

Содержание

Предисловие	15
Глава 1. Библиотека pandas и анализ данных	19
Знакомство с библиотекой pandas	19
Обработка данных, анализ, наука и библиотека pandas	21
Обработка данных.....	22
Анализ данных.....	22
Наука о данных.....	23
Предназначение библиотеки pandas	23
Процесс анализа данных.....	23
Процесс	24
Взаимосвязь между книгой и процессом анализа данных.....	28
Понятия «данные» и «анализ» в контексте нашего знакомства с библиотекой pandas	29
Типы данных	29
Временные ряды	31
Общие понятия анализа и статистики	31
Другие библиотеки Python, работающие вместе с библиотекой pandas	34
Численные и научные вычисления – NumPy и SciPy	34
Статистический анализ – StatsModels	35
Машинное обучение – scikit-learn	35
РуМС – стохастическое байесовское моделирование	35
Визуализация данных – matplotlib и seaborn.....	36
Выводы	36
Глава 2. Запуск библиотеки pandas	37
Установка Anaconda	37
IPython и Jupyter Notebook.....	39
IPython	39
Jupyter Notebook	40
Знакомство со структурами данных библиотеки pandas – Series и DataFrame	43
Импорт pandas.....	43
Объект Series.....	44
Объект DataFrame.....	48
Загрузка данных из CSV-файла в объект DataFrame.....	52
Визуализация.....	55
Выводы	56
Глава 3. Представление одномерных данных с помощью объекта Series	57
Настройка библиотеки pandas	58
Создание объекта Series	58
Создание объекта Series с помощью питоновских списков и словарей	58
Создание объекта Series с помощью функций NumPy	60
Создание объекта Series с помощью скалярного значения	61

Свойства <code>.index</code> и <code>.values</code>	61
Размер и форма объекта <code>Series</code>	62
Установка индекса во время создания объекта <code>Series</code>	63
Использование методов <code>.head()</code> , <code>.tail()</code> и <code>.take()</code> для вывода значений	64
Получение значений в объекте <code>Series</code> по метке или позиции	65
Поиск по метке с помощью оператора <code>[]</code> и свойства <code>.ix[]</code>	65
Явный поиск по позиции с помощью свойства <code>.iloc[]</code>	67
Явный поиск по меткам с помощью свойства <code>.loc[]</code>	67
Создание срезов объекта <code>Series</code>	68
Выравнивание данных по меткам индекса	73
Выполнение логического отбора	76
Переиндексация объекта <code>Series</code>	78
Модификация объекта <code>Series</code> на месте	81
Выводы	83

Глава 4. Представление табличных и многомерных данных с помощью объекта <code>DataFrame</code>	84
Настройка библиотеки <code>pandas</code>	85
Создание объектов <code>DataFrame</code>	85
Создание объекта <code>DataFrame</code> на основе результатов функций <code>NumPy</code>	85
Создание объекта <code>DataFrame</code> с помощью питонового словаря и объектов <code>Series</code>	87
Создание объекта <code>DataFrame</code> на основе CSV-файла	89
Доступ к данным внутри объекта <code>DataFrame</code>	90
Отбор столбцов в объекте <code>DataFrame</code>	91
Отбор строк в объекте <code>DataFrame</code>	92
Поиск скалярного значения по метке и позиции с помощью <code>.at[]</code> и <code>.iat[]</code>	93
Создание среза датафрейма с помощью оператора <code>[]</code>	94
Логический отбор строк	94
Одновременный отбор строк и столбцов	96
Выводы	96

Глава 5. Выполнение операций над объектом <code>DataFrame</code> и его содержимым	97
Настройка библиотеки <code>pandas</code>	97
Переименование столбцов	98
Добавление новых столбцов с помощью оператора <code>[]</code> и метода <code>.insert()</code>	99
Добавление столбцов за счет расширения датафрейма	100
Добавление столбцов с помощью конкатенации	101
Переупорядочивание столбцов	102
Замена содержимого столбца	103
Удаление столбцов	103
Присоединение новых строк	105
Конкатенация строк	107
Добавление и замена строк за счет расширения датафрейма	109
Удаление строк с помощью метода <code>.drop()</code>	109
Удаление строк с помощью логического отбора	110
Удаление строк с помощью среза	111
Выводы	111

Глава 6. Индексация данных	112
Настройка библиотеки <code>pandas</code>	112

Важность применения индексов	113
Типы индексов библиотеки pandas	115
Основной тип Index	115
Индексы Int64Index и RangeIndex, в качестве меток используются целые числа	115
Индекс Float64Index, в качестве меток используются числа с плавающей точкой	117
Представление дискретных интервалов с использованием IntervalIndex	117
Категории в качестве индекса – CategoricalIndex	118
Индексация по датам и времени с помощью DatetimeIndex	119
Индексация периодов времени с помощью PeriodIndex	119
Работа с индексами	120
Создание и использование индекса в объекте Series или объекте DataFrame	120
Отбор значений с помощью индекса	121
Преобразование данных в индекс и получение данных из индекса	123
Переиндексация объекта библиотеки pandas	124
Иерархическая индексация	125
Выводы	128
Глава 7. Категориальные данные	129
Настройка библиотеки pandas	129
Создание категориальных переменных	130
Переименование категорий	135
Добавление категорий	136
Удаление категорий	136
Удаление неиспользуемых категорий	137
Установка категорий	137
Вычисление описательных статистик для категориальной переменной	138
Обработка школьных оценок	138
Выводы	141
Глава 8. Численные и статистические методы	142
Настройка библиотеки pandas	143
Применение численных методов к объектам библиотеки pandas	143
Выполнение арифметических операций над объектами DataFrame или Series	144
Вычисление количества значений	147
Определение уникальных значений (и их встречаемости)	147
Вычисление минимума и максимума	148
Вычисление n наименьших значений и n наибольших значений	148
Вычисление накопленных значений	149
Выполнение статистических операций с объектами библиотеки pandas	150
Получение итоговых описательных статистик	150
Измерение центральной тенденции: среднее, медиана и мода	151
Вычисление дисперсии и стандартного отклонения	153
Вычисление ковариации и корреляции	154
Дискретизация и квантилизация данных	156
Вычисление ранга значений	160
Вычисление процентного изменения для каждого наблюдения серии	161
Выполнение операций со скользящим окном	161
Создание случайной выборки данных	164
Выводы	165
Глава 9. Загрузка данных	166
Настройка библиотеки pandas	166

Работа с CSV-файлами и текстовыми/табличными данными	167
Исследование CSV-файла	167
Чтение CSV-файла в датафрейм	168
Указание индекса столбца при чтении CSV-файла	168
Вывод и спецификация типа данных	169
Указание имен столбцов.....	169
Указание конкретных столбцов для загрузки.....	170
Сохранение датафрейма в CSV-файл	170
Работа с данными, в которых используются разделители полей.....	171
Обработка загрязненных данных, в которых используются разделители полей	172
Чтение и запись данных в формате Excel	174
Чтение и запись JSON-файлов	177
Чтение HTML-файлов из интернета.....	178
Чтение и запись HDF5-файлов	180
Загрузка CSV-файлов из интернета	182
Чтение из базы данных SQL и запись в базу данных SQL.....	182
Загрузка данных с удаленных сервисов	185
Загрузка базы данных по экономической статистике Федерального резервного банка Сент-Луиса	185
Загрузка данных Кеннета Френча	187
Загрузка данных Всемирного банка	188
Выводы	192
Глава 10. Приведение данных в порядок	193
Настройка библиотеки pandas	193
Что такое приведение данных в порядок?.....	194
Как работать с пропущенными данными	195
Поиск значений NaN в объектах библиотеки pandas	196
Удаление пропущенных данных.....	198
Обработка значений NaN в ходе арифметических операций	201
Заполнение пропущенных данных.....	202
Прямое и обратное заполнение пропущенных значений	203
Заполнение с помощью меток индекса.....	204
Выполнение интерполяции пропущенных значений.....	205
Обработка дублирующихся данных	207
Преобразование данных	210
Сопоставление значений другим значениям	210
Замена значений	211
Применение функций для преобразования данных	214
Выводы	218
Глава 11. Объединение, связывание и изменение формы данных	219
Настройка библиотеки pandas	219
Конкатенация данных, расположенных в нескольких объектах.....	220
Понимание семантики конкатенации, принятой по умолчанию	220
Переключение осей выравнивания	224
Определение типа соединения	225
Присоединение вместо конкатенации	226
Игнорирование меток индекса	226
Слияние и соединение данных	227
Слияние данных, расположенных в нескольких объектах.....	227

