| Годовая возвратность | 1,80% | 0,60% | 4,20% |
| Кумулятивная возвратность | 5,40% | 1,10% | 4,20% |
| Годовая волатильность | 5,80% | 6,30% | 4,60% |
| Шарп | 0,33 | 0,12 | 0,92 |
| Калмар | 0,17 | 0,06 | 1,28 |
| Стабильность | 0,49 | 0,04 | 0,75 |
| Максимальная просадка | -10,10% | -10,10% | -3,30% |
| Омега | 1,06 | 1,02 | 1,18 |
| Сортино | 0,48 | 0,18 | 1,37 |
| Асимметрия | 0,34 | 0,40 | 0,09 |
| Эксцесс | 3,37 | 3,70 | 2,59 |
| Хвостовой коэффициент | 0,91 | 0,88 | 1,03 |
| Дневная величина под риском | -0,7% | -0,8% | -0,6% |
| Валовое кредитное плечо | 0,38 | 0,37 | 0,42 |
| Ежедневный оборот | 4,70% | 4,40% | 5,10% |
| Альфа | 0,01 | 0,00 | 0,04 |
| Бета | 0,15 | 0,16 | 0,03 |

Подробнее о расчете и интерпретации метрических показателей портфельного риска и возвратности см. блокнот `pyfolio_demo.ipynb`.

Периоды просадки и влияние факторов

Функция `plot_drawdown_periods(returns)` строит график главных периодов просадки для портфеля, а несколько других графопостроительных функций показывают влияния на скользящий коэффициент Шарпа и скользящие факторы рыночных бета, или факторов размера, роста и импульса Фама — Френча:

```
fig, ax = plt.subplots(nrows=2, ncols=2, figsize=(16, 10))
axes = ax.flatten()
```

```

plot_drawdown_periods(returns=returns, ax=axes[0])
plot_rolling_beta(returns=returns,
                  factor_returns=benchmark_rets,
                  ax=axes[1])
plot_drawdown_underwater(returns=returns, ax=axes[2])
plot_rolling_sharpe(returns=returns)

```

Этот график, который выделяет подмножество визуализации, содержащейся в различных сводных таблицах, иллюстрирует, как библиотека `pyfolio` позволяет детализировать результативностные характеристики и влияние фундаментальных движущих сил риска и возвратов (рис. 5.2).

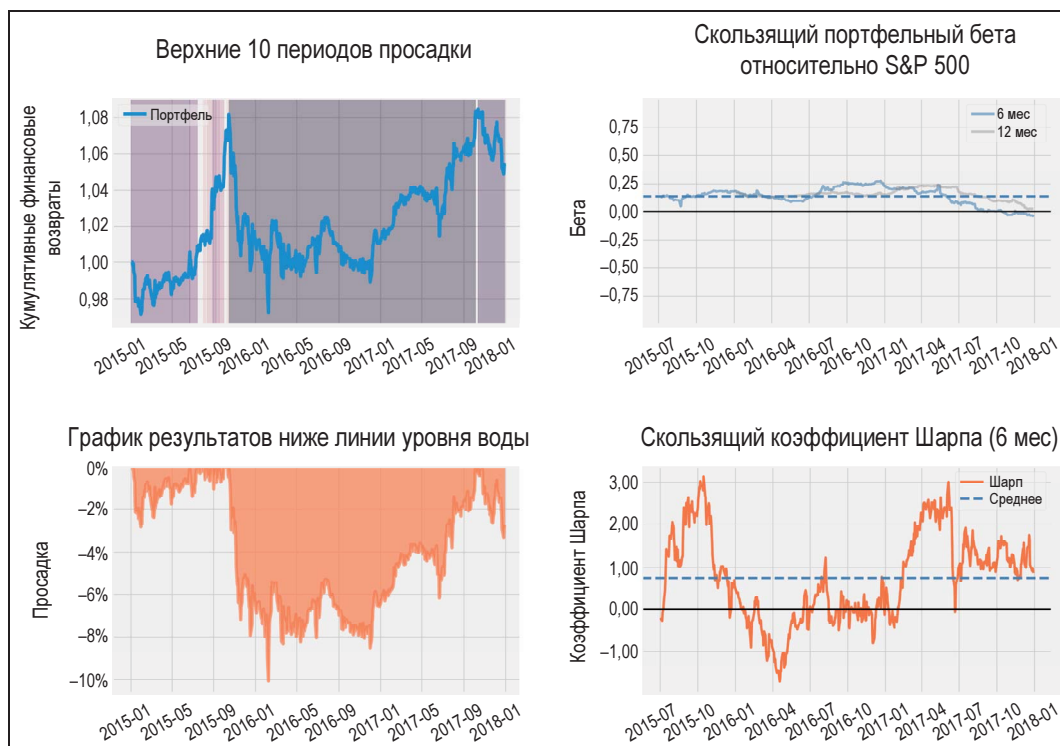


Рис. 5.2. Визуализация нескольких графиков, полученных из различных сводных таблиц библиотеки `pyfolio`

Моделирование событийного риска

Библиотека `pyfolio` также содержит временные шкалы для всевозможных событий, которые можно использовать для сравнения портфельной результативности с эталоном в течение этого периода, например, во время осенней распродажи 2015 г. после голосования по Брекситу:

```
interesting_times = extract_interesting_date_ranges(returns=returns)
(interesting_times['Fall2015'].to_frame('pf')
 .join(benchmark_rets)
 .add(1).cumprod().sub(1)
 .plot(lw=2, figsize=(14, 6),
       title='Паника после голосования по Брекситу'))
```

Результирующий график выглядит так, как показано на рис. 5.3.

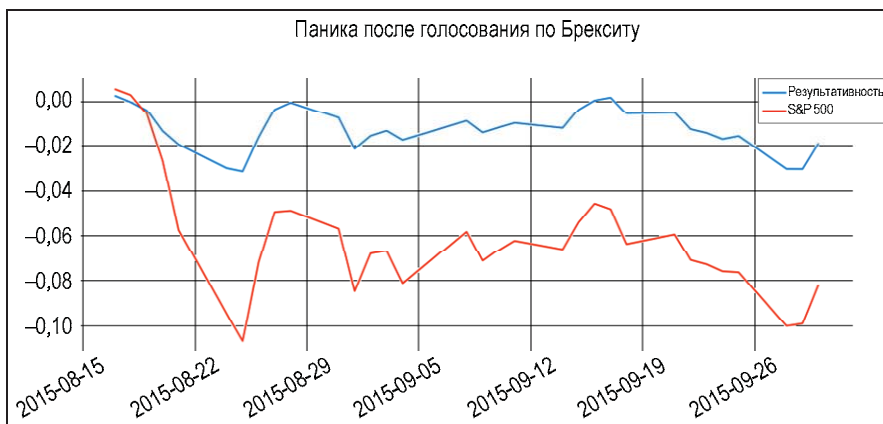


Рис. 5.3. Сравнение результативности портфеля с эталоном во время определенного события

Как избегать ловушек бэктестирования

Бэктестирование симулирует алгоритмическую стратегию, используя исторические данные с целью выявления регулярностей, которые обобщаются на новые рыночные условия. Помимо общих проблем, связанных с предсказанием неопределенного будущего в условиях меняющихся рынков, многочисленные факторы делают весьма вероятным ошибочное принятие положительной внутривыборочной результативности за обнаружение истинных регулярностей. Эти факторы включают аспекты данных, реализацию симулирования стратегии и недостатки статистических проверок и их интерпретаций. Риски ложных обнаружений умножаются с использованием большей вычислительной мощи, более крупных наборов данных и более сложных алгоритмов, которые обеспечивают выявление различных регулярностей в шуме.

Мы перечислим наиболее серьезные и распространенные методологические ошибки и сошлемся на литературу по множественному тестированию для получения более подробной информации. Мы также представим дефлированный коэффициент Шарпа, который иллюстрирует то, как можно корректировать метрики, полученные в результате повторных испытаний при использовании для анализа одного и того же набора финансовых данных.

Сложности данных

Проблемы достоверности результатов бэктестирования, обусловленные вопросами, связанными с данными, включают систематическое смещение из-за забегания вперед (look-ahead bias), систематическое смещение из-за выживших (survivorship bias) и управление выбросами.

Систематическое смещение из-за забегания вперед

Тесты торговых правил, выведенных из прошлых данных, будут выдавать смещенные результаты, когда отобранные данные, используемые для разработки правил, содержат информацию, которая фактически не имела в наличии или не была известна на момент времени, на который эти данные ссылаются.

Типичным источником смещения из-за забегания вперед является неучет общих исправлений, вносимых в финансовые отчеты после публикации. Дробления или консолидации (обратные дробления) акций также могут генерировать систематическое смещение из-за забегания вперед. При расчете отдачи от корпоративных заработков данные о заработках в расчете на долю собственности в акционерном капитале (EPS) поступают из финансовых отчетов компании с низкой частотой, в то время как рыночные цены доступны практически ежедневно. Следовательно, и данные о заработках в расчете на долю, и данные о ценах должны корректироваться с одновременным учетом дробления/консолидации.

Решение заключается в тщательном анализе временных меток, ассоциированных со всеми данными, которые вводятся в бэктекст, с целью обеспечения того, что использовались данные только с определенной точкой во времени (PIT). Высококачественные поставщики данных, такие как Compustat (www.compustat.com/), обеспечивают соблюдение этих критериев. При отсутствии данных в определенной точке во времени необходимо принимать допущения о временном сдвиге в отчетности.

Систематическое смещение из-за выживших

Систематическое смещение из-за выживших возникает, когда тест проводится на данных, которые содержат только активные в настоящее время ценные бумаги и пропускают активы, которые исчезли с течением времени, например, из-за банкротства, делистинга или приобретения. Ценные бумаги, которые больше не являются частью инвестиционного универсума, часто не показывали хорошую результативность, и включение этих случаев может положительно исказить бэктестовый результат.

Решение, естественно, состоит в перепроверке того, чтобы наборы данных включали все ценные бумаги, имевшиеся в рассматриваемый период, а не только те, которые все еще доступны при выполнении теста.

Управление выбросами

Подготовка данных перед анализом обычно включает ту или иную трактовку выбросов, например, путем винзоризации или отсека экстремальных значений.

Трудность здесь лежит в выявлении выбросов, которые действительно не отражают анализируемый период, в отличие от экстремальных значений, которые являются неотъемлемой частью рыночной среды в это время. Многие рыночные модели исходят из нормального распределения данных, когда экстремальные значения наблюдаются чаще, как это предполагается распределениями с толстыми хвостами.

Решение данной проблемы предусматривает тщательный анализ выбросов относительно вероятности возникновения экстремальных значений и приведение параметров стратегии в соответствие с этой реальностью.

Нерепрезентативный период

Бэктест не даст репрезентативного результата, который будет обобщаться на будущие периоды, если используемый период времени не отражает хорошо текущую среду, не имеет в достаточном объеме соответствующих аспектов рыночного режима и не включает достаточное число точек данных или улавливает экстремальные исторические события, которые вряд ли повторятся.

Решение этой проблемы предусматривает использование выборочных периодов, включающих важные рыночные явления, либо генерирование синтетических данных, отражающих соответствующие рыночные характеристики.

Вопросы реализации

Практические вопросы, связанные с реализацией исторического моделирования, включают неспособность выполнять рыночную переоценку, т. е. точно отражать базовые рыночные цены и учитывать просадки, нереалистичные допущения о наличии, стоимости или рыночного влияния сделок или продолжительность сигналов и сроки исполнения сделок.

Результативность на основе рыночной цены

Эта стратегия может показывать хорошую результативность в ходе бэктестирования, но со временем привести к неприемлемым убыткам или волатильности.

Решение данной проблемы предусматривает построение графика результативности во временной динамике или расчет (скользящих) рискованных метрик, таких как *стоимость под риском* (VaR) или коэффициент Сортино (подробности см. в блокаде).

Торговые издержки

Эта стратегия может допускать короткие продажи, т. е. продажи в ожидании снижения цены, требующие наличие контрагента, держать меньше ликвидных активов, которые могут двигать рынок при их продаже, либо недооценивать затраты, возникающие из-за брокерских комиссий или проскальзывания, т. е. разницы между рыночной ценой при принятии решения торговать и последующим исполнением.

Решение данной проблемы включает ограничение только высоколиквидным универсумом и принятие реалистичных параметрических допущений для торговых издержек и издержек, связанных с проскальзыванием (как показано в предыдущем примере использования библиотеки `zipline`). Это также служит гарантией от включения неустойчивых факторных сигналов с высоким затуханием и, следовательно, неустойчивой оборачиваемостью.

Определение времени сделок

Симулирование может принимать нереалистичные допущения о времени оценивания альфа-факторных сигналов и результирующих сделок. Например, сигналы могут оцениваться по ценам закрытия, когда следующая сделка доступна только по совершенно (и нередко) другим ценам открытия. Как следствие, бэктест будет значительно смещен, когда для оценивания результативности торговли используется цена закрытия.

Решение данной проблемы предусматривает тщательную оркестровку последовательности прихода сигнала, исполнения сделки и оценивания ее результативности.

Прочесывание данных и переподгонка при бэктестировании

Наиболее серьезная сложность, связанная с достоверностью результатов бэктестирования, в том числе опубликованных результатов, связана и с обнаружением мнимых регулярностей вследствие многократного тестирования во время процесса отбора стратегии⁵. Отбор стратегии после тестирования разных кандидатов на одних и тех же данных, скорее всего, сместит выбор, поскольку положительный результат, ожидаемо, будет обусловлен стохастической природой самой меры результативности. Другими словами, стратегия будет чрезмерно приспособлена, или подогнана, к имеющимся данным и будет давать обманчиво положительные результаты.

Таким образом, бэктестовая результативность не является информативной, если только не будут предоставлены сведения о числе испытаний, позволяющие оценить риск систематического смещения при отборе⁶. Это редко встречается в практических или академических исследованиях, поскольку вызывает сомнения в достоверности многих опубликованных утверждений.

Риск переподгонки во время бэктеста к конкретному набору данных возникает не только из-за непосредственного прогона многочисленных тестов, но и включает стратегии, разработанные на основе априорных (предшествующих) знаний о том,

⁵ Прочесывание данных (data snooping) — подробная ревизия данных с целью найти что-то интересное. — *Прим. перев.*

⁶ Систематическое смещение при отборе (selection bias) — систематическая ошибка, выражающаяся в появлении у изучаемой выборки признаков, не свойственных генеральной совокупности; возникает в результате применения неподходящего метода отбора. — *Прим. перев.*

что работает, а что нет, т. е. знаний о различных бэктестах, выполнявшихся другими на тех же самых данных. По этой причине трудно избежать бэктестовой переподгонки на практике.

Решения данной проблемы включают отбор планируемых тестов на основе инвестиционной или экономической теории, а не на широких усилиях по глубинному анализу данных. Они также предусматривают тестирование во всевозможных контекстах и ситуациях, в том числе, возможно, на синтетических данных.

Минимальная длина бэктеста и дефлированный коэффициент Шарпа

Маркос Лопес де Прадо⁷ (Marcos Lopez de Prado, <http://www.quantresearch.info/>) опубликовал обширный перечень статей о рисках бэктестирования и о том, как их обнаруживать или избегать. Кроме того, он разместил онлайн-симулятор бэктестовой переподгонки (<http://datagrid.lbl.gov/backtest/>).

Еще один результат включает оценку минимальной длины бэктеста, которую инвестор должен потребовать при заданном числе предпринятых попыток во избежание отбора стратегии с заданным внутривыборочным коэффициентом Шарпа во время заданного числа испытаний. Такая стратегия имеет ожидаемый вневыборочный нулевой коэффициент Шарпа. Из этого следует, что, например, если имеется только два года ежедневных бэктестовых данных, то следует опробовать не более семи вариантов стратегии, а если имеется пять лет ежедневных бэктестовых данных, то следует опробовать не более 45 вариантов стратегии. Подробности реализации см. в справочных материалах.

Лопес де Прадо и Бейли (Marcos Lopez de Prado и Bailey, 2014) также выводят дефлированный коэффициент Шарпа, вычисляя вероятность того, что коэффициент Шарпа является статистически значимым при учете инфляционного эффекта множественного тестирования, ненормальных финансовых возвратов и более коротких длин выборок (см. подпапку 03_multiple_testing, где находится Python-реализация дефлированного коэффициента Шарпа — в сценарном файле `deflated_sharpe_ratio.py` — и справочные материалы, касающиеся вывода родственных формул).

Оптимальное прекращение бэктестирования

В дополнение к ограничению бэктестирования только теми стратегиями, которые могут быть обоснованы на теоретических основаниях, в отличие от простых упражнений по глубинному анализу данных, важным является ответ на вопрос: когда прекращать выполнение дополнительных тестов?

⁷ В 2019 г. вышел перевод книги Маркоса Лопеса де Прадо под заголовком "Машинное обучение: алгоритмы для бизнеса". В оригинале книга называется "Advances in Financial Machine Learning, Marcos Lopez de Prado", ©2018 by John Wiley & Sons ("Достижения в финансовом машинном обучении"). — Прим. перев.

Основываясь на решении *задачи секретаря* из теории оптимального прекращения, рекомендуется принимать решение в соответствии со следующим эмпирическим правилом: протестировать случайную выборку из $1/e$ (примерно 37%) разумных стратегий и записать их результативность. Затем продолжать тестировать до тех пор, пока стратегия не превзойдет протестированные ранее.

Это правило применяется к тестам нескольких альтернатив с целью как можно скорее выбрать почти лучшую из них, минимизируя риск ложного утверждения.

Как управлять портфельным риском и возвратностью

Под портфельным менеджментом понимается занятие позиций в финансовых инструментах таких, которые достигают желаемого компромисса между риском и возвратностью относительно эталона. В каждом периоде менеджер отбирает позиции, которые оптимизируют диверсификацию с целью снижения рисков, при этом достигая целевого финансового возврата. В разные периоды позиции будут перебалансироваться для учета изменений в весах, возникающих в результате ценовых движений, с целью достижения или поддержания целевого рискового профиля.

Диверсификация позволяет снижать риски для заданного ожидаемого возврата за счет эксплуатации того, как ценовые движения взаимодействуют друг с другом, поскольку приросты в одном активе могут компенсировать убытки другого актива. Гарри Марковиц изобрел *теорию современного инвестиционного портфеля* (Modern Portfolio Theory, MPT) в 1952 г. и предоставил математические инструменты, которые позволяют оптимизировать диверсификацию путем выбора соответствующих портфельных весов. Марковиц показал, как портфельный риск, измеряемый как стандартное отклонение портфельных возвратов, зависит от ковариации (совместной дисперсии) возвратов от всех активов и их относительных весов. Эта взаимосвязь приводит к существованию эффективной портфельной границы, которая максимизирует портфельные возвраты при наличии максимального уровня портфельного риска.

Однако среднedisперсные границы очень чувствительны к оценкам входных данных, необходимых для их расчета, таким как ожидаемые возвраты, волатильности и корреляции. На практике среднedisперсные портфели, которые ограничивают эти входные данные с целью снижения ошибок выборки, показывают намного более высокую результативность. Эти ограниченные частные случаи включают равновзвешенные, минимально-дисперсные и риск-паритетные портфели.

Модель ценообразования капитальных активов (Capital Asset Pricing Model, CAPM) является моделью оценивания стоимости активов, которая строится на взаимосвязи между риском и возвратностью, постулируемой теорией современного инвестиционного портфеля (MPT). Она вводит концепцию рисковой премии, которую инвестор может ожидать в рыночном равновесии для удержания рискованного актива; премия компенсирует временную стоимость денег и влияние общего рыночного

риска, который не может быть устранен путем диверсификации (в отличие от идиосинкратического риска конкретных активов). Экономическое обоснование недиверсифицируемого риска заключается, например, в макродрайверах деловых рисков, сказывающихся на возвратах от долевого ценных бумаг или дефолтах по облигациям. Таким образом, математическое ожидание возврата от актива, $E[r_i]$, является суммой безрисковой процентной ставки r_f и рисковой премии, пропорциональной влиянию на актив ожидаемого избыточного возврата рыночного портфеля r_m , превышающего безрисковую ставку:

$$E[r_i] = \alpha_i + r_f + \beta_i (E[r_m] - r_f).$$

Теоретически, рыночный портфель содержит все пригодные для инвестирования активы и будет удерживаться всеми рациональными инвесторами в равновесии. На практике широкий взвешенный по стоимости индекс аппроксимирует рынок, например фондовый индекс S&P 500, для инвестиций в долевого ценные бумаги компаний США. β_i служит мерой влияния на актив повышенных возвратов рыночного портфеля. Если модель CAPM соблюдается, то компонент пересечения α_i должен быть равен нулю. На самом деле допущения, принимаемые моделью CAPM, часто не соблюдаются, и альфа улавливает возвраты, оставшиеся необъясненными влиянием широкого рынка.

Со временем исследование выявило нетрадиционные источники рисков премий, такие как импульсный эффект или эффект стоимости долевого ценной бумаги, которые давали объяснение некоторым первоначальным альфа. Экономические обоснования, такие как поведенческие тенденциозности (смещения) из-за недостаточной или чрезмерной реакции инвесторов на новую информацию, оправдывают рисковые премии за влияние этих альтернативных рисков факторов. Они эволюционировали в инвестиционные стили, призванные улавливать эти альтернативные бета, которые также стали торгуемыми в виде специализированных индексных фондов. После изоляции вкладов этих альтернативных рисков премий истинный альфа становится ограниченным идиосинкратическими возвратами от активов и способностью менеджера определять время влияния рисков.

В течение последних нескольких десятилетий гипотеза об эффективном рынке (Efficient Market Hypothesis, EMH) совершенствовалась, выправляя многие из первоначальных недостатков CAPM-модели, включая несовершенную информацию и расходы, ассоциированные с транзакциями, финансированием и агентивностью. Многие поведенческие тенденциозности (смещения) имеют тот же эффект, и некоторые трения моделируются как поведенческие тенденциозности.

МО играет важную роль в выведении новых альфа-факторов с использованием технических решений контролируемого и неконтролируемого автоматического обучения, базирующихся на рыночных, фундаментальных и альтернативных источниках данных, рассмотренных в предыдущих главах. Входы в автоматически обучающиеся модели состоят из сырых данных и признаков, вырабатываемых для улавливания информативных сигналов. Автоматически обучающиеся модели также

используются для совмещения индивидуальных предсказательных сигналов и обеспечивают высокоагрегированную предсказательную мощность.

За последние несколько десятилетий теория современного инвестиционного портфеля и ее практика значительно эволюционировали. Мы представим:

- ◆ среднedisперсную оптимизацию и ее недостатки;
- ◆ альтернативы, такие как минимально-рисковое размещение и размещение $1/n$;
- ◆ риск-паритетные подходы;
- ◆ риск-факторные подходы.

Среднedisперсная оптимизация

Теория современного инвестиционного портфеля (MPT) дает решение оптимальных портфельных весов, минимизируя волатильность для заданного ожидаемого возврата либо максимизируя возвраты для заданного уровня волатильности. Ключевыми необходимыми входами являются ожидаемые возвраты от активов, стандартные отклонения и матрица ковариаций.

Как это работает

Диверсификация работает, потому что дисперсия портфельных возвратов зависит от ковариации активов и может быть снижена до уровня ниже средневзвешенного значения дисперсий активов путем включения активов с менее идеальной корреляцией. В частности, с учетом вектора ω , портфельных весов и матрицы ковариаций Σ , портфельная дисперсия σ_{PF} определяется как:

$$\sigma_{PF} = \omega^T \Sigma \omega.$$

Марковиц показал, что задача максимизации ожидаемого портфельного возврата с учетом целевого риска имеет эквивалентное двойное представление в виде минимизации портфельного риска с учетом целевого уровня ожидаемого возврата μ_{PF} . Следовательно, задачей оптимизации становится:

$$\begin{aligned} \min_{\omega} \sigma_{PF}^2 &= \omega^T \Sigma \omega; \\ \text{при условии, что } \mu_{PF} &= \omega^T \mu; \\ \|\omega\| &= 1. \end{aligned}$$

Эффективная граница на Python

Мы можем рассчитать эффективную границу, используя функцию `scipy.optimize.minimize` и исторические оценки возвратов от активов, стандартные отклонения и матрицу ковариаций. Исходный код можно найти в подпапке `efficient_frontier` в репозитории для этой главы. В нем реализуется такая последовательность шагов:

1. Симуляция генерирует случайные веса, используя размещение Дирихле, и вычисляет среднее значение, стандартное отклонение и коэффициент Шарпа для каждого образца портфеля на основе исторических данных финансовых возвратов:

```
def simulate_portfolios(mean_ret, cov, rf_rate=rf_rate, short=True):
    alpha = np.full(shape=n_assets, fill_value=.025)
    weights = dirichlet(alpha=alpha, size=NUM_PF)
    weights *= choice([-1, 1], size=weights.shape)

    returns = weights @ mean_ret.values + 1
    returns = returns ** periods_per_year - 1
    std = (weights @ monthly_returns.T).std(1)
    std *= np.sqrt(periods_per_year)
    sharpe = (returns - rf_rate) / std

    return pd.DataFrame({'Стандартное отклонение в годовом исчислении':
                        std,
                        'Финансовые возвраты в годовом исчислении':
                        returns, 'Коэффициент Шарпа': sharpe}), weights
```

2. Ставится задача квадратической оптимизации — найти минимальное стандартное отклонение для заданного финансового возврата или максимального коэффициента Шарпа. Для этого определяются функции, которые измеряют ключевые метрики:

```
def portfolio_std(wt, rt=None, cov=None):
    """Портфельное стандартное отклонение в годовом исчислении"""
    return np.sqrt(wt @ cov @ wt * periods_per_year)

def portfolio_returns(wt, rt=None, cov=None):
    """Портфельные возвраты в годовом исчислении"""
    return (wt @ rt + 1) ** periods_per_year - 1

def portfolio_performance(wt, rt, cov):
    """ Портфельные возвраты и стандартное отклонение
        в годовом исчислении"""
    r = portfolio_returns(wt, rt=rt)
    sd = portfolio_std(wt, cov=cov)
    return r, sd
```

3. Задается целевая функция, представляющая отрицательный коэффициент Шарпа для функции `minimize` библиотеки `SciPy`, которая оптимизирует с учетом ограничений, что веса ограничиваются промежутком $[-1; 1]$ и в сумме составляют единицу в абсолютном выражении:

```
def neg_sharpe_ratio(weights, mean_ret, cov):
    r, sd = portfolio_performance(weights, mean_ret, cov)
    return -(r - rf_rate) / sd
```



```
weight_constraint = {'type': 'eq',
                    'fun': lambda x: np.sum(np.abs(x)) - 1}

def max_sharpe_ratio(mean_ret, cov, short=True):
    return minimize(fun=neg_sharpe_ratio,
                   x0=x0,
                   args=(mean_ret, cov),
                   method='SLSQP',
                   bounds=((-1 if short else 0, 1),) * n_assets,
                   constraints=weight_constraint,
                   options={'tol': 1e-10, 'maxiter': 1e4})
```

4. Эффективная граница вычисляется путем перебора диапазона целевых финансовых возвратов и нахождения соответствующих минимально-дисперсных портфелей. Оптимизационная задача и ограничения, накладываемые на портфельный риск и возвратность, как функция от весов, могут быть сформулированы следующим образом:

```
def min_vol_target(mean_ret, cov, target, short=True):

    def ret_(wt):
        return portfolio_returns(wt, mean_ret)

    constraints = [{'type': 'eq', 'fun': lambda x: ret_(x) - target},
                  weight_constraint]

    bounds = ((-1 if short else 0, 1),) * n_assets
    return minimize(portfolio_std,
                   x0=x0,
                   args=(mean_ret, cov),
                   method='SLSQP',
                   bounds=bounds,
                   constraints=constraints,
                   options={'tol': 1e-10, 'maxiter': 1e4})
```

5. Решение требует итеративного перебора интервалов допустимых значений для определения оптимальных комбинаций компромисса между риском и возвратностью:

```
def min_vol(mean_ret, cov, short=True):
    bounds = ((-1 if short else 0, 1),) * n_assets

    return minimize(fun=portfolio_std,
                   x0=x0,
                   args=(mean_ret, cov),
                   method='SLSQP',
                   bounds=bounds,
```

```
constraints=weight_constraint,
options={'tol': 1e-10, 'maxiter': 1e4})
```

```
def efficient_frontier(mean_ret, cov, ret_range, short=True):
    return [min_vol_target(mean_ret, cov, ret) for ret in ret_range]
```

Симуляция выдает подмножество допустимых портфелей, а эффективная граница определяет оптимальные внутривыборочные комбинации компромисса между риском и возвратностью, которые были достижимы с учетом исторических данных. На рис. 5.4 показан результат, включающий минимально-дисперсный портфель и портфель, максимизирующий коэффициент Шарпа, а также несколько портфелей, порожденных альтернативными оптимизационными стратегиями, которые мы обсудим в последующих разделах.

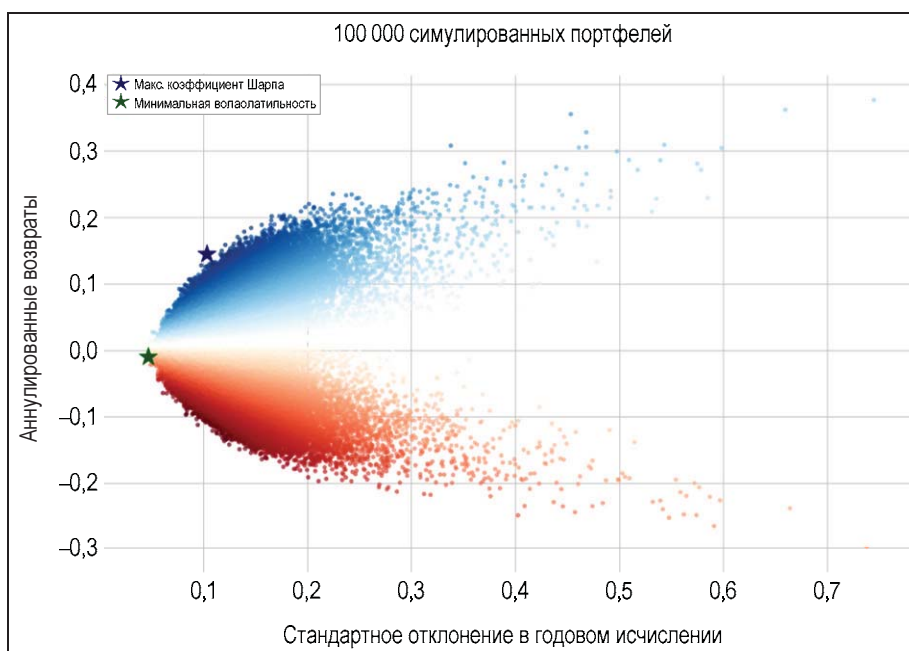


Рис. 5.4. Диаграмма эффективной границы с подмножеством допустимых портфелей, порожденных разными оптимизационными стратегиями

Портфельная оптимизация может выполняться на каждом шаге оценивания торговой стратегии с целью оптимизации позиций.

Сложности и недостатки

Предыдущий пример среднedisперсной границы иллюстрирует внутривыборочную оптимизацию в обратном направлении. На практике портфельная оптимизация требует форвардного подхода. Ожидаемые возвраты, как известно, трудно оценить точно.

Матрица ковариаций может быть оценена слегка надежнее, что привело к появлению нескольких альтернативных подходов. Однако матрицы ковариаций с коррелированными активами создают вычислительные сложности, поскольку оптимизационная задача требует инвертирования матрицы. Высокое число обусловленности индуцирует численную нестабильность, которая, в свою очередь, порождает проклятие Марковица: чем больше диверсификации требуется (коррелированной инвестиционными возможностями), тем ненадежнее являются веса, получаемые алгоритмом.

Многие инвесторы предпочитают использовать технические решения для портфельной оптимизации с менее обременительными требованиями ко входу. Далее мы предложим несколько альтернатив, призванных устранять эти недостатки, включая более современные подходы, основанные на машинном обучении.

Альтернативы среднedisперсной оптимизации

Сложности с точным входом для задачи среднedisперсной оптимизации привели к принятию нескольких практических альтернатив, которые ограничивают среднее значение, дисперсию или оба показателя, или опускают более сложные оценки возврата, такие как подход на основе паритета риска.

Портфель $1/n$

Простые портфели дают полезные эталоны для зондирования добавочной стоимости сложных моделей, которые генерируют риск переподгонки. Самая простая стратегия — равновзвешенный портфель — зарекомендовала себя как одна из лучших.

Как известно, де Мигуэль, Гарлаппи и Уппал (de Miguel, Garlappi, Uppal, 2009) сравнили вневыборочную результативность портфелей, порождаемых различными среднedisперсными оптимизаторами, включая робастные байесовы оценщики, портфельные ограничения и оптимальные комбинации портфелей, с простым правилом $1/n$. Они обнаружили, что портфель $1/n$ дает более высокий коэффициент Шарпа, чем каждая позиция класса активов, что объясняется высокой стоимостью ошибок оценивания, которые часто перевешивают преимущества изоэкренной вневыборочной оптимизации.

Портфель $1/n$ также включен в приведенную выше диаграмму эффективной границы.

Минимально-дисперсный портфель

Еще одной альтернативой является портфель с *глобально-минимальной дисперсией* (global minimum variance, GMV), в котором приоритет отдается минимизации риска. Он показан на диаграмме эффективной границы и может быть рассчитан путем минимизации портфельного стандартного отклонения с использованием среднedisперсного вычислительного каркаса следующим образом:

```
def min_vol(mean_ret, cov, short=True):
    bounds = ((-1 if short else 0, 1),) * n_assets

    return minimize(fun=portfolio_std,
                    x0=x0,
                    args=(mean_ret, cov),
                    method='SLSQP',
                    bounds=bounds,
                    constraints=weight_constraint,
                    options={'tol': 1e-10, 'maxiter': 1e4})
```

Соответствующий минимально волатильный портфель лежит на эффективной границе, как показано выше.

Глобальная портфельная оптимизация — подход Блэка — Литтермана

Подход Блэка — Литтермана (Black и Litterman, 1992) к глобальной портфельной оптимизации сочетает экономические модели со статистическим обучением и пользуется популярностью, поскольку позволяет генерировать оценки ожидаемых финансовых возвратов, которые являются правдоподобными во многих ситуациях.

Данное техническое решение отходит от допущения, что рынок является средне-дисперсным портфелем, который вытекает из равновесной модели CAPM, и строится на том факте, что наблюдаемая рыночная капитализация может рассматриваться как оптимальные веса, назначаемые рынком. Рыночные веса отражают рыночные цены, которые, в свою очередь, воплощают ожидания рынка относительно будущих возвратов.

Таким образом, этот подход может реконструировать ненаблюдаемые будущие ожидаемые возвраты из допущения о том, что рынок находится достаточно близко к равновесию, как определено моделью CAPM, и позволяет инвесторам корректировать эти оценки в соответствии с их собственными мнениями с помощью усадочного оценщика⁸. Модель может быть интерпретирована как байесов подход к портфельной оптимизации. Мы представим байесовы методы в *главе 9*.

Как определять размер ставок — правило Келли

Правило Келли имеет долгую историю в азартных играх, потому что оно дает указания о том, сколько ставить на каждую (бесконечную) последовательность ставок с варьирующимися (но благоприятными) шансами, максимизируя конечное богатство. Оно было опубликовано как новая интерпретация скорости информации

⁸ В статистике усадочный оценщик (shrinkage estimator) — это оценщик, или правило оценки, который дает улучшенную оценку величины за счет явного или неявного встраивания другой информации. В этом смысле используется для регуляризации неточно поставленных задач статистического вывода. См. https://en.wikipedia.org/wiki/Shrinkage_estimator. — *Прим. перев.*

в 1956 г. Джоном Келли (John Kelly), который был коллегой Клода Шеннона (Claude Shannon) в Bell Labs. Его заинтриговали ставки, делаемые на кандидатов в новой викторине "Вопрос на 64 000 долларов", где зритель с западного побережья использовал трехчасовую задержку с целью получения инсайдерской информации о победителях.

Келли провел взаимосвязь с информационной теорией Шеннона с целью найти решение для ставки, которая является оптимальной для долговременного роста капитала, когда шансы благоприятны, но неопределенность остается. Его правило максимизирует логарифмическое богатство как функцию от шансов на успех каждой партии игры и включает в себя неявную защиту от банкротства, поскольку $\log(0)$ является отрицательной бесконечностью, вследствие чего игрок Келли естественным образом избежит потери всего.

Оптимальный размер ставки

Келли начал с анализа игр с бинарным результатом "победа-поражение". Ключевыми переменными являются:

- ◆ b — коэффициенты задают сумму, выигрываемую для ставки 1 доллар. Шанс равен 5:1 приводит к выигрышу 5 долларов, если ставка выигрывает, плюс восстановление однодолларового капитала;
- ◆ p — вероятность задает правдоподобие благоприятного результата;
- ◆ f — доля текущего капитала для ставок;
- ◆ V — стоимость капитала в результате сделанной ставки.

Правило Келли призвано максимизировать скорость роста добавочной стоимости G бесконечно повторяющихся ставок:

$$G = \lim_{N \rightarrow \infty} \frac{1}{N} \log \frac{V_N}{V_0}.$$

Когда W и L являются числами выигрышей и проигрышей, тогда:

$$\begin{aligned} V_N &= (1 + b * f)^W (1 - f)^L V_0 \Rightarrow \\ \Rightarrow G &= \lim_{N \rightarrow \infty} \left[\frac{W}{N} \log(1 + \text{odds} * \text{share}) + \frac{L}{N} \log(1 - f) \right] = \\ &= p \log(1 + b * f) + (1 - p) \log(1 - f). \end{aligned}$$

Мы можем максимизировать скорость роста G , максимизируя G по f , как показано на примере с использованием библиотеки SymPy следующим образом:

```
from sympy import symbols, solve, log, diff
from scipy.optimize import minimize_scalar, newton, minimize
from scipy.integrate import quad
from scipy.stats import norm
```

```
share, odds, probability = symbols('share odds probability')
Value = probability * log(1 + odds * share) + (1 - probability) * log(1 - share)
solve(diff(Value, share), share)
```

```
[(odds*probability + probability - 1)/odds]
```

Мы приходим к оптимальной доле капитала для ставок (критерий Келли):

$$f^* = \frac{b * p + p - 1}{b}.$$

Оптимальная инвестиция — единственный актив

В контексте финансового рынка сложность результатов и альтернатив выше, но логика правила Келли по-прежнему применима. Оно стало популярным благодаря Эду Торпу (Ed Thorp), который сначала применил его к азартным играм (описанным в книге "Beat The Dealer" ("Победить крупье")), а затем основал успешный хедж-фонд Princeton/Newport Partners.

В случае непрерывных результатов темпы роста капитала определяются интегрированием над вероятностным распределением разных финансовых возвратов, которые могут быть оптимизированы численно:

$$E[G] = \int \log(1 * fr) P(r) dr \Leftrightarrow$$

$$\Leftrightarrow \frac{d}{df} E[G] = \int_{-\infty}^{+\infty} \frac{r}{1 * fr} P(r) dr = 0.$$

Мы можем решить это выражение для оптимального f^* с помощью программного модуля `scipy.optimize`:

```
def norm_integral(f, mean, std):
    val, er = quad(lambda s: np.log(1 + f * s) * norm.pdf(s, mean, std),
                    mean - 3 * std,
                    mean + 3 * std)

    return -val

def norm_dev_integral(f, mean, std):
    val, er = quad(lambda s: (s / (1 + f * s)) * norm.pdf(s, mean, std),
                    m-3*std, mean+3*std)

    return val

m = .058
s = .216

# Вариант 1: минимизировать математическое ожидание интеграла
sol = minimize_scalar(norm_integral, args=(m, s),
                       bounds=[0., 2.], method='bounded')
print('Оптимальная доля Келли: {:.4f}'.format(sol.x))
```

```
# Вариант 2: взять производную математического ожидания и приравнять нулю
x0 = newton(norm_dev_integral, .1, args=(m, s))
print('Оптимальная доля Келли: {:.4f}'.format(x0))
```

Оптимальная инвестиция — многочисленные активы

Мы будем использовать пример с различными долевыми ценными бумагами. Е. Чан (E. Chan, 2008) иллюстрирует, как можно прийти к применению правила Келли для случая с многочисленными активами, и что результат эквивалентен (потенциально с кредитным плечом) портфелю с максимальным коэффициентом Шарпа из среднedisперсной оптимизации.

Вычисление предусматривает скалярное произведение матрицы прецизионностей, которая является инверсией матрицы ковариаций, и матрицы возвратов:

```
mean_returns = monthly_returns.mean()
cov = monthly_returns.cov()
precision_matrix = pd.DataFrame(inv(cov), index=stocks, columns=stocks)
kelly_allocation = mean_returns.dot(precision_matrix)

kelly_allocation.sum()

46.433472367624404
```

Портфель Келли также показан на диаграмме эффективной границы (после нормализации, вследствие которой абсолютные веса в сумме составляют единицу). Многие инвесторы предпочитают снижать веса Келли с целью снижения волатильности стратегии, и в связи с этим веса полу-Келли (Half-Kelly) приобрели особую популярность.

Паритет риска

Тот факт, что предыдущие 15 лет охарактеризовались двумя крупными кризисами на глобальных рынках долевого ценного бумага, неуклонно восходящей кривой возвратности и общим снижением процентных ставок, делает паритет риска⁹ особенно привлекательным вариантом. Многие учреждения выкраивают финансовые средства для стратегических размещений с целью обеспечения паритета риска, таким образом еще более диверсифицируя свои портфели.

Простая реализация паритета риска распределяет финансовые средства среди портфельных активов в соответствии с инверсиями их дисперсий, игнорируя корреляции и, в частности, прогнозы финансовых возвратов:

⁹ Паритет риска (risk parity) — это стратегия перераспределения портфеля с использованием риска для определения размещений по различным компонентам инвестиционного портфеля. Подход к управлению портфелем на основе паритета риска основан на теории современного инвестиционного портфеля, которая стремится оптимально диверсифицировать инвестиционный портфель среди указанных активов. См. <https://www.investopedia.com/terms/r/risk-parity.asp>. — Прим. перев.

```
var = monthly_returns.var()  
risk_parity_weights = var / var.sum()
```

Риск-паритетный портфель также показан на диаграмме эффективной границы в начале этого раздела.

Риск-факторное инвестирование

Альтернативным вычислительным каркасом для оценивания входа является сведение к базовым детерминантам, или факторам, которые приводят в действие риск и возвраты от активов. Если мы поймем, каким образом факторы влияют на возвраты, то сможем конструировать более робастные портфели.

Концепция факторного инвестирования выходит за рамки разбиения на классы активов и касается базовых факторных рисков, максимизируя преимущества диверсификации. Вместо того чтобы различать инвестиционные инструменты по таким признакам, как хедж-фонды или частная (непублично торгуемая) долевая ценная бумага, факторное инвестирование призвано выявлять четкие профили риска-возвратности на основе различий во влиянии фундаментальных рисков факторов. Наивный подход к среднedisперсному инвестированию подключает (искусственные) группировки как отдельные классы активов к среднedisперсному оптимизатору. Факторное инвестирование признает, что таким группировкам присущи многие из тех же самых факторных рисков, что и у традиционных классов активов. Выгоды от диверсификации могут быть завышены, как было обнаружено инвесторами во время последнего кризиса, когда корреляции между классами рисков активов возросли из-за влияния тех же самых базовых факторных рисков.

Иерархический паритет риска

Среднedisперсная оптимизация очень чувствительна к оценкам ожидаемых финансовых возвратов и ковариации этих возвратов. Инверсия матрицы ковариаций также становится сложнее и менее точной, когда возвраты сильно коррелируют, как это часто бывает на практике. Такой результат был назван проклятием Марковица: когда диверсификация важнее, потому что инвестиции коррелируют, обычные портфельные оптимизаторы, скорее всего, приведут к нестабильному решению. Преимущества диверсификации могут быть более чем компенсированы искаженными оценками. Как уже обсуждалось, даже наивные, равновзвешенные портфели нередко превосходят среднedisперсную и рисковую оптимизацию вне выборки.

Более робастные подходы встраивали дополнительные ограничения (Clarke et al., 2002), байесовы априорные распределения (Black и Litterman, 1992) либо использовали усадочные оценщики, делающие матрицу прецизионностей более стабильной численно (Ledoit и Wolf [2003], и все они доступны в библиотеке `sklearn` (<http://scikit-learn.org/stable/modules/generated/sklearn.covariance.LedoitWolf.html>)).

Иерархический паритет риска (Hierarchical risk parity, HRP), напротив, эффективно

задействует неконтролируемое машинное обучение, достигая превосходящих вне-выборочных размещений финансовых средств среди портфельных активов.

Недавнее нововведение в портфельной оптимизации эффективно задействует теорию графов и иерархическую кластеризацию, строя портфель за три шага (Lopez de Prado, 2015):

1. Определить метрику расстояния с тем, чтобы коррелированные активы находились близко друг к другу, и применить кластеризацию с одиночной связью с целью выявления иерархических связей.
2. Использовать иерархическую корреляционную структуру для квазидиагонализации матрицы ковариаций.
3. Применить нисходящее инверсно-дисперсное взвешивание с использованием рекурсивного бисекционного поиска, обрабатывая при конструировании портфеля кластеризованные активы как дополнения, а не заменители, и снижая число степеней свободы.

Родственный метод построения *иерархически кластеризованного портфеля* (hierarchical clustering portfolio, HCP) был представлен Раффино (Raffinot, 2016). В концептуальном плане сложные системы, такие как финансовые рынки, тяготеют к обладанию структурой и часто организованы иерархическим образом, в то время как взаимодействие между элементами в иерархии формирует динамику системы. Однако матрицы корреляций не имеют представления об иерархии, что дает весам варьироваться свободно и потенциально непреднамеренно.

Как метод HRP, так и метод HCP были протестированы финансовым холдингом J. P. Morgan на различных универсумах долевого ценного бумага. Метод HRP, в частности, произвел равные или превосходящие скорректированные на риск финансовые возвраты и коэффициенты Шарпа по сравнению с наивно диверсифицированными, максимально диверсифицированными и GMV-портфелями (с глобально минимальной дисперсией).

Мы представим реализацию на Python в *главе 12*.

Резюме

В этой главе мы рассмотрели важную тему портфельного менеджмента, которая предусматривает сочетание инвестиционных позиций с целевой задачей управления компромиссами между риском и возвратностью. Мы познакомились с библиотекой *pyfolio*, предназначенной для вычисления и визуализации ключевых метрических показателей риска и возвратности, а также для сравнения результативности различных алгоритмов.

Мы стали свидетелями того, насколько для оптимизации портфельных весов и максимизации преимуществ диверсификации важны точные предсказания. Мы также развели то, как МО способно обеспечивать более эффективное конструирование

портфеля, усваивая иерархические связи из матрицы ковариаций возвратов от портфельных активов.

Теперь мы перейдем ко второй части этой книги, которая посвящена использованию автоматически обучающихся моделей. Эти модели будут производить более точные предсказания за счет более эффективного использования более разнообразной информации, улавливая более сложные регулярности, чем более простые альфа-факторы, которые были наиболее заметными до сих пор.

Мы начнем с тренировки, тестирования и настройки линейных регрессионных и классификационных моделей с использованием перекрестного (скользящего) контроля с целью достижения робастной вневыборочной результативности. Мы также вставим эти модели в каркас определения и бэк-тестирования стратегий алгоритмической торговли, которые мы рассмотрели в последних двух главах.

6

Процесс машинного обучения

В этой главе мы начнем иллюстрировать способы применения широкого спектра контролируемых и неконтролируемых автоматически обучающихся моделей для алгоритмической торговли. Сначала мы объясним допущения и примеры использования каждой модели и потом продемонстрируем соответствующие применения, использующие различные библиотеки Python. Категории моделей включают:

- ◆ линейные модели для регрессии и классификации срезовых (кросс-секционных), рядовых и панельных данных¹;
- ◆ обобщенные аддитивные модели, включая нелинейные древесные модели, такие как *деревья решений*;
- ◆ ансамблевые модели, включая случайный лес и градиентно-бустинговые машины;
- ◆ неконтролируемые линейные и нелинейные методы снижения размерности и методы кластеризации;
- ◆ нейросетевые модели, включая рекуррентные и сверточные архитектуры;
- ◆ модели подкрепляемого обучения.

Мы применим эти модели к рыночным, фундаментальным и альтернативным источникам данных, представленным в предыдущих главах. Мы будем далее опираться на материал, рассмотренный до сих пор, показывая, как вставлять эти модели в стратегию алгоритмической торговли для генерирования или комбинирования альфа-факторов или для оптимизации процесса портфельного менеджмента и оценивания их результативности.

Существует несколько аспектов, которые объединяют многие из этих моделей и их применение. Настоящая глава охватывает эти общие аспекты, что позволяет сосре-

¹ Срезовые данные (cross-sectional data) — это данные, агрегированные по разным переменным за один и тот же момент времени, например в формате "страна — ВВП за 2018 год". Панельные данные (panel data), или продольные данные — это данные, агрегированные по разным переменным за многочисленные периоды времени, например в формате "страна — ВВП 2010, ВВП 2011, ..., ВВП 2018". — *Прим. перев.*

доточиться на применении конкретной модели в последующих главах. Они нацелены на усвоение функциональных связей из данных путем оптимизации целевой функции или функции потери. Они также содержат тесно связанные методы измерения модельной результативности.

Мы проведем различие между неконтролируемым (без учителя) и контролируемым (с учителем) автоматическим обучением и контролируемыми регрессионными и классификационными задачами, а также опишем примеры использования алгоритмической торговли. Для сравнения мы приведем пример использования контролируемого обучения для статистического вывода о связях между входными и выходными данными с целью предсказания будущих выходов из будущих входов. Мы также проиллюстрируем то, как ошибки предсказания обуславливаются смещением или дисперсией модели или высоким отношением шума к сигналу в данных. И самое главное, мы представим методы диагностики источников ошибок и повышения модельной результативности.

В этой главе мы рассмотрим следующие темы:

- ◆ как работает контролируемое и неконтролируемое автоматическое обучение с использованием данных;
- ◆ как применять рабочий поток МО;
- ◆ как формулировать функции потерь для регрессии и классификации;
- ◆ как тренировать и оценивать модели контролируемого обучения;
- ◆ как компромисс между смещением и дисперсией влияет на ошибки предсказания;
- ◆ как диагностировать и устранять ошибки предсказания;
- ◆ как тренировать модель с помощью перекрестного контроля, управляя компромиссом между смещением и дисперсией;
- ◆ как реализовывать перекрестный контроль с помощью библиотеки `sklearn`;
- ◆ почему характер финансовых данных требует других подходов к вневыборочному тестированию.

Если вы уже хорошо знакомы с МО, то можете без колебаний пропустить эту главу и погрузиться прямо в изучение того, как использовать линейные модели для порождения и комбинирования альфа-факторов в стратегии алгоритмической торговли.

Усвоение регулярностей из данных

Существует целый ряд определений МО, и все они вращаются вокруг автоматического обнаружения значимых регулярностей в данных. Два ярких примера включают:

- ◆ первопроходчик ИИ Артур Самуэльсон (Arthur Samuelson) дал определение термину МО в 1959 г. как подобласти информатики, которая дает компьютерам возможность учиться, не программируя их явным образом;
- ◆ Тони Митчелл (Toni Mitchell), один из нынешних лидеров в этой области, в 1998 г. дал более конкретное определение хорошо поставленной задачи автоматического обучения: компьютерная программа учится на опыте относительно задачи и меры результативности, которая сообщает об улучшении результативности задачи вместе с опытом.

Опыт предоставляется алгоритму в виде тренировочных данных. Принципиальное отличие от предыдущих попыток построения машин, которые решают задачи, состоит в том, что правила, используемые алгоритмом для принятия решения, не программируются или жестко кодируются извне, как было в случае экспертных систем, широко известных в 1980-х годах, а извлекаются из данных.

Ключевая сложность автоматического обучения заключается в выявлении регулярностей в тренировочных данных, которые приобретают смысл при обобщении усвоенного моделью на новые данные. Существует большое число потенциальных регулярностей, которые модель может выявлять, в то время как тренировочные данные составляют лишь выборку из более крупного набора явлений, которые алгоритму нужны для того, чтобы выполнять задачу в будущем. Бесконечное число функций, которые могут генерировать определенные выходы из заданного входа, делает поисковый процесс неразрешимым без ограничений на приемлемое множество функций.

Типы регулярностей, которые алгоритм способен усваивать, лимитированы размером его пространства гипотез, с одной стороны, и количеством информации, содержащейся в выборке данных, с другой. Размер пространства гипотез существенно варьируется от алгоритма к алгоритму. С одной стороны, это ограничение хорошо тем, что обеспечивает успешный поиск, а с другой — приводит к индуктивному смещению, когда алгоритм обобщает с тренировочной выборки на новые данные.

Следовательно, ключевой проблемой становится вопрос о том, как выбрать модель с пространством гипотез, достаточно большим для того, чтобы содержать решение задачи автоматического обучения, но достаточно малым для того, чтобы обеспечить надежное обобщение с учетом размера тренировочных данных. При наличии все более и более информативных данных модель с более крупным пространством гипотез будет успешной.

Теорема об отсутствии бесплатных обедов гласит, что универсального обучающегося алгоритма не существует. Вместо этого пространство гипотез ученика должно быть приспособлено к конкретной задаче с использованием предшествующих (априорных) знаний о предметной области задачи для того, чтобы поиск значимых регулярностей был успешным. В этой главе мы уделим пристальное внимание допущениям, принимаемым моделью о связях данных для конкретной задачи,

и подчеркнем важность сочетания этих допущений с эмпирическими наблюдениями, получаемыми в ходе разведывания данных. Процесс, необходимый для усвоения задачи, можно разделить на контролируемое, неконтролируемое и подкрепляемое обучение².

Контролируемое обучение

Контролируемое (само)обучение является наиболее часто используемым типом МО. Мы посвятим подавляющую часть глав настоящей книги исследованию различных применений моделей в этой категории. Термин "контролируемый" влечет за собой выходную переменную — результат, который направляет процесс самообучения, т. е. он побуждает алгоритм заучивать правильное решение усваиваемой задачи. Контролируемое обучение призвано обобщать функциональную связь между входными и выходными данными, которая усваивается из отдельных образцов, и применять ее к новым данным.

Выходная переменная также, в зависимости от области, взаимозаменяемо называется меткой, целью, результатом, исходом, эндогенной или левосторонней переменной. Мы будем использовать y_i для наблюдений $i = 1, \dots, N$ или y в векторном обозначении. Некоторые задачи представлены несколькими результатами, также именуемыми *многометочными задачами*. Входные данные для задачи контролируемого обучения также называются признаками, экзогенными и правосторонними переменными. Вектор признаков обозначается через x_i для наблюдений $i = 1, \dots, N$ или через x в матричном обозначении.

Решение задачи контролируемого обучения является функцией \hat{f} , которая представляет то, что модель усвоила о связи между входом и выходом из выборки и аппроксимирует истинную связь, представленную как $y \approx \hat{f}(x)$. Эта функция может использоваться для выведения статистических ассоциаций или потенциально даже причинно-следственных связей между интересующими переменными за пределами выборки, либо для предсказания выходов для новых входных данных.

Обе цели сталкиваются с важным компромиссом: более сложные модели имеют больше *движущихся частей*, которые способны представлять более нюансированные связи, но их также труднее обследовать. Кроме того, они могут излишне плотно прилегать к данным и усваивать случайный шум, характерный для тренировочного образца, в отличие от систематического сигнала, который представляет собой обобщенную регулярность связи между входом и выходом. С другой стороны, слишком простые модели будут пропускать сигналы и давать смещенные результаты. Этот компромисс известен как *компромисс между смещением и дисперсией* в контролируемом обучении, но концептуально он также относится и к другим

² В данном переводе принят зарубежный подход — вместо терминов "обучение с учителем", "обучение без учителя" и "обучение с подкреплением" используются термины "контролируемое", "неконтролируемое" и "подкрепляемое обучение". — *Прим. перев.*

формам МО, где чрезмерно сложные модели могут плохо работать за пределами тренировочных данных.

Неконтролируемое обучение

Во время решения задачи неконтролируемого обучения мы только наблюдаем за признаками и не имеем предварительных статистических измерений результата. Вместо предсказания будущих результатов или вывода связей между переменными задача состоит в том, чтобы отыскивать в данных структуру без какой-либо результирующей информации, которая будет руководить поисковым процессом.

Зачастую неконтролируемые алгоритмы призваны усваивать новое представление входных данных, которое будет полезным при использовании в некоторых других задачах. Сюда входит выработка меток, которые выявляют общие черты между наблюдениями, либо итоговое описание, которое улавливает релевантную информацию, при этом требуя точек данных или признаков. Неконтролируемые обучающиеся алгоритмы также отличаются от контролируемых обучающихся алгоритмов допущениями о природе структуры, которую они стремятся обнаружить.

Применения

Существует несколько полезных применений неконтролируемого обучения, которые могут быть задействованы в алгоритмической торговле, включая следующие:

- ◆ группировка ценных бумаг с аналогичными характеристиками риска и возвратности (см. иерархический паритет риска в этой главе, где рассматривается портфельная оптимизация);
- ◆ нахождение малого числа рисковых факторов, стимулирующих результативность гораздо большего числа ценных бумаг;
- ◆ выявление торговых и ценовых регулярностей, которые различаются систематически и могут создавать более высокие риски;
- ◆ выявление латентных тем в массиве документов (например, в стенограммах телеконференций о корпоративных зарплатах), которые составляют наиболее важные аспекты этих документов.

На высоком уровне эти применения опираются на методы выявления кластеров и методы снижения размерности данных.

Кластерные алгоритмы

Кластерные алгоритмы используют меру сходства для выявления наблюдений или атрибутов данных, содержащих схожую информацию. Они резюмируют набор данных, назначая большое число точек данных меньшему числу кластеров, связывая членов кластера друг с другом теснее, чем с членами других кластеров.

Кластерные алгоритмы в основном различаются по типу кластеров, которые они будут создавать, из чего вытекают разные допущения о процессе генерации данных, которые перечислены далее.

- ◆ *Кластеризация методом k средних* — точки данных принадлежат одному из k равновеликих кластеров, которые принимают эллиптическую форму.
- ◆ *Модели гауссовых смесей* — точки данных были сгенерированы любым видом многомерных нормальных распределений.
- ◆ *Плотностные кластеры* — кластеры имеют произвольную форму и определяются только наличием минимального числа близлежащих точек данных.
- ◆ *Иерархические кластеры* — точки данных принадлежат разным надмножествам групп, которые формируются путем последовательного слияния более мелких кластеров.

Снижение размерности

В результате снижения размерности создаются новые данные, которые улавливают наиболее важную информацию, содержащуюся в исходных данных. Вместо группировки существующих данных в кластеры эти алгоритмы преобразуют существующие данные в новый набор данных, который для представления исходной информации использует значительно меньше признаков или наблюдений.

Указанные алгоритмы различаются по природе нового набора данных, который они будут создавать, как показано в списке далее.

- ◆ *Анализ главных компонент* (principal component analysis, PCA) находит линейное преобразование, которое улавливает большую часть дисперсии в существующем наборе данных.
- ◆ *Усвоение проекций в топологическом многообразии* (manifold learning)³ выявляет нелинейное преобразование, которое создает низкоразмерное представление данных.
- ◆ *Автокодировщики* (autoencoder) используют нейронные сети для сжатия данных нелинейно с минимальной потерей информации.

Мы углубимся в линейные, нелинейные и нейросетевые модели неконтролируемого обучения в нескольких следующих главах, включая важные применения *обработки естественного языка* (ЕЯ) в форме тематического моделирования и извлечения признаков на основе модели векторных вложений слов Word2vec.

³ Усвоение проекций в топологическом многообразии (manifold learning) — это попытка раскрыть структуру топологического многообразия (manifold), или многосвязной области, в наборе данных, где топологическое многообразие — это пространство, локально сходное с евклидовым. Указанный метод представляет собой нелинейный метод снижения размерности. — Прим. перев.

Подкрепляемое обучение

Подкрепляемое обучение (Reinforcement learning, RL) — это третий тип МО. Оно стремится выбирать действие, которое дает наибольшее вознаграждение с учетом набора входных данных, описывающих контекст или среду. Оно является одновременно динамичным и интерактивным: поток положительных и отрицательных вознаграждений влияет на самообучение алгоритма, а действия, предпринимаемые сейчас, могут влиять как на окружающую среду, так и на будущие вознаграждения.

Компромисс между следованием линии поведения, которая была усвоена, как приносящая определенное вознаграждение, и разведыванием новых поступков, которые могут увеличить вознаграждение в будущем, порождает метод проб и ошибок⁴. Подкрепляемое обучение оптимизирует самообучение агента с использованием теории динамических систем и, в частности, оптимального управления в рамках марковских процессов принятия решения в условиях неполной информации.

Подкрепляемое обучение отличается от контролируемого обучения, где имеющиеся тренировочные данные выкладывают перед алгоритмом контекст и правильное решение. Подкрепляемое обучение подходит для интерактивных конфигураций, где результаты становятся доступными только с течением времени, и самообучение должно проходить в онлайн-овом или непрерывном режиме по мере приобретения агентом нового опыта. Вместе с тем некоторые самые заметные достижения в области *искусственного интеллекта* (ИИ) предусматривают подкрепление, которое для аппроксимирования функциональных связей между поступками, средой и будущими вознаграждениями задействует глубокое обучение. Подкрепляемое обучение также отличается от неконтролируемого обучения, поскольку имеется обратная связь о последствиях, хотя и с задержкой.

Подкрепляемое обучение подходит для алгоритмической торговли больше всего, потому что концепция агента, максимизирующего финансовый возврат в неопределенной динамической среде, имеет много общего с инвестором или торговой стратегией, которая взаимодействует с финансовыми рынками. Этот подход был успешно применен к игровым агентам, в особенности к игре го, а также к сложным видеоиграм. Он также используется в робототехнике — например, в самодвижущихся автомобилях — или для персонализации служб, таких как рекомендации веб-сайта на основе взаимодействия с пользователем. В *главе 21* мы познакомим вас с подходами к построению стратегии алгоритмической торговли на основе подкрепляемого самообучения.

⁴ Разведывание (exploration) — это поиск новых идей или новых стратегий. Используется как антитеза для эксплуатации (exploitation), т. е. использованию существующих идей и стратегий, которые оказались успешными в прошлом. Интеллектуальная работа включает надлежащий баланс между разведыванием, сопряженным с риском пустой траты ресурсов, и эксплуатацией проверенных временем решений. — *Прим. перев.*

Рабочий поток машинного обучения

Разработка технического решения с использованием МО для стратегии алгоритмической торговли требует системного подхода, максимизирующего шансы на успех при экономии ресурсов. Также очень важно сделать процесс прозрачным и воспроизводимым с целью обеспечения сотрудничества, сопровождения и последующих усовершенствований.

На рис. 6.1 показаны ключевые шаги, начиная с определения задачи и заканчивая развертыванием предсказательного технического решения.

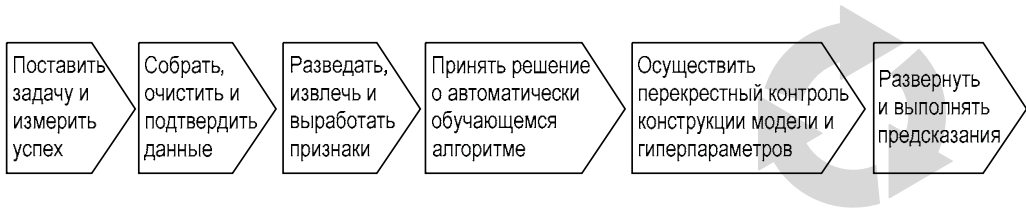


Рис. 6.1. Схема рабочего потока машинного обучения

Указанный процесс является итеративным на протяжении всей последовательности, и усилия, необходимые на разных этапах, будут варьироваться в зависимости от проекта, но этот процесс обычно должен охватывать следующие шаги:

1. Очертить рамки задачи, дать определение целей и метрик успеха.
2. Привлечь, очистить и подтвердить данные.
3. Разобраться в данных и сгенерировать информативные признаки.
4. Выбрать один или несколько автоматически обучающихся алгоритмов, подходящих для привлеченных данных.
5. Натренировать, протестировать и настроить модели.
6. Применить модель для решения исходной задачи.

В последующих разделах мы пройдемся по этим шагам, используя простой пример, иллюстрирующий некоторые из этих ключевых моментов.

Простое пошаговое руководство — *k* ближайших соседей

В блокноте `machine_learning_workflow.ipynb` в папке этой главы репозитория GitHub содержится несколько примеров, иллюстрирующих рабочий поток проекта машинного обучения с использованием набора данных с ценами на жилую недвижимость.

Мы воспользуемся довольно простым алгоритмом *k* ближайших соседей (*k*-nearest neighbors, KNN), который позволяет решать регрессионные и классификационные задачи.

Стандартная реализация в библиотеке `sklearn` определяет k ближайших точек данных (на основе евклидова расстояния) для предсказания. Она предсказывает наиболее частый класс среди соседей либо средний результат соответственно в случае классификации либо регрессии.

Очертить рамки задачи — цели и метрики

Отправной точкой для любого упражнения по отработке методов машинного обучения является принципиальный пример использования, который он призван решать. Иногда этой целью будет статистический вывод для выявления ассоциации между переменными или даже причинно-следственной связи. Однако чаще всего целью будет прямое предсказание результата, выдающего торговый сигнал.

Как в выводе, так и предсказании используются метрики для оценивания того, насколько хорошо модель достигает своей цели. Мы сосредоточимся на общих целевых функциях и соответствующих метриках ошибок для предсказательных моделей, которые могут различаться по типу переменной на выходе: непрерывные выходные переменные приводят к регрессионной задаче, категориальные переменные — к классификационной задаче, а частный случай упорядоченных категориальных переменных ведет к задаче ранжирования.

Задача может составлять эффективное сочетание нескольких альфа-факторов и быть оформлена как регрессионная задача, которая призвана предсказывать финансовые возвраты, как бинарно-классификационная задача, призванная предсказывать направление будущего движения цены, или как мультиклассовая задача, назначающая акции классам с различной результативностью. В следующем далее разделе мы представим эти цели и рассмотрим способы измерения и интерпретации связанных с ними метрик ошибок.

Предсказание против статистического вывода

Функциональная связь, производимая обучающимся под контролем алгоритмом, может использоваться для выведения заключения, т. е. для проникновения в сущность того, как генерируются результаты, либо для предсказания, т. е. для генерирования точных выходных оценок (обозначаемых через \hat{y}) для неизвестных или будущих входов (обозначаемых через x).

Применительно к алгоритмической торговле выведение заключения может использоваться для оценивания причинно-следственной или статистической зависимости возвратов от актива на рисковом факторе, тогда как предсказание может использоваться для прогнозирования рискового фактора. Совмещение обоих подходов может дать предсказание цены актива, которое в свою очередь может быть транслировано в торговый сигнал.

Статистический вывод связан с выведением заключений из выборки данных относительно параметров базового вероятностного распределения или популяции (также именуемой генеральной совокупностью). Потенциальные заключения охва-

тывают проверки статистических гипотез о характеристиках распределения индивидуальной переменной или о существовании либо силе числовых связей между переменными. Они также охватывают точечные или интервальные оценки статистических показателей.

Выведение заключения изначально зависит от допущений о процессе, который генерирует данные. Мы вернемся к этим допущениям и инструментам, которые используются для статистического вывода, при рассмотрении линейных моделей, где они хорошо отработаны. Более сложные модели делают меньше допущений о структурной связи между входом и выходом и вместо этого более открыто подходят к задаче аппроксимации функций, при этом рассматривая процесс генерирования данных как черный ящик. Эти модели, включая деревья решений, ансамблевые модели и нейронные сети, ориентированы на предсказательные задачи и часто показывают повышенную результативность при их использовании для таких задач. Тем не менее случайные леса недавно получили вычислительный каркас статистического вывода, который мы представим позже.

Причинно-следственный вывод

Причинно-следственный вывод призван выявлять причинно-следственные связи, такие, что из определенных входных значений вытекают определенные выходы, например конкретное созвездие макропеременных, побуждающих цену данного актива двигаться определенным образом при допущении, что все остальные переменные остаются постоянными.

Статистический вывод о связях между двумя или более переменными производит меры корреляции, которые могут интерпретироваться как причинно-следственная связь только при соблюдении нескольких других условий, например при исключении альтернативных объяснений или обратной причинно-следственной связи. Соблюдение этих условий требует экспериментальной ситуации, в которой все релевантные переменные, представляющие интерес, могут полностью контролироваться при изолировании причинно-следственных связей. В качестве альтернативы квазиэкспериментальные ситуации рандомизированно подвергают единицы наблюдений изменениям во входах, исключая возможность того, чтобы другие наблюдаемые или ненаблюдаемые признаки были ответственны за наблюдаемые эффекты изменения в среде.

Эти условия редко удовлетворяются, поэтому к статистическим выводам необходимо относиться с осторожностью. То же самое относится и к результативности предсказательных моделей, которые тоже опираются на статистическую ассоциативную связь между признаками и выходами, способную изменяться вместе с другими факторами, не являющимися частью модели.

Непараметрическая природа модели k ближайших соседей не поддается статистическому выводу, поэтому мы отложим этот шаг в рабочем потоке до следующей главы, где мы столкнемся с линейными моделями.

Регрессионные задачи

Регрессионные задачи призваны предсказывать непрерывную переменную. *Корень квадратный из среднеквадратической ошибки* (root-mean-square error, RMSE) является самой популярной функцией потери и метрикой ошибки, не в последнюю очередь потому, что указанная функция является дифференцируемой. Потеря является симметричной, но более крупные ошибки при расчете весят больше. Использование квадратного корня имеет преимущество измерения ошибки в единицах целевой переменной. Та же метрика в сочетании с RMSE, представленная *корнем квадратным из среднеквадратической логарифмированной ошибки* (RMSE log of the error, RMSLE), подходит, когда цель подвергается экспоненциальному росту благодаря ее асимметричному штрафу, который взвешивает отрицательные ошибки меньше, чем положительные. Кроме того, сначала можно взять логарифм цели, а затем применить RMSE, как мы сделаем в примере позже в этом разделе.

Средняя абсолютная ошибка (mean absolute error, MAE) и *медианная абсолютная ошибка* (median absolute error, MedAE) являются симметричными, но не придают более крупный вес более крупным ошибкам. Медианная абсолютная ошибка робастна к выбросам.

Показатель объясненной⁵ дисперсии вычисляет долю целевой дисперсии, которая учитывается моделью и варьируется от 0 до 1. Показатель R^2 , или коэффициент детерминации, дает тот же результат, что и среднее значение остатков, равное 0, но может отличаться в противном случае. В частности, он может быть отрицательным при расчете на вневыборочных данных (или для линейной регрессии без пересечения⁶).

В табл. 6.1 определены формулы, используемые для расчета, и соответствующие им функции библиотеки `sklearn`, которые можно импортировать из модуля `metrics`. Параметр `scoring` используется в сочетании с автоматизированными тренировочно-тестовыми функциями (такими как перекрестно-проверочная отметка `cross_val_score` и класс поиска в решетке параметров `GridSearchCV`), которые мы представим позже в этом разделе и которые проиллюстрированы в прилагаемом блокноте.

На рис. 6.2 показаны различные метрики ошибок для регрессии цен на жилую недвижимость, представленной в блокноте.

Функции библиотеки `sklearn` также поддерживают многоуровневое оценивание, т. е. закрепление многочисленных выходных значений за одним-единственным наблюдением; подробности см. в документации, ссылка на которую приведена в репозитории GitHub по адресу <https://github.com/PacktPublishing/Hands-On-Machine-Learning-for-Algorithmic-Trading/tree/master/Chapter06>.

⁵ Под термином "объяснить" здесь имеется в виду "снизить" дисперсию (остатков модели). — Прим. перев.

⁶ Пересечение (intercept) — коэффициент сдвига, который показывает, каким будет \hat{y} в случае, если все используемые в модели факторы будут равны 0; под ним подразумевается зависимость от других неописанных в модели факторов. Представляет собой сдвиг по оси y и соответствует точке на этой оси, где прямая регрессии пересекает эту ось. — Прим. перев.

Таблица 6.1. Формулы и функции модуля `metrics` библиотеки `sklearn`

| Имя | Формула | sklearn | Мерный параметр |
|---|---|--------------------------|----------------------------|
| Среднеквадратическая ошибка | $\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$ | mean_squared_error | neg_mean_squared_error |
| Среднеквадратическая логарифмированная ошибка | $\frac{1}{N} \sum_{i=1}^N (\ln(1 + y_i) - \ln(1 + \hat{y}_i))^2$ | mean_squared_log_error | neg_mean_squared_log_error |
| Средняя абсолютная ошибка | $\frac{1}{N} \sum_{i=1}^N y_i - \hat{y}_i $ | mean_absolute_error | neg_mean_absolute_error |
| Медианная абсолютная ошибка | $\text{median}(y_1 - \hat{y}_1 , \dots, y_n - \hat{y}_n)$ | median_absolute_error | neg_median_absolute_error |
| Объясненная дисперсия | $1 - \frac{\text{var}(y - \hat{y})}{\text{var}(y)}$ | explained_variance_score | explained_variance |
| Показатель R^2 | $1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$ | r2_score | r2 |

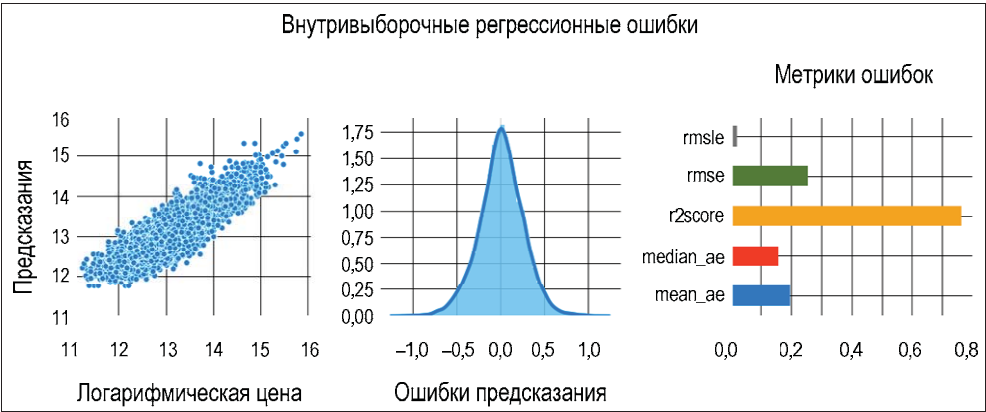


Рис. 6.2. Внутривыборочные регрессионные ошибки

Классификационные задачи

Классификационные задачи имеют выходные категориальные переменные. Большинство предсказывающих переменных, или предикторов, выдают отметку (score), обозначающую принадлежность наблюдения к определенному классу. На втором шаге эти отметки затем транслируются в фактические предсказания.

В бинарном случае, когда мы будем размечать классы как положительные и отрицательные, эта отметка, как правило, варьируется между нулем и единицей или же соответствующим образом нормализуется. После того как отметки конвертированы в предсказания 0-1, может иметься четыре результата, потому что каждый из двух существующих классов может быть предсказан правильно либо неправильно. Когда число классов больше двух, таких случаев может быть больше, если вы проводите различие между несколькими потенциальными ошибками.

Все метрики ошибок вычисляются из разбивки предсказаний на четыре поля матрицы несоответствий размером 2×2 , которая ассоциирует фактические и предсказанные классы. Метрики, представленные на рис. 6.3, такие как точность⁷, оценивают модель для заданного порога⁸.

| | | | | | |
|--------------|-------------|---------------------------|-------------------------|--|--|
| | | Факт (истина) | | Точность = $\frac{\text{Число правильных предсказаний}}{\text{Число случаев}} = \frac{TP + TN}{TP + FP + TN + FN}$ | |
| | | Утверждение Отрицание | | | |
| Предсказание | Утверждение | Истинное утверждение (TP) | Ложное утверждение (FP) | Частота истинных утверждений (чувствительность, полнота) | $\frac{\text{Число правильных положительных предсказаний}}{\text{Число положительных случаев}} = \frac{TP}{TP + FN}$ |
| | Отрицание | Ложное отрицание (FN) | Истинное отрицание (TN) | Частота ложных отрицаний (частота непопаданий) | $1 - \text{Частота истинных утверждений}$ |
| | | | | Частота истинных отрицаний (специфичность) | $\frac{\text{Число правильных отрицательных предсказаний}}{\text{Число отрицательных случаев}} = \frac{TN}{TN + FP}$ |
| | | | | Частота ложных утверждений (выпадение) | $1 - \text{Частота истинных отрицаний}$ |

Рис. 6.3. Метрики ошибок

⁷ Точность (assuracy) модели — это мера статистического смещения и описывает систематические ошибки; иными словами, это близость результатов измерений истинному значению. См. https://en.wikipedia.org/wiki/Accuracy_and_precision. — Прим. перев.

⁸ Для справки: TP (true positive) — истинное утверждение, TN (true negative) — истинное отрицание, FP (false positive) — ложное утверждение и FN (false negative) — ложное отрицание. — Прим. перев.

Классификатор не обязательно должен выдавать откалиброванные вероятности. Как раз наоборот, он должен производить отметки, которые являются относительными друг к другу при различении положительных и отрицательных случаев. Отсюда следует, что порог является величиной принятия решения, которая может и должна быть оптимизирована с учетом издержек и преимуществ от правильных и неправильных предсказаний. Более низкий порог приводит к более положительным предсказаниям с потенциально растущей частотой ложных утверждений, а для более высокого порога, скорее всего, будет верно обратное.

Рабочие характеристики приемника и площадь под кривой

Кривая рабочих характеристик приемника (receiver operating characteristics, ROC) позволяет визуализировать, организовывать и отбирать классификаторы на основе их результативности. Она вычисляет все комбинации *частот истинных утверждений* (TP rate) и *частот ложных утверждений* (FP rate), которые возникают в результате порождения предсказаний с использованием любой предсказанной отметки в качестве порога. Затем она строит эти пары в квадрате с длиной стороны, равной единице.

Классификатор, который делает случайные предсказания (с учетом несбалансированности классов), в среднем дает эквивалентную частоту истинных утверждений и частоту ложных утверждений, вследствие чего их комбинации будут лежать по диагонали, таким образом становясь эталонным случаем. Поскольку классификатор со слабой результативностью извлек бы выгоду из переназначения меток предсказаний, этот эталон также становится минимальным.

Площадь под кривой (area under the curve, AUC) определяется как площадь под кривой ROC, которая варьируется от 0,5 до максимума, равного 1. Она представляет собой сводную меру того, насколько хорошо отметки классификатора способны ранжировать точки данных по их классовой принадлежности. Более конкретно, мера AUC классификатора имеет важное статистическое свойство представлять вероятность, что классификатор будет ранжировать случайно выбранный положительный экземпляр выше, чем случайно выбранный отрицательный экземпляр, что эквивалентно статистической проверке знаковых рангов Уилкоксона. В дополнение к этому, мера AUC имеет преимущество в том, что она не чувствительна к несбалансированностям классов.

Кривые прецизионности-полноты

Когда предсказания для одного из классов представляют особый интерес, кривые прецизионности и полноты визуализируют компромисс между этими метриками ошибок для разных пороговых значений. Обе меры оценивают качество предсказаний для отдельного класса. В следующем списке указано, как они применяются к положительному классу.

- ♦ *Полнота* служит мерой доли фактических положительных членов класса, которые классификатор предсказывает как положительные для данного порога. Ука-

занная мера берет свое начало в технологии информационного поиска, где она служит мерой доли соответствующих документов, успешно идентифицированных поисковым алгоритмом.

- ♦ *Прецизионность*, напротив, служит мерой доли положительных предсказаний, которые являются правильными⁹.

Полнота, как правило, увеличивается с более низким порогом, но прецизионность может уменьшаться. Кривые прецизионности-полноты визуализируют достижимые комбинации и позволяют оптимизировать порог с учетом издержек и преимуществ пропуска многих релевантных случаев или порождения предсказаний пониженного качества.

Показатель F1 — это гармоническое среднее прецизионности и полноты для заданного порога. Он может использоваться для численной оптимизации порога с учетом относительных весов, которые эти две метрики должны принимать.

На рис. 6.4 представлены кривая ROC и соответствующий показатель — отметка AUC, наряду с кривой прецизионности-полноты и показателем F1, который, используя равные веса для прецизионности и полноты, дает оптимальный порог 0,37. Диаграмма взята из прилагаемого блокнота, где можно найти исходный код классификатора k ближайших соседей, который работает на бинаризованных ценах на жилую недвижимость:

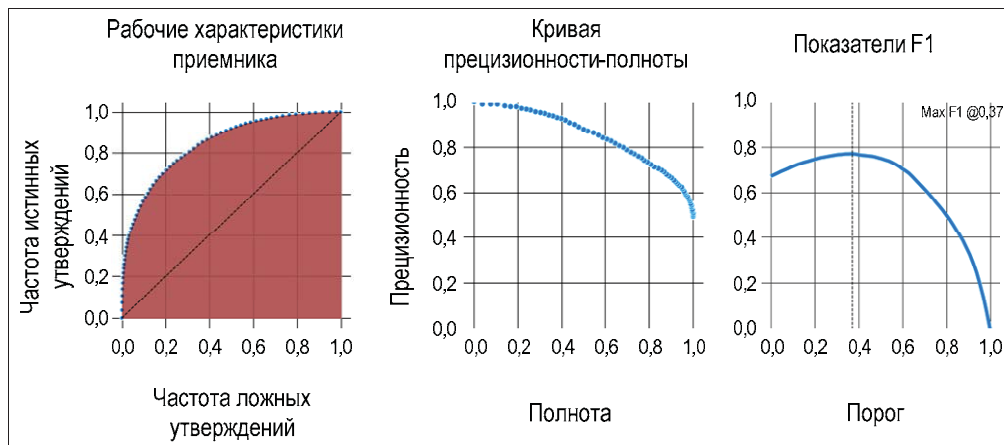


Рис. 6.4. Определение оптимального порога с использованием кривой ROC, кривой прецизионности-полноты и кривой показателей F1

⁹ Прецизионность (precision) результатов измерений — это мера статистической изменчивости, она описывает случайные ошибки или степень близости друг к другу независимых результатов измерений, полученных в конкретных установленных условиях.

См. https://en.wikipedia.org/wiki/Accuracy_and_precision. — Прим. перев.

Собрать и подготовить данные

Мы уже рассмотрели важные аспекты привлечения рыночных, фундаментальных и альтернативных данных и продолжим работу с разнообразными примерами этих источников, иллюстрируя применение различных моделей.

В дополнение к рыночным и фундаментальным данным, к которым мы будем обращаться через платформу Quantopian, мы также будем получать и преобразовывать текстовые данные, когда займемся разведкой технологии обработки естественного языка, и снимковых данных, когда мы обратимся к обработке и распознаванию изображений. Помимо получения, очистки и подтверждения данных на их связь с торговыми данными, обычно доступными в формате временного ряда, важно хранить их в формате, обеспечивающем быстрый доступ для беглого разведывательного анализа и итеративной обработки. Мы рекомендовали форматы HDF и parquet. Для более крупных объемов данных наилучшим решением является вычислительный каркас распределенной обработки Apache Spark.

Разведать, извлечь и выработать признаки

Понимание распределения отдельных величин и связей между результатами и признаками является основой для подбора подходящего алгоритма. Подбор обычно начинается с визуализаций, таких как диаграммы рассеяния, как показано в прилагаемом блокноте (и на рис. 6.5), но также предусматривает числовые оценивания, начиная с линейных метрических показателей, таких как корреляция, до нелинейных статистических показателей, таких как коэффициент ранговой корреляции Спирмена, с которым мы столкнулись при введении информационного коэффициента. Он также охватывает информационно-теоретические меры, такие как взаимная информация, как показано в следующем далее подразделе (см. рис. 6.5).

Систематический разведывательный анализ является основой того, что нередко выступает единственным самым важным компонентом успешной предсказательной

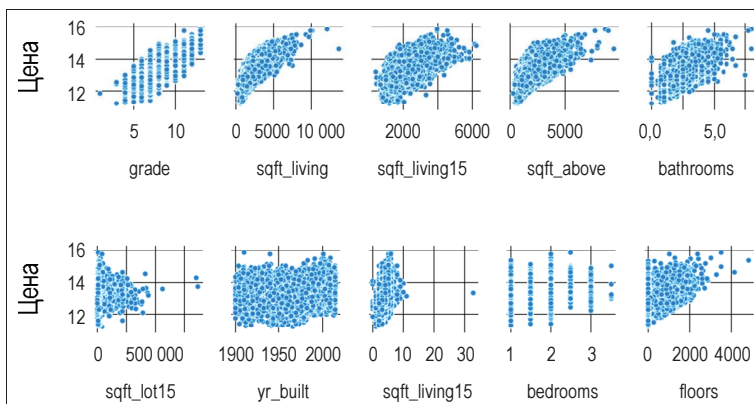


Рис. 6.5. Графики рассеяния

модели: выработка признаков¹⁰, извлекающих информацию, содержащуюся в данных, но которые не обязательно доступны алгоритму в их сырой форме. Выработка признаков выигрывает от наличия компетентных знаний предметной области, применения статистики и теории информации, а также изобретательности.

Выработка признаков опирается на изобретательный подбор преобразований данных, которые эффективно выделяют систематическую связь между входными и выходными данными. Существует много вариантов выработки признаков, которые включают обнаружение и устранение выбросов, функциональные преобразования и сочетание нескольких величин, включая неконтролируемое обучение. Мы будем приводить его примеры на всем протяжении, но подчеркнем, что этот функционал лучше всего осваивается на опыте. Онлайн-сообщество исследователей данных Kaggle (<https://www.kaggle.com/>) является отличным местом, где можно научиться у других исследователей данных, которые делятся своим опытом с сообществом Kaggle.

Использование теории информации для оценивания признаков

Взаимная информация (mutual information, MI) между признаком и результатом является мерой взаимозависимости между двумя величинами. Она расширяет понятие корреляции до нелинейных связей. Более конкретно, она квантифицирует информацию, полученную об одной случайной величине через другую случайную величину.

Понятие взаимной информации тесно связано с фундаментальным понятием энтропии случайной величины. Энтропия квантифицирует количество информации, содержащейся в случайной величине. Формально взаимная информация $I(X, Y)$ двух случайных величин X и Y определяется следующим образом:

$$I(X, Y) = \int \int p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right)$$

Библиотека `sklearn` реализует регрессию взаимной информации в функции `feature_selection.mutual_info_regression`, которая вычисляет взаимную информацию между всеми признаками и непрерывным результатом, отбирая признаки, которые с максимальной вероятностью содержат предсказательную информацию. Существует также классификационная версия данной функции (подробнее см. документацию). Блокнот `mutual_information.ipynb` содержит пример использования к финансовым данным, которые мы создали в *главе 4*.

¹⁰ Выработка признаков (feature engineering), или конструирование признаков — это процесс извлечения признаков из сырых данных и преобразования их в форматы, подходящие для автоматически обучающейся модели, используя компетентные знания предметной области данных. — *Прим. перев.*

Отобрать автоматически обучающийся алгоритм

В оставшейся части этой книги будет представлено несколько семейств моделей, начиная с линейных моделей, которые делают довольно сильные допущения о природе функциональной связи между входными и выходными переменными, вплоть до глубоких нейронных сетей, которые делают очень мало допущений. Как упоминалось во вступительном разделе, успех автоматического обучения зависит от числа принятых допущений: чем оно меньше, тем больше требуется данных со значительной информацией о связи.

По мере введения этих моделей мы будем излагать основные допущения и способы их проверки в применимых случаях.

Сконструировать и настроить модели

Процесс МО включает шаги диагностики и управления сложностью модели на основе оценок модельной ошибки обобщения. Несмещенная оценка требует статистически обоснованной и эффективной процедуры, а также метрик ошибок, которые согласуются с типом выходной переменной, что также определяет, с какой задачей мы имеем дело: регрессией, классификацией или ранжированием.

Компромисс между смещением и дисперсией

Ошибки, совершаемые автоматически обучающейся моделью во время предсказания результатов для новых входных данных, могут быть разбиты на снижаемую и неснижаемую части. Неснижаемая часть обусловлена случайной вариацией (шумом) в данных, которая не поддается измерению, как в случае с релевантными, но отсутствующими переменными или естественной вариацией. Снижаемая часть ошибки обобщения, в свою очередь, может быть разбита на смещение и дисперсию. Обе они обусловлены различиями между истинной функциональной связью и допущениями, сделанными автоматически обучающимся алгоритмом, как подробно описано в следующем списке.

- ◆ *Ошибка из-за смещения.* Гипотеза является слишком простой для того, чтобы улавливать сложность истинной функциональной связи. В результате всякий раз, когда модель пытается усвоить истинную функцию, она совершает систематические ошибки, и предсказания в среднем будут схожим образом смещены. Такая ситуация также называется *недоподгонкой*.
- ◆ *Ошибка из-за дисперсии.* Алгоритм является слишком сложным с точки зрения истинной связи. Вместо того чтобы улавливать истинные связи, он достигает перепоподгонки, т. е. чрезмерно плотно прилегает к данным, и извлекает регулярности из шума. В результате из каждого образца он заучивает совершенно разные функциональные связи, и вневыборочные предсказания будут варьироваться в широком диапазоне.

Недоподгонка против переподгонки

На рис. 6.6 проиллюстрировано явление переподгонки: косинусная функция аппроксимирована с использованием все более сложных многочленов, а внутривыборочные ошибки¹¹ измерены. Более конкретно, мы извлекаем 10 случайных выборок с некоторым добавленным шумом ($n = 30$), усваивая многочлен варьирующейся сложности (исходный код см. в прилагаемом блокноте). Модель всякий раз предсказывает новые точки данных, и для этих предсказаний мы улавливаем среднеквадратическую ошибку, а также стандартное отклонение этих ошибок.

На рис. 6.6 слева представлен многочлен первой степени; прямая линия явно не вписывается в истинную функцию, что свидетельствует о недоподгонке. Однако эта оценочная линия не будет резко отличаться от одной выборки к другой, извлеченной из истинной функции. В центре рис. 6.6 показано, что многочлен 5-й степени достаточно хорошо аппроксимирует истинную связь в промежутке $[0; 1]$. С другой стороны, многочлен 15-й степени почти идеально укладывается в небольшую выборку, но дает слабую оценку истинной связи: он достигает переподгонки к случайной вариации выборочных точек данных, и усвоенная функция будет сильно варьироваться с каждой извлеченной выборкой.

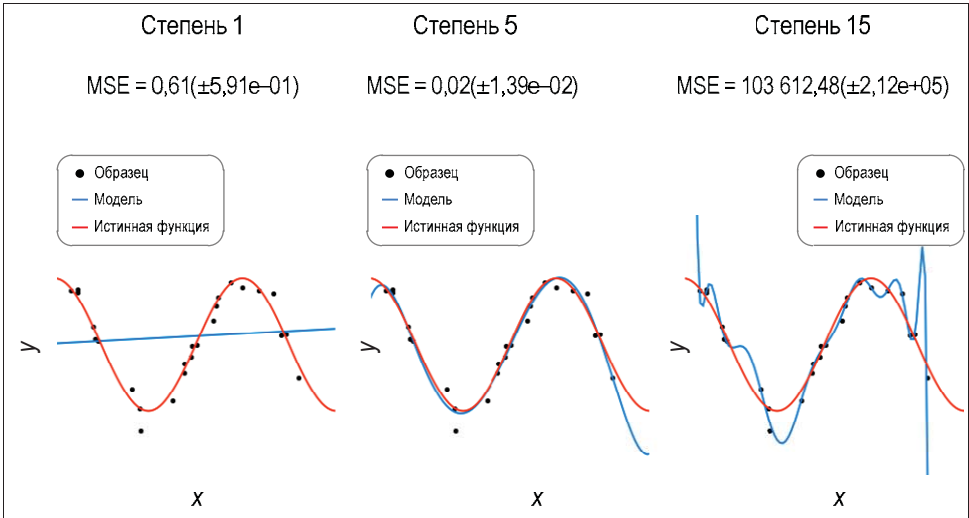


Рис. 6.6. Иллюстрация недоподгонки и переподгонки на примере многочленов разных степеней

¹¹ Суть переподгонки в том, что невольно извлекается некая остаточная дисперсия (т. е. шум), как если бы эта дисперсия представляла собой базовую структуру модели (Burnham K. P., Anderson D. R.). Недоподгонка относится к ситуациям, когда модель не способна улавливать структуру данных. — Прим. перев.

Управление компромиссом

Давайте подробнее проиллюстрируем влияние переподгонки в сопоставлении с недоподгонкой, пытаясь усвоить аппроксимацию ряда Тейлора для косинусной функции 9-й степени с некоторым дополнительным шумом. На рис. 6.7 показано, что мы извлекаем случайные образцы истинной функции и вписываем многочлены, которые переподогнаны, недоподогнаны и обеспечивают приближенно правильную степень гибкости. Затем мы выполняем предсказание вне выборки и измеряем ошибку RMSE.

Высокое смещение, но низкая дисперсия многочлена третьей степени сравнима с низким смещением, но чрезвычайно высокой дисперсией различных ошибок предсказания, видимых на первой панели. На рис. 6.7 слева показано распределение ошибок, возникающих в результате вычитания значений истинной функции. Недоподогнанный случай с прямой линией производит слабую внутривыборочную подгонку и существенно отдалится от цели вне выборки. Переподогнанная модель показывает наилучшую подгонку внутри выборки с наименьшей дисперсией ошибок, но ценой является большая дисперсия вне выборки. Надлежащая модель, которая сочетается с функциональной формой истинной модели, работает безусловно лучше вне выборки.

На рис. 6.7 справа для демонстрации того, как на практике выглядят различные типы подгонки, показаны не ошибки, а фактические предсказания.

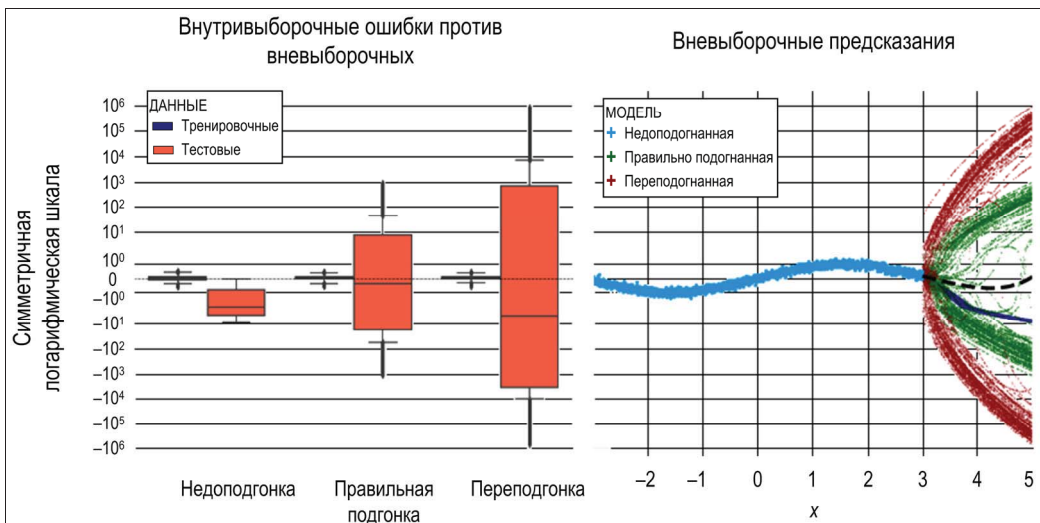


Рис. 6.7. Компромисс между смещением и дисперсией: недоподгонка против переподгонки

Кривые усвоения

Графики кривых усвоения, или графики (само)обучения, отображают эволюцию ошибок на тренировочном наборе и ошибок на тестовом наборе относительно размера набора данных, используемого для усвоения функциональной связи. Этот инструмент весьма полезен для диагностики смещения и дисперсии заданной модели, поскольку ошибки будут вести себя по-разному. Модель с высоким смещением будет иметь высокую, но похожую тренировочную ошибку, как внутри выборки, так и вне выборки, в то время как переподогнанная модель будет иметь очень низкую тренировочную ошибку.

Убывание вневыборочной ошибки иллюстрирует, что переподогнанные модели могут извлекать преимущества из дополнительных данных или инструментов, которые ограничивают сложность модели, такие как регуляризация, в то время как недоподогнанные модели должны использовать больше признаков либо иным образом увеличивать сложность модели, как показано на рис. 6.8.

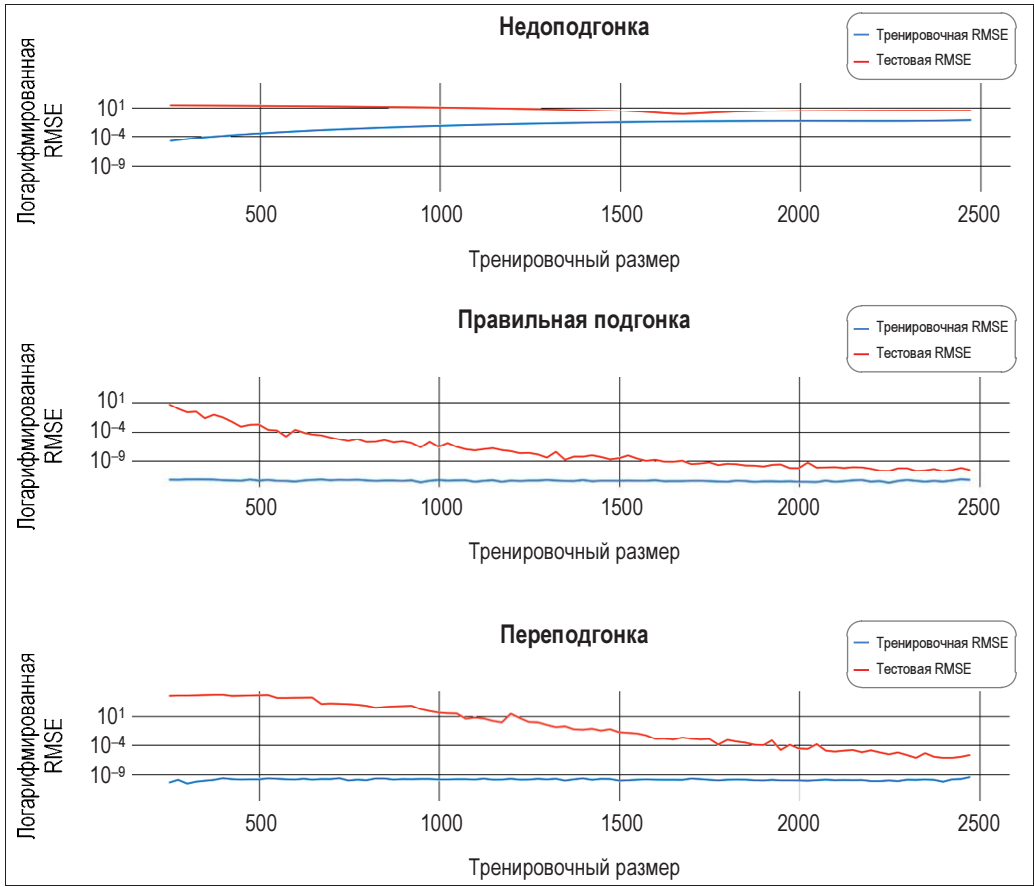


Рис. 6.8. Кривые усвоения (самообучения)

Как использовать перекрестный контроль для отбора модели

Когда для вашего случая применения имеется несколько кандидатных моделей (т. е. алгоритмов), акт выбора одной из них называется задачей *отбора модели*. Процедура отбора модели призвана выявлять модель, которая будет производить наименьшую ошибку предсказания при наличии новых данных.

Несмещенная оценка этой ошибки обобщения требует выполнения теста на данных, которые не были частью тренировки модели. Следовательно, для тренировки модели мы используем только часть имеющихся данных и оставляем другую часть данных для тестирования модели. Как показано на рис. 6.9, для получения несмещенной оценки ошибки предсказания в тренировочный набор не может просачиваться абсолютно никакая информация о тестовом наборе.

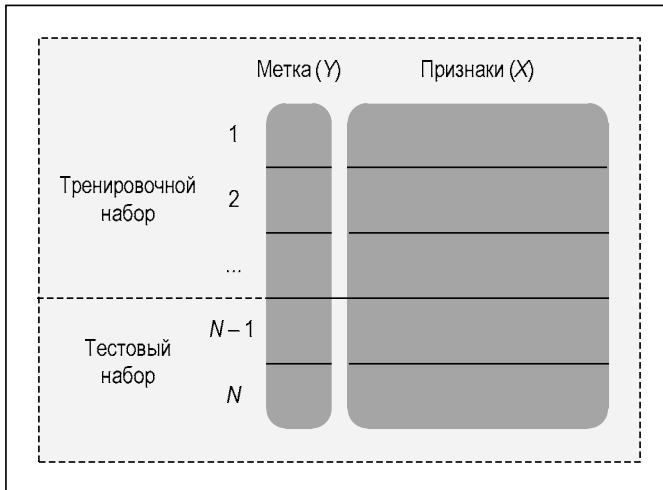


Рис. 6.9. Разбиение данных на тренировочный и тестовый наборы

Существует несколько методов, которые можно использовать для разбиения имеющихся данных, которые различаются с точки зрения объема данных, используемых для тренировки, дисперсии оценок ошибок, вычислительной интенсивности и того, учитываются ли структурные аспекты данных при разбиении данных или нет.

Как реализовать перекрестный контроль на языке Python

Мы проиллюстрируем различные варианты разбиения данных на тренировочный и тестовый наборы, показывая, как назначаются индексы макетного набора данных с десятью наблюдениями в представленном далее фрагменте кода тренировочному и тестовому наборам (подробности см. в сценарии `cross_validation.py`):


```
data = list(range(1, 11))
print(data)

[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
```

Элементарное разбиение на тренировочный и тестовый наборы

В случае одинарного разбиения данных на тренировочный и тестовый наборы следует использовать функцию `sklearn.model_selection.train_test_split`, где параметр `shuffle` по умолчанию обеспечивает рандомизированный отбор наблюдений, которые в свою очередь могут быть реплицированы путем установки параметра `random_state`. Существует также параметр `stratify`, который используется для классификационной задачи, обеспечивая, чтобы тренировочный и тестовый наборы содержали примерно одинаковые доли каждого класса, как показано в следующем фрагменте кода:

```
train_test_split(data, train_size=.8)

[[8, 7, 4, 10, 1, 3, 5, 2], [6, 9]]
```

В этом случае мы тренируем модель, используя все данные, кроме строк 6 и 9, которые потребуются для генерирования предсказаний и измерения ошибок, получаемых на метках `know`. Этот метод полезен для беглой оценки, но чувствителен к разбиению, и стандартная ошибка оценки тестовой ошибки будет выше.

Перекрестный контроль

Перекрестный контроль (или скользящий контроль, *cross-validation*, CV) является популярной стратегией отбора модели, предназначенной для перекрестного подтверждения работоспособности модели. Главная идея перекрестного контроля состоит в разбиении данных один или несколько раз таким образом, что каждый подраздел используется один раз в качестве контрольного набора, а остальная часть — в качестве тренировочного набора: часть данных (тренировочная выборка) используется для тренировки алгоритма, а оставшаяся часть (контрольная выборка) — для оценивания риска алгоритма. Затем перекрестный контроль отбирает алгоритм с наименьшим оценочным риском.

Хотя эвристика разбиения данных носит весьма общий характер, ключевым допущением перекрестной проверки является то, что данные *одинаково распределены и взаимно независимы* (*independently and identically distributed*, iid). В следующих далее разделах мы увидим, что в случае данных временного ряда это часто не так, и требуется иной подход.

Использование отложенного тестового набора

Во время отбора гиперпараметров на основе их контрольной отметки (*validation score*) имейте в виду, что эта контрольная отметка является смещенной из-за многократных тестов и больше не считается хорошей оценкой ошибки обобщения. Для несмещенной оценки частоты появления ошибок мы должны оценить эту отметку из свежего набора данных (рис. 6.10).

По этой причине, как показано на рис. 6.10, мы используем разбиение данных по трем направлениям: одна часть используется в перекрестном контроле и неоднократно разбивается на тренировочный и контрольный наборы. Оставшаяся часть отставляется в сторону как отложенный набор, который используется только после завершения перекрестного контроля для генерирования несмещенной оценки тестовой ошибки¹². Мы проиллюстрируем этот метод в следующей главе, когда начнем строить автоматически обучающиеся модели.

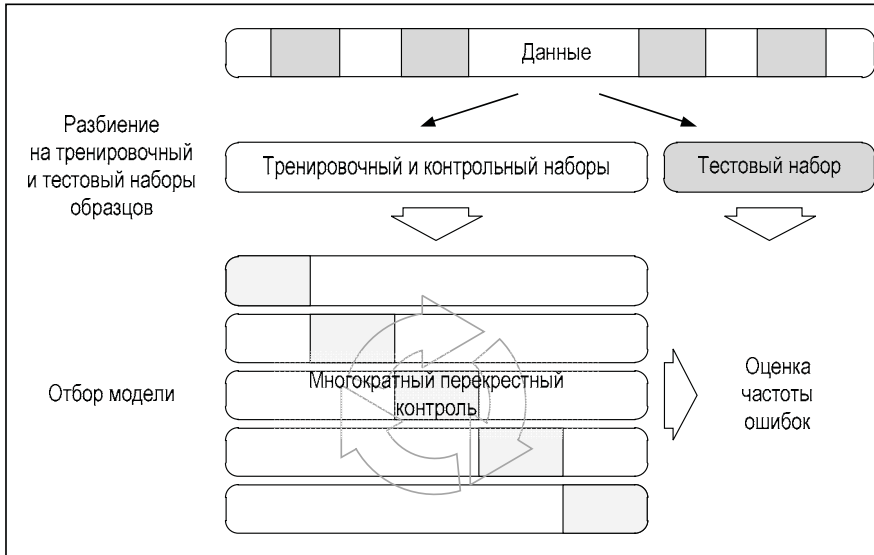


Рис. 6.10. Перекрестный контроль с участием отложенного тестового набора

Итератор k -блочного перекрестного контроля *KFold*

Итератор `sklearn.model_selection.KFold` производит несколько несвязанных между собой подразделов и назначает каждый подраздел один раз тестовому набору, как показано в следующем фрагменте кода¹³:

¹² Тренировочный набор (training set) — это набор образцов, используемых алгоритмом для самообучения, т. е. усвоения взаимосвязей и закономерностей, которое проходит в форме подгонки параметров/весов модели. Контрольный набор (validation set) — это набор образцов, используемый для регуляризации/оптимизации параметров модели и отбора окончательной модели. Тестовый набор (test set) — это набор образцов, используемый только для оценивания результативности полностью натренированной модели. После оценивания окончательной модели на тестовом наборе образцов дальнейшей настройки модели не требуется (Рикардо Гутьеррес-Окуна (Ricardo Gutierrez-Osuna)). — Прим. перев.

¹³ k -Блочный перекрестный контроль (k-fold cross validation) состоит в том, чтобы разбить полный набор данных на k одинаковых по размеру и не накладывающихся друг на друга подмножеств, именуемых блоками (fold), и затем в цикле тренировать модель на всех блоках, кроме i -го блока и оценивать метрику на i -м блоке. — Прим. перев.

```

kf = KFold(n_splits=5)
for train, validate in kf.split(data):
    print(train, validate)

```

```

[2 3 4 5 6 7 8 9] [0 1]
[0 1 4 5 6 7 8 9] [2 3]
[0 1 2 3 6 7 8 9] [4 5]
[0 1 2 3 4 5 8 9] [6 7]
[0 1 2 3 4 5 6 7] [8 9]

```

В дополнение к числу подразделов большинство объектов перекрестного контроля (CV) принимают аргумент `shuffle`, который обеспечивает рандомизацию. Для обеспечения воспроизводимых результатов следует задавать `random_state` следующим образом:

```

kf = KFold(n_splits=5, shuffle=True, random_state=42)
for train, validate in kf.split(data):
    print(train, validate)

```

```

[0 2 3 4 5 6 7 9] [1 8]
[1 2 3 4 6 7 8 9] [0 5]
[0 1 3 4 5 6 8 9] [2 7]
[0 1 2 3 5 6 7 8] [4 9]
[0 1 2 4 5 7 8 9] [3 6]

```

Перекрестный контроль с исключением по одному образцу

В исходной реализации перекрестного контроля использовался метод исключения образцов по одному `LeaveOneOut`, который задействовал каждое наблюдение один раз в качестве контрольного набора, как показано в следующем фрагменте кода:

```

loo = LeaveOneOut()
for train, validate in loo.split(data):
    print(train, validate)

```

```

[1 2 3 4 5 6 7 8 9] [0]
[0 2 3 4 5 6 7 8 9] [1]
...
[0 1 2 3 4 5 6 7 9] [8]
[0 1 2 3 4 5 6 7 8] [9]

```

Благодаря ему максимизируется число тренируемых моделей, что увеличивает вычислительные издержки. Хотя контрольные наборы друг на друга не накладываются, наложение тренировочных наборов максимизируется, тем самым подстегивая корреляцию моделей и их ошибки предсказания. В результате дисперсия ошибки предсказания будет выше для модели с большим числом блоков.

Перекрестный контроль с исключением по P образцов

Перекрестный контроль с исключением по P образцов похож на перекрестную проверку с исключением по одному образцу, и он генерирует все возможные комбинации p строк данных:

```
lpo = LeavePOut(p=2)
for train, validate in lpo.split(data):
    print(train, validate)
```

```
[2 3 4 5 6 7 8 9] [0 1]
[1 3 4 5 6 7 8 9] [0 2]
...
[0 1 2 3 4 5 6 8] [7 9]
[0 1 2 3 4 5 6 7] [8 9]
```

Класс тасованного подразделения данных *ShuffleSplit*

Объект `sklearn.model_selection.ShuffleSplit` создает независимые подразделы с потенциально накладывающимися контрольными наборами, как показано в следующем фрагменте кода:

```
ss = ShuffleSplit(n_splits=3, test_size=2, random_state=42)
for train, validate in ss.split(data):
    print(train, validate)
```

```
[4 9 1 6 7 3 0 5] [2 8]
[1 2 9 8 0 6 7 4] [3 5]
[8 4 5 1 0 6 9 7] [2 3]
```

Настройка параметров с помощью библиотеки *sklean*

Процедура отбора модели, как правило, предусматривает многократный перекрестный контроль вневыборочной результативности моделей с использованием разных алгоритмов (таких как линейная регрессия и случайный лес) или различных конфигураций. Разные конфигурации могут предусматривать изменения в гиперпараметрах либо включение или исключение разных переменных.

Библиотека `yellowbricks` расширяет API библиотеки `sklearn`, генерируя инструменты диагностической визуализации, обеспечивающие процесс отбора модели. Эти инструменты могут использоваться для исследования взаимосвязей между признаками, анализа классификационных или регрессионных ошибок, мониторинга результативности кластерного алгоритма, обследования характеристик текстовых данных и помощи в отборе модели. Далее мы продемонстрируем кривые перекрестного контроля и усвоения, которые предоставляют ценную информацию во время фазы настройки параметров — подробности реализации см. в блокноте `machine_learning_workflow.ipynb`.

Кривые перекрестного контроля с помощью библиотеки yellowbricks

Кривые перекрестного контроля (рис. 6.11, слева) визуализируют влияние одного гиперпараметра на перекрестно-контрольную результативность модели. Такая кривая широко используется для определения того, является ли модель недоподогнанной или переподогнанной к определенному набору данных.

В нашем примере регрессора k ближайших соседей с использованием класса `KNeighborsRegressor`, который имеет только один гиперпараметр, мы ясно видим, что модель недоподогнана для значений k выше 20, где контрольная ошибка падает по мере снижения числа соседей, тем самым усложняя нашу модель, потому что она делает предсказания для большего числа различных групп соседей или участков в признаковом пространстве. Модель начинает показывать симптомы переподгонки для значений ниже 20, поскольку тренировочная и контрольная ошибки расходятся, а средняя вневыборочная результативность быстро ухудшается (см. рис. 6.11).

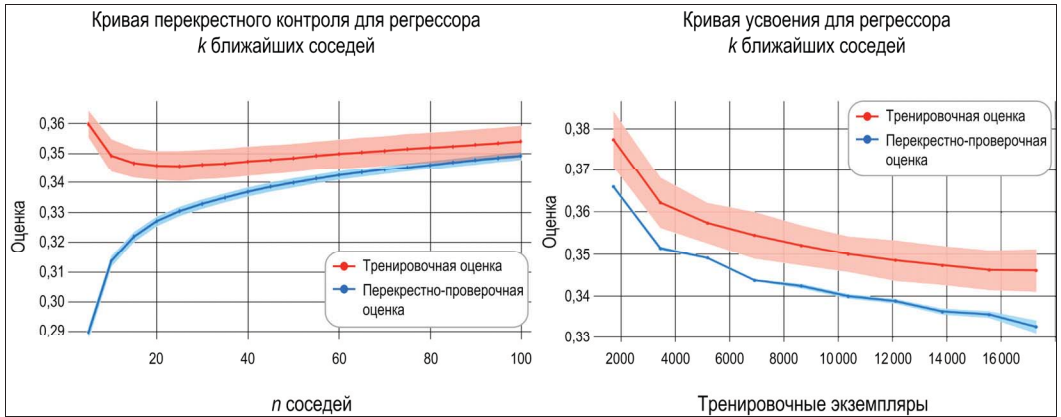


Рис. 6.11. Кривая перекрестного контроля и кривая усвоения

Кривые усвоения

Кривая усвоения (см. рис. 6.11 справа с нашим примером регрессии цен на жилую недвижимость) помогает определить, выиграет ли перекрестно-контрольная результативность модели от дополнительных данных и будут ли ошибки предсказания в большей степени стимулироваться смещением или дисперсией.

Если тренировочная и перекрестно-контрольная результативности сходятся, то дополнительные данные вряд ли улучшат результативность. В этом месте важно оценить, соответствует ли модельная результативность ожиданиям, заданным человеческим эталоном. Если это не так, то следует модифицировать настройки модельных гиперпараметров для того, чтобы лучше улавливать связь между признаками и результатом, либо выбрать другой алгоритм с большей способностью улавливать сложность.

В дополнение к этому, вариация тренировочной и тестовой ошибок, показанная затененными доверительными интервалами, дает подсказки об источниках смещения и дисперсии ошибки предсказания. Вариативность вокруг перекрестно-контрольной ошибки является свидетельством дисперсии, в то время как вариативность для тренировочного набора говорит о смещении в зависимости от размера тренировочной ошибки.

В нашем примере перекрестно-контрольная результативность продолжала падать, но инкрементные улучшения сжались, а ошибки вышли на плато, поэтому маловероятно, что будет много преимуществ от увеличения тренировочного набора. С другой стороны, данные показывают существенную дисперсию с учетом диапазона контрольных ошибок по сравнению с тем, что показано для тренировочных ошибок.

Настройка параметров с помощью интерфейса *GridSearchCV* и конвейера

Поскольку гиперпараметрическая настройка является ключевым компонентом рабочего потока машинного обучения, существуют инструменты, которые автоматизируют этот процесс. Библиотека *sklearn* содержит интерфейс *GridSearchCV*, который параллельно выполняет перекрестный контроль всех комбинаций параметров, улавливает результат и автоматически тренирует модель, используя параметрическую конфигурацию, которая работает лучше всего во время перекрестного контроля полного набора данных.

На практике тренировочный и контрольный наборы часто требуют некоторой предобработки перед перекрестным контролем. Библиотека *sklearn* предлагает конвейер *Pipeline*, который также автоматизирует все необходимые шаги предобработки признаков в автоматизированной гиперпараметрической настройке, обеспечиваемой интерфейсом *GridSearchCV*.