Настройка семантики соединения при выполнении слияния	230
Поворот данных для преобразования значений в индексы и наоборот	233
Состыковка и расстыковка данных	234
Состыковка с помощью неиерархических индексов	234
Расстыковка с помощью иерархических индексов	236
Расплавление данных для преобразования «широкого» формата в «длинный» и наоборот	239
Преимущества использования состыкованных данных	240
Выводы	241
Глава 12. Агрегирование данных	242
Настройка библиотеки pandas	242
Обзор схемы «разделение – применение – объединение»	243
Данные для примеров	244
Разделение данных	244
Группировка по значениям отдельного столбца	244
Просмотр результатов группировки	245
Группировка по нескольким столбцам	248
Группировка по уровням индекса	249
Применение агрегирующих функций, преобразований и фильтров	251
Применение агрегирующих функций к группам	251
Преобразование групп данных	253
Исключение групп из процедуры агрегирования	258
Выводы	259
Глава 13. Анализ временных рядов	260
Настройка библиотеки pandas	260
Представление дат, времени и интервалов	261
Объекты datetime, day и time	261
Создание временной метки с помощью объекта Timestamp	263
Использование объекта Timedelta для представления временного интервала	263
Введение во временные ряды	264
Индексация с помощью объекта DatetimeIndex	264
Создание временного ряда с определенной частотой	269
Вычисление новых дат с помощью смещений	271
Представление временных интервалов с помощью смещений дат	271
Привязанные смещения	274
Представление промежутков времени с помощью объектов Period	275
Создание временного интервала с помощью объекта Period	275
Индексация с помощью объекта PeriodIndex	277
Обработка праздников с помощью календарей	279
Нормализация временных меток с помощью часовых поясов	280
Операции с временными рядами	284
Опережение и запаздывание	284
Преобразование частоты временного ряда	287
Увеличение или уменьшение шага дискретизации временного ряда	289
Применение к временному ряду операций на основе скользящего окна	294
Выводы	297
Глава 14. Визуализация	298
Настройка библиотеки pandas	299

Основные инструменты визуализации	299
Создание графиков временных рядов	300
Настройка внешнего вида графика временного ряда	302
Виды графиков, часто использующиеся в статистическом анализе данных	314
Демонстрация относительных различий с помощью столбиковых диаграмм	314
Визуализация распределений данных с помощью гистограмм	316
Визуализация распределений категориальных данных с помощью ящичных диаграмм с усами	318
Отображение накопленных итогов с помощью площадных диаграмм	318
Визуализация взаимосвязи между двумя переменными с помощью диаграммы рассеяния	320
Визуализация оценок распределения с помощью графика ядерной оценки плотности	320
Визуализация корреляций между несколькими переменными с помощью матрицы диаграмм рассеяния	321
Отображение взаимосвязей между несколькими переменными с помощью тепловых карт	322
Размещение нескольких графиков на одном рисунке вручную	323
Выводы	325

Приложение 1. Советы по оптимизации вычислений

в библиотеке pandas	326
Базовое итерирование	327
Итерирование с помощью метода <code>.iterrows()</code>	328
Более лучший способ итерирования с помощью метода <code>.apply()</code>	328
Векторизация с помощью объектов Series	329
Векторизация с помощью массивов NumPy	329
Выводы	330

Приложение 2. Улучшение производительности pandas

(из официального пособия по библиотеке pandas)	331
Написание расширений на языке C для pandas	331
«Чистый» Python	331
Обычный Cython	333
Использование библиотеки Numba	333
Jit	334
Vectorize	334
Вычисление выражений с помощью функции <code>eval()</code>	335
Поддерживаемый синтаксис	336
Примеры использования функции <code>eval()</code>	336
Метод <code>DataFrame.eval()</code>	337

Приложение 3. Используем pandas для больших данных

Работаем с данными бейсбольных игр	341
Внутреннее представление датафрейма	343
Подтипы	344
Оптимизация числовых столбцов с помощью понижающего преобразования	345
Сравнение способов хранения числовых и строковых значений	347
Оптимизация типов object с помощью типа category	349
Задаем типы во время считывания данных	353
Выводы	355

Приложение 4. Пример предварительной подготовки данных в pandas (конкурсная задача Tinkoff Data Science Challenge)	356
Считывание CSV-файла в объект DataFrame	357
Преобразование типов переменных	382
Переименование категорий переменных	384
Обработка редких категорий	385
Разбиение набора данных на обучающую и контрольную.....	390
Импутация пропусков	393
Конструирование новых признаков.....	398
Создание переменной, у которой значения основаны на значениях исходной переменной	399
Создание бинарной переменной на основе значений количественных переменных.....	401
Создание переменной, у которой каждое значение – среднее значение количественной переменной, взятое по уровню категориальной переменной.....	402
Возведение в квадрат.....	403
Дамми-кодирование (One-hot Encoding)	405
Кодирование контрастами (Effect Coding)	407
Присвоение категориям в лексикографическом порядке целочисленных значений, начиная с 0 (Label Encoding)	407
Создание переменной, у которой каждое значение – частота наблюдений в категории переменных (Frequency Encoding)	409
Кодирование вероятностями зависимой переменной (Likelihood Encoding).....	410
Кодировка средним значением зависимой переменной, сглаженным через сигмоидальную функцию	412
Кодировка средним значением зависимой переменной, сглаженным через параметр регуляризации.....	415
Кодировка простым средним значением зависимой переменной по схеме leave-one-out	415
Кодировка простым средним значением зависимой переменной по схеме K-fold	416
Кодировка средним значением зависимой переменной, сглаженным через сигмоидальную функцию, по схеме K-fold	416
Присвоение категориям в зависимости от порядка их появления целочисленных значений, начиная с 1 (Ordinal Encoding).....	421
Бинарное кодирование (Binary Encoding)	422
Создание переменных-взаимодействий.....	422
Категоризация (биннинг) количественной переменной	423
Дамми-кодирование и подготовка массивов для обучения и проверки.....	429
Выбор метрики качества.....	431
Построение моделей случайного леса, градиентного бустинга и логистической регрессии	447
Математический аппарат логистической регрессии	488
Отдельная предварительная подготовка данных для логистической регрессии	494
Построение логистической регрессии в библиотеке H2O	538
Приложение 5. Пример предварительной подготовки данных в pandas (конкурсная задача предсказания отклика ОТП Банка)	563
Этап I. Построение модели на обучающей выборке – части исторической выборки и ее проверка на контрольной выборке – части исторической выборки.....	566
I.1. Считывание CSV-файла, содержащего исторические данные, в объект DataFrame.....	566

I.2. Преобразование типов переменных.....	567
I.3. Импутация пропусков, не использующая результаты математических вычислений (импутация, которую можно выполнять до/после разбиения на обучение/контроль)	570
I.4. Обработка редких категорий	572
I.5. Конструирование новых признаков, не использующее результаты математических вычислений (которое можно выполнять до/после разбиения на обучение/контроль)	575
I.6. Разбиение на обучающую и контрольную выборки.....	577
I.7. Импутация пропусков, использующая статистику – результаты математических вычислений (ее нужно выполнять после разбиения на обучение и контроль).....	577
I.8. Поиск преобразований переменных, максимизирующих нормальность распределения (дается в сокращенном виде).....	578
I.9. Биннинг как один из способов конструирования новых признаков, использующий результаты математических вычислений (нужно выполнять только после разбиения на обучение и контроль).....	582
I.10. Выполнение преобразований, исходя из информации гистограмм распределения и графиков квантиль-квантиль	587
I.11. Конструирование новых признаков	587
I.12. Стандартизация.....	589
I.13. Дамми-кодирование	589
I.14. Подготовка массивов признаков и массивов меток зависимой переменной	590
I.15. Построение логистической регрессии с помощью класса LogisticRegression библиотеки scikit-learn	590
I.16. Настройка гиперпараметров логистической регрессии с помощью класса GridSearchCV	591
I.17. Отбор признаков для логистической регрессии с помощью случайного леса (класса RFE).....	592
I.18. Отбор признаков для логистической регрессии с помощью BorutaPy.....	594
I.19. Проблема дисбаланса классов	597
Этап II. Построение модели на всей исторической выборке и применение к новым данным	605
II.1. Считывание CSV-файла, содержащего исторические данные, в объект DataFrame.....	605
II.2. Предварительная обработка исторических данных.....	605
II.3. Обучение модели логистической регрессии на всех исторических данных	611
II.4. Считывание CSV-файла, содержащего новые данные, в объект DataFrame.....	611
II.5. Предварительная обработка новых данных	612
II.6. Применение модели логистической регрессии, построенной на всех исторических данных, к новым данным.....	612
Приложение 6. Работа с датами и строками	615
Работа с датами.....	615
Работа со строками.....	618
Изменение регистра строк	618
Изменение строкового значения	620
Определение пола клиента по отчеству	620
Удаление лишних символов из строк	624
Удаление повторяющихся строк.....	627
Извлечение нужных символов из строк.....	628