Примеры реализации можно найти в прилагаемом блокноте *machine_learning_workflow.ipynb*, где можно увидеть эти инструменты в действии.

Сложности перекрестного контроля в финансах

Ключевым допущением, принимаемым в методах перекрестного контроля, обсуждавшихся до настоящего времени, является одинаковая распределенность и взаимная независимость (*iid*) образцов, располагаемых для тренировки.

Применительно к финансовым данным зачастую это не так. Напротив, финансовые данные не являются ни одинаково распределенными, ни взаимно независимыми из-за внутрирядовой корреляции и варьирующегося во времени стандартного отклонения. Этот феномен также носит название "*гетероскедастичность*" (подробнее см. следующие две главы). Класс перекрестного контроля временного ряда *TimeSeriesSplit* в модуле *sklearn.model_selection* предназначен для решения задачи линейного порядка данных временных рядов.

Перекрестный контроль временного ряда в библиотеке sklearn

Природа данных временного ряда является причиной того, что перекрестный контроль порождает ситуацию, когда данные из будущего будут использоваться для предсказания данных из прошлого. Это в лучшем случае нереалистично, а в худшем превращается в прочесывание данных, вплоть до того, что будущие данные будут отражать прошлые события.

Для того чтобы решить проблему временной зависимости, объект `sklearn.model_selection.TimeSeriesSplit` реализует форвардный (вперед направленный) тест с расширяющимся тренировочным набором, где последующие тренировочные наборы являются надмножествами прошлых тренировочных наборов:

```
tscv = TimeSeriesSplit(n_splits=5)
for train, validate in tscv.split(data):
    print(train, validate)
```

```
[0 1 2 3 4] [5]
[0 1 2 3 4 5] [6]
[0 1 2 3 4 5 6] [7]
[0 1 2 3 4 5 6 7] [8]
[0 1 2 3 4 5 6 7 8] [9]
```

Параметр `max_train_size` можно использовать для реализации форвардного перекрестного контроля, где размер тренировочного набора остается постоянным с течением времени, подобно тому, как библиотека `zipline` тестирует алгоритм торговли. Библиотека `sklearn` обеспечивает конструирование кастомизированных методов перекрестного контроля с использованием подклассирования, которое мы реализуем в последующих главах.

Прочистка, наложение эмбарго и комбинаторный перекрестный контроль

Применительно к финансовым данным метки часто выводятся из накладывающихся друг на друга точек данных, поскольку финансовые возвраты вычисляются из цен в многочисленных периодах. В контексте стратегий торговли результаты модельного предсказания, которые могут влечь за собой занятие позиции в активе, бывают известны только позже, когда это решение оценивается, например, когда позиция закрывается.

Возникающие в результате риски включают утечку информации из тестового набора в тренировочный набор, что, вероятно, приведет к искусственно инфлированной результативности, решить которую необходимо путем обеспечения того, чтобы все данные относились к определенной точке во времени, т. е. были действительно доступными и известными к тому времени, когда они будут использоваться в качестве входа в модель. В своей книге "Advances in Financial Machine Learning" ("Дос-

тижения в финансовом машинном обучении")¹⁴ Маркос Лопес де Прадо предложил ряд методов для решения этих проблем финансовых данных для перекрестного контроля.

- ◆ *Прочистка (purging)* — устранение точек тренировочных данных, где оценивание происходит после предсказания конкретно-временной точки данных в контрольном наборе во избежание смещения из-за забегания вперед.
- ◆ *Наложение эмбарго (embargoing)* — дополнительное устранение тренировочных образцов, которые следуют за тестовым периодом.
- ◆ *Комбинаторный перекрестный контроль (combinatorial cross-validation)* — форвардный перекрестный контроль строго ограничивает исторические траектории, которые могут быть протестированы. Вместо этого, с учетом T наблюдений вычисляются все возможные тренировочные/тестовые подразделы для $N < T$ групп, каждая из которых поддерживает свой порядок, и выполняются прочистка и наложение эмбарго на потенциально накладывающиеся друг на друга группы. Затем модель тренируется на всех комбинациях $N - k$ групп, при этом тестируя модель на остальных k группах. Результатом является гораздо большее число возможных исторических траекторий.

Книга Прадо "Advances in Financial Machine Learning" содержит примеры исходного кода для реализации этих подходов; исходный код этих методов также доступен посредством новой библиотеки `timeseriescv`.

Резюме

В этой главе мы познакомились с трудностями усвоения регулярностей из данных и обратились к контролируемому, неконтролируемому и подкрепляемому обучению и моделям на их основе как главным формам автоматического обучения, которые мы изучим в этой книге с целью построения стратегий алгоритмической торговли. Мы обсудили вопрос необходимости того, чтобы обучающиеся под контролем алгоритмы принимали допущения о функциональных связях, которые они пытаются усвоить, с целью ограничить поисковое пространство, при этом отметили, что они страдают от индуктивного смещения, которое может приводить к чрезмерным ошибкам обобщения.

Мы представили ключевые аспекты рабочего потока МО, ввели наиболее распространенные метрики ошибок для регрессионных и классификационных моделей, дали объяснение понятия компромисса между смещением и дисперсией и проиллюстрировали различные инструменты управления процессом отбора модели с использованием перекрестного контроля.

В следующей главе мы погрузимся в линейные регрессионные и классификационные модели и разработаем наши первые стратегии алгоритмической торговли, которые используют МО.

¹⁴ В русском переводе "Машинное обучение: алгоритмы для бизнеса". См. <https://www.piter.com/collection/soon/product/mashinnoe-obuchenie-algoritmy-dlya-biznesa>. — Прим. перев.

7

Линейные модели

Семейство линейных моделей является одним из наиболее полезных классов гипотез. Многие обучающиеся алгоритмы, которые широко применяются в алгоритмической торговле, опираются на линейные предсказательные переменные, потому что во многих случаях они могут быть эффективно натренированы, они относительно робастны в условиях шумных финансовых данных и имеют сильные связи с теорией финансов. Линейные предсказатели также интуитивны, легко интерпретируются и часто достаточно хорошо укладываются в данные или, по крайней мере, обеспечивают хорошую опорную линию.

Линейная регрессия известна уже более 200 лет, когда Лежандр и Гаусс применили ее к астрономии и начали анализировать ее статистические свойства. С тех пор для усвоения ее параметров целый ряд расширений адаптировал линейную регрессионную модель и базовый метод *обычных наименьших квадратов* (ordinary least squares, OLS).

- ◆ *Обобщенные линейные модели* (Generalized linear models, GLM) расширяют область применения, допуская результирующие переменные, отклики, которые приводят к распределению ошибок, отличному от нормального распределения. Обобщенные линейные модели включают пробит или логистические модели для *категориальных переменных отклика*, которые появляются в классификационных задачах.
- ◆ Более *робастные методы оценивания* поддерживают статистический вывод в тех случаях, когда данные нарушают исходные допущения, например, из-за корреляции во времени или между наблюдениями. Это часто имеет место в случае панельных данных, содержащих повторяющиеся наблюдения по тем же единицам измерения, таким как исторические финансовые возвраты на универсуме активов.
- ◆ *Усадочные методы* (Shrinkage method) призваны улучшать предсказательную результативность линейных моделей. Они используют штраф за сложность, который смещает усвоенные моделью коэффициенты, снижая модельную дисперсию и улучшая вневыборочную предсказательную результативность.

На практике линейные модели применяются к регрессионным и классификационным задачам в целях статистического вывода и предсказания. Многочисленные модели оценивания стоимости активов, разработанные академическими и отраслевыми исследователями, используют линейную регрессию. Ее применения включают выявление существенных факторов, которые стимулируют финансовые возвраты, получаемые от активов, например, в качестве основы для риск-менеджмента, а также предсказание возвратов на различных временных горизонтах. Классификационные задачи, с другой стороны, охватывают направленные ценовые прогнозы.

В этой главе мы рассмотрим следующие темы:

- ◆ как работает линейная регрессия и какие допущения она делает;
- ◆ как тренировать и диагностировать линейные регрессионные модели;
- ◆ как использовать линейную регрессию для предсказания будущих финансовых возвратов;
- ◆ как использовать регуляризацию для улучшения предсказательной результативности;
- ◆ как работает логистическая регрессия;
- ◆ как конвертировать регрессионную задачу в классификационную.

Примеры исходного кода, дополнительные ресурсы и справочные материалы см. в папке этой главы в онлайн-овом репозитории GitHub.

Линейная регрессия для статистического вывода и предсказания

Как следует из названия, линейные регрессионные модели исходят из допущения, что выход является результатом линейной комбинации входов. Указанная модель также допускает случайную ошибку, которая позволяет каждому наблюдению отклоняться от ожидаемой линейной связи. Причины, почему модель не описывает связь между входом и выходом идеально в детерминированном виде, включают, например, вопросы, связанные с пропущенными переменными, статистическим измерением или сбором данных.

Если мы хотим вывести статистические заключения об истинной (но не наблюдаемой) линейной связи в популяции, основываясь на регрессионных параметрах, оцененных по выборке, то нам нужно добавить допущения о статистической природе этих ошибок. Базовая регрессионная модель делает сильное допущение о том, что распределение ошибок является идентичным по всем ошибкам и что ошибки независимы друг от друга, т. е. знание об одной ошибке не помогает прогнозировать следующую ошибку. Допущение об *одинаково распределенных и взаимно независимых* (independent and identically distributed, iid) ошибках влечет за собой то, что их матрица ковариаций является единичной матрицей, умноженной на константу, представляющую дисперсию ошибки.

Эти допущения гарантируют, что оценки, которые производит метод обычных наименьших квадратов (OLS), не только не смещены, но и эффективны, т. е. имеют алгоритмы, усваивающие самые низкие ошибки выборки¹. Однако на практике эти допущения редко удовлетворяются. В финансах мы часто встречаем панельные данные с многократными наблюдениями по заданному срезу (перекрестному сечению). Попытка оценить систематическое влияние на универсум активов набора рисков факторов во временной динамике, как правило, обнаруживает корреляцию во временной или срезовой размерности, или и в обеих. Как следствие, появились альтернативные обучающиеся алгоритмы, принимающие больше матриц ковариации ошибок, и эти матрицы отличаются от кратных единичной матрице.

С другой стороны, методы, усваивающие смещенные параметры для линейной модели, могут давать оценки с более низкой дисперсией и, следовательно, улучшать прогнозные характеристики. *Усадочные методы* снижают сложность модели, применяя регуляризацию, которая добавляет в линейную целевую функцию штрафной член. Штраф положительно связан с абсолютным размером коэффициентов, вследствие чего они сжимаются относительно базового случая. Большие коэффициенты приводят к более сложной модели, которая сильнее реагирует на изменения во входах. Надлежаще откалиброванный штраф может ограничивать рост коэффициентов модели за пределами того, что предлагает оптимальный компромисс между смещением и дисперсией.

Мы представим базовые срезовые (кросс-секционные) и панельные технические решения для линейных моделей и важные усовершенствования, которые порождают точные оценки, когда ключевые допущения нарушаются. Затем мы проиллюстрируем эти решения, оценив факторные модели, которые повсеместны в разработке стратегий алгоритмической торговли. Наконец, мы сосредоточимся на методах регуляризации.

Множественная линейная регрессионная модель

Мы представим спецификацию и целевую функцию модели, методы усвоения ее параметров, статистические допущения, позволяющие выводить заключение об этих допущениях и проводить их диагностику этих допущений, а также расширения, которые адаптируют модель к ситуациям, когда эти допущения оказываются безуспешными.

¹ Ошибка выборки (sampling error) — это величина погрешности, которая возникает при оценивании значения популяционной переменной на основе значения этой переменной в выборке, взятой из этой популяции. — *Прим. перев.*

Как формулировать модель

Множественная регрессионная модель определяет линейную функциональную связь между одной непрерывной выходной переменной и p входными переменными, которые могут быть любого типа, но которые могут требовать предобработки. Многомерная регрессия, напротив, означает регрессию многочисленных выходных переменных на многочисленных входных переменных.

Линейная регрессионная модель имеет следующую форму в популяции для одного-единственного экземпляра выхода y , входного вектора $\mathbf{x}^T = [x_1 \dots x_p]$ и ошибки ε :

$$y = f(\mathbf{x}) + \varepsilon = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon = \beta_0 + \sum_{j=1}^p \beta_j x_j + \varepsilon.$$

Интерпретация коэффициентов не представляет труда: значение коэффициента β_i является частичным, средним эффектом переменной x_i на выходе при неизменности всех остальных переменных.

Указанная модель также может быть написана более компактно в матричном виде. В этом случае \mathbf{y} — это вектор из N выходных наблюдений, \mathbf{X} — матрица плана² с N строками наблюдений по P переменных плюс столбец из единиц для пересечения, и $\boldsymbol{\beta}$ — вектор, содержащий $P = p + 1$ коэффициентов β_0, \dots, β_p :

$$\begin{matrix} \mathbf{y} & = & \mathbf{X} & \boldsymbol{\beta} & + & \boldsymbol{\varepsilon}. \\ (N \times 1) & & (N \times P) & (N \times P) & & (N \times 1) \end{matrix}$$

Эта модель является линейной по своим $p + 1$ параметрам, но может моделировать нелинейные связи путем выбора или преобразования переменных соответствующим образом, например, путем включения членов разложения многочлена по базису или логарифмических членов. Она также может использовать категориальные переменные в кодировке с одним активным состоянием³ и взаимодействия между переменными, создавая новые входы в форме $x_i \cdot x_j$.

Завершая формулировку модели со статистической точки зрения, что позволяет проводить тестирование статистической гипотезы о параметрах, нам нужно принять специфические допущения о члене ошибки. Мы сделаем это после того, как сначала представим альтернативные методы усвоения параметров.

² Матрица плана регрессионной модели (design matrix), также модельная матрица, матрица регрессоров, обычно обозначаемая \mathbf{X} , — это матрица значений так называемых объясняющих переменных множества объектов. Каждая строка представляет собой отдельный объект, а последующие столбцы соответствуют переменным и их конкретным значениям для этого объекта. Преимущество матрицы плана заключается в том, что она может представлять множество видов экспериментальных проектов и статистических моделей. См. https://en.wikipedia.org/wiki/Design_matrix. — Прим. перев.

³ Словосочетание "кодирование с одним активным состоянием" (one-hot encoding) пришло из терминологии цифровых интегральных микросхем, где оно описывает конфигурацию микросхемы, в которой допускается, чтобы только один бит был положительным (активным — hot). — Прим. перев.

Как тренировать модель

Существует несколько методов усвоения модельных параметров β из данных: *обычные наименьшие квадраты* (ordinary least squares, OLS), *оценивание максимального правдоподобия* (maximum likelihood estimation, MLE) и *стохастический градиентный спуск* (stochastic gradient descent, SGD).

Наименьшие квадраты

Наименьшие квадраты — это изначальный метод усвоения параметров гиперплоскости, которая наилучшим образом аппроксимирует выход из входных данных. Как следует из названия, наилучшая аппроксимация минимизирует сумму квадратов расстояний между выходным значением и гиперплоскостью, представленной моделью.

Разница между модельным предсказанием и выходным результатом для заданной точки данных называется *остатком* (тогда как отклонение истинной модели от истинного выхода в популяции называется *ошибкой*). Следовательно, в формальных терминах метод оценивания наименьших квадратов выбирает коэффициентный вектор β , который минимизирует *сумму квадратов остатков* (residual sum of squares, RSS):

$$\text{RSS}(\beta) = \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N (y_i - f(x_i))^2 = \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^N x_{ij} \beta_j \right)^2 = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta).$$

Следовательно, коэффициенты β_{LS} наименьших квадратов вычисляются как:

$$\arg \min_{\beta_{LS}} \text{RSS} = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta).$$

Оптимальный вектор параметров, минимизирующий сумму квадратов остатков, является результатом приравнивания производных предыдущего выражения относительно β к нулю. В результате получается уникальное решение, исходя из того, что \mathbf{X} имеет полный столбцовый ранг, т. е. входные переменные не являются линейно зависимыми, как показано ниже:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Когда из \mathbf{y} и \mathbf{X} были вычтены соответствующие им средние значения, β представляет собой соотношение ковариации между входами и выходами $\mathbf{X}^T \mathbf{y}$ и выходной дисперсией $(\mathbf{X}^T \mathbf{X})^{-1}$. Существует также геометрическая интерпретация: коэффициенты, которые минимизируют сумму квадратов остатков, обеспечивают, чтобы вектор остатков $\mathbf{y} - \hat{\mathbf{y}}$ был ортогонален подпространству \mathbb{R}^N , охватываемому столбцами \mathbf{X} , а оценки $\hat{\mathbf{y}}$ были ортогональными проекциями в это подпространство.

Оценивание максимального правдоподобия

Оценивание максимального правдоподобия⁴ (MLE) является важным общим методом оценивания параметров статистической модели. Он опирается на функцию правдоподобия, которая вычисляет степень правдоподобия наблюдать выборку выходных значений для заданного набора входных данных как функцию от модельных параметров. Правдоподобие отличается от вероятностей в том, что оно не нормализовано в промежутке от 0 до 1.

Мы можем задать функцию правдоподобия для линейно-регрессионного примера, приняв как допущение распределение для члена ошибки, в частности, стандартное нормальное распределение:

$$\varepsilon_i \sim N(0, 1) \quad \forall i = 1, \dots, n.$$

Это позволит вычислить условную вероятность наблюдать заданный выход y_i при наличии соответствующего входного вектора \mathbf{x}_i и параметров $p(y_i | \mathbf{x}_i, \beta)$:

$$p(y_i | \mathbf{x}_i, \beta) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{\varepsilon_i^2}{2\sigma^2}} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i - \mathbf{x}_i\beta)^2}{2\sigma^2}}.$$

Исходя из того, что выходные значения являются условно независимыми при заданных входах, правдоподобие выборки пропорционально произведению условных вероятностей индивидуальных точек выходных данных. Поскольку легче работать с суммами, чем с произведениями, мы применяем логарифм и в результате получаем логарифмическую функцию правдоподобия:

$$\log \mathcal{L}(\mathbf{y}, \mathbf{x}, \beta) = \sum_{i=1}^n \log \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i - \mathbf{x}_i\beta)^2}{2\sigma^2}}.$$

Цель оценивания максимального правдоподобия — максимизация вероятности выходной выборки, которая фактически наблюдалась путем выбора модельных параметров, принимая наблюдаемые входы как заданные. Отсюда, оценка параметра MLE приводит к максимизации (логарифмической) функции правдоподобия:

$$\beta_{\text{MLE}} = \arg \min_{\beta} \mathcal{L}.$$

Вследствие принятого допущения о нормальном распределении максимизация логарифмической функции правдоподобия приводит к тому же решению параметра, что и в наименьших квадратах, поскольку единственное выражение, которое зависит от параметров, является квадратом остатка в экспоненте. Для других допуще-

⁴ Правдоподобие (likelihood), точнее функция правдоподобия — это совместное распределение вероятностей наблюдаемых данных, выраженное как функция статистических параметров. Оно описывает относительную вероятность или шансы получения наблюдаемых данных для всех допустимых значений параметров и используется для выявления конкретных значений параметров, наиболее правдоподобных с учетом наблюдаемых данных. См. https://en.wikipedia.org/wiki/Likelihood_function. — Прим. перев.

ний о природе распределения и моделей метод оценивания максимального правдоподобия (MLE) произведет другие результаты, и, как мы увидим позже, во многих случаях наименьшие квадраты неприменимы для логистической регрессии.

Градиентный спуск

Градиентный спуск — это универсальный оптимизационный алгоритм, который будет отыскивать стационарные точки гладких функций. Решением будет глобальный оптимум, если целевая функция является выпуклой. Варианты градиентного спуска широко используются в тренировке сложных нейронных сетей, а также для вычисления решений задач оценивания максимального правдоподобия.

Указанный алгоритм использует градиент целевой функции, содержащий ее частные производные относительно параметров. Эти производные показывают, насколько цель изменяется в направлении соответствующих параметров при бесконечно малых шагах. Выясняется, что максимальное изменение значения функции происходит в результате шага в направлении самого градиента.

Следовательно, во время минимизации функции, которая, например, описывает стоимость ошибки предсказания, данный алгоритм вычисляет градиент для текущих значений параметров с использованием тренировочных данных и модифицирует каждый параметр, исходя из отрицательного значения соответствующей ему градиентной компоненты. В результате целевая функция будет принимать меньшее значение и перемещать параметры ближе к решению. Оптимизация останавливается, когда градиент становится малым, а значения параметров меняются очень незначительно.

Размер этих шагов представляет собой темп усвоения (learning rate, темп самообучения) и является критически важным параметром, который может потребовать настройки; многие реализации включают вариант, когда темп усвоения постепенно увеличивается вместе с числом итераций. В зависимости от размера данных указанный алгоритм может итеративно перебирать весь набор данных много раз. Каждая такая итерация называется *эпохой*. Число эпох и допустимый предел (tolerance), используемый для прекращения дальнейших итераций, являются гиперпараметрами, которые можно регулировать.

Стохастический градиентный спуск случайно отбирает точку данных и в отличие от среднего значения над большей выборкой вычисляет градиент этой точки данных, тем самым достигая ускорения. Существуют также пакетные версии, которые на каждом шаге используют определенное число точек данных.

Теорема Гаусса — Маркова

Для вычисления статистики модели и проведения статистического вывода необходимо принять допущения об остатках, т. е. свойствах необъяснимой части входа. Теорема Гаусса — Маркова (Gauss — Markov theorem, GMT) определяет допущения, необходимые для того, чтобы обычные наименьшие квадраты (OLS) произво-

дили несмещенные оценки модельных параметров β , и того, когда эти оценки имеют наименьшую стандартную ошибку среди всех линейных моделей для срезовых (кросс-секционных) данных.

Базовая множественная регрессионная модель делает следующие допущения теории Гаусса — Маркова:

1. В популяции *линейность* соблюдается:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon,$$

где β_i неизвестны, но постоянны; ε — случайная ошибка.

2. Данные для входных переменных x_1, \dots, x_k представляют собой *случайную выборку* из популяции.
3. Идеальная *коллинеарность* отсутствует — нет точных линейных связей между входными переменными.
4. Ошибка имеет *условное среднее значение, равное нулю*, при любом заданном входе: $E[\varepsilon | x_1, \dots, x_k] = 0$.
5. *Гомоскедастичность*, член ошибки ε имеет постоянную дисперсию при заданных входах: $E[\varepsilon^2 | x_1, \dots, x_k] = \sigma^2$.

Четвертое допущение приводит к тому, что нет ни одной пропущенной переменной, которая коррелировала бы с любой из входных переменных. В соответствии с первыми четырьмя допущениями метод обычных наименьших квадратов (OLS) обеспечивает *несмещенные* оценки: включение нерелевантной переменной не смещает оценки пересечения и наклона⁵, но устранение релевантной переменной сместит оценки OLS. Далее, обычные наименьшие квадраты также *состоятельны*: по мере увеличения размера выборки оценки сходятся к истинному значению, поскольку стандартные ошибки становятся произвольно малыми. Обратное, к сожалению, также верно: если условное математическое ожидание ошибки не равно нулю, потому что модели не хватает релевантной переменной, или функциональная форма неверна (т. е. квадратические или логарифмические члены отсутствуют), то все оценки параметров являются смещенными. Если ошибка коррелирована с любой из входных переменных, то обычные наименьшие квадраты также не являются состоятельными, т. е. внесение дополнительных данных не устранит смещение.

Если добавить пятое допущение, то обычные наименьшие квадраты также производят наилучшие линейные, несмещенные оценки (best linear, unbiased estimates, BLUE), где "лучший" означает, что оценки имеют наименьшую стандартную ошибку среди всех линейных оценщиков. Следовательно, если пять приведенных

⁵ Наклон (slope) — угловой коэффициент, определяющий наклон линии регрессии по отношению к оси x , численно равен тангенсу угла между положительным направлением оси x и данной прямой. — Прим. перев.

выше допущений соблюдаются, и целью является статистический вывод, то оценки обычных наименьших квадратов являются именно тем путем, который решает задачу. Однако если цель заключается в предсказании, то, как мы увидим, существуют и другие оценщики, которые выторговывают некоторое смещение в обмен на более низкую дисперсию, достигая во многих средах превосходящей предсказательной результативности.

Теперь, когда мы ввели основные допущения обычных наименьших квадратов (OLS), можем взглянуть на выведение заключения в малых и больших выборках.

Как проводить статистический вывод

Вывод в контексте линейной регрессии призван делать заключения об истинной связи в популяции, исходя из выборочных данных. Оно включает проверку статистической гипотезы о значимости совокупной связи или значений конкретных коэффициентов, а также оценки доверительных интервалов.

Ключевым компонентом статистического вывода является проверочный статистический показатель (или статистика) с известным распределением. Мы используем его для принятия нулевой гипотезы как истинной и вычисления вероятности наблюдать значение этого статистического показателя в выборке, известной как p -значение. Если p -значение падает ниже порога значимости (обычно 5%), то мы эту гипотезу отклоняем, потому что такая ситуация делает фактическое значение выборки очень маловероятным. В то же время мы признаем, что p -значение отражает вероятность того, что мы ошибаемся, когда отклоняем то, что на самом деле является правильной гипотезой.

В дополнение к пяти допущениям теоремы Гаусса — Маркова классическая линейная модель исходит из *нормальности* — популяционная ошибка нормально распределена и не зависит от входных переменных. Указанное допущение влечет за собой то, что выходная переменная нормально распределена, находясь в зависимости от входных переменных. Это сильное допущение позволяет выводить точное распределение коэффициентов, которое, в свою очередь, приводит к точному распределению проверочного статистического показателя, необходимого для таких же точных проверок статистических гипотез в малых выборках. Это допущение часто оказывается безуспешным — финансовые возвраты от активов, например, не являются нормально распределенными, но, к счастью, методы, используемые в условиях нормальности, также допустимы приближенно.

Мы имеем следующие распределительные характеристики и проверочный статистический показатель, приближенно в рамках допущений 1–5 теоремы Гаусса — Маркова и как раз, когда нормальность соблюдается:

- ♦ оценки параметров подчинены многомерному нормальному распределению:

$$\hat{\beta} \sim N\left(\beta, \left(\mathbf{X}^T \mathbf{X}\right)^{-1} \sigma\right);$$

- ♦ в рамках допущений 1–5 теоремы Гаусса — Маркова оценки параметров уже являются несмещенными, и мы можем получить несмещенную оценку σ , т. е. постоянную дисперсию ошибки, используя

$$\hat{\sigma} \sim \frac{1}{N-p-1} \sum_{i=1}^N (y_i - \hat{y}_i)^2;$$

- ♦ статистический показатель t для проверок статистической гипотезы об индивидуальном коэффициенте β_j равен

$$t_j = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{v_j}} \sim t_{N-p-1}$$

и подчиняется t -распределению с $N-p-1$ степенями свободы, где v_j — это элемент j диагонали матрицы $(\mathbf{X}^T \mathbf{X})^{-1}$;

- ♦ t -распределение (Госсета, или Стьюдента) сходится к нормальному распределению, и поскольку 97,5-й квантиль нормального распределения равен 1,96, полезным эмпирическим правилом для 95%-го доверительного интервала вокруг оценки параметра является $\hat{\beta} \pm 2 \cdot \text{se}(\hat{\beta})$. Интервал, который содержит ноль, подразумевает, что мы не можем отклонить нулевую гипотезу о нулевом значении истинного параметра, и, следовательно, является несущественным для модели;
- ♦ статистический показатель F позволяет проверять ограничения на нескольких параметрах, в том числе, является ли вся регрессия значимой. Он служит мерой изменения (снижения) в сумме квадратов остатков (RSS), которое получается из добавочных переменных;
- ♦ наконец, проверка на основе *лагранжева множителя* (Lagrange multiplier, LM) является альтернативой проверке на основе статистического показателя F , лимитируя многочисленные ограничения.

Как диагностировать и устранять проблемы

Диагностика подтверждает допущения модели и не дает делать неправильные заключения во время интерпретации результата и проведения статистического вывода. Они включают меры качества подгонки и различные проверки допущений о члене ошибки, в том числе, насколько близко остатки соответствуют нормальному распределению. Более того, диагностика проверяет, является ли дисперсия остатков действительно постоянной либо проявляет гетероскедастичность⁶, и яв-

⁶ Гетероскедастичность (heteroskedasticity) — это ситуация, когда некоторые диапазоны выходного результата показывают остатки с более высокой дисперсией (что может говорить о предсказательной переменной, которая в уравнении отсутствует). — *Прим. перев.*

ляются ли ошибки условно некоррелированными либо проявляют внутрирядовую корреляцию, т. е. помогает ли знание одной ошибки предсказывать последующие ошибки.

В дополнение к описанным далее проверкам всегда важно визуально обследовать остатки с целью обнаружения возможных систематических регулярностей, поскольку они указывают на то, что в модели отсутствует один или несколько факторов, которые обуславливают выходной результат.

Качество подгонки

Меры качества подгонки оценивают степень, с которой модель объясняет вариацию в выходе. Они помогают оценивать качество спецификации модели, например, отбирать среди разных модельных конструкций. Они различаются в том, как они оценивают подгонку. Меры, обсуждаемые здесь, предоставляют внутривыборочную информацию; мы будем использовать вневыборочное тестирование и перекрестный контроль в следующем разделе, когда сосредоточимся на предсказательных моделях.

Видные меры качества подгонки включают (скорректированный) R^2 , который должен быть максимизирован и основан на оценке наименьших квадратов.

- ♦ R -квадрат (R^2) служит мерой доли вариации в выходных данных, объясненных моделью, и вычисляется как $R^2 = 1 - \text{RSS}/\text{TSS}$, где RSS — это сумма квадратов остатков, TSS — сумма квадратов отклонений выходного результата от его среднего значения. Он также соответствует квадрату коэффициента корреляции между фактическими значениями выхода и оцененными (подогнанными) моделью. Цель состоит в том, чтобы максимизировать R^2 , но он никогда не уменьшается, поскольку модель добавляет больше переменных и, следовательно, поощряет переподгонку.
- ♦ Скорректированный R^2 штрафует R^2 за добавление большего числа переменных; каждая дополнительная переменная должна значительно снижать сумму квадратов остатков (RSS), тем самым обеспечивая более высокое качество подгонки.

В качестве альтернативы должны быть минимизированы *информационный критерий Акаике* (Akaike Information Criterion, AIC) и *байесов информационный критерий* (Bayesian Information Criterion, BIC), оба основанные на оценивании максимального правдоподобия:

- ♦ информационный критерий Акаике:

$$\text{AIC} = -2 \log(\mathcal{L}^*) + 2k,$$

где \mathcal{L}^* — значение максимизированной функции правдоподобия; k — число параметров;

♦ байесов информационный критерий:

$$\text{BIC} = -2 \log(\mathcal{L}^*) + \log(N)k,$$

где N — размер выборки.

Обе метрики штрафуют за сложность, при этом BIC налагает более высокий штраф, вследствие чего он может достигать недоподгонки, тогда как AIC может достигать перепоподгонки в относительном выражении. Концептуально AIC стремится найти модель, которая наилучшим образом описывает неизвестный процесс генерирования данных, в то время как BIC пытается найти наилучшую модель среди множества кандидатов. На практике оба критерия могут использоваться совместно для ориентира при отборе модели, когда целью является внутривыборочная подгонка; в противном случае предпочтительнее перекрестный контроль и отбор на основе оценок ошибки обобщения.

Гетероскедастичность

Допущение 5 теоремы Гаусса — Маркова требует, чтобы ковариация остатков принимала форму $\Sigma = \sigma^2 \mathbf{I}$, т. е. диагональной матрицы с элементами, равными постоянной дисперсии члена ошибки. Гетероскедастичность возникает, когда дисперсия остатков, не будучи постоянной, отличается в разных наблюдениях. Если дисперсия остатков положительно коррелирует с входной переменной, т. е. когда ошибки крупнее для входных значений, которые далеки от их среднего значения, то OLS-оценки стандартных ошибок будут слишком низкими, и, следовательно, статистический показатель t будет завышен, что приведет к ложным обнаружениям связей, где их фактически нет.

Диагностика начинается с визуального осмотра остатков. Систематические регулярности в (предположительно случайных) остатках дают основание для статистических проверок нулевой гипотезы о том, что ошибки являются гомоскедастичными по отношению к различным альтернативам. Эти проверки включают статистические проверки Бройша — Пагана (Breusch — Pagan) и Уайта (White).

Существует несколько способов внесения поправок в оценки обычных наименьших квадратов на гетероскедастичность.

- ♦ Робастные стандартные ошибки (иногда именуемые "белыми" стандартными ошибками по фамилии изобретателя теста Хэлберта Уайта) принимают в расчет гетероскедастичность при вычислении дисперсии ошибок с помощью так называемого *сэндвичного оценщика*.
- ♦ Кластеризованные стандартные ошибки предполагают, что в данных есть четкие группы, которые являются гомоскедастичными, но вариация ошибок различается между группами. Эти группы могут быть разными классами финансовых активов или долевыми ценными бумагами из разных отраслей.

Несколько альтернатив обычным наименьшим квадратам (OLS) оценивают матрицу ковариации ошибок, используя другие допущения, когда $\Sigma \neq \sigma^2 \mathbf{I}$. В библиотеке StatsModels имеются следующие из них.

- ◆ *Взвешенные наименьшие квадраты* (weighted least squares, WLS) — для гетероскедастичных ошибок, где матрица ковариаций имеет только диагональные элементы данных, как для OLS, но теперь элементам данных разрешено варьироваться.
- ◆ *Допустимые обобщенные наименьшие квадраты* (generalized least squares, GLSAR) для автокоррелированных ошибок, которые подчиняются авторегрессионному процессу AR(p) (см. главу о моделях линейных временных рядов).
- ◆ *Обобщенные наименьшие квадраты* (generalized least squares, GLS) для ковариационной матрицы произвольной структуры; дает эффективные и несмещенные оценки в присутствии гетероскедастичности или внутрирядовой корреляции.

Внутрирядовая корреляция

Внутрирядовая корреляция (serial correlation) означает, что следующие подряд остатки, порожденные линейной регрессией, коррелированы, что нарушает четвертое допущение теоремы Гаусса — Маркова. Положительная внутрирядовая корреляция приводит к тому, что стандартные ошибки недооцениваются, а статистический показатель t будет инфлирован, что приведет к ложным обнаружениям, если их игнорировать. Однако существуют процедуры внесения поправок на внутрирядовую корреляцию во время расчета стандартных ошибок.

Статистический показатель Дарбина — Уотсона (Durbin — Watson) диагностирует внутрирядовую корреляцию. Он проверяет гипотезу о том, что остатки обычных наименьших квадратов не автокоррелированы, относительно альтернативы о том, что они подчиняются авторегрессионному процессу (который мы разведем в следующей главе). Проверочный статистический показатель варьируется от 0 до 4, а значения около 2 указывают на неавтокорреляцию, более низкие значения — на положительную автокорреляцию, а более высокие — на отрицательную. Точные пороговые значения зависят от числа параметров и наблюдений и должны отыскиваться в таблицах.

Мультиколлинеарность

Мультиколлинеарность возникает, когда две или более независимых переменных являются высоко коррелированными. Это создает ряд сложностей:

- ◆ трудно определить, какие факторы влияют на зависимую переменную;
- ◆ индивидуальные p -значения могут вводить в заблуждение — p -значение может быть высоким, даже если переменная является важной;
- ◆ доверительные интервалы для регрессионных коэффициентов будут чрезмерными, вероятно, даже нулевыми, что делает невозможным определение эффекта независимой переменной на выход.

Формальное или теоретическое решение, которое вносит поправки на мультиколлинеарность, отсутствует. Вместо этого следует попробовать удалить одну или несколько коррелированных входных переменных или увеличить размер выборки.

Как выполнять линейную регрессию на практике

Прилагаемый блокнот `linear_regression_intro.ipynb` иллюстрирует простую и затем множественную линейную регрессию, в последнем случае используя как обычные наименьшие квадраты OLS, так и градиентный спуск. Применительно к множественной регрессии мы генерируем две случайные входные переменные x_1 и x_2 , которые варьируются от -50 до $+50$, и выходную переменную, вычисляемую как линейную комбинацию входов плюс случайный гауссов шум для удовлетворения допущения 6 о нормальности теоремы Гаусса — Маркова:

$$y = 50 + x_1 + 3x_2 + \varepsilon, \quad \varepsilon \sim N(0, 50).$$

Обычные наименьшие квадраты из библиотеки StatsModels

Мы используем библиотеку StatsModels для оценивания множественной регрессионной модели, которая точно отражает процесс генерирования данных, следующим образом:

```
import statsmodels.api as sm

X_ols = sm.add_constant(X)
model = sm.OLS(y, X_ols).fit()

print(model.summary())
```

В результате мы получаем статистическую сводку результатов регрессии обычными наименьшими квадратами (OLS):

OLS Regression Results

| | | | |
|-------------------|------------------|---------------------|-----------|
| Dep. Variable: | Y | R-squared: | 0.789 |
| Model: | OLS | Adj. R-squared: | 0.788 |
| Method: | Least Squares | F-statistic: | 1160. |
| Date: | Tue, 14 May 2019 | Prob (F-statistic): | 1.27e-210 |
| Time: | 17:15:40 | Log-Likelihood: | -3316.5 |
| No. Observations: | 625 | AIC: | 6639. |
| Df Residuals: | 622 | BIC: | 6652. |
| Df Model: | 2 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P> t | [0.025 | 0.975] |
|-------|---------|---------|--------|-------|--------|--------|
| const | 50.2414 | 1.956 | 25.682 | 0.000 | 46.400 | 54.083 |
| X_1 | 0.9875 | 0.065 | 15.167 | 0.000 | 0.860 | 1.115 |
| X_2 | 2.9768 | 0.065 | 45.720 | 0.000 | 2.849 | 3.105 |

| | | | |
|----------------|-------|-------------------|-------|
| Omnibus: | 4.505 | Durbin-Watson: | 2.051 |
| Prob(Omnibus): | 0.105 | Jarque-Bera (JB): | 4.569 |
| Skew: | 0.192 | Prob(JB): | 0.102 |
| Kurtosis: | 2.832 | Cond. No. | 30.0 |

=====

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Верхняя часть сводки показывает характеристики набора данных, а именно метод оценивания, число наблюдений и параметров, и указывает, что оценки стандартных ошибок не учитывают гетероскедастичность. Средняя панель демонстрирует значения коэффициентов, которые точно отражают процесс генерирования искусственных данных. Мы можем подтвердить, что оценки, отображаемые в середине сводного результата, могут быть получены с использованием формулы обычных наименьших квадратов (OLS), выведенной ранее:

```
beta = np.linalg.inv(X_ols.T.dot(X_ols)).dot(X_ols.T.dot(y))

pd.Series(beta, index=X_ols.columns)

const    50.24
X_1       0.99
X_2       2.98
dtype: float64
```

На рис. 7.1 показана гиперплоскость, подогнанная этой моделью для случайно сгенерированных точек данных.

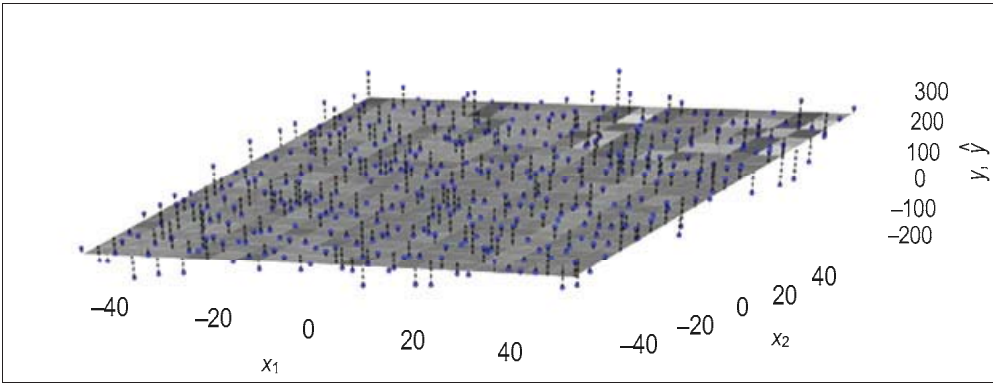


Рис. 7.1. Гиперплоскость, подогнанная моделью

Статистическая сводка результатов в правой верхней части показывает только что рассмотренные меры качества подгонки, а также F -тест, т. е. проверку со статистическим показателем F в качестве критерия, который отклоняет гипотезу о том, что

все коэффициенты равны нулю и не релевантны. Схожим образом статистический показатель t показывает, что пересечение и оба коэффициента наклона, как и следовало ожидать, имеют высокую статистическую значимость.

В нижней части сводки содержится диагностика остатков. Слева показаны асимметрия и эксцесс, которые используются для проверки гипотезы о нормальности. И обобщенная (omnibus) статистическая проверка, и проверка Харке — Бера (Jarque — Bera) не могут отклонить нулевую гипотезу о том, что остатки нормально распределены. Статистический показатель Дарбина — Уотсона (Durbin — Watson) проверяет на наличие внутрирядовой корреляции в остатках и имеет значение около 2, которое с учетом 2 параметров и 625 наблюдений не отклоняет гипотезу об отсутствии внутрирядовой корреляции.

Наконец, число обусловленности (Cond. No.) обеспечивает подтверждающее наблюдение о мультиколлинеарности: это отношение квадратных корней наибольшего и наименьшего собственных значений матрицы плана регрессии, содержащей входные данные. Значение выше 30 говорит о том, что регрессия может иметь значительную мультиколлинеарность.

Библиотека StatsModels включает дополнительные диагностические проверки, ссылки на которые приведены в блокноте.

Стохастический градиентный спуск с помощью библиотеки sklearn

Библиотека sklearn в модуле `linear_models` содержит модель `SGDRegressor`. Для усвоения параметров той же самой модели с помощью этого метода нам нужно сначала стандартизировать данные, потому что градиент чувствителен к шкале. Для этой цели мы используем класс `StandardScaler`, который для каждой входной переменной во время шага подгонки вычисляет ее среднее значение и стандартное отклонение, а затем во время шага преобразования вычитает среднее и делит на стандартное отклонение. Все это можно удобно выполнить в одной команде `fit_transform()`:

```
scaler = StandardScaler()
X_ = scaler.fit_transform(X)
```

Затем мы инстанцируем регрессор `SGDRegressor`, используя значения по умолчанию, за исключением параметра `random_state`, с целью обеспечения репликации:

```
sgd = SGDRegressor(loss='squared_loss', fit_intercept=True,
                    # Перетасовать тренировочные данные
                    # для более высоких оценок градиентов
                    shuffle=True, random_state=42,
                    # Снижать темп усвоения во временной динамике
                    learning_rate='invscaling',
                    # Параметры для траектории темпа усвоения
                    eta0=0.01, power_t=0.25)
```


Теперь можно вписать модель `sgd`, создать внутривыборочные предсказания для модели обычных наименьших квадратов и модели стохастического градиентного спуска `sgd` и вычислить корень квадратный из среднеквадратической ошибки (RMSE) для каждой из них:

```
sgd.fit(X=X_, y=y)
resids = pd.DataFrame({'sgd': y - sgd.predict(X_),
                      'ols': y - model.predict(sm.add_constant(X))})
resids.pow(2).sum().div(len(y)).pow(.5)

ols    50.06
sgd    50.06
dtype: float64
```

Как и ожидалось, обе модели дают один и тот же результат. Теперь мы возьмемся за более амбициозный проект, применив линейную регрессию для оценивания многофакторной модели ценообразования активов.

Как строить линейную факторную модель

Стратегии алгоритмической торговли используют *линейные факторные модели* для квантификации связи между возвратом от актива и источниками риска, которые являются главными движущими силами этих возвратов. Каждый факторный риск несет премию, и можно ожидать, что суммарные возвраты от активов будут соответствовать средневзвешенному значению этих рискованных премий.

В процессе портфельного менеджмента существует несколько практических применений факторных моделей от конструирования и отбора финансовых активов до риск-менеджмента и оценивания результативности. Важность факторных моделей продолжает расти, поскольку общие рискованные факторы теперь торгуемы.

- ◆ Сводка возвратов от многочисленных активов по гораздо меньшему числу факторов снижает объем данных, необходимых для оценки матрицы ковариаций во время оптимизации портфеля.
- ◆ Оценка влияния на актив или портфель этих факторов позволяет управлять результирующим риском, например, путем ввода подходящих хеджей, когда сами рискованные факторы торгуются.
- ◆ Факторная модель позволяет анализировать инкрементное сигнальное содержание новых альфа-факторов.
- ◆ Факторная модель также помогает оценивать то, за счет чего формируется результативность менеджера: обусловлена ли она его умением отбирать финансовые активы и хронометрировать рынок, либо вместо этого она может быть объяснена кренами портфелей в сторону известных движущих сил возвратности, которые сегодня могут быть воспроизведены как низкостоимостные пассивно управляемые фонды, избавляя от расходов на комиссионные за активный менеджмент.

Следующие далее примеры относятся к долевым ценным бумагам, но рисковые факторы были определены для всех классов финансовых активов (справочные материалы см. в репозитории GitHub).

От модели CAPM к пятифакторной модели Фама — Френча

Рисковые факторы являются ключевым компонентом квантитативных моделей, поскольку модель ценообразования капитальных активов (Capital Asset Pricing Model, CAPM) объяснила математическое ожидание финансовых возвратов от всех N активов r_i , $i = 1, \dots, N$, используя их соответствующую подверженность влиянию β_i единственного фактора, т. е. математического ожидания избыточного возврата от совокупного рынка, превышающего безрисковую ставку r_f . Указанная модель принимает следующую линейную форму:

$$E[r_i] = \alpha_i + r_f + \beta_i (E[r_m] - r_f).$$

Она отличается от классического фундаментального анализа в стиле Додда и Грэма, где финансовые возвраты зависят от характеристик фирмы. Обоснование заключается в том, что в совокупности инвесторы не могут устранить этот так называемый систематический риск путем диверсификации. Следовательно, в равновесии они требуют компенсации за владение активом, соизмеримой с его систематическим риском. Данная модель влечет за собой то, что с учетом эффективных рынков, где цены немедленно отражают всю публичную информацию, не должно быть превосходящих скорректированных на риск возвратов, т. е. α должен быть равен нулю.

Эмпирические тесты модели используют линейную регрессию и неуклонно оказывались безуспешными, вызывая разногласия по вопросу о том, виноваты ли в этом эффективные рынки или однофакторный аспект совместной гипотезы. Оказывается, что оба посыла, вероятно, неверны.

- ◆ Джозеф Стиглиц получил Нобелевскую премию по экономике за 2001 г. отчасти за то, что показал, что рынки, как правило, не являются идеально эффективными: если рынки являются эффективными, нет никакой ценности в сборе данных, потому что эта информация уже отражена в ценах. Однако если нет стимула для сбора информации, то трудно понять, каким образом она уже должна быть отражена в ценах.
- ◆ С другой стороны, теоретические и эмпирические улучшения в рамках CAPM-модели свидетельствуют о том, что дополнительные факторы помогают объяснять некоторые аномалии, которые заключались в превосходящих скорректированных на риск возвратах, не зависящих от влияния совокупного рынка, такие как более высокие возвраты у небольших фирм.

В 1976 г. Стивен Росс (Stephen Ross) предложил теорию арбитражного ценообразования (Arbitrage Pricing Theory, APT) в качестве альтернативы, которая учитывает несколько рисков факторов, тщательно избегая эффективности рынка. В отличие

от модели CAPM, указанная теория исходит из того, что из-за неправильной оценки стоимости могут существовать возможности для превосходящих возвратов, но они быстро задействуются за счет арбитража. В теории не уточняются, какие это факторы, но исследования автора позволяют предположить, что наиболее важными являются изменения в инфляции и промышленном производстве, а также изменения в рискованных премиях или срочной структуре процентных ставок.

Кеннет Френч (Kenneth French) и Юджин Фама (Eugene Fama), которые стали лауреатами Нобелевской премии в 2013 г., выявили дополнительные рисковые факторы, которые зависят от характеристик фирмы и широко используются сегодня. В 1993 г. трехфакторная модель Фама — Френча добавила относительные размеры и стоимость фирм в единственный источник риска по модели CAPM. В 2015 г. пятифакторная модель еще больше расширила их набор, включив в него прибыльность фирм и уровень инвестирования, которые в истекшие годы доказали свою значимость на практике. В дополнение к этому многие факторные модели включают фактор ценового импульса.

Рисковые факторы Фама — Френча рассчитываются как разница в возвратах на диверсифицированных портфелях с высокими или низкими значениями в соответствии с метриками, отражающими данный рисковый фактор. Эти возвраты получают путем сортировки акций в соответствии с указанными метриками, а затем лонгированием акций, находящихся выше определенного процентиля, и шортированием акций, находящихся ниже определенного процентиля. Метрики, ассоциированные с рисковыми факторами, определяются следующим образом:

- ◆ размер — рыночная стоимость собственного капитала (Market Equity, ME);
- ◆ стоимость — балансовая стоимость собственного капитала (Book Value of Equity, BE), деленная на ME;
- ◆ операционная прибыльность (Operating Profitability, OP) — выручка минус стоимость реализованных товаров/активов;
- ◆ инвестиции: инвестиции/активы.

Для ведомого данными обнаружения рисковых факторов также существуют технические решения неконтролируемого обучения с использованием факторов и компонентного анализа, которые мы рассмотрим в *главе 12*.

Получение рисковых факторов

Фама и Френч обеспечивает наличие обновляемых данных о рисковых факторах и исследовательских портфелях через свой веб-сайт, и для получения этих данных можно использовать библиотеку `pandas-datareader`. Дополнительную информацию по поводу такого применения см. в блокноте `fama_macbeth.ipynb`.

В частности, мы будем использовать пять факторов Фама — Френча, которые получаются в результате сортировки акций сначала на три размерные группы, а затем на две по каждому из оставшихся трех специфичных для фирмы факторов. Таким

образом, факторы предусматривают три набора взвешенных по стоимости портфелей, сформированных в виде сортировок 3×2 по размеру и балансовой стоимости относительно рыночной стоимости, размеру и операционной прибыльности, а также размеру и инвестициям. Значения рисковых факторов рассчитаны как средние возвраты портфелей (табл. 7.1).

Таблица 7.1. Рисковые факторы Фама — Френча

| Понятие | Метка | Имя | Расчет рискового фактора |
|--------------|-------|--|---|
| Размер | SMB | Малый минус большой (small minus big) | Портфель с девятью малостоимостными акциями ⁷ минус портфель с девятью высокостоимостными |
| Стоимость | HML | Высокий минус низкий (high minus low) | Портфель с двумя стоимостными акциями минус портфель с двумя ростовыми акциями (с низким значением BE/ME) |
| Прибыльность | RMW | Робастный минус слабый (robust minus weak) | Два робастных по операционной прибыли портфеля минус два слабых |
| Инвестиции | CMA | Консервативный минус агрессивный (conservative minus aggressive) | Два консервативных инвестиционных портфеля минус два агрессивных |
| Рынок | Rm-Rf | Избыточная возвратность на рынке | Взвешенная по стоимости возвратность всех фирм встроенных и внесенных в список главных бирж США с хорошими данными минус ставка по одномесячным казначейским векселям |

Мы будем использовать возвраты на месячной частоте, которую мы получаем за период 2010–2017 гг. следующим образом:

```
import pandas_datareader.data as web

ff_factor = 'F-F_Research_Data_5_Factors_2x3'
ff_factor_data = web.DataReader(ff_factor, 'famafrench', start='2010', end='2017-12')[0]
ff_factor_data.info()

<class 'pandas.core.frame.DataFrame'>
PeriodIndex: 96 entries, 2010-01 to 2017-12
```

⁷ Малостоимостные акции (small-value stock, small stock) — это акции с малой рыночной капитализацией, а также акции, которые торгуются по балансовой стоимости или ниже ее. Найти акции, соответствующие обоим этим критериям, очень трудно, но может оказаться полезным предприятием, поскольку малостоимостные акции обычно приносят высокую возвратность.

См. https://www.investopedia.com/terms/s/small_value_stock.asp. — Прим. перев.

```

Freq: M
Data columns (total 6 columns):
Mkt-RF    96 non-null float64
SMB       96 non-null float64
HML       96 non-null float64
RMW       96 non-null float64
CMA       96 non-null float64
RF        96 non-null float64
dtypes: float64(6)
memory usage: 5.2 KB

```

Фама и Френч также обеспечивают наличие многочисленных портфелей, которые мы можем использовать для иллюстрации оценивания влияния факторов, а также стоимости рискованных премий, доступных на рынке за данный период времени. Мы будем использовать панель из 17 отраслевых портфелей с ежемесячной частотой. Мы вычтем безрисковую ставку из возвратов, потому что факторная модель работает с избыточными возвратами:

```

ff_portfolio = '17_Industry_Portfolios'
ff_portfolio_data = web.DataReader(ff_portfolio, 'famafrench',
                                   start='2010', end='2017-12')[0]
ff_portfolio_data = ff_portfolio_data.sub(ff_factor_data.RF, axis=0)
ff_portfolio_data.info()

```

```

<class 'pandas.core.frame.DataFrame'>
PeriodIndex: 96 entries, 2010-01 to 2017-12
Freq: M
Data columns (total 17 columns):
Food      96 non-null float64
Mines     96 non-null float64
Oil       96 non-null float64
...
Rtail     96 non-null float64
Finan     96 non-null float64
Other     96 non-null float64
dtypes: float64(17)
memory usage: 13.5 KB

```

Теперь на основе этих панельных данных мы построим линейную факторную модель, используя метод, устраняющий безуспешность некоторых базовых линейно-регрессионных допущений.

Регрессия Фама — Макбета

При наличии данных о рискованных факторах и портфельных возвратах полезно оценить их влияние на портфель, т. е. выяснить, насколько рискованные факторы стимулируют портфельные возвраты, а также величину влияния данного фактора, т. е. какова рыночная риск-факторная премия. Рисксовая премия позволяет оценить воз-

врат для любого портфеля при условии, что влияние фактора известно или может быть принято, как допущение.

Более формально, мы будем иметь возвраты $i = 1, \dots, N$ от активов или портфелей за $t = 1, \dots, T$ периодов, и избыточные возвраты периода владения по каждому активу будут обозначаться как r_{it} . Цель состоит в том, чтобы проверить, объясняют ли $j = 1, \dots, M$ факторов F_{jt} избыточные возвраты и рисковую премию, ассоциированные с каждым фактором. В нашем случае мы имеем $N = 17$ портфелей и $M = 5$ факторов с 96 периодами данных для каждого.

Факторные модели оцениваются для многочисленных акций за определенный период. В таких срезовых (кросс-секционных) регрессиях вполне могут возникать проблемы статистического вывода, потому что фундаментальные допущения классической линейной регрессии могут не соблюдаться. Возможные нарушения охватывают измерительные ошибки, ковариацию остатков из-за гетероскедастичности и внутрирядовой корреляции и мультиколлинеарность.

Для решения проблемы статистического вывода, вызванной корреляцией остатков, Фама и Макбет предложили двухэтапную методологию для срезовой (кросс-секционной) регрессии возвратов на факторах. Двухэтапная регрессия Фама — Макбета разработана для оценивания премии, вознаграждаемой за воздействие определенного рискового фактора со стороны рынка.

♦ **Первый этап.** N -рядовая регрессия, по одной для каждого актива или портфеля, его избыточных возвратов на факторах для оценивания факторных нагрузок⁸. В матричной форме для каждого актива:

$$\mathbf{r}_i = \mathbf{F} \boldsymbol{\beta}_i + \boldsymbol{\varepsilon}_i.$$

$T \times 1 \quad T \times (m+1) \quad (m+1) \times 1 \quad T \times 1$

♦ **Второй этап.** T срезовых (кросс-секционных) регрессий, по одной для каждого периода времени, для оценивания рисковой премии. В матричной форме мы получаем вектор $\hat{\lambda}_t$ рисковых премий для каждого периода:

$$\mathbf{r}_t = \hat{\boldsymbol{\beta}} \hat{\boldsymbol{\lambda}}_t.$$

$N \times (M+1) \quad N(M+1) \quad (M+1) \times 1$

Теперь можно вычислить премии за факторные риски как среднее время и получить статистический показатель t для оценки их индивидуальной значимости, используя допущение, что оценки рисковых премий независимы во временной динамике:

$$t = \frac{\lambda_j}{\sigma(\lambda_j)/\sqrt{T}}.$$

⁸ Факторная нагрузка (factor loading) — мера влияния данного фактора на возвратность финансового актива; может интерпретироваться как стандартизированный регрессионный коэффициент и таким образом указывать на корреляцию актива с фактором. См. <https://www.theanalysisfactor.com/factor-analysis-1-introduction/>. — Прим. перев.

Если бы у нас была очень крупная и репрезентативная выборка данных о торгуемых рисков факторах, то в качестве оценки рисков премии мы могли бы использовать среднее значение выборки. Однако у нас, как правило, недостаточно длинная история, и ошибка вокруг среднего выборки может быть довольно крупной. Методология Фама — Макбета для определения факторных премий эффективно задействует ковариацию факторов с другими активами. Второй момент возвратов от активов оценить легче, чем первый момент, и получение более гранулярных данных значительно улучшает оценивание, что неверно для оценивания на основе средних.

Первый этап для получения 17 оценок факторных нагрузок можно реализовать следующим образом:

```
from statsmodels.api import OLS, add_constant

betas = []
for industry in ff_portfolio_data:
    step1 = OLS(endog=ff_portfolio_data.loc[ff_factor_data.index, industry],
                exog=add_constant(ff_factor_data)).fit()
    betas.append(step1.params.drop('const'))

betas = pd.DataFrame(betas,
                    columns=ff_factor_data.columns,
                    index=ff_portfolio_data.columns)

betas.info()

<class 'pandas.core.frame.DataFrame'>
Index: 17 entries, Food to Other
Data columns (total 6 columns):
Mkt-RF    17 non-null float64
SMB       17 non-null float64
HML       17 non-null float64
RMW       17 non-null float64
CMA       17 non-null float64
RF        17 non-null float64
dtypes: float64(6)
memory usage: 1.6+ KB
```

Относительно второго этапа мы выполняем 96 регрессий периодных возвратов для среза портфелей на факторных нагрузках:

```
lambdas = []
for period in ff_portfolio_data.index:
    step2 = OLS(endog=ff_portfolio_data.loc[period, betas.index], exog=betas).fit()
    lambdas.append(step2.params)
```

```

lambdas = pd.DataFrame(lambdas,
                        index=ff_portfolio_data.index,
                        columns=betas.columns.tolist())

lambdas.info()

```

```

<class 'pandas.core.frame.DataFrame'>
PeriodIndex: 95 entries, 2010-02 to 2017-12
Freq: M
Data columns (total 6 columns):
Mkt-RF    95 non-null float64
SMB       95 non-null float64
HML       95 non-null float64
RMW       95 non-null float64
CMA       95 non-null float64
RF        95 non-null float64
dtypes: float64(6)
memory usage: 7.7 KB

```

Наконец, мы вычисляем среднее значение для 96 периодов и получаем оценки факторной рискованной премии:

```

lambdas.mean()

Mkt-RF    1.244610
SMB       0.007380
HML      -0.696972
RMW      -0.255768
CMA      -0.308635
RF       -0.013344
dtype: float64

```

Модуль `linear_models` расширяет библиотеку `StatsModels` за счет различных моделей для панельных данных, а также реализует двухэтапную процедуру Фама — Макбета:

```

model = LinearFactorModel(portfolios=ff_portfolio_data, factors=ff_factor_data)
res = model.fit()

```

В итоге мы получим тот же результат:

```

LinearFactorModel Estimation Summary
=====
No. Test Portfolios:      17  R-squared:      0.6943
No. Factors:             6   J-statistic:    19.155
No. Observations:       95   P-value       0.0584
Date:                    Wed, Oct 31 2018  Distribution:  chi2(11)
Time:                    15:15:52
Cov. Estimator:         robust

```


Risk Premia Estimates

| | Parameter | Std. Err. | T-stat | P-value | Lower CI | Upper CI |
|--------|-----------|-----------|---------|---------|----------|----------|
| Mkt-RF | 1.2446 | 0.3928 | 3.1689 | 0.0015 | 0.4748 | 2.0144 |
| SMB | 0.0074 | 0.7055 | 0.0105 | 0.9917 | -1.3753 | 1.3901 |
| HML | -0.6970 | 0.5334 | -1.3067 | 0.1913 | -1.7424 | 0.3484 |
| RMW | -0.2558 | 0.6888 | -0.3713 | 0.7104 | -1.6057 | 1.0942 |
| CMA | -0.3086 | 0.4737 | -0.6515 | 0.5147 | -1.2371 | 0.6198 |
| RF | -0.0133 | 0.0132 | -1.0092 | 0.3129 | -0.0393 | 0.0126 |

Covariance estimator:
HeteroskedasticCovariance
See full_summary for complete results

Прилагаемый блокнот иллюстрирует использование категориальных переменных с помощью отраслевых фиктивных переменных (industry dummies) при оценивании рисков премий для более крупной панели индивидуальных акций.

Усадочные методы: регуляризация для линейной регрессии

Методы наименьших квадратов для тренировки линейной регрессионной модели будут давать наилучшие, линейные и несмещенные оценки коэффициентов, когда допущения теоремы Гаусса — Маркова соблюдаются. Такие варианты, как обобщенные наименьшие квадраты GLS, также хороши, даже когда допущения обычных наименьших квадратов OLS о матрице ковариации ошибок нарушаются. Однако существуют оценщики, которые производят смещенные коэффициенты, снижая дисперсию, достигая пониженной ошибки обобщения в целом.

Когда линейная регрессионная модель содержит много коррелированных переменных, их коэффициенты будут слабо определены, потому что эффект большого положительного коэффициента на сумму квадратов остатков (RSS) может быть нейтрализован аналогичным большим отрицательным коэффициентом на коррелированной переменной. Следовательно, из-за этого маневренного пространства коэффициентов модель будет иметь тенденцию к высокой дисперсии, что увеличивает риск того, что модель будет переподогнана к выборке.

Как хеджироваться от переподгонки

Одним из популярных методов управления переподгонкой является *регуляризация*, которая предусматривает добавление штрафного члена в функцию ошибки, противодействуя достижению коэффициентами крупных значений. Другими словами, наложение размерных ограничений на коэффициенты может смягчить результирующее потенциально негативное влияние на вневыборочные предсказания. Мы

столкнемся с регуляризационными методами для всех моделей, поскольку переопределение является воистину повсеместной проблемой.

В этом разделе мы представим усадочные методы, которые касаются двух мотивов для совершенствования подходов к линейным моделям, обсуждавшихся до сих пор.

- ◆ *Точность предсказания* — низкое смещение, но высокая дисперсия оценок по методу наименьших квадратов говорит о том, что ошибка обобщения может быть снижена за счет усадки (сжатия) или установки некоторых коэффициентов, равными нулю, тем самым выторговывая более высокое смещение в обмен на снижение в дисперсии модели.
- ◆ *Интерпретация* — большое число предсказательных переменных может усложнить интерпретацию или передачу общей картины результатов. Возможно, предпочтительнее пожертвовать некоторой детализацией, лимитируя модель меньшим подмножеством параметров с самыми сильными эффектами.

Усадочные модели (shrinkage model) ограничивают регрессионные коэффициенты, накладывая штраф на их размер. Указанные модели достигают этой цели, добавляя член в целевую функцию, вследствие чего коэффициенты усадочной модели минимизируют сумму квадратов остатков (RSS) плюс штраф, который положительно связан с (абсолютным) размером коэффициентов. Добавленный штраф превращает отыскание линейно-регрессионных коэффициентов в ограниченную минимизационную задачу, которая, в общем случае, принимает следующую лагранжеву форму:

$$\hat{\beta}^S = \arg \min_{\beta^S} \sum_{i=1}^N \left[\left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda S(\beta) \right] = \arg \min_{\beta^S} \mathbf{y} - \mathbf{X}\beta - \lambda S(\beta).$$

Регуляризационный параметр λ определяет размер штрафного эффекта, т. е. силу регуляризации. Когда λ является положительной, коэффициенты будут отличаться от неограниченных наименьших квадратов параметров, из чего вытекает смещенная оценка. Гиперпараметр λ должен выбираться адаптивно с помощью перекрестного контроля, минимизируя оценки ожидаемой ошибки предсказания.

Усадочные модели отличаются тем, как они вычисляют штраф, т. е. функциональную форму S . Наиболее распространенными версиями являются гребневая регрессия, которая использует сумму квадратов коэффициентов, тогда как модель лассо основывает штраф на сумме абсолютных значений коэффициентов.

Как работает гребневая регрессия

Гребневая регрессия сжимает (осаждает) регрессионные коэффициенты путем добавления штрафа в целевую функцию, которая равна сумме квадратов коэффициентов, что в свою очередь соответствует норме L_2 вектора коэффициентов:

$$S(\beta) = \sum_{i=1}^p \beta_i^2 = \|\beta\|^2.$$

Следовательно, гребневые коэффициенты определяются следующим образом:

$$\begin{aligned}\hat{\beta}^{\text{гребневой}} &= \arg \min_{\beta^{\text{гребневой}}} \sum_{i=1}^N \left[\left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right] = \\ &= \arg \min_{\beta^{\text{гребневой}}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}.\end{aligned}$$

Пересечение β_0 было исключено из штрафа для того, чтобы сделать процедуру независимой от начала координат, выбираемого для выходной переменной; в противном случае добавление константы ко всем выходным значениям изменит все параметры наклона, в отличие от параллельного сдвига.

Важно стандартизировать входы путем вычитания из каждого входа соответствующего среднего и деления результата на стандартное отклонение входа, потому что гребневое решение является чувствительным к шкале входных данных. Для гребневого оценщика существует также замкнутое решение, которое напоминает случай обычных наименьших квадратов (OLS):

$$\hat{\beta}^{\text{гребневой}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}.$$

Указанное решение добавляет шкалированную единичную матрицу $\lambda \mathbf{I}$ в $\mathbf{X}^T \mathbf{X}$ перед инверсией, тем самым гарантируя, что задача не является сингулярной, даже если $\mathbf{X}^T \mathbf{X}$ не имеет полного ранга. В самом начале, когда был введен указанный оценщик, это было одним из стимулов для его использования.

Гребневой штраф приводит к пропорциональной усадке всех параметров. В случае *ортонормированных входов* гребневые оценки являются просто шкалированной версией наименьших квадратов оценок, т. е.:

$$\hat{\beta}^{\text{гребневой}} = \frac{\hat{\beta}_{LS}}{1 + \lambda}.$$

Используя *сингулярное разложение* (singular value decomposition, SVD) входной матрицы \mathbf{X} , мы можем проникнуть в сущность того, как усадка влияет на входы в более общем случае, когда они не ортонормированы. Сингулярное разложение центрированной матрицы представляет главные компоненты матрицы (см. главу 11, посвященную неконтролируемому обучению), которые улавливают некоррелированные направления в столбцовом пространстве данных в порядке убывания дисперсии.

Гребневая регрессия сжимает коэффициенты на входных переменных, которые ассоциированы с направлениями в данных, имеющими меньшую дисперсию, больше, чем на входных переменных, которые коррелируют с направлениями, показывающими больше дисперсии. Следовательно, неявное допущение гребневой регрессии состоит в том, что направления в данных, которые варьируются больше всего, будут наиболее влиятельными или наиболее надежными при предсказании выхода.

Как работает регрессия лассо

Регрессия лассо, именуемая в обработке сигналов погоней за базисом, также сжимает (осаждает) коэффициенты, добавляя штраф в сумму квадратов остатков, но лассо-штраф имеет несколько иной эффект. Лассо-штраф представляет собой сумму абсолютных значений вектора коэффициентов, которая соответствует его норме L_1 . Следовательно, оценка лассо определяется по формуле:

$$\begin{aligned}\hat{\beta}^{\text{лассо}} &= \arg \min_{\beta^{\text{лассо}}} \sum_{i=1}^N \left[\left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right] = \\ &= \arg \min_{\beta^{\text{лассо}}} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) - \lambda |\beta|.\end{aligned}$$

Подобно гребневой регрессии, входы должны быть стандартизированы. Лассо-штраф делает решение нелинейным, и для коэффициентов нет выражения в замкнутой форме, как это было в гребневой регрессии. Вместо этого лассо-решение является задачей квадратичного программирования, и существуют эффективные алгоритмы, которые вычисляют всю траекторию коэффициентов, которые получаются для разных значений λ с той же вычислительной стоимостью, что и для гребневой регрессии.

Лассо-штраф имеет эффект постепенного сведения некоторых коэффициентов к нулю по мере увеличения регуляризации. По этой причине лассо может использоваться для непрерывного отбора подмножества признаков.

Как применять линейную регрессию для предсказания финансовых возвратов

Блокнот `linear_regression.ipynb` содержит примеры предсказания цен акций с использованием обычных наименьших квадратов (OLS) с помощью библиотек `StatsModels` и `sklearn`, а также гребневой и лассо-моделей. Он предназначен для работы на исследовательской платформе `Quantopian` и опирается на библиотеку `factor_library`, представленную в *главе 4*.

Подготовка данных

Нам нужно отобрать универсум долевого ценных бумаг и временной горизонт, построить и преобразовать альфа-факторы, которые мы будем использовать в качестве признаков, рассчитать форвардные возвраты, которые мы стремимся предсказать, и потенциально очистить наши данные.

Создание универсума и временной горизонт

Мы будем использовать данные долевого ценных бумаг за 2014 и 2015 гг. из кастомизированного универсума `Q100US`, который использует встроенные фильтры,

факторы и классификаторы для отбора 100 акций с самым высоким среднедолларовым объемом за последние 200 торговых дней, отфильтрованных по дополнительным критериям, принятым по умолчанию (см. документацию Quantopian, ссылка на которую приведена в репозитории GitHub). Указанный универсум динамически обновляется на основе фильтровочных критериев, вследствие чего, несмотря на то, что в любой заданной точке имеется 100 акций, в выборке может быть более 100 разных долевых ценных бумаг:

```
def Q100US():
    return filters.make_us_equity_universe(
        target_size=100,
        rankby=factors.AverageDollarVolume(window_length=200),
        mask=filters.default_us_equity_universe_mask(),
        groupby=classifiers.fundamentals.Sector(),
        max_group_weight=0.3,
        smoothing_func=lambda f: f.downsample('month_start'),
    )
```

Расчет целевого финансового возврата

Мы протестируем предсказания для различных прогнозных периодов `lookahead`, выявляя наилучшие периоды владения, которые генерируют наилучшую предсказуемость, измеряемую информационным коэффициентом. Более конкретно, мы вычисляем возвраты за 1, 5, 10 и 20 дней, используя встроенную функцию `Returns`, в результате получая более 50 тыс. наблюдений для универсума из 100 акций за два года (состоящие примерно по 252 торговых дней каждый):

```
START = '2014-01-01'
END = '2015-12-31'

UNIVERSE = Q100US()

lookahead = [1, 5, 10, 20]
returns = run_pipeline(Pipeline({'Returns': D'.format(i):
                                Returns(inputs=[USEquityPricing.close],
                                           window_length=i+1, mask=UNIVERSE)
                                for i in lookahead},
                           screen=UNIVERSE),
                      start_date=START,
                      end_date=END)
return_cols = ['Returns': D'.format(i) for i in lookahead]
returns.info()

<class 'pandas.core.frame.DataFrame'>
MultiIndex: 50362 entries, (2014-01-02 00:00:00+00:00, Equity(24 [AAPL]))
to (2015-12-31 00:00:00+00:00, Equity(47208 [GPRO]))
```

```
Data columns (total 4 columns):
Returns10D    50362 non-null float64
Returns1D     50362 non-null float64
Returns20D    50360 non-null float64
Returns5D     50362 non-null float64
dtypes: float64(4)
memory usage: 1.9+ MB
```

Отбор и преобразование альфа-факторов

Мы будем использовать более 50 признаков, которые охватывают широкий спектр факторов, основанных на рыночных, фундаментальных и альтернативных данных. Блокнот также включает кастомизированные преобразования, конвертирующие фундаментальные данные, которые, как правило, доступны с частотой ежеквартальной отчетности, в скользящие годовые итоги или средние во избежание чрезмерных сезонных колебаний.

После того как факторы были вычислены с помощью различных конвейеров, описанных в *главе 4*, мы комбинируем их с помощью метода `pd.concat()`, назначаем индексные имена и создаем категориальную переменную, отождествляющую актив для каждой точки данных:

```
data = pd.concat([returns,
                  value_factors,
                  momentum_factors,
                  quality_factors,
                  payout_factors,
                  growth_factors,
                  efficiency_factors,
                  risk_factors], axis=1).sortlevel()
data.index.names = ['date', 'asset']
data['stock'] = data.index.get_level_values('asset').map(lambda x: x.asset_name)
```

Очистка данных — пропущенные данные

На следующем этапе мы удаляем строки и столбцы, в которых не хватает более 20% наблюдений, что приводит к потере 6% наблюдений и трех столбцов:

```
rows_before, cols_before = data.shape
data = (data
        .dropna(axis=1, thresh=int(len(data)*.8))
        .dropna(thresh=int(len(data.columns) * .8)))
data = data.fillna(data.median())
rows_after, cols_after = data.shape
print('{:,d} rows and {:,d} columns dropped'
      .format(rows_before-rows_after, cols_before-cols_after))

2,985 rows and 3 columns dropped
```

В этом месте у нас имеются 51 признак и категориальный идентификатор акции:

```
data.sort_index(1).info()

<class 'pandas.core.frame.DataFrame'>
MultiIndex: 47377 entries, (2014-01-02 00:00:00+00:00, Equity(24 [AAPL]))
  to (2015-12-31 00:00:00+00:00, Equity(47208 [GPRO]))
Data columns (total 52 columns):
AssetToEquityRatio      47377 non-null float64
AssetTurnover           47377 non-null float64
CFO To Assets           47377 non-null float64
...
WorkingCapitalToAssets  47377 non-null float64
WorkingCapitalToSales   47377 non-null float64
stock                   47377 non-null object
dtypes: float64(51), object(1)
memory usage: 19.2+ MB
```

Разведывательный анализ данных

В случае линейных регрессионных моделей важно провести разведку корреляции между признаками, выявив проблемы, связанные с мультиколлинеарностью, и проверив корреляцию между признаками и целью. Указанный блокнот содержит кластерную карту, созданную средствами библиотеки Seaborn, которая показывает иерархическую структуру матрицы корреляции признаков. Она выявляет малое число высоко коррелированных кластеров.

Кодирование категориальных переменных с помощью фиктивных значений

Нам нужно конвертировать категориальную переменную `stock` в числовой формат, что позволит линейной регрессии ее обработать. Для этой цели мы используем кодирование фиктивных значений (или кодирование с одним активным состоянием), которое создает отдельные столбцы для каждого категориального уровня и отмечает наличие этого уровня в исходном категориальном столбце 1 и 0 в противном случае. Функция `pandas.get_dummies()` автоматизирует кодирование фиктивных значений. Она обнаруживает и правильно конвертирует столбцы типизированных объектов, как показано ниже. Например, если вам нужны фиктивные переменные для столбцов, содержащих целые числа, то вы можете идентифицировать их с помощью ключевого слова `columns`:

```
df = pd.DataFrame({'categories': ['A', 'B', 'C']})

categories
0      A
1      B
2      C
```

```
pd.get_dummies(df)
```

| | categories_A | categories_B | categories_C |
|---|--------------|--------------|--------------|
| 0 | 1 | 0 | 0 |
| 1 | 0 | 1 | 0 |
| 2 | 0 | 0 | 1 |

При конвертировании всех категорий в фиктивные переменные и оценивании модели с пересечением (как обычно) вы непреднамеренно создаете мультиколлинеарность: матрица теперь содержит избыточную информацию и больше не имеет полного ранга, т. е. становится сингулярной. Этого легко избежать, удалив один из новых индикаторных столбцов. Коэффициент на уровне отсутствующей категории теперь будет захвачен пересечением (которое всегда равно 1, когда все остальные категорийные фиктивные значения равны 0). Для внесения соответствующих поправок в фиктивные переменные следует использовать ключевое слово `drop_first`:

```
pd.get_dummies(df, drop_first=True)
```

| | categories_B | categories_C |
|---|--------------|--------------|
| 0 | 0 | 0 |
| 1 | 1 | 0 |
| 2 | 0 | 1 |

Для наших скомбинированных признаков и возвратов мы получаем 181 столбец, поскольку имеется более 100 акций, т. к. формулировка универсума автоматически обновляет подборку акций:

```
X = pd.get_dummies(data.drop(return_cols, axis=1))
```

```
X.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
MultiIndex: 47377 entries, (2014-01-02 00:00:00+00:00, Equity(24 [AAPL]))
```

```
to (2015-12-31 00:00:00+00:00, Equity(47208 [GPRO]))
```

```
Columns: 182 entries, DividendYield to stock_YELP INC
```

```
dtypes: float64(182)
```

```
memory usage: 66.1+ MB
```

Создание форвардных возвратов

Наша цель состоит в том, чтобы предсказать финансовые возвраты за определенный период владения. Следовательно, нам нужно выровнять признаки, содержащие значения возвратов, с точками данных соответствующего возврата в перспективе на 1, 5, 10 или 20 дней вперед для каждой доленой ценной бумаги. Мы достигаем этого путем совмещения методов библиотеки `pandas` `groupBy()` и `shift()` следующим образом:

```
y = data.loc[:, return_cols]
```

```
shifted_y = []
```



```

for col in y.columns:
    t = int(re.search(r'\d+', col).group(0))
    shifted_y.append(y.groupby(level='asset')['Returns{}'.format(t)]
                     .shift(-t).to_frame(col))
y = pd.concat(shifted_y, axis=1)
y.info()

<class 'pandas.core.frame.DataFrame'>
MultiIndex: 47377 entries, (2014-01-02 00:00:00+00:00, Equity(24 [AAPL]))
to (2015-12-31 00:00:00+00:00, Equity(47208 [GPRO]))
Data columns (total 4 columns):
Returns1D      47242 non-null float64
Returns5D      46706 non-null float64
Returns10D     46036 non-null float64
Returns20D     44696 non-null float64
dtypes: float64(4)
memory usage: 1.8+ MB

```

Для каждого ряда с финансовыми возвратами сейчас существует разное число наблюдений, поскольку сдвиг вперед создал недостающие значения в хвостовой части по каждой доле ценной бумаге.

Линейная регрессия методом OLS с использованием библиотеки StatsModels

Мы можем вычислить линейную регрессионную модель, используя обычные наименьшие квадраты (OLS) с помощью библиотеки StatsModels, как было показано ранее. Мы выбираем форвардный возврат, например за 10-дневный период владения, удаляем выбросы ниже 2,5%-го и 97,5%-го перцентилей, и соответствующим образом вписываем модель:

```

target = 'Returns10D'
model_data = pd.concat([y[[target]], X], axis=1).dropna()
model_data = model_data[model_data[target]
                        .between(model_data[target].quantile(.025),
                                model_data[target].quantile(.975))]

model = OLS(endog=model_data[target],
            exog=model_data.drop(target, axis=1))
trained_model = model.fit()
trained_model.summary()

```

Диагностическая статистика

Вся сводная информация размещена в блокноте с целью экономии пространства из-за большого числа переменных. Диагностическая статистика показывает, что

при наличии высокого p -значения в статистическом показателе Харке — Бера гипотеза о нормальном распределении остатков не может быть отклонена.

Однако статистический показатель Дарбина — Уотсона находится низко, на уровне 1,5, и поэтому мы можем спокойно отклонить нулевую гипотезу об отсутствии автокорреляции на уровне 5%. Следовательно, стандартные ошибки, по всей видимости, положительно коррелируют. Если бы наша цель заключалась в том, чтобы понять, какие факторы значимо ассоциированы с форвардными возвратами, то нам нужно было бы перезапустить регрессию с использованием робастных стандартных ошибок (параметр в методе `fit()` библиотеки `StatsModels`) или использовать другой метод, такой как панельная модель, которая допускает более сложную ковариацию ошибок.

Линейная регрессия по методу OLS с использованием библиотеки `sklearn`

Поскольку библиотека `sklearn` приспособлена к нуждам предсказания, мы вычислим линейную регрессионную модель, основываясь на ее предсказательной результативности с использованием перекрестного контроля.

Кастомизированный перекрестный контроль временного ряда

Наши данные состоят из сгруппированного временного ряда и требуют кастомизированной функции перекрестного контроля, предоставляющей тренировочные и тестовые индексы, которые обеспечивают, чтобы тестовые данные следовали сразу за тренировочными данными для каждой долевой ценной бумаги, и чтобы мы невольно не создавали смещение или утечку.

Мы можем достичь этого, используя следующую ниже функцию, которая возвращает генератор, выдающий пары тренировочных и тестовых дат. Набор тренировочных дат обеспечивает минимальную продолжительность тренировочных периодов. Число пар зависит от параметра `nfolds`. Разные тестовые периоды не накладываются и располагаются в конце периода, доступного в данных. После использования тестового периода он становится частью тренировочных данных, которые соответствующим образом увеличиваются в размере:

```
def time_series_split(d=model_data, nfolds=5, min_train=21):
    """Сгенерировать тренировочные/тестовые даты для блоков nfolds
        по крайней мере с min_train тренировочных наблюдений
    """

    train_dates = d[:min_train].tolist()
    n = int(len(dates)/(nfolds + 1)) + 1
    test_folds = [d[i:i + n]
                  for i in range(min_train, len(d), n)]
```

```

for test_dates in test_folds:
    if len(train_dates) > min_train:
        yield train_dates, test_dates
    train_dates.extend(test_dates)

```

Отбор признаков и цели

Нам нужно отобрать соответствующий ряд финансовых возвратов (мы снова будем использовать 10-дневный период владения) и удалить выбросы. Мы также конвертируем возвраты в логарифмические возвраты следующим образом:

```

target = 'Returns10D'
outliers = .01
model_data = pd.concat([y[[target]], X],
                        axis=1).dropna().reset_index('asset', drop=True)
model_data = model_data[model_data[target].between(*model_data[target]
                                                    .quantile([outliers, 1-outliers]).values)]

model_data[target] = np.log1p(model_data[target])
features = model_data.drop(target, axis=1).columns
dates = model_data.index.unique()

print(model_data.info())

```

```

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 45114 entries, 2014-01-02 to 2015-12-16
Columns: 183 entries, Returns10D to stock_YELP INC
dtypes: float64(183)
memory usage: 63.3 MB
None

```

Перекрестный контроль модели

Мы будем использовать 250 блоков, в целом предсказывая около двух дней форвардных возвратов, следующих после исторических тренировочных данных, которые будут постепенно увеличиваться в длину. Каждая итерация получает надлежащие тренировочные и тестовые даты из нашей собственной функции перекрестного контроля, отбирает соответствующие признаки и цели, а затем по мере исполнения тренирует и предсказывает. Мы улавливаем корень квадратный из среднеквадратической ошибки (RMSE), а также ранговую корреляцию Спирмена между фактическими и предсказанными значениями:

```

nfold = 250
lr = LinearRegression()

test_results, result_idx, preds = [], [], pd.DataFrame()
for train_dates, test_dates in time_series_split(dates, nfold=nfold):

```

```

X_train = model_data.loc[idx[train_dates], features]
y_train = model_data.loc[idx[train_dates], target]
lr.fit(X=X_train, y=y_train)

X_test = model_data.loc[idx[test_dates], features]
y_test = model_data.loc[idx[test_dates], target]
y_pred = lr.predict(X_test)

rmse = np.sqrt(mean_squared_error(y_pred=y_pred, y_true=y_test))
ic, pval = spearmanr(y_pred, y_test)

test_results.append([rmse, ic, pval])
preds = preds.append(y_test.to_frame('actuals')
                    .assign(predicted=y_pred))
result_idx.append(train_dates[-1])

```

Тестовые результаты — информационный коэффициент и RMSE

Мы уловили тестовые предсказания из 250 блоков и можем вычислить как совокупное, так и 21-дневное скользящее среднее:

```

test_result = pd.DataFrame(test_results, columns=['rmse', 'ic', 'pval'],
                          index=result_idx)

fig, axes = plt.subplots(nrows=2)
rolling_result = test_result.rolling(21).mean()
rolling_result[['ic', 'pval']].plot(ax=axes[0], title='Информационный коэффициент')
axes[0].axhline(test_result.ic.mean(), lw=1, ls='--', color='k')
rolling_result[['rmse']].plot(ax=axes[1],
                             title='Корень квадратный из среднеквадратической ошибки')
axes[1].axhline(test_result.rmse.mean(), lw=1, ls='--', color='k')

```

В результате мы получаем следующую ниже диаграмму на рис. 7.2, которая подчеркивает отрицательную корреляцию информационного коэффициента IC и оценки RMSE и их соответствующих значений.

На рис. 7.3 мы видим, что за весь период информационный коэффициент, измеряемый ранговой корреляцией фактических и предсказанных возвратов, слабо положителен и статистически значим:

```

preds_cleaned = preds[(preds.predicted
                      .between(*preds.predicted.quantile([.001, .999]).values))]
sns.jointplot(x='actuals', y='predicted',
              data=preds_cleaned, stat_func=spearmanr);

```



Рис. 7.2. Диаграмма, подчеркивающая отрицательную корреляцию информационного коэффициента IC и оценки RMSE

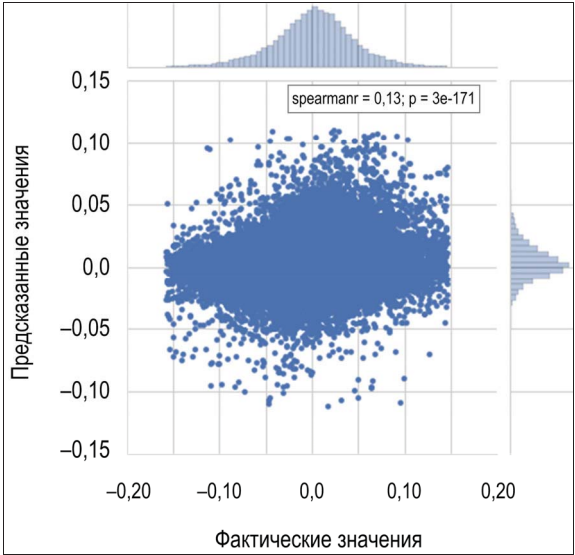


Рис. 7.3. Информационный коэффициент, измеряемый ранговой корреляцией Спирмена

Гребневая регрессия с использованием библиотеки sklearn

Для гребневой регрессии необходимо отрегулировать регуляризационный параметр с ключевым словом `alpha`, соответствующий лямбде λ , которую мы использовали ранее. Мы испытаем 21 значение между 10^{-5} и 10^5 пошагово, используя логарифмические шаги.

Чувствительность гребневой регрессии к шкале требует выполнения стандартизации входов с использованием шкалировщика `StandardScaler`. Обратите внимание, что мы всегда усваиваем среднее значение и стандартное отклонение из тренировочного набора с помощью метода `fit_transform()`, а затем применяем эти усвоенные параметры к тестовому набору с помощью метода `transform()`.

Настройка регуляризационных параметров с помощью перекрестного контроля

Затем мы приступаем к перекрестному контролю значений гиперпараметров, снова используя 250 блоков следующим образом:

```
nfolds = 250
alphas = np.logspace(-5, 5, 11)
scaler = StandardScaler()

ridge_result, ridge_coeffs = pd.DataFrame(), pd.DataFrame()
for i, alpha in enumerate(alphas):
    print i,
    coeffs, test_results = [], []
    lr_ridge = Ridge(alpha=alpha)
    for train_dates, test_dates in time_series_split(dates, nfolds=nfolds):

        X_train = model_data.loc[idx[train_dates], features]
        y_train = model_data.loc[idx[train_dates], target]
        lr_ridge.fit(X=scaler.fit_transform(X_train), y=y_train)
        coeffs.append(lr_ridge.coef_)

        X_test = model_data.loc[idx[test_dates], features]
        y_test = model_data.loc[idx[test_dates], target]
        y_pred = lr_ridge.predict(scaler.transform(X_test))

        rmse = np.sqrt(mean_squared_error(y_pred=y_pred, y_true=y_test))
        ic, pval = spearmanr(y_pred, y_test)

        test_results.append([train_dates[-1], rmse, ic, pval, alpha])
test_results = pd.DataFrame(test_results,
                             columns=['date', 'rmse', 'ic', 'pval', 'alpha'])
```

```
ridge_result = ridge_result.append(test_results)
ridge_coefs[alpha] = np.mean(coeffs, axis=0)
```

Результаты перекрестного контроля и траектории гребневых коэффициентов

Теперь мы можем построить график информационного коэффициента, полученный для значения каждого гиперпараметра, а также визуализировать то, как значения коэффициентов эволюционируют по мере увеличения регуляризации. Результаты показывают, что мы получаем наибольшее значение информационного коэффициента для значения $\lambda = 10$ (рис. 7.4). Для этого уровня регуляризации правая панель свидетельствует о том, что коэффициенты уже значительно сжаты по сравнению с (почти) неограниченной моделью с $\lambda = 10^{-5}$.

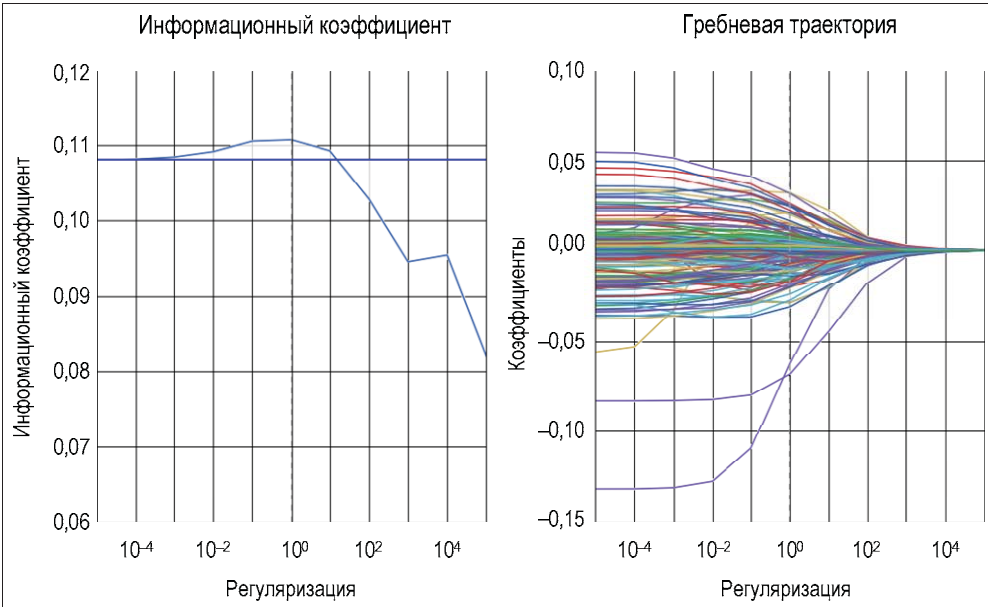


Рис. 7.4. Результаты перекрестного контроля и траектории гребневых коэффициентов

Десятка ведущих коэффициентов

Стандартизация коэффициентов позволяет сделать выводы об их относительной важности путем сравнения их абсолютной величины. Приведем 10 наиболее релевантных коэффициентов (рис. 7.5).

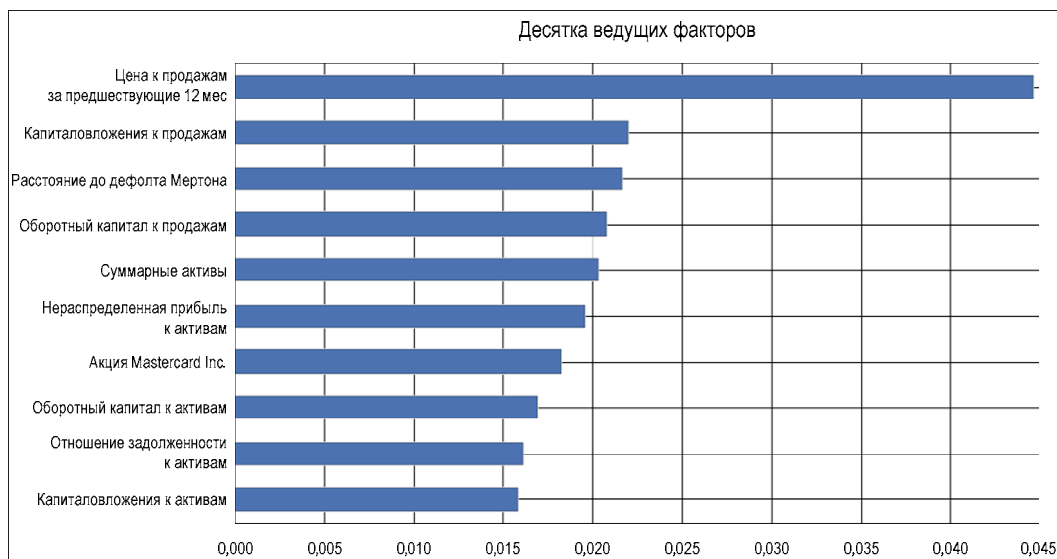


Рис. 7.5. Десятка ведущих коэффициентов

Лассо-регрессия с использованием библиотеки `sklearn`

Реализация лассо очень похожа на гребневую модель, которую мы только что выполнили. Главное отличие заключается в том, что лассо-регрессия должна прийти к решению, используя итеративный координатный спуск, тогда как гребневая регрессия может опираться на решение в замкнутой форме:

```

n folds = 250
alphas = np.logspace(-8, -2, 13)
scaler = StandardScaler()

lasso_results, lasso_coefs = pd.DataFrame(), pd.DataFrame()
for i, alpha in enumerate(alphas):
    print i,
    coefs, test_results = [], []
    lr_lasso = Lasso(alpha=alpha)
    for i, (train_dates, test_dates) in enumerate(time_series_split(dates,
                                                                    n folds=n folds)):
        X_train = model_data.loc[idx[train_dates], features]
        y_train = model_data.loc[idx[train_dates], target]
        lr_lasso.fit(X=scaler.fit_transform(X_train), y=y_train)

        X_test = model_data.loc[idx[test_dates], features]
        y_test = model_data.loc[idx[test_dates], target]
        y_pred = lr_lasso.predict(scaler.transform(X_test))

```



```

rmse = np.sqrt(mean_squared_error(y_pred=y_pred, y_true=y_test))
ic, pval = spearmanr(y_pred, y_test)

coeffs.append(lr_lasso.coef_)
test_results.append([train_dates[-1], rmse, ic, pval, alpha])
test_results = pd.DataFrame(test_results,
                             columns=['date', 'rmse', 'ic', 'pval', 'alpha'])
lasso_results = lasso_results.append(test_results)
lasso_coeffs[alpha] = np.mean(coeffs, axis=0)

```

Перекрытно-контрольный информационный коэффициент и траектория лассо-регрессии

Как и прежде, мы можем построить средний информационный коэффициент для всех тестовых наборов, используемых во время перекрытного контроля. Мы снова видим, что регуляризация улучшает информационный коэффициент над неограниченной моделью, обеспечивая наилучший вневыборочный результат на уровне $\lambda = 10^{-5}$. Оптимальное значение регуляризации сильно отличается от гребневой регрессии, поскольку вместо возведенных в квадрат значений штраф состоит из суммы абсолютных значений относительно малых значений коэффициентов. Мы также видим, что для этого уровня регуляризации коэффициенты были сжаты аналогичным образом, как и в случае гребневой регрессии (рис. 7.6).

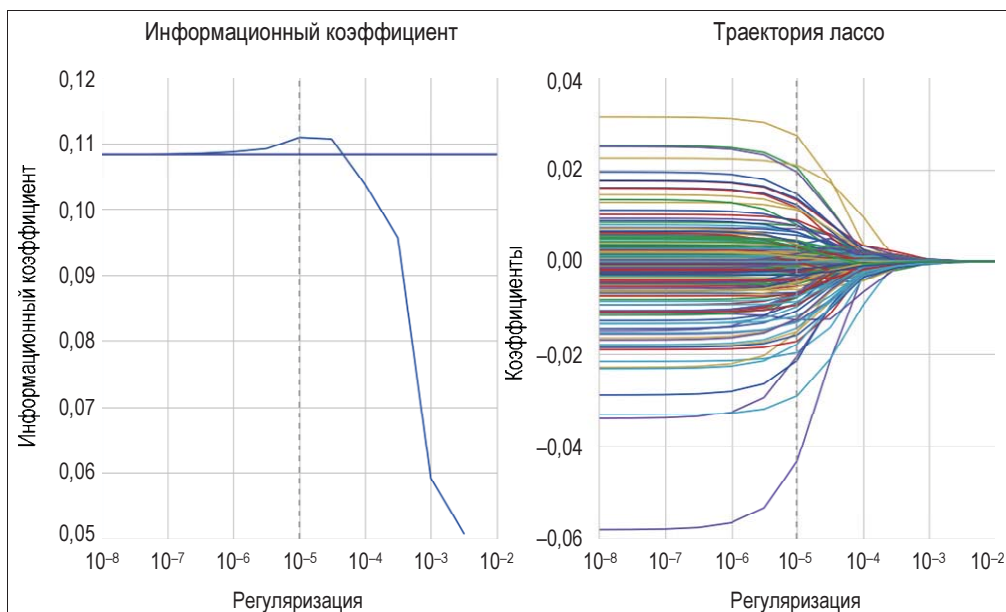


Рис. 7.6. Перекрытно-контрольный информационный коэффициент и траектория лассо-регрессии

В итоге гребневая регрессия и лассо-регрессия дадут аналогичные результаты. Гребневая регрессия часто вычисляется быстрее, но лассо-регрессия также дает непрерывный выбор подмножества признаков, постепенно сводя коэффициенты к нулю, тем самым устраняя признаки.

Линейная классификация

Обсуждавшаяся до сих пор линейная регрессионная модель предполагает количественную переменную отклика. В этом разделе мы сосредоточимся на подходах к моделированию качественных выходных переменных для статистического вывода и предсказания, т. е. процесса, который именуется *классификацией* и на практике происходит даже чаще, чем регрессия.

Предсказание качественного отклика для точки данных называется *классификацией* этого наблюдения, поскольку оно предусматривает отнесение наблюдения к категории или классу. На практике классификационные методы часто предсказывают вероятности для каждой из категорий качественной переменной, а затем используют эту вероятность для принятия решения о надлежащей классификации.

Мы могли бы подойти к классификационной задаче, игнорируя тот факт, что выходная переменная принимает дискретные значения, и применить линейную регрессионную модель, попытавшись предсказать категориальный выход с помощью нескольких входных переменных. Однако легко сконструировать примеры, где этот метод показывает очень слабые результаты. Более того, для модели не имеет интуитивного смысла производить значения больше 1 или меньше 0, когда мы знаем, что $y \in [0; 1]$.

Для предсказания качественного отклика существует целый ряд других классификационных методов или классификаторов, которые имеются в распоряжении. В этом разделе мы представим широко используемую логистическую регрессию, которая тесно связана с линейной регрессией. Более сложные методы мы рассмотрим в следующих главах, посвященных обобщенным аддитивным моделям, включающим деревья решений и случайные леса, а также градиентно-бустинговые машины и нейронные сети.

Логистическая регрессионная модель

Логистическая регрессионная модель возникает из желания моделировать вероятности выходных классов с учетом функции, линейной в x , подобно линейной регрессионной модели, в то же время обеспечивая, чтобы они в сумме составляли единицу и оставались в промежутке $[0; 1]$, как и ожидается от вероятностей.

В данном разделе мы введем целевую и функциональную формы логистической регрессионной модели и опишем метод тренировки. Затем мы проиллюстрируем способ применения логистической регрессии для статистического вывода с макроданными, используя библиотеку StatsModels, и покажем, как предсказывать цено-

вые движения с помощью регуляризованной логистической регрессии, которая реализована в библиотеке `sklearn`.

Целевая функция

Для иллюстрации мы будем использовать выходную переменную y , которая принимает значение 1, если возврат от акции является положительным в течение заданного временного горизонта d , и 0 в противном случае:

$$y_t = \begin{cases} 1, & \text{при } r_{t+d} > 0; \\ 0, & \text{в противном случае.} \end{cases}$$

Мы могли бы легко расширить y до трех категорий, где 0 и 2 отражают отрицательные и положительные ценовые движения за пределами определенного порога, и 1 в противном случае. Однако вместо моделирования выходной переменной y логистическая регрессия моделирует вероятность того, что y принадлежит к любой из категорий, заданных вектором альфа-факторов или признаков \mathbf{x}_t . Другими словами, логистическая регрессия моделирует вероятность того, что цена акции пойдет вверх в зависимости от значений переменных, включенных в модель:

$$p(\mathbf{x}_t) = \Pr(y_t = 1 | \mathbf{x}_t).$$

Логистическая функция

Для того чтобы не дать модели производить значения вне промежутка $[0; 1]$, мы должны моделировать $p(x)$ с помощью функции, которая дает выходы только между 0 и 1 по всей области x . Логистическая функция удовлетворяет этому требованию и всегда создает S-образную кривую (см. примеры из блокнота), и поэтому, независимо от значения x , мы будем получать разумное предсказание:

$$p(\mathbf{x}) = \frac{e^{\beta_0 + \sum_{i=1}^p \beta_i x_i}}{1 + e^{\beta_0 + \sum_{i=1}^p \beta_i x_i}} = \frac{e^{x\boldsymbol{\beta}}}{1 + e^{x\boldsymbol{\beta}}}.$$

Здесь вектор \mathbf{x} содержит 1 для пересечения, усвоенного первой компонентой $\boldsymbol{\beta}$ — β_0 . Мы можем преобразовать это выражение так, чтобы выделить часть, которая выглядит как линейная регрессия:

$$\underbrace{\frac{p(\mathbf{x})}{1 - p(\mathbf{x})}}_{\text{шансы}} = e^{\beta_0 + \sum_{i=1}^p \beta_i x_i} \Leftrightarrow \underbrace{\log\left(\frac{p(\mathbf{x})}{1 - p(\mathbf{x})}\right)}_{\text{логит}} = \beta_0 + \sum_{i=1}^p \beta_i x_i.$$

Величина $p(x)/[1 - p(x)]$ называется *шансами*, являясь альтернативным способом выражения вероятностей, которые мы знаем из азартных игр. Данная величина может принимать любые значения шансов между 0 и ∞ , где из низких значений

также вытекают низкие вероятности, а из высоких значений — высокие вероятности.

Логит⁹ также называется log-odds (логарифм шансов). Следовательно, логистическая регрессия представляет собой логит, линейный в x и во многом похожий на рассмотренную ранее линейную регрессию.

Оценивание максимального правдоподобия

Коэффициентный вектор β должен оцениваться с использованием имеющихся тренировочных данных. Хотя для вписывания логистической регрессионной модели мы могли бы использовать (нелинейные) наименьшие квадраты, предпочтительнее использовать более общий метод максимального правдоподобия, поскольку он имеет более качественные статистические свойства. Как мы только что рассмотрели, исходной идеей использования максимального правдоподобия для вписывания логистической регрессионной модели является поиск оценок для β таких, чтобы предсказанная вероятность $\hat{p}(x)$ максимально близко соответствовала фактическому выходному результату. Другими словами, мы пытаемся найти β , такие, что эти оценки дают число, близкое 1 для всех случаев, когда цена акции поднялась, и число, близкое к 0, в противном случае. Более формально, мы стремимся максимизировать функцию правдоподобия:

$$\max_{\beta} \mathcal{L}(\beta) = \prod_{i: y_i=1} p(x_i) \prod_{i': y_{i'}=0} (1 - p(x_{i'})).$$

Проще работать с суммами, чем с произведениями, и поэтому мы берем логарифм с обеих сторон, получив логарифмическую функцию правдоподобия и соответствующее определение логистических регрессионных коэффициентов:

$$\beta_{ML} = \arg \max \log \mathcal{L}(\beta) = \sum_{i=1}^N (y_i \log p(x_i, \beta) + (1 - y_i) \log (1 - p(x_i, \beta))).$$

Максимизация этого уравнения путем приравнивания производных \mathcal{L} по β нулю дает $p+1$ так называемых уравнений вклада¹⁰, нелинейных в параметрах, которые могут быть решены с помощью итерационных численных методов для вогнутой логарифмической функции правдоподобия.

⁹ Логит (logit), или логит-преобразование — это функция, которая увязывает вероятность принадлежности классу с интервалом $]-\infty, +\infty[$ (вместо промежутка от 0 до 1). Иногда называется ненормализованным логарифмом вероятности. — *Прим. перев.*

¹⁰ В статистике вклад (score), или уравнение вклада, указывает на степень чувствительности функции правдоподобия $\mathcal{L}(\theta; X)$ к своему параметру θ . А именно, вклад для параметра θ является градиентом логарифмической вероятности по отношению к θ .

См. [https://en.wikipedia.org/wiki/Score_\(statistics\)](https://en.wikipedia.org/wiki/Score_(statistics)). — *Прим. перев.*

Как проводить статистический вывод с помощью библиотеки StatsModels

Мы проиллюстрируем применение логистической регрессии с помощью библиотеки StatsModels на основе простого встроенного набора данных, содержащего квартальные макроданные США с 1959 по 2009 гг. (подробности см. в блокноте `logistic_regression_macro_data.ipynb`).

Переменные и их преобразования перечислены в табл. 7.2.

Таблица 7.2. Переменные квартальных макроданных

| Переменная | Описание | Преобразование |
|------------|---|--------------------|
| realgdp | Реальный валовой внутренний продукт | Годовой темп роста |
| realcons | Реальные расходы на личное потребление | Годовой темп роста |
| realinv | Реальные валовые частные внутренние инвестиции | Годовой темп роста |
| realgovt | Реальные федеральные расходы и валовые инвестиции | Годовой темп роста |
| realdpi | Реальные частные располагаемые доходы | Годовой темп роста |
| m1 | Номинальная денежная масса M1 | Годовой темп роста |
| tbilrate | Ставка по трехмесячным казначейским векселям | Уровень |
| unemp | Безработица с учетом сезонных колебаний (%) | Уровень |
| infl | Темп инфляции | Уровень |
| realint | Реальная процентная ставка | Уровень |

Для получения бинарной целевой переменной мы вычисляем 20-квартальное скользящее среднее годовых темпов роста квартального реального ВВП. Затем мы закрепляем значение 1, если текущий рост превышает скользящее среднее, и 0 в противном случае. Наконец, мы сдвигаем индикаторные переменные, выравнивая итоговый результат следующего квартала с текущим кварталом.

Мы используем пересечение, конвертируем квартальные значения в фиктивные переменные и тренируем логистическую регрессионную модель следующим образом:

```
import statsmodels.api as sm

data = pd.get_dummies(data.drop(drop_cols, axis=1),
                        columns=['quarter'], drop_first=True).dropna()
model = sm.Logit(data.target, sm.add_constant(data.drop('target',
                                                         axis=1)))

result = model.fit()
result.summary()
```

В результате для нашей модели с 198 наблюдениями и 13 переменными мы получаем статистическую сводку, включая пересечение:

Logit Regression Results

| | | | |
|----------------|------------------|-------------------|-----------|
| Dep. Variable: | target | No. Observations: | 198 |
| Model: | Logit | Df Residuals: | 185 |
| Method: | MLE | Df Model: | 12 |
| Date: | Tue, 07 May 2019 | Pseudo R-squ.: | 0.5022 |
| Time: | 12:24:25 | Log-Likelihood: | -67.907 |
| converged: | True | LL-Null: | -136.42 |
| | | LLR p-value: | 2.375e-23 |

| | coef | std err | z | P> z | [0.025 | 0.975] |
|-----------|----------|---------|--------|-------|----------|---------|
| const | -8.5881 | 1.908 | -4.502 | 0.000 | -12.327 | -4.849 |
| realcons | 130.1446 | 26.633 | 4.887 | 0.000 | 77.945 | 182.344 |
| realinv | 18.8414 | 4.053 | 4.648 | 0.000 | 10.897 | 26.786 |
| realgovt | -19.0318 | 6.010 | -3.166 | 0.002 | -30.812 | -7.252 |
| realdpi | -52.2473 | 19.912 | -2.624 | 0.009 | -91.275 | -13.220 |
| m1 | -1.3462 | 6.177 | -0.218 | 0.827 | -13.453 | 10.761 |
| tbilrate | 60.8607 | 44.350 | 1.372 | 0.170 | -26.063 | 147.784 |
| unemp | 0.9487 | 0.249 | 3.818 | 0.000 | 0.462 | 1.436 |
| infl | -60.9647 | 44.362 | -1.374 | 0.169 | -147.913 | 25.984 |
| realint | -61.0453 | 44.359 | -1.376 | 0.169 | -147.987 | 25.896 |
| quarter_2 | 0.1128 | 0.618 | 0.182 | 0.855 | -1.099 | 1.325 |
| quarter_3 | -0.1991 | 0.609 | -0.327 | 0.744 | -1.393 | 0.995 |
| quarter_4 | 0.0007 | 0.608 | 0.001 | 0.999 | -1.191 | 1.192 |

Статистическая сводка показывает, что модель была натренирована с использованием максимального правдоподобия и обеспечивает максимизированное значение логарифмической функции правдоподобия при $-67,9$.

Значение LL-Null, равное $-136,42$, является результатом максимизированной логарифмической функции правдоподобия, когда включено только пересечение. Оно формирует основу для статистического показателя псевдо- R^2 и статистической проверки отношения логарифмических правдоподобий LLR¹¹.

Статистический показатель псевдо- R^2 является заменой знакомого R^2 из метода наименьших квадратов. Он вычисляется на основе отношения максимизированных

¹¹ Статистическая проверка отношения логарифмического правдоподобия (log-likelihood ratio test, LLR) обозначает проверку значимости двух биномиальных распределений, также именуемой проверкой статистического показателя G -квадрат. См. <https://mahout.apache.org/users/clustering/llr---log-likelihood-ratio.html>. — Прим. перев.

логарифмических функций правдоподобия для нулевой модели m_0 и полной модели m_1 следующим образом:

$$\rho^2 = 1 - \frac{\log \mathcal{L}(m_1^*)}{\log \mathcal{L}(m_0^*)}.$$

Значения варьируются от 0 (когда модель не улучшает правдоподобие) до 1, где модель вписана идеально, и логарифмическое правдоподобие максимизируется при 0. Следовательно, более высокие значения указывают на более качественную подгонку.

Статистическая проверка LLR обычно сравнивает более ограниченную модель и вычисляется как:

$$\text{LLR} = -2 \log \left(\mathcal{L}(m_0^*) / \mathcal{L}(m_1^*) \right) = 2 \left(\log \mathcal{L}(m_1^*) - \log \mathcal{L}(m_0^*) \right).$$

Нулевая гипотеза состоит в том, что ограниченная модель работает лучше, но низкое p -значение свидетельствует о том, что мы можем отклонить эту гипотезу и предпочесть нулевой модели полную модель. Это похоже на F -тест для линейной регрессии (где мы также можем использовать LLR-тест при оценивании модели с использованием метода максимального правдоподобия MLE).

Статистический показатель z играет ту же роль, что и показатель t в выходе линейной регрессии, и в равной степени вычисляется как отношение оценки коэффициента и ее стандартной ошибки. P -значения также указывают на вероятность наблюдать проверочный статистический показатель при нулевой гипотезе $H_0: \beta = 0$ о том, что популяционный коэффициент равен нулю. Мы можем отклонить эту гипотезу для `intercept`, `realcons`, `realinv`, `realgovt`, `realdpi` и `unemp`.

Как применять логистическую регрессию для предсказания

Лассо-штраф L_1 и гребневой штраф L_2 могут использоваться вместе с логистической регрессией. Они имеют тот же усадочный эффект, что мы только что обсуждали, и лассо может опять же использоваться для отбора переменных с любой линейной регрессионной моделью.

Как и в случае линейной регрессии, важно простандартизировать входные переменные, поскольку регуляризованные модели чувствительны к шкале данных. Гиперпараметрическая регуляризация требует своей настройки с помощью перекрестной проверки, как в случае линейной регрессии.

Как предсказывать ценовые движения с помощью библиотеки `sklearn`

Мы продолжим пример предсказания цены, но теперь мы бинаризуем переменную результата, вследствие чего она принимает значение 1 всякий раз, когда 10-дневный возврат является положительным и 0 в противном случае; см. блокнот `logistic_regression.ipynb` в подпапке `stock_price_prediction`:

```

target = 'Returns10D'
label = (y[target] > 0).astype(int).to_frame(target)

model_data = pd.concat([label, X],
                        axis=1).dropna().reset_index('asset', drop=True)

features = model_data.drop(target, axis=1).columns
dates = model_data.index.unique()

```

С этой новой категориальной переменной выхода теперь мы можем натренировать логистическую регрессию, используя принятую по умолчанию регуляризацию L_2 . Для логистической регрессии регуляризация формулируется обратно линейной регрессии: из более высоких значений λ вытекает более низкая регуляризация, и наоборот. Мы оцениваем 11 параметрических значений с помощью перекрестного контроля следующим образом:

```

n folds = 250
Cs = np.logspace(-5, 5, 11)
scaler = StandardScaler()

logistic_results, logistic_coeffs = pd.DataFrame(), pd.DataFrame()
for C in Cs:
    coeffs = []
    log_reg = LogisticRegression(C=C)
    for i, (train_dates, test_dates) in enumerate(time_series_split(dates,
                                                                    n folds=n folds)):

        X_train = model_data.loc[idx[train_dates], features]
        y_train = model_data.loc[idx[train_dates], target]
        log_reg.fit(X=scaler.fit_transform(X_train), y=y_train)

        X_test = model_data.loc[idx[test_dates], features]
        y_test = model_data.loc[idx[test_dates], target]
        y_pred = log_reg.predict_proba(scaler.transform(X_test))[:, 1]

        coeffs.append(log_reg.coef_.squeeze())
    logistic_results = (logistic_results
                        .append(y_test
                               .to_frame('actuals')
                               .assign(predicted=y_pred, C=C)))
logistic_coeffs[C] = np.mean(coeffs, axis=0)

```

Затем мы используем отметку AUC, `roc_auc_score`, рассмотренную в предыдущей главе, для сравнения точности предсказания по различным регуляризационным параметрам:

```

auc_by_C = logistic_results.groupby('C').apply(lambda x:
                                                roc_auc_score(y_true=x.actuals.astype(int),
                                                                y_score=x.predicted))

```


Мы снова можем построить результат AUC для диапазона гиперпараметрических значений рядом с траекторией коэффициентов. Указанный диапазон демонстрирует улучшение точности предсказания, поскольку коэффициенты немного сжались при оптимальном значении регуляризации, равном 10^2 (рис. 7.7).

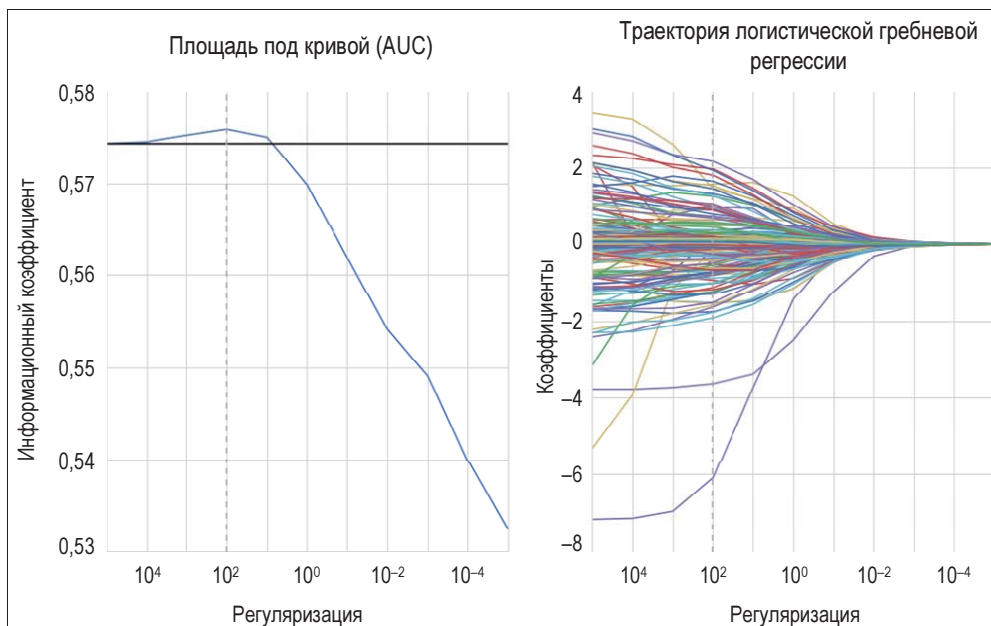


Рис. 7.7. Площадь под кривой (AUC) и траектория логистической гребневой регрессии

Резюме

В этой главе мы представили первые автоматически обучающиеся модели, использующие важный базовый случай линейных моделей для регрессионных и классификационных задач. Мы развели формулировку целевых функций для обеих этих задач, узнали о различных методах тренировки и о том, как применять указанную модель для статистического вывода и предсказания.

Мы применили эти новые методы машинного обучения для оценивания линейных факторных моделей, которые очень полезны для управления рисками, анализа новых альфа-факторов и результативности атрибутов. Мы также применили линейную регрессию и классификацию для выполнения первой предсказательной задачи по предсказанию возвратов от акций в абсолютном и направленном выражении.

В следующей главе мы рассмотрим важную тему линейных моделей временных рядов, служащих для улавливания регулярностей внутрирядовой корреляции в одномерном и многомерном случаях. Мы также узнаем о новых стратегиях торговли, разведав парную торговлю на основе концепции коинтеграции, которая улавливает динамическую корреляцию между двумя рядами с ценами на акции.

8

Модели временных рядов

В предыдущей главе мы сосредоточились на линейных моделях, приспособленных к срезовым (кросс-секционным) данным, где входные данные относятся к тому же периоду времени, что и выход, который они призваны объяснять или предсказывать. В этой главе мы сосредоточимся на данных временного ряда, где наблюдения различаются периодом, что также создает естественную упорядоченность. Наша цель будет состоять в том, чтобы выявлять исторические регулярности в данных и использовать эти регулярности для предсказания поведения временного ряда в будущем.

В предыдущей главе мы уже сталкивались с панельными данными в срезовой (кросс-секционной) и рядовой размерности и узнали, как регрессия Фама — Макбета оценивает стоимость принятия определенных факторных рисков во временной динамике и по активам. Вместе с тем связь между финансовыми возвратами во временной динамике обычно довольно низкая, поэтому указанная процедура может в значительной степени игнорировать временную размерность. Модели этой главы сосредоточены на моделях временного ряда, где прошлые значения содержат предсказательные сигналы о будущих событиях. Модели временных рядов также могут предсказывать признаки, которые затем используются в срезовых (кросс-секционных) моделях.

Более конкретно, в этой главе мы сосредоточимся на моделях, которые извлекают сигналы из ранее наблюдавшихся данных для предсказания будущих значений для того же самого временного ряда. Временной аспект торговли делает применение моделей временных рядов к рыночным, фундаментальным и альтернативным данным очень популярным. Данные временных рядов будут становиться еще более распространенными, поскольку все более широкий спектр подключенных к сети устройств собирает регулярные количественные измерения, которые могут содержать предсказательные сигналы. Ключевые сферы применения охватывают предсказание возвратов от активов, корреляций, или ковариаций, или волатильности.

В этой главе мы остановимся на линейных моделях временных рядов в качестве отправной точки для нелинейных моделей, таких как рекуррентные или сверточные нейронные сети, которые мы применим к данным временных рядов в *главах 17–21*.

Мы введем инструменты для диагностики характеристик временного ряда, включая стационарность, и извлечем признаки, которые улавливают потенциальные регулярности. Затем мы введем модели одномерных и многомерных временных рядов и применим их для предсказания регулярностей макроданных и волатильности. Мы закончим главу понятием коинтеграции и тем, как ее применять для разработки стратегии парной торговли.

В частности, мы рассмотрим следующие темы:

- ◆ как использовать анализ временного ряда для диагностики статистических показателей, информирующих процесс моделирования;
- ◆ как оценивать и диагностировать авторегрессионные модели и модели скользящего среднего временного ряда;
- ◆ как строить модели авторегрессионной условной гетероскедастичности (Autoregressive Conditional Heteroskedasticity, ARCH) для предсказания волатильности;
- ◆ как строить векторные авторегрессионные модели;
- ◆ как применять коинтеграцию для стратегии парной торговли.

Аналитические инструменты для диагностики и извлечения признаков

Данные временного ряда представляют собой последовательность значений, разделенных дискретными временными интервалами, которые, как правило, равноудалены (за исключением случаев, когда значения пропущены). Временной ряд часто моделируется как стохастический процесс, состоящий из коллекции случайных величин, $y(t_1), \dots, y(t_T)$, с одной переменной для каждой точки во времени t_i , $i = 1, \dots, T$. Одномерный временной ряд состоит из единственного значения y в каждой точке во времени, тогда как многомерный временной ряд состоит из нескольких наблюдений, которые могут быть представлены вектором.

Число периодов $\Delta t = t_i - t_j$ между разными точками во времени t_i , t_j называется *сдвигом во времени*, или *лагом*, с $T - 1$ сдвигами для каждого временного ряда. Подобно тому, как для срезовых (кросс-секционных) моделей ключевыми являются связи между разными переменными в определенной временной точке, для анализа и эксплуатации регулярностей во временных рядах основополагающими являются связи между точками данных, отделенными заданным временным сдвигом. В случае срезовых моделей мы различали входные и выходные переменные, или целевые и предсказательные переменные, соответственно метками y и x . В контексте временных рядов сдвинутые результирующие значения играют роль входных значений или значений x в срезовом (кросс-секционном) контексте.

Временной ряд называется *белым шумом*, если он представляет собой последовательность одинаково распределенных взаимно независимых (independent and

identically distributed, i.i.d.) случайных величин ε_t с выходным средним значением и дисперсией. В частности, ряд называется *гауссовым белым шумом*, если случайные величины нормально распределены с нулевым средним значением μ и постоянной дисперсией σ .

Временной ряд является линейным, если его можно записать как взвешенную сумму прошлых возмущений ε_t , которые также называются *инновациями*, и здесь принимаются, как представляющие белый шум, плюс среднее значение ряда μ :

$$y_t = \mu + \sum_{i=0}^{\infty} a_i \varepsilon_{t-i}, \quad a_0 = 1, \varepsilon \sim \text{i.i.d.}$$

Ключевой целью анализа временных рядов является понимание динамического поведения, приводимого в действие коэффициентами a_i . Анализ временных рядов предлагает методы, приспособленные к этому типу данных с целью извлечения полезных регулярностей, которые, в свою очередь, помогают строить предсказательные модели. Мы представим наиболее важные инструменты, служащие для этой цели, включая разложение на ключевые систематические элементы, анализ автокорреляции и статистику скользящего окна, такую как скользящие средние. Линейные модели временных рядов часто принимают определенные допущения о данных, такие как стационарность, и мы также представим концепцию, диагностические инструменты и типичные преобразования для достижения стационарности.

Применительно к большинству примеров этой главы мы работаем с данными, предоставляемыми Федеральной резервной системой, к которым можно получить доступ с помощью модуля библиотеки `pandas datareader`, представленного в *главе 2*. Примеры исходного кода для этого раздела имеются в блокноте `tsa_and_arma.ipynb`.

Как разложить временной ряд на регулярности

Данные временного ряда обычно содержат смесь различных регулярностей, которые можно разложить на несколько компонент, каждая из которых представляет базовую категорию регулярности. В частности, временной ряд часто состоит из систематических компонент: тренда, сезонности и циклов, а также бессистемного шума. Эти компоненты могут быть объединены в аддитивную, линейную модель, в частности, когда флуктуации не зависят от уровня ряда, или в нелинейную, мультипликативную модель.

Эти компоненты могут быть выделены автоматически. Библиотека `StatsModels` включает простой метод разбиения временного ряда на трендовую, сезонную и остаточную компоненты с использованием скользящих средних. Мы можем применить его к ежемесячным данным промышленного производства как с сильным трендом, так и с сезонностью следующим образом:

```
import pandas_datareader.data as web
import statsmodels.tsa.api as tsa
```

```

industrial_production = web.DataReader('IPGMFN', 'fred', '1988',
                                       '2017-12').squeeze().dropna()
nasdaq = web.DataReader('NASDAQCOM', 'fred', '1990',
                        '2017-12-31').squeeze().dropna()

components = tsa.seasonal_decompose(industrial_production, model='additive')
ts = (industrial_production.to_frame('Original')
      .assign(Trend=components.trend)
      .assign(Seasonality=components.seasonal)
      .assign(Residual=components.resid))
ts.plot(subplots=True, figsize=(14, 8));

```

Результирующие графики на рис. 8.1 показывают аддитивные компоненты. Остаточная компонента будет находиться в центре внимания дополнительного моделирования, принимая за основу, что компоненты тренда и сезонности более детерминированы и поддаются простой экстраполяции.

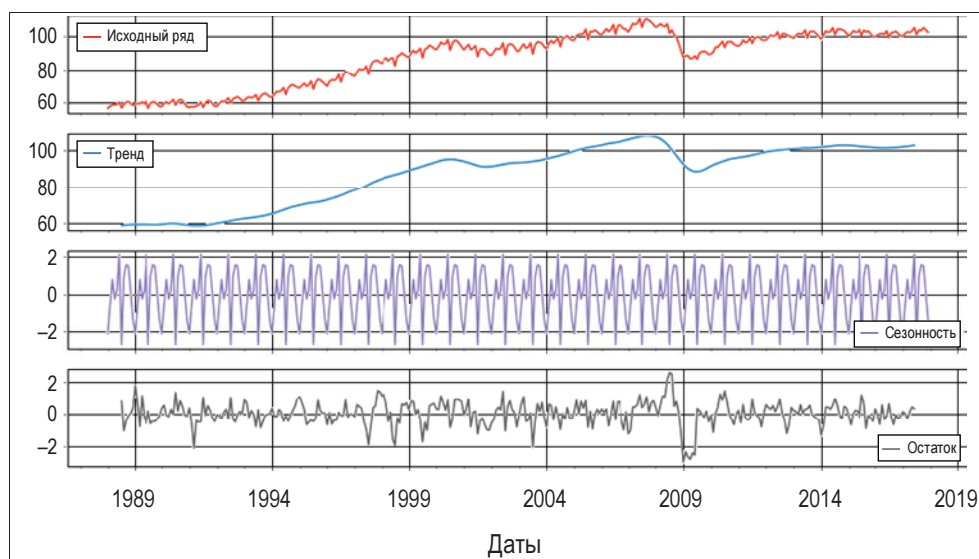


Рис. 8.1. Аддитивные компоненты

Существуют более изощренные, модельно-ориентированные подходы, которые включены в справочные материалы, размещенные в репозитории GitHub.

Как вычислять статистику скользящего окна

Учитывая упорядоченность данных временного ряда, где элементы следуют один за другим, вполне естественным является вычисление хорошо знакомой описательной статистики для периодов заданной длины с целью обнаружения стабильности или изменений в поведении и получения сглаженного представления, которое улавливает систематические аспекты, при этом отфильтровывая шум.

Статистика скользящего окна служит этому процессу: она создает новый временной ряд, в котором каждая точка данных представляет собой сводную статистику, рассчитанную за определенный период исходных данных. Скользящие средние являются самым известным примером. Исходные точки данных могут входить в вычисления с равными весами или использовать веса, например, для выделения более свежих точек данных. Экспоненциальные скользящие средние рекурсивно вычисляют веса, которые сжимаются, или затухают, для точек данных, расположенных дальше от настоящего времени. Новые точки данных обычно представляют собой сводку всех предыдущих точек данных, но они также могут быть вычислены из окружающего окна.

Библиотека `pandas` содержит очень гибкий функционал для определения различных типов окон, включая скользящие, экспоненциально взвешенные и расширяющиеся окна. На втором шаге к каждому элементу данных, захватываемых окном, можно применить вычисления. Эти вычисления включают встроенные стандартные вычисления для индивидуального ряда, такие как среднее или сумма, корреляция или ковариация для нескольких рядов, а также пользовательские функции. Примеры скользящего среднего и экспоненциального сглаживания в следующем разделе используют эти инструменты.

Скользящие средние и экспоненциальное сглаживание

Ранние прогнозные модели включали модели скользящего среднего с экспоненциальными весами, именуемые моделями экспоненциального сглаживания. Мы столкнемся со скользящими средними снова как с ключевыми строительными блоками для линейных временных рядов.

Прогнозы, которые опираются на методы экспоненциального сглаживания, используют средневзвешенные значения прошлых наблюдений, где веса экспоненциально убывают по мере старения наблюдений. Следовательно, более недавнее наблюдение получает более высокий ассоциированный вес. Эти методы популярны для временных рядов, которые не имеют очень сложных и резких регулярностей.

Экспоненциальное сглаживание представляет собой популярное техническое решение, основанное на взвешенных средних значениях прошлых наблюдений с экспоненциальным уменьшением веса по мере старения наблюдений. Другими словами, чем наблюдение свежее, тем выше ассоциированный вес. Указанный каркас генерирует надежные прогнозы быстро и в широком диапазоне временных рядов, что является большим преимуществом и имеет первостепенное значение для применения в индустрии.

Как измерять автокорреляцию

Автокорреляция (также именуемая внутрирядовой корреляцией) адаптирует понятие корреляции к контексту временных рядов: подобно тому, как корреляционный коэффициент служит мерой силы линейной связи между двумя переменными,

автокорреляционный коэффициент ρ_k служит мерой степени линейной связи между значениями временного ряда, отделенными заданным сдвигом k :

$$\rho_k = \frac{\sum_{t=k+1}^T (y_t - \bar{y})(y_{t-k} - \bar{y})}{\sum_{t=1}^T (y_t - \bar{y})^2}.$$

Следовательно, мы можем вычислить один автокорреляционный коэффициент для каждого из $T-1$ временных сдвигов во временном ряду, где T — это длина ряда. *Функция автокорреляции* (autocorrelation function, ACF) вычисляет корреляционные коэффициенты как функцию сдвига.

Автокорреляция для сдвига больше 1 (т. е. между наблюдениями на расстоянии более одного временного шага) отражает прямую корреляцию между этими наблюдениями и косвенное влияние промежуточных точек данных. Частичная автокорреляция устраняет это влияние и служит мерой только линейной зависимости между точками данных на заданном расстоянии временного сдвига. *Функция частичной автокорреляции* (partial ACF) обеспечивает все корреляции, которые возникают после устранения эффектов корреляции в более коротких сдвигах.

Существуют алгоритмы, которые вычисляют частичную автокорреляцию из выборочной автокорреляции на основе точной теоретической связи между функцией частичной автокорреляции (PACF) и функцией автокорреляции (ACF).

Коррелограмма — это просто график функции автокорреляции или функции частичной автокорреляции для последовательных временных сдвигов, $k = 0, 1, \dots, n$. Она позволяет одним взглядом осматривать корреляционную структуру по всем сдвигам. Главное применение коррелограмм заключается в обнаружении автокорреляции после удаления эффектов детерминированного тренда или сезонности. Функция автокорреляции и функция частичной автокорреляции являются ключевыми диагностическими инструментами для конструирования линейных моделей временных рядов, и мы рассмотрим примеры графиков функции автокорреляции и частичной автокорреляции в следующем разделе, посвященном преобразованиям временных рядов.

Как диагностировать и достигать стационарности

Статистические свойства, такие как среднее значение, дисперсия или автокорреляция стационарного временного ряда, не зависят от периода, т. е. они не меняются во временной динамике. Следовательно, стационарность приводит к тому, что временной ряд не имеет трендовых или сезонных эффектов и что описательные статистические показатели, такие как среднее значение или стандартное отклонение, при вычислении для разных скользящих окон, остаются постоянными или не сильно меняются с течением времени. Временной ряд разворачивается к своему среднему значению, и отклонения имеют постоянную амплитуду, тогда как краткосрочные движения в статистическом смысле всегда выглядят одинаково.

Более формально, строгая стационарность требует совместного распределения любого подмножества наблюдений временного ряда независимо от времени по всем статистическим моментам. Таким образом, помимо среднего и дисперсии, более высокие моменты, такие как асимметрия и эксцесс, также должны быть постоянными, независимо от сдвига между различными наблюдениями. В большинстве приложений мы ограничиваем стационарность первым и вторым моментами, вследствие чего временной ряд является ковариационно стационарным с постоянным средним, дисперсией и автокорреляцией.

Обратите внимание, что мы специально допускаем зависимость между наблюдениями в разных сдвигах, подобно тому, как мы хотим, чтобы входные данные для линейной регрессии коррелировались с результатом. Стационарность влечет за собой стабильность этих связей, что обеспечивает предсказание, поскольку модель может сосредотачиваться на усвоении систематических регулярностей, которые имеют место внутри стабильных статистических свойств. Это важно, поскольку классические статистические модели исходят из того, что входные данные временных рядов являются стационарными.

В следующих далее разделах представлена диагностика, помогающая обнаруживать, когда данные не являются стационарными, и преобразования, помогающие удовлетворять этим допущениям.

Преобразования временного ряда

Для того чтобы удовлетворить допущение о стационарности линейных моделей временных рядов, нам нужно преобразовать исходный временной ряд, часто в несколько шагов. Общепринятые преобразования включают применение (натурального) логарифма для конвертирования регулярности экспоненциального роста в линейный тренд и стабилизации дисперсии. Дефляция сводится к делению временного ряда на другой ряд, который вызывает трендовое поведение, например, деление номинального ряда на индекс цен с целью его конвертирования в реальную меру.

Ряд является трендово-стационарным, если он разворачивается к стабильному долгосрочному линейному тренду. Его часто можно сделать стационарным, вписав трендовую линию с использованием линейной регрессии и остатков либо включив временной индекс в регрессию или модель AR(I)MA в качестве независимой переменной (см. следующий далее раздел о моделях одномерных временных рядов), возможно, в сочетании с логарифмированием или дефляцией.

Во многих случаях для того чтобы сделать ряд стационарным, детрендрования не достаточно. Вместо этого нам нужно преобразовать исходные данные во временной ряд попериодных и/или посезонных разниц. Другими словами, мы используем результат вычитания соседних точек данных или значений в сезонных сдвигах друг из друга. Обратите внимание, что когда такое исчисление последовательных разностей применяется к логарифмически преобразованному ряду, результаты представляют моментные темпы роста или возвраты в финансовом контексте.

Если одномерный ряд становится стационарным после исчисления последовательных разностей d раз, то принято говорить, что он интегрирован с порядком d , либо если $d = 1$, то просто интегрирован. Такое поведение обусловлено так называемыми единичными корнями.

Как диагностировать и обращаться с единичными корнями

Единичные корни представляют особую проблему для определения преобразования, которое делает временной ряд стационарным. Временные ряды часто моделируются как стохастические процессы приведенной ниже авторегрессионной формы, которую мы рассмотрим подробнее в качестве строительного блока для моделей авторегрессионного интегрированного скользящего среднего (autoregressive integrated moving average, ARIMA):

$$y_t = a_1 y_{t-1} + a_2 y_{t-2} + \dots + a_p y_{t-p} + \varepsilon_t,$$

где текущее значение есть взвешенная сумма прошлых значений плюс случайное возмущение. Такой процесс имеет характеристическое уравнение следующего вида:

$$m^p + m^{p-1}a_1 - m^{p-2}a_2 - \dots - a_p = 0.$$

Если один из корней этого уравнения равен 1, то принято говорить, что процесс имеет единичный корень. Он будет нестационарным, но не обязательно должен иметь тренд. Если остальные корни характеристического уравнения в абсолютном выражении меньше 1, то первая разность процесса будет стационарной, и процесс интегрируется (с порядком 1) или $I(1)$. В случае, когда дополнительные корни в абсолютном выражении больше 1, порядок интегрирования будет выше, и потребуются дополнительные исчисления последовательных разностей.

На практике временные ряды процентных ставок или цен активов часто не являются стационарными, например, потому, что не существует ценового уровня, к которому эти ряды разворачиваются. Наиболее ярким примером нестационарного ряда является случайное блуждание для ценового временного ряда p_t , при заданной начальной цене p_0 (например, цены IPO акции) и белом шумном возмущении ε , которое удовлетворяет следующему:

$$p_t = p_{t-1} + \varepsilon_t = \sum_{s=0}^t \varepsilon_s + p_0.$$

Многократная подстановка показывает, что текущее значение p_t является суммой всех предшествующих возмущений или инноваций ε и начальной цены p_0 . Если уравнение включает постоянный член, то принято говорить, что случайное блуждание имеет дрейф. Таким образом, случайное блуждание является авторегрессионным стохастическим процессом следующей формы:

$$y_t = a_1 y_{t-1} + \varepsilon_t, \quad a_1 = 1$$

с характеристическим уравнением $m - a_1 = 0$, которое имеет единичный корень, является одновременно нестационарным и имеет порядок интегрирования, рав-

ный 1. С одной стороны, учитывая одинаково распределенную и взаимно независимую природу ε , дисперсия временного ряда равна σ^2 , которая не является второпорядково стационарной и подразумевает, что в принципе ряд мог бы со временем принимать любое значение. С другой стороны, беря первую разницу $\Delta p_t = p_t - p_{t-1}$, мы в остатке получаем $\Delta p = \varepsilon_t$, т. е. стационарность при заданном статистическом допущении о ε .

Определяющей характеристикой нестационарного ряда с единичным корнем является долгая память: поскольку текущие значения являются суммой прошлых возмущений, большие инновации сохраняются гораздо дольше, чем для средне-разворотного стационарного ряда.

В дополнение к использованию разности между соседними точками данных для удаления постоянной регулярности изменения, ее можно использовать, применяя сезонное исчисление последовательных разностей для удаления регулярностей сезонного изменения. Это предусматривает разность значений на расстоянии сдвига, которое представляет собой длину сезонной регулярности, т. е. на расстоянии четырех кварталов или 12 месяцев, с целью удалить как сезонность, так и линейный тренд.

Четко выявить правильное преобразование и, в частности, надлежащее число и сдвиги для исчисления последовательных разностей не всегда получается, в связи с чем было предложено несколько правил, которые резюмированы в следующем ниже списке.

- ◆ *Положительные автокорреляции до 10+ сдвигов*: возможно потребуется более высокопорядковое исчисление последовательных разностей.
- ◆ *Односдвиговая автокорреляция¹ близка к нулю или отрицательна, или вообще мала и безрегулярна*: в исчислении разности более высокого порядка нет необходимости.
- ◆ *Односдвиговая автокорреляция менее -0,5*: ряд может иметь чрезмерно исчисленные последовательные разности.
- ◆ Слегка избыточное или недостаточное исчисление последовательных разностей может быть исправлено с помощью членов AR или MA.
- ◆ Оптимальное исчисление последовательных разностей часто приводит к наименьшему стандартному отклонению, но не всегда.
- ◆ Модель без исчисления последовательных разностей исходит из того, что исходный ряд является стационарным, в том числе среднеразворотным. Он обычно содержит постоянный член для обеспечения ненулевого среднего значения.
- ◆ Модель с одним порядком исчисления последовательных разностей исходит из того, что первоначальный ряд имеет постоянный средний тренд и должен содержать постоянный член.

¹ То есть корреляция между элементами $x[t]$ и $x[t - 1]$ временного ряда. — Прим. перев.

- ♦ Модель с двумя порядками исчисления последовательных разностей исходит из того, что первоначальный ряд имеет варьирующий во времени тренд и не должен содержать константу.

Некоторые авторы рекомендуют дробное исчисление последовательных разностей как более гибкий подход к превращению интегрированного ряда в стационарный и вполне способный оставлять больше информации или сигнала, чем простые или сезонные разности в дискретных интервалах (см. справочные материалы в репозитории GitHub).

Единично-корневые тесты

Статистические единично-корневые тесты являются распространенным способом объективного определения необходимости (дополнительного) исчисления последовательных разностей. Они представляют собой проверки статистической гипотезы о стационарности, которые служат для определения необходимости исчисления последовательных разностей.

Расширенная проверка Дики — Фуллера (augmented Dickey — Fuller, ADF) тестирует нулевую гипотезу о том, что образец временного ряда имеет единичный корень относительно альтернативной гипотезы о стационарности. Она регрессирует временной ряд исчисленных последовательных разностей на временном тренде, первый сдвиг и все сдвинутые разности и вычисляет проверочный статистический показатель из значения коэффициента на сдвинутом временном ряду. Библиотека StatsModels обеспечивает его реализацию (см. прилагаемый блокнот).

Формально расширенная проверка Дики — Фуллера для временного ряда y_t выполняет линейную регрессию:

$$\Delta y_t = \alpha + \beta t + \gamma y_{t-1} + \delta_1 \Delta y_{t-1} + \dots + \delta_{p-1} \Delta y_{t-p+1} + \varepsilon_t,$$

где α — это константа; β — коэффициент на временном тренде; p — число сдвигов, используемых в модели. Ограничение $\alpha = \beta = 0$ влечет за собой случайное блуждание, тогда как только $\beta = 0$ влечет за собой случайное блуждание с дрейфом. Порядок сдвигов обычно определяется с использованием информационных критериев AIC и BIC, введенных в главе 7.

Статистический показатель расширенной проверки Дики — Фуллера использует выборочный коэффициент γ , который в условиях нулевой гипотезы о единично-корневой нестационарности равен нулю и отрицателен в противном случае. Он призван демонстрировать, что для интегрированного ряда значение сдвинутого ряда не должно обеспечивать полезной информации в предсказании первой разности за пределами сдвинутых разностей.

Как применять преобразования временных рядов

На рис. 8.2 показаны временные ряды фондового индекса NASDAQ и промышленного производства за 30 лет вплоть до 2017 г. в исходном виде, а также преобразованные версии соответственно после применения логарифма и последующего применения первой и сезонной разностей (с временным сдвигом 12). На диаграммах также отображается p -значение расширенной проверки Дики — Фуллера, которое в обоих случаях позволяет отклонить гипотезу о единично-корневой нестационарности после всех преобразований.

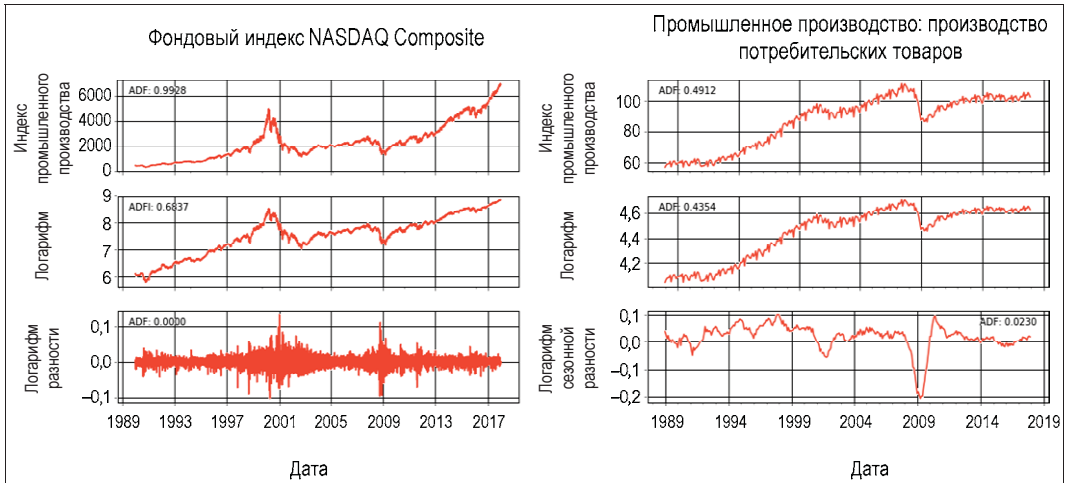


Рис. 8.2. Временные ряды фондового индекса NASDAQ Composite и промышленного производства за 30 лет вплоть до 2017–2019 гг.

Мы можем продолжить анализировать релевантные временные характеристики преобразованных временных рядов, используя квантиль-квантильный график (Q-Q), который сравнивает квантили распределения наблюдений временного ряда с квантилями нормального распределения, и коррелограммы на основе функций автокорреляции и частичной автокорреляции.

Глядя на график NASDAQ, мы замечаем, что хотя тренда нет, дисперсия не является постоянной, а скорее, показывает кластеризованные всплески вокруг периодов рыночных потрясений в конце 1980-х, 2001 и 2008 гг. Квантиль-квантильный график отмечает толстые хвосты распределения с более частыми экстремальными значениями, чем предусматривает нормальное распределение. Графики функций автокорреляции и частичной автокорреляции показывают похожие регулярности, причем автокорреляция в нескольких сдвигах выглядит значимой (рис. 8.3).

При рассмотрении месячного временного ряда промышленного производства потребительских товаров мы отмечаем большой отрицательный выброс вслед за кризисом 2008 г., а также соответствующую асимметрию (перекос) в квантильно-квантильном графике (рис. 8.4). Автокорреляция намного выше, чем для

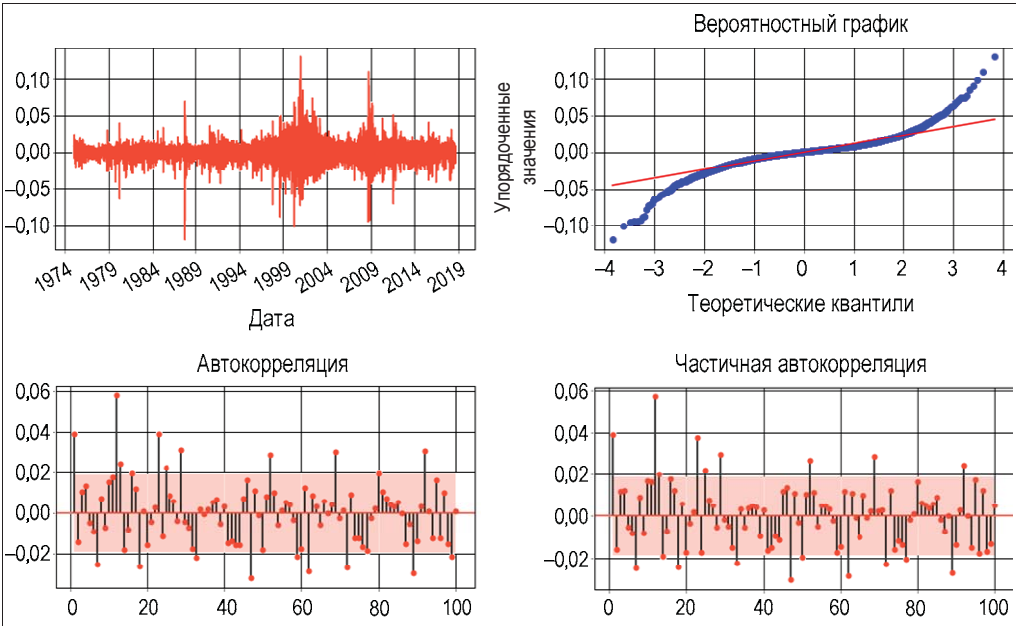


Рис. 8.3. Фондовый индекс NASDAQ Composite (логарифмированный и разностный)

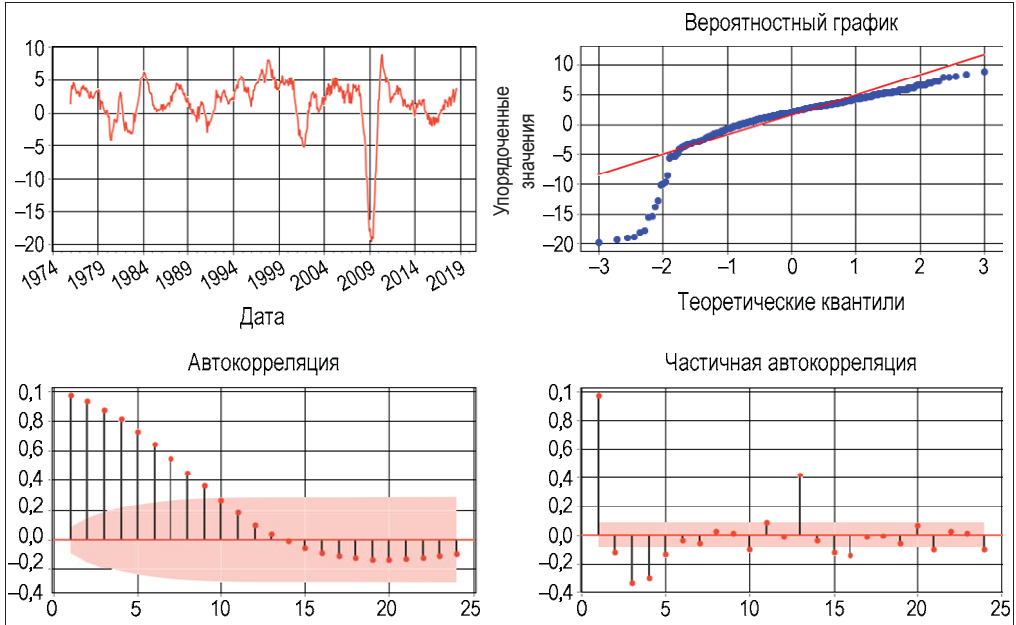


Рис. 8.4. Промышленное производство потребительских товаров (сезонные разности)

возвратности NASDAQ, и плавно снижается. График функции частичной автокорреляции показывает отчетливые положительные автокорреляционные регулярности в сдвигах 1 и 13 и значимые отрицательные коэффициенты в сдвигах 3 и 4.

Модели одномерных временных рядов

Множественные линейные регрессионные модели выражали интересующую переменную как линейную комбинацию предсказателей, или входных переменных. Модели одномерных временных рядов связывают значение временного ряда в интересующей временной точке с линейной комбинацией сдвинутых значений ряда и, возможно, членов прошлого возмущения.

В то время как модели экспоненциального сглаживания основаны на описании тренда и сезонности в данных, модели авторегрессионного интегрированного скользящего среднего ARIMA (autoregressive integrated moving average) призваны описывать автокорреляцию в данных. Модели ARIMA(p, d, q) требуют стационарности и эффективно задействуют два строительных блока:

- ◆ *авторегрессионные* (Autoregressive, AR) члены, состоящие из p -сдвинутых значений временного ряда;
- ◆ *среднескользящие* (Moving average, MA) члены, содержащие q -сдвинутых возмущений.

Буква I в аббревиатуре ARIMA означает интегрированный, потому что модель может учитывать единично-корневую нестационарность путем исчисления разности ряда d раз. Термин "авторегрессия" подчеркивает, что из модели ARIMA следует регрессия временного ряда на собственных значениях.

Мы представим строительные блоки модели ARIMA, простую авторегрессионную модель (AR) и модель скользящего среднего (MA) и объясним, как совмещать их в моделях авторегрессионного скользящего среднего (ARMA), которые могут учитывать интегрирование ряда как модели ARIMA или включать экзогенные переменные как модели AR(I)MAX. Кроме того, мы проиллюстрируем, как включать сезонные члены AR и MA с целью расширения набора инструментов за счет включения моделей SARMAX.

Как строить авторегрессионные модели

Авторегрессионная модель AR порядка p призвана улавливать линейную зависимость между значениями временного ряда в разных сдвигах и может быть записана следующим образом:

$$AR(p): y_t = \phi_0 + \phi_1 y_{t-1} + \phi_2 y_{t-2} - \dots - \phi_p y_{t-p} + \varepsilon_t, \quad \varepsilon \sim \text{i.i.d.}$$

Она очень похожа на множественную линейную регрессию на сдвинутых значениях y_t . Эта модель имеет следующее характеристическое уравнение:

$$1 - \phi_1 x - \phi_2 x^2 - \dots - \phi_p x^p = 0.$$

Инверсии решения этого уравнения в x являются характеристическими корнями, и процесс $AR(p)$ является стационарным, если все эти корни меньше 1 в абсолютных единицах, и нестабильным в противном случае. Для стационарного ряда многоступенчатые прогнозы будут сходиться к среднему значению ряда.

Мы можем вычислить модельные параметры уже хорошо знакомым методом наименьших квадратов, используя $p + 1$, ..., T наблюдений с целью обеспечения наличия данных для каждого сдвинутого члена и результата.

Как выявлять число сдвигов

На практике задача состоит в принятии решения о надлежащем порядке p сдвинутых членов. Инструменты анализа временного ряда для внутрирядовой корреляции играют ключевую роль. Функция автокорреляции оценивает автокорреляцию между наблюдениями в разных сдвигах, которая, в свою очередь, является результатом как прямой, так и косвенной линейной зависимости.

Следовательно, для авторегрессионной модели AR порядка k функция автокорреляции будет показывать значительную внутрирядовую корреляцию вплоть до временного сдвига k , и вследствие инерции, вызываемой косвенными эффектами линейной связи, будет распространяться на последующие сдвиги и в конечном итоге затухнет по мере ослабления эффекта. С другой стороны, функция частичной автокорреляции служит мерой только прямой, линейной связи между наблюдениями на расстоянии заданного сдвига, вследствие чего она не будет отражать корреляцию для сдвигов за пределами k .

Как диагностировать подгонку модели

Если модель улавливает линейную зависимость между сдвигами, то остатки должны быть похожими на белый шум.

В дополнение к обследованию функции автокорреляции с целью подтверждения отсутствия значимых автокорреляционных коэффициентов проверочный статистический показатель Q Льюнг — Бокса позволяет проверять гипотезу о том, что остаточный ряд подчиняется белому шуму. Нулевая гипотеза состоит в том, что все автокорреляционные коэффициенты равны нулю, относительно альтернативной гипотезы о том, что некоторые коэффициенты не равны. Указанный проверочный статистический показатель вычисляется из автокорреляционных коэффициентов p_k , выборки для разных сдвигов k , и подчиняется распределению X^2 :

$$Q(m) = T(T+2) \sum_{l=2}^m \frac{\rho_l^2}{T-l}.$$

Как мы увидим, библиотека StatsModels предоставляет информацию о значимости коэффициентов для разных сдвигов, и незначимые коэффициенты должны быть удалены. Если статистический показатель Q отклоняет нулевую гипотезу об отсутствии автокорреляции, то следует рассмотреть дополнительные AR -члены.

Как строить модели скользящего среднего

Как показано ниже, модель скользящего среднего (MA) порядка q вместо сдвинутых значений временного ряда использует q прошлых возмущений в модели, которая похожа на регрессию:

$$MA(p): y_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}, \quad \varepsilon \sim \text{i.i.d.}$$

Поскольку мы не наблюдаем значения белом шумных возмущений ε_t , $MA(q)$ не является регрессионной моделью, подобной тем, что мы видели до сих пор. Вместо того чтобы использовать наименьшие квадраты, модели $MA(q)$ вычисляются с использованием *максимального правдоподобия* (MLE), инициализируя или же оценивая возмущения в начале ряда, а затем рекурсивно и итеративно вычисляя остаток.

Модель $MA(q)$ берет свое название от представления каждого значения y_t как взвешенного скользящего среднего прошлых q инноваций. Другими словами, текущие оценки представляют собой коррекцию относительно прошлых ошибок, сделанных моделью. Использование скользящих средних в моделях $MA(q)$ отличается от использования экспоненциального сглаживания или вычисления сезонных компонентов временного ряда, поскольку модель $MA(q)$ в отличие от шумоподавления или вычисления трендового цикла прошлых значений призвана прогнозировать будущие значения.

Процессы $MA(q)$ всегда стационарны, потому что они являются взвешенной суммой белом шумных переменных, которые сами являются стационарными.

Как выявлять число сдвигов

Временной ряд, генерируемый процессом $MA(q)$, обусловлен остатками от q предсказаний предшествующей модели. Следовательно, функция автокорреляции для процесса $MA(q)$ покажет значительные коэффициенты для значений вплоть до сдвига q , а затем резко снизится, потому что считается, что именно так значения ряда были сгенерированы.

Связь между моделями AR и MA

Модель $AR(p)$ может быть выражена как процесс $MA(\infty)$ с использованием многократной подстановки. При наложении ограничений на размер его коэффициентов процесс $MA(q)$ становится обратимым и может быть выражен как процесс $AR(\infty)$.

Как строить модели ARIMA и их расширения

Модели авторегрессионного интегрированного скользящего среднего $ARIMA(p, d, q)$ объединяют процессы $AR(p)$ и $MA(q)$, задействуя комплементарности этих строительных блоков и упрощая разработку модели за счет более компактной формы, снижения числа параметров и, в свою очередь, снижения риска переподгонки.

Эти модели также берут на себя устранение единично-корневой нестационарности, используя d -ю разность значений временного ряда. Модель $ARIMA(p, 1, q)$ аналогична модели $ARMA(p, q)$ с первыми разностями ряда. Используя y' для обозначения исходного ряда после несезонного исчисления последовательных разностей d раз, модель $ARIMA(p, d, q)$ принимает простой вид:

$$\begin{aligned} ARIMA(p, d, q): y_t &= AR(p) + MA(q) = \\ &= \phi_0 + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q}, \quad \varepsilon \sim \text{i.i.d.} \end{aligned}$$

Модели $ARIMA$ тоже оцениваются с использованием максимального правдоподобия. В зависимости от реализации верхнепорядковые модели обычно могут включать нижнепорядковые модели. Например, библиотека `StatsModels` включает все нижнепорядковые члены p и q и не позволяет удалять коэффициенты для сдвигов ниже самого высокого значения. В этом случае верхнепорядковые модели всегда будут вписываться лучше. Следует остерегаться перепогонки модели к данным, используя слишком много членов.

Как выявлять число членов моделей AR и MA

Поскольку члены $AR(p)$ и $MA(q)$ между собой взаимодействуют, информация, предоставляемая функциями автокорреляции и частичной автокорреляции, больше не является надежной и может использоваться только в качестве отправной точки.

Традиционно использовались информационные критерии AIC и BIC, при этом обычно во время отбора модельной конструкции опирались на внутривыборочную подгонку. С другой стороны, мы можем опираться на вневыборочные тесты при перекрестном контроле многочисленных параметрических вариантов.

В следующей ниже сводке приводятся некоторые общие рекомендации по выбору порядка модели в случае рассмотрения моделей AR и MA по отдельности.

- ◆ Сдвиг, за пределами которого функция частичной автокорреляции отсекает, т. е. резко становится нулем, является индикативным числом членов AR. Если частично-автокорреляционная функция ряда исчисленных последовательных разностей резко отсекает и/или односдвиговая автокорреляция является положительной, то следует добавить один или несколько членов AR.
- ◆ Сдвиг, за пределами которого функция автокорреляции отсекает, является индикативным числом членов MA. Если автокорреляционная функция ряда исчисленных последовательных разностей показывает резкое отсечение и/или односдвиговая автокорреляция является отрицательной, то следует рассматривать возможность добавления в модель члена MA.
- ◆ Члены AR и MA могут нейтрализовывать эффекты друг друга, поэтому всегда следует стараться снижать число членов AR и MA на 1, если ваша модель содержит и то и другое, во избежание перепогонки, в особенности, если для схождения более сложной модели требуется более 10 итераций.

- ◆ Если коэффициенты AR в сумме составляют почти 1 и свидетельствуют о единичном корне в AR-части модели, то следует устранить 1 член AR и исчислить последовательные разности в модели (еще) один раз.
- ◆ Если коэффициенты MA в сумме составляют почти 1 и свидетельствуют о единичном корне в MA-части модели, то следует устранить 1 член MA и сократить порядок исчисления последовательных разностей на 1.
- ◆ Нестабильные долговременные прогнозы свидетельствуют о том, что в AR- или MA-частях модели может иметься единичный корень.

Добавление признаков — ARMAX

Модель ARMAX добавляет входные переменные или ковариацию с правой стороны ARMA-модели временных рядов (если исходить из того, что ряд является стационарным, то исчисление последовательных разностей можно пропустить):

$$\text{ARIMA}(p, d, q): y_t = \beta x_t + \text{AR}(p) + \text{MA}(q) = \\ = \beta x_t + \varphi_0 + \varphi_1 y_{t-1} + \dots + \varphi_p y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q}, \quad \varepsilon \sim \text{i.i.d.}$$

Это уравнение напоминает линейную регрессионную модель, но его довольно сложно интерпретировать, потому что влияние β на y_t не является эффектом увеличения в x_t на одну единицу, как в линейной регрессии. Вместо этого наличие сдвинутых значений y_t в правой части уравнения означает, что коэффициент может быть интерпретирован только с учетом сдвинутых значений переменной отклика, что вряд ли является интуитивно понятным.

Добавление сезонного исчисления последовательных разностей — SARIMAX

В случае временных рядов с сезонными эффектами мы можем включать члены AR и MA, которые улавливают периодичность сезонности. Например, при использовании месячных данных и когда продолжительность сезонного эффекта равна одному году, сезонные члены AR и MA будут отражать эту конкретную длину сдвига.

Тогда модель $\text{ARIMAX}(p, d, q)$ становится моделью $\text{SARIMAX}(p, d, q) \times (P, D, Q)_s$, которую немного сложнее расписать, но справочные материалы в репозитории GitHub, включая документацию библиотеки StatsModels, обеспечивают эту информацию во всех подробностях.

Для того чтобы проиллюстрировать реализацию, теперь мы построим сезонную модель ARMA, используя макроданные.

Как прогнозировать макрэкономические фундаментальные показатели

Мы построим модель SARIMAX для ежемесячных данных на временном ряде промышленного производства за период 1988–2017 гг. Как показано в первом разделе, посвященном аналитическим инструментам, данные были преобразованы логарифмически, и мы используем сезонные (сдвиг 12) разности. Мы оцениваем модель для диапазона как обыкновенных, так и конвенциональных параметров AR и MA с использованием скользящего 10-летнего окна тренировочных данных и вычисляем ошибку RMSE прогноза на 1 шаг вперед, как показано в следующем ниже упрощенном фрагменте кода (подробности см. в репозитории GitHub):

```
train_size = 120 # 10 лет тренировочных данных
test_results = {}
test_set = industrial_production_log_diff.iloc[train_size:]

for p1 in range(4): # порядок AR
    for q1 in range(4): # порядок MA
        for p2 in range(3): # сезонный порядок AR
            for q2 in range(3): # сезонный порядок MA
                preds = test_set.copy().to_frame('y_true')
                                .assign(y_pred=np.nan)

                aic, bic = [], []
                if p1 == 0 and q1 == 0:
                    continue
                print(p1, q1, p2, q2)
                convergence_error = stationarity_error = 0
                y_pred = []
                for i, T in enumerate(range(train_size,
                                            len(industrial_production_log_diff))):
                    train_set = industrial_production_log_diff.iloc[T-train_size:T]
                    try:
                        model = tsa.SARIMAX(endog=train_set, # спецификация модели
                                           order=(p1, 0, q1),
                                           seasonal_order=(p2, 0, q2, 12)).fit()
                    except LinAlgError:
                        convergence_error += 1
                    except ValueError:
                        stationarity_error += 1

                    preds.iloc[i, 1] = model.forecast(steps=1)[0] # прогноз на 1 шаг
                                                                вперед

                    aic.append(model.aic)
                    bic.append(model.bic)

                preds.dropna(inplace=True)
                mse = mean_squared_error(preds.y_true, preds.y_pred)
```

```
test_results[(p1, q1, p2, q2)] = [np.sqrt(mse),
                                   preds.y_true.sub(preds.y_pred).std(),
                                   np.mean(aic),
                                   np.std(aic),
                                   np.mean(bic),
                                   np.std(bic),
                                   convergence_error,
                                   stationarity_error]
```

Мы также собираем критерии AIC и BIC, которые показывают очень высокий коэффициент ранговой корреляции 0,94, при этом BIC отдает предпочтение моделям с немного меньшим числом параметров, чем AIC. Лучшие пять моделей по величине ошибки RMSE показаны ниже:

| | | | | RMSE | AIC | BIC |
|----|----|----|----|----------|-------------|-------------|
| p1 | q1 | p2 | q2 | | | |
| 2 | 3 | 1 | 0 | 0.009323 | -772.247023 | -752.734581 |
| 3 | 2 | 1 | 0 | 0.009467 | -768.844028 | -749.331586 |
| 2 | 2 | 1 | 0 | 0.009540 | -770.904835 | -754.179884 |
| | 3 | 0 | 0 | 0.009773 | -760.248885 | -743.523935 |
| | 2 | 0 | 0 | 0.009986 | -758.775827 | -744.838368 |

Мы повторно оцениваем модель SARIMA(2, 0, 3)×(1, 0, 0) следующим образом:

```
best_model = tsa.SARIMAX(endog=industrial_production_log_diff, order=(2, 0, 3),
                          seasonal_order=(1, 0, 0, 12)).fit()
print(best_model.summary())
```

В результате получаем следующую статистическую сводку:

Statespace Model Results

```
=====
Dep. Variable:                IPGMFN    No. Observations:   348
Model:                SARIMAX(2, 0, 3)x(1, 0, 0, 12)    Log Likelihood 1139.719
Date:                Wed, 31 Oct 2018    AIC                -2265.438
Time:                22:58:00            BIC                -2238.472
Sample:                01-01-1989        HQIC               -2254.702
                        - 12-01-2017
```

Covariance Type: opg

```
=====
              coef    std err          z      P>|z|      [0.025      0.975]
-----
ar.L1          1.4934     0.104     14.351     0.000         1.289         1.697
ar.L2         -0.5159     0.102     -5.083     0.000        -0.715        -0.317
ma.L1         -0.5499     0.114     -4.813     0.000        -0.774        -0.326
ma.L2          0.2872     0.062      4.662     0.000         0.166         0.408
ma.L3          0.1815     0.070      2.589     0.010         0.044         0.319
ar.S.L12       -0.4486     0.047     -9.533     0.000        -0.541        -0.356
sigma2         8.141e-05  5.65e-06    14.399     0.000       7.03e-05       9.25e-05
=====
```

| | | | |
|--------------------------|-------|--------------------|-------|
| Ljung-Box (Q) : | 61.58 | Jarque-Bera (JB) : | 9.97 |
| Prob(Q) : | 0.02 | Prob(JB) : | 0.01 |
| Heteroskedasticity (H) : | 1.07 | Skew: | -0.20 |
| Prob(H) (two-sided): | 0.71 | Kurtosis: | 3.73 |

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

Коэффициенты являются значимыми, и статистический показатель Q отклоняет гипотезу о дальнейшей автокорреляции. Коррелограмма также указывает на то, что мы успешно устранили автокорреляцию ряда (рис. 8.5).

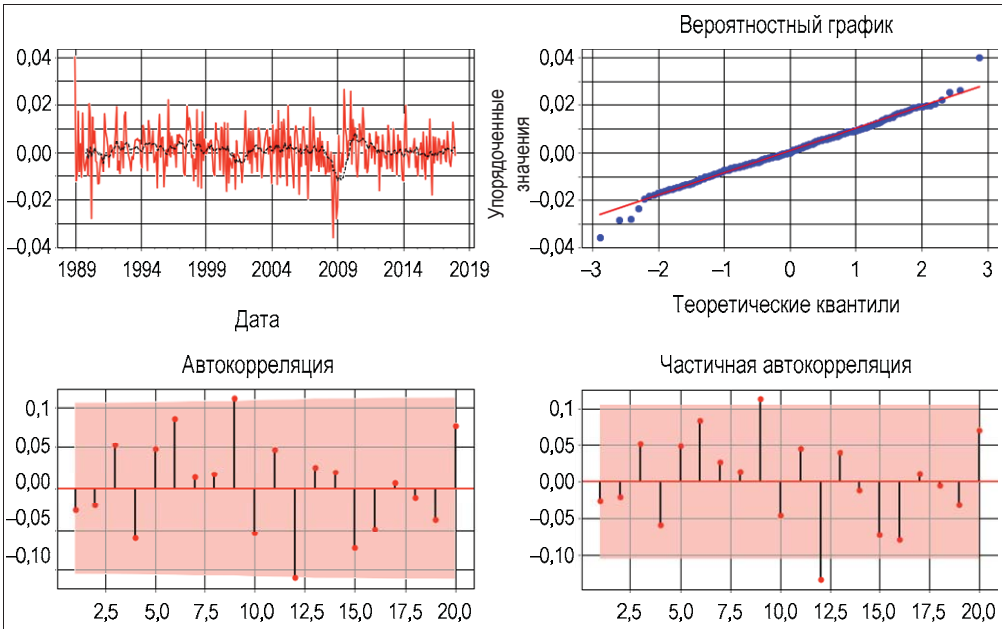


Рис. 8.5. Остатки модели SARIMAX

Как использовать модели временных рядов для прогнозирования волатильности

Особенно важной областью применения моделей одномерных временных рядов является предсказание волатильности. Волатильность финансового временного ряда, как правило, не является постоянной во временной динамике и изменяется, причем вспышки волатильности группируются вместе. Изменения в дисперсии создают сложности для прогнозирования временного ряда с использованием классических моделей ARIMA. Обращаясь к решению этой сложности, мы теперь смо-

делируем волатильность, вследствие чего сможем предсказывать изменения в дисперсии.

Гетероскедастичность — это технический термин для изменений в дисперсии переменной. Модель *авторегрессионной условной гетероскедастичности* (autoregressive conditional heteroskedasticity, ARCH) выражает дисперсию члена ошибки как функцию ошибок в предыдущих периодах. В частности, она исходит из того, что дисперсия ошибок подчиняется модели $AR(p)$.

Модель *обобщенной авторегрессионной условной гетероскедастичности* (generalized autoregressive conditional heteroskedasticity, GARCH) расширяет область применения до моделей ARMA. Процедура прогнозирования временного ряда часто сочетает модели ARIMA для математического ожидания среднего значения и модели ARCH/GARCH для математического ожидания дисперсии временного ряда. За разработку этого класса моделей Роберт Энгл (Robert Engle) и Клайв Грейнджер (Clive Granger) получили Нобелевскую премию по экономике за 2003 г. Первый из упомянутых также управляет лабораторией волатильности в школе Стерна Нью-Йоркского университета (см. справочный материал на GitHub), предоставляя многочисленные онлайн-примеры и инструменты, касающиеся моделей, которые мы будем обсуждать, и их многочисленные расширения.

Модель авторегрессионной условной гетероскедастичности (ARCH)

Модель $ARCH(p)$ — это попросту модель $AR(p)$ применительно к дисперсии остатков модели временного ряда, которая создает эту дисперсию в момент времени t в зависимости от сдвинутых наблюдений дисперсии. Более конкретно, члены ошибки ε_t являются остатками линейной модели, такой как ARIMA, на исходном временном ряду и разбиваются на зависящее от времени стандартное отклонение σ_t и возмущение z_t следующим образом:

$$ARCH(p): \text{var}(x_t) = \sigma_t^2 = \omega + \alpha_1 \varepsilon_{t-1}^2 + \dots + \alpha_p \varepsilon_{t-p}^2, \quad \varepsilon_t = \sigma_t z_t, \quad z_t \sim \text{i.i.d.}$$

Модель $ARCH(p)$ можно вычислить с помощью обычных наименьших квадратов (OLS). Энгл предложил метод выявления надлежащего порядка ARCH с использованием проверки лагранжева множителя, который соответствует F -тесту гипотезы о том, что все коэффициенты в линейной регрессии равны нулю (см. главу 7).

Одной из сильных сторон указанной модели является то, что она производит волатильность, оценивает положительный избыточный эксцесс, т. е. толстые хвосты относительно нормального распределения, что в свою очередь согласуется с эмпирическими наблюдениями о финансовых возвратах. Слабые стороны данной модели заключаются в том, что она исходит из одинакового эффекта для положительных и отрицательных шоков волатильности, поскольку она зависит от квадрата предыдущих шоков, в то время как цены активов, как известно, реагируют на положительные и отрицательные шоки по-разному. Модель ARCH также не предлагает нового понимания источника вариаций финансового временного ряда, по-

сколько она просто механически описывает условную дисперсию. Наконец, модели ARCH, скорее всего, предсказывают волатильность излишне, потому что они медленно откликаются на крупные изолированные шоки на временной ряд возвратов.

В случае правильно заданной модели ARCH стандартизированные остатки (деленные на модельную оценку для периода стандартного отклонения) должны походить на белый шум и могут быть подвергнуты проверке Лjung — Бокса на основе статистического показателя Q .

Обобщение модели ARCH — модель GARCH

Модель ARCH является относительно простой, но нередко для улавливания регулярностей волатильности временного ряда с возвратами от финансового актива требует много параметров. Обобщенная модель ARCH (generalized ARCH) применяется к временному ряду логарифмированных возвратов r_t , с возмущениями $\varepsilon_t = r_t - \mu$, которые подчиняются модели GARCH(p, q), если:

$$\varepsilon_t = \sigma_t z_t, \quad \sigma_t^2 = \omega + \sum_{i=1}^p \alpha_i \varepsilon_{t-i}^2 + \sum_{j=1}^q \beta_j \sigma_{t-j}^2, \quad z_t \sim \text{i.i.d.}$$

Модель GARCH(p, q) исходит из модели ARMA(p, q) относительно дисперсии члена ошибки ε_t .

Подобно моделям ARCH, хвостовое распределение процесса GARCH(1, 1) тяжелее, чем у нормального распределения. Данная модель сталкивается с теми же недостатками, что и модель ARCH. Например, она одинаково откликается на положительные и отрицательные шоки.

Выбор порядка сдвига

Для конфигурирования порядка сдвига моделей ARCH и GARCH используются квадраты остатков временного ряда, натренированных предсказывать среднее значение исходного ряда. Остатки центрированы на нуле, вследствие чего их квадраты также являются дисперсией. Затем необходимо обследовать графики автокорреляционной и частично-автокорреляционной функций квадратов остатков и выявить регулярности автокорреляции в дисперсии временного ряда.

Как строить модель прогнозирования волатильности

Разработка модели волатильности для временного ряда с возвратами от актива состоит из четырех шагов:

1. Построить модель ARMA временного ряда для финансового временного ряда, основываясь на внутрирядовой зависимости, выявляемой функциями автокорреляции и частичной автокорреляции.
2. Протестировать остатки модели на эффекты ARCH/GARCH, снова опираясь на функции автокорреляции и частичной автокорреляции в отношении временного ряда квадратов остатков.

3. Указать модель волатильности, если эффекты внутрирядовой корреляции значимы, и совместно вычислить уравнения среднего значения и волатильности.
4. Тщательно проверить вписанную модель и при необходимости доработать ее.



Во время применения прогнозирования волатильности к ряду финансовых возвратов внутрирядовая зависимость может быть ограничена, вследствие чего вместо модели ARMA может использоваться постоянное среднее.

Библиотека `arch` предоставляет несколько вариантов вычисления моделей прогнозирования волатильности. Она предлагает несколько вариантов моделирования математического ожидания среднего, включая постоянное среднее, модель $AR(p)$, обсуждавшуюся выше в разделе о моделях одномерных временных рядов, а также более поздние гетерогенные авторегрессионные процессы (*heterogeneous autoregressive process*, HAR), которые используют дневные (1 день), недельные (5 дней) и месячные (22 дня) сдвиги для улавливания торговых частот краткосрочных, среднесрочных и долгосрочных инвесторов.

Среднезначные модели² могут определяться и вычисляться совместно с несколькими моделями условной гетероскедастичности, которые помимо ARCH и GARCH включают экспоненциальную модель GARCH (*exponential GARCH*), допускающую асимметричные эффекты между положительными и отрицательными финансовыми возвратами, и гетерогенную модель ARCH (*heterogeneous ARCH*), которая дополняет среднезначную модель HAR.

Мы воспользуемся ежедневными финансовыми возвратами NASDAQ за 1998–2017 гг. и продемонстрируем использование модели GARCH (подробности см. в блокаде `arch_garch_models.ipynb`):

```
nasdaq = web.DataReader('NASDAQCOM', 'fred', '1998', '2017-12-31').squeeze()
```

```
# Перешкалировать для обеспечения оптимизации
nasdaq_returns = np.log(nasdaq).diff().dropna().mul(100)
```

Временной ряд с перешкалированными ежедневными финансовыми возвратами показывает только ограниченную автокорреляцию, но квадратные отклонения от среднего имеют существенную память, отраженную в медленно затухающей функции автокорреляции и высокой функции частичной автокорреляции для первых двух сдвигов и отсекающей только после первых шести сдвигов:

```
plot_correlogram(nasdaq_returns.sub(nasdaq_returns.mean()).pow(2),
                 lags=120,
                 title='NASDAQ Daily Volatility')
```

² Среднезначная модель (*mean model*) — это модель временного ряда, в которой предсказание будущих значений строится на основе выборочного среднего прошлых данных; она исходит из того, что будущие наблюдения будут взяты из того же самого распределения вероятностей. — *Прим. перев.*

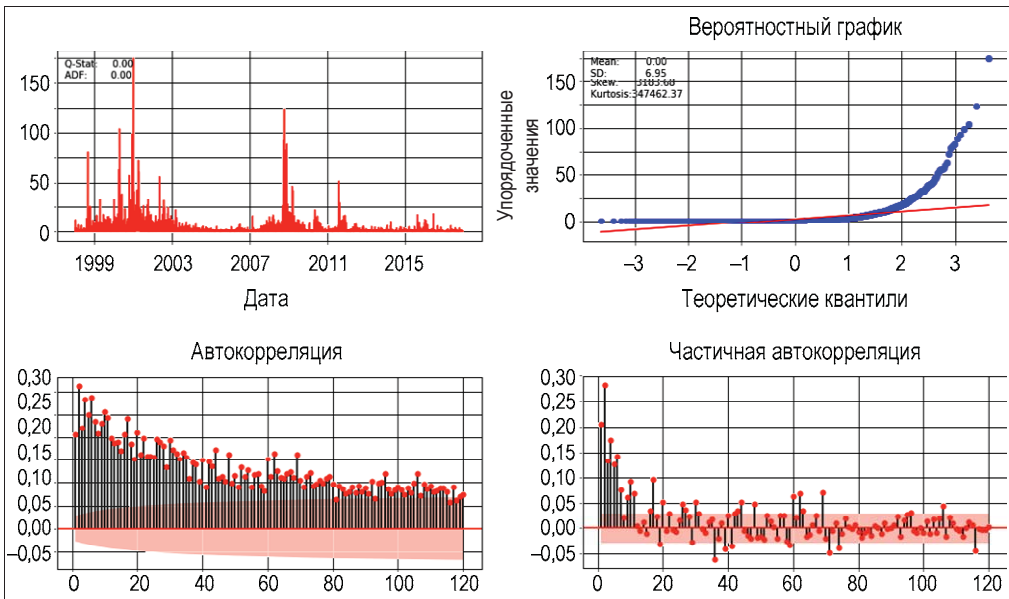


Рис. 8.6. Ежедневная волатильность биржи NASDAQ

Функция `plot_correlogram` на выходе производит следующий результат (рис. 8.6).

Следовательно, мы можем вычислить модель GARCH для улавливания линейной связи прошлых волатильностей. Мы будем использовать скользящие 10-летние окна для оценивания модели $GARCH(p, q)$ с p и q в интервале 1–4, генерируя вне-выборочные прогнозы на 1 шаг вперед. Затем мы сравним ошибку RMSE предсказанной волатильности относительно фактического квадрата отклонения финансового возврата от его среднего значения, выявляя наиболее предсказательную модель. Мы используем винсоризованные данные, лимитируя влияние экстремальных значений возврата, отраженных в очень высокой положительной асимметрии волатильности:

```
trainsize = 10 * 252 # 10 лет
data = nasdaq_returns.clip(lower=nasdaq_returns.quantile(.05),
                           upper=nasdaq_returns.quantile(.95))

T = len(nasdaq_returns)
test_results = {}
for p in range(1, 5):
    for q in range(1, 5):
        print(f'{p} | {q}')
        result = []
        for s, t in enumerate(range(trainsize, T-1)):
            train_set = data.iloc[s: t]
            test_set = data.iloc[t+1] # прогноз на 1 шаг вперед
            model = arch_model(y=train_set, p=p, q=q).fit(dispen='off')
            forecast = model.forecast(horizon=1)
```

```

mu = forecast.mean.iloc[-1, 0]
var = forecast.variance.iloc[-1, 0]
result.append([(test_set-mu)**2, var])
df = pd.DataFrame(result, columns=['y_true', 'y_pred'])
test_results[(p, q)] = np.sqrt(mean_squared_error(df.y_true, df.y_pred))

```

Модель GARCH(2, 2) достигает наименьшей ошибки RMSE (то же значение, что и GARCH(4, 2), но с меньшим числом параметров), поэтому мы идем дальше и выбираем эту модель, чтобы проинспектировать статистическую сводку:

```

am = ConstantMean(nasdaq_returns.clip(lower=nasdaq_returns.quantile(.05),
                                     upper=nasdaq_returns.quantile(.95)))
am.volatility = GARCH(2, 0, 2)
am.distribution = Normal()
model = am.fit(update_freq=5)
print(model.summary())

```

Выход показывает максимизированное логарифмическое правдоподобие, а также критерии AIC и BIC, которые обычно минимизируются при отборе моделей на основе внутривыборочной результативности (см. главу 7). Он также показывает результат для среднечисленной модели, которая в данном случае является просто постоянной оценкой, а также параметры GARCH для постоянной ω , параметры AR, α , и параметры MA, β , все из которых являются статистически значимыми:

Constant Mean - GARCH Model Results

| | | | |
|----------------|--------------------|-------------------|----------|
| Dep. Variable: | NASDAQCOM | R-squared: | -0.001 |
| Mean Model: | Constant Mean | Adj. R-squared: | -0.001 |
| Vol Model: | GARCH | Log-Likelihood: | -7484.02 |
| Distribution: | Normal | AIC: | 14980.0 |
| Method: | Maximum Likelihood | BIC: | 15019.0 |
| | | No. Observations: | 4852 |
| Date: | Wed, Oct 31 2018 | Df Residuals: | 4846 |
| Time: | 17:46:35 | Df Model: | 6 |

Mean Model

| | coef | std err | t | P> t | 95.0% Conf. Int. |
|----|--------|-----------|-------|-----------|------------------------|
| mu | 0.0521 | 1.491e-02 | 3.491 | 4.804e-04 | [2.284e-02, 8.130e-02] |

Volatility Model

| | coef | std err | t | P> t | 95.0% Conf. Int. |
|----------|--------|-----------|-------|-----------|-------------------------|
| omega | 0.0196 | 8.287e-03 | 2.365 | 1.804e-02 | [3.354e-03, 3.584e-02] |
| alpha[1] | 0.0247 | 1.470e-02 | 1.678 | 9.340e-02 | [-4.148e-03, 5.346e-02] |
| alpha[2] | 0.0627 | 2.196e-02 | 2.853 | 4.324e-03 | [1.962e-02, 0.106] |
| beta[1] | 0.5648 | 0.181 | 3.120 | 1.806e-03 | [0.210, 0.920] |
| beta[2] | 0.3337 | 0.180 | 1.853 | 6.393e-02 | [-1.932e-02, 0.687] |

Covariance estimator: robust

Теперь рассмотрим модели для многочисленных временных рядов и концепцию коинтеграции, что позволит задействовать новую торговую стратегию.

Модели многомерных временных рядов

Модели многомерных временных рядов служат для одновременного улавливания динамики многочисленных временных рядов и использования зависимостей между этими рядами для получения более надежных предсказаний.

Система уравнений

Модели одномерных временных рядов, такие как подход ARMA, который мы только что обсудили, ограничены статистическими связями между целевой переменной и ее сдвинутыми значениями или сдвинутыми возмущениями и экзогенными рядами в случае ARMAX. В отличие от них модели многомерных временных рядов также дают возможность сдвинутым значениям других временных рядов влиять на цель. Этот эффект применим ко всем рядам, приводя к сложным взаимодействиям, как показано на рис. 8.7.

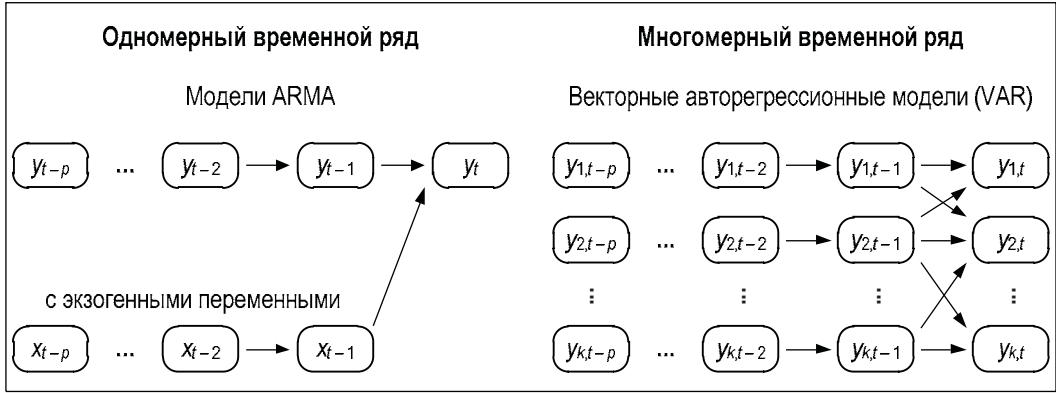


Рис. 8.7. Сложные взаимодействия модели многомерных временных рядов

В дополнение к потенциально более качественному прогнозированию многомерные временные ряды также используются для проникновения в сущность зависимостей между рядами. Например, в экономике многомерные временные ряды используются для понимания того, как изменения политики в отношении одной переменной, например процентной ставки, могут влиять на другие переменные на разных горизонтах. Функция импульсного отклика (импульсно-переходная), производимая с помощью многомерной модели, служит этой цели и позволяет симулировать то, как одна переменная откликается на внезапное изменение в других переменных. Концепция причинности Гренджера анализирует полезность одной переменной для прогнозирования другой (в смысле наименьших квадратов). Более того,

модели многомерных временных рядов допускают разложение дисперсии ошибки предсказания с целью анализа вклада других рядов.

Векторная авторегрессионная модель (VAR)

Мы увидим, как векторная авторегрессионная модель $\text{VAR}(p)$ расширяет модель $\text{AR}(p)$ до k рядов, создавая систему k уравнений, каждое из которых содержит p сдвинутых значений всех k рядов. В простейшем случае модель $\text{VAR}(1)$ для $k = 2$ принимает следующий вид:

$$\begin{aligned} y_{1,t} &= c_1 + a_{11}y_{1,t-1} + a_{12}y_{2,t-1} + \varepsilon_{1,t}; \\ y_{2,t} &= c_2 + a_{21}y_{1,t-1} + a_{22}y_{2,t-1} + \varepsilon_{2,t}. \end{aligned}$$

Эта модель может быть выражена несколько более сжато в матричном виде:

$$\begin{bmatrix} y_{1,t} \\ y_{2,t} \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} y_{1,t-1} \\ y_{2,t-1} \end{bmatrix} + \begin{bmatrix} \varepsilon_{1,t} \\ \varepsilon_{2,t} \end{bmatrix}.$$

Коэффициенты на собственных временных сдвигах предоставляют информацию о динамике самого ряда, тогда как коэффициенты в срезе переменных дают некоторое представление о сущности взаимодействий между рядами. Это обозначение распространяется на k рядов и порядок p следующим образом:

$$\mathbf{y}_t = \mathbf{c} + \mathbf{A}_1 \mathbf{y}_{t-1} + \dots + \mathbf{A}_p \mathbf{y}_{t-p} + \boldsymbol{\varepsilon}_t.$$

$k \times 1 \quad k \times 1 \quad k \times k \quad k \times 1 \quad k \times k \quad k \times 1 \quad k \times 1$

Модели $\text{VAR}(p)$ также требуют стационарности, чтобы начальные шаги со стороны моделирования одномерного временного ряда переносились дальше. Сначала следует разведать ряды и определить необходимые преобразования, затем применить расширенную проверку Дики — Фуллера, чтобы убедиться, что критерий стационарности удовлетворяется для каждого ряда, и в противном случае применить дальнейшие преобразования. Данная модель может быть вычислена с помощью обычных наименьших квадратов (OLS) в зависимости от первоначальной информации либо с помощью максимального правдоподобия, что эквивалентно для нормально распределенных ошибок, но не в противном случае.

Если некоторые или все k рядов являются единично-корневыми нестационарными, то они могут быть коинтегрированы. Данное расширение понятия единичного корня на многочисленные временные ряды означает, что линейная комбинация двух или более рядов является стационарной и, следовательно, среднеразворотной. Модель VAR не оснащена для обработки этого случая без исчисления последовательных разностей, вместо этого следует использовать векторную модель исправления ошибок (Vector Error Correction model, VECM, см. справочные материалы в репозитории GitHub). Мы продолжим разведывательный анализ коинтеграции, т. к. если она присутствует и считается устойчивой, то ее можно использовать для стратегии парной торговли.

Определение порядка временного сдвига также черпает свои подсказки из функций автокорреляции и частичной автокорреляции для каждого ряда, но ограничено тем фактом, что ко всем рядам применяется один и тот же порядок сдвига. После оценивания модели диагностика остатков также четко намекает на результат, напоминающий белый шум, и процедура отбора модели может использовать внутривыборочные информационные критерии либо, что предпочтительно, вневыборочную предсказательную результативность для перекрестного контроля альтернативных моделей, если конечной целью является применение модели для предсказания.

Как уже упоминалось в одномерном случае, предсказания исходного временного ряда требуют от нас обратных преобразований, применяемых для придания ряду стационарности перед началом тренировки модели.

Как использовать модель VAR для прогнозов макроэкономических фундаментальных показателей

Мы расширим одномерный пример с единственным временным рядом ежемесячных данных по промышленному производству и добавим ежемесячные временные ряды по потребительским настроениям, предоставляемые службой данных Федеральной резервной системы. Для получения данных с 1970 по 2017 гг. мы воспользуемся уже знакомой нам библиотекой `pandas-datareader`:

```
df = web.DataReader(['UMCSENT', 'IPGMFN'], 'fred', '1970', '2017-12').dropna()
df.columns = ['sentiment', 'ip']
```

Логарифмическое преобразование ряда с данными промышленного производства и сезонным исчислением последовательных разностей с использованием временного сдвига 12 в обоих рядах дает стационарные результаты:

```
df_transformed = pd.DataFrame({'ip': np.log(df.ip).diff(12),
                              'sentiment': df.sentiment.diff(12)}).dropna()
test_unit_root(df_transformed) # см. блокнот относительно деталей
                               # и других графиков
```

| | |
|-----------|---------|
| | p-value |
| ip | 0.03% |
| sentiment | 0.00% |

В результате мы получим ряды, приведенные на рис. 8.8.

Для того чтобы ограничить размер результата на выходе, мы оценим только модель VAR(1), применив ее реализацию `VARMAX` в библиотеке `StatsModels` (которая допускает необязательные экзогенные переменные) с постоянным трендом, используя первые 480 наблюдений:

```
model = VARMAX(df_transformed.iloc[:468], order=(1,1), trend='c').fit(maxiter=1000)
print(model.summary())
```

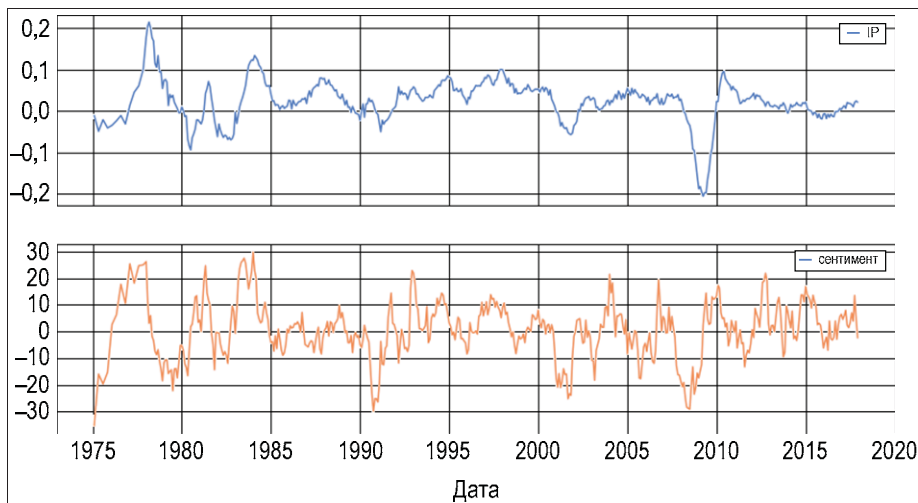


Рис. 8.8. Временные ряды промышленного производства и потребительских настроений

Это приводит к следующей статистической сводке:

Statespace Model Results

```
=====
Dep. Variable:      ['ip', 'sentiment']    No. Observations:      468
Model:              VARMA(1,1)             Log Likelihood         -71.824
                  + intercept              AIC                  169.647
Date:               Tue, 05 Feb 2019        BIC                   223.578
Time:               09:54:39               HQIC                 190.869
Sample:             0
                  - 468
Covariance Type:    opg
=====
```

```
Ljung-Box (Q):      128.08, 161.44         Jarque-Bera (JB): 130.30, 16.85
Prob(Q):            0.00, 0.00             Prob(JB):          0.00, 0.00
Heteroskedasticity (H): 0.48, 1.10         Skew:              0.20, 0.21
Prob(H) (two-sided): 0.00, 0.55           Kurtosis:          5.56, 3.83
```

Results for equation ip

```
=====
              coef    std err          z      P>|z|      [0.025      0.975]
-----
const          0.0015      0.001      2.418     0.016      0.000      0.003
L1.ip           0.9282      0.010     93.695     0.000      0.909      0.948
L1.sentiment    0.0006     6.02e-05    10.110     0.000      0.000      0.001
L1.e(ip)        0.0121      0.037      0.325     0.745     -0.061      0.085
L1.e(sentiment) -0.0001      0.000     -0.867     0.386     -0.000      0.000
=====
```

Results for equation sentiment

| | coef | std err | z | P> z | [0.025 | 0.975] |
|-------------------------|----------|---------|--------|-------|---------|--------|
| const | 0.3877 | 0.279 | 1.391 | 0.164 | -0.159 | 0.934 |
| L1.ip | -14.6492 | 5.444 | -2.691 | 0.007 | -25.320 | -3.979 |
| L1.sentiment | 0.8805 | 0.023 | 37.684 | 0.000 | 0.835 | 0.926 |
| L1.e(ip) | 39.0203 | 18.828 | 2.072 | 0.038 | 2.119 | 75.922 |
| L1.e(sentiment) | 0.0507 | 0.052 | 0.979 | 0.327 | -0.051 | 0.152 |
| Error covariance matrix | | | | | | |
| | coef | std err | z | P> z | [0.025 | 0.975] |
| sqrt.var.ip | 0.0129 | 0.000 | 40.347 | 0.000 | 0.012 | 0.014 |
| sqrt.cov.ip.sentiment | 0.0436 | 0.232 | 0.188 | 0.851 | -0.411 | 0.498 |
| sqrt.var.sentiment | 5.2755 | 0.149 | 35.506 | 0.000 | 4.984 | 5.567 |

Warnings:

[1] Covariance matrix calculated using the outer product of gradients (complex-step).

Как показано на предыдущей иллюстрации модели VAR(1), выход содержат коэффициенты для обоих уравнений временных рядов. Библиотека StatsModels предоставляет диагностические графики, позволяющие проверить соответствие остатков допущениям о белом шуме, которые в этом простом случае соблюдаются не точно (рис. 8.9).

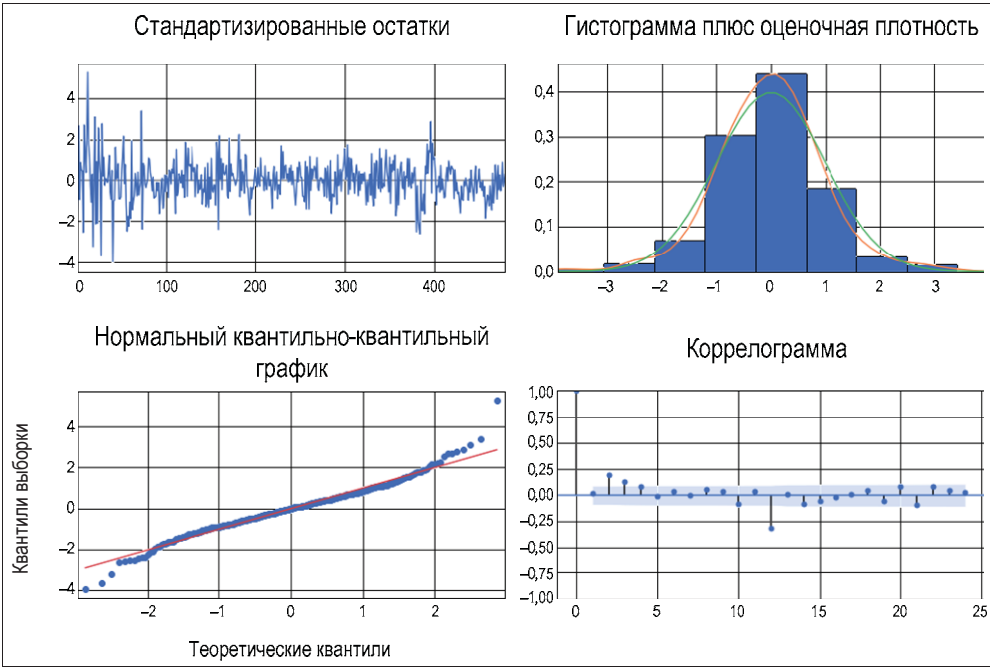


Рис. 8.9. Промышленное производство — диагностика

Вневыборочные предсказания могут генерироваться следующим образом:

```
preds = model.predict(start=480, end=len(df_transformed)-1)
```

Визуализация фактических и предсказанных значений показывает, как предсказание отстает от фактических значений и не достаточно хорошо улавливает нелинейные вневыборочные регулярности (рис. 8.10).

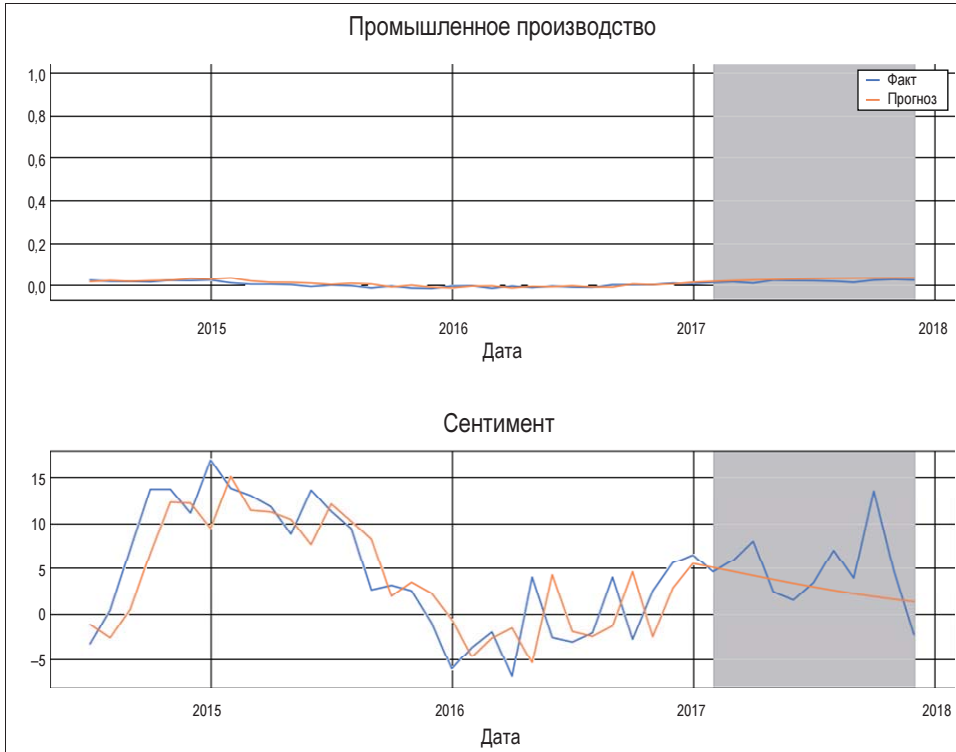


Рис. 8.10. Предсказание промышленного производства

Коинтеграция — временные ряды с общим трендом

Концепция интегрированных многомерных рядов осложняется тем, что все компонентные ряды процесса могут интегрироваться индивидуально, но процесс не является совместно интегрированным в том смысле, что существует одна или несколько линейных комбинаций рядов, которые производят новый стационарный ряд.

Другими словами, сочетание двух коинтегрированных рядов имеет стабильное среднее, к которому эта линейная комбинация разворачивается. Многомерный ряд с такой характеристикой принято считать коинтегрированным. Это также применимо, когда интегрируются индивидуальные ряды более высокого порядка, и линейная комбинация снижает совокупный порядок интеграции.

Коинтеграция отличается от корреляции: два ряда могут быть сильно коррелированы, но не обязательно будут коинтегрированы. Например, если два растущих ряда являются постоянными кратными друг другу, их корреляция будет высокой, однако любая линейная комбинация также будет расти и не будет разворачиваться к среднему значению.

Анализ VAR может по-прежнему применяться к интегрированным процессам, применяя форму исправления ошибок модели VAR, которая использует первые разности индивидуальных рядов плюс член исправления ошибок в уровнях.

Тестирование на коинтеграцию

Существуют два главных подхода к тестированию на коинтеграцию:

- ◆ двухэтапный метод Энга — Гренджера;
- ◆ процедура Йохансена.

Метод Энга — Гренджера охватывает регрессию одного ряда на другом, а затем применение единично-корневой расширенной проверки Дики — Фуллера (ADF) к остатку регрессии. Если нулевая гипотеза может быть отклонена, что позволяет допустить, что остатки являются стационарными, то ряды являются коинтегрированными. Ключевым преимуществом этого подхода является то, что коэффициент регрессии представляет собой множитель, который делает комбинацию стационарной, т. е. среднеразворотной. Мы вернемся к этому аспекту, когда будем задействовать коинтеграцию для стратегии парной торговли. С другой стороны, этот подход лимитирован выявлением коинтеграции пар временных рядов в отличие от более крупных групп рядов.

Процедура Йохансена, напротив, тестирует ограничения, налагаемые коинтеграцией на модель векторной авторегрессии (VAR), как обсуждалось в предыдущем разделе. Более конкретно, после вычитания целевого вектора из обеих сторон приведенного выше обобщенного уравнения $VAR(p)$ мы получаем формулировку модели исправления ошибок (error correction model, ECM):

$$\Delta y_t = \mathbf{c} + \mathbf{\Pi} y_{t-1} + \mathbf{\Gamma} \Delta y_{t-1} + \dots + \mathbf{\Gamma} \Delta y_{t-p} + \varepsilon_t.$$

Результирующее модифицированное уравнение $VAR(p)$ имеет только один векторный член в уровнях, т. е. не выражается как разность с помощью оператора Δ . Природа коинтеграции зависит от свойств матрицы коэффициентов $\mathbf{\Pi}$ этого члена, в частности от ее ранга. Хотя это уравнение выглядит структурно похожим на тестовые условия расширенной проверки Дики — Фуллера, теперь существует несколько потенциальных созвездий общих трендов и порядков интеграции, поскольку теперь участвуют многочисленные ряды. (Подробнее см. справочные материалы, приведенные в репозитории GitHub, в том числе касающиеся практических сложностей, связанных со шкалированием индивидуальных рядов.)

Как использовать коинтеграцию для стратегии парной торговли

Парная торговля основывается на стационарной, среднеразворотной взаимосвязи между ценами двух активов. Другими словами, соотношение или разность между двумя ценами, так называемый спред, может со временем разойтись, но в конечном счете вернуться на тот же уровень. При наличии такой пары стратегия состоит в том, чтобы лонговать (т. е. покупать) недостаточно результативный актив, т. к. для покрытия разрыва ему потребуется период повышенной результативности. Одновременно с этим следует шортить актив, который отошел от ценового якоря в положительном направлении, финансируя покупку.

Коинтеграция представляет именно такой тип стабильной взаимосвязи между двумя ценовыми рядами, заякоренными общим средним. Исходя из того, что коинтеграция будет оставаться устойчивой, в конечном счете должно последовать схождение, вызванное ростом курса акций с пониженной результативностью либо снижением курса акций с повышенной результативностью. Данная стратегия будет прибыльной независимо от того, какой из двух активов обладает добавочным преимуществом хеджирования от общих движений рынка в любом направлении.

Однако спред будет постоянно меняться, иногда расширяясь, а иногда сужаясь, или оставаться неизменным, поскольку оба актива движутся в унисон. Задача парной торговли состоит в поддержании хеджированной позиции путем корректировки относительных элементов содержимого портфеля по мере изменения спреда.

На практике при заданном универсуме активов стратегия парной торговли будет искать коинтегрированные пары путем прогона статистического теста на каждой паре. Ключевая задача здесь состоит в том, чтобы учесть смещения из-за многократного тестирования, как описано в *главе 6*. Библиотека StatsModels реализует как тест коинтеграции Энгла — Гренджера, так и тест Йохансена.

Для оценивания спреда следует выполнить линейную регрессию. Она даст коэффициент линейной комбинации двух интегрированных рядов с ценами финансовых активов, которые порождают стационарный комбинированный ряд. Как уже упоминалось, использование линейной регрессии для вычисления указанного коэффициента называется тестом коинтеграции Энгла — Гренджера.

Резюме

В этой главе мы развели линейные модели временных рядов для одномерного случая индивидуального ряда, а также многомерные модели для нескольких взаимодействующих рядов. Мы столкнулись с применениями, которые предсказывают макроэкономические фундаментальные показатели, моделями, которые прогнозируют волатильность активов или портфелей с широким применением в риск-

менеджменте, многомерными моделями VAR, которые улавливают динамику нескольких макрорядов, а также концепцией коинтеграции, которая лежит в основе популярной стратегии парной торговли.

Подобно предыдущей главе, мы стали свидетелями того, как линейные модели добавляют в модели много структуры, т. е. они делают сильные допущения, потенциально требующие преобразований и обширного тестирования с целью верификации соблюдения этих допущений. Если они соблюдаются, то тренировка и интерпретация моделей являются очень простыми, и модели обеспечивают хороший отправной случай, который способен улучшить более сложные модели, как мы увидим в последующих главах.

9

Байесово машинное обучение

В этой главе мы познакомим вас с байесовыми подходами к машинному обучению и тем, как их разные взгляды на неопределенность наращивают добавочную стоимость при разработке и оценивании стратегий алгоритмической торговли.

Байесова статистика позволяет квантифицировать неопределенность будущих событий и принципиальным образом уточнять наши оценки по мере поступления новой информации. Такой динамический подход хорошо адаптируется к эволюционирующей природе финансовых рынков. Это особенно полезно, когда релевантных данных меньше, и нам требуются методы, которые систематически интегрируют предшествующие (априорные) знания или допущения.

Мы увидим, что байесовы подходы к машинному обучению позволяют глубже проникать в сущность неопределенности вокруг статистических метрик, параметрических оценок и предсказаний. Их применения варьируются от более гранулярного риск-менеджмента до динамических обновлений предсказательных моделей, которые встраивают изменения в рыночной среде. Подход Блэка — Литтермана к размещению финансовых средств среди портфельных активов (*см. главу 5*) можно интерпретировать как байесову модель. Он вычисляет ожидаемый финансовый возврат как среднее значение рыночного равновесия и мнения инвестора, взвешенное на волатильность каждого актива, корреляции между активами и уверенность в каждом прогнозе.

Более конкретно в этой главе мы рассмотрим следующие темы:

- ◆ как байесова статистика применяется к машинному обучению;
- ◆ как использовать вероятностное программирование с помощью библиотеки PyMC3;
- ◆ как определять и тренировать автоматически обучающиеся модели;
- ◆ как выполнять современные методы генерирования выборок для выполнения приближенного вывода;
- ◆ как применять байесово машинное обучение для вычисления динамических коэффициентов Шарпа, построения байесовых классификаторов и оценивания стохастической волатильности.



Справочные материалы, ссылки на дополнительные материалы и примеры исходного кода для этой главы находятся в соответствующем каталоге репозитория GitHub. Следуйте инструкциям по установке, приведенным в *главе 1*.

Как работает байесово машинное обучение

Классическая статистика также называется частотной (фриквентистской), потому что она интерпретирует вероятность как относительную частоту события в течение продолжительного времени, т. е. после наблюдения большого числа испытаний. В контексте вероятностей событие представляет собой комбинацию одного или нескольких элементарных результатов эксперимента, таких как любой из шести равных результатов при бросании двух кубиков либо падение цены актива на 10% или более в определенный день.

Байесова статистика, напротив, рассматривает вероятность как меру уверенности или мнения о наступлении события. Байесов взгляд на вероятность оставляет больше места субъективным взглядам и, следовательно, различиям во мнениях, чем частотная интерпретация. Это различие наиболее разительно для событий, которые происходят недостаточно часто для того, чтобы приходиться к объективной мере долговременной частоты.

Иными словами, частотная статистика исходит из того, что данные являются случайной выборкой из популяции и призваны выявлять фиксированные параметры, которые сгенерировали эти данные. Байесова статистика, в свою очередь, принимает данные как данность и рассматривает параметры как случайные величины с распределением, которое можно вывести из данных. В результате частотные подходы требуют по крайней мере столько точек данных, сколько необходимо оценить параметров. Байесовы подходы, с другой стороны, совместимы с меньшими наборами данных и хорошо подходят для онлайн-ового самообучения, по одной выборке за раз.

Байесов взгляд очень полезен для многих реальных событий, которые редки или уникальны, по крайней мере в важных отношениях. Примерами могут служить результаты следующих выборов или вопрос о том, обрушатся ли рынки в течение трех месяцев. В каждом случае есть как соответствующие исторические данные, так и уникальные обстоятельства, которые разворачиваются по мере приближения события.

Сначала мы введем теорему Байеса, которая кристаллизует концепцию обновления мнений путем комбинации априорных допущений с новым эмпирическим наблюдением, и сравним результирующие параметрические оценки с их частотными аналогами. Затем мы продемонстрируем два подхода к байесовому статистическому выводу, которые позволяют проникать в сущность апостериорного распределения латентных, т. е. ненаблюдаемых, параметров, таких как их математические ожидания, при наличии разных обстоятельств:

1. Сопряженные априорные распределения обеспечивают процесс обновления, предоставляя решение в закрытой форме, но точные аналитические методы не всегда доступны.
2. Приближенный вывод симулирует распределение, которое получается в результате комбинирования допущений и данных, и использует выборки из этого распределения для вычисления статистического сущностного понимания.

Как обновлять допущения на основе эмпирического наблюдения

Теорема, которую преподобный Томас Байес сформулировал более 250 лет назад, использует фундаментальную теорию вероятностей для предписания того, каким образом вероятности или мнения должны меняться по мере поступления релевантной новой информации. Приведенная ниже цитата Джона Мейнарда Кейнса улавливает байесов образ мыслей:

"Когда факты меняются, я меняю свое мнение. А чем вы занимаетесь, сэр?"

Указанная теорема опирается на условную и полную вероятность и цепное правило (см. справочные материалы в репозитории GitHub со ссылками на источники, где дается обзор этих понятий).

Под мнением понимается единственный параметр или вектор параметров θ (имеваемых также гипотезами). Каждый параметр может быть дискретным или непрерывным. θ может быть одномерным статистическим показателем, таким как (дискретная) мода категориальной переменной или (непрерывное) среднее значение, либо более высокоразмерным множеством значений, таким как матрица ковариаций или веса глубокой нейронной сети.

Ключевым отличием от частотной статистики является то, что байесовы допущения выражаются в виде вероятностных распределений, а не значений параметров. Следовательно, в то время как частотный вывод фокусируется на точечных оценках, байесов вывод дает распределения вероятностей.

Как показано на рис. 9.1, теорема Байеса обновляет мнения о представляющих интерес параметрах путем вычисления апостериорного распределения вероятностей из следующих ниже входных данных:

- ◆ *априорное* распределение, т. е. полученное до (лат. *a priori*) наблюдений, показывает, насколько вероятной мы считаем каждую возможную гипотезу;
- ◆ *функция правдоподобия* дает вероятность наблюдать набор данных при наличии определенных значений параметров θ ;
- ◆ *наблюдения/сведения* измеряют уровень вероятности наблюдаемых данных при наличии всех возможных гипотез. Следовательно, они одинаковы для всех значений параметров и служат для нормализации числителя.

Апостериорное распределение, т. е. полученное после (лат. *a posteriori*) наблюдений, представляет собой произведение априорной вероятности и правдоподобия,

деленное на наблюдение, и отражает обновленное распределение вероятностей гипотез, принимая в расчет как априорные допущения, так и данные. Глядя на это с другой стороны, произведение априорного распределения и правдоподобия является результатом применения цепного правила, разлагающего совместное распределение данных и параметров на составляющие.

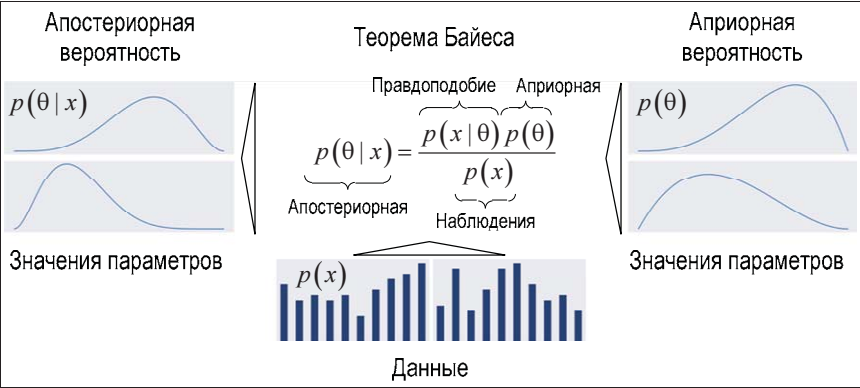


Рис. 9.1. Теорема Байеса

С более многомерными непрерывными величинами данная формулировка становится сложнее и включает (множественные) интегралы. Альтернативная формулировка использует шансы, выражая апостериорные шансы как произведение априорных шансов на коэффициент правдоподобия (подробности см. в справочных материалах).

Точный вывод: оценивание апостериорного максимума

Практические применения правила Байеса для точного вычисления апостериорных вероятностей довольно ограничены, поскольку вычислить член наблюдения в знаменателе сложно. Наблюдение (сведения) отражает вероятность наблюдаемых данных по всем возможным значениям параметров. Указанный член также называется маргинальным правдоподобием, поскольку требует маргинализации¹ распределения параметров путем сложения или интегрирования над их распределением. Это обычно возможно только в простых случаях с небольшим числом дискретных параметров, которые принимают очень мало значений.

Оценивание максимальной апостериорной вероятности (maximum a posteriori probability, MAP) задействует тот факт, что наблюдение является постоянным фак-

¹ Маргинализация (marginalizing out), или суммирование всех возможных конфигураций, состоит в том, чтобы в совместном распределении вероятностей находить распределение некой случайной величины путем исчерпания случаев ее появления вместе с другими случайными величинами. — Прим. перев.

тором, который шкалирует апостериорную вероятность, удовлетворяя требования к вероятностному распределению. Поскольку наблюдение не зависит от θ , апостериорное распределение является пропорциональным произведению правдоподобия и априорного распределения. Следовательно, оценивание максимальной апостериорной вероятности выбирает такое значение θ , которое максимизирует апостериорное распределение при наличии наблюдаемых данных и априорного мнения, т. е. моды апостериорного распределения.

Подход на основе оценивания максимальной апостериорной вероятности (MAP) контрастирует с *оцениванием максимального правдоподобия* (MLE) параметров, которое определяет распределение вероятностей. MLE подбирает такое значение параметра θ , которое максимизирует функцию правдоподобия для наблюдаемых тренировочных данных.

Если взглянуть на определения двух методов, то можно отметить, что MAP отличается от MLE тем, что включает априорное распределение. Другими словами, если априорное распределение не является постоянным, оценка θ методом MAP будет отличаться от ее аналога методом MLE:

$$\begin{aligned}\theta_{\text{MLE}} &= \arg \max_{\theta} P(X | \theta); \\ \theta_{\text{MAP}} &= \arg \max_{\theta} P(X | \theta) P(\theta).\end{aligned}$$

Решение методом MLE, как правило, отражает частотное представление о том, что оценки вероятностей должны отражать наблюдаемые соотношения. С другой стороны, влияние априорного распределения на оценку методом MAP часто соответствует добавлению в метод MLE данных, которые отражают априорные допущения. Например, сильное априорное суждение о том, что монета смещена, может быть встроено в контекст метода MLE путем добавления асимметричных данных испытаний.

Априорные распределения являются критически важным ингредиентом байесовых моделей. Теперь мы представим несколько удобных вариантов, обеспечивающих аналитический вывод.

Как отбирать априорные распределения

Априорное распределение должно отражать знание о распределении параметров, т. к. оно влияет на оценку методом MAP. Если априорное распределение достоверно не известно, то нам нужно сделать выбор, часто из нескольких разумных вариантов. В общем случае, хорошая практика состоит в том, чтобы обосновать априорное распределение и убедиться в его робастности путем проверки, что альтернативы приводят к такому же заключению.

Существует несколько типов априорных распределений.

- ◆ *Объективные априорные* распределения максимизируют влияние данных на апостериорное распределение. Если распределение параметров неизвестно, то мы можем выбрать неинформативное априорное распределение, такое как рав-

номерное распределение, также именуемое плоским априорным распределением, в соответствующем интервале значений параметров.

- ◆ Напротив, *субъективные априорные распределения* призваны встраивать в оценку информацию, внешнюю по отношению к модели.
- ◆ *Эмпирическое априорное распределение* сочетает байесов и частотный методы и использует исторические данные для устранения субъективности, например, путем оценивания разнообразных моментов, вписывая в стандартное распределение.

В контексте автоматически обучающейся модели априорное распределение можно рассматривать как регуляризатор, поскольку оно ограничивает значения, которые может принимать апостериорное распределение. Например, параметры, которые имеют нулевую априорную вероятность, не являются частью апостериорного распределения. Как правило, более хорошие данные позволяют делать более сильные заключения и снижают влияние априорного распределения.

Как не усложнять вывод — сопряженные априорные распределения

Априорное распределение является сопряженным относительно правдоподобия, когда результирующее апостериорное распределение имеет тот же тип распределения, что и априорное, за исключением разных параметров. Когда и априорное распределение, и правдоподобие являются нормально распределенными, то и апостериорное распределение тоже является нормально распределенным.

Сопряженность априорного распределения и правдоподобия приводит к решению апостериорного распределения в замкнутой форме, обеспечивая процесс обновления и позволяет избегать необходимости использования численных методов в аппроксимации апостериорного распределения. Более того, результирующее апостериорное распределение можно использовать как априорное распределение для следующего шага обновления.

Проиллюстрируем этот процесс на примере бинарной классификации движения цены акции.

Как динамически оценивать вероятности движения цены актива

Когда данные состоят из бинарных случайных величин Бернулли с определенной вероятностью успеха положительного результата, число успехов в повторных испытаниях подчиняется биномиальному распределению. Сопряженное априорное распределение — это бета-распределение с поддержкой в промежутке $[0; 1]$ и двумя параметрами формы для моделирования произвольных априорных распределений над вероятностью успеха. Следовательно, апостериорное распределение также является бета-распределением, которое мы можем вывести, непосредственно обновляя параметры.

Мы будем собирать выборки бинаризованных ежедневных финансовых возвратов фондового индекса S&P 500 разных размеров, где положительным результатом является повышение в цене. Начиная с неинформативного априорного распределения, которое выделяет равную вероятность каждой возможной вероятности успеха в промежутке $[0; 1]$, мы вычисляем апостериорное распределение для разных выборок наблюдений.

В следующем фрагменте кода показано, что обновление состоит из простого добавления наблюдаемых чисел успеха и неуспеха в параметры априорного распределения, в результате чего получается апостериорное распределение:

```
n_days = [0, 1, 3, 5, 10, 25, 50, 100, 500]

# Первоначальные 500 торговых дней
outcomes = sp500_binary.iloc[:n_days[-1]]
p = np.linspace(0, 1, 100)

# Равномерное (неинформативное) априорное распределение
a = b = 1
for i, days in enumerate(n_days):
    successes = outcomes.iloc[:days]
    theta_mle = successes.mean()
    up = successes.sum()
    down = days - up
    update = stats.beta.pdf(p, a + up, b + down)
```

На рис. 9.2 построены графики результирующих апостериорных распределений. Они иллюстрируют эволюцию от равномерного априорного распределения, которое рассматривает все вероятности успеха, как равновероятные, до все более пикообразного распределения.

После 500 испытаний вероятность сосредоточена вблизи фактической вероятности положительного движения на уровне 54,7% с 2010 по 2017 гг. Она также показывает небольшие различия между оценками MLE и MAP, где последние несколько тяготеют в сторону математического ожидания однородного априорного распределения (см. рис. 9.2).

На практике использование сопряженных априорных распределений лимитировано низкоразмерными случаями. В дополнение к этому, упрощенный подход на основе метода MAP избегает вычисления члена наблюдения, но имеет ряд недостатков, даже когда он доступен; он не возвращает распределение, которое позволило бы выводить меру неопределенности либо использовать его в качестве априорного. Следовательно, нам нужно прибегнуть не к точному выводу, а к аппроксимациям, с использованием численных методов и стохастического моделирования, которые мы введем далее.

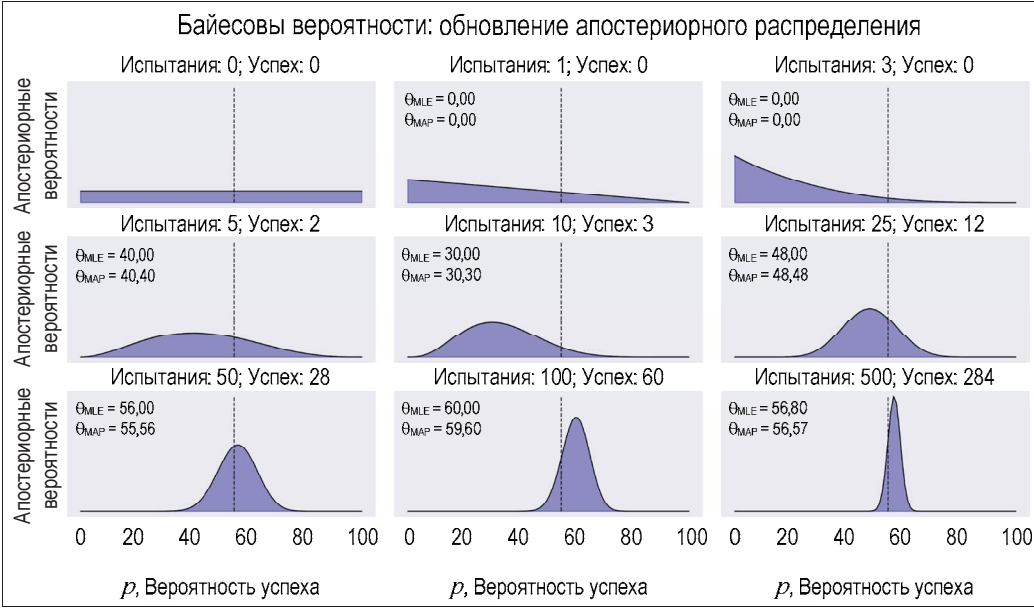


Рис. 9.2. Графики результирующих апостериорных распределений

Приближенный вывод: стохастический и детерминированный подходы

Для большинства моделей с практической значимостью получить точное апостериорное распределение аналитически и вычислить математические ожидания для латентных параметров невозможно. Модель может иметь слишком много параметров, либо апостериорное распределение может быть слишком сложным для аналитического решения. В случае непрерывных переменных их интегралы могут не иметь решений в замкнутой форме, в то время как размерность пространства и сложность интеграла могут налагать запрет на численное интегрирование. В случае дискретных переменных их маргинализации предусматривают суммирование всех возможных конфигураций скрытых величин, и хотя это всегда возможно в принципе, на практике мы часто обнаруживаем, что скрытых состояний может быть экспоненциально много, вследствие чего точный расчет является непомерно дорогим.

Хотя для некоторых приложений апостериорное распределение над ненаблюдаемыми параметрами будет представлять интерес, чаще всего в первую очередь требуется оценить математические ожидания, например, сделать предсказания. В таких ситуациях мы можем опираться на приближенный вывод.

♦ *Стохастические* технические решения, основанные на извлечении выборок методами *Монте-Карло марковских цепей* (Markov Chain Monte Carlo, MCMC), популяризировали использование байесовых методов во многих областях. Они, как правило, способны сходиться к точному результату. На практике методы из-

влечения выборок могут быть вычислительно сложными и часто ограничиваются мелкомасштабными задачами.

- ♦ *Детерминированные* методы, именуемые вариационным выводом или вариационным Байесом, основаны на аналитических аппроксимациях апостериорного распределения и способны хорошо масштабироваться до крупных приложений. Они принимают упрощенные допущения, например, о том, что апостериорное распределение разлагается определенным образом или имеет конкретную параметрическую форму, такую как гауссова. Следовательно, они не генерируют точных результатов и могут использоваться в качестве дополнения к методам извлечения выборок.

Стохастический вывод на основе выборок

Под сэмплированием подразумевается извлечение выборок $X = (x_1, \dots, x_n)$ из определенного распределения $p(x)$. Исходя из того, что выборки являются взаимно независимыми, закон больших чисел обеспечивает, что для растущего числа выборок доля конкретного экземпляра x_i в выборке (для дискретного случая) соответствует его вероятности $p(x = x_i)$. В непрерывном случае аналогичное рассуждение применимо к определенному участку пространства выборки. Следовательно, средние значения над выборками могут использоваться в качестве несмещенных оценщиков математических ожиданий параметров распределения.

Практическая задача заключается в обеспечении того, чтобы выборки при извлечении были взаимно независимыми, поскольку распределение неизвестно. Взаимно зависимые выборки по-прежнему могут быть несмещенными, но они тяготеют к увеличению дисперсии оценки, вследствие чего для такой же точной оценки, как и для взаимно независимых выборок, потребуется больше выборок.

Извлечение выборок из многомерного распределения является вычислительно затратным, поскольку число состояний экспоненциально увеличивается вместе с числом размерностей. Этот процесс обеспечивается целым рядом алгоритмов (см. справочные материалы с обзорами). Теперь мы представим несколько популярных вариантов методов, основанных на методах Монте-Карло марковских цепей (МСМС).

Генерирование выборок методами Монте-Карло марковских цепей

Цепь Маркова — это динамическая стохастическая модель, описывающая случайное блуждание по множеству состояний, связанных между собой переходными вероятностями. Марковское свойство предусматривает, что процесс не имеет памяти, а следующий шаг зависит только от текущего состояния. Другими словами, он обусловлен независимостью настоящего, прошлого и будущего, т. е. информация о прошлых состояниях не помогает предсказать будущее за пределами того, что мы знаем из настоящего.

Методы Монте-Карло опираются на многократное извлечение случайных выборок для аппроксимирования результатов, которые могут быть детерминированными, но которые не допускают аналитического, точного решения. Указанные методы были разработаны во время манхэттенского проекта для оценивания энергии на атомном уровне и получили свое устойчивое кодовое название для обеспечения секретности.

Многие алгоритмы применяют метод Монте-Карло к марковской цепи (MCMC) и обычно действуют следующим образом:

1. Начать с текущей позиции.
2. Извлечь новую позицию из распределения предлагаемых позиций.
3. Оценить вероятность новой позиции в свете данных и априорных распределений:
 - если вероятность достаточно высока, то перейти в новую позицию;
 - в противном случае остаться в текущей позиции.
4. Повторить шаг 1.
5. После заданного числа итераций вернуть все принятые позиции.

Метод MCMC стремится выявлять и разведывать интересующие участки апостериорного распределения, которые концентрируются на значительной плотности вероятности. Принято говорить, что процесс без памяти сходится, когда он неуклонно движется через близлежащие состояния высокой вероятности апостериорного распределения, где темп принятия новой позиции увеличивается. Одна из ключевых задач заключается в том, чтобы сбалансировать необходимость случайного разведывания пространства выборки с риском снижения темпа принятия.

Начальные шаги этого процесса, вероятно, будут более отражать исходную позицию, чем апостериорную, и обычно отбрасываются как образцы *обкатки* (burn-in). Ключевым свойством MCMC является то, что процесс должен забыть о своей первоначальной позиции после определенного (но неизвестного) числа итераций.

Остальные образцы называются следом процесса. Предполагая схождение, относительная частота выборок аппроксимирует апостериорное распределение и может использоваться для вычисления математических ожиданий с опорой на закон больших чисел.

Как указывалось ранее, прецизионность оценки зависит от внутрирядовой корреляции выборок, собранных случайным блужданием, каждая из которых, по своей конструкции, зависит только от предыдущего состояния. Более высокая корреляция лимитирует эффективное разведывание апостериорного распределения и должна быть подвергнута диагностическим тестам.

Общие технические решения для конструирования такой марковской цепи включают генерирование выборок по Гиббсу, генерирование выборок по Метрополису — Гастингсу и более поздние гамильтоновы методы MCMC, которые, как правило, работают лучше.

Генерирование выборок по Гиббсу

Выборка по Гиббсу сводит генерирование многомерных выборок к извлечению последовательности одномерных выборок. В исходной точке алгоритм итеративно держит $n - 1$ переменных постоянными, при этом выполняя выборку n -й величины. Он встраивает эту выборку и повторяет с начала.

Указанный алгоритм очень прост и легок в реализации, но создает сильно коррелированные выборки, которые замедляют сходжение. Его последовательная природа также не позволяет проводить параллелизацию.

Генерирование выборок по Метрополису — Гастингсу

Алгоритм Метрополиса — Гастингса случайно предлагает новые позиции на основе своего текущего состояния, эффективно разведывая пространство выборки и снижая корреляцию выборок относительно выборок, генерируемых по Гиббсу. Для того чтобы убедиться, что он генерирует выборку из апостериорного распределения, он оценивает предложенную позицию, используя произведение априорного распределения и правдоподобия, которое пропорционально апостериорному. Он принимает новую позицию с вероятностью, зависящей от результата, который относителен соответствующему значению для текущей выборки.

Ключевым преимуществом метода оценивания предложения новой позиции является то, что вместо точного оценивания апостериорного распределения он работает с пропорциональным его оцениванием. Однако для его схождения может потребоваться много времени, потому что случайные движения, которые не связаны с апостериорным распределением могут снизить темп принятия новой позиции, вследствие чего большое число шагов производит только небольшое число (потенциально коррелированных) выборок. Темп принятия может быть отрегулирован путем снижения дисперсии распределения предлагаемых позиций, однако из результирующего меньшего числа шагов следует меньшее разведывание пространства выборки.

Гамильтонов метод Монте-Карло — алгоритм NUTS

Гамильтонов метод Монте-Карло (Hamiltonian Monte Carlo, HMC) — это гибридный метод, который задействует информацию первопорядковой производной градиента правдоподобия для предложения новых состояний с целью разведывания пространства и преодоления некоторых сложностей MCMC. Кроме того, он встраивает импульс, который позволяет эффективно скакать по апостериорному распределению. В результате этого он быстрее сходится к высокоразмерному целевому распределению, чем более простые методы генерирования выборок по Метрополису или Гиббсу.

Генератор выборок NUTS (от англ. No-U-Turn — разворот на 180° запрещен) — это самонастраивающееся расширение гамильтонова метода Монте-Карло, которое адаптивно регулирует размер и число ходов по апостериорному распределению перед отбором предложения. Он хорошо работает на многомерных и сложных апостериорных распределениях и позволяет вписывать многочисленные сложные мо-

дели без специальных знаний о самом алгоритме подгонки. Как мы увидим в следующем далее разделе, данный генератор выборки используется в библиотеке PyMC3 по умолчанию.

Вариационный вывод

Вариационный вывод (variational inference, VI) — это метод МО, который аппроксимирует плотности вероятностей посредством оптимизации. В байесовом контексте он аппроксимирует апостериорное распределение следующим образом:

1. Отобрать параметризованное семейство вероятностных распределений.
2. Отыскать члена этого семейства, ближайшего к цели, измеряемого расхождением Кульбака — Лейблера².