Приложение 7. Работа с предупреждением

SettingWithCopyWarning в библиотеке pandas	631
Что представляет из себя предупреждение SettingWithCopyWarning?	632
Присваивание по цепочке (chained assignment)	633
Скрытая цепочка	635
Советы и рекомендации по работе с предупреждением SettingWithCopyWarning	637
Отключение предупреждения	637
Однотипные и многотипные объекты	638
Ложные срабатывания	639
Подробнее о присваивании по цепочке	641
Ложные пропуски	643
И вновь о скрытой цепочке	644

Приложение 8. От Pandas к Scikit-Learn – новый подход

к управлению рабочими процессами	647
Новый уровень интеграции Scikit-Learn с Pandas	647
Краткое резюме и цели статьи	647
Знакомство с классом ColumnTransformer и обновленным классом OneHotEncoder	648
Задача предсказания цен на недвижимость с Kaggle	649
Исследуем данные	649
Удаление зависимой переменной из обучающего набора	649
Кодировка отдельного столбца со строковыми значениями	650
Scikit-Learn – только двумерные данные	650
Импортируем класс, создаем экземпляр класса – модель, обучаем модель – трехэтапный процесс работы с моделью	650
У нас NumPy-массив. Где имена столбцов?	651
Проверка корректности первой строки данных	652
Используем метод .inverse_transform() для автоматизации данной операции	652
Применение преобразования к тестовому набору	653
Проблема № 1 – новые категории в тестовом наборе	654
Ошибка: Unknown Category	654
Проблема № 2 – пропущенные значения в тестовом наборе	655
Проблема № 3 – пропущенные значения в обучающем наборе	655
Необходимость импутации пропущенных значений	656
Больше о методе .fit_transform()	657
Применение нескольких преобразований к тестовому набору	657
Применение конвейера	658
Почему для тестового набора мы вызываем только метод .transform()?	659
Выполнение преобразований для нескольких столбцов со строковыми значениями	659
Обращение к отдельным этапам конвейера	659
Использование нового ColumnTransformer для отбора столбцов	660
Передаем конвейер в ColumnTransformer	660
Передаем весь объект DataFrame в ColumnTransformer	661
Извлечение названий признаков	661
Преобразование количественных переменных	662
Работа со всеми количественными признаками	662
Передаем конвейер с преобразованиями для категориальных признаков и конвейер с преобразованиями для количественных признаков в ColumnTransformer	663
Машинное обучение	664
Перекрестная проверка	665

Отбор наилучших значений гиперпараметров с помощью решетчатого поиска	665
Представление результатов решетчатого поиска в виде датафрейма pandas	666
Создание пользовательского трансформера, выполняющего основные преобразования	666
Редкие категории	666
Написание пользовательского класса	666
Применение класса BasicTransformer	669
Использование BasicTransformer в конвейере	669
Биннинг и преобразование количественных переменных с помощью нового класса KBinsDcretizer	670
Отдельная обработка всех столбцов с годами с помощью ColumnTransformer	671
Применение RobustScaler и FunctionTransformer	673
Предметный указатель	677

Предисловие

Pandas – популярная библиотека Python, применяющаяся для практического анализа данных в реальных проектах. Она предлагает воспользоваться эффективными, быстрыми и высокопроизводительными структурами данных, которые упрощают предварительную обработку и анализ информации. Это учебное пособие существенно поможет вам, предоставив в ваше распоряжение внушительный набор инструментов, предлагаемых библиотекой pandas для выполнения различных операций с данными и их анализа.

О СОДЕРЖАНИИ КНИГИ

Глава 1 «Библиотека pandas и анализ данных» – это практическое введение в основные функции библиотеки pandas. Предназначение этой главы – дать некоторое представление об использовании библиотеки pandas в контексте статистики и науки о данных. В этой главе мы рассмотрим несколько принципов, лежащих в основе науки о данных, и покажем, как они реализованы в библиотеке pandas. Эта глава задает контекст для каждой последующей главы, связанной с наукой о данных.

Глава 2 «Запуск библиотеки pandas» проинструктирует читателя по поводу того, как загрузить и установить библиотеку pandas, и познакомит его с некоторыми базовыми понятиями библиотеки pandas. Мы также рассмотрим, как можно работать с примерами с помощью iPython и тетрадок Jupyter.

Глава 3 «Представление одномерных данных с помощью объекта Series» познакомит читателя со структурой данных Series, которая используется для представления одномерных индексированных данных. Читатель узнает о том, как создавать объекты Series и как работать с данными, хранящимися внутри этих объектов. Кроме того, он узнает об индексах и выравнивании данных, а также о том, как объект Series можно использовать для создания срезов данных.

Глава 4 «Представление табличных и многомерных данных с помощью объекта DataFrame» познакомит читателя со структурой данных DataFrame, которая используется для представления и индексации многомерных данных. В этой главе читатель научится создавать объекты DataFrame, используя различные наборы статических данных, и выполнять отбор определенных столбцов и строк внутри датафрейма. Сложные запросы, операции с данными и индексация будут рассмотрены в следующей главе.

Глава 5 «Выполнение операций над объектом DataFrame и его содержимым» расширяет предыдущую главу и расскажет о том, как выполнять более сложные операции с объектом DataFrame. Мы начнем с добавления и удаления столбцов и строк, рассмотрим модификацию данных внутри объекта DataFrame (а также создание измененной копии) и выполнение арифметических операций с данными, научимся создавать иерархические индексы, а также вычислять популярные статистики по данным датафрейма.

Глава 6 «Индексация данных» расскажет об использовании различных типов индекса библиотеки pandas (Int64Index, RangeIndex, IntervalIndex, CategoricalIndex, DatetimeIndex, PeriodIndex).

Глава 7 «Категориальные данные» познакомит читателя с тем, как создавать объекты `Categorical` для представления категориальных данных и использовать их в работе.

В главе 8 «Численные и статистические методы» рассматриваются различные арифметические операции над объектами `Series` и `DataFrame`, а также вычисление статистик для объектов `pandas`.

Глава 9 «Загрузка данных» расскажет о том, как можно загрузить данные из внешних источников и записать в объекты `Series` и `DataFrame`. Кроме того, в этой главе рассматривается загрузка данных из разных источников, таких как файлы, HTTP-серверы, системы баз данных и веб-службы. Также рассматривается обработка данных в форматах `CSV`, `HTML` и `JSON`.

В главе 10 «Приведение данных в порядок» будет рассказано о том, как приводить данные в порядок, чтобы они были пригодны для анализа.

Глава 11 «Объединение, связывание и изменение формы данных» расскажет читателю о том, как можно взять несколько объектов `pandas` и объединить их с помощью операций соединения, слияния и конкатенации.

Глава 12 «Агрегация данных» расскажет о группировке и агрегации данных. В библиотеке `pandas` эти операции выполняются с помощью схемы «разделение – применение – объединение». Читатель научится использовать эту схему для различных способов группировки данных, а также применять агрегирующие функции для вычисления результатов по каждой группе данных.

Глава 13 «Анализ временных рядов» расскажет о том, как работать с временными рядами в библиотеке `pandas`. В этой главе будут освещены широкие возможности библиотеки `pandas`, существенно облегчающие анализ временных рядов.

Глава 14 «Визуализация» научит вас создавать визуализации данных на основе данных, хранящихся в объектах `Series` и `DataFrame`. Мы начнем с изучения основ, создания простой диаграммы настройки нескольких параметров диаграммы (настройки легенд, меток и цветов), рассмотрим создание нескольких распространенных типов графиков, которые используются для представления различных типов данных.

В приложении 1 «Советы по оптимизации вычислений в библиотеке `pandas`» даются некоторые рекомендации по ускорению вычислений в `pandas`.

Приложение 2 «Улучшение производительности `pandas`» представляет собой перевод одноименного раздела официального пособия по библиотеке `pandas` <https://pandas.pydata.org/pandas-docs/stable/enhancingperf.html>.