По сравнению с MCMC, вариационный Байес тяготеет к сходимости быстрее и лучше масштабируется до крупных данных. В то время как MCMC аппроксимирует апостериорное распределение с помощью выборок из цепи, которые в итоге сойдутся бесконечно близко к цели, вариационные алгоритмы аппроксимируют апостериорное распределение с помощью результата оптимизации, чье совпадение с целью не гарантируется.

Вариационный вывод лучше подходит для крупных наборов данных и для быстрого разведывания большого числа моделей. В отличие от него MCMC будет давать более точные результаты для малых наборов данных или когда время и вычислительные ресурсы накладывают меньше ограничений.

Вариационный вывод с автоматическим дифференцированием

Недостатком вариационного вывода является необходимость в специфичных для модели математических вычислениях и реализации приспособленной для них процедуры оптимизации, что замедлило широкое его распространение.

Новейший алгоритм *вариационного вывода с автоматическим дифференцированием* (Automatic Differentiation Variational Inference, ADVI) автоматизирует этот процесс таким образом, что пользователь указывает только модель, выраженную в виде программы, а ADVI автоматически генерирует соответствующий вариационный алгоритм (подробности его реализации см. в справочных материалах в репозитории GitHub).

Мы увидим, что библиотека PyMC3 поддерживает всевозможные методы вариационного вывода, в том числе ADVI.

² Расхождение Кульбака — Лейблера (Kullback–Leibler Divergence), или информационное расхождение, — это неотрицательнозначный функционал, являющийся несимметричной мерой удаленности друг от друга двух вероятностных распределений, определенных на общем пространстве элементарных событий. — Прим. перев.

Вероятностное программирование с помощью библиотеки PyMC3

Вероятностное программирование предоставляет язык для описания и подгонки вероятностных распределений, что позволяет конструировать, кодировать и автоматически оценивать сложные модели. Оно ориентировано на абстрагирование от некоторой вычислительной и аналитической сложности, позволяя сосредотачиваться на концептуально более простых и интуитивных аспектах байесова рассуждения и вывода.

С появлением новых языков эта область стала весьма динамичной. Так, компания Uber предоставила открытый доступ к исходным кодам универсального языка вероятностного программирования Pyro (на основе библиотеки PyTorch), а компания Google недавно добавила вероятностный модуль в TensorFlow (справочные материалы со ссылками на ресурсы см. в репозитории GitHub).

Вследствие этого практическая значимость и использование байесовых методов в машинном обучении, вероятно, возрастут, позволив генерировать сущностное представление о неопределенности, и будет использоваться в случаях, которые, в частности, требуют именно прозрачных моделей, а не черно-ящичных.

В следующем разделе мы представим популярную библиотеку PyMC3, которая реализует расширенное генерирование выборок методом MCMC и вариационный вывод для автоматически обучающихся моделей с использованием языка Python. Вместе с языком Stan, названным в честь Станислава Улама, который изобрел метод Монте-Карло, и разработанным Эндрю Гельманом в Колумбийском университете в 2012 г., библиотека PyMC3 — это самый популярный язык вероятностного программирования.

Байесово машинное обучение с помощью библиотеки Theano

Библиотека PyMC3 была выпущена в январе 2017 г. с добавлением гамильтоновых методов Монте-Карло в генератор выборок по Метрополису — Гастингсу, который используется в версии библиотеки PyMC2 (выпущенной в 2012 г.). В качестве программно-вычислительного процессора библиотека PyMC3 использует библиотеку Theano, выполняющую динамическую компиляцию на C и автоматическое дифференцирование. Библиотека Theano — это матрично-ориентированная оптимизационная библиотека с поддержкой GPU, которая была разработана в Монреальском институте обучающихся алгоритмов Джошуа Бенджо (Montreal Institute for Learning Algorithms, MILA) и в свое время вдохновившая на разработку платформы TensorFlow. Указанный институт недавно прекратил дальнейшее развитие библиотеки Theano из-за успеха новых библиотек глубокого обучения (подробнее см. главу 17). Библиотека PyMC4, чей выход планируется на 2019 г., вместо нее будет ис-

пользовать платформу TensorFlow, с предположительно ограниченным влиянием на API.

Рабочий поток библиотеки PyMC3

Библиотека PyMC3 ориентирована на интуитивно понятный и читаемый, но мощный синтаксис, который отражает то, как статистики описывают модели. Процесс моделирования, как правило, состоит из следующих пяти этапов:

1. Закодировать вероятностную модель, определив следующее:
 - априорные распределения, которые квантифицируют знания о латентных переменных и их неопределенность;
 - функцию правдоподобия, которая обуславливает параметры на наблюдаемых данных.
2. Проанализировать апостериорно распределение с использованием одного из вариантов, описанных в предыдущем разделе:
 - получить точечную оценку с помощью вывода методом MAP;
 - извлечь выборку из апостериорного распределения методами MCMC.
3. Аппроксимировать апостериорное распределение с использованием вариационного Байеса.
4. Проверить модель с помощью различных диагностических инструментов.
5. Сгенерировать предсказания.

Результирующая модель может использоваться для вероятностного вывода с целью получения сущностного представления о значениях параметров, а также для предсказания результатов для новых точек данных.

Мы проиллюстрируем этот рабочий процесс с помощью простой логистической регрессии (см. блокнот `bayesian_logistic_regression.ipynb`). Впоследствии мы будем использовать библиотеку PyMC3 для вычисления и сравнения байесовых коэффициентов Шарпа, оценивания динамических коэффициентов парной торговли и реализации байесовых моделей линейных временных рядов.

Определение модели — байесова логистическая регрессия

Как обсуждалось в *главе 6*, логистическая регрессия оценивает линейную связь между множеством признаков и бинарным результатом, который опосредуется сигмоидальной функцией, обеспечивающей, чтобы модель производила вероятности. Частотный подход приводил к точечным оценкам параметров, которые измеряют влияние каждого признака на вероятность того, что точка данных принадлежит положительному классу, с доверительными интервалами, основанными на допущениях о распределении параметров.

В отличие от нее байесова логистическая регрессия оценивает апостериорное распределение над самими параметрами. Апостериорное распределение допускает

более робастные оценки того, что называется байесовым убедительным (credible) интервалом для каждого параметра, с преимуществом большей прозрачности относительно неопределенности модели.

Вероятностная программа состоит из наблюдаемых и ненаблюдаемых случайных величин. Как мы уже обсуждали, мы определяем наблюдаемые случайные величины через распределения правдоподобия и ненаблюдаемые случайные величины через априорные распределения. Для этой цели библиотека PyMC3 содержит многочисленные распределения вероятностей.

Мы будем использовать простой набор данных, который классифицирует 30 тыс. человек по доходам, используя порог 50 тыс. долларов в год. Этот набор данных будет содержать информацию о возрасте, поле, отработанных часах и годах обучения. Таким образом, используя эти признаки, мы моделируем вероятность того, что человек зарабатывает более 50 тыс. долларов.

Библиотека PyMC3 значительно упрощает выполнение приближенного байесова вывода для логистической регрессии. Логистическая регрессия моделирует вероятность того, что человек i имеет высокий доход, основываясь на k признаках, как показано в левой части рис. 9.3.



Рис. 9.3. Байесова логистическая регрессия

Мы воспользуемся Python-инструкцией контекстного менеджера `with` для того, чтобы определить ручную логистическую модель `manual_logistic_model`, на которую мы позже сможем ссылаться как на вероятностную модель:

1. Случайные величины для ненаблюдаемых параметров пересечения и двух признаков выражены с использованием неинформативных априорных распределений, которые исходят из нормальности распределений со средним значением, равным 0, и стандартным отклонением, равным 100.
2. Правдоподобие совмещает параметры с данными согласно спецификации логистической регрессии.
3. Результат моделируется как бернуллиева случайная величина с вероятностью успеха, задаваемой правдоподобием:

```

with pm.Model() as manual_logistic_model:
    # Случайные величины для коэффициентов с неинформативными
    # априорными распределениями для каждого параметра
    intercept = pm.Normal('intercept', 0, sd=100)
    b1 = pm.Normal('beta_1', 0, sd=100)
    b2 = pm.Normal('beta_2', 0, sd=100)

    # Преобразовать случайные величины
    # в вектор вероятностей p(y_i=1) согласно спецификации
    # логистической регрессионной модели.
    likelihood = pm.invlogit(intercept + b1 * data.hours + b2 * data.educ)

    # Бернуллиев случайный вектор
    # с вероятностью успеха, задаваемой сигмоидальной
    # функцией и фактическими наблюдаемыми данными
    pm.Bernoulli(name='logit', p=likelihood, observed=data.income)

```

Визуализация и фигурное обозначение

Команда `pm.model_to_graphviz(manual_logistic_model)` строит фигурную схему, приведенную на рис. 9.3 справа. Эта схема показывает ненаблюдаемые параметры как светлые овалы, а наблюдаемые элементы как темные. Прямоугольник обозначает число повторений наблюдаемого модельного элемента, вытекающее из данных, включенных в определение модели.

Модуль обобщенных линейных моделей

Библиотека PyMC3 содержит многочисленные общие модели, что позволяет оставлять ручную спецификацию для кастомизированных применений. Следующий фрагмент кода определяет ту же самую логистическую регрессию в качестве члена семейства обобщенных линейных моделей (Generalized Linear Model, GLM), используя формульный формат в духе статистического языка программирования R, который портируется на Python-библиотекой patsy:

```

with pm.Model() as logistic_model:
    pm.glm.GLM.from_formula('income ~ hours + educ', data,
                           family=pm.glm.families.Binomial())

```

Вывод методом MAP

Мы получаем точечные оценки MAP для трех параметров, используя метод `find_MAP()` только что определенной модели:

```

with logistic_model:
    map_estimate = pm.find_MAP()

print_map(map_estimate)

```

```

Intercept    -6.561862
hours         0.040681
educ          0.350390

```

Библиотека PyMC3 решает оптимизационную задачу нахождения апостериорной точки с наибольшей плотностью с помощью квазиньютоновского алгоритма Бройдена — Флетчера — Гольдфарба — Шанно (Broyden — Fletcher — Goldfarb — Shanno, BFGS), но предлагает несколько альтернатив, которые предоставляются библиотекой SciPy. Результат практически идентичен соответствующей оценке в библиотеке StatsModels (подробности см. в блокноте).

Приближенный вывод — метод Монте-Карло марковской цепи

Для иллюстрации вывода методом Монте-Карло марковской цепи мы воспользуемся немного более усложненной моделью:

```
formula = 'income ~ sex + age+ I(age ** 2) + hours + educ'
```

Функция `I()` библиотеки `patsy` позволяет использовать регулярные Python-выражения для создания новых переменных на лету. Здесь мы возводим возраст в квадрат для того, чтобы уловить нелинейную связь относительно пропорциональности: чем больше опыта, тем меньше добавляется дохода в более зрелом возрасте.

Следует отметить, что переменные, измеряемые в очень отличающихся шкалах, могут замедлять генерирование выборок. Следовательно, мы сначала применяем функцию `scale()` библиотеки `sklearn` для стандартизации величин возраста `age`, часов `hours` и образования `educ`.

Как только мы определили нашу модель с новой формулой, мы готовы выполнить вывод с целью аппроксимирования апостериорного распределения. Алгоритмы MCMC генерирования выборок доступны через функцию `pm.sample()`.

По умолчанию библиотека PyMC3 автоматически отбирает наиболее эффективный генератор выборок и инициализирует процесс генерирования выборок для эффективного схождения. Для непрерывной модели библиотека PyMC3 выбирает генератор выборок NUTS, который мы обсуждали в предыдущем разделе. Она также выполняет вариационный вывод через ADVI. Это делается с целью отыскания хороших исходных параметров для генератора выборок. Одним из нескольких вариантов является использование оценки MAP.

Для того чтобы посмотреть, как выглядит схождение, мы сначала извлекаем только 100 выборок после настройки генератора на 1000 итераций. Позже они будут отброшены. Процесс генерирования выборок можно параллелизовать для многочисленных цепей с помощью аргумента `core` (за исключением тех случаев, когда используется GPU):

```
with logistic_model:
    trace = pm.sample(draws=100, tune=1000,
                      init='adapt_diag', # альтернативная инициализация
                      chains=4, cores=2,
                      random_seed=42)
```

Результирующий след содержит значения, отобранные для каждой случайной величины. Мы можем продолжить генерировать выборки, передав след предыдущего прогона в качестве входа (подробности см. в блокноте).

Убедительные интервалы

Мы можем вычислить убедительные интервалы — байесов аналог доверительных интервалов — как процентиля следа. Результирующие границы отражают уверенность в диапазоне значения параметра для заданного вероятностного порога, в отличие от числа раз, когда параметр будет находиться внутри этого диапазона для крупного числа испытаний. Упомянутый выше блокнот иллюстрирует вычисления и визуализацию.

Приближенный вывод — вариационный Байес

Интерфейс вариационного вывода очень похож на реализацию МСМС. Вместо функции `sample()` мы просто используем функцию `fit()`, с возможностью включения функции обратного вызова досрочной остановки `CheckParametersConvergence`, если процесс подгонки распределения сошелся до заданного допуска:

```
with logistic_model:
    callback = CheckParametersConvergence(diff='absolute')
    approx = pm.fit(n=100000, callbacks=[callback])
```

Мы можем извлекать выборки из аппроксимированного распределения и получать объект следа, как мы делали ранее для генератора выборок МСМС:

```
trace_advi = approx.sample(10000)
```

Обследование сводной информации о следе показывает, что результаты несколько менее точны.

Диагностика модели

Диагностика байесовской модели включает в себя верификацию того, что процесс генерирования выборок сошелся и неуклонно отбирает выборки из высоковероятностных участков апостериорного распределения, и подтверждение того, что модель хорошо представляет данные.

Схождение

Мы можем визуализировать выборки во временной динамике и обследовать их распределения, проверив качество результатов. На диаграммах рис. 9.4 показаны апостериорные распределения соответственно после первоначальных 100 и дополнительных 100 тыс. выборок, и проиллюстрировано, как схождение приводит к тому, что многочисленные цепи идентифицируют одно и то же распределение. Функция `pm.trace_plot()` также показывает эволюцию выборок (подробности см. в блокноте).

Библиотека РуМС3 производит для генератора выборок различные сводные статистические показатели (табл. 9.1). Статистические показатели доступны как отдельные функции в модуле статистики `stats`, либо путем предоставления следа в функцию `pm.summary()`.

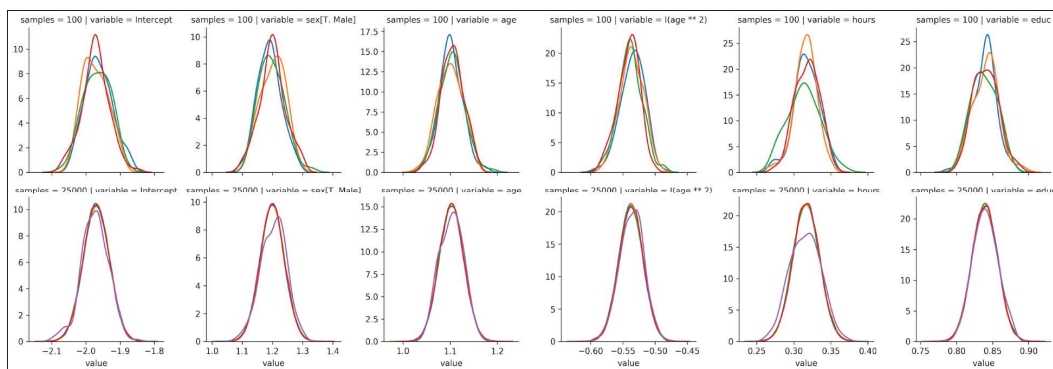


Рис. 9.4. Выборки во временной динамике и их распределения

Таблица 9.1. Сводные статистические показатели, порождаемые для генератора выборок

| statsmodels | mean | sd | mc_error | hpd_2,5 | hpd_97,5 | n_eff | Rhat |
|----------------------------|------------|----------|----------|------------|-----------|------------|----------|
| Intercept (пересечение) | −11,787778 | 1,777089 | 0,171152 | −13,785676 | −8,455286 | 14,010580 | 1,001814 |
| sex[T, Male] (пол) | 1,162391 | 0,063885 | 0,004056 | 1,061655 | 1,262571 | 63,721624 | 1,003156 |
| age (возраст) | 0,205541 | 0,071219 | 0,007049 | 0,058355 | 0,287755 | 11,705010 | 1,001347 |
| I(age ** 2) | −0,001822 | 0,000762 | 0,000075 | −0,002704 | −0,000238 | 11,544818 | 1,001402 |
| hours (часы) | 0,026945 | 0,002992 | 0,000193 | 0,023519 | 0,031850 | 171,499129 | 1,000883 |
| educ | 0,338701 | 0,030497 | 0,002267 | 0,315937 | 0,365790 | 95,228061 | 1,003296 |

Приведенная выше таблица в первом столбце содержит (вычисляемые в библиотеке StatsModels отдельно) коэффициенты `logit`, показывая, что в этом простом случае обе модели согласуются, потому что среднее значение выборки очень близко к коэффициентам.

Остальные столбцы содержат оценку *наибольшей апостериорной плотности* (highest posterior density, HPD) для минимальной ширины убедительного интервала, т. е. байесовой версии доверительного интервала, которая здесь вычисляется на уровне 95%. Статистический показатель `n_eff` подытоживает число эффективных (не отклоненных) выборок, получающееся в результате извлечения ~100 тыс. выборок.

Величина R-hat (R с крышкой), также именуемая статистическим показателем Гельмана — Рубина, проверяет схождение, сравнивая дисперсию между цепями с дисперсией внутри каждой цепи. Если генератор выборок сошелся, то эти дисперсии должны быть одинаковыми, т. е. цепи должны выглядеть одинаково. Следовательно, указанный статистический показатель должен находиться возле 1. Функция

`pm.forest_plot()` также подытоживает этот показатель для многочисленных цепей (подробности см. в блокноте).

Применительно к высокоразмерным моделям с многочисленными переменными проверка многочисленных следов становится громоздкой. При использовании алгоритма NUTS энергетический график помогает оценивать проблемы сходимости. Он резюмирует уровень эффективности, с которой случайный процесс разведывает апостериорное распределение. Такой график (рис. 9.5) показывает энергию и матрицу переходов энергии, которая должна хорошо согласовываться, как в следующем ниже примере (концептуальные подробности см. в справочных материалах).

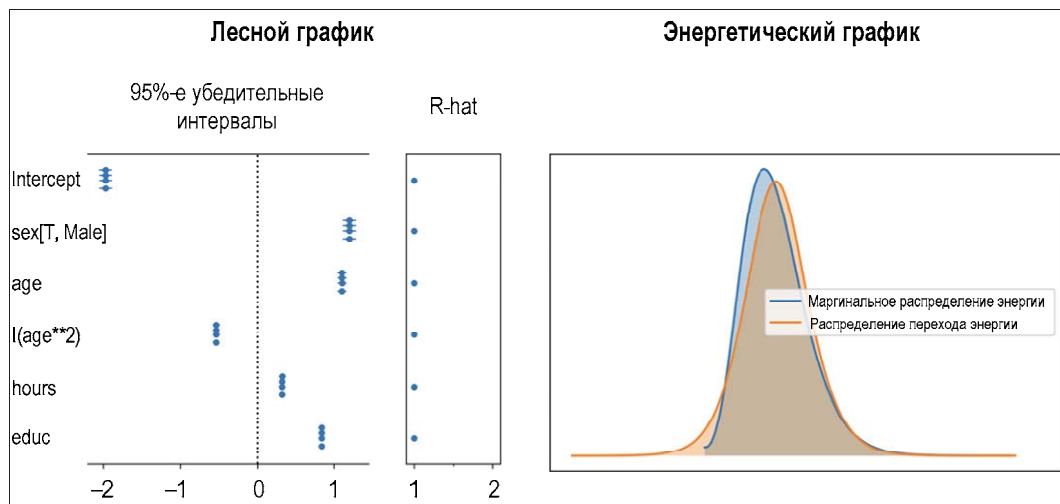


Рис. 9.5. Лесной и энергетический графики

Предсказательные проверки апостериорного распределения

Предсказательные проверки апостериорного распределения (Posterior Predictive Check, PPC) очень полезны для изучения того, насколько хорошо модель вписывается в данные. Они делают это, генерируя данные из модели, используя параметры из выборок, извлеченных из апостериорного распределения. Для этой цели мы используем функцию `pm.sample_ppc` и получаем n выборок для каждого наблюдения (модуль GLM автоматически называет результат как "y"):

```
ppc = pm.sample_ppc(trace_NUTS, samples=500, model=logistic_model)
ppc['y'].shape

(500, 29170)
```

Мы можем оценить внутривыборочную подгонку, используя отметку AUC, например, для сравнения разных моделей:

```
roc_auc_score(y_score=np.mean(ppc['y'], axis=0), y_true=data.income)

0.8294958565103577
```

Предсказание

Предсказания используют совместные переменные библиотеки Theano, заменяя тренировочные данные тестовыми перед выполнением предсказательных проверок апостериорного распределения. С целью обеспечения визуализации мы создаем переменную с одним-единственным предсказателем `hours`, создаем тренировочный и тестовый наборы данных и конвертируем первый из двух в совместную переменную. Обратите внимание, что нам нужно использовать массивы NumPy и предоставить список столбцовых меток (подробности см. в блокноте):

```
X_shared = theano.shared(X_train.values)

with pm.Model() as logistic_model_pred:
    pm.glm.GLM(x=X_shared, labels=labels,
               y=y_train, family=pm.glm.families.Binomial())
```

Затем мы запускаем генератор выборок, как и раньше, и применяем функцию `pm.sample_ppc` для результирующего следа после замены тренировочных данных тестовыми:

```
X_shared.set_value(X_test)
ppc = pm.sample_ppc(pred_trace, model=logistic_model_pred, samples=100)
```

Отметка AUC для этой модели с единственным признаком составляет 0,65. На графике рис. 9.6 показаны фактические результаты и неопределенность, окружающая предсказания для каждого выборочного значения предсказателя.

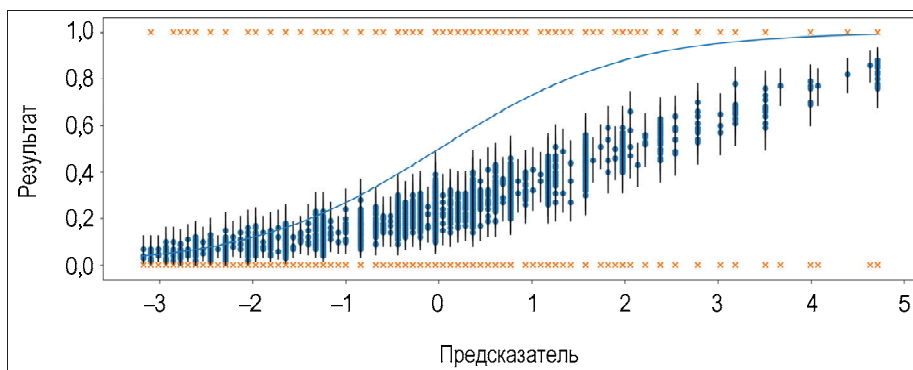


Рис. 9.6. Фактические результаты и неопределенность вокруг предсказаний

Теперь мы проиллюстрируем применение байесова анализа к примерам использования, связанным с торговлей на финансовых рынках.

Практические применения

Существует целый ряд приложений байесовых методов МО к инвестированию. Прозрачность, которую создают вероятностные оценки, естественным образом является полезной для риск-менеджмента и оценивания результативности активов.

Мы проиллюстрируем вычисление и сравнение такой метрики, как коэффициент Шарпа. Репозиторий GitHub также содержит два блокнота, на которые мы будем ссылаться далее. Эти блокноты представляют использование байесова МО для моделирования линейного временного ряда и стохастической волатильности.

Эти блокноты были адаптированы из учебных пособий, созданных в Quantopian, где Томас Вецки (Thomas Wiecki), который внес значительный вклад в популяризацию использования байесовых методов, руководит подразделением науки о данных. Справочные материалы также включают учебное руководство по использованию байесова МО для оценивания коэффициентов хеджирования парной торговли.

Байесов коэффициент Шарпа и сравнение результативности

В этом разделе мы проиллюстрируем способ определения коэффициента Шарпа как вероятностную модель и сравним результирующие апостериорные распределения для разных рядов с финансовыми возвратами. Байесово оценивание для двух групп обеспечивает полное распределение убедительных значений для размера эффекта, групповых средних и их разности, стандартных отклонений и их разности, а также нормальности данных.

Ключевые примеры использования содержат анализ разниц между альтернативными стратегиями, или между внутривыборочной финансовой возвратностью стратегии по отношению к ее вневыборочной возвратности (подробности см. в блокноте `bayesian_sharpe_ratio.ipynb`). Байесов коэффициент Шарпа также является составной частью байесовой сводной таблицы библиотеки `pyfolio`.

Определение модели

Для того чтобы смоделировать коэффициент Шарпа как вероятностную модель, нам нужны априорные сведения о распределении финансовых возвратов и параметрах, которые управляют этим распределением. *t*-Распределение Стьюдента демонстрирует толстые хвосты, которые относительно нормальному распределению для низких *степеней свободы* (*degrees of freedom*, *df*), и оно является разумным выбором для улавливания этого аспекта финансовых возвратов.

Следовательно, нам нужно смоделировать три параметра этого распределения, а именно среднее значение и стандартное отклонение возвратов, а также степени свободы. Мы будем исходить из нормального и равномерного распределений соответственно для среднего значения и стандартного отклонения и экспоненциального распределения для степени свободы с достаточно низким математическим ожиданием с целью обеспечения толстых хвостов. Возвраты основываются на этих вероятностных входах, а коэффициент Шарпа в годовом выражении является результатом стандартного вычисления, игнорируя безрисковую ставку (используя дневные возвраты):

```
mean_prior = data.stock.mean()
std_prior = data.stock.std()
std_low = std_prior / 1000
std_high = std_prior * 1000
```

```

with pm.Model() as sharpe_model:
    mean = pm.Normal('mean', mu=mean_prior, sd=std_prior)
    std = pm.Uniform('std', lower=std_low, upper=std_high)

    nu = pm.Exponential('nu_minus_two', 1 / 29, testval=4) + 2.
    returns = pm.StudentT('returns', nu=nu, mu=mean,
                          sd=std, observed=data.stock)

    sharpe = returns.distribution.mean / returns.distribution.variance ** .5 * np.sqrt(252)
    pm.Deterministic('sharpe', sharpe)

```

Блокнот содержит подробную информацию о генерировании выборок и оценивании коэффициента Шарпа для единственной акции.

Сравнение результативности

Для того чтобы сравнить результативность двух рядов возвратности, мы моделируем коэффициент Шарпа каждой группы отдельно и вычисляем размер эффекта как разницу между возвратностью с поправкой на волатильность. Визуализирование следов позволяет проникнуть в сущность распределений каждой метрики с точки зрения гранулярной результативности (рис. 9.7).

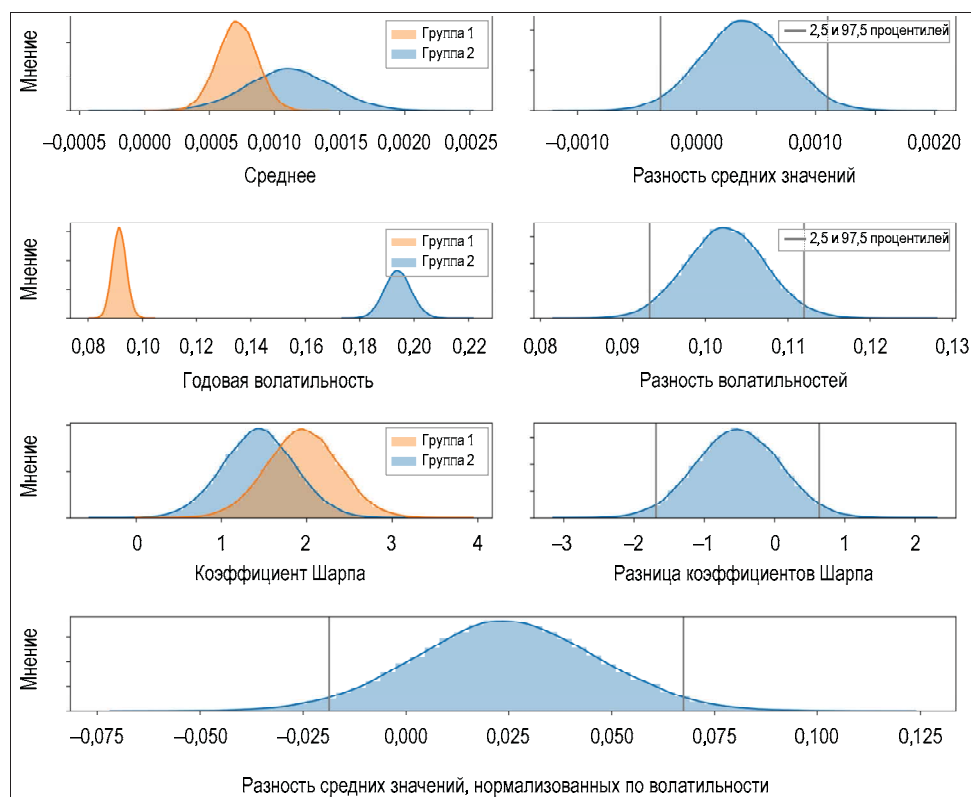


Рис. 9.7. Байесова линейная регрессия для парной торговли

В предыдущей главе мы представили парную торговлю как популярную стратегию алгоритмической торговли, основанную на коинтеграции двух или более активов. Принимая во внимание такие активы, нам нужно оценить хедж-коэффициент с целью принять решение об относительной величине длинных и коротких позиций. В базовом подходе используется линейная регрессия.

Блокнот `linear_regression.ipynb` иллюстрирует, как байесова линейная регрессия отслеживает изменения во взаимосвязи между двумя активами во временной динамике.

Байесовы модели временных рядов

Библиотека `PyMC3` включает в себя модели $AR(p)$, которые позволяют проникать в сущность неопределенности параметров точно так же, как и в предыдущих моделях. Блокнот `bayesian_time_series.ipynb` иллюстрирует модель временного ряда для одного или нескольких сдвигов.

Стохастические модели волатильности

Как обсуждалось в предыдущей главе, цены на активы имеют варьирующуюся во времени волатильность. В некоторые периоды возвраты сильно варьируются, в то время как в других они очень стабильны. Стохастические модели волатильности моделируют это с помощью латентной переменной волатильности, которая моделируется как стохастический процесс. Генератор выборок NUTS (No-U-Turn) был введен с использованием такой модели, и блокнот `stochastic_volatility.ipynb` иллюстрирует этот пример использования.

Резюме

В этой главе мы рассмотрели байесовы подходы к машинному обучению. Мы увидели, что у них есть несколько преимуществ, включая способность кодировать предварительные (априорные) знания или мнения, более глубокое проникновение в сущность неопределенности, окружающей модельные оценки и предсказания и их пригодность для онлайн-самообучения, где каждая тренировочная выборка инкрементно влияет на модельное предсказание.

Мы научились применять байесов рабочий поток, начиная со спецификации модели и вплоть до оценивания, диагностики и предсказания с помощью библиотеки `PyMC3`, и развели несколько релевантных приложений. Мы еще раз столкнемся с байесовыми моделями в *главе 14*, а также в *главе 19* в части, посвященной неконтролируемому глубокому автоматическому обучению, где мы введем вариационные автокодировщики.

В следующих двух главах будут представлены древесные, нелинейные ансамблевые модели, а именно случайные леса и градиентно-бустинговые машины.

10

Деревья решений и случайные леса

В этой главе мы познакомимся с двумя новыми классами автоматически обучающихся моделей: деревьями решений и случайными лесами. Мы увидим, как деревья решений усваивают из данных правила, которые кодируют нелинейные связи между входными и выходными переменными. Мы проиллюстрируем тренировку дерева решений и его применение для предсказания в регрессионных и классификационных задачах, визуализацию и интерпретацию правил, усвоенных моделью, и настройку модельных гиперпараметров с целью оптимизации компромисса между смещением и дисперсией и предотвращения переподгонки. Деревья решений являются не только важными автономными моделями, но и часто используются в качестве компонентов в других моделях.

Во второй части этой главы мы представим ансамблевые модели, которые совмещают многочисленные индивидуальные модели с целью получения единого агрегированного предсказания с более низкой дисперсией ошибки предсказания. Мы проиллюстрируем агрегирование бутстраповских выборок, часто именуемое бэггингом, как один из нескольких методов рандомизации конструирования индивидуальных моделей и снижения корреляции ошибок предсказания, делаемых компонентами ансамбля.

Бустинг, или ансамблевое усиление¹, — очень мощный альтернативный метод, который заслуживает своей главы для рассмотрения целого ряда новейших разработок (ему посвящена *глава 11*). Мы проиллюстрируем, как бэггинг эффективно снижает дисперсию, и научимся конфигурировать, тренировать и настраивать случайные леса. Мы увидим, как случайные леса, в качестве ансамбля из огромного числа деревьев решений, способны кардинально снижать ошибки предсказания за счет некоторой потери в интерпретации.

¹ Бустинг (boosting), или ансамблевое усиление, — это общая методика последовательной подгонки серии моделей путем предоставления для каждого последующего цикла большего веса элементам с большими остатками. — *Прим. перев.*

Иными словами, в этой главе мы рассмотрим следующие темы:

- ◆ как использовать деревья решений для регрессии и классификации;
- ◆ как проникать в суть деревьев решений и визуализировать правила принятия решения, усвоенные из данных;
- ◆ почему ансамблевые модели, как правило, обеспечивают превосходящие результаты;
- ◆ как бутстраповское агрегирование решает проблемы деревьев решений, связанные с переобучением;
- ◆ как тренировать, настраивать и интерпретировать случайные леса.

Деревья решений

Деревья решений — это автоматически обучающийся алгоритм, который предсказывает значение целевой переменной на основе правил принятия решения, усвоенных из тренировочных данных. Этот алгоритм может применяться как к регрессионным, так и к классификационным задачам путем изменения целевой функции, которая руководит тем, как указанное дерево усваивает правила принятия решения.

Мы обсудим вопросы, связанные с тем, как деревья решений используют правила для предсказаний, как их тренировать предсказывать (непрерывные) финансовые возвраты, а также (категориальные) направления ценовых движений и как их эффективно интерпретировать, визуализировать и настраивать.

Как деревья усваивают и применяют правила принятия решения

Линейные модели, которые мы изучали в *главах 7 и 8*, усваивают набор параметров для предсказания результата с использованием линейной комбинации входных переменных, возможно, после преобразования S-образной связывающей функцией в случае логистической регрессии.

Деревья решений используют другой подход: они усваивают и последовательно применяют набор правил, которые разбивают точки данных на подмножества, а затем делают одно предсказание для каждого подмножества. Предсказания основываются на выходных значениях для подмножества тренировочных образцов, которые получаются в результате применения заданной последовательности правил. Как мы увидим подробнее далее, классификационные деревья предсказывают вероятность, непосредственно оцениваемую по относительным частотам классов либо значению преобладающих классов, в то время как регрессионные модели вычисляют предсказание по среднему значению результатов для имеющихся точек данных.

Каждое такое правило опирается на один конкретный признак и использует порог для разбиения образцов на две группы со значениями ниже или выше порога по отношению к этому признаку. Бинарное дерево легко представляет логику этой мо-

дели: корень является отправной точкой для всех образцов, узлы — места применения правил принятия решения, а данные перемещаются по ребрам по мере их разбиения на меньшие подмножества до прибытия в листовой узел, где модель делает предсказание.

В линейной модели значения параметров позволяют интерпретировать влияние входных переменных на выход и модельное предсказание. Напротив, в дереве решений путь от корня к листьям создает прозрачность в отношении того, как признаки и их значения приводят к конкретным модельным решениям.

На рис. 10.1 показано, как модель усваивает правило. Во время тренировки алгоритм сканирует признаки и для каждого признака пытается найти отсечку. Эта отсечка разбивает данные, минимизируя потерю, возникающую в результате предсказаний, сделанных с использованием возникающих после разбивки подмножеств, взвешенную по числу образцов в каждом подмножестве.

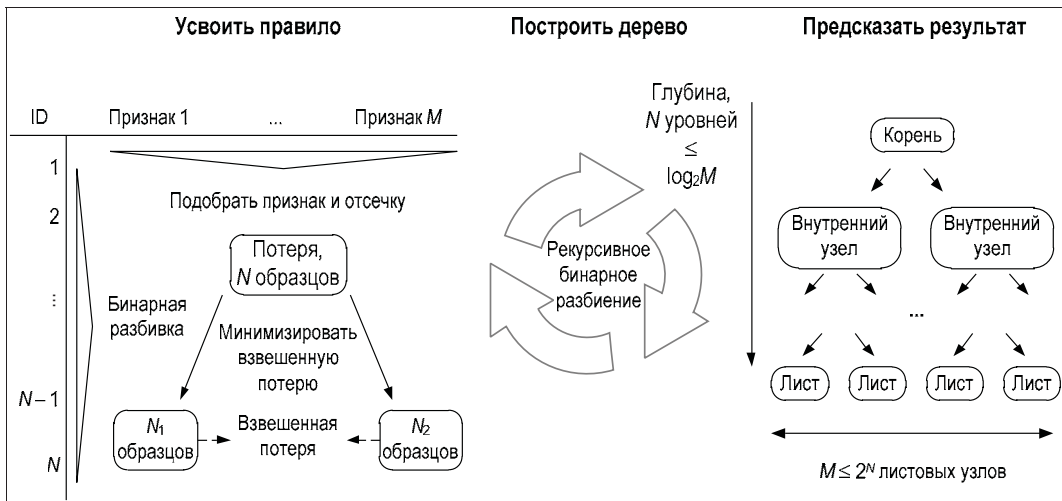


Рис. 10.1. Процесс усвоения правила принятия решения в дереве

Для построения полного дерева во время тренировки обучающийся алгоритм повторяет этот процесс деления признакового пространства, т. е. набора возможных значений p входных переменных X_1, X_2, \dots, X_p , на взаимоисключающие и коллективно исчерпывающие участки, каждый из которых представлен листовым узлом. К сожалению, алгоритм не сможет оценить все возможные разбиения признакового пространства, учитывая взрывное число возможных комбинаций последовательностей признаков и порогов. Для преодоления этого вычислительного предела древесное автоматическое обучение использует жадный нисходящий подход, именуемый *рекурсивным бинарным разбиением*.

Этот процесс является рекурсивным, поскольку он использует подмножества данных, полученных в результате предшествующих разбивок. Он является нисходящим, потому что начинается в корневом узле дерева, где все наблюдения пока при-

надлежат одному участку, а затем шаг за шагом создает две новые ветви дерева, добавляя в пространство предсказателей еще одну разбивку. Он является жадным, потому что указанный алгоритм не смотрит вперед и не оценивает потерю на несколько шагов вперед, а вместо этого подбирает лучшее правило в форме комбинации признака и порога, основываясь на непосредственном влиянии на целевую функцию. Мы вернемся к логике разбиения в более конкретном контексте регрессионных и классификационных деревьев, поскольку они существенно отличаются друг от друга.

Число тренировочных образцов продолжает сокращаться по мере того, как рекурсивные разбивки добавляют все новые узлы в дерево. Если правила разбивают образцы равномерно, давая в итоге идеально сбалансированное дерево с равным числом дочерних элементов в каждом узле, то на уровне n будет 2^n узлов с соответствующей долей от общего числа наблюдений каждый. На практике это является маловероятным, и число образцов вдоль некоторых ветвей может быстро уменьшаться, а деревья, как правило, расти до разных уровней глубины вдоль разных путей.

Для того чтобы достичь предсказания для нового наблюдения, модель использует правила, выведенные ею во время тренировки, которые принимают решение о том, какому листовому узлу должна быть назначена точка данных, а затем использует среднее (для регрессии) или моду (для классификации) тренировочных наблюдений в соответствующем участке признакового пространства. Меньшее число тренировочных образцов в заданном участке признакового пространства, т. е. в заданном листовом узле, снижает уверенность в предсказании и может свидетельствовать о переподгонке.

Рекурсивное разбиение будет продолжаться до тех пор, пока каждый листовой узел не будет содержать один-единственный образец и тренировочная ошибка не будет сведена к нулю. Мы введем несколько критериев для ограничения числа разбивок и предотвращения этой естественной тенденции деревьев решений производить крайнюю переподгонку.

Как использовать деревья решений на практике

В этом разделе мы проиллюстрируем применение древесных моделей для проникновения в сущность данных и выполнения предсказания. Для демонстрации регрессионных деревьев мы будем предсказывать финансовые возвраты, а в классификационном случае мы вернемся к примеру положительных и отрицательных движений цены актива. Примеры исходного кода для этого раздела находятся в блокноте `decision_trees.ipynb`, если не указано иначе.

Как готовить данные

Мы используем упрощенную версию набора данных, построенного в *главе 4*. Он состоит из ежедневных цен на акции, предоставляемых платформой Quandl за период 2010–2017 гг. и различных сконструированных признаков (см. блокнот

feature_engineering.ipynb из главы 4 в подпапке 00_data). Подробности см. в блокноте data_prep.ipynb настоящей главы в репозитории GitHub. Модели деревьев решений в этой главе не приспособлены для обработки пропущенных или категориальных переменных, поэтому к последним мы применим процедуру кодирования фиктивных значений (или кодирование с одним активным состоянием) после отбрасывания любой части из первых.

Как кодировать собственный класс перекрестного контроля

Мы также сконструируем собственный класс перекрестного контроля, адаптированный к формату созданных данных, который имеет мультииндекс библиотеки pandas с двумя уровнями: один для тикера и другой для данных.

```
class OneStepTimeSeriesSplit:
    """Генерирует кортежи пар (train_idx, test_idx).
        Исходит из того, что индекс содержит уровень 'date'"""

    def __init__(self, n_splits=3, test_period_length=1, shuffle=False):
        self.n_splits = n_splits
        self.test_period_length = test_period_length
        self.shuffle = shuffle
        self.test_end = n_splits * test_period_length

    @staticmethod
    def chunks(l, chunk_size):
        for i in range(0, len(l), chunk_size):
            yield l[i:i + chunk_size]

    def split(self, X, y=None, groups=None):
        unique_dates = (X.index
                        .get_level_values('date')
                        .unique()
                        .sort_values(ascending=False)[:self.test_end])

        dates = X.reset_index()[['date']]
        for test_date in self.chunks(unique_dates,
                                     self.test_period_length):
            train_idx = dates[dates.date < min(test_date)].index
            test_idx = dates[dates.date.isin(test_date)].index
            if self.shuffle:
                np.random.shuffle(list(train_idx))
            yield train_idx, test_idx

    def get_n_splits(self, X, y, groups=None):
        return self.n_splits
```


Класс `OneStepTimeSeriesSplit` обеспечивает разбивку на тренировочные и контрольные наборы и позволяет избежать смещения из-за забегания вперед, тренируя модели с использованием для каждой акции данных только вплоть до периода $t - 1$ и выполняя контроль с использованием данных за месяц T . Мы будем делать прогнозы только на один шаг вперед.

Как строить регрессионное дерево

Регрессионные деревья делают предсказания на основе среднего выходного значения для тренировочных образцов, закрепленных за конкретным узлом, и, как правило, опираются на среднеквадратическую ошибку для отбора оптимальных правил во время рекурсивного двоичного разбиения.

Получив тренировочный набор, алгоритм перебирает предсказатели X_1, X_2, \dots, X_p и возможные точки отсечения s_1, s_2, \dots, s_N , отыскивая оптимальную комбинацию. Оптимальное правило разбивает признаковое пространство на два участка $(X | X_i < s_j)$ и $(X | X_i > s_j)$ со значениями X_i -го признака ниже либо выше порога s_j , вследствие чего предсказания, основанные на тренировочных подмножествах, максимизируют снижение квадратов остатков относительно текущего узла.

В целях обеспечения визуализации начнем с упрощенного примера и будем использовать только два месяца сдвинутых финансовых возвратов для предсказания следующего месяца, в духе модели AR(2) из предыдущей главы:

$$r_t = f(r_{t-1}, r_{t-2}).$$

С помощью библиотеки `sklearn` конфигурирование и тренировка регрессионного дерева выполняются очень просто:

```
from sklearn.tree import DecisionTreeRegressor

# Сконфигурировать регрессионное дерево
regression_tree = DecisionTreeRegressor(criterion='mse', # по умолчанию
                                       max_depth=4,      # до 4 разбивок
                                       random_state=42)

# Создать тренировочные данные
y = data.returns
X = data.drop('returns', axis=1)
X2 = X.loc[:, ['t-1', 't-2']]

# Вписать древесную модель
regression_tree.fit(X=X2, y=y)

# Вписать OLS-модель (на основе обычных наименьших квадратов)
ols_model = sm.OLS(endog=y, exog=sm.add_constant(X2)).fit()
```

Сводная статистика обычных наименьших квадратов (OLS) и визуализация первых двух уровней дерева решений показывают поразительные различия между моделя-

ми. Модель OLS предоставляет три параметра для пересечения и двух признаков в соответствии с линейным допущением, которое эта модель делает о функции f .

В отличие от нее диаграмма регрессионного дерева для каждого узла первых двух уровней показывает признак и порог, используемые для разбиения данных (обратите внимание, что признаки могут использоваться повторно), а также текущее значение среднеквадратической ошибки (MSE), число образцов и предсказанное значение на основе этих тренировочных образцов (рис. 10.2).

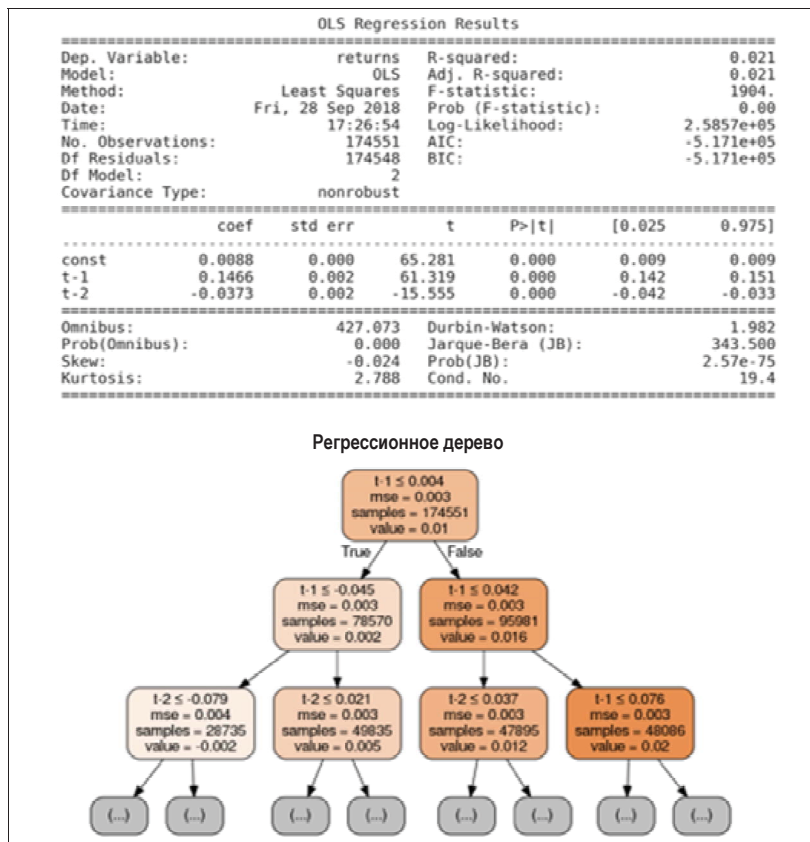


Рис. 10.2. Сводная статистика обычных наименьших квадратов (OLS) и визуализация первых двух уровней дерева решений

Древесная диаграмма также подчеркивает неравномерное распределение образцов по узлам, поскольку числа варьируются от 28 до 49 тыс. образцов всего после двух разбиений.

Для дальнейшей иллюстрации разных допущений о функциональной форме связей между входными переменными и выходом мы можем визуализировать предсказания текущих финансовых возвратов как функцию от признакового пространства, т. е. как функцию интервала значений для сдвинутых финансовых возвратов. На

рис. 10.3 показан возврат текущего периода как функция возвратов один и два периода назад для линейной регрессии и регрессионного дерева.

Результат линейной регрессионной модели на рис. 10.3 слева подчеркивает линейность связи между сдвинутым и текущим финансовыми возвратами, в то время как диаграмма регрессионного дерева справа иллюстрирует нелинейную связь, закодированную в рекурсивном разбиении признакового пространства.

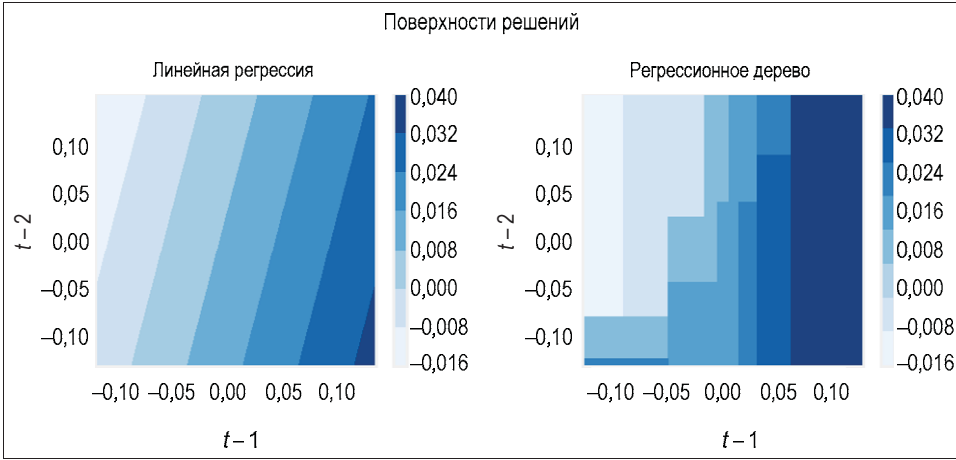


Рис. 10.3. Поверхности решений: финансовый возврат текущего периода как функция возвратов один и два периода назад для линейной регрессии и регрессионного дерева

Как строить классификационное дерево

Классификационное дерево работает точно так же, как и его регрессионная версия, за исключением того, что категориальная природа результата требует другого подхода к предсказанию и измерению потери. В то время как регрессионное дерево предсказывает отклик для наблюдения, назначенного листовому узлу, используя средний результат ассоциированных тренировочных образцов, классификационное дерево вместо этого использует моду, т. е. наиболее распространенный класс среди тренировочных образцов в соответствующем участке. Классификационное дерево также может генерировать вероятностные предсказания на основе относительных частот классов.

Как оптимизировать чистоту узлов

Во время выращивания классификационного дерева мы точно так же используем рекурсивное бинарное разбиение, но вместо оценки качества правила принятия решения, применяя снижение среднеквадратической ошибки, мы можем использовать частоту появления ошибок неправильной классификации, т. е. просто долю тренировочных образцов в определенном (листовом) узле, которые не принадлежат наиболее распространенному классу.

Однако альтернативные меры (индекс смешанности Джини или перекрестная энтропия) бывают предпочтительнее, поскольку они чувствительнее к *чистоте узла*, чем к частоте появления ошибок неправильной классификации. Под чистотой узла имеется в виду степень преобладания одного класса в узле. Узел, содержащий только образцы с результатами, принадлежащими одному классу, является чистым, и из этого следует успешное классифицирование для этого конкретного участка признакового пространства. Указанные меры рассчитывают, как показано ниже для результата классификации, выбором K значений: $0, 1, \dots, K-1$, для определенного узла m , представляющего участок R_m признакового пространства, где p_{mk} — это доля результатов класса k в узле m :

$$\text{смешанность Джини} = \sum_k p_{mk} (1 - p_{mk});$$

$$\text{перекрестная энтропия} = \sum_k p_{mk} \log(p_{mk}).$$

Обе меры: смешанность Джини и перекрестная энтропия — принимают меньшие значения, когда классовые пропорции приближаются к нулю или единице, т. е. когда в результате разбивки дочерние узлы становятся чистыми, и являются самыми высокими, когда классовые пропорции одинаковы, или равны 0,5 в бинарном случае. Диаграмма в конце данного раздела визуализирует значения, принимаемые этими двумя мерами, и частоты появления ошибок неправильной классификации по всему промежутку $[0; 1]$ пропорций.

Как тренировать классификационное дерево

Теперь мы натренируем, визуализируем и оценим классификационное дерево с пятью последовательными разбивками, используя 80% образцов для тренировки и оставшиеся 20% для предсказания. Для упрощения иллюстрации здесь мы принимаем укороченный путь и используем встроенный метод разбивки `train_test_split`, который не защищает от систематического смещения из-за заглядывания вперед, как наш собственный итератор. Из конфигурации дерева следует до $2^5 = 32$ листовых узлов, которые в сбалансированном случае в среднем будут содержать более 4300 тренировочных образцов. Взгляните на следующий фрагмент кода:

```
# Рандомизировать разбивку на тренировочные и тестовые наборы
X_train, X_test, y_train, y_test = train_test_split(X, y_binary,
                                                    test_size=0.2,
                                                    random_state=42)

# Сконфигурировать и натренировать древесного ученика
classifier = DecisionTreeClassifier(criterion='gini',
                                   max_depth=5,
                                   random_state=42)

classifier.fit(X=X_train, y=y_train)
```

```
# Выход:
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=5,
                       max_features=None, max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, presort=False,
                       random_state=42, splitter='best')
```

Выход после тренировки модели показывает все параметры классификатора `DecisionTreeClassifier`, которые мы рассмотрим подробнее в следующем разделе при обсуждении темы параметрической настройки.

Как визуализировать дерево решений

Дерево решений можно визуализировать с помощью библиотеки `graphviz` (инструкции по ее установке см. в репозитории [GitHub](#)), т. к. библиотека `sklearn` может выводить описание дерева с использованием языка описания графов DOT (файлы с расширением `dot`), применяемого этой библиотекой. Для удобства восприятия диаграммы выход можно сконфигурировать, включив метки признаков и классов и лимитировав число уровней. Это делается следующим образом:

```
dot_data = export_graphviz(classifier,
                           out_file=None, # можно конвертировать в PNG
                                       # и сохранить в файл
                           feature_names=X.columns,
                           class_names=['Down', 'Up'],
                           max_depth=3,
                           filled=True,
                           rounded=True,
                           special_characters=True)

graphviz.Source(dot_data)
```

Результат демонстрирует, что модель использует целый ряд разных признаков и показывает правила разбивки как для непрерывных, так и для категориальных (фиктивных) переменных. На диаграмме под меткой `value` показано число образцов из каждого класса, а под меткой `class` — наиболее распространенный класс (в течение выбранного периода было больше месяцев `Up`, т. е. когда цена росла) (рис. 10.4).

Как оценивать предсказания дерева решений

Для оценивания предсказательной точности нашего первого классификационного дерева мы будем использовать тестовый набор. С его помощью мы сгенерируем вероятности предсказанных классов следующим образом:

```
y_score = classifier.predict_proba(X=X_test)[: , 1] # оставлять вероятности
                                                # только для положительных классов
```

Метод `predict_proba()` порождает одну вероятность для каждого класса. В бинарном классе эти вероятности являются взаимодополняющими и в сумме составляют 1,

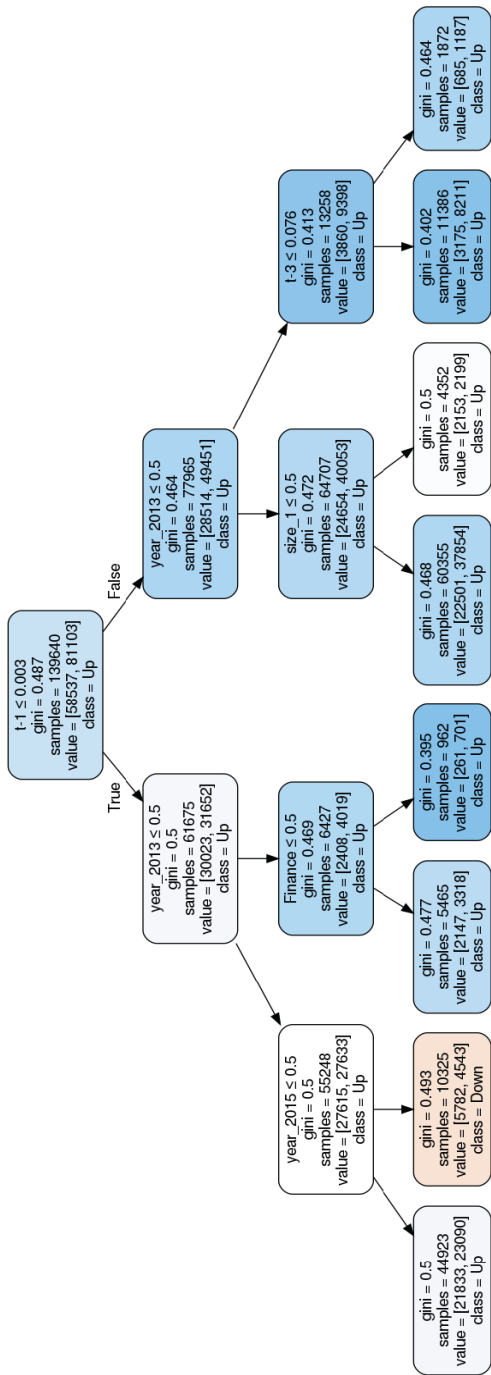


Рис. 10.4. Визуализация дерева решений

поэтому нам требуется только значение положительного класса. Для того чтобы оценить ошибку обобщения, мы будем использовать показатель площади под кривой AUC, основанный на рабочей характеристике приемника (ROC), которую мы представили в *главе 6*. Результат указывает на значительное улучшение за пределами базового значения 0,5, которое соответствует случайному угадыванию:

```
roc_auc_score(y_score=y_score, y_true=y_test)
```

0.6140

Свойство важности признаков

Деревья решений можно не только визуализировать с целью обследования пути принятия решения по данному признаку, но и предоставлять сводную меру вклада каждого признака в подгонку модели к тренировочным данным.

Свойство важности признака отражает то, насколько разбивки, порожденные признаком, помогли оптимизировать модельную метрику, используемую для оценивания качества разбивки, которая в нашем случае является индексом смещанности Джини. Важность признака вычисляется как (нормализованное) суммарное снижение этой метрики и учитывает число образцов, на которые влияет разбивка. Отсюда, признаки, использовавшиеся в дереве ранее, где узлы, как правило, содержат образцов больше, обычно считаются более важными.

На рис. 10.5 показана важность 15 верхних признаков.

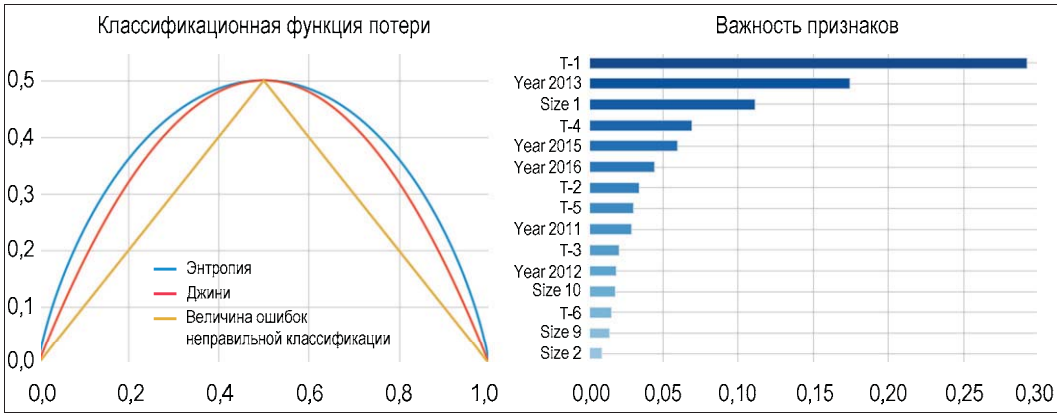


Рис. 10.5. Важность 15 верхних признаков дерева решений

Переподргонка и регуляризация

Деревья решений имеют сильную тягу к переподргонке, в особенности когда набор данных имеет большое число признаков относительно числа образцов. Как обсуждалось в предыдущих главах, переподргонка увеличивает ошибку предсказания,

потому что модель не только усваивает сигнал, содержащийся в тренировочных данных, но и шум.

Существует несколько способов устранения риска перепогонки.

- ◆ *Снижение размерности* (см. главу 12) улучшает соотношение числа признаков к числу образцов, представляя существующие признаки меньшим числом более информативных и менее шумных признаков.
- ◆ *Ансамблевые модели*, такие как случайные леса, объединяют в себе несколько деревьев, рандомизируя при этом конструирование дерева, как мы увидим во второй части этой главы.
- ◆ Деревья решений предоставляют несколько *регуляризационных* гиперпараметров, которые позволяют лимитировать рост дерева и ассоциированную с ним сложность. Хотя каждая разбивка увеличивает число узлов, она также снижает число образцов, располагаемых в расчете на узел для поддержки предсказания. Для каждого дополнительного уровня требуется удвоенное число образцов для заполнения новых узлов с одинаковой плотностью образцов.
- ◆ *Обрезка дерева* — это дополнительный инструмент снижения сложности дерева путем устранения узлов или целых частей дерева, которые создают лишь незначительное добавленное значение, но увеличивают дисперсию модели. Например, обрезка с оптимизацией совместной метрики "стоимость — сложность" начинается с большого дерева и рекурсивно снижает его размер, заменяя узлы листьями, по существу, выполняя построение дерева в обратном порядке. Различные шаги создают последовательность деревьев, которые затем можно сравнить с помощью перекрестного контроля, отобрав идеальный размер.

Как регуляризовать дерево решений

В табл. 10.1 перечислены ключевые параметры, имеющиеся в библиотеке `sklearn` для этой цели в ее реализации дерева решений. После знакомства с наиболее важными параметрами мы покажем, как применять перекрестный контроль для оптимизации гиперпараметрических настроек в отношении компромисса между смещением и дисперсией и более низких ошибок предсказания.

Параметр `max_depth` накладывает жесткие ограничения на число последовательных разбивок и является наиболее простым способом ограничения роста дерева.

Параметры `min_samples_split` и `min_samples_leaf` являются альтернативными способами ограничения роста дерева на основе данных. Вместо того чтобы устанавливать жесткие лимиты на число последовательных разбивок, эти параметры управляют минимальным числом образцов, необходимых для дальнейшей разбивки данных. Последний гарантирует определенное число образцов в расчете на лист, в то время как первый может создавать очень малые листья, если разбивка приводит к очень неравномерному распределению. Малые значения параметров обеспечивают перепогонку, в то время как большое число может помешать дереву усвоить сигнал

в данных. Значения, принятые по умолчанию, часто довольно низкие, и вам следует задействовать перекрестный контроль для разведывания диапазона потенциальных значений. Вы также можете применить тип `float` для указания процента вместо абсолютного числа.

Документация библиотеки `sklearn` содержит дополнительные сведения об использовании разнообразных параметров для разных случаев использования; см. справочные материалы в репозитории `GitHub`.

Таблица 10.1. Параметры библиотеки `sklearn` для регуляризации дерева решений

| Параметр | По умолчанию | Варианты | Описание |
|---------------------------------------|-------------------|---|--|
| <code>max_depth</code> | <code>None</code> | <code>int</code> | Максимальное число уровней: разбивать узлы до достижения глубины <code>max_depth</code> , либо все листья являются чистыми или содержат меньше образцов, чем <code>min_samples_split</code> |
| <code>max_features</code> | <code>None</code> | <code>None</code> : все признаки; <code>int</code> <code>float</code> : доля <code>auto</code> , <code>sqrt</code> : <code>sqrt(n_features)</code> <code>log2</code> : <code>log2(n_features)</code> | Число признаков, рассматриваемых для разбивки |
| <code>max_leaf_nodes</code> | <code>None</code> | <code>None</code> : неограниченное число листовых узлов; <code>int</code> | Разбивать узлы до тех пор, пока не будет создано это число листьев |
| <code>min_impurity_decrease</code> | <code>0</code> | <code>float</code> | Разбивать узлы, если смешанность уменьшается хотя бы на это значение |
| <code>min_samples_leaf</code> | <code>1</code> | <code>int</code> ; <code>float</code> (в процентах от N) | Минимальное число образцов на листовом узле. Разбивка будет рассмотрена только в том случае, если в каждой левой и правой ветвях есть хотя бы <code>min_samples_leaf</code> тренировочных образцов. Может сглаживать модель, в особенности применительно к регрессии |
| <code>min_samples_split</code> | <code>2</code> | <code>int</code> ; <code>float</code> (в процентах от N) | Минимальное число образцов, необходимых для разбиения внутреннего узла |
| <code>min_weight_fraction_leaf</code> | <code>0</code> | | Минимальная взвешенная доля от общей суммы весов всех образцов, необходимых в листовом узле. Образцы имеют одинаковый вес, если параметр <code>sample_weight</code> в методе <code>fit</code> не указан |

Обрезка дерева решений

Процедура рекурсивного бинарного разбиения, скорее всего, приведет к хорошим предсказаниям на тренировочном наборе, но будет тяготеть к перепогонке к данным и производит слабую результативность обобщения, потому что она приводит к чрезмерно сложным деревьям, что находит свое отражение в большом числе листовых узлов или подразделений признакового пространства. Меньшее число разбинок и листовых узлов в целом приводит к меньшему дереву и нередко к лучшей предсказательной результативности, а также интерпретируемости.

Один из подходов к ограничению числа листовых узлов заключается в том, чтобы избегать дальнейших разбинок, если только они не приводят к значительному улучшению целевой метрики. Недостатком этой стратегии, однако, является то, что иногда разбики, которые приводят к малым улучшениям, обеспечивают более ценные разбики позже, поскольку состав образцов продолжает меняться.

В отличие от этого обрезка деревьев начинается с выращивания очень большого дерева перед удалением или обрезкой узлов, сводя большое дерево к менее сложному и перепогонянному поддереву. Обрезка дерева с оптимизацией совместной метрики "стоимости — сложности" генерирует последовательность поддеревьев, добавляя в древесную модель штраф за дополнительные листовые узлы и регуляризационный параметр, аналогичный лассо- и гребневой линейным регрессионным моделям, который модулирует влияние штрафа. Для крупного дерева увеличивающийся штраф будет автоматически производить последовательность поддеревьев. Перекрестный контроль в отношении регуляризационного параметра может использоваться для выявления оптимального обрезанного поддерева.

Этот метод в библиотеке `sklearn` пока недоступен; дополнительные сведения и способы ручной реализации обрезки см. в справочных материалах в репозитории GitHub.

Как настраивать гиперпараметры

Деревья решений предлагают целый ряд гиперпараметров, служащих для управления и настройки результата тренировки. Перекрестный контроль является наиболее важным инструментом получения несмещенной оценки ошибки обобщения, что, в свою очередь, позволяет делать обоснованный выбор между разнообразными вариантами конфигурации. Библиотека `sklearn` предлагает несколько инструментов, которые обеспечивают процесс перекрестного контроля многочисленных конфигураций параметров, а именно вспомогательный класс `GridSearchCV`, который мы проиллюстрируем в следующем далее разделе. Диагностика также обеспечивается за счет кривых усвоения, которые оценивают потенциальные преимущества сбора дополнительных данных с целью снижения ошибки обобщения.

Класс *GridSearchCV* для деревьев решений

Библиотека `sklearn` предоставляет метод определения интервалов значений для многочисленных гиперпараметров. Он автоматизирует процесс перекрестного контроля различных комбинаций значений этих параметров с целью определения оптимальной конфигурации. Давайте пройдемся по процессу автоматической настройки вашей модели.

Первым шагом является создание экземпляра модельного объекта и определение словаря, в котором ключевые слова обозначают гиперпараметры, а значения перечисляют подлежащие проверке значения этих параметров:

```
clf = DecisionTreeClassifier(random_state=42)
param_grid = {'max_depth': range(10, 20),
              'min_samples_leaf': [250, 500, 750],
              'max_features': ['sqrt', 'auto']}
}
```

Затем мы инстанцируем объект `GridSearchCV`, предоставляя для инициализирующего метода объект-оценщик и решетку параметров `param_grid`, а также оценивающий метод `scoring`, который выставляет отметки точности предсказаний, и вариант перекрестного контроля `cv`. Для параметра `cv` мы воспользуемся нашим собственным классом `OneStepTimeSeriesSplit`, инициализированным для использования десяти блоков, и назначим параметру `scoring` значение метрики `roc_auc`. Поиск можно параллелизовать с помощью параметра `n_jobs` и автоматически получить натренированную модель с оптимальными гиперпараметрами, установив `refit=True`.

Когда все настройки будут на своем месте, мы можем выполнить подгонку объекта `GridSearchCV` так же, как и любую другую модель:

```
gridsearch_clf = GridSearchCV(estimator=clf,
                              param_grid=param_grid,
                              scoring='roc_auc',
                              n_jobs=-1,
                              cv=cv, # пользовательский класс
                                  # OneStepTimeSeriesSplit
                              refit=True,
                              return_train_score=True)
```

```
gridsearch_clf.fit(X=X, y=y_binary)
```

Процесс тренировки производит для нашего объекта `GridSearchCV` несколько новых атрибутов, и самое главное, информацию об оптимальных настройках и наилучшей перекрестно-контрольной отметке (теперь используя надлежащую конфигурацию, которая позволяет избежать систематического смещения из-за забегания вперед).

Установка параметра `max_depth` равным 13, `min_samples_leaf` равным 500 и случайную выборку только тех чисел, которые при принятии решения о разбивке соответству-

ют квадратному корню из общего числа признаков, дает наилучшие результаты со значением AUC, равным 0,5855:

```
gridsearch_clf.best_params_
{'max_depth': 13, 'max_features': 'sqrt', 'min_samples_leaf': 500}

gridsearch_clf.best_score_
0.5855
```

Указанная автоматизация довольно удобна, но мы также хотели бы проинспектировать эволюцию результативности для разных значений параметров. По завершении этого процесса объект GridSearchCV обеспечивает наличие подробных результатов перекрестного контроля, что позволяет получать более четкое представление.

Как обследовать древесную структуру

В блокноте также показано, как выполнять перекрестный контроль вручную с целью получения пользовательских атрибутов дерева, таких как общее число узлов либо листовых узлов, ассоциированных с определенными настройками гиперпараметров. Следующая ниже функция обращается к внутреннему атрибуту `tree_`, получая информацию об общем числе узлов и о том, сколько из этих узлов являются листовыми узлами:

```
def get_leaves_count(tree):
    t = tree.tree_
    n = t.node_count
    leaves = len([i for i in range(t.node_count) if t.children_left[i]==-1])
    return leaves
```

Мы можем совместить эту информацию с тренировочной и тестовой отметками, получив подробные сведения о поведении модели на протяжении всего процесса перекрестного контроля, как показано ниже:

```
train_scores, val_scores, leaves = {}, {}, {}
for max_depth in range(1, 26):
    print(max_depth, end=' ', flush=True)
    clf = DecisionTreeClassifier(criterion='gini',
                                max_depth=max_depth,
                                min_samples_leaf=500,
                                max_features='auto',
                                random_state=42)
    train_scores[max_depth], val_scores[max_depth], leaves[max_depth] = [], [], []

    for train_idx, test_idx in cv.split(X):
        X_train, y_train, = X.iloc[train_idx], y_binary.iloc[train_idx]
        X_test, y_test = X.iloc[test_idx], y_binary.iloc[test_idx]
        clf.fit(X=X_train, y=y_train)
```

```
train_pred = clf.predict_proba(X=X_train)[: , 1]
train_score = roc_auc_score(y_score=train_pred, y_true=y_train)
train_scores[max_depth].append(train_score)

test_pred = clf.predict_proba(X=X_test)[: , 1]
val_score = roc_auc_score(y_score=test_pred, y_true=y_test)
val_scores[max_depth].append(val_score)
leaves[max_depth].append(get_leaves_count(clf))
```

Результат показан на левой панели рис. 10.6. Она выделяет внутри- и вневыборочную результативность на интервале значений параметра `max_depth`, а также показывает доверительный интервал вокруг метрик ошибки. Диаграмма также демонстрирует число листовых узлов на правой логарифмической шкале и указывает на конфигурацию с наилучшей результативностью при 13 последовательных разбивках, как обозначено штриховой вертикальной линией.

Кривые усвоения

Кривая усвоения является полезным инструментом, который иллюстрирует, как контрольная и тренировочная отметки эволюционируют по мере увеличения числа тренировочных образцов.

Задача кривой усвоения состоит в том, чтобы выяснить, принесет ли модель пользу, и в какой степени она выиграет от использования большего количества данных в ходе тренировки. Она также полезна в диагностировании возможности того, что модельная ошибка обобщения обусловлена смещением либо дисперсией.

Если, например, и контрольная отметка, и тренировочная отметка сходятся к аналогично низкому значению, несмотря на увеличение размера тренировочного набора, то ошибка, скорее всего, обусловлена смещением, и дополнительные тренировочные данные вряд ли помогут.

Давайте взглянем на визуализацию на рис. 10.6.

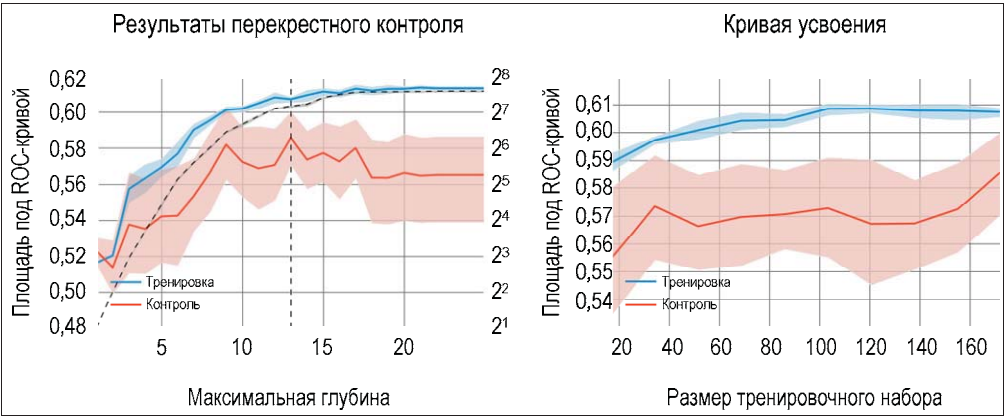


Рис. 10.6. Кривые усвоения для дерева решений

Сильные и слабые стороны деревьев решений

В отличие от линейных моделей, которые мы исследовали до сих пор, регрессионные и классификационные деревья используют совершенно иной подход к предсказанию. Как решить, какая модель больше подходит для определенной задачи? Рассмотрим следующее.

- ◆ Если связь между результатом и признаками является приближенно линейной (или может быть соответствующим образом преобразована), то результативность линейной регрессии, скорее всего, превзойдет более сложный метод, такой как дерево решений, которое не задействует эту линейную структуру.
- ◆ Если связь выглядит весьма нелинейной и более сложной, то результативность деревьев решений, скорее всего, превзойдет классические модели.

Несколько преимуществ сделали деревья решений очень популярными.

- ◆ Их довольно просто понять и интерпретировать, не в последнюю очередь потому, что их можно легко визуализировать и, следовательно, они более доступны для нетехнической аудитории. Деревья решений также называются бело-ящичными моделями, учитывая высокую степень прозрачности того, как они приходят к предсказанию. Черно-ящичные модели, такие как ансамбли и нейронные сети, могут обеспечить более высокую точность предсказания, но их логику принятия решения часто гораздо сложнее понять и интерпретировать.
- ◆ Деревья решений требуют меньше подготовки данных, чем модели, которые принимают более сильные допущения о данных или которые более чувствительны к выбросам и требуют стандартизации данных (например, регуляризованная регрессия).
- ◆ Некоторые реализации деревьев решений работают с категориальным входом, не требуют создания фиктивных переменных (повышая эффективность потребления памяти) и могут работать с пропущенными значениями, как мы увидим в *главе 11*, но это не относится к библиотеке `sklearn`.
- ◆ Предсказание делается быстро, потому что оно логарифмично числу листовых узлов (если только дерево не становится крайне несбалансированным).
- ◆ Модель может быть проверена с помощью статистических тестов, которые позволяют получить подтверждение ее надежности (см. справочные материалы в репозитории GitHub).

Деревья решений также имеют несколько ключевых недостатков.

- ◆ Деревьям решений присуща тенденция к переподгонке к тренировочному набору, и они производят высокую ошибку обобщения. Ключевыми шагами в устранении этого недостатка является обрезка (пока не поддерживается в библиотеке `sklearn`), а также регуляризация с использованием различных критериев досрочной остановки, изложенных в предыдущем разделе.
- ◆ С этим тесно связана высокая дисперсия деревьев решений, которая является результатом их способности тесно адаптироваться к тренировочному набору,

вследствие чего незначительные вариации в данных могут порождать широкие колебания в структуре деревьев решений и, следовательно, прогнозах, которые модель генерирует. Ключевым механизмом решения проблемы высокой дисперсии деревьев решений является использование ансамбля рандомизированных деревьев решений с низким смещением и некоррелированными ошибками предсказания.

- ◆ Жадный подход к автоматическому обучению дерева решений оптимизируется на основе локальных критериев, т. е. снижает ошибку предсказания на текущем узле и не гарантирует глобально оптимальный результат. Опять же, ансамбли, состоящие из рандомизированных деревьев, помогают смягчить эту проблему.
- ◆ Деревья решений также чувствительны к несбалансированным классовым весам и могут создавать смещенные деревья. Одним из вариантов является выборка с повышенной частотой недопредставленного класса или выборка с пониженной частотой более частого класса. Однако, как правило, лучше использовать класс-овые веса и напрямую настраивать целевую функцию.

Случайные леса

Деревья решений не только полезны из-за своей прозрачности и интерпретируемости, но и являются фундаментальными строительными блоками для гораздо более мощных ансамблевых моделей, которые совмещают в себе многочисленные индивидуальные деревья со стратегиями случайного варьирования их конструкции, решая проблемы переобучения и высокой дисперсии, рассмотренные в предыдущем разделе.

Ансамблевые модели

Ансамблевое обучение предусматривает совмещение нескольких автоматически обучающихся моделей в новую единую модель, которая призвана давать более качественные предсказания, чем любая индивидуальная модель. Более конкретно, ансамбль интегрирует предсказания нескольких базовых оценщиков, натренированных с использованием одного или нескольких заданных обучающихся алгоритмов, снижая ошибки обобщения, которые эти модели производят самостоятельно.

Для того чтобы ансамблевое обучение достигло этой цели, индивидуальные модели должны быть:

- ◆ *точными* — они показывают более высокую производительность по сравнению с наивным исходным уровнем (например, средним значением выборки или пропорциями классов);
- ◆ *независимыми* — их предсказания генерируются по-разному, производя разные ошибки.

Ансамблевые методы являются одними из наиболее успешных автоматически обучающихся алгоритмов, в частности для стандартных численных данных. Крупные

ансамбли очень успешны в соревнованиях по машинному обучению и могут состоять из многочисленных индивидуальных моделей, которые были объединены вручную или с использованием другого автоматически обучающегося алгоритма.

В совмещении предсказаний, делаемых разными моделями, существует несколько недостатков. К ним относятся снижение интерпретируемости, повышение сложности и стоимости тренировки, предсказания и сопровождения моделей. В результате на практике (вне соревнований) небольшие выигрыши в точности от крупномасштабного ансамблирования могут не стоить добавочных затрат.

Существует две группы ансамблевых методов, которые обычно различаются в зависимости от того, как они оптимизируют составляющие их модели, а затем интегрируют результаты для единого ансамблевого предсказания.

- ◆ *Методы пакетного усреднения* (бэггинга) тренируют несколько базовых оценщиков независимо и параллельно, а затем усредняют их предсказания. Если базовые модели не смещены и допускают разные ошибки предсказания, не имеющие высокой корреляции, то комбинированное предсказание может иметь меньшую дисперсию и быть более надежным. Это напоминает конструирование портфеля из активов с некоррелированными возвратами с целью снижения волатильности без ущерба для возвратности.
- ◆ *Методы ансамблевого усиления* (бустинга), напротив, тренируют базовых оценщиков последовательно с конкретной целью снизить смещение комбинированного оценщика. Мотив состоит в том, чтобы объединить несколько слабых моделей в мощный ансамбль, последовательно усиливая ансамбль с каждым циклом.

В оставшейся части этой главы мы сосредоточимся на методах автоматического пакетного усреднения (бэггинга), а в *главе 11* — на методах ансамблевого усиления (бустинга).

Как бэггинг снижает дисперсию модели

Мы убедились, что деревья решений, вероятно, будут делать слабые предсказания из-за высокой дисперсии. Из этого следует, что древесная структура довольно чувствительна к составу тренировочной выборки. Мы также убедились, что модель с низкой дисперсией, такая как линейная регрессия, производит аналогичные оценки, несмотря на разные тренировочные выборки при условии, что предоставлено достаточное число образцов при заданном числе признаков.

Для заданного множества независимых наблюдений, каждое из которых имеет дисперсию σ^2 , стандартная ошибка среднего значения выборки задается как σ/\sqrt{n} . Другими словами, усреднение над крупным множеством наблюдений снижает дисперсию. Таким образом, естественным способом снижения дисперсии модели и ее ошибки обобщения было бы собрать большое число тренировочных наборов из популяции, натренировать отличающуюся модель на каждом наборе данных и усреднить результирующие предсказания.

На практике мы, как правило, не можем позволить себе роскошь иметь большое число разных тренировочных наборов. Именно здесь вступает в игру агрегирование бутстраповских выборок, или *бэггинг* (от англ. *bootstrap aggregation*²). Бэггинг — это универсальный метод, снижающий дисперсию автоматически обучающейся модели, который особенно полезен и популярен при применении к деревьям решений.

Под бэггингом имеется в виду агрегирование бутстраповских выборок, т. е. случайных выборок с возмещением³. Такая случайная выборка имеет такое же число наблюдений, что и исходный набор данных, но может содержать дубликаты из-за возврата взятого образца назад в исходный набор данных.

Бэггинг повышает точность предсказания, но уменьшает интерпретируемость модели, потому что дерево больше невозможно визуализировать, чтобы понять важность каждого признака. В качестве ансамблевого алгоритма методы бэггинга (пакетного усреднения) тренируют заданное число базовых оценщиков на этих бутстрапированных выборках (пакетах), а затем агрегируют их предсказания в окончательное ансамблевое предсказание.

Бэггинг снижает дисперсию базовых оценщиков путем рандомизации, например, каждое дерево выращивается, а затем их предсказания усредняются, снижая их ошибки обобщения. Часто этот простой подход позволяет доводить модель до нужной кондиции без необходимости изменения базового алгоритма. Он лучше всего работает со сложными моделями с низким смещением и высокой дисперсией, такими как глубокие деревья решений, потому что его цель — ограничить переподгонку. Методы бустинга (ансамблевого усиления), напротив, лучше всего работают со слабыми моделями, такими как мелкие деревья решений.

Существует несколько методов бэггинга, которые отличаются процессом случайной выборки, которую они применяют к тренировочному набору:

- ♦ метод вставки (*pasting*) берет случайные образцы из тренировочных данных без возмещения, тогда как бэггинг отбирает их с возмещением (возвратом назад в тренировочный набор);

² Бэггинг (*bagging*, *bootstrap aggregating*), или агрегирование бутстраповских выборок, пакетное усреднение, — это общая методика формирования набора моделей путем взятия бутстраповских выборок из данных и их усреднения. В английском языке термин "бэггинг" имеет двойное значение и, с одной стороны, является аббревиатурой для бутстрап-агрегирования, а с другой, исходя из слова *bagging* (пакетирование), означает усреднение бутстраповских пакетов данных, формируя упакованный (*bagged*, бутстрап-агрегированный) предсказатель. Отсюда возникли термины внутripакетный (*in-bag*) и внепакетный (*out-of-bag*). Используемые для построения предсказателя выборки называются внутripакетными, а выборки, которые для этого не использовались, — внепакетными. — *Прим. перев.*

³ Бутстраповская выборка (*bootstrap sample*), или бутстрап, — это подвыборка данных, взятая с возвратом их назад в исходную совокупность наблюдаемых данных. Процесс бутстрапирования можно концептуально представить, как многократное взятие образцов с возмещением с целью получения синтетической выборки. — *Прим. перев.*

- ♦ метод случайных подпространств (random subspace) произвольно отбирает образцы из признаков (т. е. столбцов) без возмещения;
- ♦ метод случайных заплаток (random patch) тренирует базовых оценщиков путем случайной выборки как наблюдений, так и признаков.

Бутстрап-агрегированные деревья решений

Для того чтобы применить бэггинг к деревьям решений, из тренировочных данных мы создаем бутстраповские выборки, многократного отбирая образцы с возмещением, затем тренируем одно дерево решений на каждой из этих выборок и создаем ансамблевое предсказание, усредняя предсказания разных деревьев.

Бутстрап-агрегированные деревья решений обычно вырастают большими, т. е. имеют много уровней и листовых узлов и не обрезаются, вследствие чего каждое дерево имеет низкое смещение, но высокую дисперсию. И эффект усреднения их предсказаний призван снизить их дисперсии. Было показано, что бэггинг существенно улучшает предсказательную результативность, если конструировать ансамбли, которые объединяют в себе сотни или даже тысячи деревьев, натренированных на бутстраповских выборках.

Для того чтобы проиллюстрировать эффект бэггинга на дисперсию регрессионного дерева, мы можем применить метаоценщика `BaggingRegressor`, предоставляемого библиотекой `sklearn`. Он тренирует заданного пользователем базового оценщика на основе параметров, определяющих стратегию отбора:

- ♦ параметры `max_samples` и `max_features` управляют размером подмножеств, извлекаемых соответственно из строк и столбцов;
- ♦ параметры `bootstrap` и `bootstrap_features` определяют, будет ли извлекаться каждая из этих выборок с возмещением или без него.

В следующем далее примере используется экспоненциальная функция, которая генерирует тренировочные выборки для одного-единственного регрессора дерева решений `DecisionTreeRegressor` и ансамбля `BaggingRegressor`, состоящего из 10 деревьев, каждое из которых выращивается на 10 уровней в глубину. Обе модели тренируются на случайных выборках и предсказывают результаты для фактической функции с добавлением шума.

Поскольку мы знаем истинную функцию, мы можем разложить среднеквадратическую ошибку на смещение, дисперсию и шум и сравнить относительный размер этих компонентов для обеих моделей в соответствии со следующей разбивкой:

$$E[y_0 - \hat{f}(x_0)]^2 = \text{дисперсия}(\hat{f}(x_0)) + [\text{смещение}(\hat{f}(x_0))]^2 + \text{дисперсия}(\varepsilon).$$

Для 100 повторных случайных тренировочных и тестовых выборок по 250 и 500 наблюдений каждая мы находим, что дисперсия предсказаний индивидуального дерева решений почти в два раза выше, чем для небольшого ансамбля из 10 агрегированных деревьев на основе бутстраповских выборок:

```

test_size = 500
train_size = 250
reps = 100

noise = .5    # шум относительно std(y)
noise = y.std() * noise_to_signal

X_test = choice(x, size=test_size, replace=False)

max_depth = 10
n_estimators=10

tree = DecisionTreeRegressor(max_depth=max_depth)
bagged_tree = BaggingRegressor(base_estimator=tree,
                               n_estimators=n_estimators)

learners = {'Дерево решений': tree, 'Пакетноусредненный регрессор': bagged_tree}

predictions = {k: pd.DataFrame() for k, v in learners.items()}
for i in range(reps):
    X_train = choice(x, train_size)
    y_train = f(X_train) + normal(scale=noise, size=train_size)
    for label, learner in learners.items():
        learner.fit(X=X_train.reshape(-1, 1), y=y_train)
    preds = pd.DataFrame({i: learner.predict(X_test.reshape(-1, 1))}, index=X_test)
    predictions[label] = pd.concat([predictions[label], preds], axis=1)

```

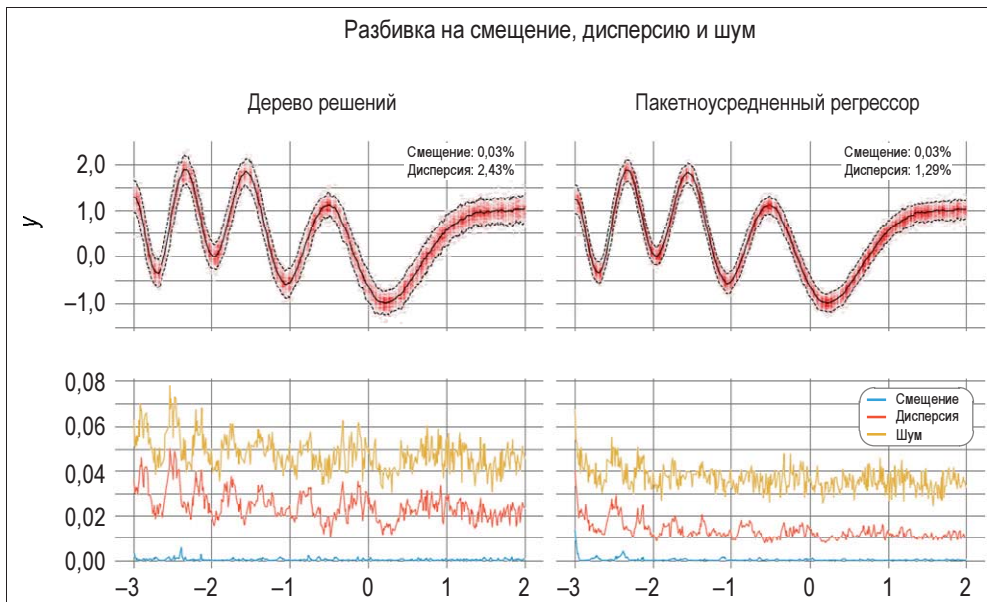


Рис. 10.7. Разбивка на смещение, дисперсию и шум, а также среднее предсказание с полосами стандартного отклонения

Для каждой модели на рис. 10.7 на верхней панели показано среднее предсказание и полоса двух стандартных отклонений вокруг среднего для обеих моделей, а на нижней панели — разбивка на смещение, дисперсию и шум на основе значений истинной функции.

Подробности реализации см. в блокноте `random_forest.ipynb`.

Как строить случайный лес

Алгоритм случайного леса опирается на рандомизацию, вносимую бутстраповскими выборками, генерируемыми бэггингом с целью еще больше снизить дисперсию и улучшить предсказательную результативность.

В дополнение к тренировке каждого члена ансамбля на бутстрапированных тренировочных данных, случайные леса также производят случайную выборку из признаков, используемых в модели (без возмещения). В зависимости от реализации, случайные выборки могут извлекаться для каждого дерева или каждой разбивки. В результате алгоритм сталкивается с разными вариантами при усвоении новых правил, на уровне дерева либо для каждой разбивки.

Размеры выборок из признаков для регрессионных и классификационных деревьев различаются:

- ◆ в случае *классификации* размер выборки, как правило, равен квадратному корню из числа признаков;
- ◆ в случае *регрессии* он может составлять от одной трети до всех признаков и должен быть выбран на основе перекрестного контроля.

На рис. 10.8 показано, как случайные леса рандомизируют тренировку индивидуальных деревьев, а затем агрегируют их предсказания в ансамблевое предсказание.

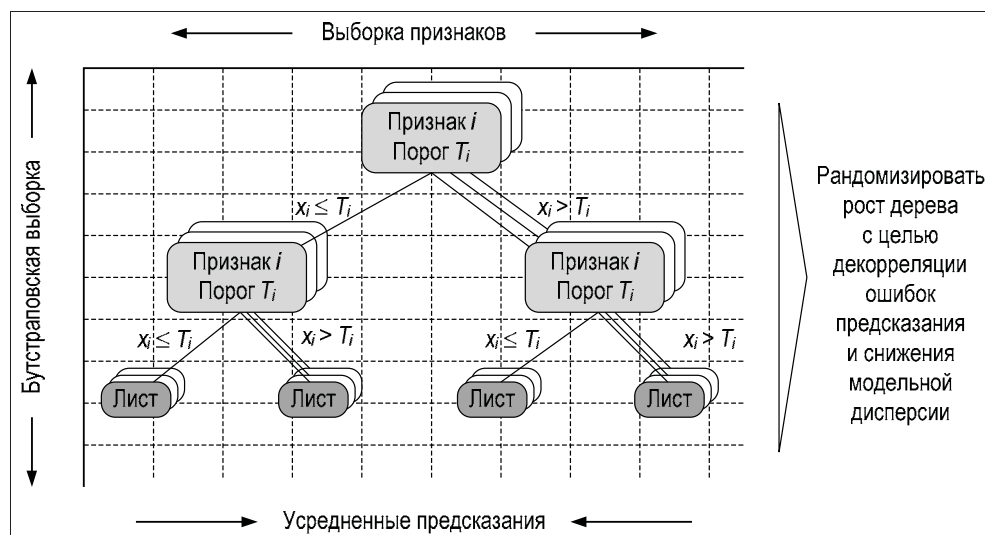


Рис. 10.8. Рандомизирование тренировки и агрегирование предсказаний

Целью рандомизации признаков в дополнение к рандомизации тренировочных наблюдений является дальнейшая декорреляция ошибок предсказания у индивидуальных деревьев. Не все признаки создаются равными, и в процессе конструирования дерева небольшое число высоко релевантных признаков будет выбираться гораздо чаще и раньше, делая деревья решений более похожими по всему ансамблю. Однако чем меньше коррелируют ошибки обобщения в индивидуальных деревьях, тем больше снижается совокупная дисперсия.

Как тренировать и настраивать случайный лес

Ключевые параметры конфигурации охватывают различные гиперпараметры для индивидуальных деревьев решений, представленных в *разд. "Как настраивать гиперпараметры"* ранее в этой главе. В табл. 10.2 перечислены дополнительные параметры для двух классов `RandomForest`.

Таблица 10.2. Дополнительные гиперпараметры

| Ключевое слово | По умолчанию | Описание |
|---------------------------|--------------------|--|
| <code>bootstrap</code> | <code>True</code> | Бутстраповские выборки во время тренировки |
| <code>n_estimators</code> | <code>10</code> | Число деревьев в лесу |
| <code>oob_score</code> | <code>False</code> | Использует внепакетные выборки для оценивания R^2 на не встречавшихся ранее данных |

Параметр `bootstrap` активируется в предыдущем схематичном описании алгоритма бэггинга, который, в свою очередь, позволяет вычислять внепакетную отметку (`oob_score` от англ. *out-of-bag score*), оценивающую точность обобщения с использованием выборок, не включенных в бутстраповскую выборку, используемую для тренировки определенного дерева (подробности см. в следующем разделе).

Параметр `n_estimators` задает число деревьев, выращиваемых в составе леса. Более крупные леса работают лучше, но также требуют больше времени для их построения. Очень важно отслеживать перекрестно-контрольную ошибку как функцию числа базовых учеников, для выявления момента, когда предельное снижение ошибки предсказания начнет ухудшаться, а стоимость дополнительной тренировки начнет перевешивать преимущества.

Параметр `max_features` управляет размером случайно отобранных подмножеств признаков, располагаемых при усвоении нового правила принятия решения и разбивки узла. Меньшее значение снижает корреляцию деревьев и, следовательно, дисперсию ансамбля, но также может увеличивать смещение. Хорошими стартовыми значениями являются `n_features` (число тренировочных признаков) для регрессионных задач и `sqrt(n_features)` для классификационных задач, но они будут зависеть от взаимосвязей между признаками и должны быть оптимизированы с помощью перекрестного контроля.

Случайные леса предназначены содержать глубокие полностью развившиеся деревья, которые могут быть созданы с использованием `max_depth=None` и `min_samples_split=2`. Однако эти значения не обязательно являются оптимальными, в особенности для высокоразмерных данных с многочисленными образцами и, следовательно, потенциально очень глубокими деревьями, которые могут стать очень вычислительно-интенсивными и емкими по потреблению памяти.

Класс `RandomForest`, предоставляемый библиотекой `sklearn`, поддерживает параллельную тренировку и предсказание путем установки параметра `n_jobs` равным числу k заданий для выполнения на разных ядрах. Значение `-1` использует все имеющиеся ядра. Накладные расходы на межпроцессный обмен могут лимитировать линейность ускорения, вследствие чего k заданий могут занимать более $1/k$ времени одного задания. Тем не менее ускорение зачастую является весьма ощутимым для крупных лесов или глубоких индивидуальных деревьев, тренировка которых может занимать значительное количество времени, когда данные являются крупными, и раздельное вычисление становится дорогостоящим.

Как всегда, наилучшая параметрическая конфигурация должна быть определена с помощью перекрестного контроля. Следующие ниже шаги иллюстрируют указанный процесс.

1. Мы воспользуемся классом `GridSearchCV` с целью выявления оптимального набора параметров для ансамбля классификационных деревьев:

```
rf_clf = RandomForestClassifier(n_estimators=10,
                               criterion='gini',
                               max_depth=None,
                               min_samples_split=2,
                               min_samples_leaf=1,
                               min_weight_fraction_leaf=0.0,
                               max_features='auto',
                               max_leaf_nodes=None,
                               min_impurity_decrease=0.0,
                               min_impurity_split=None,
                               bootstrap=True,
                               oob_score=True,
                               n_jobs=-1,
                               random_state=42,
                               verbose=1)
```

2. Далее задействуем 10-блочный перекрестный контроль собственной разработки и заполним решетку параметров значениями для ключевых конфигурационных условий:

```
cv = OneStepTimeSeriesSplit(n_splits=12, test_period_length=1, shuffle=True)
param_grid = {'n_estimators': [200, 400],
              'max_depth': [10, 15, 20],
              'min_samples_leaf': [50, 100]}
```

3. Сконфигурируем объект класса `GridSearchCV` с помощью приведенного выше входа:

```
gridsearch_clf = GridSearchCV(estimator=rf_clf,
                              param_grid=param_grid,
                              scoring='roc_auc',
                              n_jobs=-1,
                              cv=cv,
                              refit=True,
                              return_train_score=True,
                              verbose=1)
```

4. Натренируем несколько ансамблевых моделей, заданных решеткой параметров:

```
gridsearch_clf.fit(X=X, y=y_binary)
```

5. Получим наилучшие параметры, как показано ниже:

```
gridsearch_clf.bestparams

{'max_depth': 15,
 'min_samples_leaf': 100,
 'n_estimators': 400}
```

6. Лучшая оценка точности показывает малое, но значительное улучшение по сравнению с базовым уровнем одного-единственного дерева:

```
gridsearch_clf.bestscore_

0.6013
```

Свойство важности признаков для случайных лесов

Ансамбль случайного леса может содержать сотни индивидуальных деревьев, однако из бутстрап-агрегированных моделей по-прежнему можно получать совокупную сводную меру важности признаков.

Для заданного признака отметка важности определенного признака представляет собой суммарное снижение значения целевой функции, получающееся в результате разбивок на основе этого признака, усредненно по всем деревьям. Поскольку целевая функция учитывает число признаков, на которые влияет одна разбивка, эта мера неявно является средневзвешенной, вследствие чего признаки, используемые вблизи вершины дерева, получают более высокие отметки из-за большего числа наблюдений, содержащихся в гораздо меньшем числе имеющихся узлов. При усреднении по многочисленным деревьям, выращенным рандомизированным способом, отметка важности признаков теряет некоторую дисперсию и становится точнее.

Вычисление отличается для классификационных и регрессионных деревьев ввиду разных целей усвоения правил принятия решения и измеряется в терминах среднеквадратической ошибки для регрессионных деревьев и индекса Джини или энтропии для классификационных деревьев.

Библиотека `sklearn` нормализует меру важности признаков, вследствие чего она в сумме составляет 1. Вычисленная таким образом важность признаков также используется для отбора признаков в качестве альтернативы мерам взаимной информации, которые мы встречали в *главе 6* (см. класс-метатрансформер `SelectFromModel` в модуле `sklearn.feature_selection`).

В нашем примере значения важности верхних 20 признаков показаны на рис. 10.9.

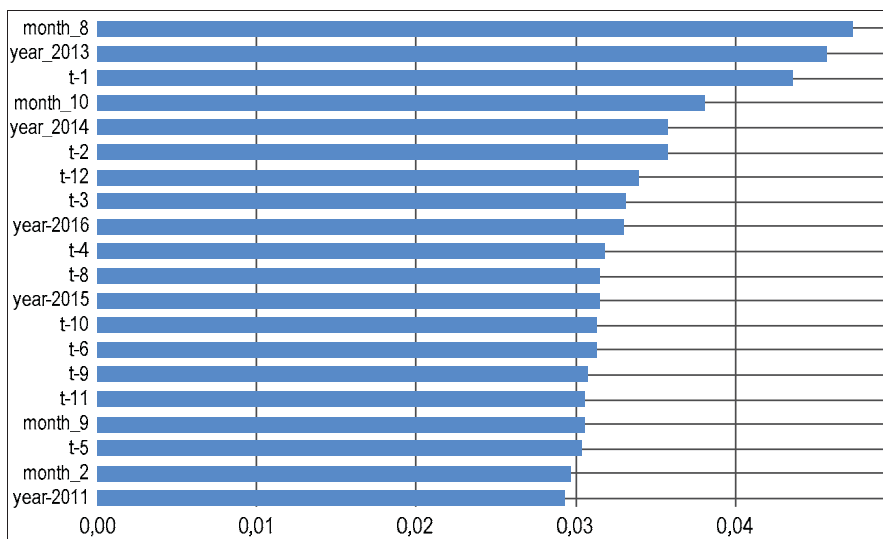


Рис. 10.9. Значения важности верхних 20 признаков

Внепакетное тестирование

Случайные леса предлагают преимущество встроенного перекрестного контроля, поскольку индивидуальные деревья тренируются на бутстрапированных версиях тренировочных данных. В результате каждое дерево использует в среднем только две трети имеющихся наблюдений. Для того чтобы понять причину, следует учесть, что бутстраповская выборка имеет тот же размер, n , что и исходная выборка, и каждое наблюдение имеет одинаковую вероятность $1/n$. Следовательно, вероятность вообще не войти в бутстраповскую выборку равна $(1 - 1/n)^n$, которая (быстро) сходится к $1/e$, или примерно одной трети.

Эта оставшаяся треть наблюдений, которые не включены в тренировочный набор, используемый для выращивания бутстрап-агрегированного (пакетно-усредняемого) дерева, называется *внепакетными наблюдениями* (out-of-bag, OOB) и может служить в качестве контрольного набора. Как и при перекрестном контроле, мы предсказываем отклик для внепакетного образца по каждому дереву, построенному без этого наблюдения, а затем усредняем предсказанные отклики (если целью является

регрессия) либо выносим решение, принимаемое большинством голосов, или берем предсказанную вероятность (если целью является классификация) для одиночного ансамблевого предсказания по каждому внепакетному образцу. Такие предсказания порождают несмещенную оценку ошибки обобщения, удобно вычисляемую во время тренировки.

Результирующая внепакетная (ООВ) ошибка является допустимой оценкой ошибки обобщения для этого наблюдения, поскольку предсказание производится с использованием правил принятия решения, усвоенных в отсутствие этого наблюдения. Как только случайный лес становится достаточно крупным, внепакетная ошибка начинает близко аппроксимировать перекрестно-контрольную ошибку с использованием метода исключения по одному образцу. Внепакетный подход к оцениванию тестовой ошибки очень эффективен для крупных наборов данных, где перекрестный контроль бывает вычислительно дорогостоящим.

Сильные и слабые стороны случайных лесов

Пакетноусредненные ансамблевые модели имеют как преимущества, так и недостатки. Преимущества случайных лесов:

- ◆ предсказательная результативность может конкурировать с лучшими алгоритмами контролируемого обучения;
- ◆ они обеспечивают надежную оценку важности признаков;
- ◆ они предлагают эффективные оценки тестовой ошибки без издержек на многократную тренировку модели, ассоциированную с перекрестным контролем.

С другой стороны, случайные леса также имеют несколько недостатков:

- ◆ ансамблевая модель по своей сути является менее интерпретируемой, чем индивидуальное дерево решений;
- ◆ тренировка большого числа глубоких деревьев может иметь высокую вычислительную стоимость (но может быть параллелизована) и потреблять очень много памяти;
- ◆ предсказания выполняются медленнее, что может создавать проблемы для приложений, требующих низкого значения задержки.

Резюме

В этой главе мы изучили новый класс моделей, способных улавливать нелинейные связи, в отличие от классических линейных моделей, которые мы исследовали до сих пор. Мы увидели, как деревья решений усваивают правила разбиения признакового пространства на участки, которые дают предсказания и, таким образом, сегментируют входные данные по определенным участкам.

Деревья решений очень полезны тем, что они позволяют уникальным образом проникать в сущность связей между признаками и целевыми переменными, и мы уви-

дели, как легко визуализировать последовательность правил принятия решения, закодированных в древесной структуре.

К сожалению, дерево решений тяготеет к переподгонке. Мы узнали, что ансамблевые модели и метод агрегирования бутстраповских выборок позволяют преодолевать некоторые недостатки деревьев решений и делать их полезными в качестве компонентов гораздо более мощных композитных моделей.

В следующей главе мы рассмотрим еще одну ансамблевую модель, которая стала считаться одним из важнейших автоматически обучающихся алгоритмов.

11

Градиентно-бустинговые машины

В предыдущей главе мы познакомились с тем, как случайные леса улучшают предсказания, выполняемые индивидуальными деревьями решений, объединяя их в ансамбль, который снижает высокую дисперсию индивидуальных деревьев. Случайные леса используют агрегирование бутстраповских выборок (пакетное усреднение), или бэггинг (от англ. *bootstrap aggregation*), с целью внесения случайности в процесс выращивания индивидуальных деревьев.

Более конкретно, бэггинг отбирает образцы из данных с их возвратом назад в набор данных, вследствие чего каждое дерево тренируется на отличающемся, но равновеликом случайном подмножестве данных (в котором некоторые наблюдения повторяются). Случайные леса также случайно отбирают подмножество признаков, вследствие чего строки и столбцы данных, используемые для тренировки каждого дерева, являются случайными версиями исходных данных. Затем ансамбль генерирует предсказания путем усреднения над выходами индивидуальных деревьев.

Индивидуальные деревья обычно выращиваются глубоко с целью обеспечения низкого смещения, опираясь на рандомизированный тренировочный процесс, получая разные некоррелированные ошибки предсказания, которые при агрегировании имеют меньшую дисперсию, чем предсказания индивидуальных деревьев. Другими словами, рандомизированная тренировка призвана декоррелировать или диверсифицировать ошибки, делаемые индивидуальными деревьями, вследствие чего ансамбль является гораздо менее восприимчивым к переподгонке, имеет меньшую дисперсию и лучше обобщается на новые данные.

В этой главе мы рассмотрим бустинг (ансамблевое усиление), альтернативный *автоматически обучающийся* алгоритм для ансамблей деревьев решений, который часто дает еще более высокие результаты. Ключевое отличие бустинга состоит в том, что он модифицирует данные, используемые для тренировки каждого дерева, на основе кумулятивных ошибок, делаемых моделью перед добавлением нового дерева. В отличие от случайных лесов, которые тренируют большое число деревьев независимо друг от друга, используя разные версии тренировочного набора, бустинг работает последовательно, используя перевзвешенные версии данных. Совре-

менные реализации бустинга также адаптируют рандомизационные стратегии случайных лесов.

В этой главе мы увидим, как за последние три десятилетия бустинг превратился в один из самых успешных автоматически обучающихся алгоритмов. На момент написания этих строк он стал доминировать в соревнованиях по машинному обучению по структурированным данным (в отличие, к примеру, от высокоразмерных снимков или речи, где связь между входом и выходом является более сложной, и где преуспевает глубокое обучение). Более конкретно, в этой главе мы рассмотрим следующие темы:

- ◆ как работает бустинг и как он соотносится с бэггингом;
- ◆ как бустинг эволюционировал от адаптивного бустинга к градиентному;
- ◆ как использовать и настраивать модели AdaBoost и градиентно-бустинговые модели с помощью библиотеки `sklearn`;
- ◆ как современные реализации градиентно-бустинговых машин (gradient boosting machine, GBM) радикально ускоряют вычисление;
- ◆ как предотвращать переподргонку градиентно-бустинговых моделей;
- ◆ как создавать, настраивать и оценивать градиентно-бустинговые модели на крупных наборах данных с помощью библиотек XGBoost, LightGBM и CatBoost;
- ◆ как интерпретировать и получить углубленные представления о данных с помощью градиентно-бустинговых моделей.

Адаптивный бустинг

Как и бэггинг (пакетное усреднение), бустинг (ансамблевое усиление) — это алгоритм ансамблевого автоматического обучения, который объединяет базовых учеников (как правило, деревья решений) в ансамбль. Бустинг был первоначально разработан для классификационных задач, но также может использоваться для регрессии и был назван одной из самых мощных идей автоматического обучения, внедренных за последние 20 лет (как описано в книге "Elements of Statistical Learning" ("Элементы статистического обучения") Тревором Хастии (Trevor Hastie) и соавт.; справочные материалы и ссылки см. в репозитории GitHub. Как и бэггинг, он является общим методом или метаметодом, который может применяться ко многим моделям статистического обучения.

Мотивом для развития бустинга был поиск метода, который объединял бы результаты многочисленных слабых моделей (предсказатель называется слабым, когда он показывает результативность чуть лучше, чем случайное угадывание) в более мощное, т. е. усиленное (бустированное) совместное предсказание. В общем случае, бустинг заучивает аддитивную гипотезу H_M формы, похожей на линейную регрессию. Однако теперь каждый из $m = 1, \dots, M$ элементов суммирования представляет

собой слабого базового ученика, именуемого h_i , который сам по себе требует тренировки. Следующая формула подытоживает этот подход:

$$H_M(x) = \sum_{m=1}^M \underbrace{h_i(x)}_{\text{слабый ученик}}.$$

Как обсуждалось в предыдущей главе, бэггинг тренирует базовых учеников на отличающихся случайных выборках из тренировочных данных. Бустинг, напротив, работает последовательно путем тренировки базовых учеников на данных, которые неоднократно модифицируются для отражения кумулятивных результатов автоматического обучения. Цель состоит в том, чтобы следующий базовый ученик компенсировал недостатки текущего ансамбля. В этой главе мы увидим, что алгоритмы бустинга отличаются тем, как они определяют недостатки. Ансамбль делает предсказания, используя средневзвешенное значение предсказаний слабых моделей.

Первый алгоритм бустинга, который сопровождался математическим доказательством того, что он усиливает результативность слабых учеников, был разработан Робертом Шапиро (Robert Schapire) и Йоавом Фройндом (Yoav Freund) примерно в 1990 г. В 1997 г. появилось практическое решение классификационных задач в виде алгоритма *адаптивного бустинга* (AdaBoost, от англ. *adaptive boosting*), получившего в 2003 г. премию Гёделя. Примерно через 5 лет этот алгоритм был расширен до произвольных целевых функций, когда Лео Брейман (Leo Breiman, который изобрел случайные леса) связал этот подход с градиентным спуском, а в 1999 г. Джером Фридман (Jerome Friedman) придумал градиентный бустинг. В последние годы появились многочисленные оптимизированные реализации, такие как XGBoost, LightGBM и CatBoost, и прочно утвердились в качестве решения для структурированных данных.

В следующих далее разделах мы кратко представим алгоритм AdaBoost, а затем сосредоточимся на градиентно-бустинговой модели, а также на нескольких современных реализациях этого очень мощного и гибкого алгоритма.

Алгоритм AdaBoost

Алгоритм AdaBoost был первым алгоритмом бустинга, который итеративно адаптировался к прогрессу кумулятивного автоматического обучения во время подгонки дополнительного члена ансамбля. В частности, алгоритм AdaBoost изменял веса на тренировочных данных, тем самым отражая кумулятивные ошибки текущего ансамбля на обучающем наборе перед подгонкой нового слабого ученика. В то время AdaBoost был самым точным классификационным алгоритмом, и на конференции NIPS в 1996 г. Лео Брейман назвал его лучшим серийным классификатором в мире.

Указанный алгоритм оказал очень заметное влияние на МО, поскольку он предоставлял теоретические гарантии результативности. Эти гарантии требуют только достаточного объема данных и слабого ученика, который надежно предсказывает

чуть лучше, чем случайное угадывание. Вследствие этого адаптивного метода, который учится поэтапно, разработка точной автоматически обучающейся модели больше не требовала точной результативности по всему признаковому пространству. Вместо этого, конструкция модели может сосредотачиваться на отыскании слабых учеников, чья результативность лишь слегка превосходит подбрасывание монеты.

Алгоритм AdaBoost является значительным отходом от бэггинга, который для снижения смещения строит ансамбли на очень глубоких деревьях. AdaBoost, напротив, в качестве слабых учеников выращивает мелкие деревья, часто производя превосходящую точность с пнями, т. е. деревьями, образованными одной-единственной разбивкой. Алгоритм начинается с равновзвешенного тренировочного набора, а затем последовательно изменяет распределение выборки. После каждой итерации алгоритм AdaBoost увеличивает веса неправильно классифицированных наблюдений и снижает веса правильно предсказанных образцов, вследствие чего последующие слабые ученики больше сосредотачиваются на особенно сложных случаях. После тренировки новое дерево решений включается в ансамбль с весом, который отражает его вклад в снижение тренировочной ошибки.

Алгоритм AdaBoost для ансамбля базовых учеников $h_m(x)$, $m = 1, \dots, M$, которые предсказывают дискретные классы, $y \in [-1; 1]$, и N тренировочных наблюдений можно резюмировать следующим образом:

1. Инициализировать выборочные веса $w_i = 1/N$ для наблюдений $i = 1, \dots, N$.
2. Для каждого базового классификатора h_m , $m = 1, \dots, M$, выполнить следующие действия:
 - вписать $h_m(x)$ в тренировочные данные, взвешенные на w_i ;
 - вычислить взвешенную частоту появления ошибок базового ученика ε_m на тренировочном наборе;
 - вычислить ансамблевый вес α_m базового ученика как функцию от его частоты появления ошибок, как показано в следующей формуле:

$$\alpha_m = \log \left(\frac{1 - \varepsilon_m}{\varepsilon_m} \right);$$

- обновить веса для неправильно классифицированных образцов в соответствии с $w_i \cdot \exp(\alpha_m)$.
3. Предсказать положительный класс, если взвешенная сумма членов ансамбля является положительной, и в противном случае — отрицательный, как показано в следующей формуле:

$$H(x) = \text{sign} \left(\sum_{m=1}^M \underbrace{\alpha_m h_m(x)}_{\text{взвешенный слабый ученик}} \right).$$

Алгоритм AdaBoost имеет целый ряд практических преимуществ, включая простоту реализации и быстрые вычисления, и его можно сочетать с любым методом выявления слабых учеников. Помимо размера ансамбля, отсутствуют гиперпараметры, которые требовали бы настройки. AdaBoost также полезен для выявления выбросов, потому что образцы, которые получают самые высокие веса, — это как раз те, которые систематически классифицируются неправильно и по своей сути неоднозначны, что также типично для выбросов.

С другой стороны, результативность алгоритма AdaBoost на конкретном наборе данных зависит от способности слабого ученика адекватно улавливать связь между признаками и результатом. Как предполагает теория, бустинг не будет работать хорошо, когда данных недостаточно или когда сложность членов ансамбля не соответствует сложности данных. Он также может быть восприимчив к шуму в данных.

AdaBoost в библиотеке sklearn

В рамках ансамблевого модуля библиотека sklearn предоставляет реализацию классификатора `AdaBoostClassifier`, поддерживающую два или более классов. Примеры исходного кода для этого раздела находятся в блокноте `gbm_baseline.ipynb`, в котором сравнивается результативность различных алгоритмов с фиктивным классификатором, всегда предсказывающим наиболее частый класс.

Сначала нам нужно определить базового оценщика `base_estimator` в качестве шаблона для всех членов ансамбля, а затем настроить сам ансамбль. Мы будем использовать классификатор, принятый по умолчанию, `DecisionTreeClassifier`, с максимальной глубиной `max_depth=1`, т. е. он будет пнем с одной разбивкой. Сложность оценщика `base_estimator` является ключевым регулировочным параметром, поскольку он зависит от природы данных. Как показано в предыдущей главе, изменения в параметре `max_depth` следует сочетать с надлежащими регуляризационными ограничениями, к примеру, используя корректировки в параметре `min_samples_split`, необходимым для разбивки внутреннего узла, `min_samples_split`, как показано в следующем фрагменте кода:

```
base_estimator = DecisionTreeClassifier(criterion='gini',
                                       splitter='best',
                                       max_depth=1,
                                       min_samples_split=2,
                                       min_samples_leaf=20,
                                       min_weight_fraction_leaf=0.0,
                                       max_features=None,
                                       random_state=None,
                                       max_leaf_nodes=None,
                                       min_impurity_decrease=0.0,
                                       min_impurity_split=None,
                                       class_weight=None,
                                       presort=False)
```

На втором шаге мы конструируем ансамбль. Параметр `n_estimators` управляет числом слабых учеников, а параметр темпа усвоения `learning_rate` определяет вклад каждого слабого ученика, как показано в следующем фрагменте кода. По умолчанию слабые ученики — это пни дерева решений:

```
ada_clf = AdaBoostClassifier(base_estimator=base_estimator,
                             n_estimators=200,
                             learning_rate=1.0,
                             algorithm='SAMME.R',
                             random_state=42)
```

Главными регулировочными параметрами, отвечающими за хорошие результаты, являются `n_estimators` и сложность базового оценщика, поскольку глубина дерева определяет степень взаимодействия между признаками.

Мы проведем перекрестный контроль ансамбля AdaBoost с использованием своей собственной 12-блочной скользящей разбивки временного ряда `OneStepTimeSeriesSplit` для предсказания на 1 месяц вперед для предыдущих 12 месяцев в выборке, используя для тренировки все имеющиеся предшествующие данные, как показано в следующем фрагменте кода:

```
cv = OneStepTimeSeriesSplit(n_splits=12, test_period_length=1, shuffle=True)

def run_cv(ada_clf, X=X_dummies, y=y, metrics=metrics, cv=cv, fit_params=None):
    return cross_validate(estimator=ada_clf,
                          X=X,
                          y=y,
                          scoring=list(metrics.keys()),
                          cv=cv,
                          return_train_score=True,
                          n_jobs=-1, # использовать все ядра
                          verbose=1,
                          fit_params=fit_params)
```

Результат показывает взвешенную тестовую точность 0,62, тестовую отметку $AUC = 0,6689$ и отрицательную логарифмическую потерю $-0,6923$, а также тестовый показатель $F1$, равный 0,5905 (рис. 11.1).

Дополнительные сведения об исходном коде для перекрестного контроля и обработки результатов см. в сопутствующем блокноте.

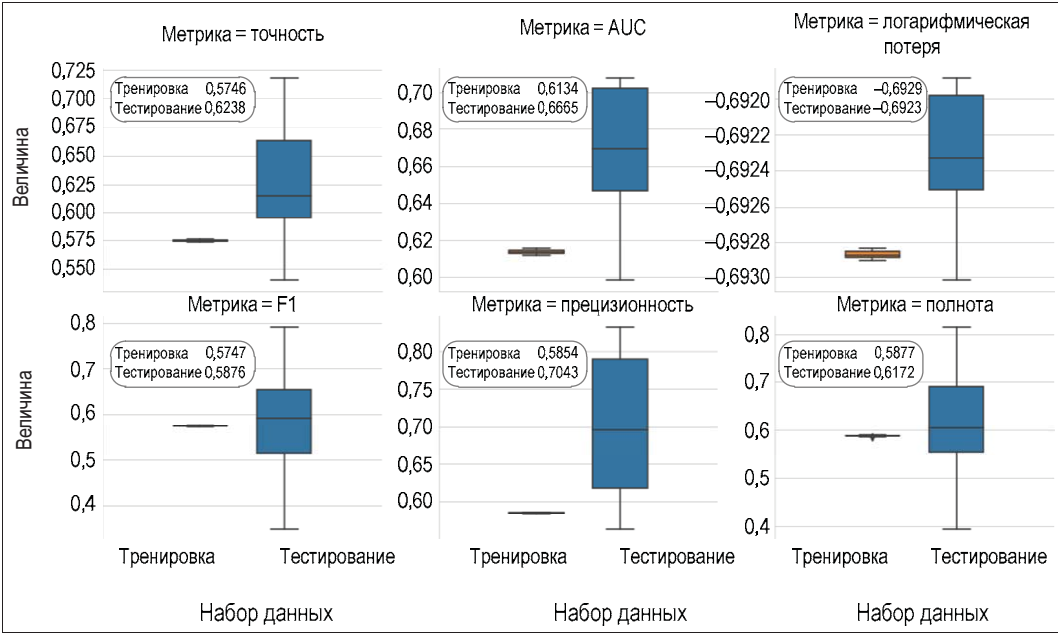


Рис. 11.1. Результаты перекрестного контроля ансамбля AdaBoost

Градиентно-бустинговые машины

Алгоритм AdaBoost также можно интерпретировать как шаговый форвардный подход к минимизации экспоненциальной функции потерь для бинарного $y \in [-1; 1]$ на каждой итерации m с целью выявления нового базового ученика h_m с соответствующим весом α_m , добавляемым в ансамбль, как показано в следующей формуле:

$$\arg \min_{\alpha, h} \sum_{i=1}^N \exp \left(\underbrace{-y_i (f_{m-1}(x_i))}_{\text{текущий ансамбль}} + \underbrace{\alpha_m h_m(x_i)}_{\text{новый член}} \right).$$

Указанная интерпретация алгоритма AdaBoost была обнаружена только через несколько лет после его публикации. Она рассматривает AdaBoost как координатный алгоритм градиентного спуска, который минимизирует конкретную функцию потерь, а именно экспоненциальную потерю.

Градиентный бустинг эффективно задействует это понимание и применяет метод бустинга к гораздо более широкому диапазону функций потерь. Указанный метод позволяет разрабатывать автоматически обучающиеся алгоритмы для решения любой задачи регрессии, классификации или ранжирования при условии, что она может быть сформулирована с использованием дифференцируемой функции потерь и, следовательно, имеет градиент. Гибкость настройки этого общего метода для многих конкретных задач предсказания является ключевой причиной популярности бустинга.

Главная идея в основе результирующего алгоритма *градиентно-бустинговых машин* (Gradient Boosting Machines, GBM) — тренировать базовых учеников усваивать отрицательный градиент текущей ансамблевой функции потери. Вследствие этого каждое дополнение в ансамбль вносит прямой вклад в снижение совокупной тренировочной ошибки с учетом ошибок, сделанных предшествующими членами ансамбля. Поскольку каждый новый член представляет новую функцию данных, также принято говорить, что градиентный бустинг оптимизирует над функциями h_m в аддитивном стиле.

Иными словами, данный алгоритм последовательно вписывает слабых учеников h_m , таких как деревья решений, в отрицательный градиент функции потери, которая оценивается для текущего ансамбля, как показано в следующей формуле:

$$H_m(x) = \underbrace{H_{m-1}(x)}_{\text{текущий ансамбль}} + \underbrace{\gamma_m h_m(x)}_{\text{новый член}} = H_{m-1}(x) + \arg \min \sum_{i=1}^n \underbrace{L(y_i, H_{m-1}(x_i) + h(x))}_{\text{функция потери}}.$$

Другими словами, на определенной итерации m алгоритм вычисляет градиент текущей потери для каждого наблюдения, а затем вписывает регрессионное дерево в эти псевдоостатки. На втором шаге он выявляет оптимальное постоянное предсказание для каждого терминального узла, которое минимизирует инкрементную потерю в результате добавления нового ученика в ансамбль.

Это отличается от автономных деревьев решений и случайных лесов, где предсказание зависит от результирующих значений тренировочных образцов, присутствующих в соответствующем терминальном или листовом узле: их среднего значения в случае регрессии или частоты положительного класса в случае бинарной классификации. Концентрация внимания на градиенте функции потери также влечет за собой то, что градиентный бустинг использует регрессионные деревья для усвоения регрессионных и классификационных правил, поскольку градиент всегда является непрерывной функцией.

Окончательная ансамблевая модель делает предсказания на основе взвешенной суммы предсказаний индивидуальных деревьев решений, каждое из которых было натренировано минимизировать ансамблевую потерю с учетом предшествующего предсказания для заданного набора значений признаков (рис. 11.2).

Градиентно-бустинговые деревья продемонстрировали передовую результативность на многочисленных эталонах классификации, регрессии и ранжирования. Они, вероятно, являются наиболее популярным алгоритмом ансамблевого обучения и как автономный предсказатель в широком круге соревнований по машинному обучению, и в реальных производственных конвейерах, например, для предсказания кликабельности для онлайн-рекламных сообщений.

Успех градиентного бустинга основан на его способности усваивать сложные функциональные связи в инкрементном стиле. Гибкость этого алгоритма требует тщательного управления риском переобучения путем настройки гиперпараметров,

которые ограничивают присущую модели тенденцию к усвоению шума в тренировочных данных в отличие от сигнала.

Мы представим ключевые механизмы управления сложностью градиентно-бустинговой древесной модели, а затем проиллюстрируем настройку модели с помощью реализации в библиотеке `sklearn`.



Рис. 11.2. Градиентный бустинг: поэтапная минимизация произвольных функций потерь

Как тренировать и настраивать модели на основе GBM

Двумя ключевыми движущими силами результативности градиентного бустинга являются размер ансамбля и сложность составляющих его деревьев решений.

Контроль за сложностью деревьев решений призван избежать усвоения узкоспецифичных правил, которые обычно охватывают очень малое число образцов в листовых узлах. В предыдущей главе мы рассмотрели наиболее эффективные ограничения, которые лимитируют способность дерева решений достигать переподгонки к тренировочным данным. Они включают следующие требования:

- ◆ минимальное число образцов для разбивки узла или принятия его в качестве терминального;
- ◆ либо минимальное улучшение качества узла, измеряемое чистотой, энтропией или среднеквадратической ошибкой в случае регрессии.

В дополнение к прямому управлению размером ансамбля существуют различные методы регуляризации, такие как усадка, с которой мы столкнулись в контексте

линейно-регрессионных моделей с гребневой и лассо-регуляризацией в главе 7. Кроме того, методы рандомизации, используемые в контексте случайных лесов, также обычно применяются к градиентно-бустинговым машинам.

Размер ансамбля и досрочная остановка

Каждая итерация бустинга стремится снизить тренировочную потерю, вследствие чего в случае крупного ансамбля тренировочная ошибка может потенциально стать очень малой, увеличивая риск перепогонки и низкой результативности на ранее не встречавшихся данных. Перекрестный контроль является лучшим подходом для отыскания оптимального размера ансамбля, который минимизирует ошибку обобщения, т. к. он зависит от применения и имеющихся данных.

Поскольку размер ансамбля должен быть указан перед тренировкой, полезно отслеживать результативность на контрольном наборе и прерывать тренировочный процесс, когда для заданного числа итераций контрольная ошибка больше не уменьшается. Это техническое решение называется досрочной остановкой и часто используется для моделей, требующих большого числа итераций и склонных к перепогонке, включая глубокие нейронные сети.

Следует иметь в виду, что использование досрочной остановки с тем же контрольным набором для большого числа испытаний также приведет к перепогонке, только в этом случае к конкретному контрольному набору, а не к тренировочному набору. При разработке стратегии торговли лучше всего избегать проведения большого числа экспериментов, т. к. риск ложных обнаружений значительно возрастает. В любом случае для получения несмещенной оценки ошибки обобщения следует иметь отложенный набор.

Усадка и темп усвоения

Усадочные методы применяют штраф к модельной функции потерь за повышенную сложность модели. Для бустинговых ансамблей усадка может быть применена путем снижения вклада каждого нового члена ансамбля на коэффициент в промежутке между 0 и 1. Этот коэффициент называется *темпом усвоения* в бустинговом ансамбле. Снижение темпа усвоения увеличивает усадку, поскольку снижает вклад каждого нового дерева решений в ансамбль.

Темп усвоения имеет эффект противоположный размеру ансамбля, тяготеющему к увеличению при более низких темпах усвоения. Было обнаружено, что более низкие темпы усвоения в сочетании с более крупными ансамблями снижают тестовую ошибку, в частности для регрессии и оценки вероятности. Большие числа итераций вычислительно обходятся дороже, но часто достижимы с быстрыми современными реализациями при условии, что индивидуальные деревья остаются мелкими. В зависимости от реализации можно также использовать адаптивные темпы усвоения, которые приспособляются к числу итераций, как правило, понижая влияние деревьев, добавляемых в процессе позже. Далее мы рассмотрим несколько примеров.

Подвыборка и стохастический градиентный бустинг

Как подробно обсуждалось в предыдущей главе, агрегирование бутстраповских выборок (бэггинг) улучшает результативность в иных случаях шумного классификатора.

Стохастический градиентный бустинг на каждой итерации использует взятие образцов без возмещения, выращивая следующее дерево на подмножестве тренировочных образцов. Преимуществом такого подхода является меньшее вычислительное усилие и нередко более высокая точность, но подвыборка должна сочетаться с усадкой.

Как вы видите, число гиперпараметров продолжает увеличиваться, подстегивая число потенциальных комбинаций, что в свою очередь увеличивает риск ложных обнаружений при выборе наилучшей модели из большого числа испытаний параметров на лимитированном объеме тренировочных данных. Лучший подход состоит в поэтапных действиях и отборе значений параметров по отдельности либо в использовании комбинации подмножеств с низкой мощностью.

Как использовать градиентный бустинг с помощью библиотеки `sklearn`

Ансамблевый модуль библиотеки `sklearn` содержит реализацию градиентно-бустинговых деревьев для регрессии и классификации, как бинарной, так и мульти-классовой. Следующий ниже фрагмент кода инициализации градиентно-бустингового классификатора `GradientBoostingClassifier` иллюстрирует ключевые регулировочные параметры, которые мы ввели ранее, в дополнение к тем, с которыми мы знакомы из просмотра моделей автономных деревьев решений. Блокнот `gbm_tuning_with_sklearn.ipynb` содержит примеры исходного кода для этого раздела.

Располагаемые функции потерь включают экспоненциальную потерю, которая приводит к алгоритму AdaBoost, и девиантность¹, которая соответствует логистической регрессии для вероятностных выходов. Мера качества узла `friedman_mse` является вариантом среднеквадратической ошибки, включающим показатель улучшения (справочные материалы со ссылками на исходные работы см. в репозитории GitHub), как показано в следующем фрагменте кода:

```
gb_clf = GradientBoostingClassifier(  
    loss='deviance',      # deviance = логистическая перп; exponential:  
                          # AdaBoost  
    learning_rate=0.1,    # сжимает вклад каждого дерева  
    n_estimators=100,     # число этапов бустинга  
    subsample=1.0,        # доля образцов для подгонки базовых учеников  
    criterion='friedman_mse', # служит мерой качества разбивки
```

¹ Девиантность (deviance) — это своего рода расстояние между двумя вероятностными моделями и сводится к двукратному логарифмическому отношению правдоподобий между двумя моделями ℓ_1/ℓ_0 , где ℓ_0 — это "меньшая" модель. — Прим. перев.

```

min_samples_split=2,
min_samples_leaf=1,
min_weight_fraction_leaf=0.0, # мин. доля суммы весов
max_depth=3,                  # необязательное значение, зависит от взаимодействия
min_impurity_decrease=0.0,
min_impurity_split=None,
init=None,
random_state=None,
max_features=None,
verbose=0,
max_leaf_nodes=None,
warm_start=False,
presort='auto',
validation_fraction=0.1,
n_iter_no_change=None,
tol=0.0001)

```

Подобно классификатору `AdaBoostClassifier`, эта модель не может справляться с пропущенными значениями. Как показано в следующем фрагменте кода, мы снова будем использовать 12-блочный перекрестный контроль, получая ошибки с целью классифицирования направленного финансового возврата для скользящих одномесячных периодов владения:

```

gb_cv_result = run_cv(gb_clf, y=y_clean, X=X_dummies_clean)
gb_result = stack_results(gb_cv_result)

```

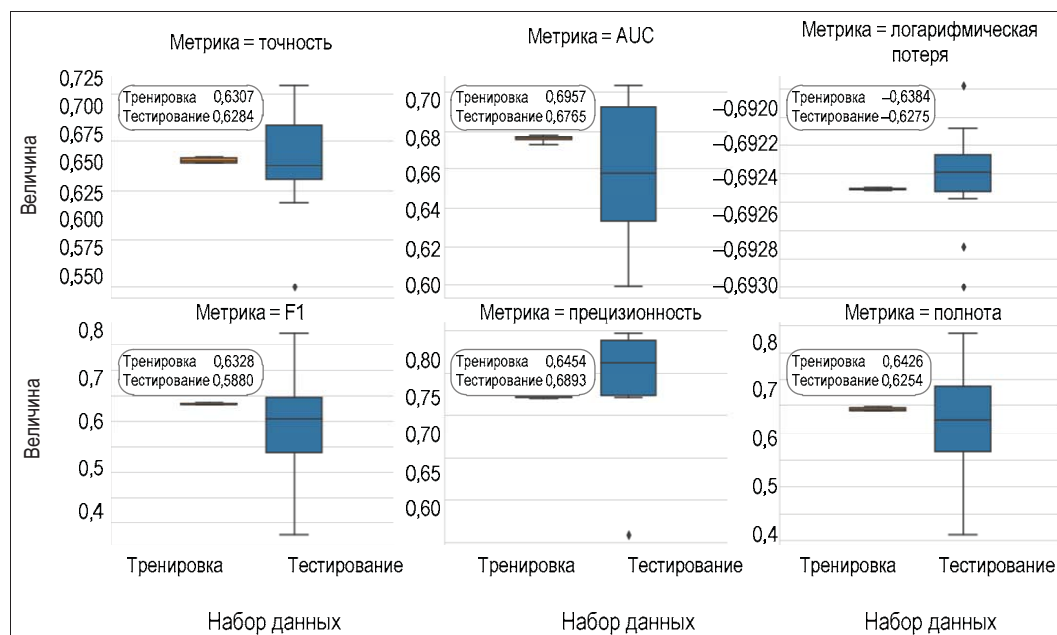


Рис. 11.3. Результаты работы градиентно-бустингового классификатора

Мы разберем и построим график результата, найдя небольшое улучшение — используя значения параметров, принятые по умолчанию — по сравнению с классификатором `AdaBoostClassifier` (рис. 11.3).

Как настраивать параметры с помощью класса *GridSearchCV*

Класс `GridSearchCV` в модуле `model_selection` обеспечивает систематическое оценивание всех комбинаций значений гиперпараметров, которые мы хотели бы протестировать. В следующем фрагменте кода мы проиллюстрируем этот функционал для семи регулировочных параметров, которые при их определении приведут к $24 \times 32 \times 4 = 576$ разным модельным конфигурациям:

```
cv = OneStepTimeSeriesSplit(n_splits=12)

param_grid = dict(
    learning_rate=[.01, .1, .2],
    max_depth=list(range(3, 13, 3)),
    max_features=['sqrt', .8, 1],
    min_impurity_decrease=[0, .01],
    min_samples_split=[10, 50],
    n_estimators=[100, 300],
    subsample=[.8, 1],
)

all_params = list(product(*param_grid.values()))
print('# Модели = :', len(all_params))

# Models = : 576
```

Метод `fit()` выполняет перекрестный контроль, используя пользовательский класс `OneStepTimeSeriesSplit` и отметку `roc_auc` для вычисления 12 блоков. Библиотека `sklearn` дает сохранить результат, как и в случае любой другой модели, используя реализацию консервации данных `joblib`, как показано в следующем фрагменте кода:

```
gs = GridSearchCV(gb_clf,
                  param_grid,
                  cv=cv,
                  scoring='roc_auc',
                  verbose=3,
                  n_jobs=-1,
                  return_train_score=True)

gs.fit(X=X, y=y)

# зафиксировать состояние результата, используя joblib с целью
# более эффективного хранения крупных массивов numpy
joblib.dump(gs, 'gbm_gridsearch.joblib')
```

После завершения работы объект `GridSearchCV` имеет несколько дополнительных атрибутов, к которым можно обратиться после загрузки законсервированного результата, чтобы узнать, какая комбинация гиперпараметров показала наилучшую результативность, и ее среднюю перекрестно-контрольную отметку AUC. Результатом будет скромное улучшение по сравнению со стандартными значениями по умолчанию. Это показано в следующем фрагменте кода:

```
pd.Series(gridsearch_result.best_params_)
```

```
learning_rate      0.0100
max_depth          9.0000
max_features       1.0000
min_impurity_decrease  0.0100
min_samples_split  10.0000
n_estimators       300.0000
subsample          0.8000
dtype: float64
```

```
gridsearch_result.best_score_
0.6853
```

Влияние параметров на тестовые отметки

Объект `GridSearchCV` сохраняет средние перекрестно-контрольные отметки, что позволяет анализировать то, как разные гиперпараметрические условия влияют на результат.

Шесть роевых графиков библиотеки `seaborn` на рис. 11.4 слева показывают распределение отметок AUC по всем значениям параметров. В данном случае самые высокие тестовые отметки AUC требовали низкого значения параметра темпа усвоения `learning_rate` и большого значения параметра максимального числа признаков `max_features`. Некоторые параметрические настройки, такие как низкий темп усвоения `learning_rate`, дают широкий спектр результатов, которые зависят от дополнительных настроек других параметров. Другие параметры совместимы с высокими отметками для всех настроек, используемых в эксперименте.

Теперь мы рассмотрим, как гиперпараметрические настройки совместно влияют на среднюю перекрестно-контрольную отметку. Для того чтобы глубже разобраться в характере взаимодействия параметрических настроек, мы можем натренировать регрессор `DecisionTreeRegressor`, в котором результатом является средняя тестовая отметка, а настройки параметров закодированы в формате фиктивных значений категориальных переменных (в кодировке с одним активным состоянием) (подробности см. в блокноте). Древесная структура подчеркивает, что использование всех признаков (`max_features_1`), низкого значения параметра `learning_rate` и параметра `max_depth`, равным более трем, привело к наилучшим результатам (рис. 11.5).

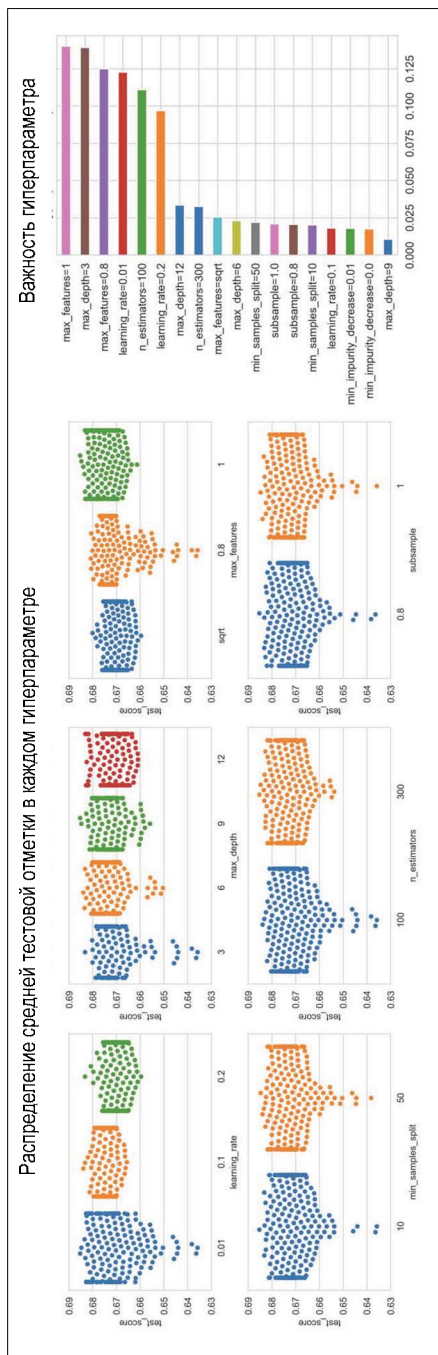


Рис. 11.4. Роевые графики с распределением отметок AUC

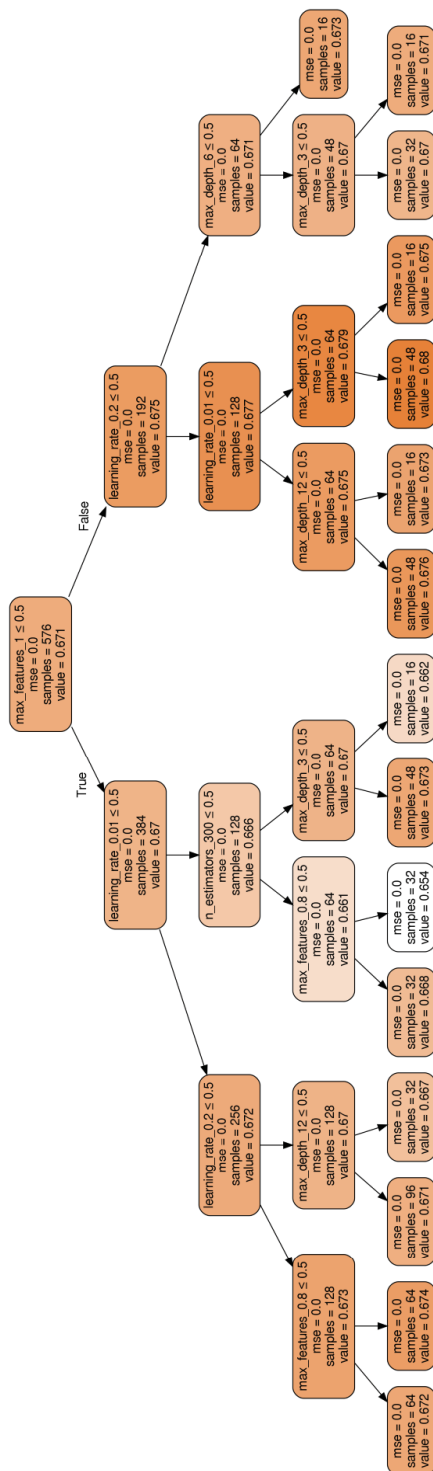


Рис. 11.5. Результаты работы древесного классификатора

Гистограмма на рис. 11.4 справа показывает влияние гиперпараметрических настроек на получение разных результатов, измеряемых важностью их признаков для дерева решений, которое выращивается до максимальной глубины. Естественно, что признаки, которые появляются вблизи вершины дерева, также аккумулируют самые высокие отметки важности.

Как тестировать на отложенном наборе данных

Наконец, мы хотели бы оценить результативность наилучшей модели на отложенном наборе, который мы исключили из упражнения по параметрическому поиску с помощью объекта `GridSearchCV`. Он содержит последние шесть месяцев периода выборки (до февраля 2018 г.; подробнее см. в блокноте). Мы получаем оценку результативности обобщения на основе отметки $AUC = 0,6622$, используя следующий фрагмент кода:

```
best_model = gridsearch_result.best_estimator_  
preds = best_model.predict(test_feature_data)  
roc_auc_score(y_true=test_target, y_score=preds)
```

0.6622

Недостатком реализации градиентного бустинга в библиотеке `sklearn` является лимитированная скорость вычисления, что затрудняет быстрое опробование разных гиперпараметрических настроек. В следующем разделе мы увидим, что за последние несколько лет появилось несколько оптимизированных реализаций, которые значительно сокращают время, необходимое для тренировки даже крупномасштабных моделей, и которые в значительной степени способствовали расширению сферы применения этого высокоэффективного алгоритма.

Быстро масштабируемые реализации градиентно-бустинговых машин

В последние годы в нескольких новых реализациях градиентного бустинга использовались различные инновации, которые ускоряют тренировку, повышают эффективность эксплуатации ресурсов и позволяют алгоритму масштабироваться до очень крупных наборов данных. Новые реализации и их источники таковы:

- ♦ библиотека `XGBoost` (extreme gradient boosting, экстремальный градиентный бустинг), впервые опубликованная в 2014 г. Тяньци Чэнем в Вашингтонском университете;
- ♦ библиотека `LightGBM`, выпущенная в январе 2017 г. корпорацией Microsoft;
- ♦ библиотека `CatBoost`, выпущенная в апреле 2017 г. компанией Яндекс.

Эти инновации решают конкретные трудности тренировки градиентно-бустинговой модели (подробные справочные материалы см. в файле `README` этой главы в репозитории `GitHub`). Реализация `XGBoost` стала первой новейшей реализацией, получившей популярность: среди 29 выигранных решений, опубликованных порта-

лом Kaggle в 2015 г., в 17 решениях использовалась реализация XGBoost. Восемь из них опирались исключительно на XGBoost, в то время как в других комбиниrowались XGBoost и нейронные сети.

Сначала мы представим ключевые инновации, которые появились с течением времени и впоследствии сошлись (вследствие чего большинство функционала доступно для всех реализаций), и после этого проиллюстрируем их реализацию.

Как алгоритмические инновации стимулируют результативность

Случайные леса можно тренировать параллельно, выращивая отдельные деревья на независимых бутстраповских выборках. Напротив, последовательный подход градиентного бустинга замедляет тренировку, что в свою очередь усложняет эксперименты с большим числом гиперпараметров, которые необходимо адаптировать к характеру задачи и набору данных.

Для того чтобы расширить ансамбль деревьями, тренировочный алгоритм инкрементно минимизирует ошибку предсказания относительно отрицательного градиента ансамблевой функции потерь, аналогично обычному оптимизатору на основе градиентного спуска. Следовательно, вычислительная стоимость во время тренировки пропорциональна времени, необходимому для оценивания влияния потенциальных точек разбиения по каждому признаку на подгонку дерева решений к текущему градиенту.

Аппроксимация второпорядковой функции потерь

Наиболее важные алгоритмические инновации снижают стоимость оценивания функции потерь, используя аппроксимации, основанные на производных второго порядка, напоминающие ньютоновский метод отыскания стационарных точек. В результате выставляющие отметки потенциальным разбиениям во время жадного расширения дерева происходит быстрее по сравнению с использованием функции полной потерь.

Как упоминалось ранее, градиентно-бустинговую модель тренируют в инкрементном стиле, минимизируя сочетания ошибки предсказания и регуляризационного штрафа для ансамбля H_M . Обозначив предсказание ансамблем результата y_i после шага m как $\hat{y}_i^{(m)}$, дифференцируемую выпуклую функцию потерь, измеряющую разницу между результатом и предсказанием, как l , и штраф, который увеличивается со сложностью ансамбля H_M , как Ω , инкрементная гипотеза h_m призвана минимизировать следующую цель:

$$\begin{aligned} L^{(m)} &= \sum_{i=1}^n \underbrace{l(y_i, \hat{y}_i^{(m)})}_{\text{потеря на шаге } m} + \sum_{i=1}^n \underbrace{\Omega(H_m)}_{\text{регуляризация}} = \\ &= \sum_{i=1}^n l\left(y_i, \hat{y}_i^{(m)} + \underbrace{h_m(x_i)}_{\text{дополнительное дерево}}\right) + \Omega(H_m). \end{aligned}$$

Регуляризационный штраф помогает избежать перепогонки, отдавая предпочтение отбору модели, которая использует простые и предсказательные регрессионные деревья. Например, в случае реализации XGBoost штраф для регрессионного дерева h зависит от числа листьев в расчете на дерево T , отметок регрессионного дерева для каждого терминального узла w и гиперпараметров γ и λ . Это отражено в следующей формуле:

$$\Omega(h) = \gamma T + \frac{1}{2} \lambda \|w\|^2.$$

Поэтому на каждом шаге алгоритм жадно добавляет гипотезу h_m , которая больше всего улучшает регуляризованную цель. Аппроксимация второпорядковой функции потери, основанная на разложении ряда Тейлора, ускоряет оценивание цели, как резюмируется в следующей формуле:

$$L^{(m)} = \sum_{i=1}^n \left[g_i f_m(x_i) + \frac{1}{2} h_i f_m^2(x_i) \right] + \Omega(h_m).$$

Здесь g_i — это первопорядковый градиент функции потери перед добавлением нового узника для заданного значения признака; h_i — соответствующее значение второпорядкового градиента (или гессииана), как показано в следующих формулах:

$$g_i = \partial_{\hat{y}_i^{(m-1)}} l(y_i, \hat{y}_i^{(m-1)});$$

$$h_i = \partial_{\hat{y}_i^{(m-1)}}^2 l(y_i, \hat{y}_i^{(m-1)}).$$

Реализация XGBoost была первой с открытыми исходными кодами, которая использовала эту аппроксимацию функции потери для вычисления оптимальных листовых отметок для заданной древесной структуры и соответствующего значения функции потери. Отметка представляет собой отношение сумм градиента и гессиианы для образцов в терминальном узле. Указанная реализация использует это значение для выставления отметки прироста информации, который получается в результате разбивки, аналогичной мерам смешанности узлов, которые мы встречали в предыдущей главе, но применимой к произвольным функциям потерь (подробный математический вывод см. в справочных материалах репозитория GitHub).

Упрощенные алгоритмы поиска разбивок

Реализация градиентного бустинга в библиотеке `sklearn` отыскивает оптимальную разбивку, которая перечисляет все варианты для непрерывных признаков. Этот точный жадный алгоритм вычислительно очень требователен, потому что перед тем, как выставить отметки потенциально очень большому числу вариантов разбивок и принять решение, он сначала должен отсортировать данные по значениям признаков. Этот подход сталкивается с проблемами, когда данные не помещаются в память или во время тренировки в распределенных условиях на нескольких компьютерах.

Приближенный алгоритм отыскания разбинок снижает число точек разбиения, закрепляя значения признаков за определяемым пользователем множеством корзин, что также может весомо снизить требования к памяти во время тренировки, поскольку для каждой корзины необходимо хранить только одну разбивку. Библиотека XGBoost внедрила алгоритм квантильного наброска², который также смог разделять взвешенные тренировочные образцы на процентильные корзины, достигая равномерного распределения. Библиотека XGBoost также представила возможность работать с разреженными данными, вызванными пропущенными значениями, частой статистикой нулевого градиента и кодированием с одним активным состоянием, а также может усваивать оптимальное направление, принимаемое по умолчанию для определенной разбивки. В результате алгоритму требуется оценивать только не пропущенные значения.

В отличие от нее в библиотеке LightGBM используется *градиентная односторонняя выборка* (gradient-based one-side sampling, GOSS), исключая значительную долю образцов с малыми градиентами, и для оценивания прироста информации и отбора надлежащего значения разбивки задействуется только оставшаяся часть. Образцы с более крупными градиентами требуют большей тренировки и тяготеют к внесению большего вклада в получение информации. В библиотеке LightGBM также используется группирование взаимоисключающих признаков, совмещая признаки, которые являются взаимоисключающими, поскольку они редко принимают ненулевые значения одновременно; это делается с целью снизить число признаков. В результате на момент выпуска реализация библиотеки LightGBM была самой быстрой.

Поуровневый рост против полистового роста

Реализация библиотеки LightGBM отличается от реализаций библиотек XGBoost и CatBoost тем, как она расставляет приоритеты в отношении того, какие узлы следует разбивать. Реализация LightGBM принимает решение о разбиении в полистовом порядке, т. е. разбивает листовой узел, который максимизирует прирост информации, даже если это приводит к несбалансированным деревьям. В отличие от нее реализации XGBoost и CatBoost расширяют все узлы поуровневым образом и сначала разбивают все узлы на заданную глубину и только после этого добавляют больше уровней. Эти два подхода расширяют узлы в разном порядке и дают разные результаты, за исключением полных деревьев. На рис. 11.6 показаны оба подхода.

Полистовые разбивки в реализации LightGBM тяготеют к увеличению сложности модели и могут ускорить схождение, но и увеличивают риск переподргонки. Дерево, выращенное в глубину с n уровнями, имеет до 2^n терминальных узлов, в то время как полистовое дерево с 2^n листьями может иметь значительно больше уровней

² Квантильный набросок (quantiles sketch) — это потоковый алгоритм оценивания распределения значений, который приблизительно отвечает на запросы о ранге значения, функции массы вероятности распределения (PMF) или гистограмме, кумулятивной функции распределения (CDF) и квантилях (медиана, минимум, максимум, 95-й процентиль и т. д.).

См. <https://datasketches.github.io/docs/Quantiles/QuantilesOverview.html>. — Прим. перев.

и содержать соответственно в некоторых листьях меньше образцов. Следовательно, настройка параметра `num_leaves` в библиотеке LightGBM требует дополнительной осторожности, и библиотека позволяет одновременно управлять параметром максимальной глубины `max_depth` с целью избежать чрезмерного дисбаланса узлов. Более недавние версии библиотеки LightGBM также предлагают поуровневый рост дерева.

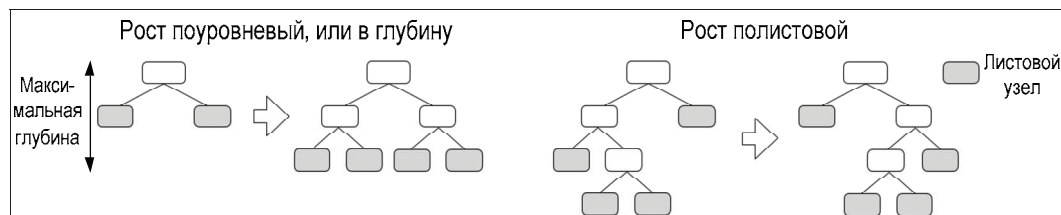


Рис. 11.6. Сравнение поуровневого роста (в глубину) и полистового роста

Тренировка на основе GPU

Все новые реализации поддерживают тренировку и предсказание на одном или нескольких графических процессорах (GPU), достигая значительных ускорений. Они совместимы с текущими GPU с поддержкой программно-аппаратной архитектуры параллельных вычислений CUDA. Требования к установке варьируются и быстро развиваются. Реализации XGBoost и CatBoost работают для нескольких текущих ее версий, но реализация LightGBM может потребовать локальной компиляции (см. репозиторий GitHub со ссылками на соответствующую документацию).

Ускорения зависят от библиотеки и типа данных и варьируются от низких, однократных до многократных. Активация GPU требует только изменения параметра задачи и никаких других гиперпараметрических модификаций.

Отсеивающая регуляризация DART — отсев для деревьев

В 2015 г. Рашми (Rashmi) и Гилад-Бахрах (Gilad-Bachrach) предложили новую модель тренировки градиентно-бустинговых деревьев, которая ориентировалась на решение проблемы, названной ими чрезмерной специализацией: деревья, добавляемые во время более поздних итераций, как правило, влияют на предсказание только нескольких экземпляров, внося незначительный вклад относительно остальных экземпляров. Однако вневыборочная результативность модели может пострадать, и она может стать чрезмерно чувствительной к вкладам малого числа деревьев, добавленных ранее в этом процессе.

Новые алгоритмы задействуют отсеивание, которое успешно использовалось для автоматического обучения более точных глубоких нейронных сетей, где отсев отключает случайную долю нейронных связей в ходе самообучения. В результате узлы в более высоких слоях не могут опираться на малое число связей при передаче информации, необходимой для предсказания. Этот метод внес значительный вклад в успех глубоких нейронных сетей в решении многих задач, а также был использо-

ван вместе с другими техническими решениями автоматического обучения, такими как логистическая регрессия, отключая случайную долю признаков. Случайные леса и стохастический градиентный бустинг также отсеивают случайное подмножество признаков.

Регуляризация DART работает на уровне деревьев и вместо индивидуальных признаков отсеивает полные деревья. Цель состоит в том, чтобы в ансамбле, сгенерированном с помощью регуляризации DART, деревья вносили более равномерный вклад в окончательное предсказание. В некоторых случаях было показано, что это дает более точные предсказания для задач ранжирования, регрессии и классификации. Указанный подход был впервые реализован в реализации LightGBM, а также доступен для реализации XGBoost.

Трактовка категориальных признаков

Реализации CatBoost и LightGBM обрабатывают категориальные переменные напрямую без необходимости кодирования фиктивных значений категориальных переменных.

Реализация CatBoost (которая так названа из-за ее обработки категориальных признаков) охватывает несколько вариантов обработки таких признаков, в дополнение к автоматическому кодированию с одним активным состоянием, и закрепляет за числовыми значениями либо категории отдельных признаков, либо комбинации категорий нескольких признаков. Другими словами, библиотека CatBoost может создавать новые категориальные признаки из комбинаций существующих признаков. Числовые значения, ассоциированные с категориальными уровнями индивидуальных признаков либо комбинаций признаков, зависят от их связи с результирующим значением. В случае классификации это связано с вероятностью наблюдать положительный класс, вычисляемый кумулятивно над выборкой, на основе априорного распределения и с коэффициентом сглаживания. Более подробные числовые примеры см. в документации.

Реализация LightGBM группирует уровни категориальных признаков, максимизируя однородность (либо минимизируя дисперсию) внутри групп по отношению к результирующим значениям.

Реализация XGBoost не работает с категориальными признаками напрямую и требует кодирования с одним активным состоянием (кодирования фиктивных значений категориальных переменных).

Дополнительные признаки и оптимизации

Реализация XGBoost оптимизировала вычисление в нескольких отношениях, обеспечив многопоточную обработку³, поддерживая данные в памяти в сжатых столб-

³ Многопоточная обработка (multithreading) — это способность CPU (или единственного ядра в многоядерном процессоре) обеспечивать несколько нитей исполнения одновременно. В русской терминологии за термином *thread* укрепился не точный перевод "поток", хотя его протогогерманский предок *thread* по форме и смыслу является родственником нашим "прясти", "прядь", т. е. процесс сплетается

цовых блоках, где каждый столбец сортируется по соответствующему значению признака. XGBoost вычисляет этот макет входных данных один раз перед тренировкой и многократно использует его для амортизации дополнительных предоперационных расходов. Поиск статистики разбивки по столбцам становится линейным сканированием при использовании квантилей и может выполняться параллельно с простой поддержкой подвыборки столбцов.

Выпущенные по очереди реализации LightGBM и CatBoost строились на основе этих инноваций, и реализация LightGBM еще больше ускорила тренировку за счет оптимизированной нитиевой обработки и снижения потребления памяти. Ввиду общедоступности исходных кодов с течением времени наблюдается тенденция сближения этих библиотек.

Библиотека XGBoost также поддерживает ограничения монотонности. Эти ограничения гарантируют, что значения для определенного признака являются только положительными или отрицательными в связи с результатом на всем его интервале. Они полезны для встраивания внешних допущений о модели, о которых известно, что они являются истинными.

Как использовать библиотеки XGBoost, LightGBM и CatBoost

Библиотеки XGBoost, LightGBM и CatBoost предлагают интерфейсы для нескольких языков, включая Python, и имеют интерфейс с библиотекой `sklearn`, совместимый с другим функционалом `sklearn`, таким как класс `GridSearchCV`, и его собственными методами тренировки градиентно-бустинговых моделей и предсказания на их основе. В блокноте `gbm_baseline.ipynb` иллюстрируется использование интерфейса с библиотекой `sklearn` для каждой реализации. Библиотечные методы часто лучше документированы, а также просты в использовании, поэтому мы ими воспользуемся для иллюстрации использования этих моделей.

Этот процесс предполагает создание специфичных для библиотек форматов данных, настройку различных гиперпараметров и оценивание результатов. Все это будет описано в последующих разделах. Сопроводительный блокнот `sklearn_gbm_tuning.ipynb` содержит файлы `gbm_utils.py` и `gbm_params.py`, которые совместно обеспечивают следующий ниже функционал и произвели соответствующие результаты.

Как создавать двоичные форматы данных

Все библиотеки имеют собственный формат данных для предвычисления признаков статистики для ускорения поиска точек разбиения, как описано выше. Кроме того, их состояние может сохраняться на диск, тем самым ускоряя начало последующей тренировки.

из нитей (ср. нить рассуждения, если брать антропоморфную аналогию). В зарубежной технической литературе данное понятие объясняется именно на нитях, которые рассматриваются как основные инфраструктурные единицы планирования в операционных системах. — *Прим. перев.*

Следующий фрагмент кода создает тренировочный и контрольный наборы двоичных данных для каждой модели, используемой с пользовательским объектом `OneStepTimeSeriesSplit`:

```
def get_datasets(features, target, kfold, model='xgboost'):
    cat_cols = ['year', 'month', 'age', 'msize', 'sector']
    data = {}
    for fold, (train_idx, test_idx) in enumerate(kfold.split(features)):
        print(fold, end=' ', flush=True)
        if model == 'xgboost':
            data[fold] = {'train':
                xgb.DMatrix(label=target.iloc[train_idx],
                           data=features.iloc[train_idx],
                           nthread=-1), # использовать имеющиеся нити
                'valid': xgb.DMatrix(label=target.iloc[test_idx],
                                    data=features.iloc[test_idx],
                                    nthread=-1)}
        elif model == 'lightgbm':
            train = lgb.Dataset(label=target.iloc[train_idx],
                               data=features.iloc[train_idx],
                               categorical_feature=cat_cols,
                               free_raw_data=False)

            # выровнять гистограммы контрольного набора
            # с тренировочным набором
            valid = train.create_valid(label=target.iloc[test_idx],
                                      data=features.iloc[test_idx])

            data[fold] = {'train': train.construct(),
                          'valid': valid.construct()}
        elif model == 'catboost':
            # получить индексы категориальных признаков
            cat_cols_idx = [features.columns.get_loc(c) for c in cat_cols]
            data[fold] = {'train': Pool(label=target.iloc[train_idx],
                                       data=features.iloc[train_idx],
                                       cat_features=cat_cols_idx),
                          'valid': Pool(label=target.iloc[test_idx],
                                       data=features.iloc[test_idx],
                                       cat_features=cat_cols_idx)}

    return data
```

Доступные варианты немного различаются:

- ◆ `xgboost` позволяет использовать все имеющиеся нити;
- ◆ `lightgbm` явно выравнивает квантили, созданные для контрольного набора, с тренировочным набором;

- ♦ реализации `catboost` требуются признаковые столбцы, идентифицируемые по индексам, а не меткам.

Как настраивать гиперпараметры

Многочисленные гиперпараметры перечислены в файле `gbm_params.py`. Каждая библиотека имеет параметрические настройки для:

- ♦ указания общих целей и обучающегося алгоритма;
- ♦ проектирования базовых учеников;
- ♦ применения различных регуляризационных методов;
- ♦ обработки досрочной остановки во время тренировки;
- ♦ активации использования GPU или параллелизации на CPU.

В документации по каждой библиотеке подробно описываются различные параметры, которые могут относиться к одной и той же концепции, но имеют различные имена в разных библиотеках. Репозиторий GitHub содержит ссылки на сайт, где выделяются параметры, относимые к библиотекам XGBoost и LightGBM.

Целевые функции и функция потерь

Указанные библиотеки поддерживают несколько алгоритмов бустинга, включая градиентный бустинг для деревьев и линейных базовых учеников, а также DART в случае библиотек LightGBM и XGBoost. Библиотека LightGBM поддерживает алгоритм GOSS, который мы описали ранее, а также случайные леса.

Привлекательность градиентного бустинга состоит в эффективной поддержке произвольных дифференцируемых функций потерь, и каждая библиотека предлагает различные варианты для задач регрессии, классификации и ранжирования. В дополнение к выбираемой функции потери можно использовать дополнительные оценочные метрики, производя мониторинг результативности во время тренировки и перекрестного контроля.

Усвоение параметров

В градиентно-бустинговых моделях, как правило, деревья решений используются для улавливания взаимодействия между признаками, и размер индивидуальных деревьев является наиболее важным регулировочным параметром. В библиотеках XGBoost и CatBoost значение максимальной глубины `max_depth` установлено по умолчанию равным 6. В отличие от них в библиотеке LightGBM по умолчанию используется значение числа листьев `num_leaves`, равное 31, которое соответствует пяти уровням для сбалансированного дерева, но не накладывает ограничений на число уровней. Во избежание перепогонки параметр `num_leaves` должен быть ниже $2^{\text{max_depth}}$. Например, для значения параметра `max_depth` с хорошей результативностью, равного 7, следует устанавливать параметр `num_leaves` равным 70–80, а не $2^7 = 128$, либо непосредственно ограничивать `max_depth`.

Число деревьев или бустинговых итераций определяет совокупный размер ансамбля. Все библиотеки поддерживают параметр досрочной остановки `early_stopping`, прерывающий тренировку, как только функция потери перестает регистрировать дальнейшие улучшения в течение заданного числа итераций. В связи с этим обычно лучше всего устанавливать большое число итераций и прекращать тренировку, основываясь на предсказательной результативности на контрольном наборе.

Регуляризация

Все библиотеки реализуют регуляризационные стратегии для базовых учеников, такие как минимальные значения числа образцов или минимальный прирост информации, необходимый для разбивок и листовых узлов.

Они также поддерживают регуляризацию на ансамблевом уровне с использованием усадки через темп усвоения, которая ограничивает вклад новых деревьев. Также имеется возможность реализации адаптивного темпа усвоения посредством функций обратного вызова, которые снижают темп усвоения по мере продвижения тренировки, что успешно используется в нейросетевом контексте. Кроме того, градиентно-бустинговая функция потери может быть регуляризована с помощью L_1 или L_2 , т. е. регуляризации, аналогичной используемой в линейных регрессионных моделях с гребневой и лассо-регуляризацией, путем модифицирования $\Omega(h_m)$ либо увеличения штрафа γ за добавление большего числа деревьев, как описано ранее.

Указанные библиотеки также позволяют использовать бэггинг или подвыборку столбцов с целью рандомизации роста случайных лесов и декоррелировать ошибки предсказания с целью снижения совокупной дисперсии. Квантизация признаков для отыскания приближенных разбивок добавляет более крупные корзины в качестве дополнительного варианта защиты от переподгонки.

Рандомизированный решеточный поиск

В целях разведывания гиперпараметрического пространства мы задаем значения ключевых параметров, которые мы хотели бы протестировать в комбинации. Библиотека `sklearn` поддерживает класс `RandomizedSearchCV` для перекрестного контроля подмножества параметрических комбинаций, которые отбираются случайно из заданных распределений. Мы реализуем собственную версию, которая позволит использовать досрочную остановку при мониторинге текущих наиболее результативных комбинаций для прерывания поискового процесса, как только будет получен удовлетворительный результат, вместо указания заданного числа итераций заранее.

С этой целью мы задаем решетку параметров в соответствии с параметрами каждой библиотеки, как и раньше, генерируем все комбинации, используя встроенный генератор декартовых произведений, предоставляемый библиотекой `itertools`, и случайно перемешиваем результат. В случае библиотеки `LightGBM` мы автоматически устанавливаем параметр `max_depth` как функцию от текущего значения параметра `num_leaves`, как показано в следующем фрагменте кода:

```

param_grid = dict(
    # общие варианты
    learning_rate=[.01, .1, .3],
    colsample_bytree=[.8, 1], # except catboost

    # lightgbm
    num_leaves=[2 ** i for i in range(9, 14)],
    boosting=['gbdt', 'dart'],
    min_gain_to_split=[0, 1, 5], # не поддерживается на GPU

all_params = list(product(*param_grid.values()))
n_models = len(all_params) # макс. число моделей для перекрестного контроля
shuffle(all_params)

```

Затем мы выполняем перекрестный контроль следующим образом:

```

results = pd.DataFrame()
GBM = 'lightgbm'
for n, test_param in enumerate(all_params, 1):
    cv_params = get_params(GBM)
    cv_params.update(dict(zip(param_grid.keys(), test_param)))
    if GBM == 'lightgbm':
        cv_params['max_depth'] = int(ceil(np.log2(cv_params['num_leaves'])))

    results[n] = run_cv(test_params=cv_params,
                        data=datasets,
                        n_splits=n_splits,
                        gb_machine=GBM)

```

Функция `run_cv` реализует перекрестный контроль для всех трех библиотек. Для примера с `light_gbm` указанный процесс выглядит следующим образом:

```

def run_cv(test_params, data, n_splits=10, gb_machine='xgboost'):
    """Настроить и проконтролировать с досрочной остановкой"""

    result = []
    cols = ['rounds', 'train', 'valid']
    for fold in range(n_splits):
        train = data[fold]['train']
        valid = data[fold]['valid']

        scores = {}
        model = lgb.train(params=test_params,
                          train_set=train,
                          valid_sets=[train, valid],
                          valid_names=['train', 'valid'],
                          num_boost_round=250,
                          early_stopping_rounds=25,

```

```

        verbose_eval=50,
        evals_result=scores)

    result.append([model.current_iteration(),
                  scores['train']['auc'][-1],
                  scores['valid']['auc'][-1]])

df = pd.DataFrame(result, columns=cols)
(df.mean()
 .append(df.std().rename({c: c + '_std' for c in cols}))
 .append(pd.Series(test_params)))

```

Метод `train()` также создает контрольные отметки, которые хранятся в словаре `scores`. Когда досрочная остановка вступает в силу, последняя итерация также является наилучшей отметкой. Дополнительные подробности см. в полной реализации в репозитории GitHub.

Как оценивать результаты

Используя GPU, мы можем тренировать модель за несколько минут и оценивать несколько сотен комбинаций параметров в течение нескольких часов, что займет много дней с использованием реализации в библиотеке `sklearn`. Для модели `LightGBM` мы исследуем факторную версию, использующую способность библиотек справляться с категориальными переменными, и фиктивную версию, использующую кодирование с одним активным состоянием.

Результаты доступны в хранилище данных HDF5 под названием `model_tuning.h5`. Примеры исходного кода оценивания модели находятся в блокноте `eval_results.ipynb`.

Результаты перекрестного контроля между моделями

При сравнении средней перекрестно-контрольной отметки AUC по четырем тестовым прогонам с тремя библиотеками мы обнаруживаем, что `CatBoost` производит немного более высокую отметку AUC для модели с самой высокой результативностью, при этом также порождает самый широкий разброс результатов (рис. 11.7).

Модель `CatBoost` с самой высокой результативностью использует следующие параметры (подробнее см. блокнот):

- ◆ `max_depth=12` и `max_bin=128`;
- ◆ `max_ctr_complexity=2`, который лимитирует число комбинаций категориальных признаков;
- ◆ `one_hot_max_size=2`, который исключает бинарные признаки из назначения признакам числовых переменных;
- ◆ `random_strength` отличается от 0 для рандомизации оценивания разбивок.

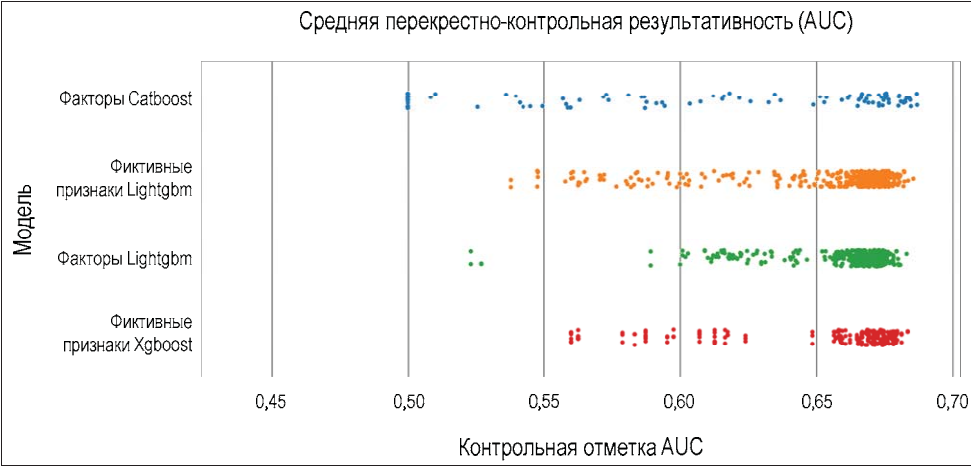


Рис. 11.7. Средняя перекрестно-контрольная результативность (AUC)

Тренировка проходит немного медленнее по сравнению с LightGBM и XGBoost (все используют GPU) со скоростью в среднем 230 секунд на модель.

Более подробный взгляд на модели с самой высокой результативностью для моделей LightGBM и XGBoost показывает, что факторная модель LightGBM достигает почти такой же хорошей результативности, как и две другие модели с гораздо меньшей сложностью. Она состоит в среднем из 41 дерева до трех уровней глубины не более восьми листьев каждый, а также с использованием регуляризации в форме `min_gain_to_split`. Она достигает значительно меньшей переподгонки на тренировочном наборе, с тренировочной отметкой AUC лишь немногим выше контрольной AUC. Она также тренируется намного быстрее, занимая всего 18 секунд на модель из-за ее меньшей сложности. На практике эта модель была бы предпочтительнее, поскольку она с большей вероятностью обеспечивает хорошую вневыборочную результативность. Подробности приведены в табл. 11.1.

Таблица 11.1. Бустинговые модели с самой высокой результативностью

| | Фиктивные признаки LightGBM | Фиктивные признаки XGBoost | Факторы LightGBM |
|--------------------------------|--------------------------------|-------------------------------|---------------------|
| Контрольная оценка AUC, % | 68,57 | 68,36 | 68,32 |
| Тренировочная оценка AUC, % | 82,35 | 79,81 | 72,12 |
| <code>learning_rate</code> | 0,1 | 0,1 | 0,3 |
| <code>max_depth</code> | 13 | 9 | 3 |
| <code>num_leaves</code> | 8192 | | 8 |
| <code>colsample_bytree</code> | 0,8 | 1 | 1 |
| <code>min_gain_to_split</code> | 0 | 1 | 0 |

Таблица 11.1 (окончание)

| | Фиктивные признаки LightGBM | Фиктивные признаки XGBoost | Факторы LightGBM |
|----------|-----------------------------|----------------------------|------------------|
| Раунды | 44,42 | 59,17 | 41,00 |
| Время, с | 86,55 | 85,37 | 18,78 |

На графике (рис. 11.8) показан эффект разных настроек параметра `max_depth` на контрольную отметку для моделей LightGBM и XGBoost: более мелкие деревья дают более широкий диапазон результатов и должны быть объединены с соответствующими настройками темпов усвоения и регуляризации для получения сильного результата, показанного в предыдущей таблице.

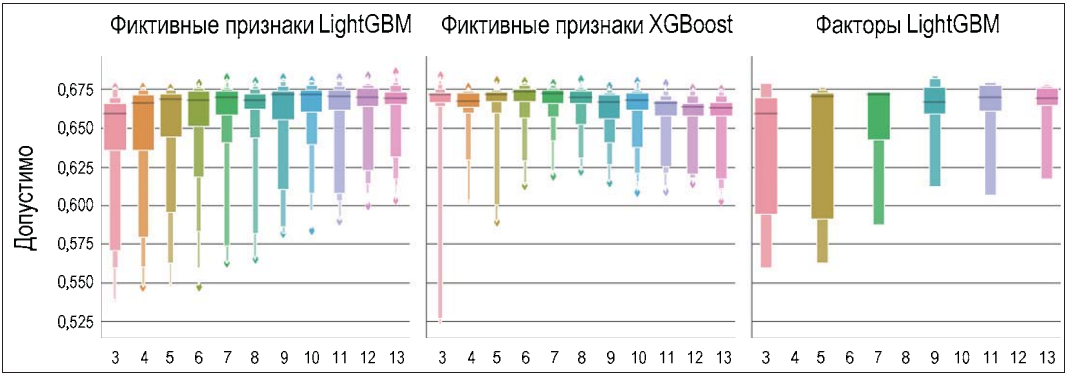


Рис. 11.8. Эффект разных настроек параметра `max_depth` на контрольную оценку для моделей LightGBM и XGBoost

Вместо регрессора `DecisionTreeRegressor`, как показано ранее, для оценивания статистической значимости разных признаков, касающихся контрольной отметки AUC, можно также использовать линейную регрессию. Для модели LightGBM с фиктивными признаками, где регрессия объясняет 68% вариации в результатах, мы обнаруживаем, что не был значимым только регуляризационный параметр `min_gain_to_split`:

Lightgbm Dummies

OLS Regression Results

| | | | |
|-------------------|------------------|---------------------|----------|
| Dep. Variable: | valid | R-squared: | 0.687 |
| Model: | OLS | Adj. R-squared: | 0.673 |
| Method: | Least Squares | F-statistic: | 26.94 |
| Date: | Sat, 20 Apr 2019 | Prob (F-statistic): | 7.92e-55 |
| Time: | 14:39:20 | Log-Likelihood: | 1018.7 |
| No. Observations: | 396 | AIC: | -2001. |
| Df Residuals: | 378 | BIC: | -1930. |
| Df Model: | 17 | | |
| Covariance Type: | HC3 | | |

| | coef | std err | z | P> z | [0.025 | 0.975] |
|----------------------|---------|-------------------|---------|---------|--------|--------|
| const | 0.6145 | 0.005 | 127.970 | 0.000 | 0.605 | 0.624 |
| boosting_gbtree | 0.0056 | 0.002 | 2.866 | 0.004 | 0.002 | 0.009 |
| learning_rate_0.1 | 0.0501 | 0.003 | 18.977 | 0.000 | 0.045 | 0.055 |
| learning_rate_0.3 | 0.0516 | 0.003 | 19.150 | 0.000 | 0.046 | 0.057 |
| max_depth_4 | 0.0060 | 0.005 | 1.094 | 0.274 | -0.005 | 0.017 |
| max_depth_5 | 0.0096 | 0.005 | 1.823 | 0.068 | -0.001 | 0.020 |
| max_depth_6 | 0.0153 | 0.005 | 3.024 | 0.002 | 0.005 | 0.025 |
| max_depth_7 | 0.0194 | 0.005 | 3.753 | 0.000 | 0.009 | 0.030 |
| max_depth_8 | 0.0196 | 0.005 | 3.733 | 0.000 | 0.009 | 0.030 |
| max_depth_9 | 0.0266 | 0.005 | 5.176 | 0.000 | 0.017 | 0.037 |
| max_depth_10 | 0.0307 | 0.005 | 5.954 | 0.000 | 0.021 | 0.041 |
| max_depth_11 | 0.0285 | 0.005 | 5.484 | 0.000 | 0.018 | 0.039 |
| max_depth_12 | 0.0312 | 0.005 | 6.178 | 0.000 | 0.021 | 0.041 |
| max_depth_13 | 0.0320 | 0.005 | 6.218 | 0.000 | 0.022 | 0.042 |
| colsample_bytree_0.8 | -0.0112 | 0.003 | -4.143 | 0.000 | -0.017 | -0.006 |
| colsample_bytree_1.0 | -0.0278 | 0.003 | -8.388 | 0.000 | -0.034 | -0.021 |
| min_gain_to_split_1 | -0.0009 | 0.003 | -0.307 | 0.759 | -0.006 | 0.005 |
| min_gain_to_split_5 | -0.0016 | 0.002 | -0.726 | 0.468 | -0.006 | 0.003 |
| ===== | | | | | | |
| Omnibus: | 11.763 | Durbin-Watson: | | 0.856 | | |
| Prob(Omnibus): | 0.003 | Jarque-Bera (JB): | | 11.104 | | |
| Skew: | -0.361 | Prob(JB): | | 0.00388 | | |
| Kurtosis: | 2.609 | Cond. No. | | 17.1 | | |
| ===== | | | | | | |

Warnings:

[1] Standard Errors are heteroscedasticity robust (HC3)

На практике проникновение в суть того, каким образом модели приходят к предсказаниям, является чрезвычайно важным, в особенности для инвестиционных стратегий, где лица, принимающие решения, часто требуют правдоподобных объяснений.

Как интерпретировать результаты GBM

Понимать, почему модель предсказывает определенный результат, очень важно по нескольким причинам, включая доверие, применимость, отчетность и отладку. Понимание нелинейной связи между признаками и результатом, обнаруженным моделью, а также взаимодействий между признаками также имеет ценность, когда цель состоит в том, чтобы узнать больше о базовых движущих силах изучаемого явления.

Общий подход к пониманию сущности предсказаний, делаемых методами древесного ансамбля, такими как градиентно-бустинговые модели или модели случайного леса, заключается в том, чтобы атрибутировать значения важности признаков каждой входной переменной. Эти значения важности признаков могут быть вычислены на индивидуальной основе для единственного предсказания либо глобально для всего набора данных целиком (т. е. для всех образцов), тем самым получая более высокий уровень перспективы того, как модель делает предсказания.

Свойство важности признаков

Существует три основных способа вычисления значений *глобальной важности признаков*.

- ◆ *Прирост* — этот классический подход, введенный Лео Брейманом в 1984 г., использует общее снижение потери или смешанности, вносимых всеми разбиивками для определенного признака. Мотивация в значительной степени является эвристической, но этот метод широко используется для отбора признаков.
- ◆ *Количество разбиивок* — это альтернативный подход, который подсчитывает, как часто признак используется для принятия решения о разбиивке, опираясь на отбор признаков для этой цели на основе результирующего прироста информации.
- ◆ *Пермутация* — этот подход случайно перетасовывает значения признаков в тестовом наборе и служит мерой величины изменения модельной ошибки, исходя из того, что важный признак должен создавать крупное увеличение ошибки предсказания. Разные варианты пермутации приводят к альтернативным реализациям этого базового подхода.

Индивидуализированные значения важности признаков, которые вычисляют релевантность признаков для единственного предсказания, встречаются реже, поскольку имеющиеся методы модельно-независимого объяснения намного медленнее, чем древесно-специфические методы.

Все реализации градиентного бустинга предоставляют отметки важности признаков после тренировки в качестве модельного атрибута. Библиотека XGBoost предоставляет три версии:

- ◆ суммарный прирост `total_gain` и прирост `gain` как его среднее в расчете на разбиивку;
- ◆ суммарное покрытие `total_cover` как число образцов в расчете на разбиивку при использовании признака;
- ◆ вес `weight` как количество разбиивок от предыдущих значений.

Эти значения доступны с помощью метода `get_score()` натренированной модели с соответствующим параметром `importance_type`. Для наиболее результативной модели XGBoost результаты будут такими, как показано на рис. 11.9 (общие меры имеют корреляцию 0,8, как и меры `cover` и `total_cover`).

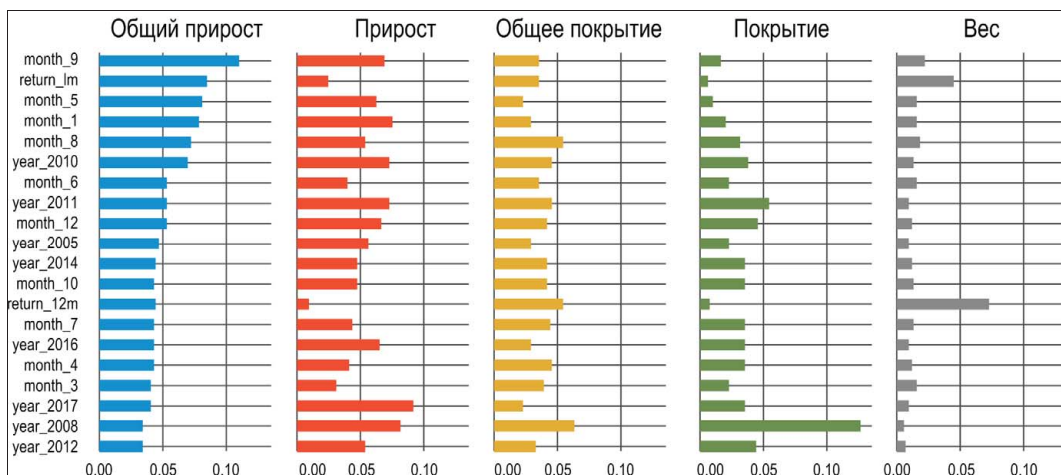


Рис. 11.9. Оценки важности признаков после тренировки модели

Хотя индикаторы для разных месяцев и лет доминируют, самый недавний финансовый возврат за 1 месяц является вторым по важности признаком с точки зрения меры `total_gain` и часто используется в соответствии с мерой `weight`, но дает низкие средние приросты, поскольку он применяется к относительно малому числу экземпляров в среднем (подробности реализации см. в блокноте).

Графики частичной зависимости

Помимо суммарного вклада отдельных признаков в модельное предсказание, графики частичной зависимости визуализируют связь между целевой переменной и множеством признаков. Нелинейная природа градиентно-бустинговых деревьев побуждает эту связь зависеть от значений всех других признаков. Следовательно, мы будем маргинализировать эти признаки. Благодаря этому мы можем интерпретировать частичную зависимость как ожидаемый целевой отклик.

Мы можем визуализировать частичную зависимость только для индивидуальных признаков или пар признаков. Последнее приводит к контурным графикам, которые показывают, как комбинации значений признаков производят разные предсказанные вероятности. Это показано в следующем фрагменте кода:

```
fig, axes = plot_partial_dependence(gbrt=best_model,
                                   X=X,
                                   features=['month_9', 'return_1m', 'return_3m',
                                             ('return_1m', 'return_3m')],
                                   feature_names=['month_9', 'return_1m', 'return_3m'],
                                   percentiles=(0.01, 0.99),
                                   n_jobs=-1,
                                   n_cols=2,
                                   grid_resolution=250)

fig.suptitle('График частичной зависимости', fontsize=18)
```

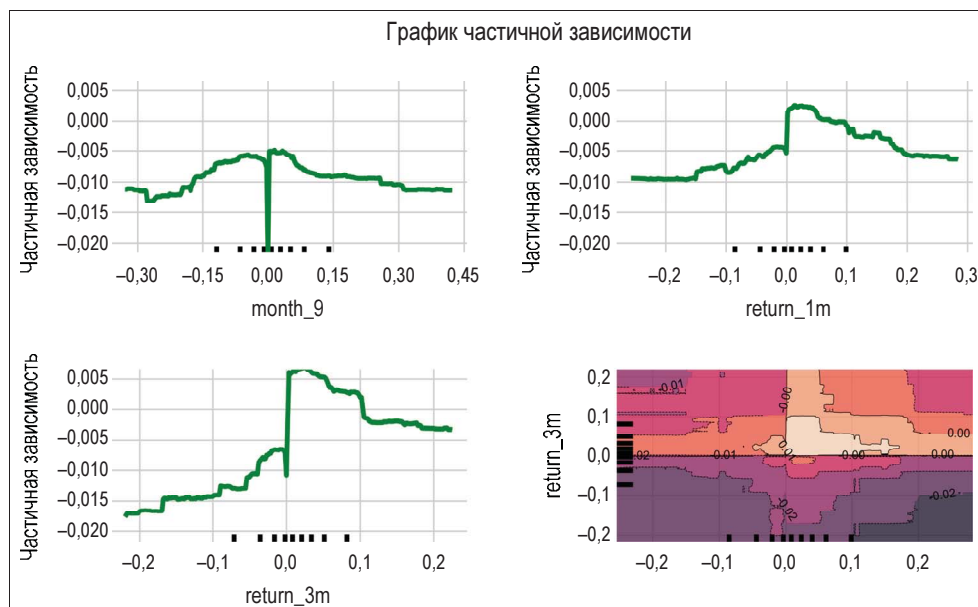


Рис. 11.10. График частичной зависимости

После некоторого дополнительного форматирования (см. сопутствующий блокнот) мы получаем график (рис. 11.10).

Правый нижний график показывает зависимость вероятности положительного финансового возврата в течение следующего месяца с учетом диапазона значений для сдвинутых одномесячного и трехмесячного возвратов после устранения выбросов в процентилях [1%; 99%]. Переменная `month_9` является фиктивной переменной, отсюда и график в виде ступенчатой (кусочно-постоянной) функции. Мы также можем визуализировать зависимость в трехмерном измерении, как показано в следующем фрагменте кода:

```
targets = ['return_1m', 'return_3m']
target_feature = [X.columns.get_loc(t) for t in targets]
pdp, axes = partial_dependence(best_model,
                               target_feature,
                               X=X,
                               grid_resolution=100)

XX, YY = np.meshgrid(axes[0], axes[1])
Z = pdp[0].reshape(list(map(np.size, axes))).T

fig = plt.figure(figsize=(14, 8))
ax = Axes3D(fig)
surf = ax.plot_surface(XX, YY, Z,
                      rstride=1,
                      cstride=1,
```

```

cmap=plt.cm.BuPu,
edgecolor='k')
ax.set_xlabel(' '.join(targets[0].split('_')).capitalize())
ax.set_ylabel(' '.join(targets[1].split('_')).capitalize())
ax.set_zlabel('Частичная зависимость')
ax.view_init(elev=22, azimuth=300)

fig.colorbar(surf)
fig.suptitle('Частичная зависимость по 1- и 3-месячному возвратам')

```

В результате получим трехмерный график частичной зависимости направления одномесячного финансового возврата от сдвинутых одномесячного и трехмесячного финансовых возвратов (рис. 11.11).

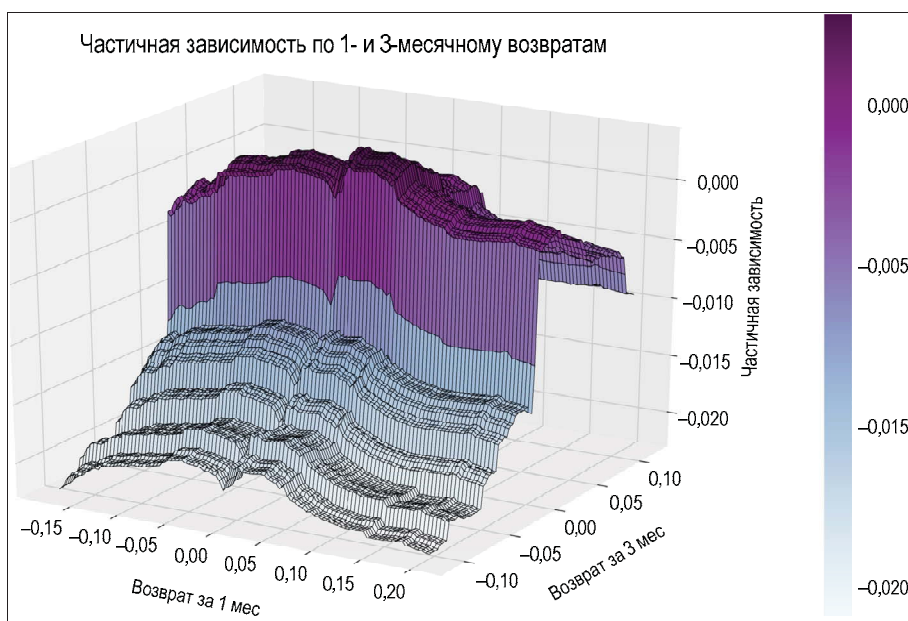


Рис. 11.11. Частичная зависимость по 1- и 3-месячному финансовым возвратам

Аддитивные объяснения Шепли (SHapley)

На конференции NIPS 2017 г. Скотт Лундберг (Scott Lundberg) и Су-Ин Ли (Su-In Lee) из Вашингтонского университета представили новый и более точный подход к объяснению вклада индивидуальных признаков в выход древесных ансамблевых моделей, именуемый аддитивными объяснениями Шепли, или значениями SHAP (SHapley Additive exPlanations, SHapley).

Этот новый алгоритм отходит от наблюдения, что методы атрибуирования признаков для древесных ансамблей, такие как те, которые мы рассматривали ранее, про-

тиворечивы, т. е. изменение в модели, которое увеличивает влияние признака на выход, может снизить значения важности для этого признака (справочные материалы с подробными иллюстрациями см. в репозитории GitHub).

Значения SHAP унифицируют идеи из коллаборативной теории игр и локальных объяснений и зарекомендовали себя как теоретически оптимальные, непротиворечивые и локально точные на основе ожиданий. Самое главное, что Лундберг и Ли разработали алгоритм, который позволяет снижать сложность вычисления этих модельно-агностических, аддитивных методов атрибуции признаков с $O(TLD^M)$ до $O(TLD^2)$, где T и M — это соответственно число деревьев и признаков, а D и L — максимальная глубина и число листьев в деревьях. Это важное нововведение позволяет объяснять предсказания из ранее не поддающихся отслеживанию моделей с тысячами деревьев и признаков за долю секунды. Реализация с открытым исходным кодом стала доступна в конце 2017 г. и совместима с древесными моделями XGBoost, LightGBM, CatBoost и sklearn.

Значения Шепли возникли в теории игр как техническое решение для закрепления значения за каждым игроком в совместной игре, которое отражает их вклад в командный успех. Значения SHAP являются адаптацией концепции теории игр к древесным моделям и рассчитываются для каждого признака и каждого образца. Они измеряют величину вносимого признаком вклада в модельный выход для определенного наблюдения. По этой причине значения SHAP дают дифференцированное представление о том, как влияние признака варьируется по разным образцам, что важно с учетом роли эффектов взаимодействия в этих нелинейных моделях.

Как резюмировать значения SHAP по признакам

Существует два способа построения графика значений SHAP для получения высокоуровневого обзора важности признаков по нескольким образцам: простое среднее значение по всем образцам, которое напоминает глобальные меры важности признаков, вычисленные ранее (как показано на левой панели следующего ниже рис. 11.12), либо график рассеяния, показывающий влияние каждого признака для каждого образца (как показано на правой панели указанного рисунка). Они очень легко создаются с использованием натренированной модели совместимой библиотеки и соответствующих входных данных, как показано в следующем фрагменте кода:

```
# Загрузить визуализационный код JS в блокнот
shap.initjs()

# Объяснить модельные предсказания, используя значения SHAP
explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test)

shap.summary_plot(shap_values, X_test, show=False)
```

График рассеяния справа на рис. 11.12 сортирует признаки по их суммарным значениям SHAP по всем образцам, а затем показывает, как каждый признак влияет на модельный выход, измеряемый значением SHAP, как функция от значения признака, представленного его цветом, где красный представляет высокие значения и синий — низкие значения относительно диапазона признака.

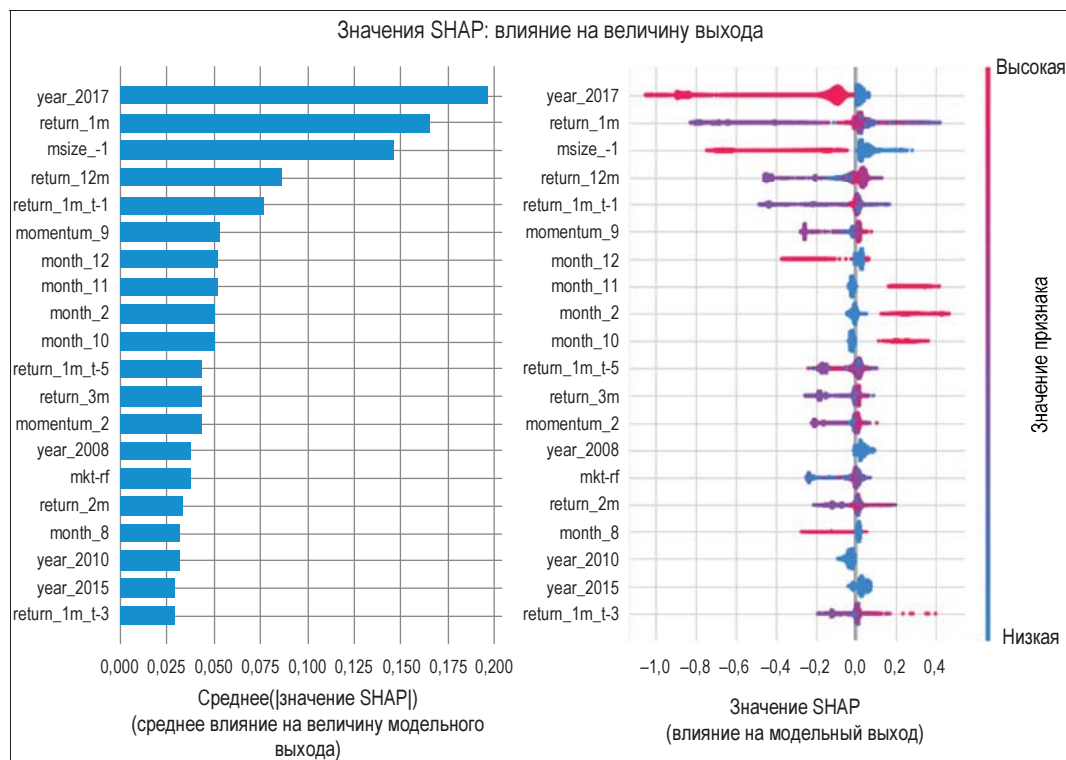


Рис. 11.12. Значения SHAP: влияние на величину модельного выхода

Как использовать графики силы для объяснения предсказания

На графике силы (см. рис. 11.12) показано кумулятивное влияние различных признаков и их значений на модельный выход, которые в данном случае составил 0,6, что значительно выше базового значения 0,13 (средний модельный выход на предоставленном наборе данных). Признаки, выделенные красным цветом, увеличивают модельный выход. Октябрь является наиболее важным признаком и увеличивает выход с 0,338 до 0,537, в то время как 2017 г. снижает выход.

Таким образом, мы получаем подробную разбивку того, как модель пришла к определенному предсказанию, как показано на рис. 11.13.

Мы также можем вычислять графики силы для многочисленных точек данных или предсказаний одновременно и использовать кластеризованную визуализацию, про-

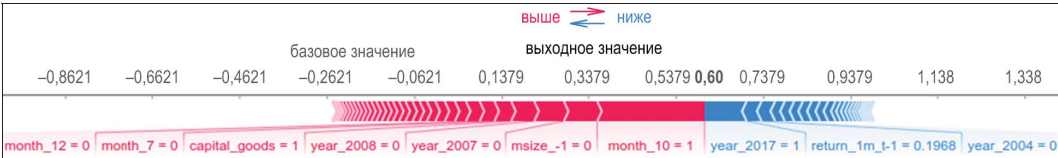


Рис. 11.13. Разбивка процесса прихода модели к определенному предсказанию

никакая в сущность того, насколько распространены определенные регулярности влияния в наборе данных. На рис. 11.14 показаны графики силы для первых 1000 наблюдений, повернутых на 90°, уложенных горизонтально и упорядоченных по влиянию разных признаков на результат для определенного наблюдения. В данной реализации для выявления этих регулярностей используется иерархическая агломеративная кластеризация точек данных по значениям SHAP признака и показывается результат в интерактивном режиме для разведывательного анализа (см. блокнот), как показано в следующем фрагменте кода:

```
shap.force_plot(explainer.expected_value, shap_values[:1000,:], X_test.iloc[:1000])
```

Он производит следующий выход (рис. 11.14).

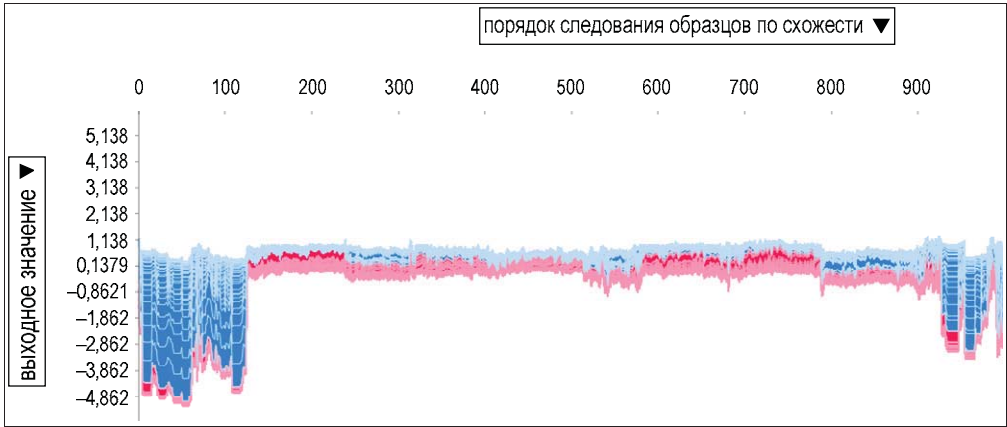


Рис. 11.14. Иерархическая агломеративная кластеризация точек данных по значениям SHAP признака с целью выявления регулярностей влияния

Как анализировать взаимодействие признаков

Наконец, значения SHAP позволяют дополнительно понять эффекты взаимодействия между разными признаками путем отделения этих взаимодействий от главных эффектов. График зависимости `shap.dependence_plot` может быть определен следующим образом:

```
shap.dependence_plot("return_1m", shap_values, X_test,
                    interaction_index=2,
                    title=Взаимодействие между 1- и 3-месячными возвратами')
```


Он показывает, как разные значения для одномесячных финансовых возвратов (на оси x) влияют на результат (значение SHAP на оси y) с исчислением последовательных разностей по трехмесячным возвратам (рис. 11.15).

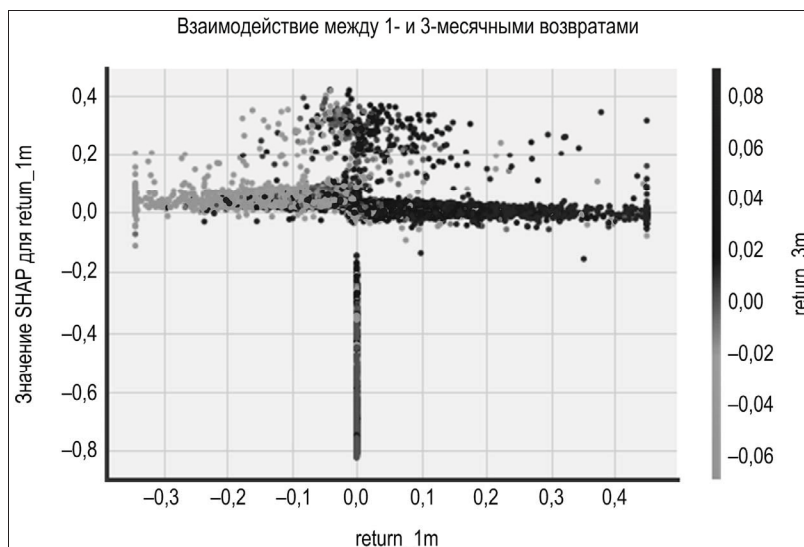


Рис. 11.15. Влияние одномесячных финансовых возвратов на результат с исчислением последовательных разностей по трехмесячным возвратам

Значения SHAP обеспечивают гранулярную атрибуцию признаков на уровне каждого отдельного предсказания и позволяют значительно расширять возможности обследования сложных моделей с помощью (интерактивной) визуализации. Сводный график рассеяния значений SHAP, показанный в начале этого раздела, предлагает гораздо более дифференцированные сущностные представления, чем глобальная гистограмма важности признаков. Силовые графики индивидуальных кластеризованных предсказаний позволяют проводить более подробный анализ, в то время как графики зависимости SHAP улавливают эффекты взаимодействия и, как следствие, дают более точные и подробные результаты, чем графики частичной зависимости.

Ограничения значений SHAP, как и в случае любой текущей меры важности признаков, касаются атрибуции влияния переменных, которые сильно коррелируют, поскольку их схожее влияние может быть подразделено произвольным образом.

Резюме

В этой главе мы рассмотрели алгоритм градиентного бустинга, который используется для построения ансамблей в последовательном стиле, добавляя мелкое дерево решений, которое для улучшения сделанных прогнозов использует только очень

небольшое число признаков. Мы увидели, как градиентно-бустинговые деревья могут очень гибко применяться к широкому спектру функций потерь и предлагают большое число возможностей для настройки модели к заданному набору данных и задаче автоматического обучения.

Недавние реализации во многом способствовали использованию градиентного бустинга, ускорив тренировочный процесс и предложив непротиворечивое и подробное проникновение в сущность важности признаков и движущих сил индивидуальных предсказаний. В следующей главе мы обратимся к подходам к МО на основе неконтролируемого автоматического обучения.

12

Неконтролируемое обучение

В *главе 6* мы рассмотрели, как неконтролируемое обучение увеличивает добавочную стоимость, раскрывая структуры в данных в отсутствии выходной переменной, такой как учитель, который направляет процесс поиска. Эта задача контрастирует со средой контролируемого обучения, на которой мы были сосредоточены в нескольких последних главах.

Неконтролируемые обучающиеся алгоритмы могут быть полезны, когда набор данных содержит только признаки и никаких измерений результата, либо когда мы хотим извлечь информацию независимо от результата. Вместо предсказания будущих результатов цель состоит в изучении информативного представления данных, которое будет полезным для решения другой задачи, в том числе разведывание набора данных.

Примеры охватывают выявление тематик с целью аннотирования документов (*см. главу 14*), уменьшения числа признаков с целью снижения риска переподгонки и вычислительных затрат на контролируемое обучение либо кластеризацию схожих наблюдений, как показано на примере использования кластеризации для размещения финансовых средств среди портфельных активов в конце этой главы.

Снижение размерности и кластеризация являются для неконтролируемого обучения главными задачами.

- ◆ *Снижение размерности* преобразовывает существующие признаки в новый, меньший набор при минимизации потери информации. Существует широкий спектр алгоритмов, отличающихся только тем, как они измеряют потерю информации, применяют ли они линейные или нелинейные преобразования или ограничения, которые они накладывают на новый набор признаков.
- ◆ *Алгоритмы кластеризации* вместо выявления новых признаков выявляют и группируют схожие наблюдения или признаки. Алгоритмы различаются тем, как они определяют сходство наблюдений, и своими допущениями о результирующих группах.

В частности, в этой главе рассматриваются следующие вопросы:

- ◆ как *анализ главных компонент* (Principal Component Analysis, PCA) и *анализ независимых компонент* (Independent Component Analysis, ICA) выполняют линейное снижение размерности;
- ◆ как применять PCA для выявления рискованных факторов и собственных портфелей из возвратов от финансовых активов;
- ◆ как применять усвоение проекций в нелинейном топологическом многообразии, которое резюмирует высокоразмерные данные, для эффективной визуализации;
- ◆ как применять алгоритмы t-SNE и UMAP для разведывания высокоразмерных альтернативных графических данных;
- ◆ как работают алгоритмы k средних, иерархической и плотностной кластеризации;
- ◆ как использовать агломеративную кластеризацию для построения робастных портфелей в соответствии с иерархическим паритетом риска.



Примеры исходного кода для каждого раздела находятся в каталоге онлайн-нового репозитория GitHub для этой главы по адресу <https://github.com/PacktPublishing/Hands-On-Machine-Learning-for-AlgorithmicTrading>.

Снижение размерности

С точки зрения линейной алгебры признаки набора данных создают векторное пространство, размерность которого соответствует числу линейно независимых столбцов (при условии, что наблюдений больше, чем признаков). Два столбца являются линейно зависимыми, когда они идеально коррелированы, вследствие чего один может быть вычислен из другого с помощью линейных операций сложения и умножения.

Другими словами, они являются параллельными векторами, которые представляют скорее одинаковые, чем разные направления или оси, и составляют только одну размерность. Схожим образом, если одна переменная является линейной комбинацией нескольких других, то она является элементом векторного пространства, скорее созданного этими столбцами, чем добавлением собственной размерности.

Число размерностей набора данных имеет важное значение, поскольку каждая новая размерность может добавлять сигнал относительно результата. Однако есть и обратная сторона, именуемая *проклятием размерности*: по мере того, как число независимых признаков растет, а число наблюдений остается постоянным, среднее расстояние между точками данных также растет, а плотность признакового пространства падает экспоненциально.

Последствия для автоматического обучения драматичны, т. к. предсказывать становится намного сложнее, когда наблюдения более отдалены; т. е. отличаются друг от

друга. В следующем далее разделе рассматриваются возникающие в связи с этим проблемы.

Снижение размерности призвано эффективнее представлять информацию в данных за счет использования меньшего числа признаков. С этой целью алгоритмы проецируют данные в низкоразмерное пространство, при этом отбрасывая изменчивость в данных, которая является неинформативной, либо выявляя низкоразмерное подпространство или топологическое многообразие, в котором или вблизи которого обитают данные.

Топологическое многообразие, или многосвязная область, — это пространство, локально напоминающее евклидово пространство¹. Одномерные топологические многообразия включают линии и круги (но не снимки восьмимерного многообразия из-за точки пересечения). Гипотеза о топологическом многообразии утверждает, что высокоразмерные данные часто находятся в низкоразмерном пространстве, которое, если его идентифицировать, позволяет точно представлять данные в этом подпространстве.

Таким образом, снижение размерности сжимает данные, находя другое, меньшее множество переменных, которые улавливают то, что является наиболее важным в исходных признаках, минимизируя потерю информации. Сжатие помогает противостоять проклятию размерности, экономит память и позволяет визуализировать существенные аспекты высокоразмерных данных, которые в иных случаях очень трудно разведывать.

Линейные и нелинейные алгоритмы

Алгоритмы снижения размерности отличаются ограничениями, которые они накладывают на новые переменные, и тем, как они призваны минимизировать потерю информации.

♦ *Линейные алгоритмы*, такие как анализ главных компонент (PCA) и анализ независимых компонент (ICA), ограничивают новые переменные линейными комбинациями исходных признаков, т. е. гиперплоскостями в более низкоразмерном пространстве. В то время как PCA требует, чтобы новые признаки были некоррелированными, ICA идет дальше и навязывает статистическую независимость — отсутствие как линейных, так и нелинейных связей. На рис. 12.1 показано, как PCA проецирует трехмерные признаки в двумерное пространство.

¹ Топологическое многообразие (manifold), или многосвязная область, — это топологическое пространство, которое является локально евклидовым (т. е. вокруг каждой точки есть окрестности, топологически такие же, как разомкнутый единичный шар в \mathbb{R}^n). Ближайшей аналогией является вера древних людей, что Земля — плоская. Это расхождение возникает главным образом из-за того, что в малых масштабах, которые мы видим, Земля действительно выглядит плоской, и топологические многообразия составляют обобщение объектов, на которых мы могли бы жить и сталкиваться с проблемой круглой/плоской Земли. <http://mathworld.wolfram.com/Manifold.html>. — Прим. перев.

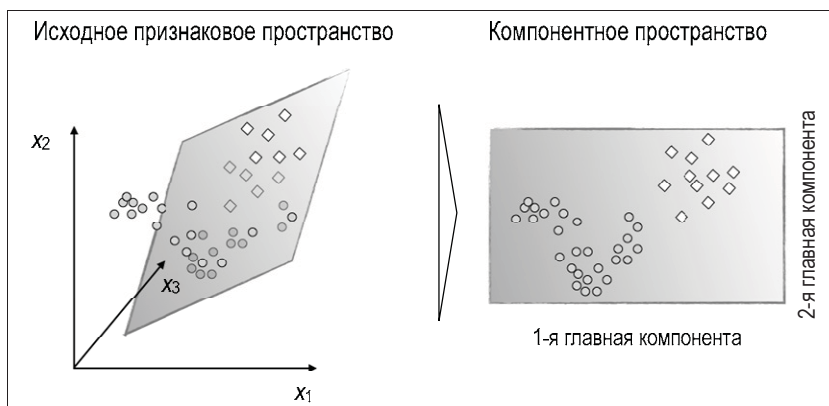


Рис. 12.1. Проецирование методом PCA трехмерных признаков в двумерное пространство

- ◆ *Нелинейные алгоритмы* не ограничиваются гиперплоскостями и могут улавливать более сложную структуру в данных. Однако, учитывая бесконечное число вариантов, для прихода к решению алгоритмы по-прежнему должны принимать допущения. В этом разделе мы покажем, в чем заключается огромная польза от алгоритмов *t-распределенного стохастического вложения соседей* (t-distributed Stochastic Neighbor Embedding, t-SNE) и *равномерной аппроксимации и проекции топологического многообразия* (Uniform Manifold Approximation and Projection, UMAP) для визуализирования высокоразмерных данных. На рис. 12.2 показано, как процедура усвоения проекции в топологическом многообразии выявляет двумерное подпространство в трехмерном признаковом пространстве (блокнот `manifold_learning.ipynb` иллюстрирует использование дополнительных алгоритмов, включая локальное линейное вложение).

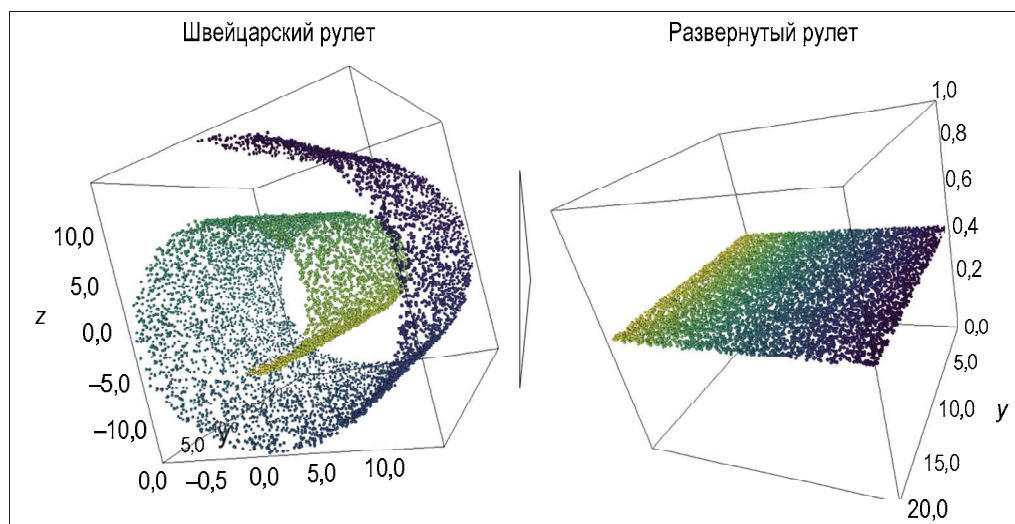


Рис. 12.2. Усвоение топологического многообразия выявляет двумерное подпространство в трехмерном признаковом пространстве

Проклятие размерности

Увеличение числа размерностей набора данных означает, что в векторе признаков, представляющем каждое наблюдение в соответствующем евклидовом пространстве, имеется больше элементов. Расстояние в векторном пространстве измеряется евклидовым расстоянием, также именуемым *нормой* L_2 , которое мы применили к вектору линейно-регрессионных коэффициентов, тренируя регуляризованную модель гребневой регрессии.

Евклидово расстояние между двумя n -мерными векторами с декартовыми координатами $\mathbf{p} = [p_1 \ p_2 \ \dots \ p_n]$ и $\mathbf{q} = [q_1 \ q_2 \ \dots \ q_n]$ вычисляется по хорошо известной формуле, разработанной Пифагором:

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}.$$

Следовательно, каждая новая размерность добавляет в сумму неотрицательный член, вследствие чего для разных векторов расстояние увеличивается вместе с числом размерностей. Другими словами, по мере роста числа признаков для данного числа наблюдений пространство признаков становится все более разреженным, т. е. менее плотным или более пустым. С другой стороны, для поддержания среднего расстояния между точками данных одинаковым более низкая плотность данных требует больше наблюдений.

На рис. 12.3 показано, сколько точек данных нужно для поддержания среднего расстояния 10 наблюдений равномерно распределенными в линейной пропорции. Оно увеличивается экспоненциально от 10^1 в одной размерности до 10^2 в двух и 10^3 в трех размерностях, поскольку данные должны расширяться в 10 раз каждый раз, когда мы добавляем новую размерность.

Блокнот `curse_of_dimensionality.ipynb` в папке репозитория GitHub для этого раздела симулирует увеличение среднего и минимального расстояний между точками данных по мере увеличения числа размерностей (рис. 12.4).

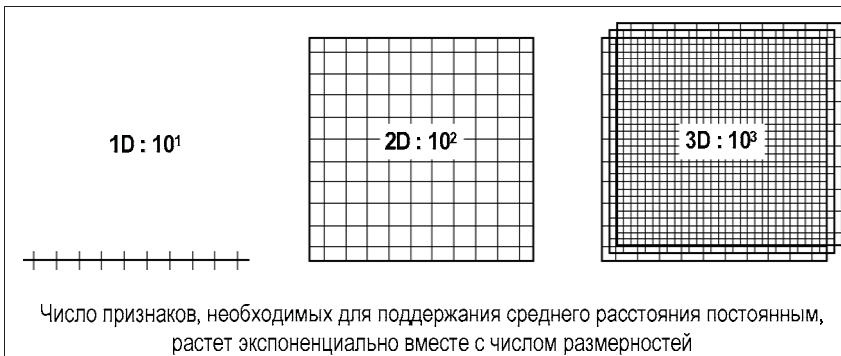


Рис. 12.3. Число точек, необходимых для поддержания среднего расстояния 10 наблюдений линейно распределенным



Рис. 12.4. Симуляция увеличения среднего и минимального расстояний между точками данных по мере увеличения числа размерностей

Указанная симуляция извлекает признаки в промежутке $[0; 1]$ из некоррелированных равномерных или коррелированных нормальных распределений и постепенно увеличивает число признаков до 2500. Среднее расстояние между точками данных увеличивается более чем в 11 раз по сравнению с признаковым диапазоном для признаков, извлеченных из нормального распределения, и более чем в 20 раз (экстремальном) случае некоррелированного равномерного распределения.

Когда расстояние между наблюдениями увеличивается, контролируемое автоматическое обучение становится труднее, поскольку предсказания для новых образцов с меньшей вероятностью будут основаны на усвоении похожих тренировочных признаков. Иными словами, число возможных уникальных строк растет экспоненциально по мере увеличения числа признаков, что значительно затрудняет эффективную выборку из такого пространства.

Точно так же сложность функций, усваиваемых гибкими алгоритмами, которые делают меньше допущений о фактической связи, растет экспоненциально вместе с числом размерностей.

Гибкие алгоритмы включают в себя древесные модели, которые мы встречали в главах 10 и 11, а также глубокие нейронные сети, которые мы рассмотрим в главе 17. Дисперсия этих алгоритмов увеличивается, поскольку они получают больше возможностей для подгонки к шуму в большем числе размерностей, что приводит к низкой результативности обобщения.

На практике признаки коррелируют, часто существенно, или не показывают большой вариативности. По этим причинам процедура снижения размерности помогает сжимать данные, не теряя большей части сигнала, и бороться с проклятием, а также экономить память. В этих случаях она дополняет собой использование регуляризации, управляя ошибкой предсказания из-за дисперсии и модельной сложности.

Тогда возникает резонный вопрос, который мы рассмотрим в следующем далее разделе: каковы наилучшие способы отыскания более низкоразмерного представления данных, которое теряет как можно меньше информации?

Линейное снижение размерности

Алгоритмы линейного снижения размерности вычисляют линейные комбинации, которые преобразуют, поворачивают и масштабируют исходные признаки, улавливая значительные вариации данных с учетом ограничений на характеристики новых признаков.

Анализ главных компонент (Principal Component Analysis, PCA), изобретенный в 1901 г. Карлом Пирсоном (Karl Pearson), находит новые признаки, которые отражают направления максимальной дисперсии данных, будучи взаимно некоррелированными или ортогональными.

Анализ независимых компонент (Independent Component Analysis, ICA), напротив, возник в сфере обработки сигналов в 1980-х годах с целью разделения разных сигналов при наложении более сильного ограничения статистической независимости.

В данном разделе представлены эти два алгоритма, а затем показано, как применять алгоритм PCA к возвратам от финансовых активов с целью усвоения рисков факторов из данных и построения так называемых характеристических (eigen) портфелей для систематических стратегий торговли.

Анализ главных компонент

Анализ главных компонент (PCA) находит главные компоненты как линейные комбинации существующих признаков и использует эти компоненты для представления исходных данных. Число компонент является гиперпараметром, которое определяет целевую размерность, и должно быть равным или меньшим числу наблюдений или столбцов в зависимости от того, что меньше.

Алгоритм PCA призван улавливать большую часть дисперсии в данных для того, упрощая восстановление исходных признаков, и так, чтобы каждая компонента добавляла информацию. Это приводит к снижению размерности, проецируя исходные данные в пространство главных компонент.

Алгоритм PCA работает путем выявления последовательности главных компонент, каждая из которых выравнивается с направлением максимальной дисперсии в данных после учета изменчивости, захваченной ранее вычисленными компонентами. Последовательная оптимизация также гарантирует, что новые компоненты не коррелируют с существующими компонентами, вследствие чего результирующее множество представляет собой ортогональный базис для векторного пространства.

Этот новый базис соответствует повернутой версии исходного базиса, вследствие чего новая ось указывает в направлении последовательно уменьшающейся дисперсии. Снижение величины дисперсии исходных данных, объясняемое каждой главной компонентой, отражает степень корреляции между исходными признаками.

Число компонент, которые улавливают, например, 95% исходной изменчивости относительно общего числа признаков, позволяет проникнуть в сущность линейно независимой информации в исходных данных.

Визуализирование PCA в двух размерностях

На рис. 12.5 показаны несколько аспектов PCA для двумерного случайного набора данных (см. блокнот `pca_key_ideas.ipynb`).

- ◆ На рис. 12.5 слева показано, как первая и вторая главные компоненты совпадают с направлениями максимальной дисперсии, будучи ортогональными.
- ◆ На рис. 12.5 в центре показано, как первая главная компонента минимизирует ошибку восстановления, измеряемую как сумма расстояний между точками данных и новой осью.
- ◆ Наконец, на рис. 12.5 справа иллюстрируется контролируемый алгоритм обычных наименьших квадратов (OLS), который аппроксимирует результирующую переменную (здесь мы выбираем x_2) (одномерной) гиперплоскостью, вычисленной из (единственного) признака. Вертикальные линии подчеркивают, как OLS минимизирует расстояние вдоль оси результата, в отличие от PCA, который минимизирует расстояния, ортогональные гиперплоскости.

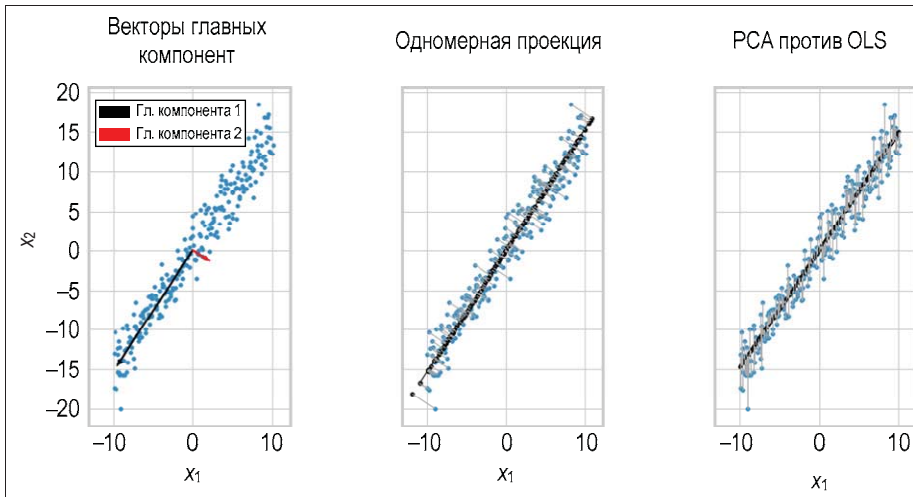


Рис. 12.5. Аспекты PCA для двумерного случайного набора данных

Допущения, принимаемые в анализе главных компонент

В анализе главных компонент делается несколько допущений, которые важно иметь в виду. К ним относятся следующие:

- ◆ высокая дисперсия приводит к высокому соотношению между сигналом и шумом;
- ◆ данные стандартизированы, вследствие чего дисперсия сопоставима по всем признакам;
- ◆ линейные преобразования улавливают релевантные аспекты данных;

- ◆ статистические показатели более высокого порядка за пределами первого и второго момента не имеют значения, из чего следует, что данные имеют нормальное распределение.

Акцент на первом и втором моментах согласуется со стандартными метриками риска/возвратности, но допущение о нормальности может противоречить характеристикам рыночных данных.

Как работает алгоритм PCA

Указанный алгоритм находит векторы, создавая гиперплоскость целевой размерности, которая минимизирует ошибку восстановления, измеряемую как сумму квадратов расстояний точек данных до плоскости. Как показано выше, эта цель соответствует нахождению последовательности векторов, которые совпадают с направлениями максимальной захваченной дисперсии с учетом других компонент, обеспечивая при этом взаимную ортогональность всех главных компонент.

На практике данный алгоритм решает задачу либо путем вычисления собственных векторов матрицы ковариаций, либо с помощью сингулярного разложения.

Мы проиллюстрируем его вычисление с помощью случайно сгенерированного трехмерного эллипса со 100 точками данных, показанного на рис. 12.6 слева, включая двумерную гиперплоскость, определенную первыми двумя главными компонентами (примеры исходного кода см. в блокноте `the_math_behind_pca.ipynb`).

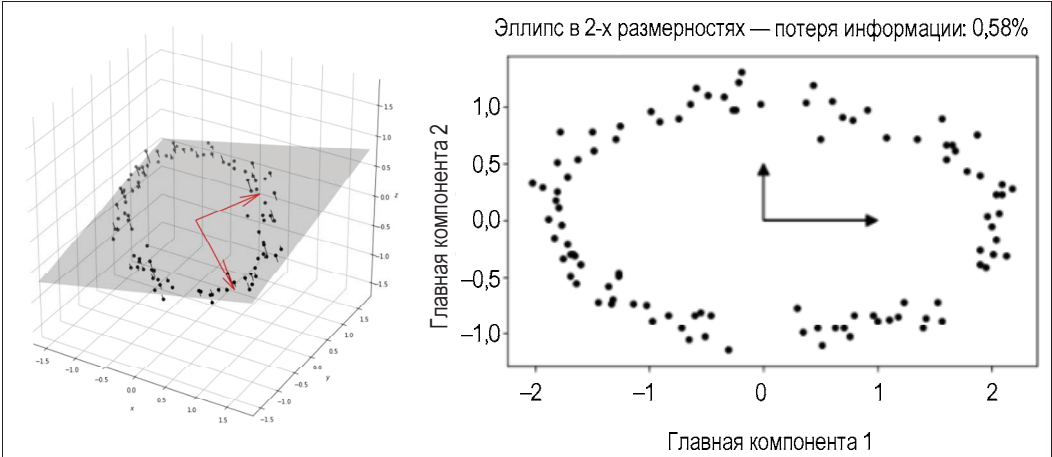


Рис. 12.6. Трехмерный эллипс и двумерная гиперплоскость

PCA на основе матрицы ковариаций

Сначала мы вычисляем главные компоненты, используя квадратную матрицу ковариаций с попарными ковариациями образцов для признаков x_i , x_j , $i, j = 1, \dots, n$ в качестве элементов в строке i и столбце j :

$$\text{cov}_{i,j} = \frac{\sum_{k=1}^N (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{N-1}.$$

Для квадратной матрицы \mathbf{M} из n размерностей мы определяем собственные векторы \mathbf{w}_i и собственные значения λ_i , $i = 1, \dots, n$ следующим образом:

$$\mathbf{M}\mathbf{w}_i = \lambda_i \mathbf{w}_i.$$

Следовательно, мы можем представить матрицу \mathbf{M} , используя собственные векторы и собственные значения, где \mathbf{W} — это матрица, содержащая собственные векторы в качестве векторов-столбцов, а \mathbf{L} — матрица, содержащая λ_i в качестве диагональных элементов (и нули в противном случае). Мы определяем спектральное разложение матрицы:

$$\mathbf{M} = \mathbf{W}\mathbf{L}\mathbf{W}^{-1}.$$

С помощью библиотеки NumPy мы реализуем это следующим образом, где кадр данных pandas содержит 100 точек данных эллипса:

```
# вычислить матрицу ковариаций:
cov = np.cov(data.T) # каждая строка представляет признак
cov.shape

(3, 3)
```

Далее вычисляются собственные векторы и собственные значения матрицы ковариаций. Собственные векторы содержат главные компоненты (где знак является произвольным):

```
eigen_values, eigen_vectors = eig(cov)
eigen_vectors

array([[ 0.71409739, -0.66929454, -0.20520656],
       [-0.70000234, -0.68597301, -0.1985894 ],
       [ 0.00785136, -0.28545725,  0.95835928]])
```

Мы можем сравнить результат с результатом, полученным из библиотеки sklearn, и обнаружим, что они совпадают в абсолютном выражении:

```
pca = PCA()
pca.fit(data)
C = pca.components_.T # столбцы = главные компоненты
C

array([[ 0.71409739,  0.66929454,  0.20520656],
       [-0.70000234,  0.68597301,  0.1985894 ],
       [ 0.00785136,  0.28545725, -0.95835928]])

np.allclose(np.abs(C), np.abs(eigen_vectors))

True
```

Мы также можем проверить спектральное разложение матрицы, начиная с диагональной матрицы **L**, содержащей собственные значения:

```
# матрица собственных значений
ev = np.zeros((3, 3))
np.fill_diagonal(ev, eigen_values)
ev # диагональная матрица

array([[1.92923132, 0.          , 0.          ],
       [0.          , 0.55811089, 0.          ],
       [0.          , 0.          , 0.00581353]])
```

Мы видим, что результат действительно соблюдается:

```
decomposition = eigen_vectors.dot(ev).dot(inv(eigen_vectors))
np.allclose(cov, decomposition)
```

True

РСА с использованием сингулярного разложения

Далее мы рассмотрим альтернативное вычисление с использованием *сингулярного разложения* (singular value decomposition, SVD). Этот алгоритм медленнее, когда число наблюдений больше, чем число признаков (типичный случай), но дает более надежную числовую стабильность, в особенности когда некоторые признаки сильно коррелируют (что нередко является изначальной причиной использования РСА).

Сингулярное разложение обобщает спектральное разложение матрицы, которую мы только что применили к квадратной и симметричной матрице ковариаций до более общего случая прямоугольных $m \times n$ -матриц. Оно имеет форму, показанную на рис. 12.7 в центре. Диагональные значения Σ являются сингулярными значениями, а транспозиция V^* содержит главные компоненты как векторы-столбцы.

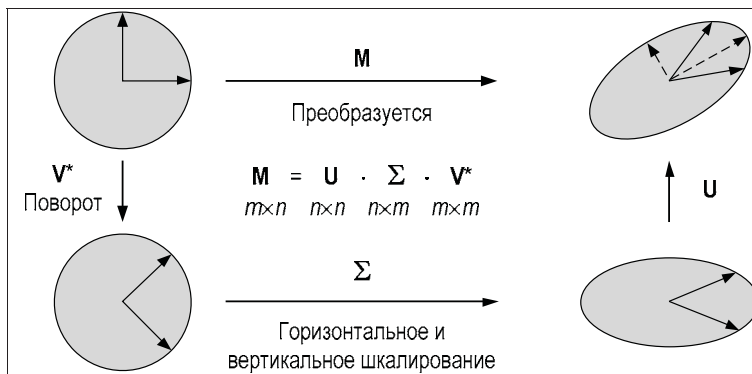


Рис. 12.7. Пошаговое сингулярное разложение

В этом случае нам нужно обеспечить, чтобы данные были центрированы со средним значением, равным нулю (предыдущее вычисление ковариации об этом позаботилось):

```
n_features = data.shape[1]
data_ = data - data.mean(axis=0)
```

С помощью центрированных данных мы вычисляем сингулярное разложение:

```
U, s, Vt = svd(data_)
U.shape, s.shape, Vt.shape
```

```
((100, 100), (3,), (3, 3))
```

Мы можем преобразовать вектор **s**, содержащий только сингулярные значения, в матрицу $m \times n$ и показать, что разложение работает:

```
# Конвертировать s из вектора в диагональную матрицу
S = np.zeros_like(data_)
S[:n_features, :n_features] = np.diag(s)
S.shape
```

```
(100, 3)
```

Мы обнаруживаем, что разложение действительно воспроизводит стандартизированные данные:

```
np.allclose(data_, U.dot(S).dot(Vt))
```

```
True
```

Наконец, мы подтверждаем, что столбцы транспозиции V^* содержат главные компоненты:

```
np.allclose(np.abs(C), np.abs(Vt.T))
```

```
True
```

В следующем разделе мы покажем, как в библиотеке **sklearn** реализован алгоритм анализа главных компонент.

PCA с помощью библиотеки **sklearn**

Реализация **sklearn.decomposition.PCA** следует стандартному API на основе методов **fit()** и **transform()**, которые соответственно вычисляют желаемое число главных компонент и проецируют данные в компонентное пространство. Вспомогательный метод **fit_transform()** выполняет все это за один шаг.