Приложение 3 «Используем `pandas` для больших данных» расскажет, как за счет применения более эффективных типов данных можно уменьшить использование памяти.

В приложениях 4 и 5 на примере конкурсной задачи Tinkoff Data Science Challenge и конкурсной задачи предсказания отклика ОТП Банка детально показаны этапы предварительной обработки данных, в частности приведение переменных к нужным типам, обработка редких категорий, импутация пропусков, конструирование признаков, также освещаются специальные процедуры предварительной обработки данных, позволяющие улучшить модель логистической регрессии.

Приложение 6 «Работа с датами и строками» посвящено таким задачам, как правильный парсинг дат различного формата, изменение регистра букв в строках, удаление лишних символов из строк, извлечение нужных символов из строк.

Приложение 7 «Работа с предупреждением `SettingWithCopyWarning` в библиотеке `pandas`» посвящено причинам появления предупреждения `SettingWithCopyWarning` и способам его устранения.

В приложении 8 «От `pandas` к `scikit-learn` – новый подход к управлению рабочими процессами» внимание уделено работе с конвейерами.

Что необходимо для чтения этой книги

Эта книга предполагает некоторое знакомство с принципами программирования, но те, кто не имеет опыта программирования или, в частности, опыта программирования на языке Python, будут довольны примерами, поскольку они в большей степени сосредоточены на конструкциях библиотеки `pandas`, нежели на языке Python или программировании вообще. Примеры приводятся для `Anaconda 5.2` для Python 3.6 и `pandas 0.23`. Если вы не установили их, в главе 2 «Запуск библиотеки `pandas`» дается инструкция относительно установки `pandas` в системах Windows, OSX и Ubuntu.

На кого рассчитана эта книга

Эта книга идеально подходит для специалистов по работе с данными, аналитиков данных и программистов Python, которые хотят погрузиться в анализ данных с использованием библиотеки `pandas`, и всех, кто интересуется анализом данных. Некоторые познания в области статистики и программирования помогут вам извлечь максимальную пользу из этой книги, но они не обязательны. Предварительный опыт работы с `pandas` также не требуется.

Соглашения

В данной книге используется несколько стилей текста для разных видов информации. Вот несколько примеров этих стилей и объяснение их смысла.

Слова в тексте, обозначающие объекты, функции, методы и другие элементы программного кода, отображаются так: «Эту информацию можно легко импортировать в датафрейм с помощью функции `pd.read_csv()` следующим образом».

Блок программного кода, введенный в интерпретаторе Python, отображается следующим образом:

```
import pandas as pd
df = pd.DataFrame.from_items([('column1', [1, 2, 3])])
print(df)
```

Любой ввод или вывод в командной строке записывается таким образом:

```
mh@ubuntu:~/Downloads$ chmod +x Anaconda-2.1.0-Linux-x86_64.sh
mh@ubuntu:~/Downloads$ ./Anaconda-2.1.0-Linux-x86_64.sh
```

Новые термины и важные слова выделены жирным шрифтом.



Предупреждения и важные примечания выглядят так.



Советы и рекомендации выглядят так.

ОТЗЫВЫ И ПОЖЕЛАНИЯ

Мы всегда рады отзывам наших читателей. Расскажите нам, что вы думаете об этой книге – что понравилось или, может быть, не понравилось. Отзывы важны для нас, чтобы выпускать книги, которые будут для вас максимально полезны.

Вы можете написать отзыв прямо на нашем сайте www.dmkpress.com, зайдя на страницу книги, и оставить комментарий в разделе «Отзывы и рецензии». Также можно послать письмо главному редактору по адресу dmkpress@gmail.com, при этом напишите название книги в теме письма.

Если есть тема, в которой вы квалифицированы, и вы заинтересованы в написании новой книги, заполните форму на нашем сайте по адресу http://dmkpress.com/authors/publish_book/ или напишите в издательство по адресу dmkpress@gmail.com.

СКАЧИВАНИЕ ИСХОДНОГО КОДА ПРИМЕРОВ

Скачать файлы с дополнительной информацией для книг издательства «ДМК Пресс» можно на сайте www.dmkpress.com или www.дмк.рф на странице с описанием соответствующей книги.

СПИСОК ОПЕЧАТОК

Хотя мы приняли все возможные меры для того, чтобы удостовериться в качестве наших текстов, ошибки все равно случаются. Если вы найдете ошибку в одной из наших книг – возможно, ошибку в тексте или в коде, – мы будем очень благодарны, если вы сообщите нам о ней. Сделав это, вы избавите других читателей от расстройств и поможете нам улучшить последующие версии этой книги.

Если вы найдете какие-либо ошибки в коде, пожалуйста, сообщите о них главному редактору по адресу dmkpress@gmail.com, и мы исправим это в следующих тиражах.

НАРУШЕНИЕ АВТОРСКИХ ПРАВ

Пиратство в интернете по-прежнему остается насущной проблемой. Издательства «ДМК Пресс» и Раскт очень серьезно относятся к вопросам защиты авторских прав и лицензирования. Если вы столкнетесь в интернете с незаконно выполненной копией любой нашей книги, пожалуйста, сообщите нам адрес копии или веб-сайта, чтобы мы могли применить санкции.

Пожалуйста, свяжитесь с нами по адресу электронной почты dmkpress@gmail.com со ссылкой на подозрительные материалы.

Мы высоко ценим любую помощь по защите наших авторов, помогающую нам предоставлять вам качественные материалы.

Глава 1

Библиотека pandas и анализ данных

Добро пожаловать на страницы книги «Изучаем pandas»! В этой книге мы отправимся в путешествие, в ходе которого вы научитесь работать с pandas, библиотекой анализа данных с открытым исходным кодом, предназначенной для языка программирования Python. Библиотека pandas предлагает высокопроизводительные и простые в использовании структуры данных и инструменты анализа, созданные с помощью языка Python. Библиотека pandas привнесла в Python массу полезных инструментов, взяв их из языка статистического программирования R, в частности объекты **data frame** (**датафрейм**), пакеты R, например `plyr` и `reshape2`, и разместила их в одной библиотеке, которую вы можете использовать в среде Python.

В первой главе мы посвятим время базовому знакомству с библиотекой pandas и тому, как она вписывается в обширную картину анализа данных. Вместо того чтобы полностью сосредоточиться на конкретных аспектах использования библиотеки pandas, эта глава призвана дать читателю ощущение своего места в обширной картине анализа данных. Цель состоит в том, чтобы при изучении библиотеки pandas вы также узнали о том, зачем нужны все эти различные функции, выполняющие задачи анализа данных.

Итак, давайте начнем. В этой главе мы рассмотрим:

- что из себя представляет библиотека pandas, почему она была создана и что она вам даст;
- как библиотека pandas связана с анализом данных и наукой о данных;
- этапы анализа данных и их поддержку в библиотеке pandas;
- общие понятия «данные» и «аналитика»;
- основные понятия анализа данных и статистического анализа;
- типы данных и их использование в библиотеке pandas;
- другие библиотеки в экосистеме Python, которые вы, вероятно, будете использовать вместе с pandas.

Знакомство с библиотекой PANDAS

Библиотека pandas – это библиотека Python, содержащая высокоуровневые структуры данных и инструменты, которые были созданы, чтобы помочь программистам Python осуществить полноценный анализ данных. Конечная цель

библиотеки pandas заключается в том, чтобы помочь вам быстро найти необходимую информацию, скрытую в данных, при этом информацию содержательно-го характера.

Разработка библиотеки pandas была начата в 2008 году Уэсом Маккинни и представлена в 2009 году как проект с открытым исходным программным кодом. В настоящее время библиотека pandas активно курируется и разрабатывается различными организациями и участниками.

Первоначально библиотека pandas предназначалась для применения в финансах, в частности благодаря ее возможностям работы с временными рядами и обработке исторической информации об акциях. Обработка финансовой информации сопряжена с массой проблем, вот некоторые из них:

- обработка данных (например, данных о котировках акций), меняющихся с течением времени;
- необходимость единого стандарта измерений нескольких потоков данных в один и тот же период времени;
- определение взаимосвязи (корреляции) между двумя и более потоками данных;
- представление дат и времени в качестве объектов первого класса;
- увеличение или уменьшение шага дискретизации временного ряда.