PCA предлагает четыре разных алгоритма, которые можно указать с помощью параметра **svd_solver**:

- ◆ **Full** вычисляет точное сингулярное разложение, используя решатель LAPACK, предоставляемый библиотекой SciPy;

- ◆ `Arpack` выполняет усеченную версию, подходящую для вычисления меньшего, чем полное, числа компонент;
- ◆ `Randomized` использует алгоритм на основе выборки, который показывает более высокую эффективность, когда набор данных имеет более 500 наблюдений и признаков, и цель состоит в том, чтобы вычислить менее 80% компонент;
- ◆ `Auto` использует алгоритм `Randomized`, где это наиболее эффективно, в противном случае он использует полное сингулярное разложение `Full`.

Справочные материалы с подробностями алгоритмической реализации см. в репозитории `GitHub`.

Другими ключевыми конфигурационными параметрами объекта `PCA` являются следующие.

- ◆ `n_components` — вычисляет все главные компоненты, передавая `None` (по умолчанию) или ограничивая число значением `int`. Для параметра `svd_solver=full` существует два дополнительных варианта: вещественное число в промежутке $[0; 1]$ вычисляет число компонент, необходимых для захвата соответствующей доли дисперсии в данных, и вариант `mle` оценивает число размерностей с использованием максимального правдоподобия.
- ◆ `whiten` — если `True`, он стандартизирует компонентные векторы в единичной дисперсии, которая в некоторых случаях может быть полезна для использования в предсказательной модели (по умолчанию `False`).

Для вычисления первых двух главных компонент трехмерного эллипса и проецирования данных в новое пространство следует использовать метод `fit_transform()` следующим образом:

```
pca = PCA(n_components=2)
projected_data = pca.fit_transform(data)
projected_data.shape
```

```
(100, 2)
```

Объясненная дисперсия первых двух компонент очень близка к 100%:

```
pca.explained_variance_ratio_

array([0.77381099, 0.22385721])
```

Рисунок в начале этого раздела показывает проекцию данных в новое двумерное пространство.

Анализ независимых компонент

Анализ независимых компонент (Independent Component Analysis, ICA) — это еще один алгоритм линейного снижения размерности, который определяет новый базис для представления исходных данных, но преследует цель, отличающуюся от анализа главных компонент (PCA).

Анализ независимых компонент (ICA) появился в сфере обработки сигналов, и задача, которую он призван решать, называется *слепым разделением источников*. Она обычно оформляется как задача коктейльной вечеринки, в которой определенное число гостей говорят одновременно, вследствие чего единственный микрофон записывает накладывающиеся друг на друга сигналы. ICA исходит из того, что существует столько же разных микрофонов, сколько и говорящих людей, и каждый из них расположен в разных местах, благодаря чему можно записывать разное сочетание сигналов. Затем ICA стремится восстановить отдельные сигналы из разных записей.

Другими словами, существует n исходных сигналов и неизвестная квадратная матрица смешивания \mathbf{A} , которая производит n -мерное множество \mathbf{M} наблюдений, такое что:

$$\underset{n \times m}{\mathbf{X}} = \underset{n \times n}{\mathbf{A}} \underset{n \times m}{\mathbf{s}}.$$

Цель состоит в том, чтобы найти матрицу $\mathbf{W} = \mathbf{A}^{-1}$, которая распутывает смешанные сигналы с целью восстановления источников.

Возможность однозначно определить матрицу \mathbf{W} зависит от негауссова распределения данных. В противном случае \mathbf{W} можно было бы поворачивать произвольно, учитывая симметрию многомерного нормального распределения при вращении.

Более того, ICA исходит из того, что смешанный сигнал является суммой его компонент и не может выявлять гауссовы компоненты, поскольку их сумма тоже является нормально распределенной.

Допущения анализа независимых компонент

Анализ независимых компонент (ICA) принимает следующие критически важные допущения:

- ◆ источники сигналов являются статистически независимыми;
- ◆ линейных преобразований достаточно для улавливания релевантной информации;
- ◆ независимые компоненты не имеют нормального распределения;
- ◆ матрица смешивания \mathbf{A} может быть инвертирована.

ICA также требует, чтобы данные были центрированы и отбелены, т. е. взаимно некоррелированы с единичной дисперсией. Предобработка данных с помощью PCA, как описано выше, обеспечивает необходимые преобразования.

Алгоритм ICA

Алгоритм FastICA, используемый в библиотеке sklearn, является алгоритмом с фиксированной точкой, который для восстановления независимых источников использует статистику более высокого порядка. В частности, в качестве индикатора

независимости он максимизирует расстояние до нормального распределения для каждой компоненты.

Альтернативный алгоритм под названием InfoMax в качестве меры статистической независимости минимизирует взаимную информацию между компонентами.

ICA в библиотеке sklearn

Реализация ICA в библиотеке sklearn использует тот же интерфейс, что и реализация PCA, поэтому мало что можно добавить. Обратите внимание, что нет никакой меры объясненной дисперсии, потому что ICA не вычисляет компоненты последовательно. Вместо этого каждая компонента призвана улавливать независимые аспекты данных.

Анализ главных компонент для алгоритмической торговли

Анализ главных компонент (PCA) полезен для алгоритмической торговли в нескольких отношениях, в том числе для ведомого данными выведения рисков факторов путем применения PCA к возвратам от финансовых активов, а также конструирования некоррелированных портфелей, основываясь на главных компонентах корреляционной матрицы возвратов от активов.

Рисковые факторы, ведомые данными

В *главе 7* мы рассмотрели модели рисков факторов, используемые в квантитативном финансировании для улавливания главных движущих сил финансовых возвратов. Эти модели объясняют разницы в возвратах от активов на основе влияния на них систематических рисков факторов и ассоциированные с ними выгоды.

В частности, мы исследовали подход Фама — Френча, который конкретизирует факторы, основываясь на предшествующих (априорных) знаниях об эмпирическом поведении средних финансовых возвратов, рассматривает эти факторы как наблюдаемые, а затем оценивает коэффициенты модели риска с использованием линейной регрессии. Альтернативный подход рассматривает рисков факторы как латентные переменные и использует методы факторного анализа, такие как PCA, для одновременного оценивания факторов и того, как они стимулируют возвраты, на основе исторических возвратов.

В этом разделе мы рассмотрим, как этот метод выводит факторы чисто статистическим или ведомым данными способом с тем преимуществом, что предварительных знаний о поведении возвратов от финансовых активов не требуется (подробнее см. модели `pca` и `risk_factor` в блокноте).

Мы воспользуемся данными платформы Quandl с ценами акций и отберем ежедневные скорректированные цены закрытия 500 акций с наибольшей рыночной капитализацией и данными за период 2010–2018 гг. Затем мы вычислим дневные финансовые возвраты следующим образом:


```

idx = pd.IndexSlice
DATA_STORE = '../data/assets.h5'
with pd.HDFStore(DATA_STORE) as store:
    stocks = store['us_equities/stocks'].marketcap.nlargest(500)
    returns = (store['quandl/wiki/prices']
               .loc[idx['2010': '2018', stocks.index], 'adj_close']
               .unstack('ticker')
               .pct_change())

```

Мы получаем 351 акцию и возвраты за более 2000 торговых дней:

```

returns.info()

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 2072 entries, 2010-01-04 to 2018-03-27
Columns: 215 entries, AAPL to GSBD
dtypes: float64(215)
memory usage: 3.4 MB

```

РСА чувствителен к выбросам, поэтому мы винзоризуем данные на 2,5%-м и 97,5%-м квантилях:

```

returns = returns.clip(lower=returns.quantile(q=.025),
                      upper=returns.quantile(q=.975),
                      axis=1)

```

РСА не допускает пропущенных данных, поэтому мы удалим акции, которые не имеют данных по крайней мере в течение 95% периода времени, и на втором шаге удалим торговые дни, которые не имеют наблюдений по крайней мере в 95% оставшихся акциях:

```

returns = returns.dropna(thresh=int(returns.shape[0] * .95), axis=1)
returns = returns.dropna(thresh=int(returns.shape[1] * .95))

```

У нас осталось 314 рядов финансовых возвратов от долевого ценных бумаг, охватывающих аналогичный период:

```

returns.info()

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 2071 entries, 2010-01-05 to 2018-03-27
Columns: 173 entries, AAPL to PLOW
dtypes: float64(173)
memory usage: 2.7 MB

```

Мы вмняем (импутируем) любые оставшиеся пропущенные значения, используя средний финансовый возврат за любой торговый день:

```

daily_avg = returns.mean(1)
returns = returns.apply(lambda x: x.fillna(daily_avg))

```

Теперь мы готовы вписать модель главных компонент в возвраты от активов, используя стандартные параметры для вычисления всех компонент с помощью алгоритма полного SVD:

```
pca = PCA()
pca.fit(returns)

PCA(copy=True, iterated_power='auto',
      n_components=None, random_state=None,
      svd_solver='auto', tol=0.0, whiten=False)
```

Мы находим, что самый важный фактор объясняет около 40% ежедневной изменчивости финансового возврата. Доминирующий фактор обычно интерпретируется как рынок, в то время как остальные факторы могут интерпретироваться как отраслевые или стилевые факторы, в соответствии с нашим обсуждением в *главах 5 и 7* в зависимости от результатов более пристального обследования (см. следующий пример).

На рис. 12.8 график справа показывает кумулятивную объясненную дисперсию и указывает, что около 10 факторов объясняют 60% возвратности этого большого среза акций.

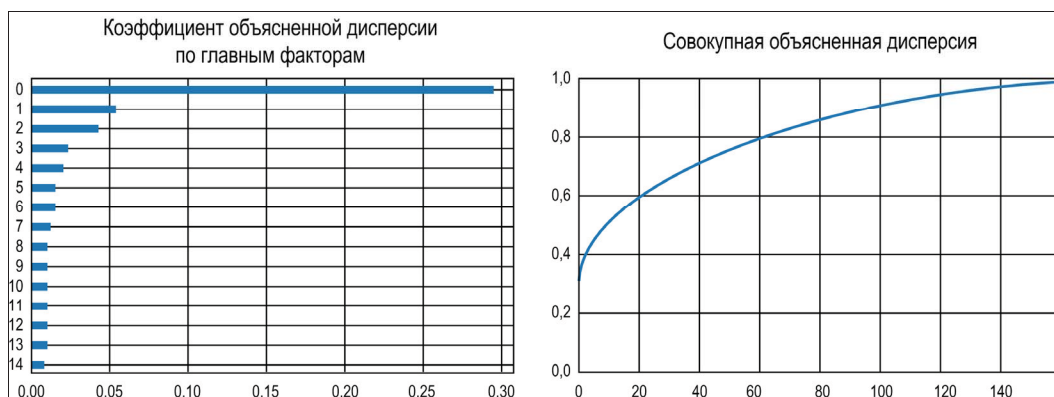


Рис. 12.8. Доли объясненной дисперсии по фактору и совокупная объясненная дисперсия

Блокнот содержит симуляцию более широкого среза акций и более продолжительного периода времени за 2000–2018 гг. Результаты показывают, что в среднем первые три компоненты объясняли 25, 10 и 5% из 500 случайно отобранных акций.

На кумулятивной диаграмме показана типичная локтевая регулярность, которая помогает определять подходящую целевую размерность, поскольку она указывает, что дополнительные компоненты добавляют меньшее объяснительное значение.

Мы можем выбрать две главных компоненты для того, чтобы убедиться, что они действительно некоррелированы:

```

risk_factors = pd.DataFrame(pca.transform(returns)[: , :2],
                             columns=['Principal Component 1',
                                       'Principal Component 2'],
                             index=returns.index)
risk_factors['Principal Component 1'].corr(risk_factors['Principal Component 2'])

```

```
-1.1270779275205213e-14
```

Кроме того, мы можем построить график временных рядов, который проиллюстрирует, как каждый фактор улавливает разные регулярности волатильности (рис. 12.9).

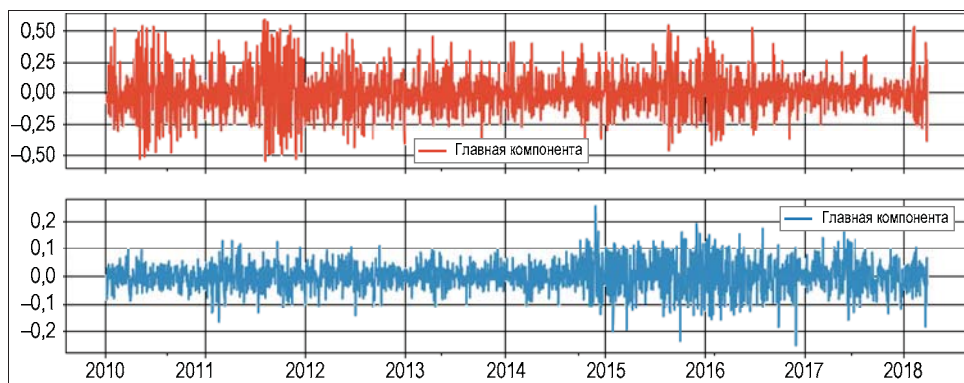


Рис. 12.9. Временные ряды, иллюстрирующие улавливание каждым фактором разных регулярностей волатильности

Модель рисков факторов будет задействовать подмножество главных компонент в качестве признаков, предсказывая будущие финансовые возвраты, подобно нашему подходу в главе 7, где рассматривались регрессия и классификация.

Характеристические портфели

Еще одно применение анализа главных компонент (PCA) предусматривает ковариационную матрицу нормализованных финансовых возвратов. Главные компоненты корреляционной матрицы улавливают подавляющую часть ковариации между активами в убывающем порядке и являются взаимно некоррелированными. Более того, стандартизированные главные компоненты можно использовать в качестве портфельных весов.

Для иллюстрации давайте возьмем 30 крупнейших акций с данными за период 2010–2018 гг.:

```

idx = pd.IndexSlice
DATA_STORE = '../data/assets.h5'
with pd.HDFStore(DATA_STORE) as store:
    stocks = store['us_equities/stocks'].marketcap.nlargest(30)

```

```
returns = (store['quandl/wiki/prices']
           .loc[idx['2010': '2018', stocks.index], 'adj_close']
           .unstack('ticker')
           .pct_change())
```

Мы снова винзоризуем и нормализуем возвраты:

```
normed_returns = scale(returns
                       .clip(lower=returns.quantile(q=.025),
                             upper=returns.quantile(q=.975),
                             axis=1)
                       .apply(lambda x: x.sub(x.mean()).div(x.std())))
```

После отбрасывания активов и торговых дней, как в предыдущем примере, у нас остается 23 актива и более 2000 торговых дней. Мы оцениваем все главные компоненты и находим, что два самых больших из них объясняют соответственно 55,9 и 15,5% ковариации:

```
pca = PCA()
pca.fit(cov)
pd.Series(pca.explained_variance_ratio_)
    .to_frame('Объясненная дисперсия') \
    .head().style.format('{:,.2%}'.format)
```

| | Объясненная дисперсия |
|---|-----------------------|
| 0 | 57.62% |
| 1 | 12.41% |
| 2 | 6.67% |
| 3 | 5.54% |
| 4 | 3.48% |

Далее мы отбираем и нормализуем четыре крупнейших компоненты, вследствие чего они в сумме составляют 1, и их можно было использовать в качестве весов для портфелей, которые мы можем сравнить с равновзвешенным портфелем, сформированным из всех акций:

```
top4 = pd.DataFrame(pca.components_[:4], columns=cov.columns)
eigen_portfolios = top4.div(top4.sum(1), axis=0)
eigen_portfolios.index = [f'Portfolio {i}' for i in range(1, 5)]
```

Веса показывают отчетливый акцент — например, портфель 3 перекладывает большие веса на Mastercard и Visa, двух платежных процессоров в выборке, в то время как на портфель 2 больше влияют технологические компании (рис. 12.10).

При сравнении результативности каждого портфеля за период выборки с рынком, состоящим из нашей небольшой выборки, мы обнаруживаем, что портфель 1 работает очень похоже, в то время как другие портфели улавливают разные регулярности финансовых возвратов (рис. 12.11).

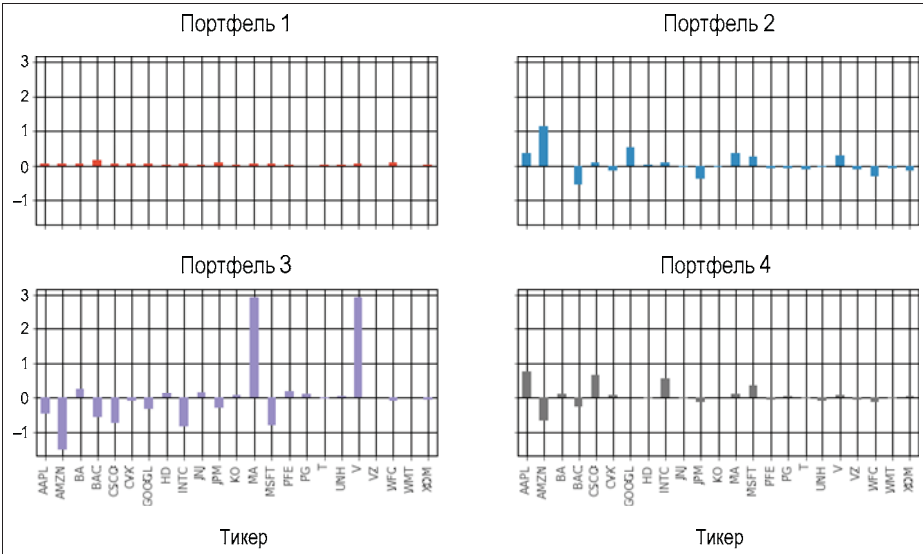


Рис. 12.10. Сравнение портфелей по финансовым активам (тикерам)

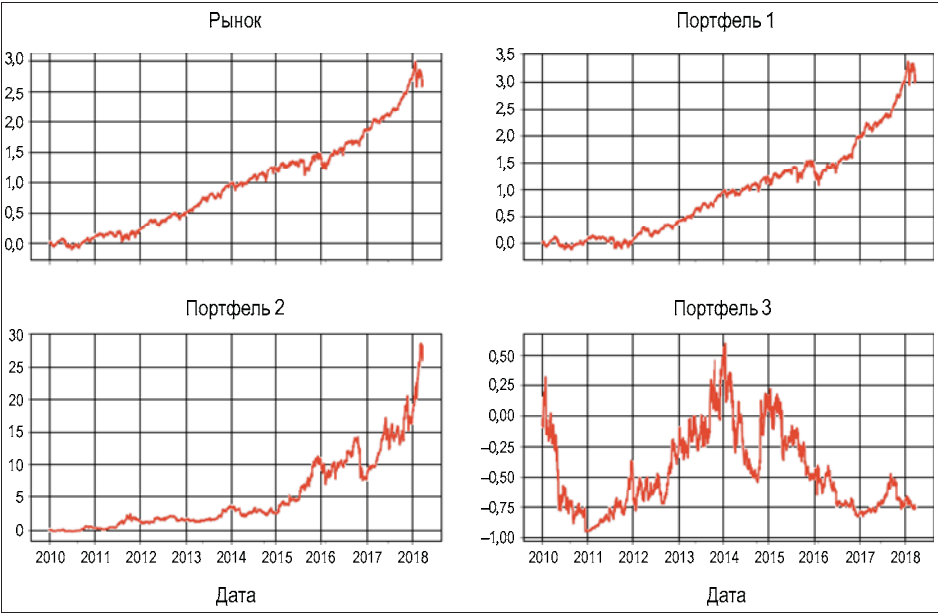


Рис. 12.11. Сравнение результативности каждого портфеля

Усвоение проекций в топологическом многообразии

Линейное снижение размерности проецирует исходные данные на низкоразмерную гиперплоскость, которая выравнивается с информативными направлениями в данных. Акцент на линейных преобразованиях упрощает вычисления и перекликается с общими финансовыми метриками, такими как цель PCA для захвата максимальной дисперсии.

Однако линейные подходы естественным образом игнорируют сигнал, отраженный в нелинейных связях в данных. Такие связи очень важны в альтернативных наборах данных, содержащих, например, изображения или текстовые данные. Обнаружение таких связей во время разведывательного анализа может давать важные подсказки о потенциальном содержимом сигнала.

В противоположность этому, гипотеза о топологическом многообразии, или многосвязной области, подчеркивает, что высокоразмерные данные часто лежат в нелинейной малоразмерной многосвязной области, встроенной в высокоразмерное пространство, или вблизи нее. Двумерный швейцарский рулет, показанный на рис. 12.2 в начале этой главы, иллюстрирует такую топологическую структуру.

Алгоритм усвоения проекции в топологическом многообразии призван отыскивать многосвязную область внутренней размерности, а затем представлять данные в этом подпространстве. В упрощенном примере дорога будет представлять одномерную топологическую область в трехмерном пространстве, а точки данных являются на основе номеров домов в качестве локальных координат.

Существует несколько технических решений, которые аппроксимируют низкоразмерное топологическое многообразие. Одним из примеров является *локально-линейное вложение* (locally-linear embedding, LLE), которое было разработано в 2000 г. Сэмом Ровейсом (Sam Roweis) и Лоуренсом Солом (Lawrence Saul) и использовалось для развертывания швейцарского рулета на рис. 12.2 (см. примеры в блокноте `manifold_learning_lle.ipynb`).

Для каждой точки данных локально-линейное вложение (LLE) определяет заданное число ближайших соседей и вычисляет веса, представляющие каждую точку как линейную комбинацию ее соседей. Указанный метод находит низкоразмерное вложение, линейно проецируя каждую окрестность на глобальные внутренние координаты на низкоразмерном топологическом многообразии, и может рассматриваться как последовательность применений PCA.

Визуализация требует снижения, по крайней мере, до трех размерностей, возможно, ниже внутренней размерности, и ставит задачу точного представления локальной и глобальной структуры. Эта задача коррелирует с увеличением расстояния, связанного с проклятием размерности. В то время как объем сферы расширяется экспоненциально вместе с числом размерностей, пространство в более низких размерностях, доступное для представления высокоразмерных данных, гораздо более ограничено.

Например, в 12 размерностях может быть 13 равноудаленных точек, но в двух размерностях могут быть только три, образуя треугольник со сторонами одинаковой длины. Следовательно, точное отражение расстояния одной точки до ее высоко-размерных соседей в более низких размерностях рискует исказить связи между всеми другими точками. В результате возникает проблема скученности: для поддержания глобальных расстояний локальные точки, возможно, придется разместить слишком близко друг к другу, и наоборот.

В следующих двух разделах рассматриваются методы, которые позволили добиться прогресса в решении проблемы скученности для визуализации сложных наборов данных. Мы будем использовать набор данных MNIST стильной одежды, более сложную альтернативу классическому эталонному набору данных MNIST рукописных цифр, используемому для компьютерного зрения. Указанный набор содержит 60 тыс. тренировочных и 10 тыс. тестовых снимков предметов стильной одежды в 10 классах (см. образцы на рис. 12.12).

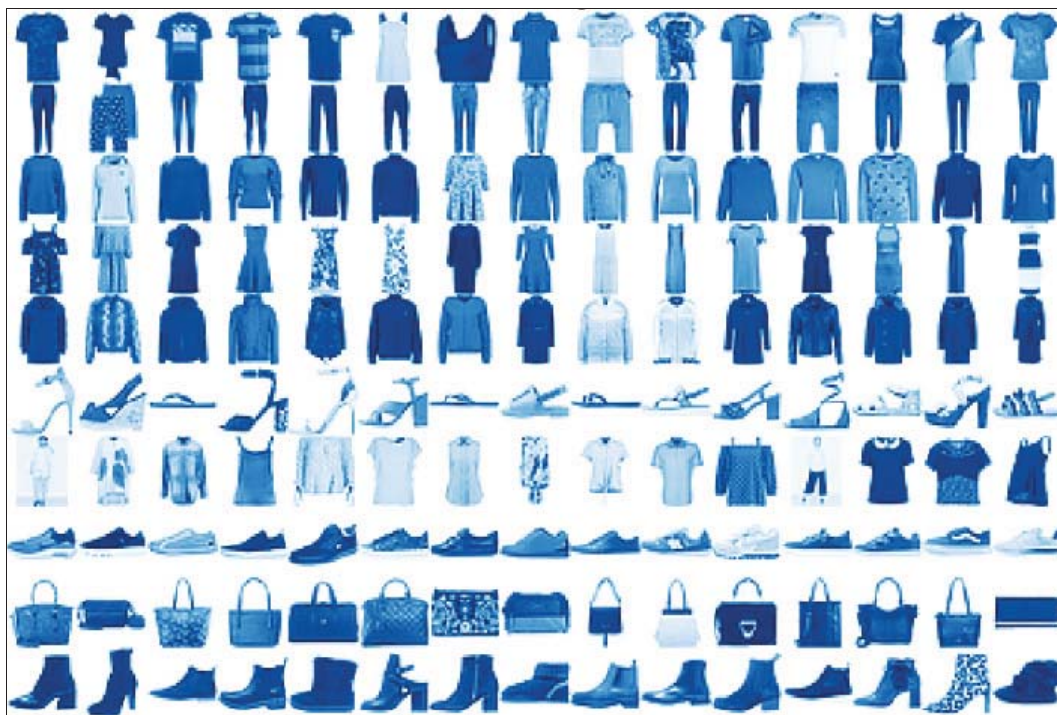


Рис. 12.12. Снимки предметов стильной одежды

Цель алгоритма усвоения проекции в топологическом многообразии для этих данных состоит в том, чтобы обнаруживать расположение классов на разных много-связных областях с целью обеспечения более легкого их распознавания и дифференцирования.

Алгоритм t-SNE

t-Распределенное стохастическое вложение соседей (t-distributed Stochastic Neighbor Embedding, t-SNE) — это удостоенный наград алгоритм, разработанный в 2010 г. Лоуренсом ван дер Маатеном (Laurens van der Maaten) и Джеффом Хинтоном (Geoff Hinton) для обнаружения регулярностей в высокоразмерных данных. В нем принят вероятностный, нелинейный подход к локализации данных на нескольких разных, но родственных низкоразмерных топологических многообразиях.

Указанный алгоритм подчеркивает поддержание вместе схожих точек в низких размерностях, в отличие от поддержания расстояния между точками, которые отделены друг от друга в высоких размерностях, являясь результатом таких алгоритмов, как PCA, которые минимизируют квадратичные расстояния.

Указанный алгоритм работает, конвертируя высокоразмерные расстояния в (условные) вероятности, где высокие вероятности влекут за собой низкое расстояние и отражают правдоподобие отбора двух точек на основе сходства. Он выполняет это, располагая нормальное распределение над каждой точкой и вычисляя плотность для точки и каждого соседа, где параметр перплексивности управляет эффективным числом соседей.

На втором шаге он упорядочивает точки в низких размерностях и использует аналогично вычисленные низкоразмерные вероятности, сочетая с высокоразмерным распределением. Он служит мерой разницы между распределениями с использованием расхождения Кульбака — Лейблера, которое накладывает высокий штраф за неправильное размещение схожих точек в низких размерностях.

Низкоразмерные вероятности используют *t*-распределение Стюдента с одной степенью свободы, поскольку оно имеет более толстые хвосты, которые снижают штраф за неправильное размещение точек, более отдаленных в высоких размерностях, тем самым управляя проблемой скученности.

Верхние изображения на рис. 12.13 показывают, как алгоритм t-SNE может различать классы снимков. Более высокое значение перплексивности увеличивает число соседей, используемых для вычисления локальной структуры, и постепенно приводит к большему акценту на глобальных взаимосвязях.

Алгоритм t-SNE в настоящее время отражает самое современное состояние дел в визуализации высокоразмерных данных. Его недостатки включают вычислительную сложность, чьи масштабы увеличиваются квадратично по числу n точек, поскольку он оценивает все попарные расстояния, но последующая древесная его реализация снизила стоимость до $n \log n$.

Алгоритм t-SNE не обеспечивает проекцию в низкоразмерное пространство новых точек данных. Сжатый выход не очень полезен для алгоритмов кластеризации, основывающихся на расстоянии или плотности, потому что указанный алгоритм рассматривает малые и большие расстояния по-разному.

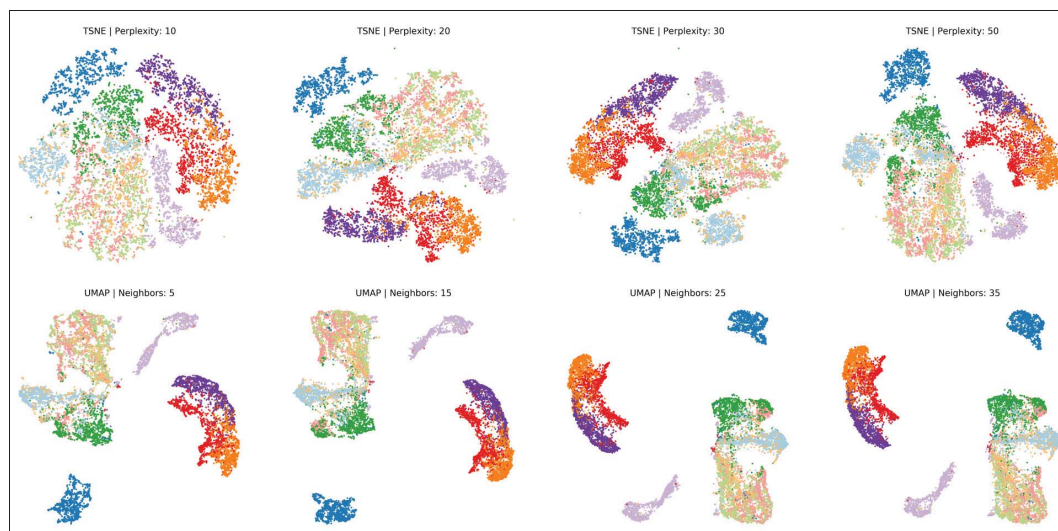


Рис. 12.13. Результаты работы алгоритмов t-SNE и UMAP

Алгоритм UMAP

Равномерная аппроксимация топологического многообразия и проекция в нем (Uniform Manifold Approximation and Projection, UMAP) является более поздним алгоритмом визуализации и общего снижения размерности. Указанный алгоритм исходит из того, что данные равномерно распределены на локально многосвязной области и ищут ближайший низкоразмерный эквивалент с использованием нечеткой топологии. Он использует параметр окрестности, который влияет на результат аналогично упомянутой выше перплексивности.

Указанный алгоритм работает быстрее и, следовательно, лучше, чем алгоритм t-SNE, масштабируется на крупные наборы данных, и иногда сохраняет глобальную структуру лучше, чем алгоритм t-SNE. Он также может работать с разными функциями расстояния, включая, например, косинусное сходство, которое используется для измерения расстояния между векторами количеств вхождений слов.

Четыре диаграммы в нижней части рис. 12.13 иллюстрируют, как алгоритм UMAP перемещает разные кластеры друг от друга дальше, в то время как алгоритм t-SNE обеспечивает более детальное представление о локальной структуре.

Блокнот также содержит интерактивные графические визуализации для каждого алгоритма, которые дают возможность разведать метки и выявить, какие объекты расположены друг к другу близко.

Кластеризация

И кластеризация, и снижение размерности резюмируют данные. Как только что подробно отмечалось, процедура снижения размерности сжимает данные и представляет их с использованием новых, меньших признаков, которые улавливают наиболее релевантную информацию. Алгоритмы кластеризации, напротив, назначают существующие наблюдения подгруппам, состоящим из похожих точек данных.

Кластеризация помогает лучше понять данные через призму категорий, усваиваемых из непрерывных переменных. Она также позволяет автоматически классифицировать новые объекты в соответствии с усвоенными критериями. Примеры связанных с кластеризацией приложений включают иерархические таксономии, медицинскую диагностику и сегментацию клиентов.

Альтернативно кластеры можно использовать для представления групп в качестве прототипов, используя (например) срединную точку кластера как наилучшего представителя усвоенной группы. Пример использования включает сжатие изображений.

Алгоритмы кластеризации различаются по стратегиям выявления групп:

- ◆ комбинаторные алгоритмы отбирают наиболее когерентные из разных групп наблюдений;
- ◆ вероятностное моделирование оценивает распределения, которые, скорее всего, стали причиной формирования кластеров;
- ◆ иерархическая кластеризация находит последовательность вложенных кластеров, которая оптимизирует когерентность на любом заданном этапе.

Алгоритмы также различаются в своем понимании того, что собой представляет полезная коллекция объектов, которая должна сочетать характеристики данных, предметную область и цель применений. Типы групп включают:

- ◆ четко разделенные группы различных форм;
- ◆ компактные кластеры на основе прототипов или центров;
- ◆ кластеры произвольной формы на основе плотности;
- ◆ кластеры, основанные на связности или графах.

Важными дополнительными аспектами кластерного алгоритма являются следующие:

- ◆ требуется ли эксклюзивное членство в кластере;
- ◆ делает ли он жесткое (бинарное) или мягкое (вероятностное) назначение точек данным кластерам;
- ◆ является ли он полным и назначает все точки данным кластерам.

В следующих далее разделах представлены ключевые алгоритмы, в том числе кластеризация k средних, иерархическая кластеризация и плотностная кластеризация,

а также модели гауссовых смесей. Блокнот `clustering_algos.ipynb` сравнивает результативность этих алгоритмов на разных помеченных наборах данных, выделяя их сильные и слабые стороны. Для измерения конгруэнтности кластерных назначений и меток в нем используется взаимная информация (см. главу 6).

Кластеризация на основе k средних

Алгоритм k средних является самым известным алгоритмом кластеризации, который был впервые предложен Стюартом Ллойдом (Stuart Lloyd) в Лаборатории Белла в 1957 г.

Указанный алгоритм находит K центроидов и назначает каждой точке данных ровно один кластер, минимизируя внутрикластерную дисперсию (именуемую инерцией). Обычно используется евклидово расстояние, но могут применяться и другие метрики. Алгоритм k средних исходит из того, что кластеры являются сферическими и имеют одинаковый размер, и игнорирует ковариацию между признаками.

Решаемая задача является вычислительно сложной (НП-трудной), потому что существует K^N способов подразделить N наблюдений на K кластеров. Стандартный итерационный алгоритм обеспечивает локальный оптимум для заданного K и действует следующим образом:

1. Случайно сформировать K кластерных центров и назначить точки ближайшему центроиду.
2. Повторять следующим образом:
 - для каждого кластера вычислить центроид как среднее значение признаков;
 - назначить каждое наблюдение ближайшему центроиду.
3. Схождение: назначения точек кластерам (или внутрикластерные вариации) не изменяются.

Блокнот `kmeans_implementation.ipynb` показывает, как кодировать указанный алгоритм на языке Python, и визуализирует итеративную оптимизацию этого алгоритма. На рис. 12.14 продемонстрировано, как результирующие центроиды разделяют признаковое пространство на участки, именуемые областями Вороного, которые очерчивают кластеры.

Результат является оптимальным для этой конкретной инициализации, но альтернативные стартовые позиции породят другие результаты. Отсюда следует, что мы вычисляем многочисленные кластеры из разных начальных значений и выбираем решение, которое минимизирует внутрикластерную дисперсию.

Алгоритм k средних требует наличия категориальных переменных, которые являются непрерывными либо кодированными с одним активным состоянием. Метрики расстояния, как правило, чувствительны к шкале, вследствие чего необходима стандартизация признаков, которая обеспечивает их одинаковый вес.

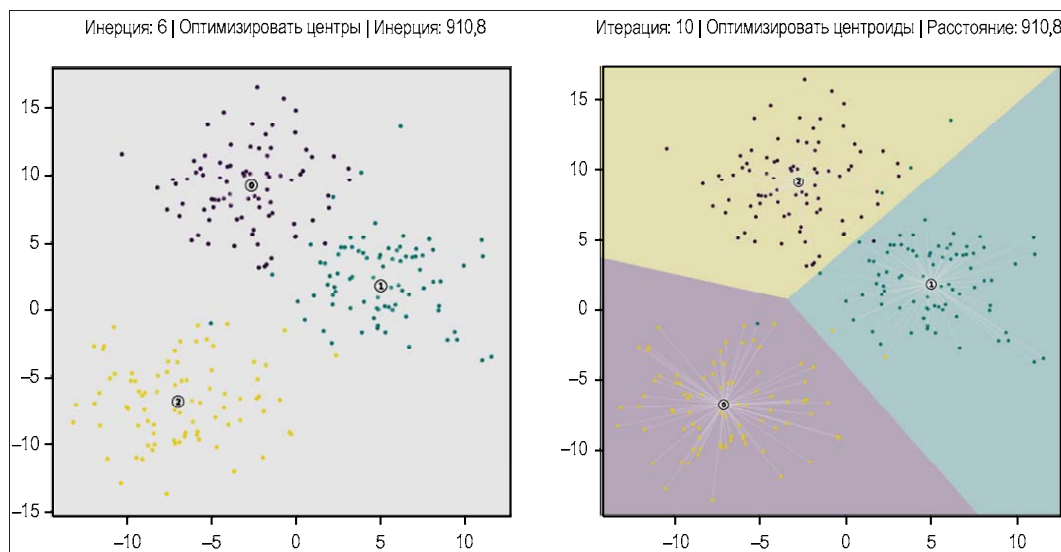


Рис. 12.14. Очерчивание кластеров областями Вороного

Сильные стороны алгоритма k средних включают его широкий диапазон применимости, быструю сходимость и линейную масштабируемость на крупные данные, при этом производя кластеры одинакового размера.

К слабым сторонам относятся:

- ◆ необходимость настройки гиперпараметра k ;
- ◆ отсутствие гарантии отыскания глобального оптимума;
- ◆ ограничительное допущение, что кластеры являются сферами, а признаки не коррелированы;
- ◆ чувствительность к выбросам.

Оценивание качества кластеров

Метрики качества кластеров помогают выбирать среди альтернативных результатов кластеризации.

В блокноте `kmeans_evaluation.ipynb` показаны следующие параметры:

1. Целевая функция k средних предлагает сравнивать эволюцию инерции или внутрикластерной дисперсии.
2. Первоначально дополнительные центры резко уменьшают инерцию, поскольку новые кластеры улучшают совокупную подгонку.
3. После того как найдено соответствующее число кластеров (если оно существует), новые центры уменьшают внутрикластерную дисперсию намного меньше, поскольку они тяготеют к разбивке естественных группировок.

4. Отсюда, когда алгоритм k средних находит хорошее кластерное представление данных, инерция склонна следовать локтевидной траектории, аналогичной коэффициенту объясненной дисперсии для РСА, как показано на рис. 12.15 (подробности реализации см. в блокноте).

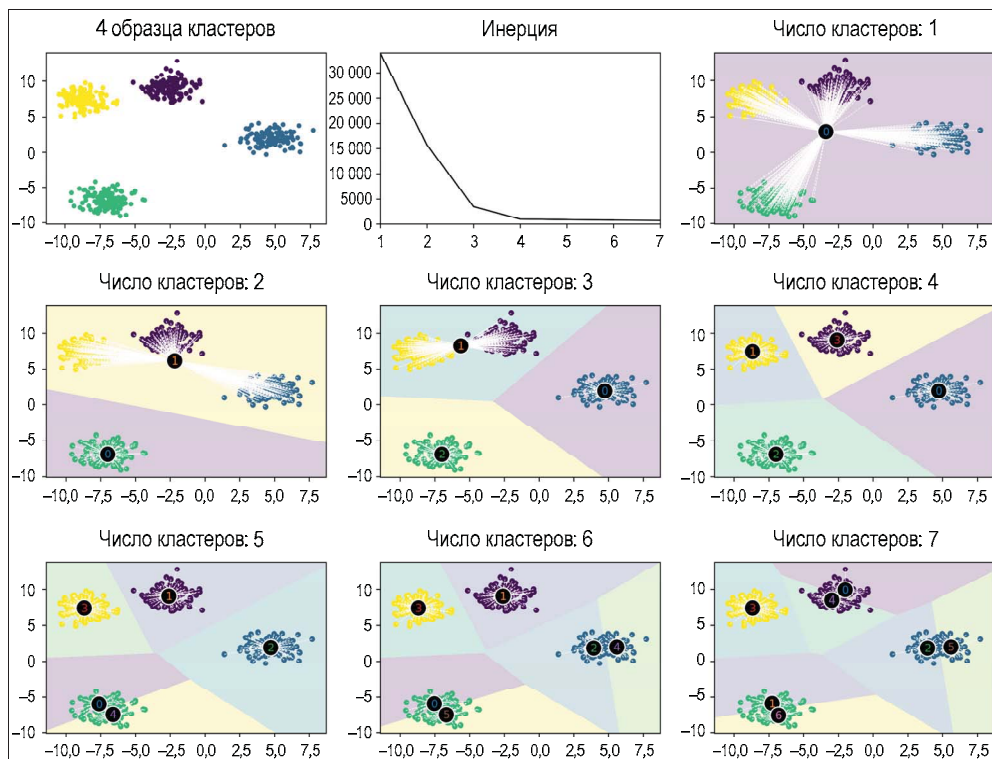


Рис. 12.15. Процесс кластеризации на разных стадиях и локтевидная траектория инерции

Силуэтный коэффициент обеспечивает более детальную картину качества кластеров. Он отвечает на вопрос: как далеко находятся точки в ближайшем кластере относительно точек в назначенном кластере?

С этой целью он сравнивает среднее внутрикластерное расстояние (a) со средним расстоянием до ближайшего кластера (b) и вычисляет следующую отметку s :

$$s = \frac{b - a}{\max(a, b)} \in [-1; 1].$$

Данная отметка может варьироваться от -1 до 1 , но отрицательные значения на практике маловероятны, поскольку из них следует, что большинство точек назначены неправильному кластеру. Полезная визуализация силуэтной отметки сравнивает значения для каждой точки данных с глобальным средним, поскольку она подчеркивает согласованность каждого кластера относительно глобальной конфигура-

ции. Эмпирическое правило состоит в том, чтобы избегать кластеров со средними отметками ниже среднего для всех образцов.

На рис. 12.16 показана выдержка из силуэтного графика для трех и четырех кластеров, где первый случай подчеркивает плохую подгонку кластера 1 по вкладам в глобальную силуэтную отметку, не дотягивающую до среднего, в то время как все четыре кластера имеют некоторые значения, которые демонстрируют отметки выше среднего.

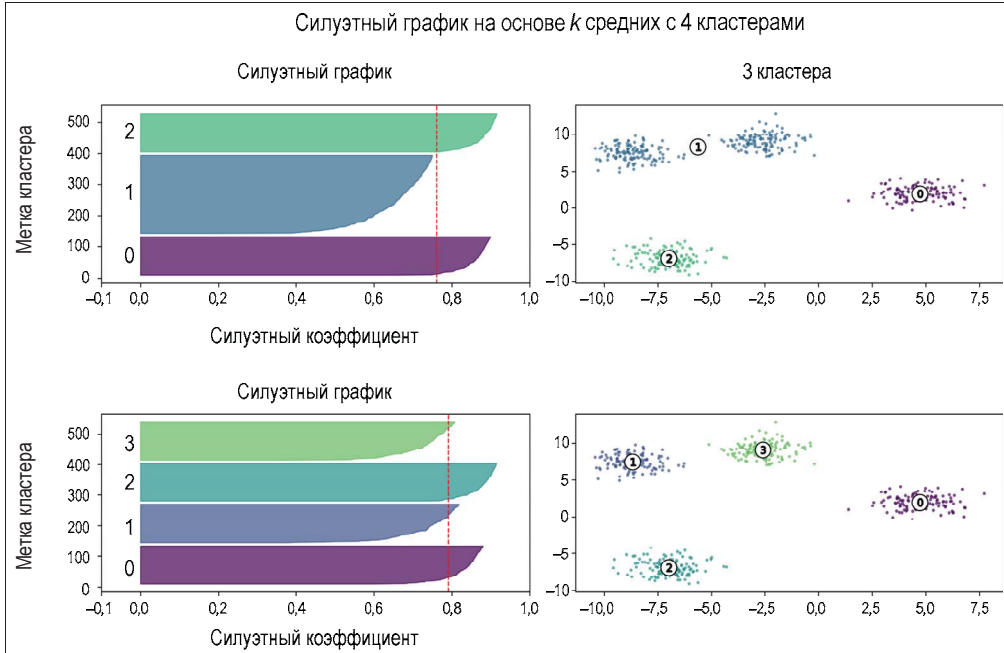


Рис. 12.16. Сравнение силуэтных отметок качества кластеров

В целом, учитывая обычно неконтролируемую природу кластеризации, необходимо варьировать гиперпараметры кластерных алгоритмов и оценивать разные результаты. Также важно калибровать шкалу признаков, в частности, когда некоторым из них должен быть дан больший вес и, таким образом, должны измеряться на более крупной шкале.

Наконец, для проверки робастности результатов следует использовать подмножества данных, с помощью которых выявляется устойчивость проявления определенных регуляризаторов.

Иерархическая кластеризация

Иерархическая кластеризация позволяет избежать необходимости указывать целевое число кластеров, поскольку она исходит из того, что данные могут последовательно объединяться во все более несхожие кластеры. Она не преследует глобаль-

ную цель, а решает инкрементно то, как создавать последовательность вложенных кластеров, которые варьируются от одного-единственного кластера до кластеров, состоящих из индивидуальных точек данных.

Существует два подхода:

- ◆ *агломеративная кластеризация* работает снизу вверх, последовательно объединяя две оставшиеся группы на основе сходства;
- ◆ *разделительная кластеризация* работает сверху вниз и последовательно разбивает оставшиеся кластеры на наиболее несовпадающие подгруппы.

Оба подхода создают $N - 1$ иерархических уровней и обеспечивают выбор кластеризации на уровне, который лучше всего подразделяет данные на однородные группы. Мы сосредоточимся на более общем подходе — агломеративной кластеризации.

Алгоритм агломеративной кластеризации идет от индивидуальных точек данных и вычисляет матрицу подобия, содержащую все взаимные расстояния. Затем он выполняет $N - 1$ шагов до тех пор, пока больше не останется несовпадающих кластеров, и каждый раз обновляет матрицу подобия, заменяя элементы, которые были объединены новым кластером, вследствие чего матрица постепенно сжимается.

Хотя иерархическая кластеризация не имеет гиперпараметров, таких как у k средних, мера несходства между кластерами (в отличие от отдельных точек данных) оказывает важное влияние на результат кластеризации. Варианты различаются следующим образом:

- ◆ *одиночная связь* — расстояние между ближайшими соседями двух кластеров;
- ◆ *полная связь* — максимальное расстояние между соответствующими членами кластера;
- ◆ *групповое среднее* — расстояние между средними значениями для каждой группы;
- ◆ *метод Уорда* — минимизирует внутрикластерную дисперсию.

Визуализация — дендрограммы

Иерархическая кластеризация позволяет получать сущностное представление о степени сходства между наблюдениями по мере того, как она продолжает сливать данные вместе. Значительное изменение метрики сходства от одного слияния к другому говорит о том, что до этой точки существовала естественная кластеризация.

Дендрограмма визуализирует последовательные слияния в виде двоичного дерева, отображая отдельные точки данных в виде листьев, а конечное слияние — в виде корня дерева. Она также показывает, как сходство монотонно уменьшается снизу вверх. Поэтому вполне естественным является выбор кластеризации путем разрезания дендрограммы.

На рис. 12.17 (подробности см. в блокноте `hierarchical_clustering.ipynb`) показана дендрограмма для классического набора данных ирисов Фишера (Iris) с четырьмя классами.

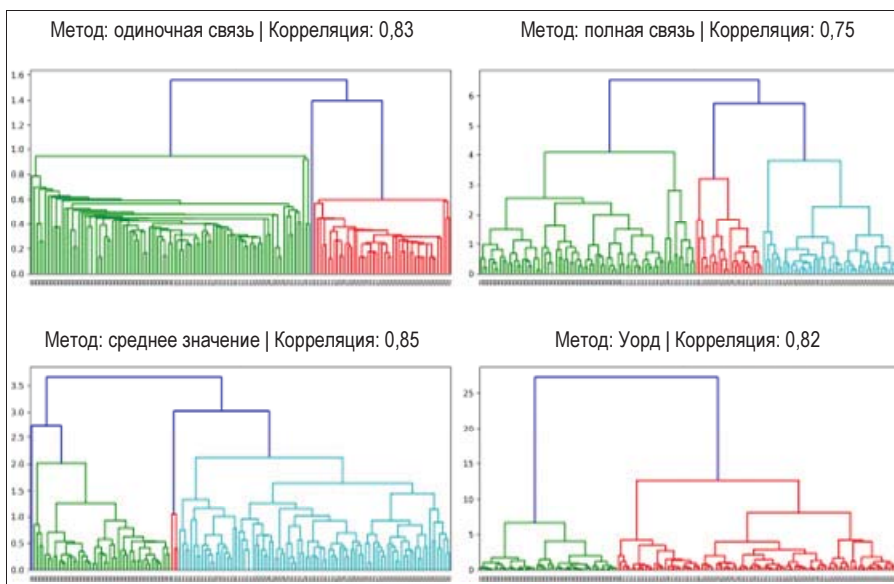


Рис. 12.17. Дендрограммы с четырьмя разными метриками расстояния на наборе данных Iris

ми и тремя признаками, использующая четыре разные метрики расстояния, введенные ранее.

Приведенная выше дендрограмма позволяет оценить подгонку иерархической кластеризации с помощью коэффициента кофенетической корреляции², сравнивающий попарные расстояния между точками и метрику подобия кластера, при которой произошло попарное слияние. Из коэффициента 1 следует, что более близкие точки всегда сливаются раньше.

Разные методы связывания производят дендрограммы разного вида, поэтому мы не сможем использовать эту визуализацию для сравнения результатов между методами. Кроме того, метод Уорда, который минимизирует внутрикластерную дисперсию, может не отражать надлежащим образом изменение в дисперсии, отражая скорее общую дисперсию, что может вводить в заблуждение. Вместо них более подходящими могут быть другие качественные показатели, такие как кофенетическая корреляция, или такие показатели, как инерция (если они согласованы с общей целью).

Сильные стороны указанной кластеризации включают следующие аспекты:

- ◆ не нужно указывать число кластеров;
- ◆ она дает представление о потенциальной кластеризации с помощью интуитивной визуализации;

² Коэффициент кофенетической корреляции (cophenetic correlation coefficient) — это мера того, насколько точно дендрограмма сохраняет попарные расстояния между исходными немоделированными точками данных. — Прим. перев.

- ♦ она создает иерархию кластеров, которые могут служить таксономией;
- ♦ ее можно совместить с k средними для уменьшения числа элементов в начале агломеративного процесса.

Недостатки иерархической кластеризации:

- ♦ высокая стоимость с точки зрения вычисления и памяти из-за многочисленных обновлений матрицы подобия;
- ♦ она не достигает глобального оптимума, потому что все слияния являются окончательными;
- ♦ проклятие размерности приводит к трудностям с шумными, высокоразмерными данными.

Плотностная кластеризация

Алгоритмы плотностной кластеризации назначают членство в кластере на основе близости к другим членам кластера. Они преследуют цель выявления плотных областей произвольных форм и размеров. Они не требуют указания определенного числа кластеров, вместо этого опираясь на параметры, определяющие размер окрестности и порог плотности (соответствующие образцы исходного кода см. в блокноте `density_based_clustering.ipynb`).

Алгоритм DBSCAN

Плотностная пространственная кластеризация приложений с шумом (density-based spatial clustering of applications with noise, DBSCAN) была разработана в 1996 г. и получила награду Test of Time на конференции KDD 2014 г. благодаря вниманию, которое на нее было обращено как в теории, так и на практике.

Указанный алгоритм призван выявлять стержневые и нестержневые образцы, где первые расширяют кластер, а последние являются частью кластера, но не имеют достаточных соседей для дальнейшего роста кластера. Другие образцы являются выбросами и никакому кластеру не назначаются.

Алгоритм использует параметр `eps` для радиуса окрестности и параметр `min_samples` для числа элементов, необходимых для стержневых образцов. Он является детерминированным и эксклюзивным и имеет трудности с кластерами разной плотности и высокоразмерными данными. Кроме того, могут возникать сложности в настройке параметров на требуемую плотность, тем более что она часто не является постоянной.

Иерархический алгоритм DBSCAN

Иерархический алгоритм DBSCAN (HDBSCAN) представляет собой более свежее усовершенствование, которое исходит из того, что кластеры являются островами потенциально различающейся плотности, преодолевая только что упомянутые проблемы сложности алгоритма DBSCAN. Он тоже призван выявлять стержневые и

нестержневые образцы. Для отбора окрестности и расширения кластера он использует параметры `min_cluster_size` и `min_samples`. Указанный алгоритм итеративно перебирает несколько значений `eps` и выбирает наиболее стабильную кластеризацию.

В дополнение к выявлению кластеров разной плотности, он позволяет проникать в сущность плотности и иерархической структуры данных.

Графики рис. 12.18 показывают, как DBSCAN и HDBSCAN могут выявлять кластеры очень разной формы.

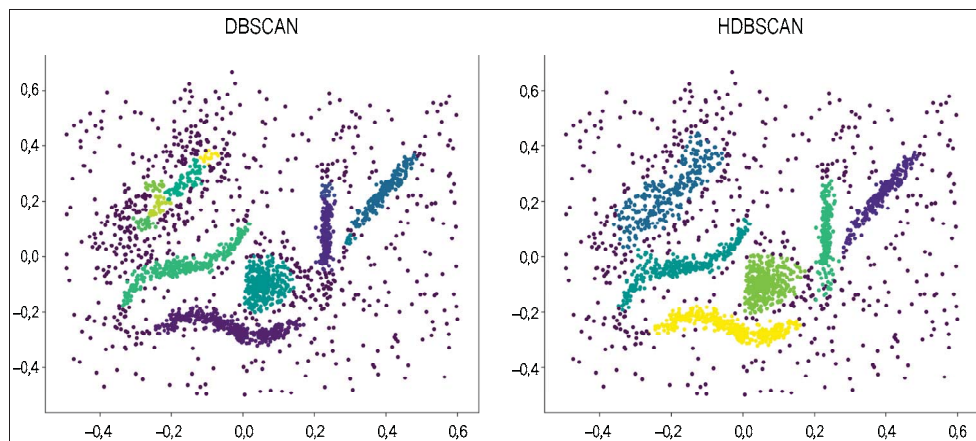


Рис. 12.18. Результаты работы алгоритмов DBSCAN и HDBSCAN

Модели гауссовых смесей

Модель гауссовой смеси (Gaussian mixture model, GMM) является генеративной моделью, которая исходит из того, что данные были сгенерированы смесью различных многомерных нормальных распределений. Данный алгоритм ориентирован на оценивание матриц средних значений и ковариаций этих распределений.

Указанный алгоритм обобщает алгоритм k средних: он добавляет ковариацию между признаками, вследствие чего кластеры могут быть эллипсоидами, а не сферами, в то время как центроиды представлены средними значениями каждого распределения. Алгоритм GMM выполняет мягкие назначения точек кластерам, поскольку каждая точка имеет вероятность быть членом любого кластера.

Алгоритм максимизации ожиданий

Алгоритм GMM использует алгоритм максимизации ожиданий (expectation-maximization algorithm), служащий для выявления компонент смеси гауссовых распределений. Цель состоит в том, чтобы усваивать параметры распределения вероятности из непомеченных данных.

Алгоритм работает итеративно следующим образом:

1. Инициализация — принять случайные центры (например, с помощью k средних).
2. Повторять следующие шаги до наступления схождения (т. е. до тех пор, пока изменения в назначениях не опустятся ниже порога):
 - *шаг ожидания*: мягкое назначение — вычислить вероятности для каждой точки из каждого распределения;
 - *шаг максимизации*: настроить параметры нормального распределения, сделав точки данных наиболее вероятными.

На рис. 12.19 показаны вероятности принадлежности кластеру GMM для набора данных Iris в виде контурных линий.

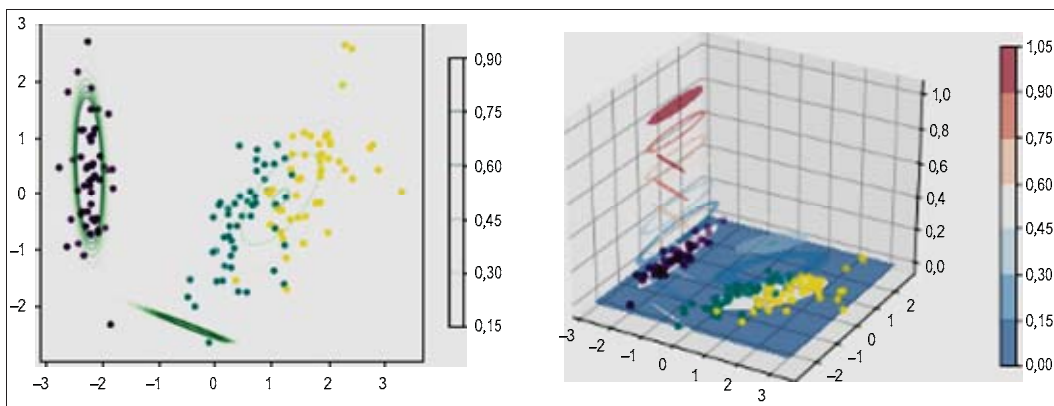


Рис. 12.19. Вероятности принадлежности кластеру GMM для набора данных Iris

Иерархический паритет риска

Ключевая идея иерархического паритета риска заключается в использовании иерархической кластеризации на матрице ковариаций для того, чтобы была возможность группировать активы с аналогичными корреляциями и уменьшать число степеней свободы, рассматривая при построении портфеля в качестве заменителей только похожие активы (подробнее см. файлы Python в подпапке `hierarchical_risk_parity`).

Первый шаг — вычислить матрицу расстояний, которая представляет близость для коррелированных активов и отвечает требованиям метрики расстояния. Полученная матрица становится входом в функцию иерархической кластеризации библиотеки SciPy, вычисляющей поочередные кластеры с помощью одного из тех методов, которые обсуждались выше:

```
def get_distance_matrix(corr):
    """Вычислить матрицу расстояний из корреляции;
    0 <= d[i,j] <= 1"""
    return np.sqrt((1 - corr) / 2)

distance_matrix = get_distance_matrix(corr)
linkage_matrix = linkage(squareform(distance_matrix), 'single')
```

Матрица `linkage_matrix` может использоваться в качестве входа в функцию `seaborn.clustermap` с целью визуализации результирующей иерархической кластеризации. Дендрограмма, отображаемая библиотекой `Seaborn`, показывает, как объединяются индивидуальные активы и кластеры активов на основе их относительных расстояний (рис. 12.20):

```
clustergrid = sns.clustermap(distance_matrix,
                             method='single',
                             row_linkage=linkage_matrix,
                             col_linkage=linkage_matrix,
                             cmap=cmap, center=0)

sorted_idx = clustergrid.dendrogram_row.reordered_ind
sorted_tickers = corr.index[sorted_idx].tolist()
```

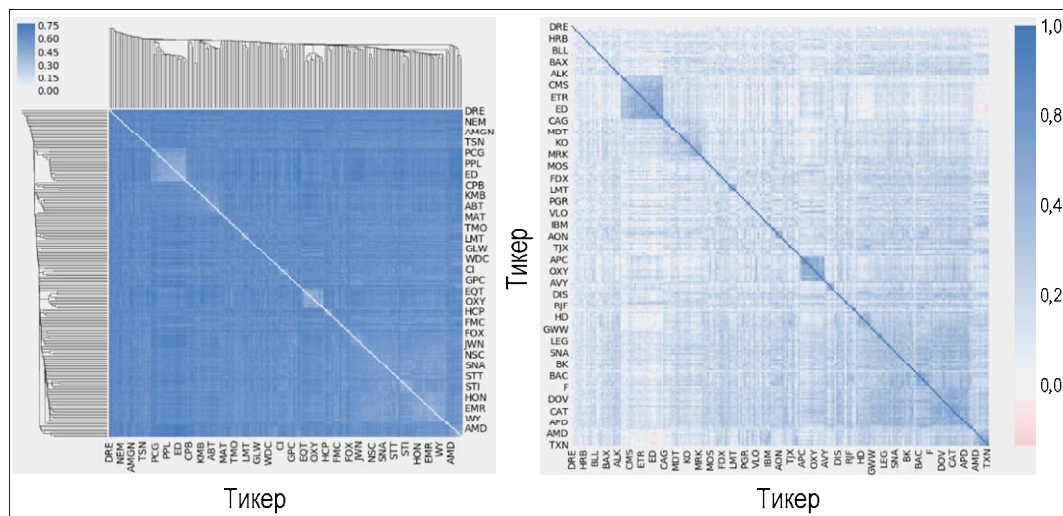


Рис. 12.20. Тепловая карта

По сравнению с тепловой картой `seaborn.heatmap` исходной матрицы корреляций теперь имеется значительно более высокая структурированность в отсортированных данных (см. рис. 12.20 справа).

Используя тикеры, отсортированные в соответствии с иерархией, индуцированной алгоритмом кластеризации, иерархический паритет риска (HRP) теперь вычисляет

нисходящее обратно-дисперсное размещение, которое поочередно регулирует веса в зависимости от дисперсии подкластеров ниже по дереву:

```
def get_cluster_var(cov, cluster_items):
    """Вычислить дисперсию в расчете на кластер"""
    cov_ = cov.loc[cluster_items, cluster_items] # срез матрицы
    w_ = get_inverse_var_pf(cov_)
    return (w_ @ cov_ @ w_).item()
```

С этой целью алгоритм использует бисекционный поиск, который позволяет распределять дисперсию кластера своим элементам на основе их относительной рискованности:

```
def get_hrp_allocation(cov, tickers):
    """Вычислить нисходящие HRP-веса"""

    weights = pd.Series(1, index=tickers)
    clusters = [tickers] # инициализировать один кластер всеми активами

    while len(clusters) > 0:
        # выполнить бисекционный поиск
        clusters = [c[start:stop] for c in clusters
                    for start, stop in ((0, int(len(c) / 2)),
                                       (int(len(c) / 2), len(c)))
                    if len(c) > 1]
        for i in range(0, len(clusters), 2): # разбивать парами
            cluster0 = clusters[i]
            cluster1 = clusters[i + 1]

            cluster0_var = get_cluster_var(cov, cluster0)
            cluster1_var = get_cluster_var(cov, cluster1)

            weight_scaler = 1 - cluster0_var / (cluster0_var +
                                                cluster1_var)

            weights[cluster0] *= weight_scaler
            weights[cluster1] *= 1 - weight_scaler

    return weights
```

Результирующее размещение финансовых средств среди портфельных активов дает веса, которые в сумме составляют 1 и отражают структуру, присутствующую в матрице корреляций (подробности см. в блокноте).

Резюме

В этой главе мы рассмотрели методы неконтролируемого обучения, которые позволяют извлекать ценный сигнал из данных, не опираясь на информацию о результатах, предоставляемую метками.

Мы увидели, как методы снижения линейной размерности, такие как анализ главных компонент (РСА) и анализ независимых компонент (ICA), могут применяться для извлечения некоррелированных или независимых компонент из данных, могущими выступать в виде рисков факторов или портфельных весов. Мы также рассмотрели передовые нелинейные технические решения по усваиванию проекций в топологическом многообразии, которые порождают ультрасовременную визуализацию сложных альтернативных наборов данных.

Во второй части мы рассмотрели несколько методов кластеризации, которые производят ведомые данными группировки в условиях различных допущений. Эти группировки могут быть полезны, например, для конструирования портфелей, в которых принципы паритета риска применяются к активам, которые были кластеризованы иерархически.

В следующих трех главах мы познакомимся с различными техническими решениями МО для ключевого источника альтернативных данных, а именно обработкой текстовых документов на естественном языке.

13

Работа

С ТЕКСТОВЫМИ ДАННЫМИ

Эта глава является первой из трех, посвященных извлечению сигналов для стратегий алгоритмической торговли из текстовых данных с использованием обработки естественного языка (natural language processing, NLP) и машинного обучения.

Текстовые данные очень богаты содержимым, но не структурированы по формату, и, следовательно, требуют дополнительной предобработки, вследствие которой автоматически обучающийся алгоритм сможет извлекать потенциальный сигнал. Ключевая проблема заключается в конвертировании текста в числовой формат для его использования алгоритмом, одновременно выражая семантику или смысл содержимого. Мы рассмотрим несколько технических решений, которые улавливают хорошо понятные людям нюансы языка, которые могут становиться входами в автоматически обучающиеся алгоритмы.

В этой главе мы введем фундаментальные технические решения для извлечения признаков, которые сосредотачиваются на индивидуальных семантических единицах, т. е. словах или коротких группах слов, именуемых *лексемами*, или токенами. Мы покажем, как представлять документы в виде векторов вхождений лексем, создав терм-документную матрицу, которая описывает частоту терминов в коллекции документов и, в свою очередь, служит входом для текстовой классификации и сентиментного анализа. Мы также введем наивный байесов алгоритм, который широко применяется для этой цели.

В следующих двух главах мы будем опираться на эти методы и использовать автоматически обучающиеся алгоритмы, такие как *тематическое моделирование* и *векторное вложение слов*, для улавливания информации, содержащейся в более широком контексте.

В частности, в этой главе мы рассмотрим следующие темы:

- ♦ как выглядит фундаментальный рабочий поток обработки естественного языка (ЕЯ);
- ♦ как строить многоязычный конвейер извлечения признаков с использованием библиотек spaCy и TextBlob;

- ◆ как выполнять задачи обработки ЕЯ, такие как *частеречная разметка* или распознавание именованных сущностей;
- ◆ как конвертировать лексемы в числа с помощью терм-документной матрицы;
- ◆ как классифицировать текст с помощью наивной байесовой модели;
- ◆ как выполнять сентиментный анализ.



Примеры исходного кода для последующих разделов находятся в репозитории GitHub для этой главы, а справочные материалы перечислены в главном файле README.

Как извлекать признаки из текстовых данных

Текстовые данные могут быть чрезвычайно ценными, учитывая те объемы информации, которыми люди обмениваются между собой и хранят с использованием естественного языка — разнообразный набор источников данных, относящихся к инвестициям, варьируется от официальных документов, таких как заявления компаний, контракты и патенты, до новостей, мнений и аналитических исследований, и даже комментариев и различных типов постов и сообщений в социальных сетях.

Многочисленные и разнообразные образцы текстовых данных доступны в режиме онлайн для разведывания возможности применения алгоритмов обработки ЕЯ, многие из которых перечислены в справочных материалах, прилагаемых к этой главе.

В целях проведения путешествия по техническим решениям и библиотекам Python, которые наиболее эффективно поддерживают реализацию данной цели, мы выделим проблемы обработки ЕЯ, представим критически важные элементы рабочего потока обработки ЕЯ и проиллюстрируем применения МО от текстовых данных до алгоритмической торговли.

Сложности обработки естественного языка

Конвертирование неструктурированного текста в машиночитаемый формат требует тщательной предобработки с целью сохранения ценных семантических аспектов данных. То, как люди выводят смысл из языка и постигают его содержимое, еще до конца не понято, и улучшение понимания языка машинами остается областью очень активных исследований.

Обработка ЕЯ представляет собой многосложную задачу, поскольку эффективное использование текстовых данных требует от МО понимания работы внутреннего механизма языка, а также знаний о мире, на который он ссылается. К числу основных сложностей относятся следующие:

- ◆ двусмысленность из-за полисемии, т. е. слово или фраза в зависимости от контекста могут иметь разные значения (например, фразу "бросившие учебу студенты сократились вдвое" можно понять по-разному);

- ◆ нестандартное и развивающееся использование языка, в особенности в социальных медиа;
- ◆ использование фразеологизмов, как, например, "выбросить полотенце на ринг" или "поднять белый флаг";
- ◆ хитроумные имена сущностей, например, "где идет „Приключения Флика“?";
- ◆ знание мира — "Мэри и Сью являются сестрами" против "Мэри и Сью являются мамами".

Рабочий поток обработки естественного языка

Ключевой целью использования автоматического обучения на текстовых данных для алгоритмической торговли является извлечение сигналов из документов. Документ представляет собой индивидуальную выборку из соответствующего источника текстовых данных, например отчета компании, новостного заголовка либо статьи или твита. Текстовый корпус, в свою очередь, представляет собой коллекцию документов.

На рис. 13.1 показаны основные шаги конвертирования документов в набор данных, который может использоваться для тренировки автоматически обучающегося алгоритма, способного выполнять пригодные на практике предсказания.

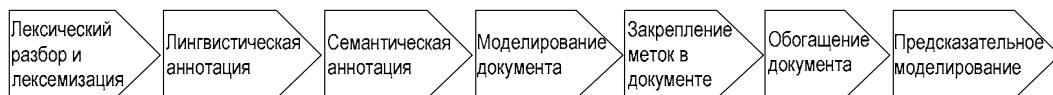


Рис. 13.1. Шаги конвертирования документов в набор данных

Фундаментальные технические решения извлекают семантические единицы, именуемые *лексемами* (токенами), в качестве текстовых признаков и используют лингвистические правила и словари для обогащения этих лексем лингвистическими и семантическими аннотациями. Модель *мешка слов* (bag-of-words, BoW) использует частоту лексем для моделирования документов в качестве векторов лексем, получая терм-документную матрицу, которая часто используется для классификации текста.

Передовые подходы используют МО для уточнения признаков, извлеченных этими фундаментальными техническими решениями, и порождают более информативные модели документов. К ним относятся тематические модели, отражающие совместное использование лексем в документах, и словарно-векторные модели, улавливающие контекст потребления лексем.

Далее мы подробнее остановимся на ключевых решениях, принимаемых на каждом шаге, и соответствующих им компромиссах, и после этого в следующем разделе проиллюстрируем их реализацию с помощью библиотеки spaCy. В табл. 13.1 обобщены ключевые задачи конвейера обработки ЕЯ.

Таблица 13.1. Ключевые задачи конвейера обработки естественного языка

| Признак | Описание |
|---|---|
| Лексемизация (токенизация) | Сегментирует текст на слова, знаки препинания и т. д. |
| Частеречная разметка (POS-теггирование) | Назначает лексемам словарные типы, такие как глагол или существительное |
| Анализ зависимостей | Разметка зависимостей синтаксических лексем, таких как субъект \Leftrightarrow объект |
| Выделение основ слов и лемматизация | Назначает базовые формы слов: was \Rightarrow be, rats \Rightarrow rat |
| Обнаружение границ предложений | Отыскивает и сегментирует индивидуальные предложения |
| Распознавание именованных сущностей | Размечает объекты реального мира, такие как люди, компании и местоположения |
| Сходство | Оценивает сходство слов, текстовых отрезков и документов |

Разбор и лексемизация текстовых данных

Лексема — это символичный экземпляр, который появляется в данном документе, и в целях дальнейшей обработки должен рассматриваться как семантическая единица. Лексикон, или вокабуляр, — это множество лексем, содержащихся в текстовом корпусе, которые считаются релевантными для дальнейшей обработки. Ключевым компромиссом в последующих решениях является точное отражение источника текста за счет расширения лексикона, что может приводить к увеличению числа признаков и усложнению модели.

Основными вариантами в этой связи является трактовка пунктуации и заглавных букв, орфографических ошибок и очень частых слов, так называемых *стоп-слов* (таких как *and* или *the*) как бессмысленного шума.

Дополнительное решение касается включения групп из n индивидуальных лексем, именуемых n -граммами, в качестве семантических единиц (индивидуальная лексема также называется *униграммой*). Примером 2-граммы (или биграммы) является "Нью-Йорк", тогда как "New York City" — это 3-грамма (или триграмма).

Цель состоит в том, чтобы создать лексемы, которые более точно отражают смысл документа. Решение может опираться на словари или сравнение относительных частот индивидуального и совместного использования. Включение n -грамм будет увеличивать число признаков, потому что число уникальных n -грамм демонстрирует тенденцию становиться выше, чем число уникальных униграмм, и с большой вероятностью будет добавлять шум, если их не отфильтровывать по значимости или частоте.

Лингвистическая аннотация

Лингвистическая аннотация охватывает применение *синтаксических* и *грамматических правил* выявления границы предложения, несмотря на неоднозначную пунктуацию и роли лексемы в предложении для частеречной разметки и анализа зависимостей. Она также позволяет выявлять общие корневые формы для выделения основ слов и лемматизации с целью группирования родственных слов.

- ◆ *Частеречные аннотации* помогают устранять неоднозначность лексем на основе их функции (это может быть необходимо, когда, например, в английском языке глагол и существительное имеют одинаковую форму), что увеличивает лексикон, но может приводить к более высокой точности.
- ◆ *Анализ зависимостей* выявляет иерархические взаимосвязи между лексемами, обычно используется для перевода с языка на язык и важен для интерактивных приложений, требующих более глубокого понимания языка, таких как виртуальные собеседники, или чатботы.
- ◆ *Выделение основ слов* использует простые правила для удаления из лексемы распространенных окончаний, таких как *s*, *ly*, *ing* и *ed* в английском языке, и сведения слова к его основе или корневой форме.
- ◆ *Лемматизация* использует более изощренные правила для получения из словоформы канонической леммы (нормальной формы) слова. Она может обнаруживать неправильные формы с окончаниями, такими как "лучше" (*-er*) и "наилучший" (*-est*), и более эффективно сжимает лексикон, но медленнее, чем выделение основ слов. Оба подхода упрощают лексикон за счет семантических нюансов.

Семантическая аннотация

Распознавание именованных сущностей (Named entity recognition, NER) призвано выявлять лексемы, которые представляют интересующие объекты, такие как люди, страны или компании. Оно может быть развито дальше, превратившись в *граф знаний*, который улавливает семантические и иерархические взаимосвязи между такими сущностями. Для приложений, которые, например, призваны предсказывать влияния новостных событий или настроений указанный ингредиент имеет критическую важность.

Закрепление меток

Многие приложения с обработкой ЕЯ учатся предсказывать результаты на основе содержательной информации, извлекаемой из текста. Контролируемое автоматическое обучение требует меток с целью научить алгоритм истинным связям между входом и выходом. Применительно к текстовым данным задача установления этой взаимосвязи может оказаться крайне сложной и потребовать явного моделирования и сбора данных.

Решения по моделированию данных включают в себя ответ на вопрос: каким образом квантифицировать настроения, коренящиеся в текстовом документе, таком как

электронное письмо, транскрибированное интервью или твит, либо какие аспекты исследовательской работы или новостного отчета закреплять за конкретным результатом?

Примеры использования

Использование МО с текстовыми данными для алгоритмической торговли опирается на извлечение значимой информации в виде признаков, которые прямо или косвенно предсказывают будущие ценовые движения. Приложения варьируются от использования краткосрочного влияния новостей на рынок до долгосрочного фундаментального анализа драйверов переоценки стоимости активов. Приведем несколько примеров:

- ♦ оценивание настроений в отзывах о продукции с целью выявления конкурентной позиции компании или отраслевых трендов;
- ♦ обнаружение аномалий в кредитных договорах с целью предсказания вероятности или последствий дефолта;
- ♦ предсказание последствий влияния новостей с точки зрения направленности, величины и затронутых сущностей.

Так, американский финансовый холдинг JP Morgan, например, разработал предсказательную модель, основанную на 250 тыс. аналитических отчетах, результативность которой превзошла несколько базовых индексов и произвела некоррелированные сигналы относительно сентиментных факторов, сформированных из консенсусных оценок по показателю корпоративных заработков в расчете на долю собственности в акционерном капитале (EPS) и изменений в рекомендациях.

От текста к лексемам — конвейер обработки естественного языка

В этом разделе мы покажем, как строить конвейер обработки ЕЯ с использованием Python-библиотеки `sraCy` с открытым исходным кодом. Библиотека `textacy` основывается на библиотеке `sraCy` и обеспечивает легкий доступ к атрибутам `sraCy` и дополнительный функционал.

Подробности относительно следующих далее примеров исходного кода, а также инструкции по установке и дополнительные сведения см. в блокноте `nlp_pipeline_with_sraCy.ipynb`.

Конвейер по обработке естественного языка с помощью библиотек `sraCy` и `textacy`

Библиотека `sraCy` — это широко используемая библиотека Python с полным комплектом функционала для быстрой обработки текста на нескольких языках. Использование механизмов лексемизации и аннотации требует установки языковых

моделей. Функционал, который мы будем использовать в этой главе, требует только малых моделей; более крупные модели также включают словарные векторы, которые мы рассмотрим в *главе 15*.

После установки и связывания мы можем создать экземпляр языковой модели библиотеки spaCy и затем вызвать ее на документе. В результате spaCy произведет объект doc, который разбивает текст на лексемы и обрабатывает его в соответствии с конфигурируемыми конвейерными компонентами, которые по умолчанию состоят из разметчика (tagger), лексического анализатора (parser) и распознавателя именovaných сущностей (ner), как на рис. 13.2.

```
nlp = spacy.load('en')
nlp.pipe_names

['tagger', 'parser', 'ner']
```

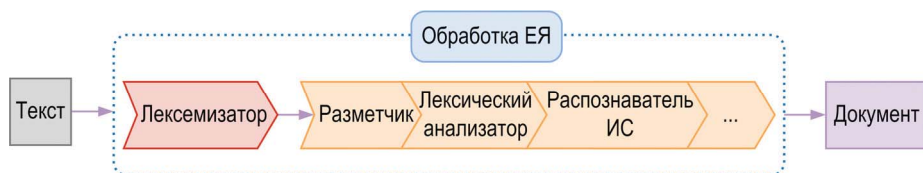


Рис. 13.2. Конвейер библиотеки SpaCy

Давайте проиллюстрируем указанный конвейер с помощью простого предложения на английском языке (в переводе: Apple планирует приобрести английский стартап за 1 млрд долларов):

```
sample_text = 'Apple is looking at buying U.K. startup for $1 billion'
doc = nlp(sample_text)
```

Разбор, лексемизирование и аннотирование предложения

Содержимое разобранного документа поддается итеративной обработке, причем каждый элемент имеет многочисленные атрибуты, порождаемые обрабатывающим конвейером. Показанный ниже пример иллюстрирует, каким образом можно обращаться к следующим атрибутам:

- ◆ text — исходный текст;
- ◆ lemma_ — корень слова;
- ◆ pos_ — простая частеречная метка;
- ◆ tag_ — подробная частеречная метка;
- ◆ dep_ — синтаксическая взаимосвязь или зависимость между лексемами;
- ◆ shape_ — форма слова относительно выделения заглавными буквами, пунктуации или цифр;

- ◆ `is_alpha` — проверка, является ли лексема буквенно-цифровой;
- ◆ `is_stop` — проверка, находится ли лексема в списке частых слов для данного языка.

Мы перебираем каждую лексему и назначаем ее атрибуты кадру данных `pd.DataFrame`:

```
pd.DataFrame([[t.text, t.lemma_, t.pos_, t.tag_,
               t.dep_, t.shape_, t.is_alpha, t.is_stop]
              for t in doc],
             columns=['text', 'lemma', 'pos', 'tag', 'dep', 'shape', 'is_alpha', 'is_stop'])
```

И на выходе получаем следующий результат (табл. 13.2).

Таблица 13.2. Результат обработки предложения

| text | lemma | pos | tag | dep | shape | is_alpha | is_stop |
|---------|---------|-------|-----|----------|-------|----------|---------|
| Apple | apple | PROPN | NNP | nsubj | Xxxxx | TRUE | FALSE |
| is | be | VERB | VBZ | aux | xx | TRUE | TRUE |
| looking | look | VERB | VBG | ROOT | xxxx | TRUE | FALSE |
| at | at | ADP | IN | prep | xx | TRUE | TRUE |
| buying | buy | VERB | VBG | pcomp | xxxx | TRUE | FALSE |
| U.K. | u.k. | PROPN | NNP | compound | X.X. | FALSE | FALSE |
| startup | startup | NOUN | NN | dobj | xxxx | TRUE | FALSE |
| for | for | ADP | IN | prep | xxx | TRUE | TRUE |
| \$ | \$ | SYM | \$ | quantmod | \$ | FALSE | FALSE |
| 1 | 1 | NUM | CD | compound | d | FALSE | FALSE |
| billion | billion | NUM | CD | pobj | xxxx | TRUE | FALSE |

Мы можем визуализировать синтаксическую зависимость в браузере или блокноте с помощью следующей строки кода:

```
displacy.render(doc, style='dep', options=options, jupyter=True)
```

Результатом является дерево зависимостей на рис. 13.3.

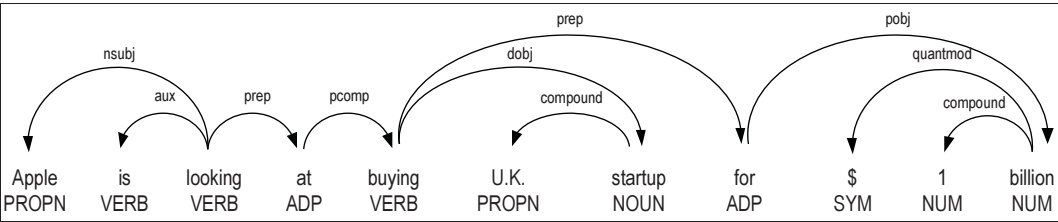


Рис. 13.3. Дерево зависимостей

Дополнительная информация о значении атрибутов может быть получена с помощью инструкции `spacy.explain()`, как здесь:

```
spacy.explain("VBZ")
```

```
verb, 3rd person singular present
```

Пакетная обработка документов

Теперь мы прочитаем более крупный набор из 2963 новостных статей вещательной корпорации BBC (подробнее об источнике данных см. в репозитории GitHub), которые принадлежат к пяти категориям и хранятся в отдельных текстовых файлах. Мы должны сделать следующее:

1. Вызвать метод `glob()` объекта `Path` библиотеки `pathlib`.
2. Перебрать полученный список путей.
3. Прочитать все строки новостной статьи, за исключением заголовка в первой строке.
4. Добавить очищенный результат в список:

```
files = (DATA_DIR / 'bbc').glob('**/*.txt')
bbc_articles = []
for i, file in enumerate(files):
    with file.open(encoding='latin1') as f:
        lines = f.readlines()
        body = ' '.join([l.strip() for l in lines[1:]])
        body = body.strip()
        bbc_articles.append(body)
```

```
len(bbc_articles)
```

```
2963
```

Определение границ предложений

Мы проиллюстрируем обнаружение предложений, вызвав объект `nlp` на первой статье:

```
doc = nlp(bbc_articles[0])
type(doc)
```

```
spacy.tokens.doc.Doc
```

Библиотека `spaCy` вычисляет границы предложений из дерева синтаксического разбора, вследствие чего пунктуация и выделение заглавными буквами играют важную, но не решающую роль. В результате границы будут совпадать с границами высказывания, даже в случае плохо размеченного знаками препинания текста.

Мы можем обратиться к разобранным предложениям с помощью атрибута `sents`:

```
sentences = [s for s in doc.sents]
sentences[:3]
```

```
[Voting is under way for the annual Bloggies which recognize the best web
blogs - online spaces where people publish their thoughts - of the year. ,
Nominations were announced on Sunday, but traffic to the official site was
so heavy that the website was temporarily closed because of too many
visitors.,
Weblogs have been nominated in 30 categories, from the top regional blog,
to the best-kept-secret blog.]
```

Распознавание именованных сущностей

Библиотека `sraCy` позволяет распознавать именованные сущности с помощью атрибута `ent_type_`:

```
for t in sentences[0]:
    if t.ent_type_:
        print('{} | {} | {}'.format(t.text, t.ent_type_, spacy.explain(t.ent_type_)))
```

```
annual | DATE | Absolute or relative dates or periods
the | DATE | Absolute or relative dates or periods
year | DATE | Absolute or relative dates or periods
```

Библиотека `textacy` обеспечивает доступ к именованным сущностям, которые появляются в первой статье:

```
from textacy.extract import named_entities

entities = [e.text for e in named_entities(doc)]
pd.Series(entities).value_counts()
```

```
year          4
US            2
annual        2
Tsunami Blog  2
South-East Asia Earthquake  2
dtype: int64
```

N-граммы

N -граммы объединяют N расположенных подряд лексем. N -граммы могут быть полезны для модели мешка слов (BoW), потому что, в зависимости от текстового контекста, обработка чего-то наподобие "data scientist" (исследователь данных), как единой лексемы, может быть более значимой, чем обработка двух разных лексем: data и scientist.

Библиотека `textacy` позволяет легко просматривать N -граммы заданной длины n , встречающиеся с частотой по крайней мере `min_freq` раз:

```
from textacy.extract import ngrams

pd.Series([n.text for n in ngrams(doc, n=2, min_freq=2)]).value_counts()

Tsunami Blog          2
annual Bloggies       2
East Asia              2
Asia Earthquake       2
dtype: int64
```

Потоковый API библиотеки `sraCy`

Для прогона более крупного числа документов через обрабатывающий конвейер можно применить потоковый API библиотеки `sraCy` следующим образом:

```
iter_texts = (bbc_articles[i] for i in range(len(bbc_articles)))
for i, doc in enumerate(nlp.pipe(iter_texts, batch_size=50, n_threads=8)):
    if i % 100 == 0:
        print(i, end = ' ')
    assert doc.is_parsed

0 100 200 300 400 500 600 700 800 900 1000 1100 1200 1300 1400 1500 1600 1700 1800 1900
2000 2100 2200
```

Многоязычная обработка естественного языка

Библиотека `sraCy` содержит натренированные языковые модели для английского, немецкого, испанского, португальского, французского, итальянского и голландского языков¹, а также многоязычную модель для распознавания именованных сущностей (NER). Перекрестно-языковое использование является простым, т. к. API не меняется.

Мы проиллюстрируем модель испанского языка, используя параллельный корпус субтитров видеозаписей TED Talk (справочные материалы со ссылкой на источник данных см. в репозитории GitHub). Для этого мы создаем экземпляры обеих языковых моделей:

```
model = {}
for language in ['en', 'es']:
    model[language] = spacy.load(language)
```

¹ В настоящее время библиотека `sraCy` поддерживает модели 9 языков. Модели отсутствующих языков можно подгрузить отдельно. См. документацию по адресу <https://spacy.io/usage/models>. — Прим. перев.

Затем мы читаем небольшие соответствующие образцы текста в каждой модели:

```
text = {}
path = Path('data/TED')
for language in ['en', 'es']:
    file_name = path / 'TED2013_sample.{}'.format(language)
    text[language] = file_name.read_text()
```

Процедура обнаружения границ предложений использует ту же логику, но находит другую разбивку:

```
parsed, sentences = {}, {}
for language in ['en', 'es']:
    parsed[language] = model[language](text[language])
    sentences[language] = list(parsed[language].sents)
print('Предложения:', language, len(sentences[language]))
```

```
Предложения: en 19
Предложения: es 22
```

Частеречная разметка работает таким же образом:

```
pos = {}
for language in ['en', 'es']:
    pos[language] = pd.DataFrame([[t.text, t.pos_,
                                   spacy.explain(t.pos_)]
                                  for t in sentences[language][0]],
                                  columns=['Token', 'POS Tag', 'Meaning'])
```

Результатом являются параллельные аннотации лексем для документов на английском и испанском языках (табл. 13.3).

Таблица 13.3. Параллельные аннотации лексем

| Лексема | Частеречная метка | Смысл | Лексема | Частеречная метка | Смысл |
|---------|-------------------|--------------------|--------------|-------------------|----------------|
| There | ADV | наречие | Existe | VERB | глагол |
| s | VERB | глагол | una | DET | определитель |
| a | DET | определитель | estrecha | ADJ | прилагательное |
| tight | ADJ | прилагательное | y | CONJ | союз |
| and | CCONJ | сочинительный союз | sorprendente | ADV | прилагательное |

В следующем разделе показано, как использовать разобранные и аннотированные лексемы для построения терм-документной матрицы, которая может использоваться для классификации текста.

Обработка естественного языка с помощью библиотеки TextBlob

TextBlob — это Python-библиотека, которая предоставляет простой API для общих задач обработки ЕЯ и основывается на естественно-языковом инструментарии библиотек *NLTK* (Natural Language Toolkit) для обработки ЕЯ и *Pattern* для глубинного анализа Всемирной паутины. Библиотека TextBlob обеспечивает частеречную разметку, извлечение именных групп, sentimentный анализ, классификацию, машинный перевод и многое другое.

Для иллюстрации использования библиотеки TextBlob мы отбираем спортивные статьи BBC с заголовком "Ronaldo considering new contract" (Рональдо подумывает о новом контракте). Подобно библиотеке spaCy и другим библиотекам, первым шагом является прогон документа через конвейер, представляемый объектом TextBlob с целью назначения аннотаций, необходимых для различных задач (см. блокнот `nlp_with_textblob.ipynb` для этого раздела):

```
from textblob import TextBlob

DATA_DIR = Path('../data')
path = Path(DATA_DIR, 'bbcsport')
files = path.glob('**/*.txt')
doc_list = []
for i, file in enumerate(files):
    topic = file.parts[-2]
    article = file.read_text(encoding='latin1').split('\n')
    heading = article[0].strip()
    body = ' '.join([l.strip() for l in article[1:]]).strip()
    doc_list.append([topic, heading, body])

docs = pd.DataFrame(doc_list, columns=['topic', 'heading', 'body'])
article = docs.sample(1).squeeze()
parsed_body = TextBlob(article.body)
```

Выделение основ слов и лемматизация

Для выделения основ слов мы создаем экземпляр стеммера *SnowballStemmer* для выделения основ слов из библиотеки *NLTK*, на каждой лексеме вызываем его метод `stem()` и показываем модифицированные лексемы:

```
from nltk.stem.snowball import SnowballStemmer

stemmer = SnowballStemmer('english')
[(word, stemmer.stem(word))
 for i, word in enumerate(parsed_body.words)
 if word.lower() != stemmer.stem(parsed_body.words[i])]
```

```
[('Manchester', 'manchest'),
 ('United', 'unit'),
 ('agreeing', 'agre'),
 ('Portugal', 'portug'),
 ('joined', 'join'),
 ...
 ('nobody', 'nobodi'),
 ('knows', 'know'),
 ('future', 'futura')]
```

И затем лемматизируем слова, вызывая библиотечный метод `lemmatize()`:

```
[(word, word.lemmatize())
 for i, word in enumerate(parsed_body.words)
 if word != parsed_body.words[i].lemmatize()]
```

```
[('has', 'ha'),
 ('sides', 'side'),
 ('knows', 'know'),
 ('knows', 'know'),
 ('clubs', 'club'),
 ('knows', 'know')]
```

Полярность и субъективность настроений

Библиотека `TextBlob` предоставляет отметки полярности и субъективности для разобранных документов с помощью словарей, предоставляемых библиотекой `Pattern`. Эти словари увязывают прилагательные, часто встречающиеся в отзывах о продукции, с отметками полярности настроений в промежутке от -1 до $+1$ (отрицательный \leftrightarrow положительный) и аналогичной отметкой субъективности (объективный \leftrightarrow субъективный).

Атрибут `sentiment` предоставляет среднее значение для каждого из них на соответствующих лексемах, тогда как атрибут `sentiment_assessments` показывает базовые значения для каждой лексемы (см. блокнот):

```
parsed_body.sentiment
```

```
Sentiment(polarity=0.2487603305785124, subjectivity=0.38013350286077563)
```

От лексем к числам — терм-документная матрица

В этом разделе мы сначала рассмотрим то, как модель мешка слов (BoW) приводит текстовые данные к представлению в виде числового векторного пространства, которое позволяет сравнивать документы, используя их расстояние. Затем мы покажем, как создавать терм-документную матрицу с помощью библиотеки `sklearn`.

Модель мешка слов

Модель мешка слов (BoW) представляет документ, основываясь на частоте содержащихся в нем терминов, или лексем. Каждый документ становится вектором с одним элементом для каждой лексемы в лексиконе, служа отражением релевантности лексемы документу.

При наличии лексикона терм-документная матрица вычисляется очень просто. Вместе с тем она также является грубым упрощением, поскольку абстрагируется от порядка слов и грамматических взаимосвязей. Тем не менее она зачастую быстро дает хорошие результаты в задачах текстовой классификации и поэтому является очень полезной отправной точкой.

На рис. 13.4 слева показано, как эта модель документа конвертирует текстовые данные в матрицу с числовыми элементами, где каждая строка соответствует документу, а каждый столбец — лексеме в лексиконе. Результирующая матрица обычно является одновременно очень высокоразмерной и разреженной, т. е. содержит много нулевых элементов, потому что большинство документов содержит только небольшую часть совокупного лексикона.

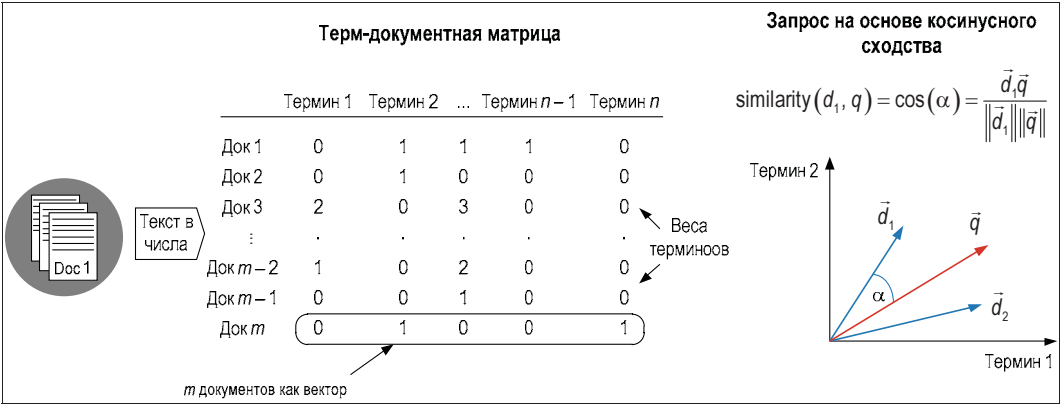


Рис. 13.4. Результирующая терм-документная матрица

Для улавливания релевантности лексемы существует несколько способов взвешивания векторного элемента лексемы документа. Мы проиллюстрируем применение библиотеки `sklearn` с использованием бинарных флагов, которые указывают на наличие или отсутствие, а также количества вхождений и взвешенные количества вхождений, которые отражают различия в частотах терминов во всех документах; т. е. в текстовом корпусе.

Измерение сходства документов

Представление документов в виде векторов слов назначает каждому документу местоположение в векторном пространстве, созданном лексиконом. Интерпретируя векторные элементы как декартовы координаты в этом пространстве, мы можем

использовать угол между двумя векторами для измерения их сходства, т. к. векторы, которые указывают в одном и том же направлении, содержат одни и те же члены с одинаковыми частотными весами.

Предыдущий рис 13.4 (справа) иллюстрирует — упрощенно в двух размерностях — вычисление расстояния между документом, представленным вектором \mathbf{d}_i и вектором запроса (множеством поисковых терминов либо еще одним документом), представленным вектором \mathbf{q} .

Косинусное сходство равно косинусу угла между двумя векторами. Оно транслирует размер угла в число в промежутке $[0; 1]$, т. к. все векторные элементы являются неотрицательными весами лексем. Значение 1 означает, что оба документа идентичны относительно их лексем, тогда как значение 0 означает, что два документа содержат только несовпадающие лексем.

Как показано на рисунке, косинус угла равен скалярному произведению векторов, т. е. сумме произведений их координат, деленной на произведение длин, измеренных евклидовыми нормами каждого вектора.

Терм-документная матрица с помощью библиотеки `sklearn`

Предобрабатывающий модуль библиотеки `sklearn` предлагает два инструмента для создания терм-документной матрицы. Класс `CountVectorizer` использует двоичные или абсолютные количества вхождений, с помощью их измеряя *частоту термина* $\text{tf}(d, t)$ для каждого документа d и лексемы t .

Класс `TfidfVectorizer`, напротив, взвешивает (абсолютную) частоту термина на *обратную документную частоту* (inverse document frequency, *idf*). В результате термин, который появляется в большем числе документов, получит меньший вес, чем термин с той же самой частотой для определенного документа, но с меньшей частотой во всех документах. Более конкретно, используя параметры по умолчанию, элементы $\text{tf-idf}(d, t)$ для терм-документной матрицы вычисляются как:

$$\text{tf-idf}(d, t) = \text{tf}(d, t) \cdot \text{idf}(t);$$

$$\text{idf}(t) = \log \frac{1 + n_d}{1 + \text{df}(d, t)} + 1.$$

Здесь n_d — это число документов, а $\text{df}(d, t)$ — документная частота термина t . Результирующие векторы tf-idf для каждого документа нормализуются относительно их абсолютных или квадратичных итогов (подробнее см. документацию библиотеки `sklearn`). Мера tf-idf первоначально использовалась в технологии информационного поиска для ранжирования результатов поисковой системы и впоследствии оказалась полезной для текстовой классификации и кластеризации.

Оба инструмента используют один и тот же интерфейс и выполняют лексемизацию и дальнейшую предобработку списка документов перед векторизацией текста

путем генерирования количеств вхождений лексем для заполнения терм-документной матрицы.

Ключевые параметры, влияющие на размер лексикона, перечислены ниже:

- ◆ `stop_words` — использовать встроенный или предоставить собственный список (частых) слов, подлежащих исключению;
- ◆ `ngram_range` — включать n -граммы в интервале для n , определенном кортежем (n_{\min}, n_{\max}) ;
- ◆ `lowercase` — конвертировать символы соответственно в нижний или верхний регистр (по умолчанию `True`);
- ◆ `min_df/max_df` — игнорировать слова, которые появляются в меньшем/большем числе (`int`) документов или меньшей/большей доле документов (если `float`, то промежуток равен $[0, 0; 1, 0]$);
- ◆ `max_features` — соответствующим образом лимитировать число лексем в лексиконе;
- ◆ `binary` — установить ненулевые частоты, равными 1 (`True`).

Следующие ниже примеры исходного кода и дополнительные сведения см. в блокноте `document_term_matrix.ipynb`. В качестве иллюстрации мы снова используем 2963 новостных статей вещательной корпорации BBC.

Применение класса *CountVectorizer*

Указанный выше блокнот содержит интерактивную визуализацию, которая разведывает влияние значений параметров `min_df` и `max_df` на размер лексикона. Мы считываем статьи в кадр данных, конфигурируем экземпляр класса `CountVectorizer` для создания бинарных флагов, используем все лексемы и вызываем его метод `fit_transform()`, получая терм-документную матрицу:

```
binary_vectorizer = CountVectorizer(max_df=1.0,
                                   min_df=1,
                                   binary=True)

binary_dtm = binary_vectorizer.fit_transform(docs.body)

binary_dtm
<2963x30652 sparse matrix of type '<class 'numpy.int64''>'
  with 581397 stored elements in Compressed Sparse Row format>
```

На выходе получится разреженная матрица `scipy.sparse` в строчном формате, которая эффективно хранит небольшую долю (<0,5%) из 581 397 ненулевых элементов в 2963 (документных) строках и 30 652 (лексемных) столбцах.

Визуализирование распределения лексикона

Визуализация показывает, что требование к лексемам может появляться по крайней мере в 1% и менее 50% документов ограничивает лексикон примерно до 10% из почти 30 тыс. лексем.

В результате на рис. 13.5 слева мы получаем моду чуть более 100 уникальных лексем на документ, а справа показана гистограмма документных частот для остальных лексем.

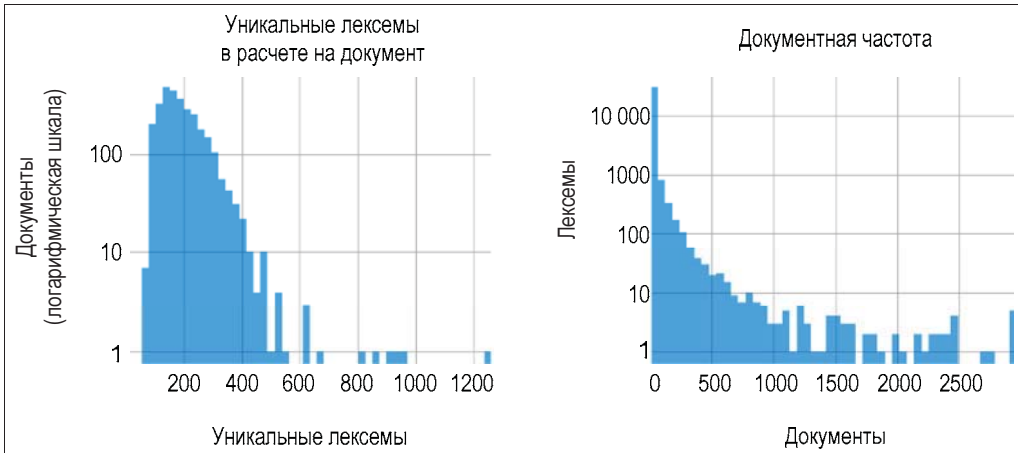


Рис. 13.5. Распределение терм-документных частот

Отыскание наиболее похожих документов

Результат векторизатора `CountVectorizer` позволяет отыскивать наиболее похожие документы с использованием функции `pdist()` для попарных расстояний, предоставляемых модулем `scipy.spatial.distance`. Она возвращает сжатую матрицу расстояний, элементы которой соответствуют верхнему треугольнику квадратной матрицы. Мы используем функцию `np.triu_indices()`, транслируя индекс, минимизирующий расстояние в строчные и столбцовые индексы, которые в свою очередь соответствуют векторам ближайших лексем (функция `np.triu_indices()` возвращает индексы верхнего треугольника массива):

```
m = binary_dtm.todense() # pdist не принимает разряженный формат
pairwise_distances = pdist(m, metric='cosine')

closest = np.argmin(pairwise_distances) # индекс, минимизирующий расстояние
rows, cols = np.triu_indices(n_docs) # получить строково-столбцовые индексы
rows[closest], cols[closest]
```

(11, 75)

Статьи № 11 и 75 являются наиболее близкими по косинусному сходству (табл. 13.4), потому что они совместно используют 58 лексем (см. блокнот).

Таблица 13.4. Статьи с наиболее близким косинусным сходством

| Тема | технологии | технологии |
|-----------|--|--|
| заголовок | Software watching while you work | BT program to beat dialer scams |
| тело | Software that can not only monitor every keystroke and action performed at a PC but can also be used as legally binding evidence of wrongdoing has been unveiled. Worries about cyber-crime and sabotage have prompted many employers to consider monitoring employees | BT is introducing two initiatives to help beat rogue dialer scams, which can cost dial-up net users thousands. From May, dial-up net users will be able to download free software to stop computers using numbers not on a user's pre-approved list |
| заголовок | Программа наблюдает, как вы работаете | Программа BT побьет мошеннических автонаборщиков номеров |
| тело | Представлена программа, которая не только способна контролировать каждое нажатие клавиши и действие, выполняемое на ПК, но и может использоваться в качестве юридически обязательного доказательства неправомерных действий. Беспокойство по поводу киберпреступности и саботажа побудило многих работодателей рассматривать возможность мониторинга сотрудников | BT представляет две инициативы, которые помогают победить мошеннические программы, автоматического набора номеров, которые могут стоить тысячи долларов пользователям с коммутируемым доступом к сети. С мая пользователи коммутируемых соединений смогут скачивать бесплатное программное обеспечение, которое прекратит использование компьютерами номеров, не включенных в предварительно утвержденный пользователем список |

Оба векторизатора, и `CountVectorizer`, и `TfidfVectorizer`, можно использовать с библиотекой `spaCy`. Например, для выполнения лемматизации и исключения определенных символов во время лексемизации мы используем следующее:

```
nlp = spacy.load('en')
def tokenizer(doc):
    return [w.lemma_ for w in nlp(doc)
            if not w.is_punct | w.is_space]
vectorizer = CountVectorizer(tokenizer=tokenizer, binary=True)
doc_term_matrix = vectorizer.fit_transform(docs.body)
```

Дополнительные сведения и примеры см. в блокноте.

Классы *TfidfTransformer* и *TfidfVectorizer*

Класс `TfidfTransformer` вычисляет веса `tf-idf` из терм-документной матрицы лексемных частот, в частности той, которая была произведена векторизатором `CountVectorizer`.

Класс `TfidfVectorizer` выполняет оба вычисления за один шаг. Он добавляет несколько параметров в API векторизатора `CountVectorizer`, который управляет сглаживающим поведением.

Вычисление меры `tf-idf` для небольшого образца текста (в переводе: позвоню тебе завтра; вызови мне такси; пожалуйста, позвони мне..., ну, пожалуйста) работает следующим образом:

```
sample_docs = ['call you tomorrow',
               'Call me a taxi',
               'please call me... PLEASE!']
```

Мы вычисляем лексемную частоту, как мы только что сделали:

```
vectorizer = CountVectorizer()
tf_dtm = vectorizer.fit_transform(sample_docs).todense()
tokens = vectorizer.get_feature_names()
term_frequency = pd.DataFrame(data=tf_dtm, columns=tokens)
print(term_frequency)
```

| | call | me | please | taxi | tomorrow | you |
|---|------|----|--------|------|----------|-----|
| 0 | 1 | 0 | 0 | 0 | 1 | 1 |
| 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| 2 | 1 | 1 | 2 | 0 | 0 | 0 |

Документная частота — это число документов, содержащих лексему:

```
vectorizer = CountVectorizer(binary=True)
df_dtm = vectorizer.fit_transform(sample_docs).todense().sum(axis=0)
document_frequency = pd.DataFrame(data=df_dtm, columns=tokens)
print(document_frequency)
```

| | call | me | please | taxi | tomorrow | you |
|---|------|----|--------|------|----------|-----|
| 0 | 3 | 2 | 1 | 1 | 1 | 1 |

Веса `tf-idf` — это соотношение этих значений:

```
tfidf = pd.DataFrame(data=tf_dtm/df_dtm, columns=tokens)
print(tfidf)
```

| | call | me | please | taxi | tomorrow | you |
|---|------|------|--------|------|----------|------|
| 0 | 0.33 | 0.00 | 0.00 | 0.00 | 1.00 | 1.00 |
| 1 | 0.33 | 0.50 | 0.00 | 1.00 | 0.00 | 0.00 |
| 2 | 0.33 | 0.50 | 2.00 | 0.00 | 0.00 | 0.00 |

Эффект сглаживания

Во избежание деления на ноль векторизатор `TfidfVectorizer` использует сглаживание для документных и лексемных частот:

- ◆ `smooth_idf` — добавить 1 к документной частоте, как если бы дополнительный документ содержал каждую лексему в лексиконе, тем самым предотвращая деления на ноль;

- ◆ `sublinear_tf` — применить сублинейное шкалирование меры `tf`; другими словами, заменить `tf` на $1 + \log(tf)$.

В сочетании с нормированными весами результаты отличаются незначительно:

```
vect = TfidfVectorizer(smooth_idf=True,
                      norm='l2', # квадратичные веса
                          # в сумме составляют 1 по документу
                      sublinear_tf=False, # если True, то использовать 1+log(tf)
                      binary=False)

print(pd.DataFrame(vect.fit_transform(sample_docs).todense(),
                  columns=vect.get_feature_names()))
```

| | call | me | please | taxi | tomorrow | you |
|---|------|------|--------|------|----------|------|
| 0 | 0.39 | 0.00 | 0.00 | 0.00 | 0.65 | 0.65 |
| 1 | 0.43 | 0.55 | 0.00 | 0.72 | 0.00 | 0.00 |
| 2 | 0.27 | 0.34 | 0.90 | 0.00 | 0.00 | 0.00 |

Как обобщать новостные статьи с помощью класса *TfidfVectorizer*

Из-за способности векторов `tf-idf` назначать лексемам значимые веса они также используются для обобщения текстовых данных. Например, функция `autotldr` на социально-новостном сайте Reddit основана на аналогичном алгоритме. Пример использования новостных статей BBC см. в блокноте.

Предобработка текста — краткий обзор

Большое число технических приемов обработки естественного языка для его использования в автоматически обучающихся моделях, которые мы представили в этом разделе, необходимы для решения многосложной природы этого высоко неструктурированного источника данных. Выработка хороших языковых признаков является трудным и полезным и, возможно, самым важным шагом в разблокировании семантического значения, скрытого в текстовых данных.

На практике опыт помогает отбирать преобразования, которые в большинстве своем удаляют шум, чем сигнал, но, так или иначе, необходимость в перекрестном контроле и сравнении результативности разных комбинаций вариантов предобработки остается.

Классификация текста и сентиментный анализ

После конвертирования текстовых данных в числовые признаки с использованием методов обработки ЕЯ, рассмотренных в предыдущих разделах, процедура классификации текста работает так же, как и любая другая классификационная задача.

В этом разделе мы применим это техническое решение предобработки текста к новостным статьям, отзывам о продукции и данным Twitter и научим вас работе с раз-

личными классификаторами с целью предсказания дискретных категорий новостей, отметок отзывов и полярности настроений.

Прежде всего, мы введем наивную байесову модель, т. е. алгоритм вероятностной классификации, который хорошо работает с текстовыми признаками, производимыми моделью мешка слов.



Примеры исходного кода для этого раздела находятся в блокноте `text_classification.ipynb`.

Наивный байесов классификатор

Наивный байесов алгоритм очень популярен для классификации текста, потому что его низкая вычислительная стоимость и потребности в памяти обеспечивают тренировку на очень крупных, высокоразмерных наборах данных. Его предсказательная результативность может конкурировать с более сложными моделями, обеспечивает хорошую отправную точку и наиболее известна своим успешным обнаружением спама.

Указанная модель основана на теореме Байеса (см. главу 9) и допущении, что различные признаки взаимно независимы при заданном классе результатов. Другими словами, для заданного результата сведения о значении одного признака (например, о наличии лексемы в документе) не дает никакой информации о значении другого признака.

Памятка по теореме Байеса

Теорема Байеса выражает условную вероятность одного события (например, что письмо является спамным в отличие от неспамного) при условии, что наступает другое событие (например, что письмо содержит определенные слова), следующим образом:

$$\begin{aligned}
 & \underbrace{P(\text{является спамом} \mid \text{имеет слова})}_{\text{апостериорная вероятность}} = \\
 & = \frac{\overbrace{P(\text{имеет слова} \mid \text{является спамом})}^{\text{правдоподобие}} \overbrace{P(\text{является спамом})}^{\text{априорная вероятность}}}{\underbrace{P(\text{имеет слова})}_{\text{наблюдение}}}.
 \end{aligned}$$

Апостериорная (после опыта) вероятность того, что письмо на самом деле является спамом при условии, что оно содержит определенные слова, зависит от взаимодействия трех факторов:

- ◆ *априорная* (до опыта) вероятность того, что письмо является спамом;
- ◆ *правдоподобие* (правдоподобная вероятность) встретить эти слова в спамном электронном письме;
- ◆ *наблюдение*, т. е. вероятность увидеть эти слова в электронном письме.

Для того чтобы вычислить апостериорную вероятность, мы можем проигнорировать наблюдение, потому что оно является одинаковым для всех результатов (спамные сообщения против неспамных), и безусловная априорная вероятность может быть легко вычислена.

Тем не менее, правдоподобие создает непреодолимые сложности для лексикона значительного размера и реального корпуса электронных писем. Причина — комбинаторный взрыв числа слов, которые появлялись или не появлялись совместно в разных документах и которые не дают выполнить оценивание, необходимое для вычисления таблицы вероятностей и закрепления значения за правдоподобием.

Допущение об условной независимости

Допущение, которое делает модель податливой и оправдывает ее название как наивной, заключается в том, что признаки не зависят от результата. В качестве иллюстрации давайте отклассифицируем электронное письмо с тремя словами "отправляй деньги сейчас", вследствие чего теорема Байеса принимает следующий вид:

$$P(\text{спам} | \text{отправляй деньги сейчас}) = \frac{P(\text{отправляй деньги сейчас} | \text{спам}) P(\text{спам})}{P(\text{отправляй деньги сейчас})}.$$

Формально допущение, что эти три слова являются условно независимыми, означает, что вероятность наблюдать термин "отправляй" не зависит от наличия других терминов при условии, что сообщение является спамным; другими словами, $P(\text{отправляй} | \text{деньги, сейчас, спам}) = P(\text{отправляй} | \text{спам})$. В результате мы можем упростить функцию правдоподобия:

$$\begin{aligned} P(\text{спам} | \text{отправляй деньги сейчас}) &= \\ &= \frac{P(\text{отправляй} | \text{спам}) P(\text{деньги} | \text{спам}) P(\text{сейчас} | \text{спам}) P(\text{спам})}{P(\text{отправляй деньги сейчас})}. \end{aligned}$$

Используя наивное допущение об условной независимости, каждый член числителя можно вычислять как относительные частоты из тренировочных данных. Знаменатель является постоянным для всех классов и может игнорироваться, когда апостериорные вероятности необходимо сравнивать, а не калибровать. Априорная вероятность становится менее актуальной по мере увеличения числа факторов, т. е. признаков.

Подведем итог: преимущества наивной байесовой модели заключаются в быстрой тренировке и быстром предсказании, поскольку число параметров является линейным по числу признаков, а их оценивание вместо дорогостоящей итерационной оптимизации имеет решение в замкнутой форме (на основе частот тренировочных данных). Она также интуитивно понятна и до некоторой степени интерпретируема, не требует гиперпараметрической настройки и относительно робастна к нерелевантным признакам при наличии достаточного сигнала.

Однако, когда допущение о независимости не соблюдается, и классификация текста зависит от комбинаций признаков либо корреляции признаков, указанная модель будет работать плохо.

Классифицирование новостных статей

Мы начинаем с иллюстрации наивной байесовой модели для классифицирования новостных статей с использованием статей BBC, которые мы считываем, как и раньше, и получаем кадр данных с 2963 статьями в пяти категориях:

```
docs = pd.DataFrame(doc_list, columns=['topic', 'heading', 'body'])
docs.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2963 entries, 0 to 2962
Data columns (total 3 columns):
topic      2963 non-null object
heading    2963 non-null object
body       2963 non-null object
dtypes: object(3)
memory usage: 69.5+ KB
```

Тренировка и оценивание мультиномного наивного байесова классификатора

Мы разбиваем данные на тренировочный и тестовый наборы по умолчанию в пропорции 75:25, обеспечивая, чтобы классы тестового набора близко отражали тренировочный набор:

```
y = pd.factorize(docs.topic)[0] # создать целочисленные значения классов
X = docs.body
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    random_state=1,
                                                    stratify=y)
```

Мы переходим к заучиванию лексикона из тренировочного набора и преобразованию обоих наборов данных с помощью векторизатора `CountVectorizer` с настройками по умолчанию, в результате получая почти 26 тыс. признаков:

```
vectorizer = CountVectorizer()
X_train_dtm = vectorizer.fit_transform(X_train)
X_test_dtm = vectorizer.transform(X_test)
X_train_dtm.shape, X_test_dtm.shape

((1668, 25919), (557, 25919))
```

Тренировка и предсказание следуют стандартному интерфейсу подгонки/предсказания библиотеки `sklearn`:

```
nb = MultinomialNB()
nb.fit(X_train_dtm, y_train)
y_pred_class = nb.predict(X_test_dtm)
```

Мы оцениваем многоклассовые предсказания с использованием метрики точности и обнаруживаем, что классификатор с настройками по умолчанию достиг почти 98%:

```
accuracy_score(y_test, y_pred_class)

0.97666068222621
```

Сентиментный анализ

Сентиментный анализ, или анализ настроений, представляет собой одно из самых популярных применений технологии обработки ЕЯ и машинного обучения в торговле на финансовых рынках, поскольку положительные или отрицательные перспективы финансовых активов или других ценовых движущих сил могут влиять на возвратность.

Как правило, модельные подходы к сентиментному анализу опираются на словари, такие как в библиотеке TextBlob, либо модели, натренированные на результатах, специфичных для конкретной области. Последнее предпочтительнее, поскольку позволяет выполнять более целенаправленное закрепление меток, например, привязывая текстовые признаки не к косвенным сентиментным отметкам, а к последующим ценовым изменениям.

Мы проиллюстрируем автоматическое обучение для анализа настроений на основе набора данных социальной сети Twitter с бинарными метками полярности и крупного набора данных отзывов о деятельности предприятий с веб-сайта Yelp с пятибалльной шкалой результатов.