Для выполнения этих операций необходим инструмент, который позволяет нам извлекать, индексировать, очищать и приводить в порядок, изменять и объединять данные, создавать срезы данных и выполнять различные виды анализа как для одномерных, так и для многомерных данных, включая данные разного типа, которые автоматически выравниваются по набору общих индексных меток. И вот как раз здесь на помощь приходит библиотека pandas, которая предлагает множество полезных и мощных функций, например:

- быстрые и эффективные объекты Series и DataFrame для обработки данных с помощью встроеной индексации;
- интеллектуальное выравнивание данных с помощью индексов и меток;
- интегрированная обработка пропущенных данных;
- инструменты для приведения данных в порядок;
- встроенные инструменты для чтения и записи данных для обмена между объектами Series и DataFrame в памяти, файлами, базами данных и веб-службами;
- возможность обработки данных, хранящихся в различных популярных форматах, таких как CSV, Excel, HDF5 и JSON;
- изменение формы и поворот данных;
- интеллектуальное создание срезов данных на основе меток, сложная индексация и отбор из больших наборов данных подмножеств по определенному критерию;
- удаление и вставка столбцов из объектов Series и DataFrame для изменения размера;
- агрегирование или преобразование данных с помощью мощного инструмента «разделение – применение – объединение»;
- иерархическая индексация, облегчающая работу с высокоразмерными данными в низкоразмерной структуре данных;
- высокопроизводительное слияние и соединение наборов данных;

- разнообразные функции для работы с временными рядами, включая создание диапазона дат и преобразование частоты временного ряда, вычисление скользящих статистик, скользящих линейных регрессий, смещение дат и сдвиг временного ряда с запаздыванием;
- оптимизация для достижения более высокой производительности, включающая программный код, написанный на Cython или C.

Мощный набор функций в сочетании с бесшовной интеграцией с Python и другими инструментами экосистемы Python позволил библиотеке pandas найти широкое применение во многих областях. Она используется в самых разных академических и коммерческих областях, включая финансы, нейробиологию, экономику, статистику, рекламу и веб-аналитику. Она стала одним из наиболее предпочтительных инструментов для специалистов по работе с данными.

Python долгое время широко использовался для сбора данных и подготовки, но при этом в меньшей степени был предназначен для анализа данных и моделирования. Библиотека pandas помогает заполнить этот пробел, позволяя вам выполнить весь рабочий процесс анализа данных в среде Python, не переходя на такой более специализированный язык, как R. Это очень важно, поскольку люди, знакомые с языком Python, являющимся более универсальным языком программирования, чем R (язык, ориентированный в большей степени на статистиков), получают в свое распоряжение массу функций по представлению и обработке данных, имеющихся в R, и при этом полностью остаются в невероятно богатой экосистеме Python.

В сочетании с IPython, тетрадками Jupyter и широким выбором других библиотек среда Python в плане выполнения анализа данных превосходит по производительности, эффективности и возможности совместной работы многие другие инструменты. Все это привело к тому, что многие пользователи широко применяют библиотеку pandas в самых различных отраслях.

ОБРАБОТКА ДАННЫХ, АНАЛИЗ, НАУКА И БИБЛИОТЕКА PANDAS

Мы живем в мире, в котором каждый день генерируются и записываются огромные объемы данных. Эти данные поступают из множества информационных систем, устройств и датчиков. Почти все, что вы делаете, и все то, что вы используете в рамках своей деятельности, генерирует данные, которые можно собрать или которые уже собраны.

Во многом такая ситуация была обусловлена повсеместным распространением услуг, связанных с информационными сетями, и стремительно возросшими возможностями хранения данных. Все это в сочетании с постоянно снижающейся стоимостью хранения повысило эффективность сбора и хранения даже самых тривиальных данных.

В итоге было накоплено огромное количество данных, готовые для загрузки. Но эти данные сосредоточились по всему киберпространству, и на самом деле их еще нельзя назвать **информацией (information)**. Данные представляют собой коллекции зарегистрированных событий, будь то финансовые данные, или ваше общение в социальных сетях, или ваш персональный трекер здоровья, отслеживающий сердцебиение в течение дня. Эти данные хранятся в различных форматах, расположены в разных местах, и исследование сырых данных может дать ценную

информацию.

Логично, что весь процесс можно разбить на три основные дисциплины:

- обработка данных;
- анализ данных;
- наука о данных.

Эти три дисциплины часто пересекаются. Вопросы о том, где заканчивается одна дисциплина и начинается другая, остаются открытыми. В следующих разделах мы дадим определения этим дисциплинам.

Обработка данных

Данные разбросаны по всей планете. Они хранятся в разных форматах. Они имеют разный уровень качества. Поэтому существует потребность в инструментах и процессах сбора данных, дающих такое представление данных, которое можно использовать для принятия решений. Инструмент, который используется для работы с данными на этапе подготовки к анализу, должен уметь решать различные задачи. Функционал такого инструмента включает в себя:

- программируемость для повторного использования и совместного использования;
- загрузку данных из внешних источников;
- локальное сохранение данных;
- индексацию данных для их эффективного извлечения;
- выравнивание данных в разных наборах на основе атрибутов;
- объединение данных, расположенных в разных наборах;
- преобразование данных в другое представление;
- очистку данных от «мусора»;
- эффективную обработку загрязненных данных;
- группировку данных;
- агрегирование данных по схожим характеристикам;
- применение функций, вычисляющих среднее или выполняющих преобразования;
- выполнение запросов или создание срезов для исследования подмножеств данных;
- изменение формы данных;
- создание отдельных категорий данных;
- изменение частоты временного ряда.

Существует масса инструментов для обработки данных. Каждый из них отличается поддержкой элементов этого списка, способами их развертывания и способами их использования. Эти инструменты включают в себя реляционные базы данных (SQL Server, Oracle), электронные таблицы (Excel), системы обработки событий (такие как Spark) и более общие инструменты, такие как R и pandas.

Анализ данных

Анализ данных – это процесс извлечения смысла из данных. Данные, представленные в количественном виде, часто называют **информацией (information)**. Анализ данных – это процесс получения информации из данных путем создания моделей и применения математического аппарата для поиска закономерностей.

Он часто переключается с обработкой данных, и не всегда можно четко провести различие между ними. Многие инструменты обработки данных также содержат аналитические функции, а инструменты анализа данных нередко предлагают возможности обработки данных.

Наука о данных

Наука о данных – это процесс использования статистики и анализа данных для понимания **явлений (phenomena)**, скрытых в данных. Наука о данных обычно начинается с информации и применяет к ней более сложный анализ на основе знаний, относящихся к разным предметным областям. К этим предметным областям относятся математика, статистика, информатика, компьютерные науки, машинное обучение, классификация, кластерный анализ, интеллектуальный анализ данных, базы данных и визуализация. Наука о данных носит многодисциплинарный характер. Ее методы анализа могут сильно отличаться друг от друга и зависеть от конкретной предметной области.

Предназначение библиотеки pandas

В первую очередь библиотека pandas – превосходный инструмент обработки данными. Все потребности, описанные ранее, будут рассмотрены в этой книге с использованием библиотеки pandas. На решение этих задач и направлен основной функционал библиотеки pandas, и именно на решении большей части таковых задач мы и сосредоточимся в этой книге.

Стоит отметить, что основное предназначение библиотеки pandas – это подготовка данных. Однако библиотека pandas также предоставляет несколько функций для выполнения анализа данных. Эти возможности подразумевают вычисление описательных статистик и функций, необходимых для финансового анализа, например вычисления корреляции.

Поэтому сама по себе библиотека pandas не является инструментом для научных исследований. Это скорее инструмент обработки данных с некоторыми возможностями анализа. Библиотека pandas явно оставляет за скобками сложный статистический, финансовый анализ, предлагая его выполнить другим библиотекам Python, таким как SciPy, NumPy, scikit-learn, а для визуализации данных использует графические библиотеки, например **matplotlib** и **ggvis**.

Мы сосредоточимся на преимуществе библиотеки pandas над другими языками, такими как R, поскольку приложения на основе библиотеки pandas могут использовать обширную сеть надежных фреймворков Python, уже созданных и протестированных Python-сообществом.

ПРОЦЕСС АНАЛИЗА ДАННЫХ

Основная цель этой книги – научить вас использовать библиотеку pandas для обработки данных. Однако есть второстепенная и, возможно, не менее важная цель – показать, как библиотека pandas встроена в те процессы, которые специалист по работе с данными/аналитик данных выполняет в повседневной жизни.