Данные социальной сети Twitter

Мы используем набор данных, который содержит 1,6 млн тренировочных и 350 тестовых твитов за 2009 г. с алгоритмически назначенными бинарными положительными и отрицательными сентиментными отметками, которые довольно равномерно разделены (подробности разведывательного анализа данных см. в соответствующем блокноте).

Мультиномный наивный Байес

Мы создаем терм-документную матрицу с 934 лексемами следующим образом:

```
vectorizer = CountVectorizer(min_df=.001, max_df=.8, stop_words='english')
train_dtm = vectorizer.fit_transform(train.text)
```

```
<1566668x934 sparse matrix of type '<class 'numpy.int64'>'
  with 6332930 stored elements in Compressed Sparse Row format>
```

Затем мы тренируем классификатор `MultinomialNB`, как и раньше, и предсказываем тестовый набор:

```
nb = MultinomialNB()
nb.fit(train_dtm, train.polarity)
predicted_polarity = nb.predict(test_dtm)
```

Результат имеет точность более 77,5%:

```
accuracy_score(test.polarity, y_pred_class)
```

```
0.7768361581920904
```

Сравнение с sentimentными отметками объекта *TextBlob*

Мы также получаем sentimentные отметки модели *TextBlob* для твитов и отмечаем (см. диаграмму на рис. 13.6 слева), что положительные тестовые твиты получают значительно более высокую sentimentную отметку. Затем мы используем модель *MultinomialNB* и метод `predict_proba()` для вычисления предсказываемых вероятностей и сравнения обеих моделей, используя соответствующие кривые ROC (см. диаграмму на рис. 13.6 справа).

В данном случае наивная байесова модель превосходит модель *TextBlob*.

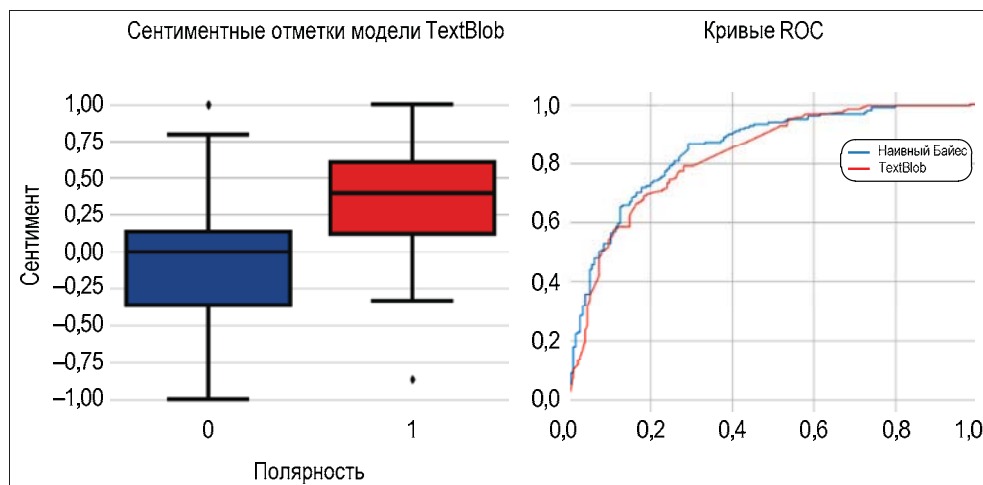


Рис. 13.6. Сентиментные отметки двух моделей в сравнении

Отзывы о деятельности предприятий — конкурсные данные веб-сайта Yelp

Наконец, мы применим sentimentный анализ к значительно более крупному набору данных отзывов о деятельности предприятий с веб-сайта Yelp с пятью классами результатов. Данные состоят из нескольких файлов с информацией о предприятии, пользователе, отзыве и других аспектах, которые веб-сайт Yelp предоставляет для поощрения инноваций в области науки о данных.

Мы будем использовать около 6 млн отзывов, подготовленных за период 2010–2018 гг. (подробнее см. соответствующий блокнот). На диаграммах рис. 13.7 показаны число отзывов и среднее число звезд в год.

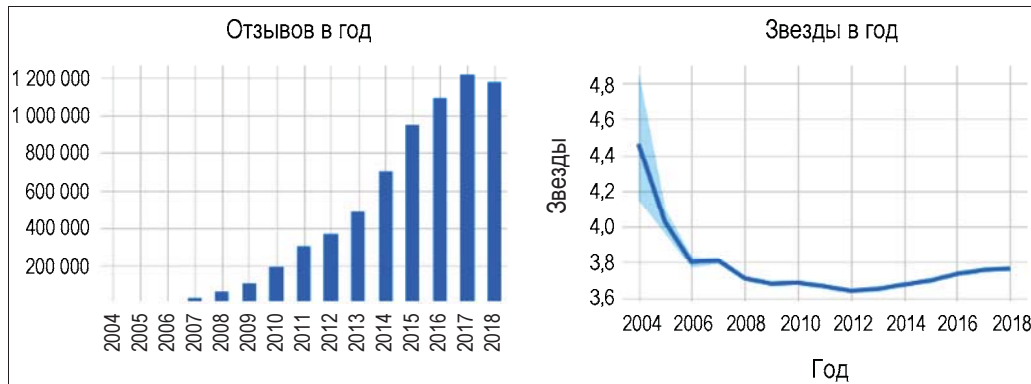


Рис. 13.7. Графики, показывающие число отзывов и среднее число звезд в год

В дополнение к текстовым признакам, получаемым из текстов отзывов, мы также будем использовать другую информацию, представляемую с отзывом или о пользователе.

Мы натренируем различные модели на данных за 2017 г. и будем использовать данные за 2018 г. в качестве тестового набора.

Эталонная точность

Используя наиболее частое число звезд ($=5$) для предсказания тестового набора, мы достигаем точности, близкой к 52%:

```
test['predicted'] = train.stars.mode().iloc[0]
accuracy_score(test.stars, test.predicted)
```

```
0.5196950594793454
```

Мультиномная наивная байесова модель

Затем мы тренируем наивный байесов классификатор, используя терм-документную матрицу, созданную векторизатором `CountVectorizer` с настройками по умолчанию:

```
nb = MultinomialNB()
nb.fit(train_dtm, train.stars)
predicted_stars = nb.predict(test_dtm)
```

Предсказание дает точность 64,7% на тестовом наборе, 24,4%-е улучшение по сравнению с эталоном:

```
accuracy_score(test.stars, predicted_stars)
```

```
0.6465164206691094
```

Логистическая регрессия по схеме "один против всех"

Мы приступаем к тренировке логистической регрессии по схеме "один против всех", которая тренирует одну модель для каждого класса, рассматривая остальные классы как отрицательный класс, и предсказывает вероятности для каждого класса с использованием разных моделей.

Используя только текстовые признаки, мы тренируем и оцениваем модель следующим образом:

```
logreg = LogisticRegression(C=1e9)
logreg.fit(X=train_dtm, y=train.stars)
y_pred_class = logreg.predict(test_dtm)
```

Модель достигает значительно более высокой точности на уровне 73,6%:

```
accuracy_score(test.stars, y_pred_class)

0.7360498864740219
```

Комбинирование текстовых и числовых признаков

Набор данных содержит различные числовые признаки (подробности реализации см. в соответствующем блокноте).

Векторизаторы производят метрики из модуля `scipy.sparse`. Для того чтобы объединить векторизованные текстовые данные с другими признаками, нужно сначала конвертировать их в разреженные матрицы; многие объекты библиотеки `sklearn` и других библиотек, таких как `LightGBM`, способны обрабатывать эти очень эффективные по потреблению памяти структуры данных. Конвертирование разреженной матрицы в плотный массив `NumPy` рискует достичь переполнения памяти.

Большинство переменных являются категориальными, и поэтому мы используем кодирование с одним активным состоянием, т. к. у нас достаточно крупный набор данных, чтобы уместить в памяти увеличившееся число признаков.

Мы конвертируем кодированные числовые признаки и объединяем их с терм-документной матрицей:

```
train_numeric = sparse.csr_matrix(train_dummies.astype(np.int8))
train_dtm_numeric = sparse.hstack((train_dtm, train_numeric))
```

Мультиномная логистическая регрессия

Логистическая регрессия также обеспечивает мультиномный вариант тренировки, который быстрее и точнее, чем реализация по схеме "один против всех". Мы используем решатель `lbfgs` (подробнее см. документацию библиотеки `sklearn`, ссылки на которую можно найти в репозитории [GitHub](#)):

```
multi_logreg = LogisticRegression(C=1e9,
                                   multi_class='multinomial',
                                   solver='lbfgs')
multi_logreg.fit(train_dtm_numeric.astype(float), train.stars)
y_pred_class = multi_logreg.predict(test_dtm_numeric.astype(float))
```

Эта модель доводит результативность до 74,6%-й точности:

```
accuracy_score(test.stars, y_pred_class)
0.7464488070176475
```

В этом случае настройка регуляризационного параметра `c` не привела к очень значительным улучшениям (см. блокнот).

Градиентно-бустинговая машина

Для целей иллюстрации мы также тренируем градиентно-бустинговый древесный ансамбль LightGBM с настройками по умолчанию и целью `multiclass`:

```
param = {'objective': 'multiclass', 'num_class': 5}
booster = lgb.train(params=param,
                    train_set=lgb_train,
                    num_boost_round=500,
                    early_stopping_rounds=20,
                    valid_sets=[lgb_train, lgb_test])
```

Базовые настройки не улучшают мультиномную логистическую регрессию, но дальнейшая параметрическая настройка остается неиспользуемым вариантом возможного улучшения точности:

```
y_pred_class = booster.predict(test_dtm_numeric.astype(float))
accuracy_score(test.stars, y_pred_class.argmax(1) + 1)
0.738665855696524
```

Резюме

В этой главе мы рассмотрели целый ряд технических решений и вариантов обработки неструктурированных данных с целью извлечения семантически значимых числовых признаков для их использования в автоматически обучающихся моделях.

Мы рассмотрели базовый конвейер лексемизации и аннотатирования и проиллюстрировали его реализацию для нескольких естественных языков с использованием библиотек `sraCu` и `TextBlob`. На основе этих результатов мы создали модель документа, опираясь на модель мешка слов, представив документы в виде числовых векторов. Мы научились совершенствовать конвейер предобработки, а затем использовать векторизованные текстовые данные для классифицирования и сентиментного анализа.

В оставшихся двух главах, посвященных альтернативным текстовым данным, мы научимся резюмировать текст, воспользовавшись в следующей главе неконтролируемым автоматическим обучением с целью выявления скрытых (латентных) тем, и изучим технические решения для представления слов в виде векторов, которые отражают контекст словопотребления и были очень успешно применены с целью получения более богатых текстовых признаков для различных классификационных задач.

14

Тематическое моделирование

В предыдущей главе мы конвертировали неструктурированные текстовые данные в числовой формат, используя модель мешка слов. Эта модель абстрагируется от порядка слов и представляет документы в виде словарных векторов, где каждый элемент представляет релевантность лексемы документу.

Результирующая терм-документная матрица (document-term matrix, DTM) (также нередко встречается транспонированная терм-документная матрица) полезна для сравнения документов друг с другом или с вектором запроса на основе их лексического содержимого и соответственно быстрого отыскания иголки в стоге сена или классифицирования документов.

Однако эта модель документа является одновременно высокоразмерной и очень разреженной. В результате она мало что делает для резюмирования содержимого или приближения к пониманию того, о чем это содержимое. В этой главе мы будем использовать неконтролируемое машинное обучение в форме тематического моделирования для извлечения из документов скрытых тем. Эти темы могут проникать в сущность информации большого объема документов в автоматическом режиме. Они очень полезны для понимания самого стога сена и позволяют выполнять краткую разметку документов, потому что они используют степень ассоциации тем и документов.

Тематические модели позволяют извлекать изолированные, интерпретируемые текстовые признаки, которые могут использоваться разнообразными способами для извлечения торговых сигналов из больших коллекций документов. Они ускоряют выполнение обзора документов, помогают выявлять и кластеризовать схожие документы и могут быть аннотированы в качестве основы для предсказательного моделирования. Их применение охватывает выявление ключевых тем в раскрываемой корпоративной информации или стенограммах телеконференций о корпоративных заработках, клиентских отзывах или контрактах, аннотированных с использованием, например, сентиментного анализа либо прямой разметки последующими возвратами от активов.

Более конкретно в этой главе мы рассмотрим следующие темы:

- ◆ что достигается с помощью тематического моделирования, в чем его важность, и как оно развивалось;
- ◆ как *латентно-семантическое индексирование* (Latent Semantic Indexing, LSI) снижает размерность терм-документной матрицы (DTM);
- ◆ как *вероятностный латентно-семантический анализ* (probabilistic LSA, pLSA) использует генеративную модель для извлечения тем;
- ◆ как *латентное размещение Дирихле* (Latent Dirichlet Allocation, LDA) уточняет вероятностный латентный семантический анализ и почему эта тематическая модель является самой популярной;
- ◆ как визуализировать и оценивать результаты тематического моделирования;
- ◆ как реализовывать латентное размещение Дирихле (LDA) с помощью библиотек `sklearn` и `gensim`;
- ◆ как применять тематическое моделирование к коллекциям стенограмм телеконференций о корпоративных заработках и отзывов о деятельности предприятий Yelp.



Примеры исходного кода для следующих далее разделов находятся в каталоге репозитория GitHub для этой главы, а справочные материалы перечислены в главном файле README.

Усвоение скрытых тем: цели и подходы

Тематическое моделирование призвано обнаруживать скрытые темы, или тематики, в документах, которые улавливают семантическую информацию, лежащую за пределами отдельных слов. Оно направлено на решение ключевой сложности в построении автоматически обучающегося алгоритма, который способен учиться на текстовых данных, выходя за пределы лексического уровня того, что было написано, и поднимаясь на семантический уровень того, что подразумевалось. Полученные темы можно использовать для аннотирования документов на основе их взаимосвязи с разнообразными темами.

Другими словами, тематическое моделирование призвано автоматически обобщать крупные коллекции документов с целью обеспечения организации и управления, а также поиска и выдачи рекомендаций. В то же время оно может обеспечивать понимание документов в той мере, в какой люди могут интерпретировать описания тем.

Тематические модели стремятся устранять проклятие размерности, которое может поражать модель мешка слов. Представление документа, основанное на высоко-размерных разреженных векторах, нередко делает меры сходства шумными, что приводит к неточному измерению расстояния и переподгонке моделей текстовой классификации.

Более того, модель мешка слов игнорирует порядок слов и теряет контекст, а также семантическую информацию, потому что она не способна улавливать синонимию (несколько слов имеют одно и то же значение) и полисемию (одно слово имеет несколько значений). Вследствие этого документный поиск или поиск по сходству может пропустить момент, когда документы не будут проиндексированы терминами, используемыми для поиска или сравнения.

Эти недостатки порождают следующий вопрос: как моделировать и усваивать смысловые темы, обеспечивающие более продуктивное взаимодействие с текстовыми данными?

От линейной алгебры к иерархическим вероятностным моделям

В первоначальных попытках тематических моделей улучшить модель векторного пространства (разработанную в середине 1970-х годов) применялась линейная алгебра. Это делалось с целью снижения размерности документной матрицы. Указанный подход аналогичен алгоритму, который мы обсуждали в *главе 12*, посвященной неконтролируемому автоматическому обучению. Несмотря на эффективность этих моделей, трудно оценить их результаты без эталонной модели.

В ответ появились вероятностные модели, которые принимают вид процесса явного генерирования документов и предоставляют алгоритмы для реконструкции этого процесса и восстановления базовых тем.

В табл. 14.1 выделены ключевые этапы эволюции модели, которую мы рассмотрим более подробно в последующих разделах.

Таблица 14.1. Ключевые этапы эволюции модели

| Модель | Год | Описание |
|--|------|--|
| Латентно-семантическое индексирование (LSI) | 1988 | Снижает размерность словарного пространства, улавливая семантические взаимосвязи между документом и термином |
| Вероятностный латентно-семантический анализ (pLSA) | 1999 | Реконструирует процесс, который исходит из того, что слова генерируют тему, а документы — это смеси тем |
| Латентное размещение Дирихле (LDA) | 2003 | Добавляет генеративный процесс для документов: трехуровневую иерархическую байесову модель |

Латентно-семантическое индексирование

Латентно-семантическое индексирование (Latent Semantic Indexing, LSI; так называемый латентно-семантический анализ — Latent Semantic Analysis, LSA) берется улучшить результаты запросов, в которых опущены релевантные документы, со-

держащие синонимы терминов запроса. Оно призвано моделировать взаимосвязи между документами и терминами, чтобы иметь возможность предсказывать, что некий термин следует ассоциировать с документом, хотя из-за вариативности в словопотреблении такая ассоциация не наблюдалась.

Метод LSI использует линейную алгебру, отыскивая заданное число k латентных тематик путем разложения терм-документной матрицы (DTM). Более конкретно, он использует *сингулярное разложение* (SVD), отыскивая наилучшую нижнеранговую аппроксимацию терм-документной матрицы с использованием k сингулярных значений и векторов. Латентно-семантическое индексирование является (LSI) применением технических решений неконтролируемого обучения по снижению размерности, с которыми мы столкнулись в *главе 12*, к текстовому представлению, которое мы рассмотрели в *главе 13*. Авторы экспериментировали с иерархической кластеризацией, но сочли, что моделировать взаимосвязи "документ — тема" и "тема — термин" явным образом или улавливать ассоциации документов или терминов с несколькими темами было бы слишком ограничительным подходом.

В этом контексте сингулярное разложение служит для выявления набора некоррелированных индексирующих переменных или коэффициентов, которые позволяют представлять каждый термин и документ своим вектором коэффициентных значений.

На рис. 14.1 показано, как сингулярное разложение раскладывает терм-документную матрицу на три матрицы, две из которых содержат ортогональные сингулярные векторы, а третья — это диагональная матрица с сингулярными значениями, служащими шкалирующими коэффициентами. Благодаря некоторой корреляции в исходных данных, сингулярные значения снижаются, вследствие чего выборка только самых больших T сингулярных значений производит нижнеразмерную аппроксимацию исходной терм-документной матрицы, которая теряет относительно мало информации. Следовательно, в сокращенной версии строки или столбцы, которые имели N элементов, имеют только $T < N$ элементов.

M
документов

$\begin{pmatrix} \dots & \dots & \dots \\ \vdots & \Sigma & \vdots \\ \dots & \dots & \dots \end{pmatrix}$

Терм-документная матрица

N терминов

$\begin{pmatrix} \dots & \dots & \dots \\ \vdots & U & \vdots \\ \dots & \dots & \dots \end{pmatrix}$

Сингулярные векторы ($M \times N$)

Схожесть по схеме "документ — тема"

$\begin{bmatrix} \dots & \dots & \dots \\ \Sigma & & \\ \dots & \dots & \dots \end{bmatrix}$

Сингулярные значения ($N \times N$)

Сила темы

$\begin{pmatrix} \dots & \dots & \dots \\ \vdots & V^T & \vdots \\ \dots & \dots & \dots \end{pmatrix}$

Сингулярные векторы ($N \times M$)

Схожесть по схеме "термин — тема"

1. Снизить размерность, используя $T < N$ сингулярных значений.

2. Оценить терм-документную матрицу, используя $U_T \Sigma_T$.

Рис. 14.1. Процедура разложения терм-документной матрицы методом сингулярного разложения

Это сокращенное разложение может быть интерпретировано, как показано на рис. 14.1, где первая матрица $M \times T$ представляет взаимосвязи между документами и темами, диагональная матрица шкалирует темы по их силе в корпусе, а третья матрица моделирует терм-документную взаимосвязь.

Строки матрицы, которая получается в результате произведения первых двух матриц $U_T \Sigma_T$, соответствуют местоположениям исходных документов, спроецированных в латентное тематическое пространство.

Как реализовать LSI с помощью библиотеки sklearn

Мы проиллюстрируем применение латентно-семантического индексирования (LSI), используя данные новостных статей BBC, которые мы ввели в предыдущей главе, потому что они достаточно малы для того, чтобы позволить провести быструю тренировку и сравнить отнесения тем к меткам категорий. Дополнительные сведения о реализации см. в блокноте `latent_semantic_indexing.ipynb`:

1. Мы начинаем с загрузки документов и создания тренировочного и (стратифицированного) тестового наборов с 50 статьями.
2. Затем мы векторизуем данные с помощью класса `TfidfVectorizer`, получив взвешенные частоты терм-документной таблицы (DTM), и фильтруем слова, которые появляются менее чем в 1% или более чем в 25% документов, а также частые стоп-слова, получив лексикон примерно из 2900 слов:

```
vectorizer = TfidfVectorizer(max_df=.25, min_df=.01,
                             stop_words='english',
                             binary=False)

train_dtm = vectorizer.fit_transform(train_docs.article)
test_dtm = vectorizer.transform(test_docs.article)
```

3. Мы используем класс `TruncatedSVD` библиотеки `sklearn`, который вычисляет только самые большие сингулярные значения, снижая размерность терм-документной матрицы. Детерминированный алгоритм `arpack` обеспечивает точное решение, однако принятая по умолчанию рандомизированная реализация эффективнее работает с большими матрицами.
4. Мы вычисляем пять тем для сопоставления с пятью категориями, которые объясняют только 5,4% общей дисперсии терм-документной матрицы, поэтому более высокие значения будут разумными:

```
svd = TruncatedSVD(n_components=5, n_iter=5, random_state=42)
svd.fit(train_dtm)
svd.explained_varianceratio

array([0.00187014, 0.01559661, 0.01389952, 0.01215842, 0.01066485])
```

5. Латентно-семантическое индексирование (LSI) выявляет новый ортогональный базис для терм-документной матрицы, который снижает ранг до числа желаемых тем.

6. Метод `transform()` натренированного объекта `svd` проецирует документы в новое тематическое пространство, т. е. результат снижения размерности документных векторов, что соответствует приведенному ранее преобразованию $U_T Z_T$:

```
train_doc_topics = svd.transform(train_dtm)
train_doc_topics.shape
```

```
(2175, 5)
```

7. Мы можем выбрать статью и посмотреть ее расположение в тематическом пространстве. Мы извлекаем политическую статью (с категорией `Politics`), которая наиболее (положительно) ассоциирована с темами 1 и 2:

```
i = randint(0, len(train_docs))
train_docs.iloc[i, :2].append(pd.Series(doc_topics[i],
                                         index=topic_labels))
```

```
Category                                Politics
Heading    UK heading wrong way, says Howard
Topic 1                                0.28
Topic 2                                0.25
Topic 3                                0.19
Topic 4                                0.04
Topic 5                                0.04
dtype: object
```

8. Отнесения тем для этого примера совпадают со средними тематическими весами для каждой проиллюстрированной далее категории (категория `Politics` — самая левая). Они иллюстрируют, как LSI выражает k тем как направления в k -мерном пространстве (блокнот содержит проекцию в двумерное пространство средних отнесений тем в расчете на категорию).
9. Каждая категория четко определена, и тестовые отнесения сочетаются с тренировочными отнесениями. Однако веса являются как положительными, так и отрицательными, что затрудняет интерпретацию тем (рис. 14.2).
10. Мы также можем показать слова, которые наиболее тесно ассоциированы с каждой темой (в абсолютном выражении). Темы, по всей видимости, улавливают некоторую семантическую информацию, но не дифференцированы (рис. 14.3).

Сильные и слабые стороны

Преимущества латентно-семантического индексирования (LSI) включают удаление шума и смягчение проклятия размерности, а также улавливание некоторой семантики и кластеризацию как документов, так и терминов.

Однако результаты LSI трудно интерпретировать, поскольку темы являются векторами слов с положительными и отрицательными элементами. Кроме того, не

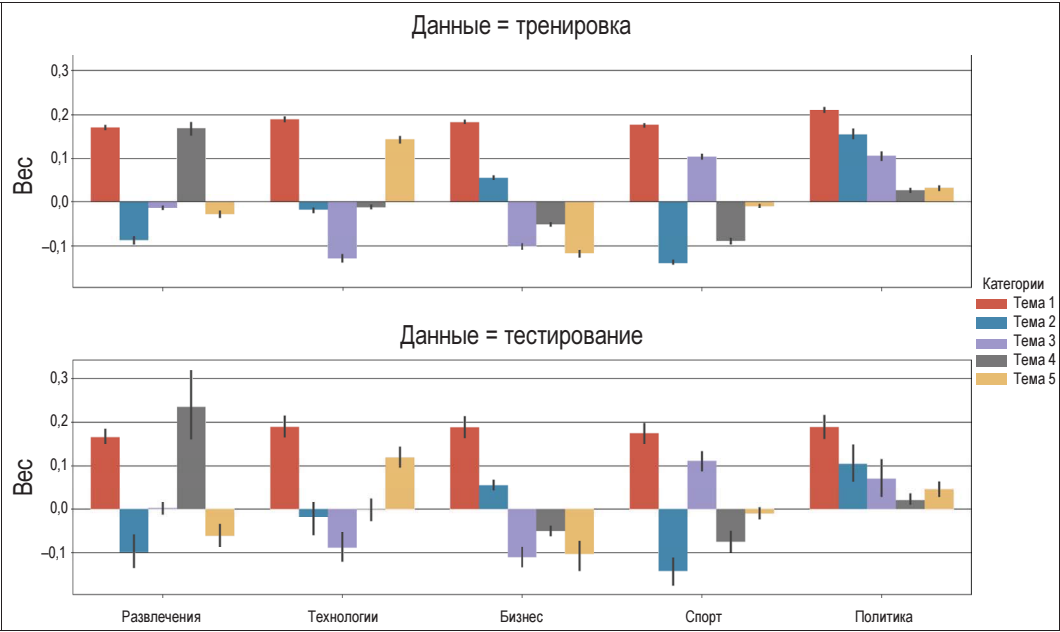


Рис. 14.2. Отношения тем к категориям на тестовом и тренировочном наборах данных методом латентно-семантического индексирования (LSI)



Рис. 14.3. Слова, наиболее тесно ассоциированные с каждой темой (в абсолютном выражении) методом латентно-семантического индексирования (LSI)

существует базовой модели, которая позволяла бы проводить оценивание соответствия и давать руководящие указания во время отбора числа размерностей или тем.

Вероятностный латентный семантический анализ

Вероятностный латентно-семантический анализ (Probabilistic Latent Semantic Analysis, pLSA) принимает статистический взгляд на латентный семантический анализ и создает генеративную модель для решения проблемы нехватки теоретических обоснований латентного семантического анализа.

pLSA явным образом моделирует вероятность каждого совместного появления документов d и слов w , описываемых терм-документной матрицей (DTM) как смесь условно взаимно независимых мультиномиальных распределений, которые содержат темы t .

Симметричная формулировка этого генеративного процесса совместных появлений слов и документов исходит из того, что как слова, так и документы генерируются латентным тематическим классом, тогда как асимметричная модель исходит из того, что темы отбираются с учетом данного документа, а слова являются результатом второго шага с учетом данной темы:

$$P(w, d) = \underbrace{\sum_t P(d|t)P(w|t)}_{\text{симметричная}} = \underbrace{P(d) \sum_t P(t|d)P(w|t)}_{\text{асимметричная}}.$$

Число тем является гиперпараметром, выбираемым перед тренировкой, и не усваивается из данных.

В вероятностных моделях часто используется приведенное ниже фигурное обозначение, которое служит для выражения зависимостей. На рис. 14.4 кодируются взаимосвязи, только что описанные для асимметричной модели. Каждый прямоугольник представляет несколько элементов, таких как M документов во внешнем блоке и N слов каждого документа во внутреннем блоке. Мы наблюдаем только за документами и их содержимым, и модель выводит скрытое, или латентное, распределение тем.

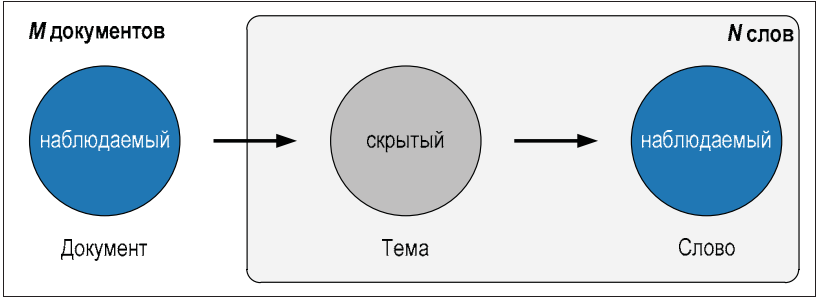


Рис. 14.4. Взаимосвязи в асимметричной модели

Преимущество использования вероятностной модели заключается в том, что теперь мы можем сравнивать модели, оценивая вероятность, которую они назначают новым документам с учетом имеющихся параметров, усвоенных в ходе тренировки.

Как реализовать pLSA с помощью библиотеки sklearn

Вероятностный латентный семантический анализ (pLSA) эквивалентен неотрицательному разложению матрицы с использованием целевой меры расхождения Кульбака — Лейблера (см. справочный материал в репозитории GitHub по адресу <https://github.com/PacktPublishing/Hands-On-Machine-Learning-for-Algorithmic-Trading>). Следовательно, для реализации этой модели мы можем применить класс `sklearn.decomposition.NMF`, следуя примеру латентного семантического анализа (LSA).

Используя ту же самую тренировочно-тестовую разбивку терм-документной матрицы, произведенную векторизатором `TfidfVectorizer`, мы выполняем подгонку pLSA следующим образом:

```
nmf = NMF(n_components=n_components,
          random_state=42,
          solver='mu',
          beta_loss='kullback-leibler',
          max_iter=1000)
nmf.fit(train_dtm)
```

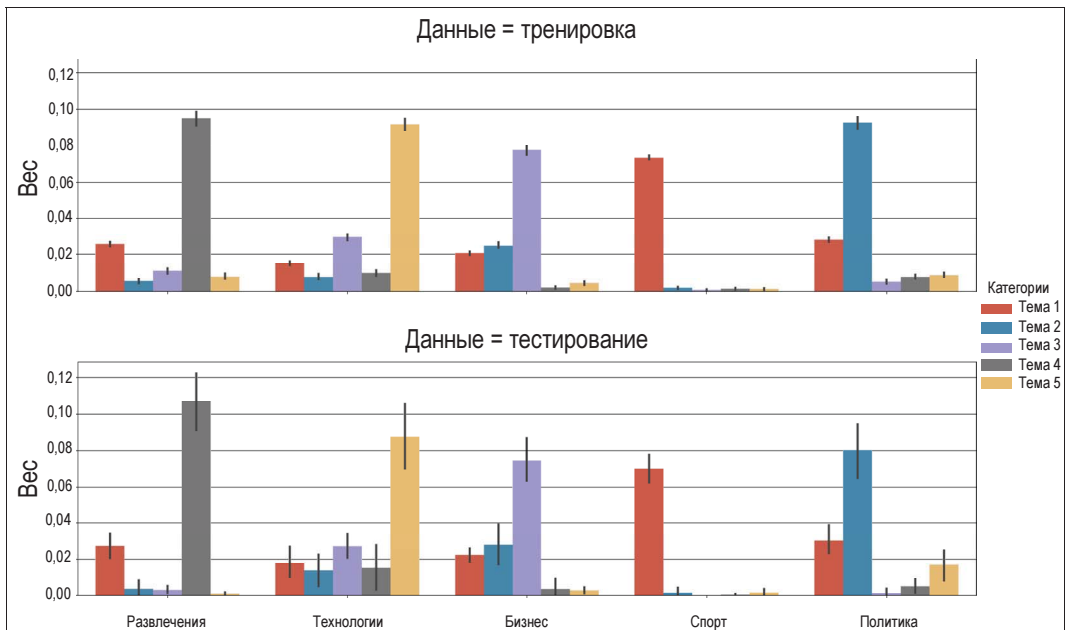


Рис. 14.5. Отнесения тем к категориям на тестовом и тренировочном наборах данных методом вероятностного латентного семантического анализа (pLSA)

Мы получаем меру ошибки восстановления, которая заменяет ранее рассматривавшуюся меру объясненной дисперсии:

```
nmf.reconstruction_err_  
316.2609400385988
```

По своей вероятностной природе pLSA производит только положительные тематические веса, которые приводят к более простым взаимосвязям "тема — категория" для тестовых и тренировочных наборов (рис. 14.5).

Мы также видим, что словарные списки, описывающие каждую тему, начинают приобретать больше смысла; например, категория развлечений наиболее непосредственно связана с темой 4, которая включает слова film (фильм), star (звезда) и т. д.



Рис. 14.6. Слова, наиболее тесно ассоциированные с каждой темой (в абсолютном выражении) методом вероятностного латентного семантического анализа (pLSA)

Латентное размещение Дирихле

Латентное размещение Дирихле (Latent Dirichlet allocation, LDA) расширяет вероятностный латентный семантический анализ, добавляя генеративный процесс для тем.

Это самая популярная тематическая модель, поскольку она тяготеет к порождению значимых тем, которые людям знакомы, может назначить темы новым документам и расширяться. Варианты моделей LDA могут включать метаданные, такие как авторы или графические данные, или усваивать иерархические темы.

Как работает LDA

Латентное размещение Дирихле (LDA) представляет собой иерархическую байесовую модель, которая исходит из того, что темы являются распределениями вероятностей над словами, а документы — распределениями над темами. Более конкретно,

указанная модель исходит из того, что темы подчиняются разреженному распределению Дирихле, из чего следует, что документы охватывают только малое множество тем, а темы используют часто только малое множество слов.

Распределение Дирихле

Распределение Дирихле порождает векторы вероятностей, которые могут быть использованы с дискретными распределениями, т. е. оно случайно генерирует заданное число величин, которые являются положительными и в сумме составляют единицу, что и ожидается от вероятностей. Оно имеет параметр с положительным, действительным значением, который контролирует концентрацию вероятности. Его значения ближе к нулю означают, что только несколько величин будут положительными и получат наибольшую массу вероятности. На рис. 14.7 показаны три выборки размера 10 для $\alpha = 0,1$ (блокнот `dirichlet_distribution.ipynb` содержит симуляцию, благодаря которой вы можете поэкспериментировать с различными значениями параметров).

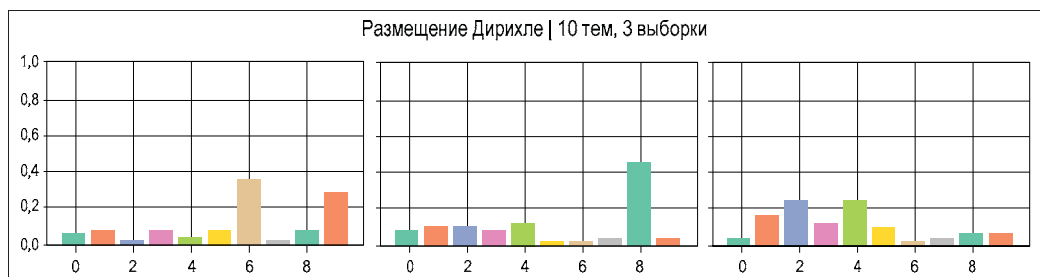


Рис. 14.7. Размещение Дирихле

Генеративная модель

Распределение Дирихле занимает видное место в тематической модели LDA, которая принимает вид следующего генеративного процесса, когда автор добавляет статью в массив документов:

1. Случайно смешать небольшое подмножество совместных тем K в соответствии с тематическими вероятностями.
2. Для каждого слова выбрать одну из тем в соответствии с документно-тематическими вероятностями (документ — тема).
3. Выбрать слово из словарного списка темы в соответствии с тематико-словарными вероятностями (тема — слово).

В результате содержание статьи зависит от веса каждой темы и от терминов, составляющих каждую тему. Распределение Дирихле регулирует выборку тем для документов и слов для тем и кодирует идею о том, что документ охватывает только

несколько тем, в то время как каждая тема часто использует только малое число слов.

Фигурное обозначение для модели LDA на рис. 14.8 резюмирует эти взаимосвязи.

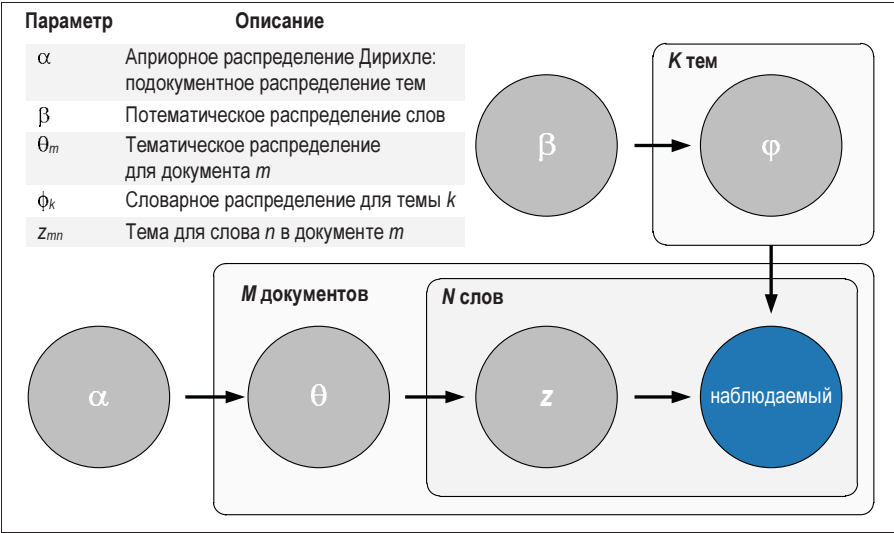


Рис. 14.8. Взаимосвязи модели LDA

Реконструирование процесса

Генеративный процесс является воображаемым, но оказывается полезным, потому что позволяет восстанавливать различные распределения. Алгоритм латентного размещения Дирихле (LDA) реконструирует работу воображаемого автора и приходит к сводке взаимосвязей по схеме "документ — тема — слово", которая кратко описывает следующее:

- ◆ процентный вклад каждой темы в документ;
- ◆ вероятностная ассоциация каждого слова с темой.

Латентное размещение Дирихле решает задачу байесова вывода о восстановлении распределений из корпуса документов и слов, которые они содержат, путем реконструирования (обратного конструирования) предполагаемого процесса генерирования содержимого. В оригинальной исследовательской работе для аппроксимации апостериорного распределения используется вариационный Байес (variational Bayes, VB). Альтернативы включают выборку по Гиббсу и распространение ожиданий. Позже мы проиллюстрируем его реализации с использованием библиотек sklearn и gensim.

Как оценивать темы LDA

Неконтролируемые тематические модели не гарантируют, что результат будет осмысленным или интерпретируемым, а объективная метрика, которой можно было бы оценивать результат, как в контролируемом обучении, отсутствует. Оценивание на основе человеческой темы считается золотым стандартом, но обходится потенциально дорого и не легко доступно в широком масштабе.

Два варианта более объективного оценивания результатов включают перплексивность, которая оценивает модель на не встречавшихся ранее документах, и метрики тематической когерентности, которые ориентированы на оценивание семантического качества обнаруженных регулярностей.

Перплексивность

Перплексивность применительно к LDA служит мерой того, насколько хорошо восстановленное моделью тематико-словарное вероятностное распределение предсказывает образец, например, не встречавшиеся текстовые документы. Она основана на энтропии $H(p)$ этого распределения p и вычисляется относительно множества лексем w :

$$2^{H(p)} = 2^{-\sum_w p(w) \log_2 p(w)}.$$

Чем значение данной меры ближе к нулю, тем распределение лучше предсказывает образец.

Тематическая когерентность

Тематическая когерентность служит мерой семантической непротиворечивости результатов тематической модели, т. е. будут ли люди воспринимать слова и их вероятности, ассоциированные с темами, как значимые.

С этой целью она выставляет отметку каждой теме, измеряя степень семантического сходства между словами, наиболее релевантными для темы. Более конкретно, меры когерентности основаны на вероятности наблюдать множество слов W , которые все вместе определяют тему.

Мы используем две меры когерентности, которые были разработаны для LDA и на практике показали, что они согласуются с человеческим суждением о качестве темы, а именно меры UMass и UCI.

Метрика UCI определяет отметку пары слов как сумму точечной взаимной информации (Pointwise Mutual Information, PMI) между двумя различными парами (лучших) тематических слов $w_i, w_j \in w$ и сглаживающим коэффициентом ε :

$$\text{когерентность} = \sum_{(w_i, w_j) \in W} \log \frac{p(w_i, w_j) + \varepsilon}{p(w_i)p(w_j)}.$$

Вероятности вычисляются из частот совместных появлений слов в скользящем окне над внешним корпусом документов, таким как Википедия, вследствие чего эту метрику можно рассматривать как внешнее сравнение с данными непосредственно-го семантического наблюдения (ground truth).

Напротив, метрика UMass использует совместные появления в определенном числе документов D из тренировочного корпуса для вычисления отметки когерентности:

$$\text{когерентность}_{\text{UMass}} = \sum_{(w_i, w_j) \in W} \log \frac{D(w_i, w_j) + \varepsilon}{D(w_j)}.$$

Вместо сравнения с внешними данными непосредственных наблюдений эта мера отражает внутренне присущую когерентность. Обе меры получили признание, как хорошо согласующиеся с человеческими суждениями. В обоих случаях значения, близкие к нулю, означают, что тема является более когерентной, т. е. непротиворечивой.

Как реализовать LDA с помощью библиотеки sklearn

Используя новостные данные BBC, как и раньше, мы применяем класс `sklearn.decomposition.LatentDirichletAllocation` для тренировки модели LDA с пятью темами (подробности относительно параметров см. в документации sklearn, детали реализации можно найти в блокноте `lda_with_sklearn.ipynb`):

```
lda = LatentDirichletAllocation(n_components=5,
                               n_jobs=-1,
                               max_iter=500,
                               learning_method='batch',
                               evaluate_every=5,
                               verbose=1,
                               random_state=42)
ldat.fit(train_dtm)

LatentDirichletAllocation(batch_size=128, doc_topic_prior=None,
                          evaluate_every=10, learning_decay=0.7,
                          learning_method='batch', learning_offset=10.0,
                          max_doc_update_iter=100, max_iter=500, mean_change_tol=0.001,
                          n_components=5, n_jobs=None, n_topics=None, perp_tol=0.1,
                          random_state=42, topic_word_prior=None,
                          total_samples=1000000.0, verbose=1)
```

Указанная модель отслеживает перплексивность в образце во время тренировки и прекращает итеративную обработку, как только эта мера перестает улучшаться. Результат, как обычно, можно сохранить и загрузить с объектами sklearn:

```
joblib.dump(lda, model_path / 'lda.pkl')
lda = joblib.load(model_path / 'lda.pkl')
```

Как визуализировать результаты LDA с помощью библиотеки pyLDAvis

Визуализирование темы способствует оцениванию качества темы с помощью человеческого суждения. Библиотека pyLDAvis — это перенесенная на Python библиотека LDAvis, разработанная в R и D3.js. Мы представим ее ключевые концепции; каждый блокнот с реализацией LDA содержит примеры.

Библиотека pyLDAvis показывает глобальные взаимосвязи между темами, при этом обеспечивая их семантическое оценивание путем обследования терминов, наиболее тесно ассоциированных с каждой темой и, наоборот, темы, ассоциированные с каждым термином. В ней также решается проблема, состоящая в том, что часто встречающиеся в корпусе термины, как правило, доминируют в мультиномиальном распределении над словами, определяющими тему. LDAvis вводит релевантность r термина w теме t , порождая гибкое ранжирование ключевых терминов с использованием весового параметра $0 \leq \lambda \leq 1$.

Обозначив через ϕ_{kw} модельную вероятностную оценку наблюдать термин w для темы t и через w маргинальную вероятность в корпусе, мы имеем:

$$r(w, k | \lambda) = \lambda \log(\phi_{kw}) + (1 - \lambda) \log \frac{\phi_{kw}}{p_w}.$$

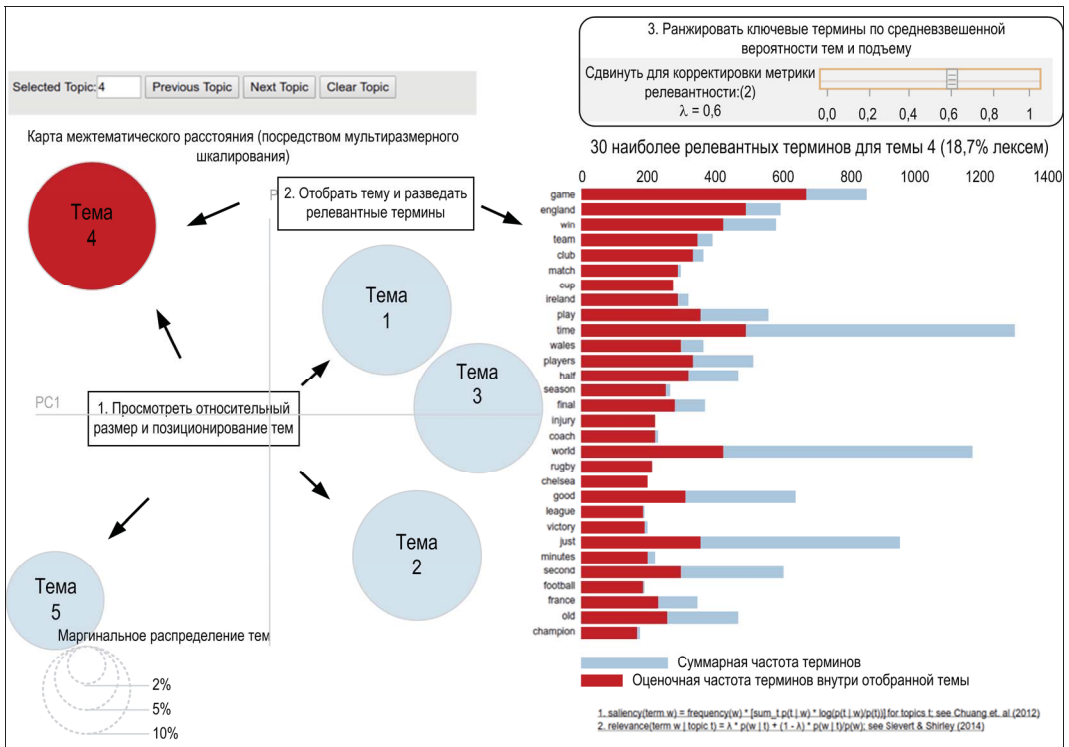


Рис. 14.9. Результаты работы инструмента LDAvis

Первый член служит мерой степени ассоциации термина t с темой w , а второй член — мерой подъема (lift), или выпуклости, т. е. насколько более вероятен термин для темы, чем в корпусе (рис. 14.9).

Данный инструмент позволяет пользователю интерактивно изменять λ , корректируя релевантность, которая обновляет ранжирование терминов. Пользовательские исследования показали, что $\lambda = 0,6$ порождает наиболее убедительные результаты.

Как реализовать LDA с помощью библиотеки `gensim`

Библиотека `gensim` представляет собой специализированную библиотеку обработки ЕЯ с быстрой реализацией LDA и многочисленным дополнительным функционалом. Мы также будем использовать ее в следующей главе, посвященной словарным векторам (подробнее см. в блокноте `latent_dirichlet_allocation_gensim.ipynb`).

Она обеспечивает преобразование терм-документной матрицы, произведенной библиотекой `sklearn`, в структуры данных библиотеки `gensim` следующим образом:

```
train_corpus = Sparse2Corpus(train_dtm, documents_columns=False)
test_corpus = Sparse2Corpus(test_dtm, documents_columns=False)
id2word = pd.Series(vectorizer.get_feature_names()).to_dict()
```

Алгоритм LDA библиотеки `gensim` содержит многочисленные настройки, которые имеют следующий вид:

```
LdaModel(corpus=None,
          num_topics=100,
          id2word=None,
          distributed=False,
          chunksize=2000, # число документов в тренировочном блоке
          passes=1,       # число проходов по корпусу во время тренировки
          update_every=1, # число документов для перебора при обновлении
          alpha='symmetric',
          eta=None,       # априорное мнение о вероятности слов
          decay=0.5,      # % лямбды, забываемой во время
                        # исследования нового документа
          offset=1.0,     # управляет замедлением нескольких
                        # первых итераций
          eval_every=10,  # как часто оценивать лог. перплексивность
                        # (дорого)
          iterations=50,  # макс. число итераций по корпусу
          gamma_threshold=0.001, # мин. изменение в гамме
                        # для продолжения
          minimum_probability=0.01, # фильтровать темы
                        # с низкой вероятностью

          random_state=None,
          ns_conf=None,
          minimum_phi_value=0.01, # нижняя граница вероятностей терминов
```

```
per_word_topics=False,      # вычислять большинство
                             # вероятностей слово - тема
callbacks=None,
dtype=<class 'numpy.float32'>)
```

Библиотека `gensim` также предоставляет модель `LdaMulticore` для параллельной тренировки, которая способна ускорить тренировку с использованием многопроцессорного функционала Python для параллельных вычислений.

Тренировка модели требует лишь создания экземпляра объекта `LdaModel` следующим образом:

```
lda = LdaModel(corpus=train_corpus,
               num_topics=5,
               id2word=id2word)
```

Тематическая когерентность служит мерой совместного появления слов в теме. Она добавляет отметку для каждой отдельной верхнеранговой пары слов. Отметка представляет собой логарифм вероятности, что документ, содержащий по крайней мере один экземпляр верхнерангового слова также содержит по крайней мере один экземпляр нижнерангового слова.

Большие отрицательные значения указывают на то, что слова встречаются не часто; значения, близкие к нулю, указывают на то, что слова встречаются чаще. Библиотека `gensim` позволяет производить оценивание тематической когерентности и показывает наиболее важные слова в каждой теме:

```
coherence = lda_gensim.top_topics(corpus=train_corpus, coherence='u_mass')
```

Результаты можно показать следующим образом:

```
topic_coherence = []
topic_words = pd.DataFrame()

for t in range(len(coherence)):
    label = topic_labels[t]
    topic_coherence.append(coherence[t][1])
    df = pd.DataFrame(coherence[t][0], columns=[(label, 'prob'),
                                              (label, 'term')])
    df[(label, 'prob')] = df[(label, 'prob')].apply(lambda x: '{:.2%}'.format(x))
    topic_words = pd.concat([topic_words, df], axis=1)

topic_words.columns = pd.MultiIndex.from_tuples(topic_words.columns)
pd.set_option('expand_frame_repr', False)
topic_words.head().to_csv('topic_words.csv', index=False)
print(topic_words.head())

pd.Series(topic_coherence, index=topic_labels).plot.bar();
```

На выходе будут показаны лидирующие слова по каждой теме (табл. 14.2).

Таблица 14.2. Лидирующие слова по каждой теме

| Тема 1 | | Тема 2 | | Тема 3 | | Тема 4 | | Тема 5 | |
|--------------|--------|--------------|--------|--------------|--------|--------------|---------|--------------|----------|
| вероят-ность | термин | вероят-ность | термин | вероят-ность | термин | вероят-ность | термин | вероят-ность | термин |
| 0,55% | online | 0,90% | best | 1,04% | mobile | 0,64% | market | 0,94% | labour |
| 0,51% | site | 0,87% | game | 0,98% | phone | 0,53% | growth | 0,72% | blair |
| 0,46% | game | 0,62% | play | 0,51% | music | 0,52% | sales | 0,72% | brown |
| 0,45% | net | 0,61% | won | 0,48% | film | 0,49% | economy | 0,65% | election |
| 0,44% | used | 0,56% | win | 0,48% | use | 0,45% | prices | 0,57% | united |

Соответствующие отметки когерентности подчеркивают ухудшение качества темы (по крайней мере, частично из-за относительно небольшого набора данных) (рис. 14.10).

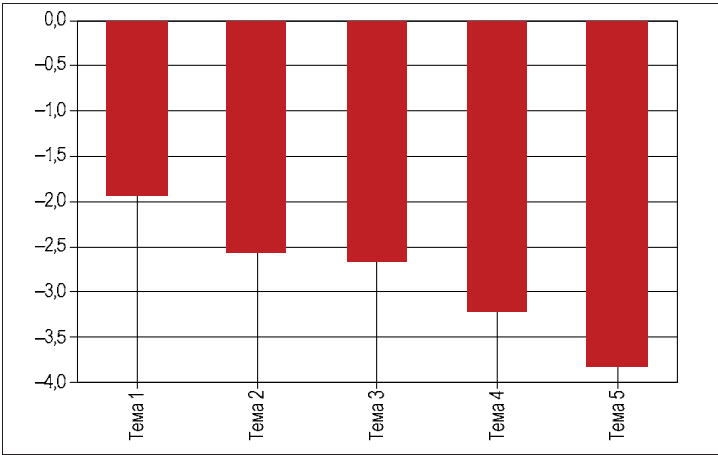


Рис. 14.10. Ухудшение качества темы

Тематическое моделирование применительно к телеконференциям о корпоративных зарплатах

В главе 3 мы познакомились с тем, как выскробывать данные телеконференций о корпоративном заработке с веб-сайта SeekingAlpha. В этом разделе мы проиллюстрируем тематическое моделирование с использованием этого источника. В примере будет использована выборка примерно со 500 стенограммами телеконференций о корпоративном заработке со второй половины 2018 г. Для практического

применения было бы очень желательно иметь более крупный набор данных. Каталог `earnings_calls` содержит несколько файлов с приведенными далее примерами.

Подробные сведения о загрузке, разведывании и предобработке данных, а также тренировке и оценивании индивидуальных моделей см. в блокноте `lda_earnings_calls.ipynb`. Файл `run_experiments.py` для экспериментов описан ниже.

Предобработка данных

Стенограммы состоят из отдельных высказываний представителя компании, модератора и, как правило, серии вопросов и ответов при общении с аналитиками. Мы будем рассматривать каждое высказывание как отдельный документ, игнорируя высказывания модератора, и в результате получим 26 314 элементов соответственно со средними и медианными количествами вхождений слов 144 и 64:

```
documents = []
for transcript in earnings_path.iterdir():
    content = pd.read_csv(transcript / 'content.csv')
    documents.extend(content.loc[(content.speaker!='Operator')
                                & (content.content.str.len() > 5),
                                'content'].tolist())

len(documents)
```

26314

Для обработки этих документов и сохранения очищенного и аннотированного текста в новый текстовый файл мы используем библиотеку `sraCu`, как показано в *главе 13* (см. блокнот).

В ходе разведывания данных выявляются специфические для предметной области стоп-слова, такие как `year` (год) и `quarter` (квартал), которые мы удаляем на втором шаге, где мы также отфильтровываем высказывания с менее чем десятью словами, вследствие чего остается около 16 150 слов.

Тренировка и оценивание модели

Для иллюстрации мы создадим терм-документную матрицу, содержащую термины, появляющиеся между 0,5 и 50% документов примерно для 1560 признаков. Тренировка 15-тематической модели с использованием 25 проходов по корпусу занимает чуть более двух минут на четырехъядерном процессоре `i7`.

Лидирующие 10 слов по каждой теме выявляют несколько разных тем, которые варьируются от очевидной финансовой информации до клинических испытаний (тема 4) и вопросов снабженческой цепочки (12) (рис. 14.11).

После использования метрики релевантности из библиотеки `pyLDavis` с 0,6-м взвешиванием безусловной частоты относительно подъема формулировки тем становятся более интуитивными. Это хорошо видно по теме 14 о результативности продаж (рис. 14.12).

| | | | | | | | | | | | | | | | |
|---|-------------|---------------|------------|---------------|-------------|------------|-------------|-------------|----------|------------|-----------|----------|------------|--------|---------|
| 0 | statement | expense | cloud | basis | patient | channel | focus | brand | want | project | lot | yes | price | maybe | bit |
| 1 | financial | total | technology | adjust | study | brand | acquisition | category | right | capital | thing | kind | production | kind | little |
| 2 | release | period | service | guidance | program | launch | investment | comp | price | investment | right | little | demand | okay | china |
| 3 | risk | income | solution | ebitda | clinical | experience | improve | team | thing | asset | mean | bit | volume | guess | loan |
| 4 | gaap | loss | platform | tax | trial | marketing | deliver | inventory | need | debt | actually | half | low | guy | service |
| 5 | officer | approximately | datum | low | phase | online | strategy | traffic | contract | portfolio | yes | pretty | capacity | sort | bank |
| 6 | chief | non | large | billion | datum | platform | invest | performance | lot | value | people | guidance | fleet | want | credit |
| 7 | conference | month | team | earning | development | digital | value | weather | say | balance | way | say | order | just | mention |
| 8 | measure | gaap | industry | gross | fda | consumer | progress | great | sure | return | different | thing | vessel | follow | tier |
| 9 | information | decrease | provide | approximately | cancer | user | performance | assortment | great | flow | obviously | low | supply | wonder | card |
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |

Рис. 14.11. Лидирующие 10 слов по каждой теме при тематическом моделировании применительно к телеконференциям о корпоративном заработке

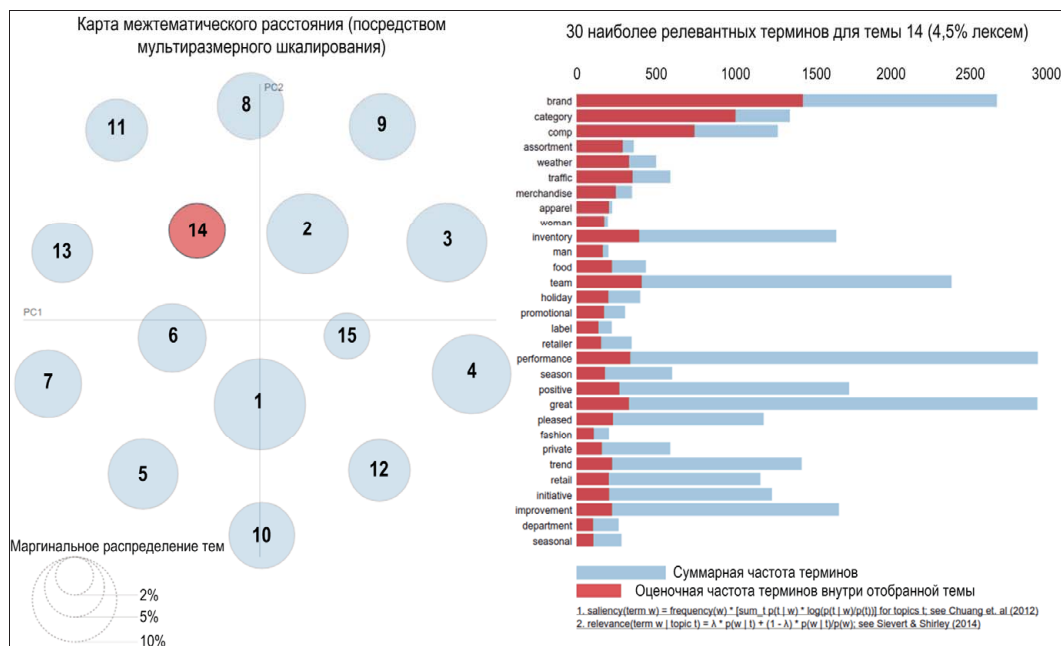


Рис. 14.12. Результативность продаж для темы 14

В указанном блокноте также иллюстрируется поиск документов по их тематической ассоциации. В этом случае аналитик может просмотреть соответствующие высказывания на наличие нюансов, применить сентиментный анализ для дальнейшей обработки текстовых данных по конкретной теме или назначить метки, выведенные из рыночных цен.

Проведение экспериментов

В качестве иллюстрации влияния разных параметрических настроек мы провели несколько сотен экспериментов для различных ограничений на терм-документную матрицу (DTM) и модельные параметры. Конкретнее, мы даем параметрам `min_df` и `max_df` варьироваться соответственно от 50 до 500 слов и от 10 до 100% документов, используя поочередно бинарные и абсолютные количества вхождений. Затем мы натренировали модели LDA с темами в количестве от 3 до 50, используя 1 и 25 проходов по корпусу.

Диаграммы рис. 14.13 иллюстрируют результаты с точки зрения когерентности темы (чем выше, тем лучше) и перплексивности (чем ниже, тем лучше). Когерентность падает после 25–30 тем, а перплексивность схожим образом увеличивается.

Блокнот включает результаты регрессии, которые квантифицируют связи между параметрами и результатами. Мы обычно получаем лучшие результаты, используя абсолютные количества вхождений и меньший лексикон.

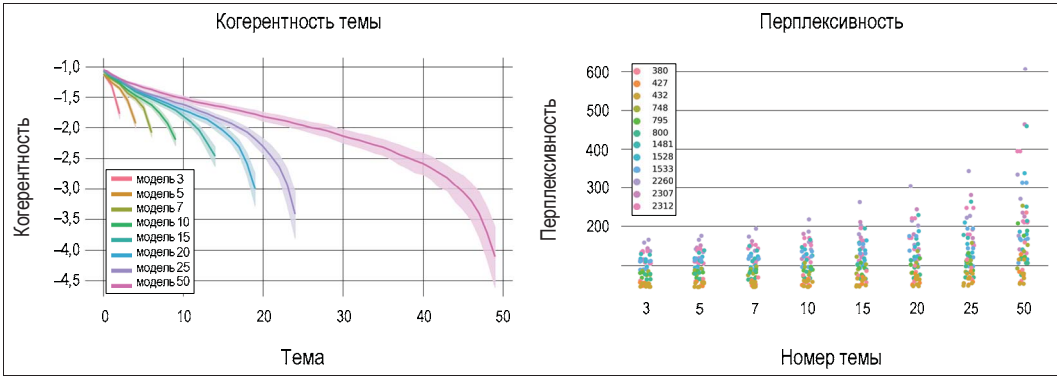


Рис. 14.13. Результаты с точки зрения когерентности и перплексивности

Тематическое моделирование применительно к отзывам о деятельности предприятий на веб-сайте Yelp

Блокнот `lda_yelp_reviews.ipynb` содержит пример LDA, применительно к 6 млн отзывов о деятельности предприятий на веб-сайте Yelp. Отзывы являются более однородными по длине, чем высказывания, извлеченные из стенограмм телеконференций о корпоративных заработках. После очистки, как и раньше, 10-й и 90-й процентиля варьируются от 14 до 90 лексем.

Мы показываем результаты для одной модели, используя лексикон в размере 3800 лексем на основе `min_df = 0,1%` и `max_df = 25%` с единственным проходом по корпусу для того, чтобы избежать длительного времени тренировки для 20 тем. Для вычисления значений релевантности можно использовать атрибут `pyldavis topic_info` с `lambda = 0,6`, в результате чего будет сгенерирован список слов (подробнее см. в блокноте), представленный на рис. 14.14.

Библиотека `gensim` предоставляет реализацию класса `LdaMultiCore`, которая предусматривает выполнение параллельной тренировки с использованием Python-модуля `multiprocessing` и повышает производительность на 50% при использовании четырех рабочих процессов. Однако большее число рабочих процессов не приводит к добавочному снижению времени тренировки из-за узких мест, связанных с вводом-выводом.

| | | | | | | | | | | | | | | | | | | | | |
|---|----------|----------|------------|------------|----------|---------|---------|----------|-----------|--------------|------------|----------|--------|--------|-------------|----------|---------|----------|----------|------|
| 0 | say | burger | price | bar | order | rice | car | room | coffee | thank | pizza | store | line | kid | office | location | har | review | tea | de |
| 1 | tell | cheese | sushi | beer | table | taco | call | hotel | breakfast | recommend | friendly | buy | wait | year | massage | dog | nail | yelp | drink | la |
| 2 | know | fry | pretty | wine | server | chicken | fix | stay | cream | highly | staff | shop | hour | fun | patient | parking | cut | read | coupon | et |
| 3 | bad | sandwich | quality | drink | waitress | soup | phone | vegas | chocolate | work | love | find | long | class | doctor | park | salon | write | shake | le |
| 4 | go | salad | buffet | menu | wait | spicy | company | pool | ice | professional | amazing | item | early | play | appointment | street | color | birthday | water | pour |
| 5 | ask | chicken | restaurant | atmosphere | minute | shrimp | tell | floor | rake | job | super | sell | minute | game | rare | downtown | wash | star | smoothie | in |
| 6 | customer | sauce | decent | night | ask | noodle | charge | casino | egg | experience | definitely | product | late | old | staff | area | job | mom | boba | con |
| 7 | want | bread | eat | patio | take | roll | credit | strip | brunch | wedding | awesome | purchase | open | son | question | south | stylist | sister | milk | pas |
| 8 | not | meat | average | dinner | seat | thai | day | club | donut | guy | favorite | sale | ticket | child | feel | drive | haircut | mother | ayce | que |
| 9 | rude | potato | high | cocktail | waiter | dish | repair | bathroom | cookie | amazing | nice | tire | movie | school | pain | north | paint | daughter | green | du |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |

Рис. 14.14. Лидирующие 10 слов по каждой теме при тематическом моделировании применительно к отзывам о деятельности предприятий Yelp

Резюме

В этой главе мы рассмотрели использование тематического моделирования для проникновения в сущность содержимого большой коллекции документов. Мы рассмотрели латентный семантический анализ, который использует процедуру снижения размерности терм-документной матрицы (DTM), проецируя документы в латентное тематическое пространство. Хотя этот метод является эффективным в устранении проклятия размерности, порождаемого высокоразмерными словарными векторами, он не улавливает значительный объем семантической информации. Вероятностные модели принимают явные допущения о взаимодействии документов, тем и слов, давая алгоритмам возможность реконструировать процесс генерирования документов и оценивать подгонку модели на новых документах. Мы увидели, что метод латентного размещения Дирихле (LDA) способен извлекать убедительные темы, позволяя получать высокоуровневое понимание крупных объемов текста в автоматическом режиме, а также выявлять релевантные документы целенаправленным образом.

В следующей главе мы научимся тренировать нейронные сети, которые встраивают индивидуальные слова в высокоразмерное векторное пространство, улавливая важную семантическую информацию и позволяя использовать результирующие словарные векторы в качестве высококачественных текстовых признаков.

15

Векторное вложение слов

В двух предыдущих главах мы применили модель мешка слов для конвертирования текстовых данных в числовой формат. В результате получались разреженные векторы фиксированной длины, которые представляют документы в высокоразмерном словарном пространстве. Это позволяет оценивать сходство документов и создавать признаки для тренировки автоматически обучающегося алгоритма и классифицирования содержимого документа или оценивания выраженных в нем настроений. Однако эти векторы игнорируют контекст, в котором используется термин, вследствие чего, например, другое предложение, содержащее те же слова, будет кодироваться точно таким же вектором.

В этой главе мы представим альтернативный класс алгоритмов, в которых для усвоения векторного представления индивидуальных семантических единиц, таких как слово или абзац, используются нейронные сети. Эти векторы являются скорее плотными, чем разреженными, и вместо десятков тысяч двоичных или дискретных элементов имеют несколько сотен вещественных значений. Они называются *векторными вложениями*¹, потому что назначают каждой семантической единице местоположение в непрерывном векторном пространстве.

Векторные вложения являются результатом тренировки модели соотносить лексемы с их контекстом и обладают таким преимуществом, что из сходного словопотребления следует сходный вектор. Более того, мы увидим, как вложения кодируют семантические аспекты, такие как взаимосвязи между словами посредством их относительного расположения. По этой причине они являются мощным техническим решением для использования в глубоко обучающихся моделях, которые мы представим в следующих главах.

¹ Векторное вложение слова (word embedding), или вложение слова в векторное пространство — это векторное представление слова как часть распределенного представления текста в n -мерном пространстве, идея которого состоит в том, чтобы создать плотный вектор для каждого слова, выбранный таким образом, чтобы он был похож на векторы слов, которые появляются в схожем контексте. См. <https://hackernoon.com/word-embeddings-in-nlp-and-its-applications-fab15eaf7430>. — Прим. перев.

Более конкретно в этой главе мы рассмотрим следующие темы:

- ◆ что такое векторные вложения слов, как они работают и улавливают семантическую информацию;
- ◆ как использовать натренированные словарные векторы;
- ◆ какие сетевые архитектуры широко используются для тренировки моделей векторных вложений слов Word2vec;
- ◆ как тренировать модель Word2vec с помощью библиотек Keras, gensim и TensorFlow;
- ◆ как визуализировать и оценивать качество словарных векторов;
- ◆ как тренировать модель Word2vec на основе финансовой отчетности, подаваемой в комиссию SEC;
- ◆ как модель векторных сложений документов Doc2vec расширяет модель векторных сложений слов Word2vec.

Как векторные вложения слов кодируют семантику

Модель мешка слов представляет документы как векторы, а те отражают лексемы, которые они содержат. Векторные вложения слов представляют лексемы как векторы более низких размерностей, вследствие чего их относительное расположение отражает их взаимосвязь с точки зрения того, как они используются в контексте. Они воплощают дистрибутивную гипотезу из лингвистики, которая утверждает, что слова лучше всего определяются компанией, которая их окружает.

Словарные векторы способны улавливать многочисленные семантические аспекты; схожие слова не только близки друг к другу, но и могут иметь несколько степеней сходства, например, слово *driver* (водитель) может быть похоже на слово *motorist* (автомобилист) или *to cause* (обуславливать). Кроме того, как мы проиллюстрируем позже в этом разделе, векторные вложения отражают взаимосвязи между парами слов, такими как аналогии (Токио для Японии то же самое, что Париж для Франции, или "лучше" для "хорошо" то же самое, что "хуже" для "плохо").

Векторные вложения, точнее вложения слов в векторное пространство, являются результатом тренировки автоматически обучающейся модели предсказывать слова из их контекста, или наоборот. В следующем разделе мы расскажем, как эти нейронные языковые модели работают, и представим успешные подходы, включая модели Word2vec, Doc2vec и fastText.

Как нейронно-языковые модели усваивают словопотребление в контексте

Векторные вложения слов получаются в результате тренировки мелкой нейронной сети предсказывать слово в условиях его контекста. В то время как традиционные

языковые модели определяют контекст как слова, предшествующие цели, модели векторных вложений слов используют слова, содержащиеся в симметричном окне, окружающем цель. В отличие от этого модель мешка слов в качестве контекста использует все документы целиком и для улавливания совместного появления слов вместо предсказательных векторов применяет (взвешенные) количества вхождений слов.

Ранее использовавшиеся нейронные языковые модели включали нелинейные скрытые слои, которые увеличивали вычислительную сложность. Двухслойная нейронная сеть Word2vec и ее расширения упростили архитектуру, обеспечив тренировку на крупных наборах данных (Википедия, например, содержит более 2 млрд лексем; дополнительную информацию о нейронных сетях прямого распространения см. в главе 17.)

Модель Word2vec — усвоение векторных вложений в широком масштабе

Модель векторных сложений слов Word2vec представляет собой двухслойную нейронную сеть, которая на входе принимает текстовый корпус, а на выходе выводит множество векторных вложений для слов в этом корпусе. Для эффективного усвоения словарных векторов с использованием мелких нейронных сетей, изображенных на рис. 15.1, существует две отличающиеся архитектуры.

- ♦ Модель *непрерывного мешка слов* (Continuous-Bag-Of-Words, CBOW) предсказывает целевое слово, используя в качестве входа среднее значение векторов контекстных слов, вследствие чего их порядок не имеет значения. Модель CBOW тренируется быстрее и демонстрирует тенденцию быть немного точнее для частых терминов, но уделяет меньше внимания нечастым словам.

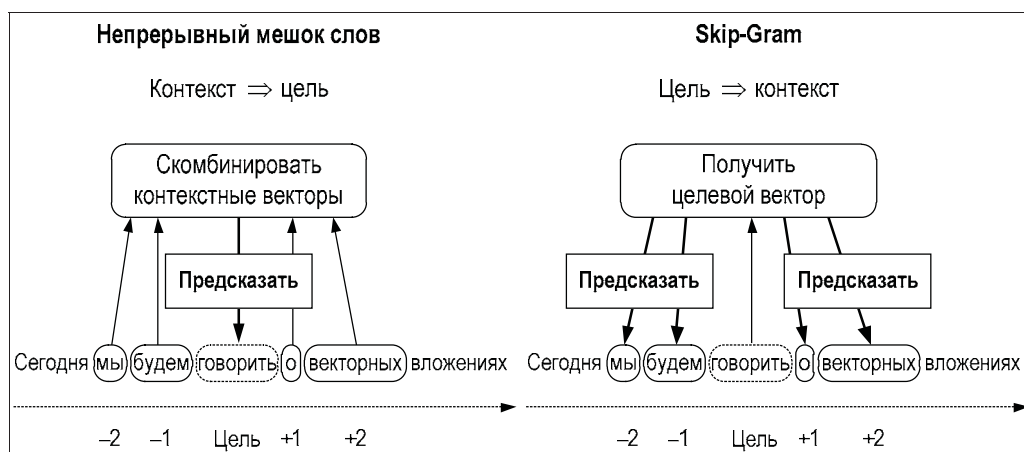


Рис. 15.1. Две архитектуры для эффективного усвоения словарных векторов с использованием мелких нейронных сетей

- ♦ Модель *Skip-Gram* (SG), напротив, использует целевое слово для предсказания слов, отобранных из контекста. Она хорошо работает с малыми наборами данных и находит хорошие представления даже для редких слов или фраз.

Следовательно, модель Word2vec получает вектор вложения в качестве входа и вычисляет скалярное произведение с другим вектором вложения. Обратите внимание, что, исходя из нормированности векторов, скалярное произведение максимизируется (в абсолютных единицах), когда векторы равны, и минимизируется, когда они ортогональны.

Затем она использует алгоритм обратного распространения, корректируя веса вложения в ответ на потерю, вычисляемую целевой функцией из-за любых классификационных ошибок. В следующем разделе мы увидим, как модель Word2vec вычисляет потерю.

Тренировочный процесс осуществляется путем перемещения контекстного окна над документами, обычно сегментированными в предложения. Каждая полная итерация по корпусу называется *эпохой*. В зависимости от данных для схождения качества вектора может потребоваться несколько десятков эпох.

Технически на практике было показано, что модель SG раскладывает словарно-контекстную матрицу, неявно содержащую точечную взаимную информацию соответствующих пар слов и контекста (справочные материалы см. в репозитории GitHub).

Модельная цель — упрощенная активационная функция softmax

Модели векторных сложений слов Word2vec призваны предсказывать одно-единственное слово из потенциально очень большого лексикона. В нейронных сетях часто используется активационная функция softmax, которая отображает любое число действительных значений в равное число вероятностей для воплощения соответствующей мультиклассовой цели, где h соответствует вложению, v — входным векторам, c — это контекст слова w :

$$p(w|c) = \frac{\exp(h^T v'_w)}{\sum_{w_i \in V} \exp(h^T v'_{w_i})}.$$

Однако сложность активационной функции softmax масштабируется вместе с числом классов, т. к. для стандартизации вероятностей знаменатель требует вычисления скалярного произведения для всех слов в лексиконе. Модели Word2vec повышают эффективность за счет использования упрощенной версии активации softmax либо подходов, основанных на выборке (подробнее см. справочные материалы).

- ♦ *Иерархическая активация softmax* организует лексикон в виде бинарного дерева со словами в виде листовых узлов. Уникальный путь к каждому узлу можно использовать для вычисления вероятности слова.

- ◆ **Контрастивное к шуму оценивание** (Noise-contrastive estimation, NCE) отбирает внеконтекстные "шумные слова" и аппроксимирует мультиклассовую задачу с помощью бинарной классификационной задачи². Производная NCE приближается к градиенту softmax по мере увеличения числа образцов, но только 25 образцов могут дать сходимость, аналогичную активации softmax, с темпом, который в 45 раз быстрее.
- ◆ **Отрицательная выборка** (Negative sampling, NEG) отбрасывает шумные образцы слов для аппроксимации NCE и непосредственно максимизирует вероятность целевого слова. Следовательно, вместо точности на тестовом наборе NEG оптимизирует семантическое качество векторов вложения (аналогичные векторы для аналогичного использования). Однако это может привести к более плохим представлениям для нечастых слов, чем целевая иерархическая активационная функция softmax.

Автоматическое обнаружение фраз

Предобработка, как правило, предусматривает обнаружение фраз, т. е. выявление лексем, которые широко используются вместе и должны получать одно векторное представление (например, "New York City", см. обсуждение n -грамм в главе 13).

Авторы исходной модели Word2vec используют простой метод оценивания лифта, который выявляет два слова w_i , w_j как биграмму, если их совместное появление превышает заданный порог относительно индивидуального внешнего вида каждого слова, скорректированного коэффициентом дисконтирования δ :

$$\text{отметка}(w_i, w_j) = \frac{\text{количество}(w_i, w_j) - \delta}{\text{количество}(w_i) \text{количество}(w_j)}.$$

Данное правило выставления отметки может применяться многократно для выявления все более длинных фраз.

Альтернативой является отметка нормализованной точечной взаимной информации (normalized point-wise mutual information, NPMI), которая, будучи точнее, обходится для вычисления дороже. В ней используется относительная частота слова $P(w)$, и она варьируется в промежутке от +1 до -1:

$$\text{NPMI} = \frac{\ln(P(w_i, w_j) / P(w_i)P(w_j))}{-\ln(P(w_i, w_j))}.$$

² Контрастивное к шуму оценивание (Noise-contrastive estimation, NCE) сводит задачу оценивания языковой модели к задаче оценивания параметров вероятностного бинарного классификатора, использующего те же параметры для различения выборок из эмпирического распределения от выборок, генерируемых распределением шума. См. <https://arxiv.org/pdf/1410.8251.pdf>. — Прим. перев.

Как оценивать векторные вложения — векторная арифметика и аналогии

Модель мешка слов создает документные векторы, которые отражают наличие и релевантность лексем документу. *Латентно-семантический анализ* в процессе снижает размерность этих векторов и выделяет то, что может быть интерпретировано как латентные понятия. *Латентное размещение Дирихле* представляет документы и термины в виде векторов, содержащих веса латентных тем.

Размерности словарных и фразовых векторов не имеют явного смысла. Однако векторные вложения кодируют схожее словопотребление как близость в латентном пространстве таким образом, что оно переносится на семантические взаимосвязи. Это приводит к интересным свойствам, когда аналогии могут выражаться путем добавления и вычитания словарных векторов.

На рис. 15.2 показано, как вектор, соединяющий векторы "Париж" и "Франция" (т. е. разницу их векторных вложений), отражает взаимосвязь "столица". Аналогичная взаимосвязь, "Лондон: Великобритания", соответствует тому же вектору, т. е. "Великобритания" находится очень близко к местоположению, полученному путем добавления вектора "столица" в вектор "Лондон".

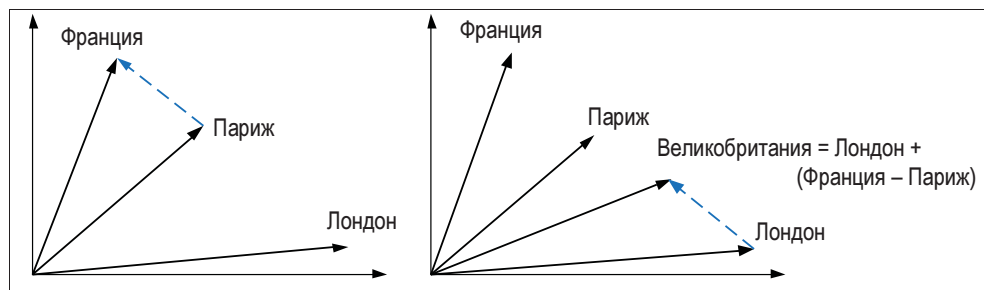


Рис. 15.2. Отражение взаимосвязей с помощью векторов вложений

Подобно тому, как слова могут использоваться в различных контекстах, они могут быть связаны с другими словами по-разному, и эти взаимосвязи соответствуют разным направлениям в латентном пространстве. Соответственно, существует несколько типов аналогий, которые векторные вложения должны отражать, если позволяют тренировочные данные.

Авторы модели Word2vec предоставляют список из нескольких тысяч взаимосвязей, охватывающих аспекты географии, грамматики и синтаксиса, а также семейных отношений, оценивая качество векторов вложения. Как показано выше, данный тест подтверждает, что целевое слово (Великобритания) находится ближе всего к результату добавления вектора, представляющего аналогичную связь (Париж: Франция), в целевое дополнение (Лондон).

На рис. 15.3 спроецированы 300-мерные векторные вложения наиболее тесно связанных аналогий для модели Word2vec, натренированной на корпусе Википедии

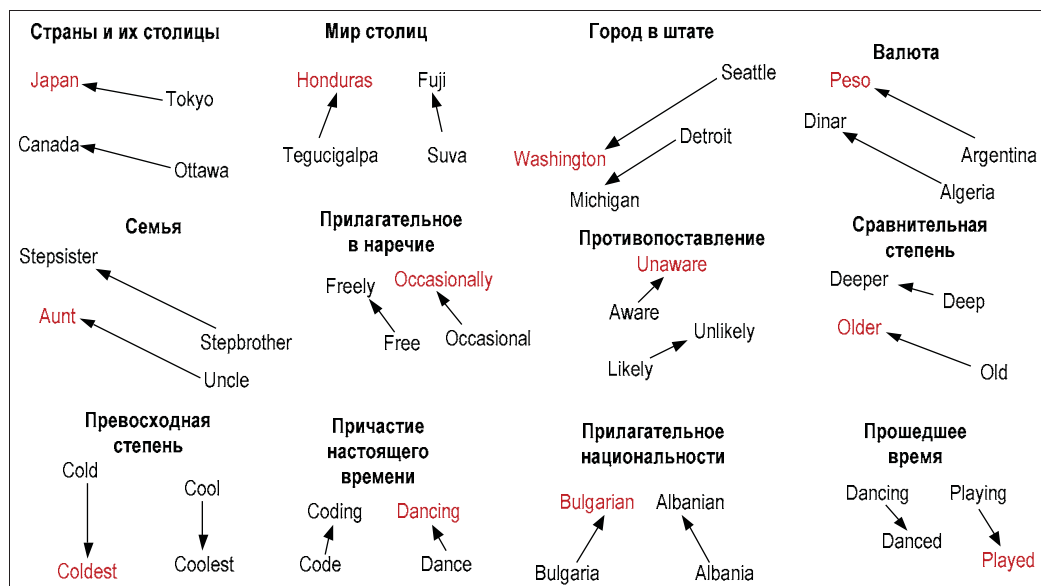


Рис. 15.3. Работа с моделями вложения слов в векторное пространство

с более чем 2 млрд лексем, в две размерности с использованием *анализа главных компонент* (PCA). Тест из более чем 24 400 аналогий из приведенных ниже категорий достиг точности более 73,5% (см. блокнот).

Подобно другим методам неконтролируемого автоматического обучения, цель усвоения векторов вложения — генерировать признаки для других задач, таких как классификация текста или сентиментный анализ.

Существует несколько вариантов получения векторов вложения для определенного корпуса документов:

- ♦ использовать вложения, усвоенные из общего крупного корпуса, такого как Википедия или Google Новости;
- ♦ натренировать собственную модель, используя документы, которые отражают интересующую область.

Чем менее общим и более специализированным является содержимое последующей задачи моделирования текста, тем предпочтительнее второй подход. Вместе с тем качественные векторные вложения слов требуют большого количества данных и информативных документов, содержащих сотни миллионов слов.

Как использовать предварительно натренированные словарные векторы

Существует несколько источников предварительно натренированных векторных вложений слов. Популярными вариантами являются глобальные векторы для представления слов GloVe (Global Vectors for Word Representation) и встроенные векторы библиотеки spaCy (подробности см. в блокноте `using_trained_vectors.ipynb`).

GloVe — глобальные векторы для представления слов

GloVe — это неконтролируемый алгоритм, разработанный в Стэнфордской лаборатории обработки ЕЯ, который усваивает векторные представления слов из агрегированной глобальной статистики совместного появления слов (см. справочные материалы). Векторы предварительно натренированы и доступны из следующих веб-масштабных источников:

- ◆ хранилище данных веб-обходчика Common Crawl (<http://commoncrawl.org/>) с 42 миллиардами или 840 миллиардами лексем и лексиконом на 1,9 или 2,2 млн лексем;
- ◆ корпус Wikipedia 2014 + Gigaword 5 с 6 миллиардами лексем и лексиконом на 400 тыс. лексем;
- ◆ корпус Twitter с 2 миллиардами твитов, 27 миллиардами лексем и лексиконом на 1,2 млн лексем.

Для конвертирования и загрузки векторных текстовых файлов в объект KeyedVector можно использовать библиотеку gensim:

```
from gensim.models import Word2vec, KeyedVectors
from gensim.scripts.glove2Word2vec import glove2Word2vec

glove2word2vec(glove_input_file=glove_wiki_file,
               word2vec_output_file=word2vec_wiki_file)
model = KeyedVectors.load_Word2vec_format(word2vec_wiki_file, binary=False)
```

Авторы Word2vec предоставляют текстовые файлы, содержащие более 24 тыс. тестов на аналогию, которые библиотека gensim использует для оценивания словарных векторов.

Словарные векторы, натренированные на корпусе Wikipedia, охватывают все аналогии и достигают общей точности 75,5% с некоторыми вариациями по категориям (табл. 15.1).

Таблица 15.1. Словарные векторы, натренированные на корпусе Wikipedia

| Категория | Образцы | Точность, % | Категория | Образцы | Точность, % |
|--------------------------|---------|-------------|-----------------------|---------|-------------|
| capital-common-countries | 506 | 94,86 | comparative | 1332 | 88,21 |
| capital-world | 8372 | 96,46 | superlative | 1 056 | 74,62 |
| city-in-state | 4242 | 60,00 | present-participle | 1056 | 69,98 |
| currency | 752 | 17,42 | nationality-adjective | 1640 | 92,50 |
| family | 506 | 88,14 | past-tense | 1560 | 61,15 |
| adjective-to-adverb | 992 | 22,58 | plural | 1332 | 78,08 |
| opposite | 756 | 28,57 | plural-verbs | 870 | 58,51 |

Векторы хранилища Common Crawl на 100 тыс. наиболее частых лексем охватывают около 80% аналогий и достигают несколько большей точности в 78%, в то время как векторы корпуса Twitter охватывают только 25% с точностью 62%.

Как тренировать собственные векторные вложения слов

Многие задачи требуют векторные вложения или специфичный для предметной области лексикон, который предварительно натренированные модели, основанные на обобщенном корпусе, представляют не совсем хорошо или не могут этого делать вообще. Стандартные модели Word2vec не могут назначать векторы словам вне лексикона и вместо этого используют вектор по умолчанию, что снижает их предсказательную значимость.

Например, во время работы с отраслевыми документами словарь или его использование могут со временем меняться по мере появления новых технологий или продукции. По этой причине векторные вложения тоже должны эволюционировать. Кроме того, в сообщениях о корпоративных заработках используются нюансы языка, не полностью отраженные в векторах GloVe, предварительно натренированных на статьях из Википедии.

Мы проиллюстрируем архитектуру модели Word2vec с помощью библиотеки Keras, которую представим более подробно в следующей главе, и более результативную адаптацию кода в библиотеке gensim, предоставленную авторами модели Word2vec. Блокнот word2vec.ipynb содержит дополнительные сведения о реализации, включая справочный материал со ссылкой на реализацию в библиотеке TensorFlow.

Архитектура модели Skip-Gram в библиотеке Keras

Для иллюстрации сетевой архитектуры модели Word2vec мы используем набор данных TED Talk с выровненными английскими и испанскими субтитрами, который мы впервые представили в *главе 13*.

Блокнот содержит исходный код для лексемизации документов и закрепления уникального идентификатора за каждым элементом лексикона. Нам требуется, по крайней мере, пять появлений в корпусе, при этом мы поддерживаем лексикон из 31 300 лексем.

Оценивание контрастивное к шуму

Библиотека Keras содержит метод `make_sampling_table`, который позволяет создавать тренировочный набор в виде пар контекстных и шумных слов с соответствующими метками, отобранными в соответствии с их частотами в корпусе.

В результате получается 27 млн положительных и отрицательных примеров контекстных и целевых пар.

Компоненты модели

Модель Skip-Gram содержит 200-мерный вектор вложений по каждому элементу лексикона, в результате давая 31 300×200 тренируемых параметров плюс два для сигмоидного выхода.

На каждой итерации модель вычисляет скалярное произведение контекстного и целевого векторов вложений, передает результат через сигмоиду, производя вероятности, и корректирует вложение на основе градиента потери.

Для визуализации модели используется TensorBoard — инструмент, который позволяет проецировать векторы вложений в три размерности для разведывания местоположения слов и фраз.

Словарные векторы из финансовой отчетности SEC с использованием библиотеки gensim

В этом разделе мы усвоим векторы слов и фраз из ежегодных документов финансовой отчетности, подаваемых в комиссию по ценным бумагам и биржам США (SEC), используя библиотеку `gensim` для того, чтобы проиллюстрировать потенциальную ценность векторных вложений слов для алгоритмической торговли. В следующих далее разделах мы объединим эти векторы как признаки с ценовыми финансовыми возвратами для того, чтобы натренировать нейронную сеть предсказывать цены долевого ценных бумаг по содержанию документов деловой отчетности.

В частности, мы используем набор данных, содержащий более 22 тыс. годовых отчетов по форме 10-K за период 2013–2016 гг., которые были поданы зарегистрированными на бирже компаниями и содержат как финансовую информацию, так и комментарии руководства (см. главу 3). Примерно для половины отчетов по форме 11-K компаний у нас есть цены на акции, которые можно использовать для разметки данных с целью предсказательного моделирования (подробности со справочными материалами об источниках данных и блокноты см. в папке `sec-filings`).

Предобработка

Каждый отчет представляет собой отдельный текстовый файл, а главный индекс содержит метаданные. Мы выделяем наиболее информативные разделы, а именно:

- ♦ *пункты 1 и 1A* — предпринимательская деятельность и рискованные факторы;
- ♦ *пункты 7 и 7A* — комментарии руководства и раскрытая информация о рыночных рисках.

Блокнот `preprocessing.ipynb` показывает, как разбирать и лексимизировать текст с помощью библиотеки `sraSu` аналогично подходу, принятому в главе 14. Мы не лемматизируем лексемы с целью сохранения нюансов словопотреблений.

Автоматическое обнаружение фраз

Мы используем библиотеку `gensim` для обнаружения фраз, как было показано ранее. Модуль `Phrases` выставляет отметки лексемам, а класс `Phraser` преобразует текстовые данные соответствующим образом. В блокноте показано, как повторять процесс для создания длинных фраз:

```
def create_ngrams(max_length=3):
    """Применение библиотеки для создания n-грамм"""

    n_grams = pd.DataFrame()
    start = time()
    for n in range(2, max_length + 1):
        print(n, end=' ', flush=True)

        sentences = LineSentence(f'ngrams_{n - 1}.txt')
        phrases = Phrases(sentences=sentences,
                           min_count=25, # игнорировать термины с низким количеством
                           threshold=0.5, # принимать фразы с более высокой отметкой
                           max_vocab_size=40000000, # обрезать менее частые слова
                                                # ради экономии памяти
                           delimiter=b' ', # как объединять n-граммные лексемы
                           progress_per=50000, # регистрировать продвижение каждые n
                           scoring='npmi')

        s = pd.DataFrame([[k.decode('utf-8'), v]
                           for k, v in phrases.export_phrases(sentences)]
                           , columns=['phrase', 'score']).assign(length=n)

        n_grams = pd.concat([n_grams, s])
        grams = Phraser(phrases)
        sentences = grams[sentences]
        Path(f'ngrams_{n}.txt').write_text('\n'.join([' '.join(s) for s in sentences]))

    n_grams = n_grams.sort_values('score', ascending=False)
    n_grams.phrase = n_grams.phrase.str.replace('_', ' ')
    n_grams['ngram'] = n_grams.phrase.str.replace(' ', '_')

    with pd.HDFStore('vocab.h5') as store:
        store.put('ngrams', n_grams)

create_ngrams()
```

Наиболее частые биграммы включают `common_stock` (обыкновенные акции), `united_states` (Соединенные Штаты), `cash_flows` (потоки денежных средств), `real_estate` (недвижимость) и `interest_rates` (процентные ставки).

Тренировка модели

Класс `gensim.models.Word2vec` реализует представленные ранее архитектуры SG и CBOW. Блокнот `Word2vec` содержит дополнительные сведения о реализации.

С целью обеспечения эффективного по памяти потребления текста класс `LineSentence` создает генератор из индивидуальных предложений, содержащихся в предоставленном текстовом файле:

```
sentence_path = Path('data', 'ngrams', f'ngrams_2.txt')
sentences = LineSentence(sentence_path)
```

Класс `Word2vec` предлагает ранее введенные параметры конфигурации:

```
model = Word2vec(sentences,
                  sg=1,          # 1=skip-gram; в противном случае CBOW
                  hs=0,          # иерархич. softmax, если 1;
                                # отрицательный отбор, если 0
                  size=300,      # размерность вектора
                  window=3,      # макс. расстояние между целевым
                                # и контекстным словом
                  min_count=50,   # игнорировать слова с более
                                # низкой частотой
                  negative=10,    # количество шумных слов для отрицательной
                                # выборки
                  workers=8,      # без нитей
                  iter=1,         # без эпох = итерации по корпусу
                  alpha=0.025,   # первоначальный темп усвоения
                  min_alpha=0.0001 # окончательный темп усвоения
                  )
```

Блокнот показывает, как сохранять состояние и перезагружать модели для продолжения тренировки или как хранить векторы вложения отдельно, например, для использования в автоматически обучающихся моделях.

Оценивание модели

Базовый функционал включает выявление похожих слов:

```
sims = model.wv.most_similar(positive=['iphone'],
                             restrict_vocab=15000)
print(pd.DataFrame(sims, columns=['term', 'similarity']))
```

| | term | similarity |
|---|---------------------|------------|
| 0 | android | 0.600454 |
| 1 | smartphone | 0.581685 |
| 2 | app | 0.559129 |
| 3 | smartphones | 0.533848 |
| 4 | smartphones_tablets | 0.526129 |

| | | |
|---|--------------|----------|
| 5 | handsets | 0.514813 |
| 6 | smart_phones | 0.512868 |
| 7 | apple | 0.507795 |
| 8 | apps | 0.505517 |
| 9 | handset | 0.491526 |

Мы также можем проверить индивидуальные аналогии, используя соответственно положительные и отрицательные вклады:

```
analogy = model.wv.most_similar(positive=['france', 'london'],
                                negative=['paris'],
                                restrict_vocab=15000)
print(pd.DataFrame(analogy, columns=['term', 'similarity']))
```

| | term | similarity |
|---|----------------|------------|
| 0 | united_kingdom | 0.606630 |
| 1 | germany | 0.585644 |
| 2 | netherlands | 0.578868 |
| 3 | italy | 0.547168 |
| 4 | india | 0.545213 |
| 5 | spain | 0.539029 |
| 6 | singapore | 0.535106 |
| 7 | australia | 0.525464 |
| 8 | belgium | 0.523677 |
| 9 | sweden | 0.510462 |

Влияние параметрических настроек на результативность

Аналогии могут использоваться для оценивания влияния разных параметрических настроек. Следующие ниже результаты выделяются особо (см. подробные результаты в папке `models`):

- ◆ отрицательная выборка превосходит по результативности иерархическую активацию `softmax` и при этом тренируется быстрее;
- ◆ архитектура `Skip-Gram` превосходит по результативности `CBOW` при заданной целевой функции;
- ◆ разные значения параметра `min_count` имеют меньшее влияние, причем середина точка, равная 50, показывает наилучшие результаты.

Дальнейшие эксперименты с самой результативной моделью `SG` с использованием отрицательной выборки и параметром `min_count`, равным 50, показывают следующее:

- ◆ контекстные окна меньше пяти снижают результативность;
- ◆ более высокий темп отрицательной выборки повышает результативность за счет более медленной тренировки;
- ◆ более крупные векторы улучшают результативность, причем размер 600 дает наилучшую точность на уровне 38,5%.

Сентиментный анализ с помощью модели Doc2vec

Классифицирование текста требует совмещения нескольких векторных вложений слов. Общепринятым подходом является усреднение векторов вложения для каждого слова в документе. В указанном подходе используется информация из всех векторных вложений и эффективно задействуется векторное сложение, при этом получается другая точка локализации в пространстве вложения. Однако релевантная информация о порядке слов теряется.

В отличие от этого, современный подход к генерированию векторных вложений для фрагментов текста, таких как абзац или отзыв о продукции, должен использовать модель векторных вложений документов Doc2vec. Эта модель была разработана авторами модели векторных вложений слов Word2vec вскоре после публикации их первоначального инновационного вклада.

Подобно модели Word2vec, существуют две разновидности модели Doc2vec.

- ◆ *Модель распределенного мешка слов* (distributed bag of words, DBOW) соответствует архитектуре CBOW модели Word2vec. Документные векторы являются результатом тренировки сети синтетической задаче предсказания целевого слова и основываются на векторах контекстных слов и векторе doc документа.
- ◆ *Модель распределенной памяти* (distributed memory, DM) соответствует архитектуре Skip-Gram модели Word2vec. Векторы doc являются результатом тренировки нейронной сети для предсказания целевого слова с использованием вектора doc полного документа.

Класс Doc2vec библиотеки gensim реализует этот алгоритм.

Тренировка модели Doc2vec на сентиментных данных Yelp

Мы используем случайную выборку из 500 тыс. отзывов о деятельности предприятий с веб-сайта Yelp (см. главу 13) с соответствующими звездными рейтингами (см. блокнот `yelp_sentiment.ipynb`):

```
df = pd.read_parquet('combined.parquet', engine='fastparquet').loc[:, ['stars', 'text']]
stars = range(1, 6)
sample = pd.concat([df[df.stars==s].sample(n=100000) for s in stars])
sample.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 500000 entries, 52085 to 3365007
Data columns (total 2 columns):
stars      500000 non-null int64
text       500000 non-null object
```

```
dtypes: int64(1), object(1)
memory usage: 11.4+ MB
```

Мы применяем простую предобработку, удаляя стоп-слова и знаки препинания с помощью лексемизатора библиотеки NLTK, и отбрасываем отзывы с менее 10 лексемами:

```
import nltk
nltk.download('stopwords')
from nltk import RegexpTokenizer
from nltk.corpus import stopwords

tokenizer = RegexpTokenizer(r'\w+')
stopword_set = set(stopwords.words('english'))

def clean(review):
    tokens = tokenizer.tokenize(review)
    return ' '.join([t for t in tokens if t not in stopword_set])

sample.text = sample.text.str.lower().apply(clean)
sample = sample[sample.text.str.split().str.len()>10]

sample.sample(n=10)

      stars  text
3713191 1    called 938 placed order informer ian manager a...
3632813 3    ok best tip sell stuff buffalo exchange sharin...
1414414 5    afford rooms well worth money absolutely amazi...
...
1354353 5    rita must love custard black cherry little bit...
2760    1    drittes goa pfaffing erlebt absolut nix unterh...
1118726 1    visited week ago im finally writing review pla...
```

Создание входных данных

Класс `gensim.models.doc2vec` обрабатывает документы в формате `TaggedDocument`, содержащем лексемизированные документы и уникальный тег, позволяющий обращаться к документным векторам после тренировки:

```
sentences = []
for i, (_, text) in enumerate(sample.values):
    sentences.append(TaggedDocument(words=text.split(), tags=[i]))
```

Тренировочный интерфейс работает аналогично модели `Word2vec` с дополнительными параметрами, относящимися к спецификации задачи для алгоритма `Doc2vec`:

```
model = Doc2Vec(documents=sentences,
                dm=1,           # алгоритм: использовать распределенную память
                size=300,
```

```

window=5,
alpha=0.025, # первоначальный темп усвоения
min_count=0,
workers=8,
epochs=5,
negative=5,
dm_concat=0, # 1: конкатенир., не сумм./усредн.
               # контекстные векторы
dbow_words=0) # 1: натренир. словарные векторы,
               # 0: только векторы doc

```

```
model.save('test.model')
```

Вы также можете использовать метод `train()` для того, чтобы продолжать процесс усвоения и, например, итеративно снижать темп усвоения:

```

for _ in range(10):
    alpha *= .9
    model.train(sentences,
                total_examples=model.corpus_count,
                epochs=model.epochs,
                alpha=alpha)

```

В результате мы можем обращаться к документным векторам как признакам, тренируя сентиментный классификатор:

```

X = np.zeros(shape=(len(sample), size))
y = sample.stars.sub(1) # модели требуется [0, 5) меток
for i in range(len(sample)):
    X[i] = model[i]

```

```
X.shape
```

```
(485681, 300)
```

Мы тренируем градиентно-бустинговую машину `lightgbm` следующим образом:

1. Создаем объекты `lightgbm` класса `Dataset` из тренировочного и тестового наборов:

```

train_data = lgb.Dataset(data=X_train, label=y_train)
test_data = train_data.create_valid(X_test, label=y_test)

```

2. Задаем тренировочные параметры для мультиклассовой модели с пятью классами (в противном случае используя значения по умолчанию):

```

params = {'objective' : 'multiclass',
          'num_classes': 5}

```

3. Тренируем модель в течение 250 итераций и отслеживаем ошибку на контрольном наборе:

```
lgb_model = lgb.train(params=params,
                      train_set=train_data,
                      num_boost_round=250,
                      valid_sets=[train_data, test_data],
                      verbose_eval=25)
```

4. Градиентно-бустинговая машина `lightgbm` предсказывает вероятности для всех пяти классов. Мы получаем предсказания классов с помощью функции `np.argmax()`, которая выдает индекс столбца с наибольшей предсказанной вероятностью:

```
y_pred = np.argmax(lgb_model.predict(X_test), axis=1)
```

5. Мы вычисляем отметку точности результата, причем видим улучшение более чем на 100% по сравнению с базовым уровнем 20% для пяти сбалансированных классов:

```
accuracy_score(y_true=y_test, y_pred=y_pred)
```

```
0.44955063467061984
```

6. Наконец, мы подробнее рассматриваем предсказания для каждого класса, используя матрицу несоответствий:

```
cm = confusion_matrix(y_true=y_test, y_pred=y_pred)
cm = pd.DataFrame(cm / np.sum(cm), index=stars, columns=stars)
```

7. И визуализируем результат в виде тепловой карты из библиотеки `Seaborn` (рис. 15.4):

```
sns.heatmap(cm, annot=True, cmap='Blues', fmt='.1%')
```

В итоге метод `doc2vec` позволил нам достичь очень существенного улучшения тестовой точности по сравнению с наивным эталоном без какой-либо настройки. Если

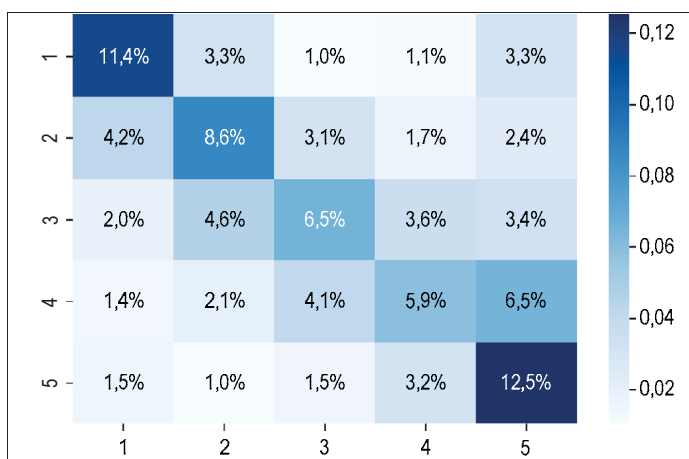


Рис. 15.4. Тепловая карта с результатами

мы выберем только верхние и нижние отзывы (соответственно с пятью и одной звездой) и натренируем бинарный классификатор, то метрика AUC достигнет более 0,86, используя 250 тыс. образцов из каждого класса.

Бонус — модель Word2vec для машинного перевода

Блокнот `translation.ipynb` демонстрирует, что взаимосвязи, закодированные на одном языке, часто соответствуют аналогичным взаимосвязям на другом языке.

Он иллюстрирует, как словарные векторы могут использоваться для перевода слов и фраз путем проецирования словарных векторов из пространства векторного вложения одного языка в пространство другого языка, используя для этого трансляционную матрицу.

Резюме

Эта глава началась с описания того, насколько эффективнее векторные вложения слов кодируют семантику индивидуальных лексем, чем модель мешка слов, которую мы использовали в *главе 13*. Мы также увидели, как выполнять оценивание векторного вложения, подтверждая правильность представления семантических взаимосвязей между словами с помощью линейной векторной арифметики.

Для усвоения векторных сложений слов мы используем мелкие нейронные сети, которые раньше тренировались медленно в масштабе веб-данных, содержащих миллиарды лексем. Модель векторных вложений слов Word2vec сочетает несколько алгоритмических инноваций, которые позволяют значительно ускорить тренировку, и установила новый стандарт для генерирования текстовых признаков. Мы увидели, как использовать предварительно натренированные словарные векторы, воспользовавшись библиотеками `sraCy` и `gensim`, и научились тренировать собственные векторные вложения слов. Затем мы применили модель Word2vec к документам финансовой отчетности, подаваемым в комиссию SEC. Наконец, мы рассмотрели расширенную модель векторных вложений документов Doc2vec, которая усваивает векторные представления документов аналогично словарным векторам, и применили ее к отзывам о деятельности предприятий с веб-сайта Yelp.

Теперь можно приступить к оставшимся главам (доступным в электронном архиве, как было упомянуто в *предисловии*), посвященным глубокому обучению, начиная с введения в сети прямого распространения, популярные каркасы глубокого обучения и технические приемы для эффективной тренировки в широком масштабе.

16

Дальнейшие действия

Цель этой книги заключалась в том, чтобы дать вам возможность применять машинное обучение (МО) к различным источникам данных и извлекать сигналы, полезные для разработки и исполнения инвестиционной стратегии. С этой целью мы ввели МО как важный элемент развития стратегии торговли на финансовых рынках. Мы увидели, что МО может добавлять стоимость на нескольких этапах в процессе конструирования, тестирования, исполнения и оценивания стратегии.

Стало ясно, что стержневое ценностное предложение машинного обучения состоит в способности извлекать действенную информацию из гораздо более крупных объемов данных более систематически, чем когда-либо могли делать эксперты-люди. С одной стороны, это ценностное предложение действительно приобрело популярность благодаря взрывному росту цифровых данных, что сделало его более перспективным и необходимым для задействования вычислительных мощностей по обработке данных. С другой стороны, применение МО по-прежнему требует значительного вмешательства человека и компетентных знаний для определения целей, отбора и обработки данных, разработки и оптимизации модели и надлежащего использования результатов.

В этой заключительной главе мы кратко изложим ключевые инструменты, приложения и уроки, извлеченные из этой книги для того, чтобы после стольких подробностей не упустить из виду общую картину. Затем мы определим области, которые мы не охватили и на которых было бы целесообразно сосредоточиться по мере того, как вы будете стремиться расширять многие введенные нами технические решения МО и будете становиться продуктивными в ежедневном их применении. Кроме того, мы выделим ряд навыков, которые представляют ценность для повышения индивидуальной результативности.

Итак, в этой главе мы рассмотрим следующие темы:

- ◆ краткий обзор основных результатов и извлеченных уроков;
- ◆ определение следующих шагов, которые необходимо предпринять на основе методов, описанных в этой книге;
- ◆ предложение способов встраивания МО в собственный инвестиционный процесс.

Ключевые итоги и извлеченные уроки

Важные сущностные идеи, которые следует иметь в виду, когда вы переходите к практическому применению МО для торговли на финансовых рынках, включают следующее:

- ◆ данные — это единственный самый важный ингредиент;
- ◆ компетентные знания предметной области помогают реализовывать потенциальную ценность данных, в особенности в области финансов;
- ◆ МО предлагает инструменты для многих случаев использования, которые должны быть доработаны и объединены с целью создания решений новых задач с использованием данных;
- ◆ широкий выбор модельных целевых задач и диагностики результативности является ключевым для продуктивного последовательного продвижения в направлении оптимальной системы;
- ◆ переподгонка во время бэктестирования на исторических данных является серьезным вызовом, требующим значительного внимания;
- ◆ прозрачность вокруг черно-ящичных моделей помогает укреплять доверие и обеспечивать их принятие.

Мы подробнее остановимся на каждой из этих идей.

Данные — единственный самый важный ингредиент

Рост популярности МО в торговле на финансовых рынках и во всех остальных сферах в значительной степени дополняет взрывной рост данных, который мы рассмотрели очень подробно. В *главе 2* мы проиллюстрировали, как получать доступ и работать с этими источниками данных, исторически являющимися основой количественных инвестиций. В *главе 3* мы представили каркас с критериями определения потенциальной ценности альтернативных наборов данных.

Ключевое сущностное осознание заключается в том, что современные методы МО, такие как глубокие нейронные сети, являются успешными по той причине, что их предсказательная результативность продолжает улучшаться вместе с большим количеством данных. С другой стороны, сложность модели и данных должны сочетаться так, чтобы уравнивать компромисс между смещением и дисперсией. Управление качеством данных и интеграция наборов данных являются ключевыми шагами в реализации их потенциальной ценности.

Контроль качества

Как и нефть, являющаяся популярным в наши дни сравнением, данные проходят по конвейеру в несколько этапов от сырой формы до рафинированного продукта, который может служить топливом для стратегии торговли. Крайне важно уделять

пристальное внимание качеству конечного продукта с целью получения из него желаемого пробега.

Иногда вы получаете данные в сыром виде и управляете многочисленными преобразованиями, необходимыми для ваших целей. Но чаще всего вы имеете дело с промежуточным продуктом, и от вас требуется получить ясность о том, что именно данные измеряют на текущий момент.

В отличие от нефти, объективного стандарта качества данных зачастую не существует, т. к. источники данных продолжают расширяться. Вместо этого качество зависит от содержимого их сигнала, который, в свою очередь, зависит от ваших инвестиционных целей. Затратоэффективное оценивание новых наборов данных требует продуктивного рабочего потока, включая соответствующую инфраструктуру, о которой мы поговорим в следующем разделе.

Интеграция данных

Ценность данных для инвестиционной стратегии часто зависит от сочетания дополнительных источников рыночных, фундаментальных и альтернативных данных. Мы увидели, что предсказательная мощь автоматически обучающихся алгоритмов, таких как древесные ансамбли или нейронные сети, частично обусловлена их способностью обнаруживать нелинейные связи, в частности эффекты взаимодействия между переменными.

Способность модулировать влияние переменной как функции других признаков модели процветает на тех входных данных, которые улавливают различные аспекты целевого результата. Сочетание цен на активы с макроэкономическими фундаментальными данными, данными социальных настроений, кредитно-карточных платежей и спутниковых снимков, вероятно, даст значительно более надежные предсказания в различных экономических и рыночных режимах, чем каждый источник сам по себе (при условии, что данные достаточно велики для того, чтобы иметь возможность усваивать скрытые связи).

Работа с данными из многочисленных источников увеличивает проблемы правильного закрепления меток. Очень важно назначать точные временные метки во избежание систематического смещения из-за забегания вперед, когда алгоритм тестируется с данными, которые фактически стали доступными позже. Например, данные могут иметь временные метки, назначенные поставщиком, которые требуют внесения поправки с целью отражения определенной точки во времени, когда они были бы доступны для живого алгоритма.

Компетентное знание предметной области помогает разблокировать ценность данных

Мы подчеркнули, что не только данные являются необходимой движущей силой успешных приложений МО, но и компетентное знание предметной области, которое имеет решающее значение для стратегического направления, выработки признаков, отбора данных и конструирования модели.

В любой области у практиков есть теории о движущих силах ключевых результатов и взаимосвязях между ними. Финансы выделяются объемом соответствующих количественных исследований, как теоретических, так и эмпирических. Маркос Лопес де Прадо и другие (справочные материалы см. в репозитории GitHub <https://github.com/PacktPublishing/Hands-On-Machine-Learning-for-Trading>) критикуют большинство эмпирических результатов, учитывая повсеместный глубокий анализ данных, который может делать находки недействительными. Тем не менее, четкое понимание работы финансовых рынков должно лежать в основе отбора и использования данных, а также обоснования стратегий, основанных на машинном обучении. Мы изложили эти ключевые идеи в *главах 4 и 5*.

С другой стороны, новые технические решения в области МО, вероятно, породят новые гипотезы о движущих силах финансовых результатов, которые наполнят теорию МО, а затем должны быть независимо протестированы.

Выработка признаков и исследование альфа-факторов

В большей степени, чем сырые данные, выработка признаков часто является ключом к тому, чтобы сделать сигнал полезным для алгоритма. Мобилизация десятилетий исследований рискованных факторов, которые стимулируют возвратность на теоретических и эмпирических основаниях, является хорошей отправной точкой для приоритизации преобразований данных, которые с большей вероятностью будут отражать релевантную информацию.

Вместе с тем стоит отметить, что только изобретательная выработка признаков приведет к инновационным стратегиям, которые смогут конкурировать на рынке в течение длительного времени. Убедительное изложение, в особенности в отношении новых альфа-факторов, объясняющее, как они работают, с учетом устоявшихся идей о динамике рынка и поведении инвесторов, обеспечит более высокую уверенность при размещении капитала.

Риски ложных обнаружений и переподгонки к историческим данным делают еще более необходимым не позволять данным говорить самим за себя, а приоритизировать стратегии перед началом тестирования. Мы рассмотрели вопрос дефлирования коэффициента Шарпа в свете числа экспериментов.

Машинное обучение — это комплект инструментов для решения задач с помощью данных

МО предлагает алгоритмические методы и технические решения, которые могут применяться ко многим случаям использования. Все главы, начиная со второй (как упоминалось в *главе 1*), представили МО как набор разнообразных инструментов, которые могут добавлять стоимость на различных этапах развития торговой стратегии, включая:

- ◆ генерирование идей и исследование альфа-факторов;
- ◆ агрегирование сигналов и оптимизирование портфеля;

- ◆ тестирование стратегии;
- ◆ исполнение сделок;
- ◆ оценивание стратегии.

Более того, автоматически обучающиеся алгоритмы предназначены для дальнейшего их развития, адаптации и комбинирования при решении новых задач в различных контекстах. По этим причинам важно понимать ключевые концепции и идеи, лежащие в основе этих алгоритмов, в дополнение к способности применять их к данным для продуктивных экспериментов и исследований, как описано в *главе 6*.

Кроме того, наилучшие результаты часто достигаются путем совмещения людей-экспертов с инструментами МО. В *главе 1* мы рассмотрели квантоментальный инвестиционный стиль, в котором сходятся дискреционная торговля и алгоритмическая торговля. Этот подход, вероятно, будет и далее приобретать все большее значение и будет зависеть от гибкого и изобретательного применения охваченных нами основополагающих инструментов и их расширений на всевозможные наборы данных.

Диагностика модели помогает ускорить оптимизацию

В *главе 6* мы изложили некоторые из наиболее важных концепций. Автоматически обучающиеся алгоритмы усваивают связи между входными данными и целью, принимая допущения о функциональной форме. Если самообучение алгоритма основано на шуме, а не на сигнале, то страдает предсказательная результативность.

Разумеется, мы пока не знаем, как отделять сигнал и шум от перспективы завтрашних результатов. Диагностика модели, например с помощью кривых усвоения и верификационного теста оптимизации, помогает облегчать эту фундаментальную трудность и калибровать выбор или конфигурацию алгоритма для данных или поставленной задачи. Эта задача может быть упрощена путем определения сфокусированных модельных целей, а в случае многосложных моделей, путем проведения грани между недостаточной результативностью из-за проблем, связанных с оптимизационным алгоритмом, или недостаточной результативностью из-за самой цели.

Принятие решений без бесплатного обеда

Ни одна система, ни компьютерная, ни человеческая, не имеет оснований для надежного предсказания результатов относительно новых примеров, помимо тех, которые она наблюдала во время тренировки. Единственный выход — иметь некоторые дополнительные предшествующие знания или принимать допущения, которые выходят за рамки тренировочных примеров. Мы рассмотрели широкий спектр алгоритмов, начиная с *глав 7* и *8* и заканчивая нелинейными ансамблями в *главах 10* и *11*, а также нейронными сетями в различных главах из электронного архива настоящей книги.