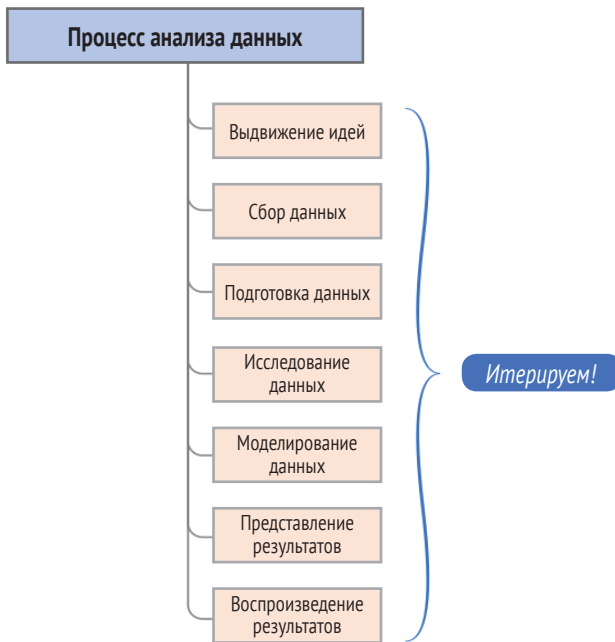
Описание шагов, связанных с процессом анализа данных, дается на веб-странице библиотеки pandas:

- обработка и очистка данных;
- анализ/моделирование данных;
- преобразование данных в удобную форму.


Этот небольшой список является хорошим исходным определением, но он не охватывает процесса анализа в целом и не может ответить на вопрос, почему в библиотеке pandas реализовано так много инструментов. В следующем разделе данный процесс будет рассмотрен подробнее.

Процесс

Предлагаемый процесс – это последовательность операций, которую называют процессом анализа данных, и он представлен на следующей диаграмме:



Эта схема задает структуру для определения логических шагов, которые будут предприняты в ходе работы с данными. Сейчас давайте кратко рассмотрим каждый этап этого процесса, а также некоторые задачи, которые вы как аналитик данных будете выполнять, используя библиотеку pandas.

 Важно понимать, что это не простой линейный процесс. Лучше всего этот процесс осуществлять интерактивно и гибким/итеративным способом.

Выдвижение идей

Первый этап процесса анализа данных – определить, что вы хотите выяснить. Он называется **выдвижение идей (ideation)**, мы формулируем идеи по поводу того, что мы хотим сделать и доказать. Идея в целом связана с гипотезой о наличии тех или иных структур в данных, которые можно использовать для принятия решений.

Эти решения часто принимаются в контексте бизнеса, а также в других дисциплинах, например в научно-исследовательской деятельности. Сейчас модно применять анализ данных для решения бизнес-задач, поскольку глубокое понимание данных может повысить прибыль компаний.

Однако какое решение мы обычно хотим принять? Ниже приводятся некоторые наиболее часто задаваемые вопросы:

- почему это происходит?
- можем ли мы предсказать будущее с использованием исторических данных?
- как я могу оптимизировать бизнес-операции в будущем?

Этот список ни в коем случае не является исчерпывающим, но он охватывает значительную часть причин, по которым проводится анализ данных. Чтобы получить ответы на эти вопросы, нужно осуществить сбор и анализ данных, связанных с решаемой проблемой. Необходимо определить, какие данные мы будем изучать, какова польза от исследования, как будут получены данные, каковы критерии успешности решения и как в конечном итоге информация будет представлена.

Сама по себе библиотека `pandas` не предлагает инструментов, позволяющих формулировать идеи. Но как только вы научитесь использовать библиотеку `pandas`, вы, естественно, поймете, как `pandas` может помочь вам в формулировании идей. Ведь к этому моменту вы вооружитесь мощным инструментом, который можно использовать для выдвижения множества сложных гипотез.

Сбор данных

Как только вы сформулировали идею, вы должны получить данные, чтобы попытаться проверить свою гипотезу. Эти данные могут быть собраны вашей организацией или получены от внешних поставщиков данных. Эти данные обычно поставляются в виде архивных данных или могут быть предоставлены в режиме реального времени (хотя `pandas` не так широко известен, как инструмент обработки данных в режиме реального времени).

Часто данные являются «сырыми», даже если они взяты из источников данных, созданных вами или вашей организацией. Понятие «сырые» означает, что данные могут быть неправильно сформированы, записаны в разных форматах или содержать ошибки. Они могут быть неполными, и может потребоваться аугментация данных¹.

В мире имеется множество бесплатных данных. Многие данные не являются бесплатными и фактически требуют значительных денежных сумм. Некоторые из них есть в свободном доступе и имеют публичные API-интерфейсы, а другие можно получить по подписке. Часто данные, которые вы покупаете, являются более чистыми, но это не всегда так.

В любом случае, библиотека `pandas` предлагает надежный и простой в использовании набор инструментов для сбора данных из разных источников, при этом данные могут иметь разный формат. Кроме того, библиотека `pandas` позволяет не только извлекать данные, но и выполнять первоначальную структуризацию данных с помощью объектов `Series` и `DataFrame`, не прибегая к написанию сложного

¹ Аугментация данных (data augmentation) – это создание дополнительных обучающих данных из имеющихся данных. – *Прим. перев.*

программного кода, который может потребоваться в других инструментах или языках программирования.

Подготовка данных

На этом этапе мы подготавливаем исходные данные для проведения дальнейшего исследования. Часто подготовка является весьма интересным этапом анализа данных. Как правило, все проблемы, связанные с данными, связаны с качеством этих данных. Вы, вероятно, потратите много времени на устранение проблем, связанных с качеством данных.

Зачем? Ну, есть ряд причин:

- данные просто некорректны;
- какая-то информация в наборе данных отсутствует;
- данные записаны в единицах измерения, не подходящих для вашего анализа;
- данные записаны в форматах, не удобных для вашего анализа;
- данные находятся на том уровне детализации, который не подходит для вашего анализа;
- не все интересующие вас поля можно получить из одного источника;
- представление данных отличается в зависимости от поставщика.

Процесс подготовки данных сосредоточен на решении вышеприведенных проблем. Библиотека pandas предлагает множество отличных возможностей для процесса подготовки данных, который часто называют **приведением данных в порядок (data tidying)**. Ее арсенал включает в себя интеллектуальные средства обработки пропущенных данных, преобразование типов данных, преобразование форматов, изменение частоты измерений, объединение данных, расположенных в нескольких наборах данных, кодировку/преобразование символов и группировку данных среди прочих операций. Мы подробно рассмотрим все эти операции.

Исследование данных

Исследование данных подразумевает изучение ваших данных, чтобы попытаться открыть какие-то закономерности в данных. Исследование может включать в себя различные задачи, например:

- изучение того, как переменные связаны друг с другом;
- анализ распределения данных;
- обнаружение и исключение выбросов;
- создание быстрых визуализаций;
- быстрое создание новых представлений данных или моделей для последующего использования в ходе моделирования.

Исследование данных является одной из самых сильных сторон библиотеки pandas. Хотя большинство языков программирования также позволяет выполнить исследование данных, каждый из них требует подготовительных действий, не связанных непосредственно с исследованием данных.

Используя принцип **read-eval-print-loop (REPL)**, реализованный в IPython и/или тетрадках Jupyter, библиотека pandas создает исследовательскую среду, в которой объем этих подготовительных действий сведен к минимуму. Выразительность синтаксиса библиотеки pandas позволяет кратко описать сложные операции с данными, а немедленно появляющийся результат выполненной операции

помогает быстро проверить правильность ваших действий, не прибегая к повторной компиляции и полному переписыванию своего программного кода.

Моделирование данных

На этапе моделирования вы формализуете свои открытия, найденные в ходе исследования данных. Для этого закономерностям, обнаруженным в данных и позволяющим извлечь из данных содержательный смысл, нужно дать ясную интерпретацию. Итогом становится **модель (model)**, позволяющая выразить обнаруженные закономерности в математическом виде и содержащая программный код, позволяющий преобразовать данные в полезную информацию.

Процесс моделирования является итеративным, когда с помощью исследования данных вы выбираете переменные, необходимые для вашего анализа, подаете переменные на вход модели, строите модель и определяете, насколько хорошо модель подтверждает ваши исходные предположения. Он может включать формальное моделирование структуры данных, а также может сочетать методы из различных областей, например из статистики, машинного обучения и исследования операций.

Библиотека `pandas` предлагает мощные инструменты для моделирования данных. Именно на этом этапе вы формализуете модель данных, используя объекты `DataFrame` и стараясь сделать процесс построения модели максимально компактным. Кроме того, вы можете воспользоваться всеми возможностями языка `Python`, чтобы полностью автоматизировать процесс создания модели.

С аналитической точки зрения библиотека `pandas` предлагает в первую очередь интегрированную поддержку описательной статистики, которая понадобится вам при решении разных задач. Впрочем, библиотека `pandas` использует язык `Python`, и поэтому если вам нужны более продвинутые аналитические возможности, вы очень легко можете интегрировать `pandas` с другими библиотеками обширной научной среды `Python`.

Представление результатов

Предпоследний этап процесса анализа данных – представление результатов, как правило, в виде отчета или презентации. Вам нужно дать убедительное и подробное объяснение своего решения. Часто это можно сделать с помощью различных инструментов визуализации в `Python` и затем вручную создать презентацию.

Тетрадки (ноутбуки) Jupyter – это мощный инструмент для создания презентаций по итогам вашего анализа. Эти тетрадки представляют собой инструмент как для выполнения кода, так и для подробного аннотирования выполняемых действий. Их можно использовать для создания высокоэффективных презентаций, включающих фрагменты программного кода, стилизованный текст и графику.



Мы кратко рассмотрим тетрадки `Jupyter` в главе 2 «Запуск библиотеки `pandas`».

Воспроизведение результатов анализа

Важный этап процесса анализа данных – обмен и воспроизведение результатов. Часто говорят, что если другие исследователи не смогут воспроизвести ваш эксперимент и полученные результаты, то вы ничего не доказали.

К счастью для вас, вы легко сможете сделать свой анализ воспроизводимым, воспользовавшись библиотекой pandas и Python. Это можно сделать, поделившись с другими исследователями программным кодом Python, который лежит в основе программного кода библиотеки pandas, а также данными.

Тетрадки Jupyter также являются удобным инструментом для хранения программного кода и проектов, которыми можно легко поделиться со всеми, у кого установлен Jupyter Notebook. В интернете существует много бесплатных и безопасных сайтов обмена, которые позволяют вам создавать или размещать свои тетрадки Jupyter для совместного использования.

По поводу итеративности и гибкости

Очень важно понимать, что обработка данных, анализ и наука – это итеративный процесс. Несмотря на то что ранее мы рассмотрели этапы, выстроив прямую последовательность, в конечном итоге вы будете двигаться как в прямом, так и в обратном порядке. Например, на этапе исследования вы можете выявить аномалии в данных, которые связаны с проблемами чистоты данных, они решаются на этапе подготовки, и поэтому необходимо вернуться, чтобы исправить эти проблемы.

Это обычная ситуация, возникающая в процессе анализа данных. Вы словно совершаете путешествие, пытаетесь решить свою первоначальную задачу, все время получая новое, дополнительное представление о данных, с которыми работаете. Эта информация, возможно, побудит вас задать новые вопросы, более точные вопросы или осознать, что ваши первоначальные вопросы не являлись теми вопросами, на которые вы хотели бы найти ответы. Процесс анализа данных – это и в самом деле путешествие, а не конкретный пункт назначения.

ВЗАИМОСВЯЗЬ МЕЖДУ КНИГОЙ И ПРОЦЕССОМ АНАЛИЗА ДАННЫХ

Ниже приведено краткое описание этапов анализа данных, о которых вы узнаете в этой книге. Не волнуйтесь, что последовательность этапов не совпадает с последовательностью изложения материала в книге. Книга в логическом порядке знакомит вас с библиотекой pandas, и вы можете вернуться к любой главе, которая соответствует интересующему вас этапу анализа данных.

Этап процесса анализа данных	Глава книги
Выдвижение идей	Выдвижение идей – это самая творческая часть науки о данных. У вас должна быть идея. Она должна подтверждаться фактами (здесь вам понадобятся данные), и должна быть возможность наблюдать эти факты после проведения анализа (на основе фактов из прошлого строим модель, прогнозируя их в будущем)
Сбор данных	Сбор данных главным образом освещается в главе «Загрузка данных»
Подготовка данных	Подготовка данных преимущественно освещается в главе 10 «Приведение данных в порядок», но это довольно общая тема, которая встречается в большинстве глав данной книги
Исследование данных	Тема исследования данных охватывает материал, начинающий с главы 3 «Представление одномерных данных с помощью объекта Series» и заканчивающийся главой 14 «Визуализация». Таким образом, большая часть книги посвящена исследованию данных. Но наиболее подробное освещение данного этапа дается в главе 14 «Визуализация»

Этап процесса анализа данных	Глава книги
Моделирование данных	Моделирование данных подробно освещается в главе 3 «Представление одномерных данных с помощью объекта Series» и главе 4 «Представление табличных и многомерных данных с помощью объекта DataFrame», а также в главе 11 «Объединение, связывание и изменение формы данных» и главе 13 «Анализ временных рядов»
Представление результатов	Представление результатов – это главная тема главы 14 «Визуализация»
Воспроизведение результатов	Тема воспроизведения результатов проходит по всей книге, поскольку примеры представлены в виде тетрадок Jupyter. Работая в тетрадках, вы по умолчанию используете инструмент для воспроизведения результатов, и у вас есть возможность поделиться этими тетрадками различными способами

ПОНЯТИЯ «ДААННЫЕ» И «АНАЛИЗ» В КОНТЕКСТЕ НАШЕГО ЗНАКОМСТВА С БИБЛИОТЕКОЙ PANDAS

По мере изучения библиотеки `pandas` и анализа данных вы будете сталкиваться с разными понятиями из области данных, моделирования и анализа. Давайте рассмотрим некоторые из этих понятий и выясним, как они связаны с библиотекой `pandas`.

Типы данных

При работе с сырыми данными вы столкнетесь с несколькими типами данных, которые нужно будет представить с помощью объектов `pandas`. Они важны для понимания, поскольку инструменты, необходимые для работы с каждым типом, отличаются друг от друга.

Библиотека `pandas`, по сути, используется для обработки структурированных данных. Однако, помимо этого, она предлагает несколько инструментов, позволяющих преобразовать неструктурированные данные в такое представление, которое мы уже сможем обработать.

Структурированные данные

Структурированные данные – это тип данных, при котором последние упорядочены в вертикальные столбцы (поля) и горизонтальные строки (записи или наблюдения). Примером таких данных являются данные в реляционных базах и электронных таблицах. Структурированные данные зависят от модели данных, которая представляет собой определенную структуру, содержательного смысла данных и часто от способа обработки данных. Сюда входят идентификация типа данных (целое число, число с плавающей точкой, строка и т. д.) и ограничения, накладываемые на данные, например количество символов, максимальное и минимальное значения, или ограничение на определенный набор значений.

Структурированные данные – это тип данных, который предназначен для использования в библиотеке `pandas`. Как мы увидим сначала на примере объекта `Series`, а затем на примере объекта `DataFrame`, `pandas` помещает структурированные данные в один или несколько столбцов данных (при этом каждый столбец имеет один конкретный тип данных) и одну или несколько строк данных.

Неструктурированные данные

Неструктурированные данные – это данные, которые не имеют определенной структуры, не предполагают наличия заранее определенных столбцов определенного типа. Примером неструктурированных данных могут быть такие виды информации, как фотографии и графические изображения, видеоролики, потоковые данные датчиков, веб-страницы, PDF-файлы, презентации PowerPoint, электронные письма, записи в блогах, страницы Википедии и документы Word.

Хотя библиотека pandas не может обработать неструктурированные данные напрямую, она предлагает ряд инструментов для сбора структурированных данных из неструктурированных источников. В качестве конкретного примера мы рассмотрим инструменты pandas, позволяющие извлекать веб-страницы, определенные фрагменты контента и записывать в объект DataFrame.

Полуструктурированные данные

Полуструктурированные данные занимают промежуточное положение между неструктурированными и структурированными данными. Их можно рассматривать как тип структурированных данных, у которых нет строгой структуры моделей данных. Формат JSON – это пример полуструктурированных данных. Хотя хороший JSON-файл будет иметь определенный формат, здесь нет определенной схемы данных, которая всегда строго соблюдается. В большинстве случаев данные будут записаны в виде определенной схемы, которую легко преобразовать в структурированный тип данных, например в объект DataFrame, но для этого, возможно, потребуется задать определенный тип данных.

Переменные

Моделируя данные в библиотеке pandas, мы будем исследовать одну или несколько переменных и искать статистический смысл, анализируя значения этой переменной или значения нескольких переменных. При этом термин «переменная» не тождествен понятию «переменная» в программировании, здесь речь идет о статистических переменных.