Мы увидели, что линейная модель делает сильное допущение о том, что связь между входами и выходами имеет очень простую форму, в то время как модели, обсуждаемые позже, призваны усваивать более сложные функции. Хотя вполне очевидно, что простая модель потерпит неудачу в большинстве случаев, сложная модель оказывается не всегда лучше. Если истинная связь является линейной, но данные являются шумными, то сложная модель заучит шум как часть сложной связи, из существования которой она исходит. В этом заключается базовая идея, лежащая в основе теоремы об отсутствии бесплатных обедов, которая постулирует, что не существует алгоритма, который был бы универсально превосходящим для любой задачи. Хорошая подгонка в некоторых случаях происходит за счет низкой результативности в других местах.

Ключевыми инструментами для приспособливания выбранного алгоритма к данным являются разведывательный анализ данных и эксперименты, основанные на понимании допущений, которые делает модель.

Управление компромиссом между смещением и дисперсией

Другой взгляд на проблему адаптации алгоритма к данным — это компромисс между смещением и дисперсией, которые вызывают ошибки предсказания, выходящие за пределы естественной шумности данных. Простая модель, которая не улавливает связи в данных адекватно, будет недостаточно подогнана и проявлять смещение, т. е. делать систематически неправильные предсказания. Модель, которая слишком сложна, будет переподогнана и усваивать шум в дополнение к любому сигналу, вследствие чего результат покажет много дисперсии для разных образцов.

Ключевым инструментом для диагностирования этого компромисса на любой итерации процесса отбора и оптимизации модели является кривая усвоения (самообучения). Она показывает, как тренировочные и контрольные ошибки зависят от размера выборки. Это позволяет выбирать между разными вариантами повышения результативности: отрегулировать сложность модели или получить больше точек данных.

Чем ближе ошибка на тренировочных данных к человеческим ошибкам или другим эталонам, тем больше вероятность того, что модель переподогнана. Низкая ошибка на контрольных данных говорит о том, что нам повезло, и мы нашли хорошую модель. Если контрольная ошибка является высокой, то нам не повезло. Однако если она продолжает снижаться с размером тренировочного набора, то, вероятно, поможет больший объем данных. Если тренировочная ошибка является высокой, то большее количество данных вряд ли поможет, и вместо этого мы должны добавить признаки либо применить более гибкий алгоритм.

Определение адресных модельных целей

Одним из первых шагов в процессе автоматического обучения является определение оптимизационной цели алгоритма. Иногда выбор является простым, например, в регрессионной задаче. Классификационная задача может быть сложнее, напри-

мер, когда нас интересуют прецизионность и полнота. Консолидация конфликтующих целей в единую метрику, такую как показатель F1, помогает сосредоточить оптимизационные усилия. Мы также можем включить условия, которые должны не оптимизироваться, а соблюдаться. Мы убедились, что суть подкрепляемого обучения заключается в определении правильной функции вознаграждения, которая направляет процесс самообучения агента.

Верификационная проверка оптимизации

Эндрю Нг (Andrew Ng) (см. справочные материалы в репозитории GitHub: <https://github.com/PacktPublishing/Hands-On-Machine-Learning-for-Algorithmic-Trading>) подчеркивает различие между недостаточной результативностью, вызванной проблемами, связанными с обучающимся и оптимизационным алгоритмами. Многосложные модели, такие как нейронные сети, исходят из нелинейности связей, и поисковый процесс оптимизационного алгоритма может закончиться оптимумом, который является локальным, а не глобальным.

Если, например, модели не удастся правильно перевести фразу, то указанный тест сравнивает отметки за правильное предсказание и решение, обнаруженное поисковым процессом оптимизационного алгоритма. Если обучающийся алгоритм получает за правильное решение более высокую отметку, то поисковый алгоритм требует улучшений. В противном случае обучающийся алгоритм оптимизируется с неправильным целевым объектом.

Остерегайтесь перепогонки к историческим данным

На протяжении всей книги мы неоднократно освещали риски ложных обнаружений из-за бэктестовой перепогонки, т. е. перепогонки к историческим данным. В *главе 5* изложены основные движущие силы и потенциальные рецептуры. Низкое отношение шума к сигналу и относительно малые наборы данных (по сравнению со снимковыми или текстовыми веб-масштабными данными) делают эту проблему особенно серьезной в области торговли на финансовых рынках. Осведомленность имеет решающее значение, поскольку легкость доступа к данным и инструментам с целью применения МО увеличивает риски экспоненциально.

Избежать бэктестовой перепогонки невозможно, потому что не существует метода ее предотвращения. Однако мы представили методы корректировки бэктестовых метрик с целью учета многократных испытаний, таких как коэффициент Шарпа. Во время работы над живой торговой стратегией частью процесса реализации стратегии должна быть стадия бумажной (виртуальной) торговли и внимательного мониторинга результативности во время исполнения на рынке.

Как проникать в сущность черно-ящичных моделей

Глубокие нейронные сети и многосложные ансамбли могут вызывать подозрение, когда они считаются непроницаемыми черно-ящичными моделями, в частности,

в свете риска бэктекстовой переподгонки. В *главе 11* мы представили несколько методов, позволяющих проникать в суть того, как эти модели делают предсказания.

В дополнение к обычным мерам важности признаков недавнее теоретико-игровое нововведение *аддитивных объяснений Шепли* (SHapley Additive exPlanations, SHAP) является значительным шагом к пониманию механики многосложных моделей. Значения SHAP позволяют точно относить признаки и их значения к предсказаниям, вследствие чего становится легче подтверждать логику модели в свете конкретных теорий о поведении рынка для определенной инвестиционной цели. Помимо обоснования, точные отметки важности признаков и атрибуция предсказаний позволяют глубже понимать движущие силы интересующего инвестиционного результата.

С другой стороны, существуют некоторые разногласия в отношении того, до какой степени прозрачность вокруг модельных предсказаний должна быть самоцелью. Джеффри Хинтон (Geoffrey Hinton), один из изобретателей глубокого автоматического обучения, утверждает, что человеческие решения также нередко туманны, и машины также должны оцениваться по их результатам, чего мы ожидаем от инвестиционных менеджеров.

Машинное обучение для торговли на финансовых рынках на практике

По мере интегрирования многочисленных инструментов и технических решений в ваш инвестиционный и торговый процесс, вы можете сконцентрировать свои усилия на целом ряде вещей. Если ваша цель состоит в том, чтобы принимать оптимальные решения, то вы должны выбирать проекты, которые реалистичны, но амбициозны, с учетом вашего текущего набора навыков. Это поможет вам разрабатывать эффективный рабочий поток, подкрепленный производительными инструментами, и получать практический опыт.

Мы кратко перечислим некоторые инструменты, которые полезны для расширения экосистемы языка Python, описанной в этой книге, и обратимся к справочным материалам, перечисленным в репозитории GitHub для более глубокого погружения в тему. К ним относятся технологии больших данных, которые в итоге будут необходимы для реализации МО для стратегий торговли в широком масштабе. Мы также перечислим некоторые платформы, которые позволяют реализовывать стратегии торговли с использованием языка Python, часто с доступом к источникам данных и алгоритмам и библиотекам МО.

Технологии управления данными

Центральная роль данных в процессе МО требует знакомства с рядом технологий по хранению, преобразованию и анализу данных в широком масштабе, включая использование облачных служб, таких как Amazon Web Services, Azure и Google Cloud.

Системы управления базами данных

Хранение данных подразумевает использование баз данных, в которых исторически доминируют *реляционные системы управления базами данных* (RDBMS, РСУБД), использующие SQL для хранения и извлечения данных в четко определенном табличном формате с коммерческими поставщиками, такими как Oracle и Microsoft, и реализациями с открытым исходным кодом, такими как PostgreSQL и MySQL. В последнее время появились альтернативы, которые часто коллективно обозначаются как NoSQL, но являются довольно разнообразными.

- ◆ *Хранилище в формате ключ-значение* — быстрый доступ операций чтения/записи к объектам. Мы рассмотрели формат HDF5 в *главе 2*, который обеспечивает быстрый доступ к кадру данных библиотеки pandas.
- ◆ *Столбчатое хранилище* задействует однородность данных в столбце, обеспечивая сжатие и более быстрые операции на основе столбцов, такие как агрегация. Используется в популярных технических решениях в виде хранилищ данных Amazon Redshift, Apache Parquet, Cassandra или Google Big Table.
- ◆ *Документное хранилище* предназначено для хранения данных, которые не соответствуют жестко заданной схеме, требуемой в СУБД. Популяризируется веб-приложениями, использующими формат JSON или XML, с которыми мы столкнулись в *главе 4*. Используется, например, в MongoDB.
- ◆ *Графовая база данных* предназначена для хранения сетей с узлами и ребрами и специализируется на запросах о сетевых метриках и взаимосвязях. Применяется в Neo4J и Apache Giraph.

В хорошо зарекомендовавших себя системах реляционных баз данных был принят определенный переход к указанным выше способам хранения. Экосистема языка Python обеспечивает взаимодействие со многими стандартными источниками данных и обеспечивает быстрые форматы HDF5 и Parquet, как было продемонстрировано на протяжении всей книги.

Технологии больших данных — Hadoop и Spark

Управление данными в широком масштабе, т. е. в сотнях гигабайтах и более, требует использования нескольких машин, образующих кластер для параллельного выполнения операций чтения, записи и вычисления, т. е. он распределен по разным машинам.

Экосистема Hadoop стала открытым программным каркасом для распределенного хранения и обработки больших данных с использованием разработанной Google модели программирования MapReduce. Указанная экосистема диверсифицировалась под крышей Apache Foundation и сегодня включает в себя многочисленные проекты, которые охватывают различные аспекты управления данными в широком масштабе. Ключевыми инструментами в среде Hadoop являются:

- ◆ *Apache Pig* — язык обработки данных для реализации крупномасштабных конвейеров "извлечение — преобразование — загрузка" (extract-transform-load, ETL) с использованием технологии MapReduce, разработанный в Yahoo;

- ◆ *Apache Hive* — стандарт де-факто для интерактивных SQL-запросов на петабайтах данных, разработанный в Facebook;
- ◆ *Apache HBASE* — база данных NoSQL для реально-временного доступа операций чтения/записи, которая линейно масштабируется до миллиардов строк и миллионов столбцов и может объединять источники данных с использованием разных схем.

Платформа Apache Spark стала самой популярной в сфере интерактивной аналитики в кластере. Каркас MapReduce допускал параллельные вычисления, но требовал многократных операций чтения/записи с диска для обеспечения избыточности данных. Платформа Spark кардинально ускорила ход вычисления в широком масштабе вследствие *отказоустойчивой распределенной структуры данных* RDD (Resilient Distributed Data), которая допускает сильно оптимизированное вычисление прямо в оперативной памяти. Сюда входят итеративные вычисления, необходимые для оптимизации, например, для градиентного спуска для многочисленных автоматически обучающихся алгоритмов.

Инструменты машинного обучения

В этой книге мы рассмотрели многочисленные библиотеки экосистемы Python. Язык Python эволюционировал и стал излюбленным языком для науки о данных и МО, и набор библиотек с открытым исходным кодом продолжает диверсифицироваться и созреть, опираясь на прочный фундамент научно-вычислительных библиотек NumPy и SciPy. Планируется, что появится версия 1.0 популярной библиотеки pandas, которая внесла значительный вклад в популяризацию использования Python для науки о данных. Интерфейс библиотеки Scikit-learn стал стандартом для современных библиотек МО, таких как XGBoost или LightGBM, которые часто взаимодействуют с различными инструментами автоматизации рабочих процессов, такими как GridSearchCV и Pipeline, которые мы неоднократно использовали на протяжении всей книги.

Существует несколько поставщиков, которые стремятся обеспечивать рабочий процесс МО.

- ◆ Компания H₂O.ai (<https://www.h2o.ai/>) предлагает демократизирующую ИИ платформу H₂O, которая интегрирует в себе облачные вычисления с автоматизацией МО. Она позволяет пользователям выполнять подгонку тысяч потенциальных моделей к своим данным для разведывания регулярностей в данных. Она имеет интерфейсы на Python, а также R и Java.
- ◆ Компания DataRobot ставит своей целью автоматизировать процесс разработки моделей, предоставляя платформу для быстрого построения и развертывания предсказательных моделей в облаке или в локальной системе.
- ◆ Компания Dataiku предоставляет коллаборативную платформу науки о данных, призванную помогать аналитикам и инженерам разведывать, прототипировать, строить и поставлять собственные продукты данных.

Существует также несколько инициатив с открытым исходным кодом, возглавляемых компаниями, которые развивают и расширяют экосистему Python:

- ◆ квантитативный хедж-фонд Two Sigma вносит инструменты квантитативного анализа в среду блокнота Jupyter в рамках проекта beakerx;
- ◆ агрегатор данных Bloomberg интегрировал блокнот Jupyter в свой терминал с целью обеспечения интерактивного анализа своих финансовых данных.

Онлайновые торговые платформы

Главными средствами разработки стратегий торговли с использованием машинного обучения являются онлайн-платформы, которые часто разыскивают успешные стратегии торговли и размещают в них капитал. Популярные решения включают платформы Quantopian, QuantConnect, QuantRocket и более позднюю Alpha Trading Labs, которая сосредоточена на высокочастотной торговле.

Кроме того, платформа Interactive Brokers (IB) предлагает API Python, который можно использовать для разработки собственного торгового решения.

Платформа Quantopian

Мы представили платформу Quantopian и продемонстрировали использование ее исследовательской и торговой среды для анализа и тестирования стратегий торговли на основе исторических данных. Quantopian использует язык Python и предлагает целый ряд учебных материалов.

На платформе Quantopian проводятся ежедневные соревнования по рекрутированию алгоритмов для своего хедж-фондового портфеля с прямым вовлечением людей через Интернет. Quantopian предоставляет капитал победившему алгоритму. Живая торговля была прекращена в сентябре 2017 г., но платформа по-прежнему предоставляет большой спектр исторических данных и привлекает активное сообщество разработчиков и биржевых торговцев, что является хорошей отправной точкой для обсуждения идей и обучения у других.

Платформа QuantConnect

QuantConnect — это еще одна ведомая сообществом алгоритмическая торговая платформа с открытым исходным кодом, которая конкурирует с Quantopian. Она также предоставляет интерактивную среду разработки (IDE) для бэктестирования и живой торговли с помощью алгоритмических стратегий с использованием Python и других языков.

Платформа QuantConnect также имеет динамичное глобальное сообщество со всего мира и предоставляет доступ к многочисленным классам активов, включая долевые ценные бумаги, фьючерсы, форекс и криптовалюты. Она предлагает интеграцию с живой торговлей через различных брокеров, таких как IB, OANDA и GDAX.

Платформа QuantRocket

QuantRocket — это платформа на основе Python для исследования, бэктестирования и запуска автоматизированных квантитативных стратегий торговли. Она предоставляет инструменты сбора данных, несколько поставщиков данных, исследовательскую среду, несколько бэктестирующих и живую и бумажную торговлю через брокера IB. Она гордится своей поддержкой международной торговли долевыми ценными бумагами и отличается своей гибкостью.

Платформа QuantRocket поддерживает несколько механизмов — собственный Moonshot, а также сторонние механизмы, выбираемые пользователем. Хотя платформа QuantRocket не имеет традиционной интегрированной среды разработки, она хорошо интегрирована с блокнотами Jupyter, производя подобные. Вместе с тем платформа QuantRocket не является бесплатной, и на момент написания этих строк цены начинались с 19 долларов в месяц.

Заключение

Мы начали с того, что отметили взрывной рост цифровых данных и появления МО как стратегического потенциала для стратегий инвестирования и торговли. Эта динамика отражает глобальные деловые и технологические тренды, выходящие за рамки финансирования, и она, скорее всего, сохранится, нежели затормозится или обратится вспять. Многие инвестиционные компании только начинают задействовать ряд инструментов искусственного интеллекта, равно как и люди, которые начинают приобретать соответствующие навыки, в то время как производственные процессы адаптируются к этим новым возможностям для создания стоимости, как описано во вступительной главе.

На горизонте появляются все новые интересные возможности для применения МО к торговле на финансовых рынках, которые, вероятно, еще больше подтолкнут текущий импульс. В ближайшие годы они, скорее всего, останутся актуальными и будут включать автоматизацию процесса МО, генерирование синтетических тренировочных данных и появление квантовых вычислений. Необычайный динамизм в указанной сфере подразумевает, что только это одно может наполнить целую книгу, и путешествие продолжит оставаться захватывающим.

Глоссарий

Baltic Dry Index (BDI) — торговый индекс, ежедневно рассчитываемый Балтийской биржей (Baltic Exchange). Индекс отражает стоимость перевозок сухого груза (уголь, руда, зерно и т. п.) морем по двадцати основным торговым маршрутам.

BATS Global Markets (Better Alternative Trading System) — компания, управляющая тремя частными фондовыми биржами, расположенная в пригороде Канзас-Сити, штат Миссури, США.

Bloomberg L.P. — один из двух ведущих поставщиков финансовой информации для профессиональных участников финансовых рынков.

BSE India Limited (The Bombay Stock Exchange Limited) — Бомбейская фондовая биржа.

BWIC (bid wanted in competition, требуется заявка на конкурентной основе) — ситуация, когда институциональный инвестор представляет свой список заявителей на облигации различным дилерам ценных бумаг, дилеры делают свои заявки на указанные бумаги, и с теми, кто предлагает самые высокие заявки, заключаются контракты. См. <https://www.investopedia.com/terms/b/bwic.asp>. — *Прим. перев.*

Chicago Board Options Exchange (CBOE) — Чикагская биржа опционов.

Deutsche Börse AG — Немецкая биржа, созданная в форме акционерного общества и осуществляющая организацию работы различных рынков ценных бумаг.

Euronext — Панъевропейская фондовая биржа, имеющая филиалы в Бельгии, Франции, Нидерландах и Португалии.

EURO STOXX 50 — фондовый индекс ценных бумаг в еврозоне, разработанный поставщиком индекса STOXX, принадлежащим Deutsche Börse Group.

Hong Kong Exchanges and Clearing Limited — крупный финансовый холдинг, образованный в 2000 г. путем объединения компаний Гонконгской фондовой биржи, Гонконгской фьючерсной биржи (Hong Kong Futures Exchange) и Гонконгской расчетно-клиринговой компании по сделкам с ценными бумагами (Hong Kong Securities Clearing Company).

i.i.d. (independent and identically distributed) — одинаково распределенный и взаимно независимый.

International Securities Exchange (ISE, Международная биржа ценных бумаг) — дочерняя компания американской многонациональной корпорации финансовых услуг Nasdaq, Inc. Является электронной биржей опционов.

Investors Exchange (IEX, биржа инвесторов) — фондовая биржа, базирующаяся в США, основанная в 2012 г. и запущена в качестве национальной биржи ценных бумаг в сентябре 2016 г. 24 октября 2017 г. IEX получила нормативное одобрение SEC на включение компаний в список. IEX зарегистрировала свою первую публичную компанию Interactive Brokers 5 октября 2018 г.

Japan Exchange Group, Inc. (Группа японских бирж) — японская корпорация финансовых услуг, которая управляет несколькими фондовыми биржами, включая Токийскую фондовую биржу и фондовую биржу Осаки.

k -блочный перекрестный контроль (k -fold cross validation) состоит в том, чтобы разбить полный набор данных на k одинаковых по размеру и не накладывающихся друг на друга подмножеств, именуемых блоками (fold), и затем в цикле тренировать модель на всех блоках, кроме i -го блока, и оценивать метрику на i -м блоке.

London Stock Exchange (LSE) Group — Британская фондовая биржа и финансовая информационная компания.

NASDAQ (National Association of Securities Dealers Automated Quotation, Автоматизированные котировки Национальной ассоциации дилеров по ценным бумагам) — американская биржа, специализирующаяся на акциях высокотехнологичных компаний (производство электроники, программного обеспечения и т. п.).

National Stock Exchange of India Limited — Национальная фондовая биржа Индии.

NYSE (New York Stock Exchange) — Нью-Йоркская фондовая биржа, расположенная на Уолл-стрит, главная фондовая биржа США, крупнейшая в мире по обороту.

p -значение (p -value) — при заданной случайной модели, которая воплощает в себе нулевую гипотезу, p -значение является вероятностью получения столь же необычных или предельных результатов, что и наблюдаемые результаты.

PnL (net mark-to-market value of profits and losses) — чистая стоимость прибылей и убытков по текущим рыночным ценам.

R -квадрат (R -squared) — доля дисперсии, объясненная моделью, со значениями в интервале от 0 до 1.

S&P 500 — фондовый индекс, в корзину которого включено 500 избранных акционерных компаний США, имеющих наибольшую капитализацию. Список принадлежит компании Standard & Poor's и ею же составляется. Индекс публикуется с 4 марта 1957 г.

***t*-значение** (*t*-value) — стандартизованная версия проверочной статистики, используемой в качестве критерия при проверке статистической гипотезы.

***t*-распределение** (*t*-distribution) — эталонное распределение (в частности, полученное из нулевой гипотезы), с которым может быть сопоставлена наблюдаемое *t*-значение.

Активы под управлением (*assets under management, AUM*) служат мерой общей рыночной стоимости всех финансовых активов, которыми финансовое учреждение, такое как взаимный фонд, венчурная компания или брокерский дом, управляет от имени своих клиентов и себя.

Акционерный капитал (*shareholders equity* или просто *equity*), или собственный капитал, — остаточное требование владельцев корпорации после погашения долгов. Собственный капитал равен совокупным активам фирмы за вычетом ее совокупных обязательств.

Альтернативная гипотеза (*alternative hypothesis*) — обратная нулевой гипотезе (то, что вы надеетесь доказать).

Альфа (*alpha*) — показатель, рассчитываемый для ценной бумаги или портфеля ценных бумаг, связывающий возвратность ценной бумаги (портфеля) с возвратностью близкого фондового индекса. Высокое значение альфа означает, что ценная бумага (портфель) работает лучше, чем ожидалось при его заданном бета (волатильности). Термин заимствован из статистики, где альфа (уровень значимости) — это вероятностный порог "необычности", который случайные результаты должны превзойти, чтобы фактические результаты считались статистически значимыми. См. <https://ru.wikipedia.org/wiki/Альфа-коэффициент>.

Арбитраж (*arbitrage*) — практика использования преимуществ разницы цен между двумя или более рынками и заключения комбинации совпадающих сделок, которые извлекают выгоду из таких дисбалансов. Например, под валютным арбитражем понимается одновременная покупка и продажа валюты на разных валютных рынках с целью получения прибыли от разницы обменных курсов в двух местах.

См. <https://www.thefreelibrary.com/Currency+arbitrage+as+a+tool+of+corporate+financial+management-a0469641512>.

Асимметрия (*skew*), или скошенность, — состояние, когда один хвост распределения длиннее другого.

Базисный пункт (*basis point, bps* или *bips* во мн. ч.) — единица измерения, используемая в финансах для описания процентного изменения стоимости или ставки финансового инструмента. Один базисный пункт эквивалентен сотой доле процента, 0,01%, или 0,0001 в десятичной форме. Схожим образом дробный базисный пункт, такой как 1,5 базисных пункта, эквивалентен 0,015% или 0,00015 в десятичной форме.

См. <https://www.investopedia.com/ask/answers/what-basis-point-bps/>.

Биномиальное испытание (*binomial trial*) — испытание с двумя результатами.

Биномиальное распределение (binomial distribution) — распределение набора успехов в x испытаниях. Синоним: бернуллиево распределение.

Биржевой фонд (exchange-traded fund, ETF), или торгуемый на бирже фонд, — ликвидная ценная бумага, которая отслеживает фондовый индекс, товар, облигации или корзину финансовых активов.

См. <https://www.investopedia.com/terms/e/etf.asp>.

Бустинг (boosting), или ансамблевое усиление, — общая методика последовательной подгонки серии моделей путем предоставления для каждого последующего цикла большего веса элементам с большими остатками.

Бутстраповская выборка (bootstrap sample), или бутстрап, — подвыборка данных, взятая с возвратом данных назад в исходную совокупность наблюдаемых данных. Процесс бутстрапирования можно концептуально представить, как многократное взятие образцов с возмещением с целью получения синтетической выборки.

Бэггинг (bagging, bootstrap aggregating), или агрегирование бутстраповских выборок, пакетное усреднение, — общая методика формирования набора моделей путем взятия бутстраповских выборок из данных и их усреднения. В английском языке термин "бэггинг" имеет двоякое значение, который, с одной стороны, является аббревиатурой для бутстрап-агрегирования, а с другой, исходя из слова bagging (пакетирование), означает усреднение бутстраповских пакетов данных, формируя упакованный (bagged, бутстрап-агрегированный) предсказатель. Отсюда возникли термины "внутрипакетный" (in-bag) и "внепакетный" (out-of-bag). Используемые для построения предсказателя выборки называются внутрипакетными, а выборки, которые для этого не использовались, — внепакетными.

Бэктестирование (backtesting), или ретроактивное тестирование, ретротестирование, — тестирование на исторических данных с целью получения результатов и анализа риска и возвратности инвестиций, прежде чем рисковать любым фактическим капиталом.

Векторное вложение слова (word embedding), или вложение слова в векторное пространство, — векторное представление слова как части распределенного представления текста в n -мерном пространстве, идея которого состоит в том, чтобы создать плотный вектор для каждого слова, выбранный таким образом, чтобы он был похож на векторы слов, которые появляются в схожем контексте.

См. <https://hackernoon.com/word-embeddings-in-nlp-and-its-applications-fab15eaf7430>, а также <https://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf>.

Взаимный фонд (mutual fund), или фонд взаимных инвестиций, — финансовый механизм, состоящий из денежных средств, собранных от многочисленных инвесторов с целью инвестирования в ценные бумаги, такие как акции, облигации, инструменты денежного рынка и другие финансовые активы. Представляет собой портфель активов, тщательно отобранных и приобретенных профессиональными финансистами на вложения многих тысяч мелких вкладчиков.

См. <https://www.investopedia.com/terms/m/mutualfund.asp>.

Винзоризация (winsorizing) — преобразование статистического показателя путем лимитирования экстремальных значений в статистических данных с целью уменьшения эффекта возможно мнимых выбросов. Данный метод назван в честь инженера-биостатистика Чарльза П. Уинзора (1895–1951). Эффект тот же, что и при амплитудном ограничении в обработке сигналов.

Влияние на финансовый рынок (market impact) означает влияние, которое участник рынка оказывает на цену финансового актива, когда он покупает или продает его.

Вмешательство учителя (teacher forcing) — это метод ускоренной и эффективной тренировки RNN-сети, когда в качестве входа в следующий шаг вместо выхода модели используется непосредственное наблюдение из предыдущего шага.

См. обзорную статью по указанному методу по адресу:

<https://machinelearningmastery.com/teacher-forcing-for-recurrent-neural-networks/>.

Внебиржевой (over-the-counter, OTC) относится к процессу торговли ценными бумагами для компаний, которые не котируются на официальной бирже, такой как Нью-Йоркская фондовая биржа (NYSE). Ценные бумаги, которые торгуются вне биржи, торгуются не на централизованной бирже, а через дилерскую сеть. Брокерско-дилерские компании торгуют ценными бумагами вне биржи, ведя переговоры непосредственно друг с другом через компьютерные сети и по телефону. Дилеры действуют как маркет-мейкеры, используя доску объявлений OTC, представляющую собой междилерскую котировочную систему, которая предоставляет торговую информацию.

См. <https://www.investopedia.com/terms/o/otc.asp>.

Внутри выборки и вне выборки (in-sample и out-of-sample) — тестирование обычно проводится путем разбиения набора данных на внутривыборочный (in-sample) период, используемый для первоначального оценивания параметров и отбора модели, и вневыборочный (out-of-sample) период, используемый для оценивания результативности предсказания. В частности, если имеются данные, скажем, за 3 года, необходимые для расчета волатильности, то модель, используемая в течение этого периода, будет "внутри выборки". Но если использовать исторические данные для предсказания на перспективу, то оценивание будет выполняться за период времени, для которых нет данных (вне выборки). Таким образом, обычно понятие "вне выборки" относится к предсказанию там, где у нас нет данных. Технически, даже использование модели для оценивания сегодняшней волатильности на основе исторической выборки является предсказанием вне выборки, потому что у нас нет сиюминутной волатильности.

См. [https://ec.europa.eu/eurostat/statistics-explained/index.php/Glossary:](https://ec.europa.eu/eurostat/statistics-explained/index.php/Glossary:In-sample_vs._out-of-sample_forecasts)

[In-sample_vs._out-of-sample_forecasts](https://ec.europa.eu/eurostat/statistics-explained/index.php/Glossary:In-sample_vs._out-of-sample_forecasts).

Внутрирядовая корреляция (serial correlation), или автокорреляция, — взаимосвязь между наблюдениями одной и той же переменной в течение определенных периодов времени. Если величина внутрирядовой корреляции переменной равна

нулю, то это означает, что корреляция отсутствует, и все наблюдения друг от друга не зависят.

См. <https://www.investopedia.com/terms/s/serial-correlation.asp>.

Выборка (sample) — подмножество большей по размеру совокупности данных либо сам процесс отбора образцов из нее.

Выезд на кривой отдачи (riding the yield curve) — торговая стратегия, которая включает в себя покупку долгосрочной облигации и ее продажу до ее созревания с целью получения прибыли от снижения доходности, которое происходит в течение срока действия облигации.

См. <https://www.investopedia.com/terms/r/ridingtheyieldcurve.asp>.

Выработка признаков (feature engineering), или конструирование признаков, — процесс извлечения признаков из сырых данных и преобразования их в форматы, подходящие для автоматически обучающейся модели, с использованием компетентных знаний предметной области данных.

Гетероскедастичность (heteroskedasticity) — ситуация, когда некоторые диапазоны выходного результата показывают остатки с более высокой дисперсией (что может говорить о предсказательной переменной, которая в уравнении отсутствует).

Гиперпараметры (hyperparameters) — параметры, которые необходимо установить перед подгонкой алгоритма.

Главная компонента (principal component) — линейная комбинация предсказательных переменных.

Глубина рынка (market depth) — способность рынка конкретного актива поддерживать крупные ордера на данный актив без существенного движения цены актива. Чем больше открытых лимитных ордеров на обеих сторонах книги ордеров биржи, тем больше глубина этой книги. Глубина также тесно связана с ликвидностью: чем глубже книга ордеров для актива, тем больше ликвидности книга ордеров предоставляет этому активу.

Градиентный бустинг (gradient boosting) — более общая форма бустинга, которая создана с точки зрения минимизации функции стоимости.

Гэп (gap), или скачок, — разрыв в цене, в частности между ценой старого контракта и ценой нового контракта.

Дефицит реализации (implementation shortfall) — разница между ценой решения и окончательной ценой исполнения (включая комиссии, налоги и т. д.) для торговой сделки. Он также известен под названием "проскальзывания". Агентская торговля в значительной степени связана с минимизацией дефицита реализации и поиском ликвидности.

Дисперсия (variance) — сумма квадратов отклонений от среднего, деленная на $n - 1$, где n — число значений данных. Синоним: среднеквадратическое отклонение.

Длинная позиция (long), или лонг, — покупка ценной бумаги в ожидании, что она вырастет в цене. Разница между ценой продажи и ценой покупки приносит при-

быль или убыток. Инвесторы входят в длинную позицию (покупают, лонгуют) в ожидании, что ценные бумаги возрастут в цене, и входят в короткую позицию (продают, шортят) в ожидании, что ценные бумаги упадут в цене.

Длинно-короткая инвестиционная стратегия (long-short strategy) — это инвестиционная стратегия, обычно связанная с хедж-фондами и с некоторыми прогрессивными традиционными менеджерами активами, которая заключается в покупке длинных акций, т. е. с ожидаемым ростом цены, и продаже коротких акций, т. е. с ожидаемым снижением цены. Прилагательные "длинный" и "короткий" обозначают направление движения цены соответственно вверх и вниз, не обязательно обозначая продолжительность.

Долевая ценная бумага (equity) — это ценная бумага, обеспечивающая право собственности в компании через покупку обыкновенных либо привилегированных акций (долей). Термин equity также означает нетто-стоимость или чистую стоимость компании.

См. http://www.morningstar.com/InvGlossary/equity_definition_what_is.aspx.

Долларовый финансовый возврат (dollar return) — возврат на портфельные инвестиции в течение любого периода оценивания, включающего изменение рыночной стоимости портфеля. Он также включает любые распределения, накопленные из портфеля в течение этого периода.

См. http://www.investorwords.com/16372/dollar_return.html#ixzz5ZB1R24cm.

Заполнение (fill) — действие по исполнению или удовлетворению ордера на ценную бумагу или биржевой товар. Это основной ордер при сделках с акциями, облигациями или любым другим видом ценных бумаг. Аннулирование (kill) — это запрос на отмену сделки между ее размещением и ее исполнением.

Защитный опцион put (protective put) — стратегия риск-менеджмента, которую инвесторы используют для защиты от потери нерезализованной прибыли в акциях или других активах. Опцион put действует как страховой полис (он стоит денег), который уменьшает потенциальную прибыль инвестора от владения ценной бумагой, но и сокращает риск потери денег, если стоимость ценной бумаги снижается. См. <https://www.investopedia.com/terms/p/protective-put.asp>.

Идиосинкратический риск (idiosyncratic risk) — риск, который относится к факторам, влияющим на отдельный актив, по сравнению с систематическим риском, который относится к более широким трендам, сказывающимся на всем рынке в целом. См. <https://www.investopedia.com/search?q=idiosyncratic+asset>.

Инвестиционная стратегия (investment strategy) — систематический план размещения инвестируемых активов среди инвестиционных вариантов, таких как облигации, депозитные сертификаты, товары, недвижимость, акции. Эти планы учитывают такие факторы, как экономические тенденции, инфляция и процентные ставки. Другие факторы включают возраст инвестора, уровень толерантности к риску и краткосрочные или долгосрочные цели роста. Корпоративные инвестиционные стратегии определяют средства, необходимые для достижения

конкурентного преимущества, и денежные результаты (прибыль), ожидаемые от таких решений.

См. <http://www.businessdictionary.com/definition/investment-strategy.html>.

Индикатор чистого дисбаланса ордеров (Net Order Imbalance Indicator, NOII) — информация о кроссах на открытии и закрытии на фондовом рынке NASDAQ, предоставляемая пользователям рынка перед исполнением кроссов. Индикатор NOII показывает истинное предложение и спрос на акции на основе фактических ордеров на покупку и продажу за 10 минут до закрытия рынка и за 5 минут до открытия рынка. См. <https://www.investopedia.com/terms/n/net-order-imbalance-indicator-noii.asp>.

Информационный коэффициент (information coefficient) описывает корреляцию между прогнозируемой и фактической возвратностью активов, варьируясь в диапазоне -1 до 1 ; иногда используется для измерения вклада финансового аналитика. См. <https://www.investopedia.com/terms/i/information-coefficient.asp>.

Информационное соотношение (information ratio) — показатель возвратности портфеля, выходящий за пределы возвратности эталонного показателя, обычно индекса, по сравнению с волатильностью этой возвратности; вычисляется по формуле $(R_p - R_b) / \text{ошибка прослеживания}$, где R_p — это возвратность инвестиционного портфеля, R_b — эталонная возвратность, ошибка прослеживания — стандартное отклонение избыточной возвратности по отношению к эталонной возвратности. IR пытается определить стабильность результативности портфеля путем включения в расчет риска отклонения или компоненты стандартного отклонения.

См. <https://www.investopedia.com/terms/i/informationratio.asp>.

Квант (quant) — биржевой специалист по квантитативным методам.

См. <https://habr.com/company/itcapital/blog/307854/>.

Квантиль (quantile) — такое значение, что p процент значений принимает данное значение или меньшее и $(100 - p)$ процент значений принимает данное значение или большее.

Квантильный набросок (quantiles sketch) — потоковый алгоритм оценивания распределения значений, который приближенно отвечает на запросы о ранге значения, функции массы вероятности распределения (PMF) или гистограмме, кумулятивной функции распределения (CDF) и квантилях (медиана, минимум, максимум, 95-й процентиль и т. д.).

См. <https://datasketches.github.io/docs/Quantiles/QuantilesOverview.html>.

Квантоментальный (quantamental) — неологизм, который образован из двух терминов "квантитативный подход" и "фундаментальный подход" и означает сочетание новейших квантитативных методов, в том числе на основе машинного обучения и ИИ, с классическими методами на основе фундаментальных величин.

Корпоративные заработки (earnings) обычно относятся к чистому доходу после уплаты налогов, иногда именуемому нижней строкой, или прибылями компании (по формуле: валовая выручка *минус* стоимость реализованной продукции/услуг = чистая выручка *плюс* другие доходы *минус* операционные расходы = валовая прибыль (EBITDA, корпоративный заработок) *минус* налоги, проценты, износ и амортизация = чистая прибыль (нераспределенная прибыль, это *тоже* корпоративный заработок). Заработки являются основным определяющим фактором долевой цены компании, поскольку заработки и связанные с ними обстоятельства могут указывать на прибыльность и успешность предприятия в долгосрочной перспективе. Заработки являются, пожалуй, самым важным и наиболее изученным показателем в финансовой отчетности компании, поскольку они показывают прибыльность по сравнению с оценками со стороны аналитиков и руководства компании. См. <https://www.investopedia.com/terms/e/earnings.asp>.

Кэрри-трейд (carry-trade, carry) — стратегия, в соответствии с которой высокодоходная валюта финансирует сделку с низкодоходной валютой. Иными словами, это получение прибыли на валютном рынке за счет разной величины процентных ставок. См. <https://www.investopedia.com/terms/c/currencycarrytrade.asp>.

Ковариация (covariance), или совместная дисперсия, — метрический показатель, который показывает степень, с которой переменная варьируется совместно с другой (т. е. имеет аналогичную величину и направление).

Компромисс между риском и возвратностью (risk-return trade-off) гласит, что потенциальная возвратность возрастает с увеличением риска. Используя этот принцип, люди связывают низкий уровень неопределенности с низкой потенциальной возвратностью, а высокий уровень неопределенности или риска с высокой потенциальной возвратностью. Согласно компромиссу "риск — возвратность" вложенные деньги могут давать более высокую возвратность только в том случае, если инвестор примет более высокую вероятность убытков. См. <https://www.investopedia.com/terms/r/riskreturntradeoff.asp>.

Контрастивное к шуму оценивание (Noise-contrastive estimation, NCE) сводит задачу оценивания языковой модели к задаче оценивания параметров вероятностного бинарного классификатора, использующего те же параметры для различения выборок из эмпирического распределения от выборок, генерируемых распределением шума. См. <https://arxiv.org/pdf/1410.8251.pdf>.

Контрольный набор (validation set) — набор образцов, используемый для регуляции/оптимизации параметров модели и отбора окончательной модели.

Короткая позиция (short), или шорт — продажа ценной бумаги в ожидании, что она упадет в цене. Короткая позиция сначала продается, а затем покупается по более низкой цене. Разница между ценой продажи и ценой покупки приносит прибыль или убыток. Инвесторы входят в длинную позицию (покупают, лонгуют) в ожидании, что ценные бумаги возрастут в цене, и входят в короткую позицию (продают, шортят) в ожидании, что ценные бумаги упадут в цене.

Котировка (quote) — последняя цена, по которой торгуется биржевой актив — ценная бумага или товар, при этом имеется в виду самая недавняя цена, с которой согласились покупатель и продавец и по которой была проведена сделка с определенной суммой актива.

См. <http://www.businessdictionary.com/definition/quote.html>.

Коэффициент кофенетической корреляции (cophenetic correlation coefficient) — мера того, насколько точно дендрограмма сохраняет попарные расстояния между исходными немоделированными точками данных.

Кредитное плечо (leverage), или финансовое плечо, — дополнительная покупательная способность, создаваемая маржинальной торговлей, позволяющая эффективно платить меньше полной цены за актив, используя заемные средства. Кредитное плечо обычно представлено в виде соотношения: например, если на торговом счете имеется 10 тыс. долларов, и вы занимаете еще 10 тыс. долларов, то ваше кредитное плечо равно 2:1.

Кросс (cross) — данные, которые биржа NASDAQ собирает и публикует по всем заинтересованностям на покупку и продажу за 2 минуты до своего открытия или закрытия и называются кроссом на открытии (opening cross) или кроссом на закрытии (closing cross). Такое распространение информации о заинтересованности в цене помогает ограничить сбои в ликвидности на открытии биржи. Кросс на закрытии биржи сопоставляет запросы и предложения по данной акции с целью создать окончательную цену дня. См. <https://www.investopedia.com/terms/n/net-order-imbalance-indicator-noii.asp>.

Лимитный ордер — поручение клиента брокеру исполнить заявку в заданном интервале цен.

Логит (logit), или логит-преобразование, — функция, которая увязывает вероятность принадлежности классу с интервалом $(-\infty; +\infty)$ (вместо интервала от 0 до 1). Иногда называется ненормализованным логарифмом вероятности.

Лонговать (go long) — покупать активы, открывать длинную позицию на покупку в надежде на рост цены актива, причем это вовсе не означает, что позиция будет удерживаться продолжительное время. Торговец "лонгует", когда цена растет, давая ему возможность получать прибыль от покупки.

Лямбда (lambda) — интенсивность (в расчете на единицу времени или пространства), с которой события происходят.

Малостоимостные акции (small-value stock, small stock) — акции с малой рыночной капитализацией, а также акции, которые торгуются по балансовой стоимости или ниже ее. Найти акции, соответствующие обоим этим критериям, очень трудно, но может оказаться полезным предприятием, поскольку малостоимостные акции обычно приносят высокую возвратность.

См. https://www.investopedia.com/terms/s/small_value_stock.asp.

Маргинализация (marginalizing out), или суммирование всех возможных конфигураций, состоит в том, чтобы в совместном распределении вероятностей находить

распределение некой случайной величины путем исчерпания случаев ее появления вместе с другими случайными величинами.

Маргинальное распределение (marginal distribution), или частное распределение — распределение вероятностей подмножества совокупности случайных величин. Например, при заданных двух случайных величинах X и Y , совместное распределение которых известно, частное распределение X — это простое распределение вероятности X , усредненное по информации об Y .

См. [http://cyclowiki.org/wiki/Маргинальное распределение](http://cyclowiki.org/wiki/Маргинальное_распределение).

Маржинколл (margin call) — сигнал для уведомления инвестора о приближении к критическому уровню и не обязывает вносить средства для поддержания кредита и покрытия дальнейших возможных убытков. В отличие от него, стопаут (stop out) — это более низкий чем "маржинколл" уровень, помогающий избежать потери финансовых средств, предоставленных через кредитное плечо. В этом случае, если размер убытков продолжает увеличиваться, то сделка, или даже несколько открытых позиций, закрывается принудительно, причем основные убытки в этом случае несет только инвестор.

Маркетмейкер (market maker) — брокер, дилер или инвестиционная компания, которые принимают на себя рыночный риск (системный риск), вступая во владение финансовым активом и торгуя им в качестве принципала. Проще говоря, это организации, которые предоставляют ликвидность бирже путем размещения лимитных ордеров в ее книге заказов для того, чтобы сделки могли совершаться в диапазоне цен. Маркетмейкеры обязаны постоянно выдавать котировки цен спроса и предложения, а также гарантировать полную продажу или поглощение финансового актива по определенной цене.

См. <http://www.businessdictionary.com/definition/market-maker.html>.

Мартингейл (martingale) — тип инвестиционной стратегии, используемой торговцами для капитализации убытков. По мере того как цены на акции снижаются, инвестор покупает больше инвестиций, чтобы расширить свой портфель по более низкой цене. Инвестор покупает больше акций по более низкой цене с верой в то, что цена в конечном итоге увеличится и принесет чистую прибыль. Например, инвестор может приобрести акции Google по цене 500 долларов за акцию, а затем развернуться и купить больше акций, если она упадет до 450 долларов за акцию. Когда цена восстанавливается, инвестор получает большой финансовый возврат на инвестиции. См. http://www.investorwords.com/8639/martingale_system.html#ixzz5ZR4W4YYN.

Матрица несоответствий (confusion matrix), или ошибок, — отображение в табличной форме (2×2 в бинарном случае) количеств записей по их предсказанному и фактическому состоянию, или результату, классификации.

Матрица плана регрессионной модели (design matrix), также модельная матрица, матрица регрессоров, обычно обозначаемая как X , — это матрица значений так называемых объясняющих переменных множества объектов. Каждая строка представляет собой отдельный объект, а последующие столбцы соответствуют

переменным и их конкретным значениям для этого объекта. Преимущество матрицы плана заключается в том, что она может представлять множество видов экспериментальных проектов и статистических моделей.

См. https://en.wikipedia.org/wiki/Design_matrix.

Медиана (median) — такое значение, что половина отсортированных данных находится выше и ниже его.

Ментальный учет (mental accounting) — экономическая концепция, созданная экономистом Ричардом Талером (Richard H. Thaler, теория рациональных расчетов), которая утверждает, что люди классифицируют личные финансовые средства по-разному и поэтому склонны к иррациональному принятию решений в своем поведении по расходованию и инвестированию финансовых средств. См. <https://www.investopedia.com/terms/m/mentalaccounting.asp>.

Многонитиевая обработка (multithreading) — способность CPU (или единственного ядра в многоядерном процессоре) обеспечивать несколько нитей исполнения одновременно. В русской терминологии за термином *thread* укрепился неточный перевод "поток", хотя его протогерманский предок *thread* по форме и смыслу является родственником нашим "прясти", "прядь", т. е. процесс сплетается из нитей (ср. нить рассуждения, если брать антропоморфную аналогию). В зарубежной технической литературе данное понятие объясняется именно на нитях, которые рассматриваются как основные инфраструктурные единицы планирования в операционных системах.

Среднезначная модель (mean model) — это модель временного ряда, в которой предсказание будущих значений строится на основе выборочного среднего прошлых данных; она исходит из того, что будущие наблюдения будут взяты из того же самого распределения вероятностей.

Модель ценообразования капитальных активов (Capital Asset Pricing Model, CAPM) описывает связь между систематическим риском и ожидаемой доходностью от активов, в особенности от акций. Указанная модель является однофакторной, и единственным фактором служит бета, который согласно модели показывает, насколько хорошо акция движется относительно рынка. Акции, которые двигались выше рынка, имели более высокий бета и, следовательно, более высокий риск и возврат. См. <https://www.investopedia.com/terms/c/capm.asp>.

Мультиколлинеарность (multicollinearity) — ситуация, при которой предикторные переменные имеют идеальную или почти идеальную корреляцию, а регрессия может быть нестабильной, либо ее невозможно вычислить.

Наклон (slope) — угловой коэффициент, определяющий наклон линии регрессии по отношению к оси x , численно равен тангенсу угла между положительным направлением оси x и данной прямой.

Накопления (accruals) — заработанные доходы и понесенные расходы, которые оказывают общее влияние на отчет о прибылях и убытках.

Нулевая гипотеза (null hypothesis) — общее утверждение или принимаемая по умолчанию позиция относительно того, что между двумя наблюдаемыми событиями взаимосвязи не существует, или что причиной всему является случайность.

Обрезка дерева решения (decision tree pruning) с оптимизацией совместной метрики стоимости-сложности призвана снижать сложность и стоимость дерева решения. Данный метод подрезает деревья решений, стремясь минимизировать метрику, которая сочетает в себе два фактора: размер (сложность) и точность (или, что эквивалентно, частоту ошибок) дерева. См. <https://pdfs.semanticscholar.org/d5df/17fc1534094da809034e404f66d057717e14.pdf>.

Опционы call — это соглашения, которые дают покупателю опциона право, но не обязанность, купить акции, облигации, товар или другой инструмент по оговоренной цене в течение определенного периода времени. Акции, облигации или товары называются базовым активом. Покупатель опциона call получает прибыль, когда базовый актив увеличивается в цене.

См. <https://www.investopedia.com/terms/c/calloption.asp>.

Опцион put (put option) — опционный контракт, дающий владельцу право, но не обязательство, продать определенное количество базового актива по оговоренной цене в течение определенного периода времени. Поскольку базовый актив обесценивается по стоимости, сам опцион put растет в цене; поэтому покупка опционов put является методом его шортирования (занятия короткой позиции). Является противоположностью опциона call, который дает держателю право купить базовую ценную бумагу по указанной цене до истечения срока действия опциона. См. <https://www.investopedia.com/terms/p/putoption.asp>.

Ордер (order), или заявка, — подтвержденная заявка одной стороны другой стороне на покупку, продажу, доставку или получение товаров или услуг в соответствии с указанными условиями. Когда ордер принимается принимающей стороной, он становится юридически обязательным договором.

См. <http://www.businessdictionary.com/definition/order.html>.

Ордерная книга (order book) — это электронный список заявок на покупку и продажу конкретной ценной бумаги или финансового инструмента, организованный по уровню цен. Она имеет и другие названия — стакан заявок, книга приказов или глубина рынка. См. <https://www.investopedia.com/terms/o/order-book.asp>.

Ордер не по установившемуся рынку (not-held order) — рыночный или лимитный ордер, в котором клиент не желает автоматически совершать сделки на внутреннем рынке (установившемся рынке, market held) и дает брокеру или торговцу время и свободу выбора цены с целью получить наилучшую возможную цену.

Относительный финансовый возврат (relative return) — соотношение между финансовым возвратом, реализованном на фонде или активе, и финансовым возвратом, реализованном заданным эталоном. Относительные финансовые возвраты часто используются для оценивания менеджеров фондов, от которых ожида-

ется, что они достигнут финансовых возвратов, превышающих движения индексов неуправляемых активов.

См. http://www.investorwords.com/15640/relative_return.html.

Оценщик (estimator) — правило (формула) для расчета оценки заданной величины на основе наблюдаемых данных. Различаются правило (оценщик), интересующая величина (эстиманд) и результат применения правила (оценка).

См. <https://en.wikipedia.org/wiki/Estimator>.

Ошибка (error) — разница между точкой данных и предсказанным либо средним значением.

Ошибка выборки (sampling error) — величина погрешности, которая возникает при оценивании значения популяционной переменной на основе значения этой переменной в выборке, взятой из этой популяции.

Ошибка прослеживания (tracking error), иногда риск отклонения, — стандартное отклонение возвратности портфеля от эталонной возвратности за определенный период времени. См. <https://financetrain.com/tracking-error-tracking-risk/>.

Панельные данные (panel data), или продольные данные, — данные, агрегированные по разным переменным за многочисленные периоды времени, например, в формате "страна — ВВП 2010, ВВП 2011, ..., ВВП 2018".

Паритет риска (risk parity) — стратегия перераспределения портфеля с использованием риска для определения размещений по различным компонентам инвестиционного портфеля. Подход к управлению портфелем на основе паритета риска основан на теории современного инвестиционного портфеля, которая стремится оптимально диверсифицировать инвестиционный портфель среди указанных активов. См. <https://www.investopedia.com/terms/r/risk-parity.asp>.

Паритет цен put и call (put-call parity) — связь между ценой колл и ценой пут для опциона с одинаковыми характеристиками (цена исполнения, дата истечения срока действия, базовый). Она используется в теории арбитражного ценообразования. Если разные портфели, состоящие из колл и пут, имеют одинаковое значение во время истечения срока действия, то из этого следует, что они будут иметь одинаковое значение вплоть до истечения срока действия. Таким образом, значения портфелей перемещаются строго в параллельно-шаговом режиме. См. http://www.investorwords.com/6631/put_call_parity.html#ixzz5ZADB6MWrr.

Первичное публичное размещение (initial public offering, IPO) — самое первое предложение акций частной корпорации общественности.

Перенос вперед (roll forward), или накат, — продление срока действия фьючерсного контракта путем закрытия первоначального контракта и открытия нового контракта с более долгим сроком на тот же базовый актив по текущей рыночной цене. Позволяет торговцу поддерживать позицию после первоначального истечения контракта, так как фьючерсные контракты имеют конечные даты истечения. Обычно проводится незадолго до истечения первоначального контракта и

требует урегулирования прибыли или убытка по первоначальному контракту. См. <https://www.investopedia.com/terms/r/rollforward.asp#ixzz5YCTu343D>.

Перенос назад (roll backward), или откат, — укорочение срока действия фьючерсного контракта путем выхода из одной позиции и входа в новую позицию с более близким сроком действия. Обычно проводится незадолго до истечения первоначального контракта и требует урегулирования прибыли или убытка по первоначальному контракту.

См. <https://www.investopedia.com/terms/r/rollforward.asp#ixzz5YCTu343D>.

Пересечение (intercept) — коэффициент сдвига регрессионной модели, который показывает, каким будет \hat{y} в случае, если все используемые в модели факторы будут равны 0; под ним подразумевается зависимость от других неописанных в модели факторов. Представляет собой смещение по оси y и соответствует точке на этой оси, где прямая регрессии пересекает эту ось.

Повторный отбор (resampling) — процесс многократного взятия выборок из наблюдаемых данных.

Подстановочный оценщик энтропии (plug-in entropy estimator) — формулировка энтропии распределения, где вероятности символов или блоков были заменены их относительными частотами в выборке. См. <https://arxiv.org/abs/1601.06014>.

Позиционер (position-taker), или покупатель позиций, — физическое или юридическое лицо, которое должно принимать преобладающие позиции на рынке, не располагая долей рынка, чтобы влиять на рыночную позицию самостоятельно.

Позиция (position) — сумма ценной бумаги, товара или валюты, которыми владеет физическое лицо, дилер, учреждение или другая налогооблагаемая сущность. Они бывают двух типов: короткие позиции, которые заимствуются, а затем продаются, и длинные позиции, которыми владеют, а которые затем продаются. В зависимости от рыночных тенденций, движений и колебаний, позиция может быть прибыльной или убыточной. Пересчет стоимости позиции для отражения ее фактической текущей стоимости на открытом рынке в отрасли называется маркировкой по рынку. См. <https://www.investopedia.com/terms/p/position.asp>.

Покупатель цен (price-taker) — физическое или юридическое лицо, которое должно принимать преобладающие цены на рынке, не располагая долей рынка, чтобы влиять на рыночную цену самостоятельно. В большинстве конкурентных рынков фирмы являются покупателями цен. Если фирмы устанавливают на свою продукцию более высокие цены, чем преобладающие рыночные цены, то потребители просто покупают ее у другого продавца по более низкой цене. На фондовом рынке индивидуальные инвесторы считаются покупателями цен, а маркет-мейкеры — те, кто устанавливает цену на ценную бумагу и предлагает его на рынке.

См. <https://www.investopedia.com/terms/p/pricetaker.asp#ixzz5YmH5MrO8>.

Полураспад (half-life) в физике — это время, затрачиваемое определенным количеством вещества на распад до половины его массы. Указанное понятие связано

с измерением скорости этого процесса. Применительно к данной теме полураспад показывает медленность процесса либо время достижения ожидаемого значения. См. <http://marcoagd.usuarios.rdc.puc-rio.br/half-life.html>.

Портфель (portfolio), как правило, представляет собой группу финансовых активов, таких как акции, облигации, товары, валюты и эквиваленты денежных средств, а также их фондовые аналоги, включая взаимные, биржевые и закрытые фонды. См. <https://www.investopedia.com/terms/p/portfolio.asp>.

Портфельное страхование с постоянной пропорцией (Constant Proportion Portfolio Insurance, CPPI) — тип портфельного страхования, в котором инвестор устанавливает нижний уровень долларовой стоимости своего портфеля, а затем структурирует размещение активов вокруг этого решения. В CPPI используется два класса активов: рискованный актив (обычно акционерные капиталы или взаимные фонды) и консервативный актив, состоящий из денежных средств, их эквивалентов или казначейских облигаций. Процент, выделенный каждому из них, зависит от значения "подушки", определенной как текущая стоимость портфеля минус значение нижнего уровня, и мультипликативного коэффициента, где большее число обозначает более агрессивную стратегию. См. <https://www.investopedia.com/terms/c/cppi.asp>.

Правдоподобие (likelihood), точнее, функция правдоподобия — это совместное распределение вероятностей наблюдаемых данных, выраженное как функция статистических параметров. Оно описывает относительную вероятность или шансы получения наблюдаемых данных для всех допустимых значений параметров и используется для идентификации конкретных значений параметров, наиболее правдоподобных с учетом наблюдаемых данных. См. https://en.wikipedia.org/wiki/Likelihood_function.

Признак (feature) — индивидуальное измеряемое свойство или характеристика наблюдаемого явления. В матричной форме признак обычно представляется столбцом матрицы. Строки матрицы образуют наблюдения (также именуемые образцами, экземплярами или примерами). В данном контексте выборками, как правило, называется подмножество строк.

Прецизионность (precision), или точность результатов измерений, — мера статистической изменчивости, которая описывает случайные ошибки или степень близости друг к другу независимых результатов измерений, полученных в конкретных установленных условиях. См. https://en.wikipedia.org/wiki/Accuracy_and_precision.

Проверочный статистический показатель (test statistic) — метрический показатель целевой разницы или эффекта, используемый в качестве критерия при проверке статистической гипотезы.

Проскальзывание (slippage) — разница в цене, по которой брокер получает указание от принципала выполнить ордер, и цене, по которой ордер фактически исполнен. См. <http://www.businessdictionary.com/definition/slippage.html>.

Протокол FIX (financial information exchange, FIX) — протокол обмена сообщениями, разработанный для обмена актуальной информацией по сделкам. Находится в свободном доступе и использовании (<http://www.fixprotocol.org>).

Прочесывание данных (data snooping) — подробная ревизия данных с целью найти что-то интересное.

Разведывание (exploration) — это поиск новых идей или новых стратегий. Применяется как антитеза для эксплуатации (exploitation), т. е. использованию существующих идей и стратегий, которые оказались успешными в прошлом. Интеллектуальная работа включает в себя надлежащий баланс между разведыванием, сопряженным с риском пустой траты ресурсов и эксплуатацией проверенных временем решений.

Разворот (reversal) — изменение направления тренда, то есть цена, вероятно, будет продолжать двигаться в этом разворотном направлении в течение продолжительного периода времени.

Размещение финансовых средств среди портфельных активов (asset allocation) — инвестиционная стратегия, призванная балансировать риск и вознаграждение путем выделения и перераспределения средств по портфельным активам в соответствии с целями, рискоустойчивостью и инвестиционным горизонтом человека. Три основных класса активов — долевые ценные бумаги (акции), активы с фиксированным доходом, денежные средства и их эквиваленты — имеют разные уровни риска и возвратности, поэтому каждый из них будет вести себя по-разному с течением времени.

См. <https://www.investopedia.com/terms/a/assetallocation.asp>.

Расхождение Кульбака-Лейблера (Kullback-Leibler Divergence), или информационное расхождение, — неотрицательнозначный функционал, являющийся несимметричной мерой удаленности друг от друга двух вероятностных распределений, определенных на общем пространстве элементарных событий.

Решеточный поиск (grid search) предусматривает, что для каждого гиперпараметра заранее подбирается список значений, которые могут оказаться для него хорошими, затем пишется вложенный цикл for, который пробует все комбинации этих значений с целью найти их контрольные точности, и отслеживает те, которые показывают наилучшую результативность.

Рисковая премия (risk premium) — финансовый возврат, превышающий безрисковый финансовый возврат, ожидаемый от инвестиций.

См. <https://www.investopedia.com/search?q=Risk+premium>.

Рисковый фактор (risk factor) — измеримая характеристика или элемент, изменение которого может повлиять на стоимость актива, например, обменный курс, процентная ставка и рыночная цена.

См. <http://www.businessdictionary.com/definition/risk-factor.html>.

Рынок долевых ценных бумаг (equity market), или фондовый рынок, — это рынок, где выпускаются в обращение и торгуются акции, то есть доли собственности в компании, через биржу или внебиржевые рынки.

См. <https://www.investopedia.com/terms/s/shareholdersequity.asp>.

Рыночный ордер (market order) — поручение клиента брокеру немедленно купить или продать товар по текущей лучшей цене. Наиболее распространенный способ исполнения ордеров.

Рыночный портфель (market portfolio) — теоретическая диверсифицированная группа инвестиций, каждый актив которой взвешен пропорционально его суммарному присутствию на рынке.

См. <https://www.investopedia.com/terms/m/market-portfolio.asp>.

Рыночный режим (market regime) — расширенный период, когда специфическая комбинация активов доминирует над ценовой динамикой целевого инструмента.

Сведение на основе максимума (max pooling) — операция в CNN-сети, когда берется максимальное значение участка/окна окрестных признаков с целью снижения числа активаций; дает инвариантность сдвига в изображении.

Скачки цен спроса-предложения (bid-ask bounce) относятся к конкретной ситуации, когда цена акции или другого актива быстро скачет туда-сюда в очень ограниченном интервале между ценой bid и ценой ask.

См. <https://www.investopedia.com/ask/answers/013015/whats-difference-between-bidask-spread-and-bidask-bounce.asp>.

Скользкая возвратность (roll return, rolling return) представляет собой среднегодовую возвратность за период, заканчивающийся указанным годом.

См. <https://www.investopedia.com/terms/r/rollingreturns.asp>.

Случайный отбор (random sampling) — взятие элементов в выборку в произвольном порядке.

Систематическое смещение из-за выживших (survivorship bias) — ситуация, при которой в качестве инвестиционного универсума используется текущий универсум и, следовательно, игнорируется тот факт, что некоторые компании могли обанкротиться, и ценные бумаги в результате были исключены из списка.

Систематическое смещение из-за забегания вперед (look-ahead bias) — систематическая ошибка, возникающая из-за использования информации, которая не была публичной в момент принятия симулируемого решения.

Систематическое смещение из-за репрезентативности (representativeness bias), или систематическая ошибка из-за стереотипизации, означает, что менеджеры оценивают вероятность события на основе его близости с другими событиями. Горизонтальное смещение означает, что люди склонны классифицировать предмет по его аналогам и прогнозировать его в будущем в соответствии с его сходством. Вертикальное смещение подразумевает, что на финансовых рынках люди легко склонны судить или прогнозировать акции в соответствии с их собственной предысторией. См. <https://www.hindawi.com/journals/mpe/2014/686201/>.

Систематическое смещение при отборе образцов (selection bias) — систематическая ошибка, выражающаяся в появлении у изучаемой выборки признаков, не свойственных генеральной совокупности; возникает в результате применения неподходящего метода отбора.

Смешанность (impurity) — степень разнородности классов в подразделе данных.

Спред (spread) — разность между лучшими ценами заявок на продажу и на покупку в один и тот же момент времени на какой-либо актив. Словом "спред" также называют разность цен двух различных сходных товаров, торгуемых на открытых рынках.

Спред между ценой спроса и ценой предложения (bid-ask spread) — сумма, на которую цена предложения (ask) превышает цену спроса (bid) на ценную бумагу на рынке. Спред между ценами спроса и предложения представляет собой, по сути, разницу между самой высокой ценой, которую покупатель готов заплатить за актив, и самой низкой ценой, которую продавец готов принять, чтобы продать его. См. <https://www.investopedia.com/terms/b/bid-askspread.asp>.

Срединная цена (mid-price) — цена между лучшей ценой продавцов актива, или ценой предложения (ask), и лучшей ценой покупателей актива, или ценой спроса (bid). См. https://en.wikipedia.org/wiki/Mid_price.

Срезовые данные (cross-sectional data) — данные, агрегированные по разным переменным за один и тот же момент времени, например в формате "страна — ВВП за 2018 год".

Стандартное нормальное распределение (standard normal) — нормальное распределение со средним, равным 0, и стандартным отклонением, равным 1.

Стандартное отклонение (standard deviation) — квадратный корень из дисперсии.

Статистическая проверка отношения логарифмического правдоподобия (log-likelihood ratio test, LLR) обозначает проверку значимости двух биномиальных распределений, также именуется проверкой статистического показателя G -квадрат.

См. <https://mahout.apache.org/users/clustering/llr---log-likelihood-ratio.html>.

Стопаут (stop out), или принудительная остановка, — более низкий чем "маржин-колл"-уровень, помогающий избежать потери финансовых средств, предоставленных через кредитное плечо. Если количество убытков продолжает увеличиваться, то сделка, или даже несколько открытых позиций, закрывается принудительно, причем основные убытки в этом случае несет только инвестор.

См. <https://profitgid.ru/margin-call-i-stop-out.html>.

Сторона покупки (buy side) и **сторона продажи** (sell side) относятся к фирмам, которые соответственно приобретают и торгуют товарами и услугами. Применительно к финансам, сторона покупки относится к покупке ценных бумаг и включает в себя инвестиционных менеджеров, пенсионные фонды и хедж-фонды; а сторона продажи относится к фирмам, которые выпускают, продают

или торгуют ценными бумагами, и включает в себя инвестиционные банки, консультативные фирмы и корпорации.

Сторона ставки (side of the bet) показывает, в какую сторону пойдет движение цены — в длинную (вверх) или короткую (вниз). Этот термин имеет синоним "позиция", т. е. соответственно длинная или короткая позиция.

Стоимость под риском (value at risk, VaR) — статистический метод измерения портфельного риска: максимальная величина стоимости, которую можно ожидать потерять за данный временной горизонт. Например, если портфель имеет 95% дневной VaR в размере 100 тыс. долларов США, то это означает, что, статистически говоря, существует 95% уверенность в том, что портфель не потеряет более 100 тыс. долларов США стоимости в течение следующего дня.

Стратегия "стоп-лосс" — стратегия торговли с преобладанием ордеров стоп-лосс, т. е. размещением сделок на бирже для немедленного их исполнения, в случае если актив достигает определенной ценовой точки. Как следует из названия, этот вид ордера предназначен для ограничения убытков.

Структурный сдвиг (structural break), или структурное изменение или разрыв, — неожиданный сдвиг во временном ряде, который может привести к огромным ошибкам предсказания и ненадежности модели в целом.

См. https://en.wikipedia.org/wiki/Structural_break.

Телеконференция о корпоративных заработках (earnings call) — встреча или веб-трансляция, в которой публичная компания обсуждает финансовые результаты отчетного периода. Данный термин происходит от понятия "заработки в расчете на пай (долю) в акционерном капитале" (earnings per share, EPS), т. е. число в нижней строке отчета о прибылях и убытках, разделенное на число долей (паев) в обращении. Базирующийся в США Национальный институт по связям с инвесторами (NIRI) сообщает, что 92% компаний, представленных своими членами, проводят такие телеконференции, и что практически все они транслируются через Интернет. См. https://en.wikipedia.org/wiki/Earnings_call.

Тестовый набор (test set) — набор образцов, используемый только для оценивания результативности полностью натренированной модели. После оценивания окончательной модели на тестовом наборе образцов дальнейшей настройки модели не требуется (Рикардо Гутьеррес-Окуна (Ricardo Gutierrez-Osuna)).

Теория современного инвестиционного портфеля (Modern portfolio theory, MPT), или портфельная теория Марковица, — теория, которая описывает то, как инвесторы, не склонные к риску, могут создавать портфели для оптимизации или максимизации ожидаемой возвратности на основе заданного уровня рыночного риска, подчеркивая, что риск является неотъемлемой частью более высокого вознаграждения.

См. <https://www.investopedia.com/terms/m/modernportfoliotheory.asp>.

Тик (tick) — мера минимального восходящего или нисходящего движения цены ценной бумаги. Тик также может относиться к изменению цены ценной бумаги

от сделки к сделке. С 2001 г., с появлением децимализации, минимальный размер тика для торговли акциями выше 1 доллара составляет 1 цент. Тик представляет собой стандарт, на котором стоимость ценной бумаги может колебаться. Тик обеспечивает определенное приращение цены, отраженное в местной валюте, связанной с рынком, на котором торгуется ценная бумага, на которое может измениться общая цена ценной бумаги.

Топологическое многообразие (manifold), или многосвязная область, — это топологическое пространство, которое является локально евклидовым (т. е. вокруг каждой точки есть окрестности, топологически такие же, как разомкнутый единичный шар в \mathbb{R}^n). Ближайшей аналогией является вера древних людей, что Земля плоская. Это расхождение возникает главным образом из-за того, что в малых масштабах, которые мы видим, Земля действительно выглядит плоской, и топологические многообразия составляют обобщение объектов, на которых мы могли бы жить и сталкиваться с проблемой круглой/плоской Земли. См. <http://mathworld.wolfram.com/Manifold.html>.

Точность (accuracy) модели — мера статистического смещения. Описывает систематические ошибки; иными словами, это близость результатов измерений истинному значению. См. https://en.wikipedia.org/wiki/Accuracy_and_precision.

Точность top-1 (top-1 accuracy) — условная точность, когда модельный ответ (с наибольшей вероятностью) должен быть строго ожидаемым ответом. Точность top-5 означает, что любой из 5 модельных ответов с наибольшей вероятностью должен соответствовать ожидаемому ответу.

Трансакционные издержки (transaction cost) — стоимость осуществления любой экономической сделки, сопровождающая участие в рынке. В инвестировании это затраты, понесенные при покупке или продаже активов, такие как комиссии и спред. См. https://en.wikipedia.org/wiki/Transaction_cost.

Тренировочный набор (training set) — набор образцов, используемых алгоритмом для самообучения, т. е. усвоения взаимосвязей и регулярностей, которое проходит в форме подгонки параметров/весов модели.

Умное бета (smart beta) инвестирование сочетает в себе преимущества пассивного инвестирования и преимущества активных инвестиционных стратегий. Цель умного бета-инвестирования — получить альфа, снизить риск или увеличить диверсификацию за стоимость ниже, чем традиционный активный менеджмент, и слегка выше, чем прямое индексное инвестирование. Оно стремится отыскать наилучшую структуру оптимально диверсифицированного портфеля.

См. <https://www.investopedia.com/terms/s/smart-beta.asp>.

Усадочный оценщик (shrinkage estimator) — оценщик, или правило оценки, который дает улучшенную оценку величины за счет явного или неявного встраивания другой информации. В этом смысле используется для регуляризации неточных поставленных задач статистического вывода.

См. https://en.wikipedia.org/wiki/Shrinkage_estimator.

Уровень коротких продаж (short interest), т. е. продаж в надежде на падение цены — индикатор рыночного сентимента, представляющий собой число акций (долевых ценных бумаг), которые инвесторы продали, но еще не покрыли. Он используется для определения того, надеются ли инвесторы на понижение (медвежий рынок) или повышение цены (бычий рынок), и иногда используется в качестве встречного индикатора.

См. <https://www.investopedia.com/terms/s/shortinterest.asp>.

Учет заслуг (credit assignment) — процесс выявления среди множества действий, выбранных в эпизоде, тех, которые отвечают за конечный результат. И более того, это попытка определить лучшие и худшие решения, выбранные во время эпизода, для того чтобы лучшие решения подкреплялись, а худшие штрафовались. Используется в подкрепляемом обучении.

См. <https://www.igi-global.com/dictionary/credit-assignment/40159>.

Усвоение топологического многообразия (manifold learning) — попытка раскрыть структуру топологического многообразия (manifold), или многосвязной области, в наборе данных, где топологическое многообразие — это пространство, локально сходное с евклидовым. Данный метод представляет собой нелинейный подход к снижению размерности.

Фактор (factor) — это характерная, поддающаяся квантификации особенность актива с существенной информацией о риске и возвратности. В случае долевых ценных бумаг наиболее известные и наиболее документированные факторы включают стоимость, размер, импульс, низкую волатильность и качество.

Факторное инвестирование (factor investing) — стратегия, которая выбирает ценные бумаги по атрибутам, связанным с более высокой возвратностью. Существует два основных типа факторов, определяющих возвратность активов: макроэкономические факторы и стилевые факторы. Первые охватывают широкие риски по всем классам активов, в то время как вторые призваны объяснять возвратности и риски в рамках классов активов. Некоторые общие макроэкономические факторы включают кредит, инфляцию и ликвидность, в то время как стилевые факторы среди прочих охватывают стиль, стоимость и импульс.

См. <https://www.investopedia.com/terms/f/factor-investing.asp>.

Факторная нагрузка (factor loading) — мера влияния данного фактора на возвратность финансового актива; может интерпретироваться как стандартизированный регрессионный коэффициент и таким образом указывать на корреляцию актива с фактором.

См. <https://www.theanalysisfactor.com/factor-analysis-1-introduction/>.

Финансовый возврат (return), или возврат на инвестицию, или финансовая отдача, — деньги, сделанные или потерянные на инвестиции. Возврат выражается номинально как изменение денежной стоимости инвестиции с течением времени, либо в процентах из соотношения прибыли к инвестициям. В этом случае он именуется возвратностью. См. <https://www.investopedia.com/terms/r/return.asp>.

Форвардная оптимизация (walk-forward optimization) — торговая стратегия оптимизируется с помощью внутривыборочных данных внутри временного окна во временном ряде финансовых данных. Остальная часть данных резервируется для тестирования вне выборки. Небольшая порция зарезервированных данных, следующих после внутривыборочных, тестируется с использованием зафиксированных результатов. Затем внутривыборочное временное окно переносится вперед на период, охватываемый вневыборочным тестом, и процесс повторяется. В конце все зафиксированные результаты используются для выявления торговой стратегии. См. также "вне выборки и внутри выборки".

См. https://en.wikipedia.org/wiki/Walk_forward_optimization.

Форвардное тестирование (forward testing) — моделирование реальных рыночных данных виртуально, только на бумаге, т. е. хотя вы и двигаетесь по рынкам вживую, на самом деле вы не вкладываете реальные деньги и занимаетесь виртуальной торговлей с целью лучше понять движения рынков.

См. https://en.wikipedia.org/wiki/Walk_forward_optimization.

Фундаментальные финансовые показатели акционерного капитала (equity fundamentals) включают в себя денежный поток, доходность активов, соотношение собственного/заемного капитала (консервативное лeverажное соотношение), историю удержания прибыли для финансирования будущего роста, обоснованность управления капиталом для максимизации заработков и финансовых возвратов акционеров.

См. <https://www.investopedia.com/articles/fundamental/03/022603.asp>.

Фьючерсный контракт (futures contract) — юридическое соглашение о покупке или продаже определенного товара или актива по заранее определенной цене в конкретное время в будущем. Фьючерсные контракты стандартизированы по качеству и количеству для облегчения торговли на фьючерсной бирже. Покупатель фьючерсного контракта берет на себя обязательство купить базовый актив по истечении срока действия фьючерсного контракта. Продавец фьючерсного контракта берет на себя обязательство предоставить базовый актив на дату истечения срока действия.

См. <https://www.investopedia.com/terms/f/futurescontract.asp>.

Хвост (tail) — длинная узкая часть частотного распределения, где относительно предельные значения встречаются с низкой частотой.

Хеджирование (hedging) — действие по снижению риска торговли. Например, если свечной график предполагает, что рынок биткойна очень нерешителен, и торговец считает, что цена собирается подняться, то он может купить еще несколько BTC, одновременно торгуя коррелированным активом (например, ETH).

Цена предложения (ask) — цена, которую продавец готов принять за ценную бумагу или другой финансовый инструмент. Также упоминается как предложение (offer). См. <http://www.businessdictionary.com/definition/ask.htm>.

Цена спроса (bid), или заявка, — самая высокая цена, которую любой покупатель готов заплатить за данную ценную бумагу в данный момент времени. Как пра-

вило, заявка (bid) ниже, чем предложения (ask), и разница между ними называется спредом между ценой спроса и ценой предложения.

См. <http://www.investorwords.com/469/bid.htm>.

Ценные бумаги (securities) — финансовые или инвестиционные инструменты (некоторые — оборотные, другие — нет), купленные и проданные на финансовых рынках, такие как облигации, долговые обязательства, облигации, опционы, акции и купоны. См. <http://www.businessdictionary.com/definition/securities.html>.

Частота ошибок top-5 (top-5 error rate) — процент тестовых примеров, для которых правильный класс не был среди верхних 5 предсказанных классов.

Ценностная функция (value function) — функция состояния (или пары "состояние — действие"), которая оценивает, насколько хорошо агенту находиться в данном состоянии (или насколько хорошо выполнять данное действие в данном состоянии). См. <http://www.incompleteideas.net/book/ebook/node34.html>.

Ценовая коррекция (retracement) — временный ценовой откат, который происходит в рамках более крупного тренда. Ключевой момент здесь в том, что этот ценовой откат является кратковременным и не указывает на изменение более крупного тренда. Разворот (reversal), с другой стороны, это ситуация, когда тренд меняет направление, т. е. цена, вероятно, будет продолжать двигаться в этом разворотном направлении в течение продолжительного периода времени. См. <https://www.investopedia.com/articles/trading/06/retracements.asp>.

Шкалирование (scaling) — сплющивание или расширение данных, обычно для приведения многочисленных переменных к одинаковой шкале измерения.

Шортить (go short) — продавать активы, открывать короткую позицию на продажу в надежде на снижение цены актива, причем это вовсе не означает, что позиция будет удерживаться непродолжительное время. Торговец "шортит", когда цена падает, давая ему возможность получать прибыль от продажи.

Электронные коммуникационные сети ECN (electronic communication networks) — электронная система осуществления сделок купли-продажи биржевых товаров.

Предметный указатель



Ссылки на номера страниц типа "A-123" указывают на главы электронного архива.

A

alphalens, библиотека 154
Applied Quantitative Research (AQR) 52

B

BeautifulSoup, библиотека 120

C

CatBoost, библиотека 383
ConvNet, сверточная нейронная сеть A-41
CountVectorizer, класс 454
◇ визуализирование распределения
лексикона 455
◇ отыскание наиболее похожих
документов 455
CSV, стандартный формат данных 102

D

Dataiku, компания 518
Datamir, компания по поиску информации
117
DataRobot, компания 518

G

gensim, библиотека 482
◇ автоматическое обнаружение фраз 501
◇ словарные векторы из финансовой
отчетности SEC 500
Gnip, компания 117
GridSearchCV, класс:
◇ влияние параметров на тестовые отметки
375
◇ донастройка параметров 374
◇ тестирование на отложенном наборе 377

H

H2O, компания 518
HDF5, иерархический формат данных 102
HTML:
◇ извлечение данных 120
◇ разбор с использованием регулярных
выражений 126

K

KNeighborsRegressor, класс 220

L

LightGBM, библиотека 383, 466

N

NLTK, библиотека 450
 NumPy, библиотека 144
 n-грамма 441, 447

P

pandas, библиотека 144
 pandas-datareader, библиотека 101
 ◇ биржа IEX 90
 Parquet, двоичный столбчатый формат данных 102
 pyfolio, библиотека
 ◇ влияние факторов 173
 ◇ моделирование событийного риска 174
 ◇ периоды просадки 173
 ◇ получение входа:
 ▫ из библиотеки alphalens 170
 ▫ из бэктеста библиотеки zipline 171
 ◇ применительно:
 ▫ к внутри- и вневыборочной результативности 170
 ▫ к измерению результативности 167
 ◇ моделирование форвардного тестирования 171
 ◇ статистика результативности 172
 PyMC3, библиотека:
 ◇ применительно к вероятностному программированию 319
 ◇ рабочий поток 320
 ▫ байесова логистическая регрессия 320
 ▫ вариационный Байес 324
 ▫ диагностика модели 324
 ▫ метод Монте-Карло марковской цепи 323
 ▫ предсказание 327
 Python, реализация перекрестного контроля 215

Q

Quandl, платформа 93
 QuantConnect, платформа 519
 Quantopian, платформа 90, 519
 ◇ справочные материалы и аккаунт 91
 QuantRocket, платформа 520

R

RavenPack, поставщик аналитики Больших данных 118
 requests, библиотека 120
 RS Metrics, поставщик приложений Больших данных на основе ИИ 118
 RSS-канал 96

S

Scrapy, библиотека 124
 Selenium, библиотека 121
 spaCy, библиотека 443
 splash, библиотека 124
 statsmodels, библиотека:
 ◇ оценивание методом обычных наименьших квадратов 256
 ◇ применительно к проведению статистических рассуждений 268
 StockTwits, социальная сеть 118

T

TA-Lib, библиотека 147
 TensorBoard, инструмент визуализации 500
 textacy, библиотека 443, 447
 TextBlob, библиотека 450
 TfIdFTransformer, класс 456
 TfIdFVectorizer, класс 456
 ◇ обобщение новостных статей 458
 ◇ сглаживание 457
 Two Sigma, квантитативный хеджевый фонд 519
 t-распределенное стохастическое вложение соседей (t-SNE), алгоритм 404

X

XBRL, расширяемый язык деловой отчетности 95
 XGBoost, библиотека 383

Z

zipline, библиотека:
 ◇ применение 148
 ▫ для построения портфеля 165
 ▫ для тестирования портфеля 165

А

Автокодировщик 199, А-97
 Автокорреляция, измерение 277
 Автоматизация браузерная 121
 Авторегрессия векторная (VAR) 299
 Актив под управлением (AUM) 51
 Алгебра линейная, мигрирование в сторону иерархических вероятностных моделей 469
 Алгоритм:
 ◇ GloVe 498
 ◇ t-распределенного стохастического вложения соседей (t-SNE) 404
 ◇ автоматически обучающийся 66, 362, А-2
 ◇ адаптивного бустинга (AdaBoost) 363, 364
 ▫ применение с библиотекой sklearn 368, 372
 ◇ аппроксимации ожиданий 433
 ◇ Бройдена — Флетчера — Гольдфарба — Шанно (BFGS) квазиньютоновский 323
 ◇ кластерный 401
 ▫ иерархические кластеры 198, 199
 ▫ кластеризация k-средних 198, 199
 ▫ модели гауссовых смесей 198, 199
 ▫ плотностные кластеры 198, 199
 ◇ линейный 403
 ◇ Метрополиса — Хастингса 317
 ◇ равномерной аппроксимации и проекции топологического многообразия (UMAP) 404
 ◇ торговый, привлечение людей через Интернет 56
 Альтернативные системы торговли (ATS) 43
 Альтернативы средне-дисперсной оптимизации:
 ◇ инвестирование риск-факторное 191
 ◇ оптимизация глобальная портфельная 187
 ◇ паритет риска 190
 ▫ иерархический 191
 ◇ портфель:
 ▫ минимально-дисперсный 186
 ▫ 1/n 186
 ◇ правило Келли 187
 Альфа 129, 168, 172

Альфа-фактор 129, 130
 ◇ единственный из рыночных данных 149
 ◇ исследование 129, 512
 ◇ ресурсы 162
 Анализ:
 ◇ зависимостей 442
 ◇ латентно-семантический (LSA) 469, 496
 ▫ вероятностный (pLSA) 468, 474, 475
 ◇ независимых компонент (ICA) 402, 407, 413, А-8
 ▫ алгоритм 414
 ▫ допущения 414
 ▫ с помощью библиотеки Sklearn 415
 ◇ сентиментный 458, 462
 ▫ с помощью модели Doc2Vec 504
 Анализ главных компонент (PCA) 199, 402, 407, 409, 415, А-8
 ◇ визуализирование в двух размерностях 408
 ◇ для алгоритмической торговли 415
 ▫ рисковые факторы, ведомые данными 415
 ▫ характеристические портфели 418
 ◇ допущения 408
 ◇ матрица ковариаций 409
 ◇ применение с сингулярным разложением (SVD) 411
 ◇ с помощью библиотеки Sklearn 412
 Аннотация:
 ◇ лингвистическая 442
 ◇ семантическая 442
 ◇ частеречная 442

Б

Балтийский фрахтовый индекс сухогрузного тоннажа (BDI) 107
 Бар:
 ◇ временной 84, 85
 ◇ долларовый 87
 ◇ объемный 86
 ◇ тиковый 83
 Бета 172
 Библиотека данных Кеннета Френча 101
 Библиотеки Python с открытыми исходными кодами для алгоритмической торговли 162
 Биграмма 441
 Биржа NASDAQ 71

Бустинг градиентный:

- ◇ настройка 466
- ◇ применение с библиотекой Sklearn 372

Бэкстест 113

Бэкестирование, ловушки:

- ◇ вопросы реализации 177
- ◇ длина бэктеста 179
- ◇ издержки торговые 177
- ◇ координация сделок временная 178
- ◇ остановка оптимальная 179
- ◇ переподгонка бэктестовая 178
- ◇ предотвращение 175
- ◇ прочесывание данных 178
- ◇ результативность на основе рыночной цены 177
- ◇ сложности данных 176

В

Важность признаков глобальная 392

Векторы словарные 492

- ◇ из финансовой отчетности SEC:
 - автоматическое обнаружение фраз 501
 - влияние параметрических настроек на результативность 503
 - оценивание модели 502
 - предобработка 500
 - тренировка модели 502
- ◇ предварительно натренированные:
 - GloVe 498
 - применение 497

Вероятность:

- ◇ апостериорная 460
- ◇ априорная 459

Вложение:

- ◇ векторное, оценивание 496
- ◇ локально-линейное (LLE) 421
- ◇ слов векторное 438, 491
 - кодирование семантики 492
 - тренировка 499
- ◇ соседей t-распределенное стохастическое (t-SNE) 423

Возврат финансовый:

- ◇ на вложенный капитал 46, 251
- ◇ одинаково распределенный и взаимно независимый (iid) 168
- ◇ форвардный, создание 155

Возвратность скользящая 48

Всемирный банк 102

Вход ортонормированный 250

Выборка:

- ◇ градиентная односторонняя (GOSS) 380
- ◇ отрицательная (NEG) 495
- ◇ по Гиббсу 317
- ◇ по Метрополису — Хастингсу 317
- Выделение основ слов 442, 450
- Выезд на кривой отдачи 48
- Выработка признаков 512
- Выражения регулярные применительно к разбору HTML 126
- Выскабливание данных 119

Г

Генерирование выборок методом Монте-Карло марковских цепей 315

- ◇ выборка:
 - по Гиббсу 317
 - по Метрополису — Гастингсу 317
- ◇ гамильтонов Монте-Карло (HMC) 317
- Гетероскедастичность 221, 235
- Гиперпараметр:
 - ◇ класс GridsearchCV для деревьев решений 346
 - ◇ кривые усвоения 348
 - ◇ настройка 345, 374
 - ◇ обследование древесной структуры 347
- Гипотеза об эффективном рынке (ЕМН) 47, 181
- Гомоскедастичность 231
- Горизонт данных временной 114
- Граф знаний 442

Д

Данные 510

- ◇ Twitter, сентиментный анализ:
 - мультиномный наивный Байес 462
 - сентиментные отметки объекта TextBlob 463
- ◇ альтернативные:
 - источники:
 - датчики 107, 109
 - предприятия 107, 108
 - физические лица 107
 - критерии оценивания 111
 - обработка 119
 - оценивание 110
- ◇ геолокационные 110, 119

Данные (*prod.*):

- ◇ извлечение из HTML 120
- ◇ контроль качества 510
- ◇ ордерной книги 72
 - TotalView-ITCH биржи NASDAQ 73
 - лексический разбор двоичных сообщений ITCH 74, 76
 - реконструирование ордерной книги 79, 80
 - реконструирование сделок 79, 80
 - протокол FIX 72
 - регуляризация тиковых данных 83
- ◇ преобразование в факторы 143
- ◇ разблокирование ценности с помощью компетенции в предметной области 511
- ◇ рыночные 67
- ◇ с привязкой к определенной точке во времени (point-in-time, PIT) 149
- ◇ сбор данных на основе бронирования ресторанов 122
- ◇ службы OpenTable, выскабливание 119
- ◇ социальных настроений 117
- ◇ спутниковые 118
- ◇ текстовые:
 - извлечение признаков 439
 - лексемизация 441
- ◇ усваивание регуляризаторов из данных 195
- ◇ финансовой отчетности 95
- ◇ фундаментальные 95
- ◇ электронных почтовых квитанций 119

Датчик:

- ◇ геолокационные данные 110
- ◇ спутник 109

Движение цены акции, динамическое оценивание вероятностей 312, 313

Дерево:

- ◇ классификационное:
 - оптимизация в целях чистоты узлов 338
 - построение 338
 - тренировка 339
- ◇ решений 194, 332
 - визуализирование 340
 - класс поиска по решетке параметров GridsearchCV 335
 - кодирование класса кастомизированного перекрестного контроля 335
 - обрезка 343
 - оценивание предсказаний 340
 - подготовка 334

- построение регрессионного дерева 336
- правила принятия решений 332, 334
- применение 334
- регуляризация 342
- сильные стороны 349
- слабые стороны 349

Диагностика модели 324, 513

- ◇ верификационная проверка оптимизации 515
- ◇ определение адресных модельных целей 514
- ◇ предсказательные проверки апостериорного распределения 326
- ◇ схождение 324
- ◇ теорема об отсутствии бесплатных обедов 513
- ◇ управление компромиссом между смещением и дисперсией 514

Дисперсия глобально-минимальная (GMV) 186

Дневная стоимость под риском (VaR) 172

Документы, измерение сходства 452

Долгая краткосрочная память (LSTM) A-70

Допущение об одинаковой распределенности и взаимной независимости 225

Доступ:

- ◇ к данным дистанционный с помощью библиотеки pandas 87
 - модуль pandas-datareader для рыночных данных 88
 - чтение HTML-таблиц 87
- ◇ через API к рыночным данным 87

Дробление акций 99

Е

Евростат 102

Естественно-языковой инструментарий (NLTK) 450

З

Задача многометочная 197

Закрепление меток 442

Заработки:

- ◇ в расчете на долю собственности в акционерном капитале (EPS) 176
- ◇ на долю собственности в акционерном капитале разбавленные (Diluted EPS) 99

И

- Импульс А-17
- Инвестиции в стратегический потенциал 54
- Индексирование латентно-семантическое (LSI) 468 *См. Анализ латентно-семантический (LSA):*
 - ◇ реализация с помощью библиотеки sklearn 475
 - ◇ сильные стороны 472
 - ◇ слабые стороны 472
- Индикатор:
 - ◇ сентиментный 135
 - ◇ чистого дисбаланса ордеров (NOII) 74
- Инновации алгоритмические для улучшения результативности:
 - ◇ алгоритмы поиска разбивок 379
 - ◇ аппроксимация функции потери второго порядка 378
 - ◇ поуровневый рост против полистового роста 380
 - ◇ признаки:
 - и оптимизации 382
 - категориальные 382
 - ◇ тренировка на основе GPU 381
- Инструменты аналитические для диагностики и извлечения признаков 274
- Интеграция данных 511
- Интернет вещей (IoT) 109
- Интерфейс межмашинной связи «компьютер — компьютер» 73
- Информация взаимная (MI), информационная теория применительно к оцениванию признаков 210
- Искусственный интеллект (AI) 53, 200, А-1

К

- Квантиль факторный 155
- Классификатор наивный байесов 459
 - ◇ допущение об условной независимости 460
 - ◇ мультиномный 464
 - оценивание 461
 - тренировка 461
- Классификация 265, 458
 - ◇ стандартная промышленная 96
- Классифицирование новостных статей 461

- Кластеризация 425
 - ◇ агломеративная 430
 - ◇ иерархическая 199, 429
 - дендограммы 430
 - инвестиционных портфелей (НСР) 192
 - сильные стороны 431
 - слабые стороны 432
 - ◇ на основе k-средних 199, 426
 - иерархическая 429
 - модель гауссовой смеси (GMM) 433
 - оценивание качества 427
 - паритет риска иерархический 434
 - плотностная 432
 - ◇ на основе модели гауссовой смеси 199
 - ◇ плотностная 199, 432
 - пространственная приложений с шумом (DBSCAN) 432
 - ◇ разделительная 430
- Когерентность тематическая 479
- Коинтеграция 303
 - ◇ применение для стратегии парной торговли 305
 - ◇ тестирование 304
- Коллинеарность 231
- Комиссия по ценным бумагам и биржам США (SEC) 40, 41, 43, 44, 95
- Компетентное знание предметной области для разблокирования ценности данных 511
- Компромисс между смещением и дисперсией 197
- Конвейер по обработке ЕЯ 440, 443
 - ◇ с помощью библиотек spaCy/textacy 443
 - API потоковый 448
 - n-граммы 447
 - аннотирование предложений 444
 - лексемизация 444
 - обработка документов пакетная 446
 - обработка ЕЯ многоязычная 448
 - определение границ предложений 446
 - разбор лексический 444
 - распознавание именованных сущностей 447
- Конкурсный набор данных веб-сайта Yelp 463
 - ◇ комбинирование текстовых признаков с числовыми признаками 465
 - ◇ логистическая регрессия по схеме "один против всех" 465

Конкурсный набор данных веб-сайта Yelp (прод.):

- ◇ машина градиентно-бустинговая 466
- ◇ модель мультиномиальная (мультиномная) наивная байесова 464
- ◇ регрессия мультиномная логистическая 465

◇ точность эталонная 464

Конструирование модели 211

- ◇ компромисс между смещением и дисперсией 211
- ◇ кривые усвоения 214
- ◇ недоподгонка против переподгонки 212
- ◇ управление компромиссом 213

Конструирование признаков

См. Выработка признаков

Корень квадратный:

- ◇ из среднеквадратической логарифмированной ошибки (RMSLE) 204
- ◇ из среднеквадратической ошибки (RMSE) 204

Коррелограмма 278

Корреляция внутрирядовая

См. Автокорреляция

Коэффициент:

- ◇ информационный (IC) 34, 159, 160, 169
- ◇ Калмар (коэффициент просадки) 172
- ◇ Омега 172
- ◇ передаточный (TC) 169
- ◇ просадки (Калмар) 172
- ◇ Сортино 172
- ◇ хвостовой 172
- ◇ Шарпа (SR) 167, 168
 - моделирование как вероятностная модель 328
 - сравнение с результативностью 329

Кривая:

- ◇ прецизионности — полноты 207
- ◇ усвоения (самообучения) 214

Критерий:

- ◇ информационный Акаике (AIC) 234
- ◇ информационный байесов (BIC) 234
- ◇ оценивания качества альтернативных наборов данных 113
 - горизонт временной 114
 - надежность 114
 - риски правовые 113
 - риски репутационные 113

▫ частота 114

▫ эксклюзивность 114

- ◇ оценивания качества содержимого сигнала альтернативных наборов данных 111

▫ качество альфа 112

▫ классы активов 111

▫ премии за риск 112

▫ содержимое альфа 112

▫ стиль инвестиционный 112

- ◇ оценивания технических аспектов альтернативных наборов данных
 - задержка 115
 - формат 115

Кэри-трейд 49

Л

Латентное размещение Дирихле (LDA) 468, 476, 477, 496

- ◇ когерентность тематическая 479

- ◇ модель генеративная 477

- ◇ оценивание 479

- ◇ перплексивность 479

- ◇ процесс реконструирования 478

- ◇ работа с размещением 476

- ◇ реализация с помощью библиотеки:
 - genism 482
 - sklearn 480

- ◇ результаты визуализирования с помощью библиотеки pyLDAvis 481

Лексема 441

Лексикон 441

Лемматизация 442

Лес случайный 350

- ◇ модели ансамблевые 350

- ◇ настройка 356

- ◇ построение 355

- ◇ применение бэггинга к деревьям решений 353

- ◇ свойство важности признаков 358

- ◇ сильные стороны 360

- ◇ слабые стороны 360

- ◇ снижение модельной дисперсии с помощью бэггинга 351

- ◇ тестирование внепакетное 359

- ◇ тренировка 356

Линейность 231

М

Материальная непубличная информация (MNPI) 113

Матрица терм-документная (DTM) 467, 468, 470, 474, 487, 490

◊ с помощью библиотеки sklearn 453

Машина градиентно-бустинговая (GBM) 363, 369

◊ реализации 377

- алгоритмические инновации для стимулирования результативности 378

◊ результаты:

- аддитивные объяснения Шепли (SHAP) 395
- глобальная важность признаков 392
- графики частичной зависимости 393
- интерпретация 391

Машинное обучение (ML) 34, 66, 165, 509, 516

◊ байесово:

- вероятность максимальная апостериорная (MAP) 310
- обновление предварительного допущения на основе эмпирического наблюдения 309
- применение библиотеки Theano 319
- работа 308

◊ в инвестиционной индустрии 42

- инвестирование факторное 46
- первопроходцы алгоритмические 51
- торговля:
 - высокочастотная 44
 - электронная 43
- умные бета-фонды 46

◊ данные альтернативные 54, 56

◊ инструменты 512, 518

◊ поток рабочих 201

- k-ближайших соседей (KNN) 201
- анализ разведывательный 209
- выработка признаков 209
- задачи:
 - классификационные 205
 - регрессионные 204
- извлечение признаков 209
- конструирование модели 211
- настройка параметров с помощью библиотеки scikit-learn 219

- отбор автоматически обучающегося алгоритма 211
- перекрестный контроль для отбора модели 215
- подготовка данных 209
- предсказание против статистического рассуждения 202
- разведывание признаков 209
- сбор данных 209
- сложности финансового перекрестного контроля 221
- цели и метрики 202

◊ применительно к решению задач с данными 512

◊ с текстовыми данными, примеры 443

Метод:

◊ ансамблевого усиления (бустинг), 351

◊ гамильтонов Монте-Карло (HMC) 317

◊ Монте-Карло марковских цепей (MCMC) 314

- приближенное рассуждение 324
- убедительные интервалы 324

◊ обычных наименьших квадратов (OLS) 224

◊ оценивания робастный 224

◊ пакетного усреднения (бэггинг) 351

◊ усадочный 224, 226

Микроструктура рынка 68

Множитель лагранжей (LM) 233

Модель:

◊ ARIMA:

- ARMAX-модель 289
- SARIMAX-модель 289
- выявление количества:
 - MA-членов 288
 - AR-членов 288
- построение 287

◊ Doc2Vec:

- анализ сентиментный 504
- мешок слов распределенный (DBOW) 504
- память распределенная (DM) 504
- создание входных данных 505
- тренировка на сентиментных данных веб-сайта Yelp 504

◊ Skip-Gram (SG) 494

- архитектура в библиотеке Keras 499
 - компоненты модели 500
 - оценивание контрастивное к шуму 499

Модель (*прод.*):

- ◇ Word2Vec 493
 - автоматическое обнаружение фраз 495
 - для машинного перевода 508
 - упрощение активационной функции softmax 494
 - ◇ авторегрессионная 285
 - векторная (VAR) 299
 - число временных сдвигов 286
 - диагностика подгонки 286
 - ◇ авторегрессионной условной гетероскедастичности (ARCH) 293
 - ◇ ансамблевая 343, 350
 - ◇ временного ряда байесова 330
 - ◇ волатильности стохастические 330
 - ◇ гауссовой смеси (GMM) 433
 - алгоритм максимизации ожиданий 433
 - ◇ генеративная 477
 - ◇ гетероскедастичности обобщенной авторегрессионной условной (GARCH) 293, 294
 - выбор порядка сдвига 294
 - ◇ градиентно-бустинговых машин (GBM) 368
 - бустинг стохастический градиентный 372
 - донастройка 370
 - остановка досрочная 371
 - подвыборка 372
 - размер ансамбля 371
 - темп усвоения 371
 - тренировка 370
 - усадка 371
 - ◇ линейная обобщенная 224
 - ◇ логистическая регрессионная 265
 - оценивание максимального правдоподобия 267
 - функция:
 - логистическая 266
 - целевая 266
 - ◇ мешка слов (BOW) 440, 447, 452
 - непрерывного (CBOW) 493
 - распределенного (DBOW) 504
 - ◇ многомерного временного ряда 298
 - векторная авторегрессионная модель (VAR) 299
 - применение модели VAR для макроэкономических фундаментальных прогнозов 300
 - система уравнений 298
 - ◇ множественной линейной регрессии 226
 - градиентный спуск 230
 - диагностирование проблем 233
 - исполнение 237
 - наименьшие квадраты 228
 - оценивание максимального правдоподобия 229
 - статистическое рассуждение 232
 - теорема Гаусса — Маркова (GMT) 230
 - тренировка 228
 - устранение проблем 233
 - формулирование 227
 - ◇ нейронно-языковая, усвоение словоупотребления 492
 - ◇ одномерных временных рядов 285
 - построение:
 - авторегрессионной модели 285
 - модели ARIMA 287
 - модели скользящего среднего 287
 - прогнозирование:
 - волатильности 292
 - макроэкономических фундаментальных показателей 290, 292
 - ◇ прогнозирования волатильности 294, 296, 297
 - ◇ распределенной памяти (DM) 504
 - ◇ скользящего среднего 287
 - выявление числа временных сдвигов 287
 - связь моделей AR и MA 287
 - ◇ факторная линейная 240
 - ценообразования активов (CAPM) 241
 - получение рисков факторов 242
 - ◇ ценообразования капитальных активов (CAPM) 46, 140, 164, 180, 241
 - ◇ черно-ящичная 515
- Моделирование тематическое 438, 467
- ◇ для телеконференций о корпоративных заработках 484
 - отзывы о деятельности предприятий веб-сайта Yelp 488
 - предобработка данных 485
 - проведение экспериментов 487
 - тренировка модели 485
- Мультиколлинеарность 236

Н

Наблюдение 459

◊ внепакетное (OOB) 359

◊ сведения 309

Надежность данных 114

Наименьшие квадраты:

◊ взвешенные (WLS) 236

◊ обобщенные (GLS) 236

◊ обычные (OLS) 228, 230, 256

- применение библиотeki statsmodels 237, 238

Наложение эмбарго 223

Настройка:

◊ гиперпараметров:

- поиск рандомизированный решеточный 386
- регуляризация 386
- усвоение параметров 385
- целевые функции и функция потери 385

◊ параметров с помощью библиотеки scikit-learn 219

- кривые перекрестного контроля с помощью библиотеки yellowbricks 220
- кривые усвоения 220
- применение класса GridSearchCV 221

Насыщение A-12

Нейронная сеть:

◊ прямого распространения (FFNN) A-70

◊ рекуррентная (RNN) 64, 66, A-1, A-70

◊ сверточная (CNN), A-1 A-41, A-70

Номер учетный adsh 98

Норма L2 405

Нормальность ошибки 232

Нью-Йоркская фондовая биржа (NYSE) 68

О

Обнаружение:

◊ ликвидности *См.* Предвосхищение заявок

◊ фраз автоматическое 495

Оборачиваемость фактора 161

Обработка естественного языка (NLP) 199

◊ рабочий поток 440

◊ с помощью библиотеки TextBlob 450

- выделение основ слов 450

▫ лемматизация 450

▫ полярность настроений 451

▫ субъективность 451

◊ сложности 439

Обрезка дерева 343

Общепринятые принципы бухгалтерского учета США (GAAP) 95

Общий регламент по защите данных (GDPR) 113

Объяснение Шепли (SHapley) аддитивное 516

◊ анализ взаимодействия признаков 398

◊ применение графиков силы 397

◊ резюмирование 396

Одинаковая распределенность и взаимная независимость (iid) 168, 221, 225, 275

Оптимизация среднedisперсная 182

◊ альтернативы 186

◊ процедура 182

◊ сложности и недостатки 185

◊ эффективная граница на Python 182, 185

Ордер:

◊ "все или ничего" 71

◊ "заполнить или уничтожить" 72

◊ "исполнить немедленно или аннулировать" 72

◊ лимитный 71

◊ на открытии/закрытии рынка 72

◊ не по установившемуся рынку 72

◊ рыночный 71

◊ стоп (стоп-ордер) 71

Остановка досрочная A-15

Отбор модели 215

Отношения логарифмических правдоподобий (LLR) 269

Отсев A-15

Оценивание:

◊ альтернативных наборов данных 110

◊ контрастивное к шуму (NCE) 495, 499

◊ максимального правдоподобия (MLE) 228, 229, 267

◊ максимальной апостериорной вероятности (MAP) 310

◊ результатов 388

- перекрестного контроля по всем моделям 388

Оценка несмещенная 231

Оценщик сэндвичный 235

Ошибка 228

- ◇ медианная абсолютная (MedAE) 204
- ◇ среднеквадратическая (MSE) A-12
- ◇ средняя абсолютная (MAE) 204

П

Память распределенная См. Модель распределенной памяти

Паритет риска иерархический (HRP) 191, 434

Первичное публичное размещение (IPO) 79

Перекрестный контроль (CV) 216

- ◇ в финансах, сложности 221
 - комбинаторный 223
 - временного ряда 222
 - с наложением эмбарго 223
 - с прочисткой 223

◇ отложенный тестовый набор 216

◇ итератор KFold 217

◇ комбинаторный 223

◇ с исключением:

- по P образцов 219
- по одному образцу 218

◇ тасованное подразделение ShuffleSplit 219

Переподгонка:

- ◇ бэктестовая, риски 515
- ◇ интерпретация 249
- ◇ точность предсказания 249
- ◇ управление 248
- ◇ устранение риска 342

Период нерепрезентативный 177

Пермутация 392

Перплексивность 479

Персонально идентифицируемая информация (PII) 113

Платформа торговая онлайн-овая 519

- ◇ QuantConnect 519
- ◇ Quantopian 519
- ◇ QuantRocket 520

Плотностная пространственная кластеризация приложений с шумом (DBSCAN), алгоритм 432

◇ иерархический алгоритм DBSCAN 432

Плотность апостериорная наибольшая (HPD) 325

Площадка торговая 68

Площадь под кривой (AUC) 207

Подготовка данных:

- ◇ линейная регрессия 251
 - кодирование категориальных переменных 254
 - отбор альфа-факторов 253
 - очистка данных 253
 - преобразование альфа-факторов 253
 - разведывательный анализ данных 254
 - расчет целевого финансового возврата 252
 - создание универсума ценных бумаг и временной горизонт 251
 - создание форвардных финансовых возвратов 255

Подход:

- ◇ Блэка — Литтермана 187, 307
- ◇ статистический 314
 - по сравнению с детерминированным подходом 314

Полнота 207

Полосы Боллинджера 149

Портфель:

- ◇ инвестиционный:
 - иерархически кластеризованный (HCP) 192
 - перебалансировка 165
 - построение с помощью библиотеки zipline 165
 - тестирование с помощью библиотеки zipline 165
 - управление риском и возвратностью 180

◇ характеристический 407

Портфельное страхование с постоянной пропорцией (CPPI) 133

Поставщик рыночных данных 93

Правдоподобие 459

Правила синтаксические и грамматические 442

Правило Келли 187

Предвосхищение заявок 45

Предобработка текста 458

Преобразование аффинное A-10

Прецизионность 208

Прибыль операционная 242

Признаки, извлечение из текстовых данных 439

Примеры использования МО для торговли на финансовых рынках 62

- ◇ глубокий анализ данных для извлечения признаков 62

- ◇ контролируемое самообучение для создания альфа-факторов 63
- ◇ подкрепляемое самообучение 65
- ◇ распределение финансовых средств среди портфельных активов 64
- ◇ тестирование торговых идей 64
- Проверка предсказательная апостериорного распределения 326
- Проклятие размерности 402
- Просадка максимальная 172
- Протокол:
 - ◇ FIX электронного обмена финансовой информацией 72
 - ◇ TCP, базовый управляющий протокол передачи 73
 - ◇ передачи гипертекста (HTTP) 120
- Прочистка 223
- Прямой доступ к рынку (DMA) 44
- Публичная служба распространения информации (PDS) 96

Р

- Рабочие характеристики приемника (ROC) 207
- Рабочий поток машинного обучения 201
- Равномерная аппроксимация:
 - ◇ и проекция топологического многообразия (UMAP), алгоритм 404
 - ◇ топологического многообразия и проекция в нем (UMAP) 424
- Разбиение рекурсивное бинарное 345
- Разбор лексический 444
- Разложение сингулярное (SVD) 250, 411, 470
- Разметка частеречная 439
- Разнообразные источники данных, сочетание факторов 152
- Распознавание именованных сущностей (NER) 442
- Распределение:
 - ◇ апостериорное 309
 - ◇ априорное 309
 - объективное 311
 - отбор 311
 - сопряженность 312
 - субъективное 312
 - эмпирическое 312
- Рассуждение:
 - ◇ вариационное (VI) 318
 - с автоматическим дифференцированием (ADVI) 318
 - ◇ на основе причинно-следственных взаимосвязей 202
 - ◇ приближенное 323
 - вариационное 318
 - генерирование выборок методом Монте-Карло марковских цепей 315
 - детерминированные методы 314
 - стохастические технические решения 314
 - стохастическое рассуждение на основе выборок 315
- Расширяемый язык деловой отчетности (XBRL) 67, 95
- Революция альтернативных данных 105, 107
- Регрессия гребневая 249
- ◇ применение библиотеки sklearn 261
 - коэффициенты 262
 - настройка регуляризационных параметров с помощью перекрестного контроля 261
 - результаты перекрестного контроля 262
 - траектории гребневых коэффициентов 262
- ◇ работа с регрессией 249
- Регрессия:
 - ◇ лассо 251
 - применение библиотеки sklearn 263
 - перекрестно-проверенный информационный коэффициент 264
 - траектория 264
 - ◇ линейная:
 - для статистического рассуждения и предсказания 225
 - применение для предсказания финансового возврата 251
 - регуляризация 248
 - ◇ линейная обычными наименьшими квадратами:
 - с помощью библиотеки:
 - диагностическая статистика 256
 - кастомизированный перекрестный контроль временного ряда 257
 - отбор признаков и цели 258
 - оценивание 256
 - перекрестный контроль модели 258
 - тестовые предсказания 259
 - с помощью библиотеки sklearn 257

Регрессия (*прод.*):

- ◇ логистическая:
 - байесова 320
 - визуализация и фигурное обозначение 322
 - модуль обобщенных линейных моделей 322
 - рассуждение на основе MAP 322
- мультиномиальная 465
- по схеме "один против всех" 465
- предсказание ценовых движений с помощью библиотеки `sklearn` 270
- применительно к предсказанию 270

◇ Фама — Макбета 244, 246, 247

Регрессор k ближайших соседей 220

Результативность предсказательная по факторным квантилям 156

Ретротест *См.* Бэкtest

Решение задачи:

- ◇ внутрирядовая корреляция 236
- ◇ гетероскедастичность 235
- ◇ меры качества подгонки 234
- ◇ мультиколлинеарность 236

Риск аттестационный 169

Рынок альтернативных данных 115

- ◇ геолокационные данные 119
- ◇ данные социальных настроений 117
- ◇ данные электронных почтовых квитанций 119
- ◇ поставщики данных 116
- ◇ примеры использования 116
- ◇ спутниковые данные 118

Ряд временной:

- ◇ применение преобразований 283
- ◇ разложение на регулярности 275
- ◇ с фундаментальными данными:
 - извлечение набора данных с примечаниями 97
 - извлечение финансовых отчетов 96
 - получение ежеквартальной отчетности компании Apple 98
 - построение 96
 - временного ряда цена/зарботки 99, 100

С

Самообучение:

- ◇ глубокое (DL) A-1
- ◇ контролируемое 197, 200

- ◇ неконтролируемое 198, 401
 - кластерные алгоритмы 198
 - применения 198
 - снижение размерности 199

◇ подкрепляемое 200

◇ предсказательное A-97

Свертка A-41

Сент-луисские данные ФРС 101

Сеть состязательная генеративная (GAN) A-2

Система электронного сбора, анализа и извлечения данных (EDGAR) 95

Систематическое смещение:

- ◇ из-за выживших 176
 - ◇ из-за забегания вперед 176
- Системы управления базами данных 517
- ◇ графовая база данных 517
 - ◇ хранилище:
 - в формате "ключ — значения" 517
 - документное 517
 - столбчатое 517

Смешанность Джини 339

Смешанность узла *См.* Чистота узла

Снижение размерности 343, 401, 402

◇ автокодировщики 199

◇ алгоритм:

- линейный 403
 - нелинейный 404
- ◇ анализ главных компонент (PCA) 199
 - ◇ линейное 407
 - ◇ проклятие размерности 402, 405, 406
 - ◇ усвоение проекций в топологическом многообразии 199, 421

Соотношение:

- ◇ информационное (IR) 34, 159, 168, 169
- ◇ цены к зарботкам (P/E) 99

Спуск градиентный 230

◇ стохастический (SGD) 228, A-16

- применение с помощью библиотеки `sklearn` 239

Среда разработки интерактивная (IDE) 519

Среднее значение условное 231

Ставка:

- ◇ единственный актив 189
- ◇ многочисленные активы 190
- ◇ определение размера 187
- ◇ оптимальный размер 188

Статистика скользящего окна:

- ◇ вычисление 276
- ◇ скользящие средние 277
- ◇ экспоненциальное сглаживание 277

Стационарность:

- ◇ диагностирование единичных корней 280, 281
- ◇ достижение 278
- ◇ обращение с единичными корнями 280, 281
- ◇ преобразования временного ряда 279
- ◇ проверки на единичные корни 281

Стенограмма телеконференций

- о корпоративных зарплатах 125

Степени свободы 328

Стоимость под риском (VaR) 177

Стоп-слова 441

Стратегия:

- ◇ алгоритмической торговли 61
- ◇ квантитативная 61
- ◇ торговли на финансовых рынках:
 - бэкестирование стратегии 60
 - исполнение 58
 - исследование альфа-факторов 58
 - конструирование 57
 - менеджмент данных 58
 - оптимизация портфеля 60
 - оценивание 58
 - привлечение источников данных 58
 - риск-менеджмент 60

Сумма квадратов остатков (RSS) 228

Сходство косинусное 424, 453

Т

Таксономии XBRL 95

Темные пулы 44

Темп усвоения (самообучения) A-17

Теорема:

- ◇ Байеса 309, 459
- ◇ Гаусса — Маркова (GMT) 230, 232
- ◇ об отсутствии бесплатных обедов 196

Теория:

- ◇ арбитражного ценообразования (APT) 241
- ◇ современного инвестиционного портфеля (MPT) 46, 180

Технологии управления данными 516

- ◇ технологии Больших данных:
 - Hadoop 517
 - Spark 517

Торговля:

- ◇ высокочастотная (HFT) 44, 66
- ◇ по плану 165

Точность предсказания 249

Триграмма 441

У

Униграмма 441

Управление выбросами 176

Усвоение проекции в топологическом многообразии 199, 421

- ◇ алгоритм t-SNE 423
- ◇ алгоритм UMAP 424

Ф

Фактор 131

- ◇ волатильностный 140
 - ключевые метрики 141
 - обоснование 140
- ◇ встроенный в платформу Quantopian 147, 148
- ◇ импульсный 132
 - ключевые метрики 134
 - обоснование 132
- ◇ качественный 141
 - ключевые метрики 142
 - обоснование 142
- ◇ обрачиваемость 161
- ◇ отобранный, вычисление из сырых фоновых данных:
 - вычисление:
 - импульсных факторов 145
 - факторных бета 146
 - загрузка данных 144
 - использование разных периодов владения 146
 - передискретизация с дневной частоты на месячную 144
 - применение сдвинутых во времени финансовых возвратов 146
- ◇ преобразование данных в факторы 143
- ◇ размерный 140
 - ключевые метрики 141
 - обоснование 140
- ◇ сентиментный 132
 - ключевые метрики 134
 - обоснование 132
- ◇ создание факторных квантилей 155
- ◇ сочетание из разнообразных источников данных 152, 153

Фактор (прод.):

- ◇ стоимостной 136
 - ключевые метрики 138
 - обоснование 137

Финансовые отчеты и примечания (FSN) 96

Фонд:

- ◇ биржевой (ETF) 43
- ◇ взаимный 43
- ◇ квантоментальный 53

Формат двоичных данных 383

Формула Пифагора 405

Функция:

- ◇ автокорреляции (ACF) 278
 - частичной (PACF) 278
- ◇ активационная softmax иерархическая 494
- ◇ правдоподобия 309

Х

Хранение данные эффективное с помощью библиотеки pandas 102

Ц

Цена средневзвешенная по объему (VWAP) 83, 86

Ценообразование:

- ◇ средневзвешенное по времени (TWAP) 62
 - ◇ средневзвешенное по объему (VWAP) 62
- Центральный индексный ключ 98

Ч

Частота:

- ◇ данных 114
 - ◇ документная обратная 453
 - ◇ истинных утверждений (TP rate) 207
 - ◇ ложных утверждений (FP rate) 207
 - ◇ появления ошибок 339
 - ◇ термина 453
- Чикагская товарная биржа 66
-
- Чистота узла 339

Ш

Штраф по норме параметра A-14

Э

Эксклюзивность данных 114
Электронные коммуникационные сети (ECN) 43
Электронный обмен финансовой информацией (FIX) 72
Элемент вентильный рекуррентный (Gated Recurrent Units, GRU) A-70
Энтропия перекрестная 339