Переменной является любая характеристика, число или количество, которое можно измерить или сосчитать. Название «переменная» обусловлено тем, что значение характеристики может меняться от наблюдения к наблюдению, а также может меняться со временем. Значение цены акции, возраст, пол, доходы и расходы в бизнесе, страна рождения, капитальные затраты, школьные оценки, цвет глаз и тип транспортного средства являются примерами переменных.

Существует несколько общих типов статистических переменных, с которыми мы столкнемся при работе с библиотекой pandas:

- категориальные переменные;
- непрерывные переменные.

Категориальные переменные

Категориальная переменная (categorical variable) – это переменная, которая может принимать одно значение из ограниченного и обычно фиксированного набора возможных значений. Каждое из возможных значений часто называется **уровнем (level)**. Примерами категориальных переменных являются пол, штат, социальный класс, группа крови, гражданство, время наблюдения или рейтинг, например шкала Лайкерта. Все категориальные переменные являются качествен-

ными характеристиками, которые могут описать продукт, но при этом не измеряются в непрерывной шкале. Например, возьмем переменную *Штат*. Допустим, он имеет уровни Аризона, Северная Каролина и Висконсин. Человек либо проживает в штате Аризона, либо в штате Северная Каролина, либо в штате Висконсин. Не существует золотой середины между штатами, нельзя вычислить среднее значение штата, и нет естественного способа упорядочить эти категории, нельзя сказать, что штат Аризона хуже/лучше штата Северная Каролина, а штат Северная Каролина хуже/лучше штата Висконсин. Категориальные переменные в библиотеке *pandas* представлены объектами *Categorical*, специальным типом данных, который соответствует категориальным переменным в статистике.

Непрерывные переменные

Непрерывная переменная (continuous variable) – это переменная, которая может принимать бесконечное (неисчислимо) количество значений. Все непрерывные переменные являются характеристиками, которые количественно описывают продукт и измеряются в непрерывной шкале. Примерами непрерывных переменных являются возраст, высота, время и температура. Мы можем вычислить средний возраст, среднюю высоту, среднее время и среднюю температуру. Мы можем упорядочить значения. 20-летний младше 30-летнего, а 30-летний младше 40-летнего. Непрерывные переменные в библиотеке *pandas* представлены либо типом *float*, либо типом *integer* (нативными питоновскими типами данных).

Временные ряды

В библиотеке *pandas* временные ряды – это объект первого класса. Время добавляет важное дополнительное измерение. Часто переменные не зависят от времени их регистрации, то есть момент времени, в который их фиксировали, не имеет значения. Но во многих случаях это не так. Временной ряд представляет собой переменную, у которой значения зарегистрированы в определенные временные интервалы, таким образом, наблюдения изначально упорядочены по времени.

Стохастическая модель для временного ряда обычно отражает тот факт, что наблюдения, близкие друг к другу во времени, будут более тесно связаны, чем наблюдения, которые далеко отстоят друг от друга по времени. Модели временных рядов часто используют естественное одностороннее упорядочение времени, поэтому значения для заданного периода будут выражаться как полученные из прошлого, а не из будущего.

Чаще всего в библиотеке *pandas* работают с финансовыми данными, где переменная представляет собой стоимость акции, изменяющуюся через равные интервалы времени в течение дня. Нам часто нужно определить скорость изменения цены в определенные интервалы времени. Кроме того, мы можем скорректировать цены нескольких акций по определенным интервалам времени.

Это настолько важная и мощная возможность библиотеки *pandas*, что мы посвятим ей целую главу.

Общие понятия анализа и статистики

В этой книге мы только коснемся общих понятий статистики и технической стороны анализа данных. Однако стоит раскрыть несколько понятий, некоторые из которых реализованы непосредственно в библиотеке *pandas*. Остальные понятия

связаны с другими библиотеками, например с библиотекой SciPy, однако вы также можете встретить их во время работы с библиотекой pandas, поэтому мы тоже о них расскажем.

Количественный и качественный анализ

Качественный анализ – это научное исследование данных, которые можно наблюдать, но нельзя количественно измерить. Основное внимание уделяется качественным характеристикам. Примерами качественных данных могут быть:

- мягкость кожи;
- изящность бега.

Количественный анализ – это исследование данных, которые можно измерить количественно. Примерами количественных данных могут быть:

- количество;
- цена;
- высота.

Библиотека pandas главным образом работает с количественными данными, предлагая разнообразные инструменты для их обработки. Библиотека pandas не предназначена для проведения качественного анализа, но позволяет вам представить качественную информацию в различных видах.

Одномерный и многомерный анализ

В каком-то смысле статистика представляет собой практику изучения переменных и, в частности, наблюдение за этими переменными. Статистика преимущественно опирается на анализ одной переменной, которая называется одномерным анализом. **Одномерный анализ (univariate analysis)** – это простейшая форма анализа данных. Он не имеет отношения к анализу причин или взаимосвязей и обычно используется для описания или подытоживания данных, а также поиска закономерностей в данных.

Многомерный анализ (multivariate analysis) – это метод моделирования, в котором участвуют две или более переменных, влияющих на результат эксперимента. Многомерный анализ часто связан с такими понятиями, как корреляция и регрессия, которые помогают нам понять взаимосвязь между несколькими переменными, а также влияние этой взаимосвязи на результат.

Библиотека pandas в основном предлагает инструменты для проведения одномерного анализа. Главным образом речь идет о вычислении описательных статистик, хотя можно вычислить такой показатель, как корреляция (поскольку ее часто используют в финансах и других областях).

Остальные более сложные статистические процедуры можно выполнить с помощью библиотеки StatsModels. Опять же это не недостаток библиотеки pandas, а конкретное решение, направленное на то, чтобы более сложные операции выполнять с помощью специализированных библиотек Python.

Описательные статистики

Описательные статистики – это показатели, которые подытоживают данные, обычно в тех случаях, где набор данных представляет собой генеральную совокупность или выборку из одной переменной (одномерные данные). Эти показатели описывают набор данных и являются мерами центральной тенденции, а также мерами изменчивости и дисперсии.

Например, следующие показатели являются описательными статистиками:

- распределение (например, нормальное, пуассоновское);
- центральная тенденция (например, среднее, медиана и мода);
- дисперсия (например, дисперсия, стандартное отклонение).

Позже мы увидим, что объекты `Series` и `DataFrame` поддерживают вычисление большей части описательных статистик.

Индуктивные статистики

Индуктивные статистики отличаются от описательных тем, что с их помощью мы пытаемся сделать выводы о данных, а не просто подытожить данные. Примерами индуктивных статистик являются:

- t-тест;
- хи-квадрат;
- ANOVA;
- бутстреп.

Эти индуктивные методы были перенесены из библиотеки `pandas` в другие инструменты типа `SciPy` и `StatsModels`.

Стохастические модели

Стохастические модели представляют собой вид статистического моделирования, который включает в себя одну или несколько случайных величин, а также подразумевает использование временных рядов. Цель стохастической модели – оценить вероятность того, что результат будет находиться в пределах определенного интервала, и спрогнозировать условия для разных ситуаций.

Примером стохастического моделирования являются симуляции Монте-Карло. Симуляции Монте-Карло применяются для расчета стоимости финансовых портфелей, которые подвержены влиянию различных риск-факторов и имеют нетривиальные распределения доходности.

Библиотека `pandas` предлагает фундаментальную структуру данных `DataFrame`, ее часто используют при работе с временными рядами, чтобы на их основе потом создать и запустить стохастическую модель. Хотя можно создать свои собственные стохастические модели и провести анализ с помощью `pandas` и `Python`, во многих случаях более удобным инструментом для стохастического моделирования будут специализированные библиотеки типа `RyMC`.

Вероятность и байесовская статистика

Байесовская статистика – это подход к статистическому оцениванию, основанный на теореме Байеса, математическом уравнении, использующем простые аксиомы теории вероятностей. Она позволяет аналитику вычислить любую условную вероятность интересующего события. Условная вероятность – это просто вероятность события A , при условии что наступило событие B .

Поэтому если использовать термины теории вероятностей, события уже произошли и были зафиксированы (поскольку мы знаем вероятность). Используя теорему Байеса, мы можем вычислить вероятность различных интересующих нас событий, учитывая уже имеющиеся данные.

Байесовское моделирование выходит за рамки этой книги, но опять же соответствующие модели данных прекрасно обрабатываются с помощью библиотеки `pandas`, и затем их можно анализировать с помощью таких библиотек, как `RyMC`.